

# Least Squares Fitting of Data by Linear or Quadratic Structures

David Eberly, Geometric Tools, Redmond WA 98052

<https://www.geometrictools.com/>

This work is licensed under the Creative Commons Attribution 4.0 International License. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/> or send a letter to Creative Commons, PO Box 1866, Mountain View, CA 94042, USA.

Created: July 15, 1999

Last Modified: October 30, 2018

## Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>The General Formulation for Nonlinear Least-Squares Fitting</b>	<b>4</b>
<b>3</b>	<b>Affine Fitting of Points Using Height Fields</b>	<b>5</b>
3.1	Fitting by a Line in 2 Dimensions . . . . .	5
3.1.1	Pseudocode for Fitting by a Line . . . . .	6
3.2	Fitting by a Plane in 3 Dimensions . . . . .	7
3.2.1	Pseudocode for Fitting by a Plane . . . . .	8
3.3	Fitting by a Hyperplane in $n + 1$ Dimensions . . . . .	9
3.3.1	Pseudocode for Fitting a Hyperplane . . . . .	10
<b>4</b>	<b>Affine Fitting of Points Using Orthogonal Regression</b>	<b>11</b>
4.1	Fitting by a Line [1 Dimension] . . . . .	11
4.1.1	Pseudocode for the General Case . . . . .	12
4.2	Fitting by a Hyperplane [ $(n - 1)$ Dimensions] . . . . .	13
4.2.1	Pseudocode for the General Case . . . . .	13
4.3	Fitting by a Flat [ $k$ Dimensions] . . . . .	14
4.3.1	Pseudocode for the General Case . . . . .	16
<b>5</b>	<b>Fitting a Hypersphere to Points</b>	<b>17</b>
5.1	Fitting Using Differences of Lengths and Radius . . . . .	17

5.1.1	Pseudocode for the General Case . . . . .	18
5.2	Fitting Using Differences of Squared Lengths and Squared Radius . . . . .	19
5.2.1	Pseudocode for the General Case . . . . .	20
5.2.2	Pseudocode for Circles . . . . .	21
5.2.3	Pseudocode for Spheres . . . . .	22
5.3	Fitting the Coefficients of a Quadratic Equation . . . . .	23
5.3.1	Pseudocode for the General Case . . . . .	25
5.3.2	Pseudocode for Circles . . . . .	26
5.3.3	Pseudocode for Spheres . . . . .	27
<b>6</b>	<b>Fitting a Hyperellipsoid to Points</b>	<b>29</b>
6.1	A Tractable Nonlinear Least-Squares Error Function . . . . .	29
6.1.1	Angle Derivative for an Ellipse . . . . .	31
6.1.2	Angle Derivatives for an Ellipsoid . . . . .	31
6.2	Initial Choice for the Parameters of the Error Function . . . . .	33
6.2.1	Initial Choice for the Parameters of an Ellipse . . . . .	33
6.2.2	Initial Choice for the Parameters of an Ellipsoid . . . . .	34
6.2.3	Initial Choice for the Parameters of a Hyperellipsoid . . . . .	35
6.3	Fitting Using Derivative-Based Methods . . . . .	37
6.4	Fitting by Reduction to Polynomial Equations . . . . .	37
6.4.1	Reduction for an Ellipse . . . . .	38
6.4.2	Reduction for an Ellipsoid . . . . .	39
<b>7</b>	<b>Fitting a Cylinder to 3D Points</b>	<b>39</b>
7.1	Representation of a Cylinder . . . . .	40
7.2	The Least-Squares Error Function . . . . .	40
7.3	An Equation for the Radius . . . . .	41
7.4	An Equation for the Center . . . . .	41
7.5	An Equation for the Direction . . . . .	43
7.6	Fitting for a Specified Direction . . . . .	45
7.7	Pseudocode and Experiments . . . . .	45

<b>8</b>	<b>Fitting a Cone to 3D Points</b>	<b>52</b>
8.1	Initial Choice for the Parameters of the Error Function . . . . .	53
8.1.1	Simple Attempt to Reconstruct Height Extremes . . . . .	56
8.1.2	Attempt to Reconstruct the Cone Axis Direction . . . . .	58
8.1.3	Attempt to Reconstruct the Cone Vertex . . . . .	58
<b>9</b>	<b>Fitting a Paraboloid to 3D Points of the Form <math>(x, y, f(x, y))</math></b>	<b>62</b>

# 1 Introduction

This document describes least-squares minimization algorithms for fitting point sets by linear structures or quadratic structures. The organization is somewhat different from that of the previous version of the document. Modifications include the following.

- A section on the general formulation for nonlinear least-squares fitting is now available. The standard approach is to estimate parameters using numerical minimizers (Gauss–Newton or Levenberg–Marquardt).
- A new algorithm for fitting points by a circle, sphere or hypersphere is provided. The algorithm is non-iterative, so the computation time is bounded and small.
- In the previous version, the sections about fitting of points by ellipses or ellipsoids were severely lacking details and not useful for developing algorithms. Several algorithms are now provided for such fitting, including a general approach for fitting points by hyperellipsoids.
- The document for fitting points by a cylinder has been moved to this document. The website hyperlink to the cylinder document has been redirected to this document.
- A section has been added for fitting points by a single-sided cone.
- Pseudocode is now provided for each of the algorithms. Hyperlinks still exist for those algorithms implemented in the GTEngine source code.

Other documents using least-squares algorithms for fitting points with curve or surface structures are available at the website. The document for fitting points with a torus is new to the website (as of August 2018).

- [Least-Squares Fitting of Data with Polynomials](#)
- [Least-Squares Fitting of Data with B-Spline Curves](#)
- [Least-Squares Reduction of B-Spline Curves](#)
- [Fitting 3D Data with a Helix](#)
- [Least-Squares Fitting of Data with B-Spline Surfaces](#)
- [Fitting 3D Data with a Torus](#)

The document [Least-Squares Fitting of Segments by Line or Plane](#) describes a least-squares algorithm where the input is a set of line segments rather than a set of points. The output is a line (segments in  $n$  dimensions) or a plane (segments in 3 dimensions) or a hyperplane (segments in  $n$  dimensions).

## 2 The General Formulation for Nonlinear Least-Squares Fitting

Let  $\mathbf{F}(\mathbf{p}) = (F_0(\mathbf{p}), F_1(\mathbf{p}), \dots, F_{n-1}(\mathbf{p}))$  be a vector-valued function of the parameters  $\mathbf{p} = (p_0, p_1, \dots, p_{m-1})$ . The nonlinear least-squares problem is to minimize the real-valued error function  $E(\mathbf{p}) = |\mathbf{F}(\mathbf{p})|^2$ .

Let  $J = d\mathbf{F}/d\mathbf{p} = [dF_r/dp_c]$  denote the Jacobian matrix, which is the matrix of first-order partial derivatives of the components of  $\mathbf{F}$ . The matrix has  $n$  rows and  $m$  columns, and the indexing  $(r, c)$  refers to row  $r$  and column  $c$ . A first-order approximation is

$$\mathbf{F}(\mathbf{p} + \mathbf{d}) \doteq \mathbf{F}(\mathbf{p}) + J(\mathbf{p})\mathbf{d} \tag{1}$$

where  $\mathbf{d}$  is an  $m \times 1$  vector with small length. Consequently, an approximation to the error function is

$$E(\mathbf{p} + \mathbf{d}) = |\mathbf{F}(\mathbf{p} + \mathbf{d})|^2 = |\mathbf{F}(\mathbf{p}) + J(\mathbf{p})\mathbf{d}|^2 \tag{2}$$

The goal is to choose  $\mathbf{d}$  to minimize  $|\mathbf{F}(\mathbf{p}) + J(\mathbf{p})\mathbf{d}|^2$  and, hopefully, with  $E(\mathbf{p} + \mathbf{d}) < E(\mathbf{p})$ . Choosing an initial  $\mathbf{p}_0$ , the hope is that the algorithm generates a sequence  $\mathbf{p}_i$  for which  $E(\mathbf{p}_{i+1}) < E(\mathbf{p}_i)$  and, in the limit,  $E(\mathbf{p}_j)$  approaches the global minimum of  $E$ . The algorithm is referred to as Gauss–Newton iteration.

For a single Gauss–Newton iteration, we need to choose  $\mathbf{d}$  to minimize  $|\mathbf{F}(\mathbf{p}) + J(\mathbf{p})\mathbf{d}|^2$  where  $\mathbf{p}$  is fixed. This is a linear least-squares problem which can be formulated using the normal equations

$$J^T(\mathbf{p})J(\mathbf{p})\mathbf{d} = -J^T(\mathbf{p})\mathbf{F}(\mathbf{p}) \tag{3}$$

The matrix  $J^T J$  is positive semidefinite. If it is invertible, then

$$\mathbf{d} = -(J^T(\mathbf{p})J(\mathbf{p}))^{-1}\mathbf{F}(\mathbf{p}) \tag{4}$$

If it is not invertible, some other algorithm must be used to choose  $\mathbf{d}$ ; one option is to use gradient descent for the step. A Cholesky decomposition can be used to solve the linear system.

During Gauss–Newton iteration, if  $E$  does not decrease for a step of the algorithm, one can modify the algorithm to Levenberg–Marquardt iteration. The idea is to smooth the linear system to

$$(J^T(\mathbf{p})J(\mathbf{p}) + \lambda I) \mathbf{d} = -J^T(\mathbf{p})\mathbf{F}(\mathbf{p}) \tag{5}$$

where  $I$  is the identity matrix of appropriate size and  $\lambda > 0$  is the smoothing factor. The strategy for choosing the initial  $\lambda$  and how to adjust it as you compute iterations depends on the problem at hand.

For a more detailed discussion, see [Gauss–Newton algorithm](#) and [Levenberg–Marquardt algorithm](#). Implementations of the Cholesky decomposition, Gauss–Newton method and Levenberg–Marquardt method in GTEngine can be found in [GteCholeskyDecomposition.h](#), [GteGaussNewtonMinimizer.h](#) and [GteLevenbergMarquardtMinimizer.h](#).

### 3 Affine Fitting of Points Using Height Fields

We have a set of measurements  $\{(\mathbf{X}_i, h_i)\}_{i=1}^m$  for which  $\mathbf{X}_i \in \mathbb{R}^n$  are sampled independent variables and  $h_i \in \mathbb{R}$  is a sampled dependent variable. The hypothesis is that  $h$  is related to  $\mathbf{X}$  via an affine transformation  $h = \mathbf{A} \cdot \mathbf{X} + b$ , where  $\mathbf{A}$  is an  $n \times 1$  vector of constants and  $b$  is a scalar constant. The goal is to estimate  $\mathbf{A}$  and  $b$  from the samples. The choice of name  $h$  stresses that the measurement errors are in the direction of *height* above the plane containing the  $\mathbf{X}$  measurements.

#### 3.1 Fitting by a Line in 2 Dimensions

The measurements are  $\{(x_i, h_i)\}_{i=1}^m$  where  $x$  is an independent variable and  $h$  is a dependent variable. The affine transformation we want to estimate is  $h = ax + b$ , where  $a$  and  $b$  are scalars. This defines a line that

best fits the samples in the sense that the sum of the squared errors between the  $h_i$  and the line values  $ax_i + b$  is minimized. Note that the error is measured only in the  $h$ -direction.

Define the error function for the least-squares minimization to be

$$E(a, b) = \sum_{i=1}^m [(ax_i + b) - h_i]^2 \quad (6)$$

This function is nonnegative and its graph is a paraboloid whose vertex occurs when the gradient satisfies  $\nabla E(a, b) = (\partial E/\partial a, \partial E/\partial b) = (0, 0)$ . This leads to a system of two linear equations in  $a$  and  $b$  which can be easily solved. Precisely,

$$\begin{aligned} 0 &= \partial E/\partial a = 2 \sum_{i=1}^m [(ax_i + b) - h_i]x_i \\ 0 &= \partial E/\partial b = 2 \sum_{i=1}^m [(ax_i + b) - h_i] \end{aligned} \quad (7)$$

and so

$$\begin{bmatrix} \sum_{i=1}^m x_i^2 & \sum_{i=1}^m x_i \\ \sum_{i=1}^m x_i & \sum_{i=1}^m 1 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^m x_i h_i \\ \sum_{i=1}^m h_i \end{bmatrix} \quad (8)$$

The system is solved by standard numerical algorithms. If implemented directly, this formulation can lead to an ill-conditioned linear system. To avoid this, you should first compute the averages  $\bar{x} = (\sum_{i=1}^m x_i)/m$  and  $\bar{h} = (\sum_{i=1}^m h_i)/m$  and subtract them from the data. The fitted line is of the form  $h - \bar{h} = \bar{a}(x - \bar{x}) + \bar{b}$ . The linear system of equations that determines the coefficients is

$$\begin{bmatrix} \sum_{i=1}^m (x_i - \bar{x})^2 & 0 \\ 0 & m \end{bmatrix} \begin{bmatrix} \bar{a} \\ \bar{b} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^m (x_i - \bar{x})(h_i - \bar{h}) \\ 0 \end{bmatrix} \quad (9)$$

and has solution

$$\bar{a} = \frac{\sum_{i=1}^m (x_i - \bar{x})(h_i - \bar{h})}{\sum_{i=1}^m (x_i - \bar{x})^2}, \quad \bar{b} = 0 \quad (10)$$

In terms of the original inputs,  $a = \bar{a}$  and  $b = \bar{h} - \bar{a}\bar{x}$ .

### 3.1.1 Pseudocode for Fitting by a Line

Listing 1 contains pseudocode for fitting a height line to points in 2 dimensions.

---

**Listing 1.** Pseudocode for fitting a height line to points in 2 dimensions. The number of input points must be at least 2. The returned Boolean value is true as long as the numerator of equation (10) is positive; that is, when the points are not all the same point. An implementation in a slightly more general framework is [GteApprHeightLine2.h](#).

```
bool FitHeightLine(int numPoints, Vector2 points[],
                  Real& barX, Real& barH, Real& barA)
{
    // Compute the mean of the points.
    Vector2 mean = { 0, 0 };
    for (int i = 0; i < numPoints; ++i)
    {
```

```

    mean += points[i];
}
mean /= numPoints;

// Compute the linear system matrix and vector elements.
Real xxSum = 0, xhSum = 0;
for (int i = 0; i < numPoints; ++i)
{
    Vector2 diff = points[i] - mean;
    xxSum += diff[0] * diff[0];
    linear += diff[0] * diff[1];
}

// Solve the linear system.
if (xxSum > 0)
{
    // Compute the fitted line h(x) = barH + barA * (x - barX).
    barX = mean[0];
    barH = mean[1];
    barA = linear / xxSum;
    return true;
}
else
{
    // The output is invalid. The points are all the same.
    barX = 0;
    barH = 0;
    barA = 0;
    return false;
}
}

```

### 3.2 Fitting by a Plane in 3 Dimensions

The measurements are  $\{(x_i, y_i, h)\}_{i=1}^m$  where  $x$  and  $y$  are independent variables and  $h$  is a dependent variable. The affine transformation we want to estimate is  $h = a_0x + a_1y + b$ , where  $a_0$ ,  $a_1$  and  $b$  are scalars. This defines a plane that best fits the samples in the sense that the sum of the squared errors between the  $h_i$  and the plane values  $a_0x_i + a_1y_i + b$  is minimized. Note that the error is measured only in the  $h$ -direction.

Define the error function for the least-squares minimization to be

$$E(a_0, a_1, b) = \sum_{i=1}^m [(a_0x_i + a_1y_i + b) - h_i]^2 \quad (11)$$

This function is nonnegative and its graph is a hyperparaboloid whose vertex occurs when the gradient satisfies  $\nabla E(a_0, a_1, b) = (\partial E/\partial a_0, \partial E/\partial a_1, \partial E/\partial b) = (0, 0, 0)$ . This leads to a system of three linear equations in  $a_0$ ,  $a_1$  and  $b$  which can be easily solved. Precisely,

$$\begin{aligned}
 0 &= \partial E/\partial a_0 = 2 \sum_{i=1}^m [(Ax_i + By_i + C) - z_i]x_i \\
 0 &= \partial E/\partial a_1 = 2 \sum_{i=1}^m [(Ax_i + By_i + C) - z_i]y_i \\
 0 &= \partial E/\partial b = 2 \sum_{i=1}^m [(Ax_i + By_i + C) - z_i]
 \end{aligned} \quad (12)$$

and so

$$\begin{bmatrix} \sum_{i=1}^m x_i^2 & \sum_{i=1}^m x_i y_i & \sum_{i=1}^m x_i \\ \sum_{i=1}^m x_i y_i & \sum_{i=1}^m y_i^2 & \sum_{i=1}^m y_i \\ \sum_{i=1}^m x_i & \sum_{i=1}^m y_i & \sum_{i=1}^m 1 \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ b \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^m x_i h_i \\ \sum_{i=1}^m y_i h_i \\ \sum_{i=1}^m h_i \end{bmatrix} \quad (13)$$

The solution is solved by standard numerical algorithms. If implemented directly, this formulation can lead to an ill-conditioned linear system. To avoid this, you should first compute the averages  $\bar{x} = (\sum_{i=1}^m x_i)/m$ ,  $\bar{y} = (\sum_{i=1}^m y_i)/m$  and  $\bar{h} = (\sum_{i=1}^m h_i)/m$  and subtract them from the data. The fitted plane is of the form  $h - \bar{h} = \bar{a}_0(x - \bar{x}) + \bar{a}_1(y - \bar{y}) + \bar{b}$ . The linear system of equations that determines the coefficients is

$$\begin{bmatrix} \ell_{00} & \ell_{01} & 0 \\ \ell_{01} & \ell_{11} & 0 \\ 0 & 0 & m \end{bmatrix} \begin{bmatrix} \bar{a}_0 \\ \bar{a}_1 \\ \bar{b} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^m (x_i - \bar{x})^2 & \sum_{i=1}^m (x_i - \bar{x})(y_i - \bar{y}) & 0 \\ \sum_{i=1}^m (x_i - \bar{x})(y_i - \bar{y}) & \sum_{i=1}^m (y_i - \bar{y})^2 & 0 \\ 0 & 0 & m \end{bmatrix} \begin{bmatrix} \bar{a}_0 \\ \bar{a}_1 \\ \bar{b} \end{bmatrix} \quad (14)$$

$$= \begin{bmatrix} \sum_{i=1}^m (h_i - \bar{h})(x_i - \bar{x}) \\ \sum_{i=1}^m (h_i - \bar{h})(y_i - \bar{y}) \\ 0 \end{bmatrix} = \begin{bmatrix} r_0 \\ r_1 \\ 0 \end{bmatrix}$$

and has solution

$$\bar{a}_0 = \frac{\ell_{11}r_0 - \ell_{01}r_1}{\ell_{00}\ell_{11} - \ell_{01}^2}, \quad \bar{a}_1 = \frac{\ell_{00}r_1 - \ell_{01}r_0}{\ell_{00}\ell_{11} - \ell_{01}^2}, \quad \bar{b} = 0 \quad (15)$$

In terms of the original inputs,  $a_0 = \bar{a}_0$ ,  $a_1 = \bar{a}_1$  and  $b = \bar{h} - \bar{a}_0\bar{x} - \bar{a}_1\bar{y}$ .

### 3.2.1 Pseudocode for Fitting by a Plane

Listing 2 contains pseudocode for fitting a height plane to points in 3 dimensions.

---

**Listing 2.** Pseudocode for fitting a height plane to points in 3 dimensions. The number of input points must be at least 3. The returned Boolean value is `true` as long as the matrix of the linear system has nonzero determinant. An implementation in a slightly more general framework is [GteApprHeightPlane3.h](#).

```
bool FitHeightPlane(int numPoints, Vector3 points[],
    Real& barX, Real& barY, Real& barH, Real& barA0, Real& barA1)
{
    // Compute the mean of the points.
    Vector3 mean = { 0, 0, 0 };
    for (int i = 0; i < numPoints; ++i)
    {
        mean += points[i];
    }
    mean /= numPoints;

    // Compute the linear system matrix and vector elements.
    Real xxSum = 0, xySum = 0, xhSum = 0, yySum = 0, yhSum = 0;
    for (int i = 0; i < numPoints; ++i)
    {
        Vector3 diff = points[i] - mean;
        xxSum += diff[0] * diff[0];
        xySum += diff[0] * diff[1];
    }
}
```



```

    xhSum += diff[0] * diff[2];
    yySum += diff[1] * diff[1];
    yhSum += diff[1] * diff[2];
}

// Solve the linear system.
Real det = xxSum * yySum - xySum * xySum;
if (det != 0)
{
    // Compute the fitted plane  $h(x,y) = \text{barH} + \text{barA0} * (x - \text{barX}) + \text{barA1} * (y - \text{barY})$ .
    barX = mean[0];
    barY = mean[1];
    barH = mean[2];
    barA0 = (yySum * xhSum - xySum * yhSum) / det;
    barA1 = (xxSum * yhSum - xySum * xhSum) / det;
    return true;
}
else
{
    // The output is invalid. The points are all the same or they are collinear.
    barX = 0;
    barY = 0;
    barH = 0;
    barA0 = 0;
    barA1 = 0;
    return false;
}
}

```

### 3.3 Fitting by a Hyperplane in $n + 1$ Dimensions

The measurements are  $\{(\mathbf{X}_i, h_i)\}_{i=1}^m$  where the  $n$  components of  $\mathbf{X}$  are independent variables and  $h$  is a dependent variable. The affine transformation we want to estimate is  $h = \mathbf{A} \cdot \mathbf{X} + b$ , where  $\mathbf{A}$  is an  $n \times 1$  vector of constants and  $b$  is a scalar constant. This defines a hyperplane that best fits the samples in the sense that the sum of the squared errors between the  $h_i$  and the hyperplane values  $\mathbf{A} \cdot \mathbf{X}_i + b$  is minimized. Note that the error is measured only in the  $h$ -direction.

Define the error function for the least-squares minimization to be

$$E(\mathbf{A}, b) = \sum_{i=1}^m [(\mathbf{A} \cdot \mathbf{X}_i + b) - h_i]^2 \quad (16)$$

This function is nonnegative and its graph is a hyperparaboloid whose vertex occurs when the gradient satisfies  $\nabla E(\mathbf{A}, b) = (\partial E / \partial \mathbf{A}, \partial E / \partial b) = (\mathbf{0}, 0)$ . This leads to a system of  $n + 1$  linear equations in  $\mathbf{A}$  and  $b$  which can be easily solved. Precisely,

$$\begin{aligned} \mathbf{0} &= \partial E / \partial \mathbf{A} = 2 \sum_{i=1}^m [(\mathbf{A} \cdot \mathbf{X}_i + b) - h_i] \mathbf{X}_i \\ 0 &= \partial E / \partial b = 2 \sum_{i=1}^m [(\mathbf{A} \cdot \mathbf{X}_i + b) - h_i] \end{aligned} \quad (17)$$

and so

$$\begin{bmatrix} \sum_{i=1}^m \mathbf{X}_i \mathbf{X}_i^T & \sum_{i=1}^m \mathbf{X}_i \\ \sum_{i=1}^m \mathbf{X}_i^T & \sum_{i=1}^m 1 \end{bmatrix} \begin{bmatrix} \mathbf{A} \\ b \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^m h_i \mathbf{X}_i \\ \sum_{i=1}^m h_i \end{bmatrix} \quad (18)$$

The solution is solved by standard numerical algorithms. If implemented directly, this formulation can lead to an ill-conditioned linear system. To avoid this, you should first compute the averages  $\bar{\mathbf{X}} = (\sum_{i=1}^m \mathbf{X}_i)/m$  and  $\bar{h} = (\sum_{i=1}^m h_i)/m$  and subtract them from the data. The fitted hyperplane is of the form  $h - \bar{h} = \bar{\mathbf{A}} \cdot (\mathbf{X} - \bar{\mathbf{X}}) + \bar{b}$ . The linear system of equations that determines the coefficients is

$$\begin{bmatrix} \sum_{i=1}^m (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})^\top & \mathbf{0} \\ \mathbf{0}^\top & m \end{bmatrix} \begin{bmatrix} \bar{\mathbf{A}} \\ \bar{b} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^m (h_i - \bar{h})(\mathbf{X}_i - \bar{\mathbf{X}}) \\ 0 \end{bmatrix} \quad (19)$$

and has solution

$$\bar{\mathbf{A}} = \left( \sum_{i=1}^m (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})^\top \right)^{-1} \left( \sum_{i=1}^m (h_i - \bar{h})(\mathbf{X}_i - \bar{\mathbf{X}}) \right), \quad \bar{b} = 0 \quad (20)$$

In terms of the original inputs,  $\mathbf{A} = \bar{\mathbf{A}}$  and  $b = \bar{h} - \bar{\mathbf{A}} \cdot \bar{\mathbf{X}}$ .

### 3.3.1 Pseudocode for Fitting a Hyperplane

Listing 3 contains pseudocode for fitting a height hyperplane to points in  $n + 1$  dimensions.

---

**Listing 3.** Pseudocode for fitting a height hyperplane to points in  $n + 1$  dimensions. The number of input points must be at least  $n$ . The returned Boolean value is true as long as the matrix of the linear system has nonzero determinant.

```
bool FitHeightHyperplane(int numPoints, Vector<n + 1> points [],
    Vector<n>& barX, Real& barH, Vector<n>& barA)
{
    // Compute the mean of the points.
    Vector<n + 1> mean = Vector<n>::ZERO;
    for (int i = 0; i < numPoints; ++i)
    {
        mean += points[i];
    }
    mean /= numPoints;

    // Compute the linear system matrix and vector elements. The function
    // Vector<n> Head<n>(Vector<n + 1> V) returns (V[0], ..., V[n-1]).
    Matrix<n, n> L = Matrix<n, n>::ZERO;
    Vector<n> R = Vector<n>::ZERO;
    for (int i = 0; i < numPoints; ++i)
    {
        Vector<n + 1> diff = points[i] - mean;
        Vector<n> XminusBarX = Head<n>(diff);
        Real HminusBarH = diff[n];
        L += OuterProduct(XminusBarX, XminusBarX); // (X[i]-barX[i])*(X[i]-barX[i])^T
        R += HminusBarH * XminusBarX;
    }

    // Solve the linear system.
    Real det = Determinant(L);
    if (det != 0)
    {
        // Compute the fitted plane h(X) = barH + Dot(barA, X - barX).
        barX = Head<n>(mean);
        barH = mean[n];
        barA = SolveLinearSystem(L, R); // solve L*A = R
        return true;
    }
}
```

```

}
else
{
    // The output is invalid. The points appear not to live on a
    // hyperplane; they might live in an affine subspace of dimension
    // smaller than n.
    barX = Vector<n>::ZERO;
    barH = 0;
    barA = Vector<n>::ZERO;
    return false;
}
}
}

```

---

## 4 Affine Fitting of Points Using Orthogonal Regression

We have a set of measurements  $\{\mathbf{X}_i\}_{i=1}^m$  for which  $\mathbf{X}_i \in \mathbb{R}^n$  are sampled independent variables. The hypothesis is that the points are sampled from a  $k$ -dimensional affine subspace in  $n$ -dimensional space. Such a space is referred to as a *k-dimensional flat*. The classic cases include fitting a line to points in  $n$  dimensions and fitting a plane to points in 3 dimensions. The latter is a special case of fitting a hyperplane, an  $(n - 1)$ -dimensional flat, to points in  $n$  dimensions.

In the height-field fitting algorithms, the least-squares errors were measured in a specified direction (the height direction). An alternative is to measure the errors in the perpendicular direction to the purported affine subspace. This approach is referred to as *orthogonal regression*.

### 4.1 Fitting by a Line [1 Dimension]

The algorithm may be applied to sample points  $\{\mathbf{X}_i\}_{i=1}^m$  in any dimension  $n$ . Let the line have origin  $\mathbf{A}$  and unit-length direction  $\mathbf{D}$ , both  $n \times 1$  vectors. Define  $\mathbf{Y}_i = \mathbf{X}_i - \mathbf{A}$ , which can be written as  $\mathbf{Y}_i = (\mathbf{D} \cdot \mathbf{Y}_i)\mathbf{D} + \mathbf{D}_i^\perp$  where  $\mathbf{D}_i^\perp$  is the perpendicular vector from  $\mathbf{X}_i$  to its projection on the line. The squared length of this vector is  $|\mathbf{D}_i^\perp|^2 = |\mathbf{Y}_i - d_i\mathbf{D}|^2$ . The error function for the least-squares minimization is  $E(\mathbf{A}, \mathbf{D}) = \sum_{i=1}^m |\mathbf{D}_i^\perp|^2$ . Two alternate forms for this function are

$$E(\mathbf{A}, \mathbf{D}) = \sum_{i=1}^m \left( \mathbf{Y}_i^\top (I - \mathbf{D}\mathbf{D}^\top) \mathbf{Y}_i \right) \quad (21)$$

and

$$E(\mathbf{A}, \mathbf{D}) = \mathbf{D}^\top \left( \sum_{i=1}^m \left( (\mathbf{Y}_i \cdot \mathbf{Y}_i)I - \mathbf{Y}_i \mathbf{Y}_i^\top \right) \right) \mathbf{D} = \mathbf{D}^\top M \mathbf{D} \quad (22)$$

where  $M$  is a positive semidefinite symmetric matrix that depends on  $\mathbf{A}$  and the  $\mathbf{Y}_i$  but not in  $\mathbf{D}$ .

Compute the derivative of equation (21) with respect to  $\mathbf{A}$  to obtain

$$\frac{\partial E}{\partial \mathbf{A}} = -2 \left[ I - \mathbf{D}\mathbf{D}^\top \right] \sum_{i=1}^m \mathbf{Y}_i \quad (23)$$

At a minimum value of  $E$ , it is necessary that this derivative is zero, which it is when  $\sum_{i=1}^m \mathbf{Y}_i = 0$ , implying  $\mathbf{A} = (1/m) \sum_{i=1}^m \mathbf{X}_i$ , the average of the sample points. In fact there are infinitely many solutions,  $\mathbf{A} + s\mathbf{D}$ , for any scalar  $s$ . This is simply a statement that  $\mathbf{A}$  is a point on the best-fit line, but any other point on the line may serve as the origin for that line.

Equation (22) is a quadratic form  $\mathbf{D}^T M \mathbf{D}$  whose minimum is the smallest eigenvalue of  $M$ , computed using standard eigensystem solvers. A corresponding unit length eigenvector  $\mathbf{D}$  completes our construction of the least-squares line. The covariance matrix of the input points is  $C = \sum_{i=1}^m \mathbf{Y}_i \mathbf{Y}_i^T$ . Defining  $\delta = \sum_{i=1}^m \mathbf{Y}_i^T \mathbf{Y}_i$ , we see that  $M = \delta I - C$ , where  $I$  is the identity matrix. Therefore,  $M$  and  $C$  have the same eigenspaces. The eigenspace corresponding to the minimum eigenvalue of  $M$  is the same as the eigenspace corresponding to the maximum eigenvalue of  $C$ . In an implementation, it is sufficient to process  $C$  and avoid the additional cost to compute  $M$ .

#### 4.1.1 Pseudocode for the General Case

Listing 4 contains pseudocode for fitting a line to points in  $n$  dimensions with  $n \geq 2$ .

---

**Listing 4.** Pseudocode for fitting a line to points in  $n$  dimensions using orthogonal regression. The number of input points must be at least 2. The returned Boolean value is `true` as long as the covariance matrix of the linear system has a 1-dimensional eigenspace for the maximum eigenvalue of the covariance matrix.

```
bool FitOrthogonalLine(int numPoints, Vector<n> points[],
    Vector<n>& origin, Vector<n>& direction)
{
    // Compute the mean of the points.
    Vector<n> mean = Vector<n>::ZERO;
    for (int i = 0; i < numPoints; ++i)
    {
        mean += points[i];
    }
    mean /= numPoints;

    // Compute the covariance matrix of the points.
    Matrix<n,n> C = Matrix<n,n>::ZERO;
    for (int i = 0; i < numPoints; ++i)
    {
        Vector<n> diff = points[i] - mean;
        C += OuterProduct(diff, diff); // diff * diff^T
    }

    // Compute the eigenvalues and eigenvectors of C, where the eigenvalues are sorted
    // in nondecreasing order (eigenvalues[0] <= eigenvalues[1] <= ...).
    Real eigenvalues[n];
    Vector<n> eigenvectors[n];
    SolveEigensystem(C, eigenvalues, eigenvectors);

    // Set the output information.
    origin = mean;
    direction = eigenvectors[n-1];

    // The fitted line is unique when the maximum eigenvalue has multiplicity 1.
    return eigenvalues[n-2] < eigenvalues[n-1];
}
```

---

Specializations for 2 and 3 dimensions are simple, computing only the upper-triangular elements of  $C$  and passing them to specialized eigensolvers for 2 and 3 dimensions. Implementations are [GteApprOrthogonalLine2.h](#) and [GteApprOrthogonalLine3.h](#).

## 4.2 Fitting by a Hyperplane [( $n - 1$ ) Dimensions]

The algorithm may be applied to sample points  $\{\mathbf{X}_i\}_{i=1}^m$  in any dimension  $n$ . Let the hyperplane be defined implicitly by  $\mathbf{N} \cdot (\mathbf{X} - \mathbf{A}) = 0$ , where  $\mathbf{N}$  is a unit-length normal to the hyperplane and  $\mathbf{A}$  is a point on the hyperplane. Define  $\mathbf{Y}_i = \mathbf{X}_i - \mathbf{A}$ , which can be written as  $\mathbf{Y}_i = (\mathbf{N} \cdot \mathbf{Y}_i)\mathbf{N} + \mathbf{N}_i^\perp$  where  $\mathbf{N}_i^\perp$  is a vector that is perpendicular to  $\mathbf{N}$ . The squared length of the projection of  $\mathbf{Y}_i$  onto the normal line for the hyperplane is  $(\mathbf{N} \cdot \mathbf{Y}_i)^2$ . The error function for the least-squares minimization is  $E(\mathbf{A}, \mathbf{N}) = \sum_{i=1}^m (\mathbf{N} \cdot \mathbf{Y}_i)^2$ . Two alternate forms for this function are

$$E(\mathbf{A}, \mathbf{N}) = \sum_{i=1}^m \left( \mathbf{Y}_i^\top (\mathbf{N}\mathbf{N}^\top) \mathbf{Y}_i \right) \quad (24)$$

and

$$E(\mathbf{A}, \mathbf{N}) = \mathbf{N}^\top \left( \sum_{i=1}^m \mathbf{Y}_i \mathbf{Y}_i^\top \right) \mathbf{N} = \mathbf{N}^\top C \mathbf{N} \quad (25)$$

where  $C = \sum_{i=1}^m \mathbf{Y}_i \mathbf{Y}_i^\top$  is the covariance matrix of the  $\mathbf{Y}_i$ .

Compute the derivative of equation (24) with respect to  $\mathbf{A}$  to obtain

$$\frac{\partial E}{\partial \mathbf{A}} = 2 \left( \mathbf{N}\mathbf{N}^\top \right) \sum_{i=1}^m \mathbf{Y}_i \quad (26)$$

At a minimum value of  $E$ , it is necessary that this derivative is zero, which it is when  $\sum_{i=1}^m \mathbf{Y}_i = 0$ , implying  $\mathbf{A} = (1/m) \sum_{i=1}^m \mathbf{X}_i$ , the average of the sample points. In fact there are infinitely many solutions,  $\mathbf{A} + \mathbf{W}$ , where  $\mathbf{W}$  is any vector perpendicular to  $\mathbf{N}$ . This is simply a statement that the average is on the best-fit hyperplane, but any other point on the hyperplane may serve as the origin for that hyperplane.

Equation (25) is a quadratic form  $\mathbf{N}^\top C \mathbf{N}$  whose minimum is the smallest eigenvalue of  $C$ , computed using standard eigensystem solvers. A corresponding unit-length eigenvector  $\mathbf{N}$  completes our construction of the least-squares hyperplane.

### 4.2.1 Pseudocode for the General Case

Listing 5 contains pseudocode for fitting a hyperplane to points in  $n$  dimensions with  $n \geq 3$ .

---

**Listing 5.** Pseudocode for fitting a line to points in  $n$  dimensions using orthogonal regression. The number of input points must be at least  $n$ . The returned Boolean value is `true` as long as the covariance matrix of the linear system has a 1-dimensional eigenspace for the minimum eigenvalue of the covariance matrix.

```
bool FitOrthogonalHyperplane(int numPoints, Vector<n> points [],
    Vector<n>& origin, Vector<n>& normal)
{
    // Compute the mean of the points.
```

```

Vector<n> mean = Vector<n>::ZERO;
for (int i = 0; i < numPoints; ++i)
{
    mean += points[i];
}
mean /= numPoints;

// Compute the covariance matrix of the points.
Matrix<n,n> C = Matrix<n,n>::ZERO;
for (int i = 0; i < numPoints; ++i)
{
    Vector<n> diff = points[i] - mean;
    C += OuterProduct(diff, diff); // diff * diff^T
}

// Compute the eigenvalues and eigenvectors of M, where the eigenvalues are sorted
// in nondecreasing order (eigenvalues[0] <= eigenvalues[1] <= ...).
Real eigenvalues[n];
Vector<n> eigenvectors[n];
SolveEigensystem(C, eigenvalues, eigenvectors);

// Set the output information.
origin = mean;
normal = eigenvectors[0];

// The fitted hyperplane is unique when the minimum eigenvalue has multiplicity 1.
return eigenvalues[0] < eigenvalues[1];
}

```

A specialization for 3 dimensions is simple, computing only the upper-triangular elements of  $C$  and passing them to a specialized eigensolver for 3 dimensions. An implementation is [GteApprOrthogonalPlane3.h](#).

### 4.3 Fitting by a Flat [ $k$ Dimensions]

Orthogonal regression to fit  $n$ -dimensional points by a line or by a hyperplane can be generalized to fitting by an affine subspace called a  $k$ -dimensional flat, where you may choose  $k$  such that  $1 \leq k \leq n - 1$ . A line is a 1-flat and a hyperplane is an  $(n - 1)$ -flat.

For dimension  $n = 3$ , we fit with flats that are either lines ( $k = 1$ ) or planes ( $k = 2$ ). For dimensions  $n \geq 4$  and flat dimensions  $1 \leq k \leq n - 1$ , the generalization of orthogonal regression is the following. The bases mentioned here are for the linear portion of the affine subspace; that is, the basis vectors are relative to an origin at a point  $\mathbf{A}$ . The flat has an orthonormal basis  $\{\mathbf{F}_j\}_{j=1}^k$  and the orthogonal complement has an orthonormal basis  $\{\mathbf{P}_j\}_{j=1}^{n-k}$ . The union of the two bases is an orthonormal basis for  $\mathbb{R}^n$ . Any input point  $\mathbf{X}_i$ ,  $1 \leq i \leq m$ , can be represented by

$$\mathbf{X}_i = \mathbf{A} + \sum_{j=1}^k f_{ij} \mathbf{F}_j + \sum_{j=1}^{n-k} p_{ij} \mathbf{P}_j = \left( \mathbf{A} + \sum_{j=1}^k f_{ij} \mathbf{F}_j \right) + \left( \sum_{j=1}^{n-k} p_{ij} \mathbf{P}_j \right) \quad (27)$$

The left-parenthesized term is the portion of  $\mathbf{X}_i$  that lives in the flat and the right-parenthesized is the portion that is the deviation of  $\mathbf{X}_i$  from the flat. The least-squares problem is about choosing the two bases so that the sum of squared lengths of the deviations is as small as possible.

Define  $\mathbf{Y}_i = \mathbf{X}_i - \mathbf{A}$ . The basis coefficients are  $f_{ij} = \mathbf{F}_j \cdot \mathbf{Y}_i$  and  $p_{ij} = \mathbf{P}_j \cdot \mathbf{Y}_i$ . The squared length of the deviation is  $\sum_{j=1}^{n-k} p_{ij}^2$ . The error function for the least-squares minimization is the sum of the squared

lengths for all inputs,  $E = \sum_{i=1}^m \sum_{j=1}^{n-k} p_{ij}^2$ . Two alternate forms for this function are

$$E(\mathbf{A}, \mathbf{P}_1, \dots, \mathbf{P}_{n-k}) = \sum_{i=1}^m \mathbf{Y}_i^\top \left( \sum_{j=1}^{n-k} \mathbf{P}_j \mathbf{P}_j^\top \right) \mathbf{Y}_i \quad (28)$$

and

$$E(\mathbf{A}, \mathbf{P}_1, \dots, \mathbf{P}_{n-k}) = \sum_{j=1}^{n-k} \mathbf{P}_j^\top \left( \sum_{i=1}^m \mathbf{Y}_i \mathbf{Y}_i^\top \right) \mathbf{P}_j = \sum_{j=1}^{n-k} \mathbf{P}_j^\top \mathbf{C} \mathbf{P}_j \quad (29)$$

where  $\mathbf{C} = \sum_{i=1}^m \mathbf{Y}_i \mathbf{Y}_i^\top$ .

Compute the derivative of equation (28) with respect to  $\mathbf{A}$  to obtain

$$\frac{\partial E}{\partial \mathbf{A}} = 2 \left( \sum_{j=1}^{n-k} \mathbf{P}_j \mathbf{P}_j^\top \right) \sum_{i=1}^m \mathbf{Y}_i \quad (30)$$

At a minimum value of  $E$ , it is necessary that this derivative is zero, which it is when  $\sum_{i=1}^m \mathbf{Y}_i = 0$ , implying  $\mathbf{A} = (1/m) \sum_{i=1}^m \mathbf{X}_i$ , the average of the sample points. In fact there are infinitely many solutions,  $\mathbf{A} + \mathbf{W}$ , where  $\mathbf{W}$  is any vector in the orthogonal complement of the subspace spanned by the  $\mathbf{P}_j$ ; this subspace is the one spanned by the  $\mathbf{F}_j$ . This is simply a statement that the average is on the best-fit flat, but any other point on the flat may serve as the origin for that flat. Because we are choosing  $\mathbf{A}$  to be the average of the input points, the matrix  $\mathbf{C}$  is the covariance matrix for the input points.

The last term in equation (29) is a sum of quadratic forms involving the matrix  $\mathbf{C}$  and the vectors  $\mathbf{P}_j$  that are a basis (unit length, mutually perpendicular). The minimum value of the quadratic form is the smallest eigenvalue  $\lambda_1$  of  $\mathbf{C}$ , so we may choose  $\mathbf{P}_1$  to be a unit-length eigenvector of  $\mathbf{C}$  corresponding to  $\lambda_1$ . We must choose  $\mathbf{P}_2$  to be unit length and perpendicular to  $\mathbf{P}_1$ . If the eigenspace for  $\lambda_1$  is 1-dimensional, the next smallest value we can attain by the quadratic form is the smallest eigenvalue  $\lambda_2$  for which  $\lambda_1 < \lambda_2$ .  $\mathbf{P}_2$  is chosen to be a corresponding unit-length eigenvector. However, if  $\lambda_1$  has an eigenspace of dimension larger than 1, we can choose  $\mathbf{P}_2$  in that eigenspace but which is perpendicular to  $\mathbf{P}_1$ .

Generally, let  $\{\lambda_\ell\}_{\ell=1}^r$  be the  $r$  distinct eigenvalues of the covariance matrix  $\mathbf{C}$ ; we know  $1 \leq r \leq n$  and  $\lambda_1 < \lambda_2 < \dots < \lambda_r$ . Let the dimension of the eigenspace for  $\lambda_\ell$  be  $d_\ell \geq 1$ ; we know that  $\sum_{\ell=1}^r d_\ell = n$ . List the eigenvalues and eigenvectors in order of increasing eigenvalue, including repeated values,

$$\underbrace{\begin{matrix} \lambda_1 & \cdots & \lambda_1 \\ \mathbf{V}_1^1 & \cdots & \mathbf{V}_{d_1}^1 \end{matrix}}_{d_1 \text{ terms}} \quad \underbrace{\begin{matrix} \lambda_2 & \cdots & \lambda_2 \\ \mathbf{V}_1^2 & \cdots & \mathbf{V}_{d_2}^2 \end{matrix}}_{d_2 \text{ terms}} \quad \cdots \quad \underbrace{\begin{matrix} \lambda_r & \cdots & \lambda_r \\ \mathbf{V}_1^r & \cdots & \mathbf{V}_{d_r}^r \end{matrix}}_{d_r \text{ terms}} \quad (31)$$

The list has  $n$  items. The eigenvalue  $\lambda_\ell$  has an eigenspace with orthonormal basis  $\{\mathbf{V}_j^\ell\}_{j=1}^{d_\ell}$ . In this list, choose the first  $k$  eigenvectors to be  $\mathbf{P}_1$  through  $\mathbf{P}_k$  and choose the last  $n - k$  eigenvectors to be  $\mathbf{F}_1$  through  $\mathbf{F}_{n-k}$ .

It is possible that one (or more) of the  $\mathbf{P}_j$  and one (or more) of the  $\mathbf{F}_j$  are in the eigenspace for the same eigenvalue. In this case, the fitted flat is not unique, and one should re-examine the choice of dimension  $k$  for the fitted flat. This is analogous to the following situations in dimension  $n = 3$ :

- The input points are nearly collinear but you are trying to fit those points with a plane. The covariance matrix likely has two distinct eigenvalues  $\lambda_1 < \lambda_2$  with  $d_1 = 2$  and  $d_2 = 1$ . The basis vectors are  $\mathbf{P}_1 = \mathbf{V}_1^1$ ,  $\mathbf{F}_1 = \mathbf{V}_2^1$ , and  $\mathbf{F}_2 = \mathbf{V}_1^2$ . The first two of these are from the same eigenspace.
- The input points are spread out over a portion of a plane (and are not nearly collinear) but you are trying to fit those points with a line. The covariance matrix likely has two distinct eigenvalues  $\lambda_1 < \lambda_2$  with  $d_1 = 1$  and  $d_2 = 2$ . The basis vectors are  $\mathbf{P}_1 = \mathbf{V}_1^1$ ,  $\mathbf{P}_2 = \mathbf{V}_1^2$ , and  $\mathbf{F}_1 = \mathbf{V}_2^2$ . The last two of these are from the same eigenspace.
- The input points are not well fit by a flat of any dimension. For example, your input points are uniformly distributed over a sphere. The covariance matrix likely has one distinct eigenvalue  $\lambda_1$  of multiplicity  $d_1 = 3$ . Neither a line nor a plane is a good fit to the input points—in either case, a  $\mathbf{P}$ -vector and an  $\mathbf{F}$ -vector are in the same eigenspace.

The computational algorithm is to compute the average  $\mathbf{A}$  and covariance matrix of the points. Use an eigensolver whose output eigenvalues are sorted in nondecreasing order. Choose the  $\mathbf{F}_j$  to be the last  $n - k$  eigenvectors output by the eigensolver.

#### 4.3.1 Pseudocode for the General Case

Listing 6 contains pseudocode for fitting a  $k$ -dimensional flat to points in  $n$ -dimensions where  $n \geq 2$  and  $1 \leq k \leq n - 1$ .

---

**Listing 6.** Pseudocode for fitting a  $k$ -dimensional flat to points in  $n$  dimensions using orthogonal regression. The number of input points must be at least  $k + 1$ . The returned Boolean value is `true` as long as the covariance matrix of the linear system has a basis of eigenvectors sorted by nondecreasing eigenvalues for which the following holds. The subbasis that spans the linear space of the flat and the subbasis that spans the orthogonal complement of the linear space of the flat do not both contain basis vectors from the same eigenspace.

```

bool FitOrthogonalFlat(int numPoints, Vector<n> points[],
    Vector<n>& origin, Vector<k>& flatBasis, Vector<n-k>& complementBasis)
{
    // Compute the mean of the points.
    Vector<n> mean = Vector<n>::ZERO;
    for (int i = 0; i < numPoints; ++i)
    {
        mean += points[i];
    }
    mean /= numPoints;

    // Compute the covariance matrix of the points.
    Matrix<n,n> C = Matrix<n,n>::ZERO;
    for (int i = 0; i < numPoints; ++i)
    {
        Vector<n> diff = points[i] - mean;
        C += OuterProduct(diff, diff); // diff * diff^T
    }

    // Compute the eigenvalues and eigenvectors of M, where the eigenvalues are sorted
    // in nondecreasing order (eigenvalues[0] <= eigenvalues[1] <= ...).
    Real eigenvalues[n];
    Vector<n> eigenvectors[n];
    SolveEigensystem(C, eigenvalues, eigenvectors);
}

```



```

// Set the output information. The basis for the fitted flat corresponds
// to the largest variances and the basis for the complement of the fitted
// flat corresponds to the smallest variances.
origin = mean;
for (int i = 0; i < n - k; ++i)
{
    complementBasis[i] = eigenvectors[i];
}
for (int i = 0; i < k; ++i)
{
    flatBasis[i] = eigenvectors[n - k + i];
}

// The fitted flat and its complement do not have vectors from the same eigenspace.
return eigenvalues[n - k - 1] < eigenvalues[n - k];
}

```

In the special case of fitting with a line ( $k = 1$ ), the line direction is the only element in `flatBasis`. In the special case of fitting with a hyperplane ( $k = n - 1$ ), the hyperplane normal is the only element in `complementBasis`.

## 5 Fitting a Hypersphere to Points

The classic cases are fitting 2-dimensional points by circles and 3-dimensional points by spheres. In  $n$ -dimensions, the objects are called hyperspheres, defined implicitly by the quadratic equation  $|\mathbf{C} - \mathbf{X}|^2 = r^2$ , where  $\mathbf{C}$  is the center and  $r$  is the radius. Three algorithms are presented for fitting points by hyperspheres.

### 5.1 Fitting Using Differences of Lengths and Radius

The sample points are  $\{\mathbf{X}_i\}_{i=1}^m$ . The least-squares error function involves the squares of differences between lengths and radius,

$$E(\mathbf{C}, r) = \sum_{i=1}^m (|\mathbf{C} - \mathbf{X}_i| - r)^2 \quad (32)$$

The minimization is based on computing points where the gradient of  $E$  is zero. The partial derivative with respect to  $r$  is

$$\frac{\partial E}{\partial r} = -2 \sum_{i=1}^m (|\mathbf{C} - \mathbf{X}_i| - r) \quad (33)$$

Setting the derivative equal to zero and solving for the radius,

$$r = \frac{1}{m} \sum_{i=1}^m |\mathbf{C} - \mathbf{X}_i| \quad (34)$$

which says that the radius is the average of the distances from the sample points to the center  $\mathbf{C}$ . The partial derivative with respect to  $\mathbf{C}$  is

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{C}} &= 2 \sum_{i=1}^m (|\mathbf{C} - \mathbf{X}_i| - r) \frac{\partial |\mathbf{C} - \mathbf{X}_i|}{\partial \mathbf{C}} \\ &= 2 \sum_{i=1}^m (|\mathbf{C} - \mathbf{X}_i| - r) \frac{\mathbf{C} - \mathbf{X}_i}{|\mathbf{C} - \mathbf{X}_i|} \\ &= 2 \sum_{i=1}^m \left( (\mathbf{C} - \mathbf{X}_i) - r \frac{\mathbf{C} - \mathbf{X}_i}{|\mathbf{C} - \mathbf{X}_i|} \right) \end{aligned} \quad (35)$$

Setting the derivative equal to zero and solving for the center,

$$\mathbf{C} = \frac{1}{m} \sum_{i=1}^m \mathbf{X}_i + r \frac{1}{m} \sum_{i=1}^m \frac{\mathbf{C} - \mathbf{X}_i}{|\mathbf{C} - \mathbf{X}_i|} = \frac{1}{m} \sum_{i=1}^m \mathbf{X}_i + \left( \frac{1}{m} \sum_{i=1}^m |\mathbf{C} - \mathbf{X}_i| \right) \left( \frac{1}{m} \sum_{i=1}^m \frac{\mathbf{C} - \mathbf{X}_i}{|\mathbf{C} - \mathbf{X}_i|} \right) \quad (36)$$

The average of the samples is  $\bar{\mathbf{X}}$ . Define length  $L_i = |\mathbf{C} - \mathbf{X}_i|$ ; the average of the lengths is  $\bar{L}$ . Define unit-length vector  $\mathbf{U}_i = (\mathbf{C} - \mathbf{X}_i)/|\mathbf{C} - \mathbf{X}_i|$ ; the average of the unit-length vectors is  $\bar{\mathbf{U}}$ . Equation (36) becomes

$$\mathbf{C} = \bar{\mathbf{X}} + \bar{L}\bar{\mathbf{U}} =: \mathbf{F}(\mathbf{C}) \quad (37)$$

where the last equality defines the vector-valued function  $\mathbf{F}$ . The function depends on the independent variable  $\mathbf{C}$  because both  $\bar{L}$  and  $\bar{\mathbf{U}}$  depend on  $\mathbf{C}$ . Fixed-point iteration can be applied to solve equation (37),

$$\mathbf{C}_0 = \bar{\mathbf{X}}; \quad \mathbf{C}_{i+1} = \mathbf{F}(\mathbf{C}_i), \quad i \geq 0 \quad (38)$$

Depending on the distribution of the samples, it is possible to choose a different initial guess for  $\mathbf{C}_0$  that (hopefully) leads to faster convergence.

### 5.1.1 Pseudocode for the General Case

Listing 7 contains pseudocode for fitting a hypersphere to points. The case  $n = 2$  is for circles and the case  $n = 3$  is for spheres.

---

**Listing 7.** Fitting a hypersphere to points using least squares based on squared differences of lengths and radius. If you want the incoming hypersphere center to be the initial guess for the center, set `inputCenterIsInitialGuess`; otherwise, the initial guess is computed to be the average of the samples. The maximum number of iterations is also specified. The returned function value is the number of iterations used.

```
int FitHypersphere(int numPoints, Vector<n> X[], int maxIterations, bool inputCenterIsInitialGuess,
    Vector<n>& center, Real& radius)
{
    // Compute the average of the data points.
    Vector<n> averageX = X[0];
    for (int i = 1; i < numPoints; ++i)
    {
        averageX += X[i];
    }
    averageX /= numPoints;

    // The initial guess for the center is either the incoming center of the
    // average of the sample points.
    if (!inputCenterIsInitialGuess)
    {
        center = averageX;
    }

    int iteration;
    for (iteration = 0; iteration < maxIterations; ++iteration)
    {
        // Update the estimate for the center.
        Vector<n> previousCenter = center;

        // Compute average L and average U.
        Real averageL = 0;
        Vector<n> averageU = Vector<n>::ZERO;
        for (int i = 0; i < numPoints; ++i)
```

```

{
  Vector<n> CmXi = center - X[i];
  Real length = Length(CmXi);
  if (length > 0)
  {
    averageL += length;
    averageU -= CmXi / length;
  }
}
averageL /= numPoints;
averageU /= numPoints;

center = averageX + averageL * averageU;
radius = averageL;

// Test for convergence.
if (center == previousCenter)
{
  break;
}
}
return ++iteration;
}

```

The convergence test uses an exact equality, that the previous center  $\mathbf{C}'$  and the current center  $\mathbf{C}$  are the same. In practice you might want to specify a small  $\varepsilon > 0$  and instead exit when  $|\mathbf{C} - \mathbf{C}'| \leq \varepsilon$ .

Specializations for 2 and 3 dimensions have the same implementation as the general case. Implementations are [GteApprCircle2.h](#) and [GteApprSphere3.h](#).

## 5.2 Fitting Using Differences of Squared Lengths and Squared Radius

The sample points are  $\{\mathbf{X}_i\}_{i=1}^m$ . The least-squares error function involves the squares of differences between squared lengths and squared radius,

$$E(\mathbf{C}, r^2) = \sum_{i=1}^m (|\mathbf{C} - \mathbf{X}_i|^2 - r^2)^2 \quad (39)$$

The minimization is based on computing points where the gradient of  $E$  is zero. The partial derivative with respect to  $r^2$  is

$$\frac{\partial E}{\partial r^2} = -2 \sum_{i=1}^m (|\mathbf{C} - \mathbf{X}_i|^2 - r^2) \quad (40)$$

Define  $\Delta_i = \mathbf{C} - \mathbf{X}_i$ . Setting the derivative to zero and solving for the squared radius,

$$r^2 = \frac{1}{m} \sum_{i=1}^m |\mathbf{C} - \mathbf{X}_i|^2 = \frac{1}{m} \sum_{i=1}^m \Delta_i^\top \Delta_i \quad (41)$$

which says that the squared radius is the average of the squared distances from the sample points to the center  $\mathbf{C}$ . The partial derivative with respect to  $\mathbf{C}$  is

$$\frac{\partial E}{\partial \mathbf{C}} = 4 \sum_{i=1}^m (|\mathbf{C} - \mathbf{X}_i|^2 - r^2) (\mathbf{C} - \mathbf{X}_i) = 4 \sum_{i=1}^m (\Delta_i^\top \Delta_i - r^2) \Delta_i \quad (42)$$

Setting this to zero, the center and radius must satisfy

$$\sum_{i=1}^m \left( \Delta_i^\top \Delta_i - r^2 \right) \Delta_i = \mathbf{0} \quad (43)$$

Expanding the squared lengths,

$$\Delta_i^\top \Delta_i = |\mathbf{C}|^2 - 2\mathbf{C}^\top \mathbf{X}_i + |\mathbf{X}_i|^2 \quad (44)$$

Substituting into equation (41),

$$r^2 = |\mathbf{C}|^2 - 2\mathbf{C}^\top \left( \frac{1}{m} \sum_{i=1}^m \mathbf{X}_i \right) + \frac{1}{m} \sum_{i=1}^m |\mathbf{X}_i|^2 = |\mathbf{C}|^2 - 2\mathbf{C}^\top \mathbf{A} + \frac{1}{m} \sum_{i=1}^m |\mathbf{X}_i|^2 \quad (45)$$

where  $\mathbf{A} = (\sum_{i=1}^m \mathbf{X}_i) / m$  is the average of the samples. Define  $\mathbf{Y}_i = \mathbf{X}_i - \mathbf{A}$ . Some algebra will show that

$$\begin{aligned} \Delta_i^\top \Delta_i - r^2 &= |\mathbf{C}|^2 - 2\mathbf{C}^\top \mathbf{X}_i + |\mathbf{X}_i|^2 - \left( |\mathbf{C}|^2 - 2\mathbf{C}^\top \mathbf{A} + \frac{1}{m} \sum_{j=1}^m |\mathbf{X}_j|^2 \right) \\ &= -2\mathbf{C}^\top \mathbf{Y}_i + |\mathbf{X}_i|^2 - \frac{1}{m} \sum_{j=1}^m |\mathbf{X}_j|^2 \\ &= -2(\mathbf{C} - \mathbf{A})^\top \mathbf{Y}_i + |\mathbf{Y}_i|^2 - \frac{1}{m} \sum_{j=1}^m |\mathbf{Y}_j|^2 \\ &= -2(\mathbf{C} - \mathbf{A})^\top \mathbf{Y}_i + B_i \end{aligned} \quad (46)$$

where the last equality defines  $B_i$ . Equation (43) becomes

$$\begin{aligned} \mathbf{0} &= \sum_{i=1}^m \left( \Delta_i^\top \Delta_i - r^2 \right) \Delta_i \\ &= \sum_{i=1}^m \left( -2(\mathbf{C} - \mathbf{A})^\top \mathbf{Y}_i + B_i \right) ((\mathbf{C} - \mathbf{A}) - \mathbf{Y}_i) \\ &= ((\mathbf{C} - \mathbf{A})^\top \sum_{i=1}^m \mathbf{Y}_i) (\mathbf{C} - \mathbf{A}) + 2 \left( \sum_{i=1}^m \mathbf{Y}_i \mathbf{Y}_i^\top \right) (\mathbf{C} - \mathbf{A}) + \left( \sum_{i=1}^m B_i \right) (\mathbf{C} - \mathbf{A}) - \sum_{i=1}^m B_i \mathbf{Y}_i \end{aligned} \quad (47)$$

It is easily shown that  $\sum_{i=1}^m \mathbf{Y}_i = \mathbf{0}$  and  $\sum_{i=1}^m B_i = 0$ ; therefore,

$$\mathbf{0} = 2 \left( \sum_{i=1}^m \mathbf{Y}_i \mathbf{Y}_i^\top \right) (\mathbf{C} - \mathbf{A}) - \sum_{i=1}^m \left( \mathbf{Y}_i^\top \mathbf{Y}_i \right) \mathbf{Y}_i \quad (48)$$

The least-squares center is obtained by solving the previous equation,

$$\mathbf{C} = \mathbf{A} + \frac{1}{2} \left( \sum_{i=1}^m \mathbf{Y}_i \mathbf{Y}_i^\top \right)^{-1} \sum_{i=1}^m \left( \mathbf{Y}_i^\top \mathbf{Y}_i \right) \mathbf{Y}_i \quad (49)$$

### 5.2.1 Pseudocode for the General Case

Listing 8 contains pseudocode for fitting a hypersphere to points. The case  $n = 2$  is for circles and the case  $n = 3$  is for spheres.

**Listing 8.** Fitting a hypersphere to points using least squares based on squared differences of squared lengths and square radius. The algorithm requires inverting the covariance matrix. If the matrix is invertible, the output center and radius are valid and the function returns `true`. If the matrix is not invertible, the function returns `false`, and the center and radius are invalid (but set to zero so at least they are initialized).

```

bool FitHypersphere(int numPoints, Vector<n> X[], Vector<n>& center, Real& radius)
{
    // Compute the average of the data points.
    Vector<n> A = Vector<n>::ZERO;
    for (int i = 0; i < numPoints; ++i)
    {
        A += X[i];
    }
    A /= numPoints;

    // Compute the covariance matrix M of the Y[i] = X[i]-A and the right-hand side R of the linear
    // system M*(C-A) = R.
    Matrix<n, n> M = Matrix<n, n>::ZERO;
    Vector<n> R = Vector<n>::ZERO;
    for (int i = 0; i < numPoints; ++i)
    {
        Vector<n> Y = X[i] - A;
        Matrix<n, n> YYT = OuterProduct(Y, Y); // Y*Transpose(Y)
        Real YTY = Dot(Y, Y); // Transpose(Y)*Y
        M += YYT;
        R += YTY * Y;
    }
    R /= 2;

    // Solve the linear system M*(C-A) = R for the center C. The function 'bool Solve(M, R, S)' tries
    // to solve the linear system M*S = R. If M is invertible, the function returns true and S is the
    // solution. If M is not invertible, the function returns false and S is invalid.
    Vector<n> CmA;
    if (Solve(M, R, CmA))
    {
        center = A + CmA;
        Real rsqr = 0;
        for (int i = 0; i < numPoints; ++i)
        {
            Vector<n> delta = X[i] - center;
            rsqr += Dot(delta, delta);
        }
        rsqr /= numPoints;
        radius = sqrt(rsqr);
        return true;
    }
    else
    {
        center = Vector<n>::ZERO;
        radius = 0;
        return false;
    }
}

```

---

## 5.2.2 Pseudocode for Circles

Listing 9 is a specialization of Listing 8 for circles in 2 dimensions. The matrix inversion requires only computing the upper-triangular part of the covariance matrix and uses cofactors for inversion.

---

**Listing 9.** Fitting a circle to points using least squares based on squared differences of squared lengths and square radius. The algorithm requires inverting the covariance matrix. If the matrix is invertible, the output center and radius are valid and the function returns true. If the matrix is not invertible, the function returns false, and the center and radius are invalid (but set to zero so at least they are initialized).

```

bool FitCircle(int numPoints, Vector2 X[], Vector2& center, Real& radius)
{
    // Compute the average of the data points.
    Vector2 A = { 0, 0 };
    for (int i = 0; i < numPoints; ++i)
    {
        A += X[i];
    }
    A /= numPoints;

    // Compute the covariance matrix M of the Y[i] = X[i]-A and the right-hand side R of the linear
    // system M*(C-A) = R.
    Real M00 = 0, M01 = 0, M11 = 0;
    Vector2 R = { 0, 0 };
    for (int i = 0; i < numPoints; ++i)
    {
        Vector2 Y = X[i] - A;
        Real Y0Y0 = Y[0] * Y[0], Y0Y1 = Y[0] * Y[1], Y1Y1 = Y[1] * Y[1];
        M00 += Y0Y0; M01 += Y0Y1; M11 += Y1Y1;
        R += (Y0Y0 + Y1Y1) * Y;
    }
    R /= 2;

    // Solve the linear system M*(C-A) = R for the center C.
    Real det = M00 * M11 - M01 * M01;
    if (det != 0)
    {
        center[0] = A[0] + (M11 * R[0] - M01 * R[1]) / det;
        center[1] = A[1] + (M00 * R[1] - M01 * R[0]) / det;
        Real rsqr = 0;
        for (int i = 0; i < numPoints; ++i)
        {
            Vector2 delta = X[i] - center;
            rsqr += Dot(delta, delta);
        }
        rsqr /= numPoints;
        radius = sqrt(rsqr);
        return true;
    }
    else
    {
        center = { 0, 0 };
        radius = 0;
        return false;
    }
}

```

---

### 5.2.3 Pseudocode for Spheres

Listing 10 is a specialization of Listing 8 for spheres in 3 dimensions. The matrix inversion requires only computing the upper-triangular part of the covariance matrix and uses cofactors for inversion.

---

**Listing 10.** Fitting a sphere to points using least squares based on squared differences of squared lengths and square radius. The algorithm requires inverting the covariance matrix. If the matrix is invertible, the

output center and radius are valid and the function returns true. If the matrix is not invertible, the function returns false, and the center and radius are invalid (but set to zero so at least they are initialized).

```

bool FitSphere(int numPoints, Vector3 X[], Vector3& center, Real& radius)
{
    // Compute the average of the data points.
    Vector3 A = { 0, 0, 0 };
    for (int i = 0; i < numPoints; ++i)
    {
        A += X[i];
    }
    A /= numPoints;

    // Compute the covariance matrix M of the Y[i] = X[i]-A and the right-hand side R of the linear
    // system M*(C-A) = R.
    Real M00 = 0, M01 = 0, M02 = 0, M11 = 0, M12 = 0, M22 = 0;
    Vector3 R = { 0, 0, 0 };
    for (int i = 0; i < numPoints; ++i)
    {
        Vector2 Y = X[i] - A;
        Real Y0Y0 = Y[0] * Y[0], Y0Y1 = Y[0] * Y[1], Y0Y2 = Y[0] * Y[2];
        Real Y1Y1 = Y[1] * Y[1], Y1Y2 = Y[1] * Y[2], Y2Y2 = Y[2] * Y[2];
        M00 += Y0Y0; M01 += Y0Y1; M02 += Y0Y2;
        M11 += Y1Y1; M12 += Y1Y2; M22 += Y2Y2;
        R += (Y0Y0 + Y1Y1 + Y2Y2) * Y;
    }
    R /= 2;

    // Solve the linear system M*(C-A) = R for the center C.
    Real cof00 = M11 * M22 - M12 * M12;
    Real cof01 = M02 * M12 - M01 * M22;
    Real cof02 = M01 * M12 - M02 * M11;
    Real det = M00 * cof00 + M01 * cof01 + M02 * cof02;
    if (det != 0)
    {
        Real cof11 = M00 * M22 - M02 * M02;
        Real cof12 = M01 * M02 - M00 * M12;
        Real cof22 = M00 * M11 - M01 * M01;
        center[0] = A[0] + (cof00 * R[0] + cof01 * R[1] + cof02 * R[2]) / det;
        center[1] = A[1] + (cof01 * R[0] + cof11 * R[1] + cof12 * R[2]) / det;
        center[2] = A[2] + (cof02 * R[0] + cof12 * R[1] + cof22 * R[2]) / det;
        Real rsqr = 0;
        for (int i = 0; i < numPoints; ++i)
        {
            Vector3 delta = X[i] - center;
            rsqr += Dot(delta, delta);
        }
        rsqr /= numPoints;
        radius = sqrt(rsqr);
        return true;
    }
    else
    {
        center = { 0, 0, 0 };
        radius = 0;
        return false;
    }
}

```

### 5.3 Fitting the Coefficients of a Quadratic Equation

The general quadratic equation that represents a hypersphere in  $n$  dimensions is

$$b_0 + \mathbf{b}_1 \cdot \mathbf{X} + b_2 |\mathbf{X}|^2 = 0 \quad (50)$$

where  $b_0$  and  $b_2 \neq 0$  are scalar constants and  $\mathbf{b}_1$  is an  $n \times 1$  vector of scalars. Define

$$\mathbf{b} = \begin{bmatrix} b_0 \\ \mathbf{b}_1 \\ b_2 \end{bmatrix}, \quad \mathbf{V} = \begin{bmatrix} 1 \\ \mathbf{X} \\ |\mathbf{X}|^2 \end{bmatrix} \quad (51)$$

which are both  $(n+2) \times 1$  vectors. The quadratic equation is  $\mathbf{b} \cdot \mathbf{V} = 0$ . Because  $\mathbf{b}$  is not zero, we can remove a degree of freedom by requiring  $|\mathbf{b}| = 1$ .

Given samples  $\{\mathbf{X}_i\}_{i=1}^m$ , we can estimate the constants using a least-squares algorithm where the error function is

$$E(\mathbf{b}) = \sum_{i=1}^m (\mathbf{b} \cdot \mathbf{V}_i)^2 = \mathbf{b}^\top \left( \sum_{i=1}^m \mathbf{V}_i \mathbf{V}_i^\top \right) \mathbf{b} = \mathbf{b}^\top V \mathbf{b} \quad (52)$$

where as a tuple,  $\mathbf{V}_i = (1, \mathbf{X}_i, |\mathbf{X}_i|^2)$ , and where  $V = \sum_{i=1}^m \mathbf{V}_i \mathbf{V}_i^\top$  is a positive semidefinite matrix (symmetric, eigenvalues are nonnegative). The error function is minimized by a unit-length eigenvector  $\mathbf{b}$  that corresponds to the minimum eigenvalue of  $V$ . The equation (50) is factored into

$$\left| \mathbf{X} + \frac{\mathbf{b}_1}{2b_2} \right|^2 = \left| \frac{\mathbf{b}_1}{2b_2} \right|^2 - \frac{b_0}{b_2} \quad (53)$$

from which we see that the hypersphere center is  $\mathbf{C} = -\mathbf{b}_1/(2b_2)$  and the radius is  $r = \sqrt{(|\mathbf{b}_1|^2 - 4b_0b_2)/(4b_2^2)}$ .

As is typical of these types of problems, it is better to subtract translate the samples by their average to obtain numerical robustness when computing with floating-point arithmetic. Define  $\mathbf{A} = (\sum_{i=1}^m \mathbf{X}_i)$  and  $\mathbf{Y}_i = \mathbf{X}_i - \mathbf{A}$ . Equation (50) becomes

$$f_0 + \mathbf{f}_1 \cdot \mathbf{Y} + f_2 |\mathbf{Y}|^2 = 0 \quad (54)$$

where

$$\mathbf{b} = \begin{bmatrix} b_0 \\ \mathbf{b}_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} f_0 - \mathbf{f}_1 \cdot \mathbf{A} + f_2 |\mathbf{A}|^2 \\ \mathbf{f}_1 - 2f_2 \mathbf{A} \\ f_2 \end{bmatrix} = \begin{bmatrix} 1 & -\mathbf{A}^\top & |\mathbf{A}|^2 \\ \mathbf{0} & I & -2\mathbf{A} \\ 0 & \mathbf{0}^\top & 1 \end{bmatrix} \begin{bmatrix} f_0 \\ \mathbf{f}_1 \\ f_2 \end{bmatrix} = M \mathbf{f} \quad (55)$$

with  $I$  the  $n \times n$  identity matrix and  $\mathbf{0}$  the  $n \times 1$  zero vector. The last equality defines the  $(n+2) \times (n+2)$  upper-triangular matrix  $M$  and the  $(n+2) \times 1$  vector  $\mathbf{f}$ . The error function is now

$$E(\mathbf{f}) = \mathbf{b}^\top V \mathbf{b} = (R\mathbf{f})^\top V (R\mathbf{f}) = \mathbf{f}^\top (R^\top V R) \mathbf{f} = \mathbf{f}^\top W \mathbf{f} \quad (56)$$

where the last equation defines the positive semidefinite matrix  $W$ . Observe that

$$W = R^\top \left( \sum_{i=1}^m \mathbf{V}_i \mathbf{V}_i^\top \right) R = \sum_{i=1}^m (R^\top \mathbf{V}_i) (R^\top \mathbf{V}_i)^\top = \sum_{i=1}^m \mathbf{W}_i \mathbf{W}_i^\top \quad (57)$$

where

$$\mathbf{W}_i = R^\top \mathbf{V}_i = \begin{bmatrix} 1 & \mathbf{0}^\top & 0 \\ -\mathbf{A} & I & \mathbf{0} \\ |\mathbf{A}|^2 & -2\mathbf{A}^\top & 1 \end{bmatrix} \begin{bmatrix} 1 \\ \mathbf{X}_i \\ |\mathbf{X}_i|^2 \end{bmatrix} = \begin{bmatrix} 1 \\ \mathbf{X}_i - \mathbf{A} \\ |\mathbf{X}_i - \mathbf{A}|^2 \end{bmatrix} \quad (58)$$



Therefore, we can subtract  $\mathbf{A}$  from the samples, compute the matrix  $W$ , extract its eigenvector  $\mathbf{f}$  corresponding to its minimum eigenvalue, compute  $\mathbf{b} = R\mathbf{f}$  and then the center and radius using equation (53). The matrix  $W$  written as a block matrix *and* divided by the number of points  $m$  to keep the numerical intermediate values on the order of the sample values is

$$W = \begin{bmatrix} 1 & \mathbf{0}^T & \frac{1}{m} \sum_{i=1}^m |\mathbf{Y}_i|^2 \\ \mathbf{0} & \frac{1}{m} \sum_{i=1}^m \mathbf{Y}_i \mathbf{Y}_i^T & \frac{1}{m} \sum_{i=1}^m |\mathbf{Y}_i|^2 \mathbf{Y}_i \\ \frac{1}{m} \sum_{i=1}^m |\mathbf{Y}_i|^2 & \frac{1}{m} \sum_{i=1}^m |\mathbf{Y}_i|^2 \mathbf{Y}_i & \frac{1}{m} |\mathbf{Y}_i|^4 \end{bmatrix} \quad (59)$$

### 5.3.1 Pseudocode for the General Case

Listing 11 contains pseudocode for fitting a hypersphere to points. The case  $n = 2$  is for circles and the case  $n = 3$  is for spheres.

---

**Listing 11.** Fitting a hypersphere to points using least squares to fit the coefficients of a quadratic equation defining the hypersphere. The algorithm requires computing the squared radius in terms of the components of an eigenvector. For samples not really distributed well on a hypersphere, the purported squared radius might be negative. The function returns true when that computation is nonnegative or false when it is negative.

```
bool FitHypersphere(int numPoints, Vector<n> X[], Vector<n>& center, Real& radius)
{
    // Compute the average of the data points and the squared length of the average.
    Vector<n> A = Vector<n>::ZERO;
    for (int i = 0; i < numPoints; ++i)
    {
        A += X[i];
    }
    A /= numPoints;
    Real sqrLenA = Dot(A, A);

    // Compute the components of W. Block(r, c, rsize, csize) is an accessor to the block
    // whose upper-left location is (r, c) and whose size is rsize-by-csize. The indices
    // r and c are zero-based.
    Matrix<n + 2, n + 2> W = Matrix<n + 2, n + 2>::ZERO;
    for (int i = 0; i < numPoints; ++i)
    {
        Vector<n> Y = X[i] - A;
        Matrix<n, n> YYT = OuterProduct(Y, Y); // Y*Transpose(Y)
        Real YTY = Dot(Y, Y); // Transpose(Y)*Y
        Vector YTYT = YTY * Y;
        Real YTYTY = YTY * YTY;
        W.Block(0, 0, 1, 1) += 1;
        W.Block(0, n + 1, 1, 1) += YTY;
        W.Block(1, 1, n, n) += YYT;
        W.Block(1, n + 1, n, 1) += YTYT;
        W.Block(n + 1, 0, 1, 1) += YTY;
        W.Block(n + 1, 1, 1, n) += YTYT;
        W.Block(n + 1, n + 1, 1, 1) += YTYTY;
    }
    W /= numPoints;

    // Compute the eigenvalues and eigenvectors of M, where the eigenvalues are sorted
    // in nondecreasing order (eigenvalues[0] <= eigenvalues[1] <= ...).
    Real eigenvalues[n + 2];
    Vector<n> eigenvectors[n + 2];
    SolveEigensystem(W, eigenvalues, eigenvectors);
}
```

```

// Block(i, isize) is an accessor to the block whose initial location is i and whose
// size is isize.
Real f0 = eigenvectors[0].Block(0, 1);
Vector<n> f1 = eigenvectors[0].Block(1, n);
Real f2 = eigenvectors[0].Block(n + 1, 1);
Real b0 = f0 - Dot(f1, A) + f2 * sqrLenA;
Vector<n> b1 = f1 - 2 * f2 * A;
Real b2 = f2;

if (b2 != 0)
{
    Real discr = Dot(b1, b1) - 4 * b0 * b2;
    if (discr >= 0)
    {
        center = -b1 / (2 * b2);
        radius = sqrt(discr / (4 * b2 * b2));
        return true;
    }
}

center = Vector<n>::ZERO;
radius = 0;
return false;
}

```

### 5.3.2 Pseudocode for Circles

Listing 12 contains pseudocode for fitting a circle to points.

**Listing 12.** Fitting a circle to points using least squares to fit the coefficients of a quadratic equation defining the hypersphere. The algorithm requires computing the squared radius in terms of the components of an eigenvector. For samples not really distributed well on a hypersphere, the purported squared radius might be negative. The function returns true when that computation is nonnegative or false when it is negative.

```

bool FitCircle(int numPoints, Vector2 X[], Vector2& center, Real& radius)
{
    // Compute the average of the data points and the squared length of the average.
    Vector2 A = { 0, 0 };
    for (int i = 0; i < numPoints; ++i)
    {
        A += X[i];
    }
    A /= numPoints;
    Real sqrLenA = Dot(A, A);

    // Compute the upper-triangular components of W.
    Matrix4x4 W = Matrix4x4::ZERO;
    for (int i = 0; i < numPoints; ++i)
    {
        Vector2 Y = X[i] - A;
        Real Y0Y0 = Y[0] * Y[0], Y0Y1 = Y[0] * Y[1], Y1Y1 = Y[1] * Y[1];
        Real RR = Y0Y0 + Y1Y1, RRRR = RR * RR;
        Real Y0RR = Y0 * RR, Y1RR = Y1 * RR;
        W(0, 3) += RR;
        W(1, 1) += Y0Y0;
        W(1, 2) += Y0Y1;
        W(1, 3) += Y0RR;
    }
}

```

```

        W(2, 2) += Y1Y1;
        W(2, 3) += Y1RR;
        W(3, 3) += RRRR;
    }
    W /= numPoints;
    W(0, 0) = 1;

    // Fill in the lower-triangular components of W.
    W(3, 0) = W(0, 3);
    W(2, 1) = W(1, 2);
    W(3, 1) = W(1, 3);
    W(3, 2) = W(2, 3);

    // Compute the eigenvalues and eigenvectors of M, where the eigenvalues are sorted
    // in nondecreasing order
    // eigenvalues[0] <= eigenvalues[1] <= eigenvalues[2] <= eigenvalues[3]
    Real eigenvalues[4];
    Vector4 eigenvectors[4];
    SolveEigensystem(W, eigenvalues, eigenvectors);

    Real f0 = eigenvectors[0][0];
    Vector2 f1 = { eigenvectors[0][1], eigenvectors[0][2] };
    Real f2 = eigenvectors[0][3];
    Real b0 = f0 - Dot(f1, A) + f2 * sqrLenA;
    Vector2 b1 = f1 - 2 * f2 * A;
    Real b2 = f2;

    if (b2 != 0)
    {
        Real discr = Dot(b1, b1) - 4 * b0 * b2;
        if (discr >= 0)
        {
            center = -b1 / (2 * b2);
            radius = sqrt(discr / (4 * b2 * b2));
            return true;
        }
    }

    center = Vector2::ZERO;
    radius = 0;
    return false;
}

```

---

### 5.3.3 Pseudocode for Spheres

Listing 13 contains pseudocode for fitting a sphere to points.

---

**Listing 13.** Fitting a sphere to points using least squares to fit the coefficients of a quadratic equation defining the hypersphere. The algorithm requires computing the squared radius in terms of the components of an eigenvector. For samples not really distributed well on a hypersphere, the purported squared radius might be negative. The function returns `true` when that computation is nonnegative or `false` when it is negative.

```

bool FitSphere(int numPoints, Vector3 X[], Vector3& center, Real& radius)
{
    // Compute the average of the data points and the squared length of the average.
    Vector3 A = { 0, 0, 0 };
    for (int i = 0; i < numPoints; ++i)
    {

```

```

    A += X[i];
}
A /= numPoints;
Real sqrLenA = Dot(A, A);

// Compute the upper-triangular components of W.
Matrix5x5 W = Matrix5x5::ZERO;
for (int i = 0; i < numPoints; ++i)
{
    Vector3 Y = X[i] - A;
    Real Y0Y0 = Y[0] * Y[0], Y0Y1 = Y[0] * Y[1], Y1Y1 = Y[1] * Y[1];
    Real Y1Y1 = Y[1] * Y[1], Y1Y2 = Y[1] * Y[2], Y2Y2 = Y[2] * Y[2];
    Real RR = Y0Y0 + Y1Y1 + Y2Y2, RRRR = RR * RR;
    Real Y0RR = Y0 * RR, Y1RR = Y1 * RR, Y2RR = Y2 * RR;
    W(0, 4) += RR;
    W(1, 1) += Y0Y0;
    W(1, 2) += Y0Y1;
    W(1, 3) += Y0Y2;
    W(1, 4) += Y0RR;
    W(2, 2) += Y1Y1;
    W(2, 3) += Y1Y2;
    W(2, 4) += Y1RR;
    W(3, 3) += Y2Y2;
    W(3, 4) += Y2RR;
    W(4, 4) += RRRR;
}
W /= numPoints;
W(0, 0) = 1;

// Fill in the lower-triangular components of W.
W(4, 0) = W(0, 4);
W(2, 1) = W(1, 2);
W(3, 1) = W(1, 3);
W(4, 1) = W(1, 4);
W(3, 2) = W(2, 3);
W(4, 2) = W(2, 4);
W(4, 3) = W(3, 4);

// Compute the eigenvalues and eigenvectors of M, where the eigenvalues are sorted
// in nondecreasing order
// eigenvalues[0] <= eigenvalues[1] <= eigenvalues[2] <= eigenvalues[3] <= eigenvalues[4]
Real eigenvalues[5];
Vector5 eigenvectors[5];
SolveEigensystem(W, eigenvalues, eigenvectors);

Real f0 = eigenvectors[0][0];
Vector3 f1 = { eigenvectors[0][1], eigenvectors[0][2], eigenvectors[0][3] };
Real f2 = eigenvectors[0][4];
Real b0 = f0 - Dot(f1, A) + f2 * sqrLenA;
Vector3 b1 = f1 - 2 * f2 * A;
Real b2 = f2;

if (b2 != 0)
{
    Real discr = Dot(b1, b1) - 4 * b0 * b2;
    if (discr >= 0)
    {
        center = -b1 / (2 * b2);
        radius = sqrt(discr / (4 * b2 * b2));
        return true;
    }
}

center = Vector3::ZERO;
radius = 0;
return false;
}

```

## 6 Fitting a Hyperellipsoid to Points

The classic cases are fitting 2-dimensional points by ellipses and 3-dimensional points by ellipsoids. In  $n$ -dimensions, the objects are called hyperellipsoids, defined implicitly by the quadratic equation

$$(\mathbf{X} - \mathbf{C})^\top RDR^\top (\mathbf{X} - \mathbf{C}) = 1 \quad (60)$$

The center of the hyperellipsoid is  $\mathbf{C}$ . The matrix  $R$  is an  $n \times n$  rotation matrix with  $n \times 1$  column vectors  $\mathbf{U}_i$  for  $0 \leq i < n$ . These are unit-length vectors that are mutually perpendicular. The matrix  $D$  is an  $n \times n$  diagonal matrix with diagonal elements  $d_i = 1/e_i^2$  for  $0 \leq i < n$ . The extent  $e_i$  is the distance from  $\mathbf{C}$  to the two extreme points of the hyperellipsoid in the direction  $\mathbf{U}_i$ . The matrix  $M = RDR^\top$  is positive definite and can be expressed as

$$M = \sum_{i=0}^{n-1} \frac{\mathbf{U}_i \mathbf{U}_i^\top}{e_i^2} = \sum_{i=0}^{n-1} \frac{\mathbf{V}_i \mathbf{V}_i^\top}{e_i^2 |\mathbf{V}_i|^2} \quad (61)$$

The second equality uses vectors  $\mathbf{V}_i$  that are not necessarily unit length, which is useful when robust computations require arbitrary precision arithmetic for exact results. The length of  $\mathbf{V}_i$  is not necessary to compute; rather, we compute  $|\mathbf{V}_i|^2 = \mathbf{V}_i \cdot \mathbf{V}_i$  to eliminate floating-point rounding errors that occur when normalizing  $\mathbf{V}_i$  to obtain  $\mathbf{U}_i$ .

Let the samples be  $\{\mathbf{X}_i\}_{i=1}^m$ . If we choose to mimic the approach for fitting points by hyperspheres using differences of lengths and squared radius, the equivalent is to minimize the squared distance from points to the purported hyperellipsoid  $H$ . The nonlinear least-squares error function is

$$E(\mathbf{C}, R, D) = \sum_{i=1}^m L(\mathbf{X}_i; \mathbf{C}, R, D)^2 \quad (62)$$

where  $L(\mathbf{X}; \mathbf{C}, R, D)$  is the distance from point  $\mathbf{X}$  to the hyperellipsoid defined by equation (60). This problem is more difficult to approach than that for fitting hyperspheres. The function  $L(\mathbf{X}; \mathbf{C}, R, D)$  is nontrivial to compute, as shown by the mathematical details in [Distance from a Point to an Ellipse, an Ellipsoid, or a Hyperellipsoid](#). If one were to use Gauss–Newton or Levenberg–Marquardt algorithms for minimizing  $E(\mathbf{C}, R, D)$ , we would need to compute derivatives of  $L$  with respect to  $\mathbf{C}$ ,  $R$  or  $D$ . This is a complicated endeavor for general dimensions, but it is feasible in dimensions 2 and 3. An alternative is using Powell’s direction-set method to search for a minimum, an algorithm that does not require computing derivatives. Even so, the iteration involves a large number of evaluations of  $L$ , which is computationally expensive. Implementations for fitting ellipses and ellipsoids to points using this approach are [GteApprEllipse2.h](#) and [GteApprEllipsoid3.h](#).

### 6.1 A Tractable Nonlinear Least-Squares Error Function

A more tractable approach uses the nonlinear least-squares error function

$$E(\mathbf{C}, R, D) = \sum_{i=1}^m \left[ (\mathbf{X}_i - \mathbf{C})^\top RDR^\top (\mathbf{X}_i - \mathbf{C}) - 1 \right]^2 \quad (63)$$

Various algorithms are presented in this section for locating the parameters that minimize  $E$ , all requiring derivative computations. We can parameterize  $D$  by its  $n$  diagonal components  $d_i$  listed as a vector  $\mathbf{d}$ . We can also parameterize  $R$  by its  $n(n-1)/2$  Rodrigues angles  $\theta_{rc}$  for  $0 \leq r < n$  and  $r < c < n$ , listed as

a vector  $\boldsymbol{\theta}$ . These angles are the upper-triangular entries of an  $n \times n$  skew-symmetric matrix  $S$ , and  $R$  is generated from them by the exponential map  $R = \exp(S)$ . The error function can then be thought of as a function  $E(\mathbf{C}, \boldsymbol{\theta}, \mathbf{d})$ , which depends on  $2n + n(n-1)/2$  scalar parameters.

In the following discussion, define  $\boldsymbol{\Delta}_i = \mathbf{X}_i - \mathbf{C}$ . The first-order derivative of  $E$  with respect to the center is

$$\frac{\partial E}{\partial \mathbf{C}} = -4 \sum_{i=1}^m \left( \boldsymbol{\Delta}_i^\top R D R^\top \boldsymbol{\Delta}_i - 1 \right) R D R^\top \boldsymbol{\Delta}_i \quad (64)$$

The first-order derivative of  $E$  with respect to the component  $d_j$  of  $\mathbf{d}$  is

$$\frac{\partial E}{\partial d_j} = 2 \sum_{i=1}^m \left( \boldsymbol{\Delta}_i^\top R D R^\top \boldsymbol{\Delta}_i - 1 \right) \boldsymbol{\Delta}_i^\top \mathbf{U}_j \mathbf{U}_j^\top \boldsymbol{\Delta}_i \quad (65)$$

where  $E_j$  is the matrix whose elements are all 0 except for the  $j$ -th diagonal entry which is one. The vector  $\mathbf{U}_j$  is the unit-length  $j$ -th column of  $R$ .

The first-order derivatives with respect to the  $\theta_{rc}$  are complicated to compute in high dimensions. However, we can derive simplified equations for the derivatives of  $E$  with respect to those angles,

$$\frac{\partial E}{\partial \theta_{rc}} = 2 \sum_{i=1}^m \left( \boldsymbol{\Delta}_i^\top R D R^\top \boldsymbol{\Delta}_i - 1 \right) \boldsymbol{\Delta}_i^\top \left( R D \frac{\partial R^\top}{\partial \theta_{rc}} + \frac{\partial R}{\partial \theta_{rc}} D R^\top \right) \boldsymbol{\Delta}_i \quad (66)$$

In the matrix-valued expression involving  $D$ ,  $R$  and  $\partial R / \partial \theta_{rc}$ , multiply by  $R^\top$  on the left and  $R$  on the right to obtain

$$D \left( \frac{\partial R^\top}{\partial \theta_{rc}} R \right) + \left( R^\top \frac{\partial R}{\partial \theta_{rc}} \right) D = \left[ (d_k - d_j) \mathbf{U}_j^\top \frac{\partial \mathbf{U}_k}{\partial \theta_{rc}} \right] = [M_{rc, jk}] = M_{rc} \quad (67)$$

where the right-hand side defines the matrix  $M_{rc}$  and its entry  $M_{rc, jk}$  in row  $j$  and column  $k$ . The diagonal entries of  $M$  are all zero ( $M_{rc, jj} = 0$ ) because the differences  $d_k - d_j = 0$  for  $k = j$ .  $M_{rc}$  is symmetric ( $M_{rc, kj} = M_{rc, jk}$ ) because of the mutual orthogonality of the unit-length columns of  $R$ ; that is,  $\mathbf{U}_j^\top \mathbf{U}_k = 0$  for  $j \neq k$  which implies  $\mathbf{U}_j^\top \partial \mathbf{U}_k / \partial \theta_{rc} = -\partial \mathbf{U}_j^\top / \partial \theta_{rc} \mathbf{U}_k$  and because  $(d_j - d_k) = -(d_k - d_j)$ .

Define the vector  $\mathbf{P}_{jk}$  for  $j < k$  to have all 0 entries except at indices  $j$  and  $k$  where the entries are 1. Define the vector  $\mathbf{N}_{jk}$  for  $j < k$  to have all 0 entries except at index  $j$  where the entry is 1 and at index  $k$  where the entry is  $-1$ . The matrix  $M$  can be factored into a linear combination of outer products of vectors,

$$M_{rc} = \sum_{j=0}^{n-1} \sum_{k=j+1}^{n-1} \frac{1}{2} M_{rc, jk} \left( \mathbf{P}_{jk} \mathbf{P}_{jk}^\top - \mathbf{N}_{jk} \mathbf{N}_{jk}^\top \right) \quad (68)$$

Multiplying by  $R$  on the left and  $R^\top$  on the right,

$$R D \frac{\partial R^\top}{\partial \theta_{rc}} + \frac{\partial R}{\partial \theta_{rc}} D R^\top = R M_{rc} R^\top = \sum_{j=0}^{n-1} \sum_{k=j+1}^{n-1} M_{rc, jk} \left( \mathbf{U}_j \mathbf{U}_k^\top + \mathbf{U}_k \mathbf{U}_j^\top \right) \quad (69)$$

where I have used  $R \mathbf{P}_{jk} = \mathbf{U}_j + \mathbf{U}_k$  and  $R \mathbf{N}_{jk} = \mathbf{U}_j - \mathbf{U}_k$ . Substituting this into equation (66),

$$\frac{\partial E}{\partial \theta_{rc}} = 2 \sum_{j=0}^{n-1} \sum_{k=j+1}^{n-1} M_{rc, jk} \sum_{i=1}^m \left( \boldsymbol{\Delta}_i^\top R D R^\top \boldsymbol{\Delta}_i - 1 \right) \boldsymbol{\Delta}_i^\top \left( \mathbf{U}_j \mathbf{U}_k^\top + \mathbf{U}_k \mathbf{U}_j^\top \right) \boldsymbol{\Delta}_i \quad (70)$$

Defining  $\theta_{rr} = d_r$  and  $M_{rr,jj} = 1$ , we can extend equation (70) to represent the derivatives with respect to the diagonal terms of  $D$  and the angles that generate  $R$ ,

$$\frac{\partial E}{\partial \theta_{rc}} = 2 \sum_{j=0}^{n-1} \sum_{k=j}^{n-1} M_{rc,jk} \mathbf{U}_j^\top \left( \sum_{i=1}^m \left( \Delta_i^\top R D R^\top \Delta_i - 1 \right) \Delta_i \Delta_i^\top \right) \mathbf{U}_k \quad (71)$$

for  $0 \leq r \leq c < n$ .

### 6.1.1 Angle Derivative for an Ellipse

For 2-dimensional rotations, the only Rodrigues angle is  $\theta_{01}$ . I will omit the subscript for ease of reading. The columns of  $R$  written as 2-tuples are  $\mathbf{U}_0 = (\cos \theta, \sin \theta)$  and  $\mathbf{U}_1 = (-\sin \theta, \cos \theta)$ . The derivatives are  $d\mathbf{U}_0/d\theta = (-\sin \theta, \cos \theta) = \mathbf{U}_1$  and  $d\mathbf{U}_1/d\theta = (-\cos \theta, -\sin \theta) = -\mathbf{U}_0$ . The matrix  $M_{01}$  has the off-diagonal entry  $M_{01,01} = (d_1 - d_0) \mathbf{U}_0^\top d\mathbf{U}_1/d\theta = (d_0 - d_1)$ . The angle derivative of  $E$  is

$$\frac{\partial E}{\partial \theta} = (d_0 - d_1) \mathbf{U}_0^\top \left( \sum_{i=1}^m \left( \Delta_i^\top R D R^\top \Delta_i - 1 \right) \Delta_i \Delta_i^\top \right) \mathbf{U}_1 \quad (72)$$

### 6.1.2 Angle Derivatives for an Ellipsoid

Rodrigues angles are used to parameterize rotation matrices via the exponential map. Specifically,  $R = \exp(S)$  for some skew-symmetric matrix

$$S = \begin{bmatrix} 0 & -s_2 & s_1 \\ s_2 & 0 & -s_0 \\ -s_1 & s_0 & 0 \end{bmatrix} = \text{Skew}(\mathbf{s}) \quad (73)$$

where the last equality defines the skew-symmetric matrix corresponding to the input 3-tuple. Define

$$\mathbf{s} = (s_0, s_1, s_2), \quad \theta = \sqrt{s_0^2 + s_1^2 + s_2^2}, \quad \alpha(\theta) = \frac{\sin \theta}{\theta}, \quad \beta(\theta) = \frac{1 - \cos \theta}{\theta^2} \quad (74)$$

The rotation matrix is

$$R(\mathbf{s}) = I + \alpha S + \beta S^2 = \begin{bmatrix} 1 - \beta(s_1^2 + s_2^2) & -\alpha s_2 + \beta s_0 s_1 & +\alpha s_1 + \beta s_0 s_2 \\ +\alpha s_2 + \beta s_0 s_1 & 1 - \beta(s_0^2 + s_2^2) & -\alpha s_0 + \beta s_1 s_2 \\ -\alpha s_1 + \beta s_0 s_2 & +\alpha s_0 + \beta s_1 s_2 & 1 - \beta(s_0^2 + s_1^2) \end{bmatrix} \quad (75)$$

where  $R(\mathbf{s})$  indicates that the rotation matrix is parameterized by the components of  $\mathbf{s}$ . When  $\mathbf{s} \neq \mathbf{0}$ , a unit-length rotation axis is  $\mathbf{s}/\theta$ .

Given a  $3 \times 3$  rotation matrix  $R = [r_{ij}]$ , the 9 elements of the matrix are dependent; specifically, they are related by equation (75). Observe that  $\partial \theta / \partial s_\ell = s_\ell / \theta$  for any  $\ell$ . Define

$$\gamma(\theta) = \frac{\alpha'(\theta)}{\theta} = \frac{\theta \cos \theta - \sin \theta}{\theta^3}, \quad \delta(\theta) = \frac{\beta'(\theta)}{\theta} = \frac{\theta \sin \theta - 2(1 - \cos \theta)}{\theta^4} \quad (76)$$

where  $\alpha(\theta)$  and  $\beta(\theta)$  are defined in equation (74). The first-order partial derivatives of the rotation matrix components with respect to the components of  $\mathbf{s}$  are

$$\frac{\partial R}{\partial s_0} = \begin{bmatrix} -\delta s_0(s_1^2 + s_2^2) & -\gamma s_0 s_2 + (\beta + \delta s_0^2)s_1 & \gamma s_0 s_1 + (\beta + \delta s_0^2)s_2 \\ \gamma s_0 s_2 + (\beta + \delta s_0^2)s_1 & -2\beta s_0 - \delta s_0(s_0^2 + s_2^2) & -(\alpha + \gamma s_0^2) + \delta s_0 s_1 s_2 \\ -\gamma s_0 s_1 + (\beta + \delta s_0^2)s_2 & (\alpha + \gamma s_0^2) + \delta s_0 s_1 s_2 & -2\beta s_0 - \delta s_0(s_0^2 + s_1^2) \end{bmatrix} \quad (77)$$

and

$$\frac{\partial R}{\partial s_1} = \begin{bmatrix} -2\beta s_1 - \delta s_1(s_1^2 + s_2^2) & -\gamma s_1 s_2 + (\beta + \delta s_1^2)s_0 & (\alpha + \gamma s_1^2) + \delta s_0 s_1 s_2 \\ \gamma s_1 s_2 + (\beta + \delta s_1^2)s_0 & -\delta s_1(s_0^2 + s_2^2) & -\gamma s_0 s_1 + (\beta + \delta s_1^2)s_2 \\ -(\alpha + \gamma s_1^2) + \delta s_0 s_1 s_2 & \gamma s_0 s_1 + (\beta + \delta s_1^2)s_2 & -2\beta s_1 - \delta s_1(s_0^2 + s_1^2) \end{bmatrix} \quad (78)$$

and

$$\frac{\partial R}{\partial s_2} = \begin{bmatrix} -2\beta s_2 - \delta s_2(s_1^2 + s_2^2) & -(\alpha + \gamma s_2^2) + \delta s_0 s_1 s_2 & \gamma s_1 s_2 + (\beta + \delta s_2^2)s_0 \\ (\alpha + \gamma s_2^2) + \delta s_0 s_1 s_2 & -2\beta s_2 - \delta s_2(s_0^2 + s_2^2) & -\gamma s_0 s_2 + (\beta + \delta s_2^2)s_1 \\ -\gamma s_1 s_2 + (\beta + \delta s_2^2)s_0 & \gamma s_0 s_2 + (\beta + \delta s_2^2)s_1 & -\delta s_2(s_0^2 + s_1^2) \end{bmatrix} \quad (79)$$

The columns of the matrices are written as

$$R = \begin{bmatrix} \mathbf{U}_0 & \mathbf{U}_1 & \mathbf{U}_2 \end{bmatrix}, \quad \frac{\partial R}{\partial s_k} = \begin{bmatrix} \frac{\partial \mathbf{U}_0}{\partial s_k} & \frac{\partial \mathbf{U}_1}{\partial s_k} & \frac{\partial \mathbf{U}_2}{\partial s_k} \end{bmatrix} \quad (80)$$

Define

$$\varepsilon(\phi) = \frac{\phi - \sin \phi}{\phi^3} = \frac{1 - \alpha(\phi)}{\phi^2} = \beta(\phi) + \gamma(\phi) \quad (81)$$

then some algebra leads to

$$\begin{aligned} \mathbf{U}_0^\top \frac{\partial \mathbf{U}_1}{\partial s_0} &= -s_1 \beta - s_0 s_2 \varepsilon, & \mathbf{U}_0^\top \frac{\partial \mathbf{U}_2}{\partial s_0} &= -s_2 \beta + s_0 s_1 \varepsilon, & \mathbf{U}_1^\top \frac{\partial \mathbf{U}_2}{\partial s_0} &= -1 + (s_1^2 + s_2^2) \varepsilon, \\ \mathbf{U}_0^\top \frac{\partial \mathbf{U}_1}{\partial s_1} &= +s_0 \beta - s_1 s_2 \varepsilon, & \mathbf{U}_0^\top \frac{\partial \mathbf{U}_2}{\partial s_1} &= +1 - (s_0^2 + s_2^2) \varepsilon, & \mathbf{U}_1^\top \frac{\partial \mathbf{U}_2}{\partial s_1} &= -s_2 \beta - s_0 s_1 \varepsilon, \\ \mathbf{U}_0^\top \frac{\partial \mathbf{U}_1}{\partial s_2} &= -1 + (s_0^2 + s_1^2) \varepsilon, & \mathbf{U}_0^\top \frac{\partial \mathbf{U}_2}{\partial s_2} &= +s_0 \beta + s_1 s_2 \varepsilon, & \mathbf{U}_1^\top \frac{\partial \mathbf{U}_2}{\partial s_2} &= +s_1 \beta - s_0 s_2 \varepsilon \end{aligned} \quad (82)$$

and

$$R^\top \left( RD \frac{\partial R}{\partial s_k}^\top + \frac{\partial R}{\partial s_k} DR^\top \right) R = \begin{bmatrix} 0 & (d_1 - d_0) \mathbf{U}_0^\top \frac{\partial \mathbf{U}_1}{\partial s_k} & (d_2 - d_0) \mathbf{U}_0^\top \frac{\partial \mathbf{U}_2}{\partial s_k} \\ (d_1 - d_0) \mathbf{U}_0^\top \frac{\partial \mathbf{U}_1}{\partial s_k} & 0 & (d_2 - d_1) \mathbf{U}_1^\top \frac{\partial \mathbf{U}_2}{\partial s_k} \\ (d_2 - d_0) \mathbf{U}_0^\top \frac{\partial \mathbf{U}_2}{\partial s_k} & (d_2 - d_1) \mathbf{U}_1^\top \frac{\partial \mathbf{U}_2}{\partial s_k} & 0 \end{bmatrix} \quad (83)$$

The angle derivatives of  $E$  are

$$\begin{aligned} \frac{\partial E}{\partial s_k} &= \left( (d_1 - d_0) \mathbf{U}_0^\top \frac{\partial \mathbf{U}_1}{\partial s_k} \right) \mathbf{U}_0^\top \left( \sum_{i=1}^m \left( \mathbf{\Delta}_i^\top R D R^\top \mathbf{\Delta}_i - 1 \right) \mathbf{\Delta}_i \mathbf{\Delta}_i^\top \right) \mathbf{U}_1 + \\ &\quad \left( (d_2 - d_0) \mathbf{U}_0^\top \frac{\partial \mathbf{U}_2}{\partial s_k} \right) \mathbf{U}_0^\top \left( \sum_{i=1}^m \left( \mathbf{\Delta}_i^\top R D R^\top \mathbf{\Delta}_i - 1 \right) \mathbf{\Delta}_i \mathbf{\Delta}_i^\top \right) \mathbf{U}_2 + \\ &\quad \left( (d_2 - d_1) \mathbf{U}_1^\top \frac{\partial \mathbf{U}_2}{\partial s_k} \right) \mathbf{U}_1^\top \left( \sum_{i=1}^m \left( \mathbf{\Delta}_i^\top R D R^\top \mathbf{\Delta}_i - 1 \right) \mathbf{\Delta}_i \mathbf{\Delta}_i^\top \right) \mathbf{U}_2 \end{aligned} \quad (84)$$



## 6.2 Initial Choice for the Parameters of the Error Function

Without any application-specific knowledge about the sample points, the simplest initial choice for the center, axis directions and axis extents is based on fitting the sample points with a Gaussian distribution. The mean of the distribution is used for the center. The covariance matrix is computed from the sample points and mean. The eigenvectors of this matrix are used for the axis directions. If  $M$  is the covariance matrix, the eigendecomposition produces  $M = RDR^T$ , where  $R$  is a rotation matrix and  $D$  is a diagonal matrix with nonnegative (but in practice usually positive) diagonal entries  $d_i$  for  $0 \leq i < n$ . You might choose the axis extents to be  $1/\sqrt{d_i}$ , but in the limit as the number of samples becomes dense enough to cover a full hyperellipsoid, this choice of axis extents is not good.

In general dimensions, the initial center is chosen to be the mean of the sample points,

$$\mathbf{C} = \frac{1}{m} \sum_{i=1}^m \mathbf{X}_i \quad (85)$$

The covariance matrix is

$$M = \frac{1}{m} \sum_{i=1}^m (\mathbf{X}_i - \mathbf{C})(\mathbf{X}_i - \mathbf{C})^T \quad (86)$$

Construct a right-handed orthonormal basis of eigenvectors  $\{\mathbf{U}_j\}_{j=0}^{n-1}$  and create the rotation matrix whose columns are these eigenvectors,

$$R = \begin{bmatrix} \mathbf{U}_0 & \mathbf{U}_1 & \cdots & \mathbf{U}_{n-1} \end{bmatrix} \quad (87)$$

Let the corresponding eigenvalues be  $\{\lambda_j\}_{j=0}^{n-1}$ . We need to choose initial values for the axis extents  $\{e_j\}_{j=0}^{n-1}$  based on the eigenvalues are other information specific to hyperellipsoids.

### 6.2.1 Initial Choice for the Parameters of an Ellipse

The ellipse is parameterized by  $\mathbf{P}(\theta) = \mathbf{C} + e_0 \cos(\theta)\mathbf{U}_0 + e_1 \sin(\theta)\mathbf{U}_1$  for  $\theta \in [0, 2\pi)$ . The algebraic equation that implicitly defines the ellipse is  $(\mathbf{P}(\theta) - \mathbf{C})^T RDR^T (\mathbf{P}(\theta) - \mathbf{C}) = 1$ , where  $R = [\mathbf{U}_0 \ \mathbf{U}_1]$  is a rotation matrix and  $D = \text{Diag}(1/e_0^2, 1/e_1^2)$ .

The mean of the ellipse points is the average of those points in the sense of integration over the ellipse,

$$\begin{aligned} \boldsymbol{\mu} &= \frac{1}{2\pi} \int_0^{2\pi} \mathbf{P}(\theta) d\theta \\ &= \frac{1}{2\pi} \int_0^{2\pi} (\mathbf{C} + e_0 \cos(\theta)\mathbf{U}_0 + e_1 \sin(\theta)\mathbf{U}_1) d\theta \\ &= \mathbf{C} + \frac{e_0}{2\pi} \left( \int_0^{2\pi} \cos(\theta) d\theta \right) \mathbf{U}_0 + \frac{e_1}{2\pi} \left( \int_0^{2\pi} \sin(\theta) d\theta \right) \mathbf{U}_1 \\ &= \mathbf{C} \end{aligned} \quad (88)$$

It is natural to choose the initial center  $\mathbf{C}$  to be the mean  $\boldsymbol{\mu}$  of the sample points.

The covariance matrix for the ellipse is

$$\begin{aligned}
\mathcal{C} &= \frac{1}{2\pi} \int_0^{2\pi} (\mathbf{P}(\theta) - \mathbf{C})(\mathbf{P}(\theta) - \mathbf{C})^\top d\theta \\
&= \frac{1}{2\pi} \int_0^{2\pi} (e_0 \cos(\theta) \mathbf{U}_0 + e_1 \sin(\theta) \mathbf{U}_1)(e_0 \cos(\theta) \mathbf{U}_0 + e_1 \sin(\theta) \mathbf{U}_1)^\top d\theta \\
&= \frac{e_0^2}{2\pi} \left( \int_0^{2\pi} \cos^2(\theta) d\theta \right) \mathbf{U}_0 \mathbf{U}_0^\top + \frac{e_1^2}{2\pi} \left( \int_0^{2\pi} \sin^2(\theta) d\theta \right) \mathbf{U}_1 \mathbf{U}_1^\top \\
&\quad + \frac{e_0 e_1}{2\pi} \left( \int_0^{2\pi} \cos(\theta) \sin(\theta) d\theta \right) (\mathbf{U}_0 \mathbf{U}_1^\top + \mathbf{U}_1 \mathbf{U}_0^\top) \\
&= \frac{e_0^2}{2} \mathbf{U}_0 \mathbf{U}_0^\top + \frac{e_1^2}{2} \mathbf{U}_1 \mathbf{U}_1^\top \\
&= R \text{Diag} \left( \frac{e_0^2}{2}, \frac{e_1^2}{2} \right) R^\top
\end{aligned} \tag{89}$$

Computing  $\mathcal{C}$  to be the covariance matrix for the sample points, factor  $\mathcal{C} = R \text{Diag}(\lambda_0, \lambda_1) R^\top$  using an eigendecomposition. The estimates for the axis extents are then  $e_j = \sqrt{2\lambda_j}$  for  $j = 0, 1$ .

In summary, compute

$$\boldsymbol{\mu} = \frac{1}{m} \sum_{i=1}^m \mathbf{X}_i, \quad \mathcal{C} = \frac{1}{m} \sum_{i=1}^m (\mathbf{X}_i - \boldsymbol{\mu})(\mathbf{X}_i - \boldsymbol{\mu})^\top \tag{90}$$

Choose the initial center  $\mathbf{C} = \boldsymbol{\mu}$ . Factor  $\mathcal{C} = R \text{Diag}(\lambda_0, \lambda_1) R^\top$  where  $R$  is a rotation matrix. Choose the initial orientation to be  $R$ . Choose the initial diagonal matrix  $D = \text{Diag}(1/(2\lambda_0), 1/(2\lambda_1))$ .

### 6.2.2 Initial Choice for the Parameters of an Ellipsoid

The ellipsoid is parameterized by  $\mathbf{P}(\theta, \phi) = \mathbf{C} + e_0 \cos(\theta) \sin(\phi) \mathbf{U}_0 + e_1 \sin(\theta) \sin(\phi) \mathbf{U}_1 + e_2 \cos(\phi) \mathbf{U}_2$  for  $\theta \in [0, 2\pi)$  and  $\phi \in [0, \pi]$ . The algebraic equation that implicitly defines the ellipsoid is  $(\mathbf{P}(\theta, \phi) - \mathbf{C})^\top R D R^\top (\mathbf{P}(\theta, \phi) - \mathbf{C}) = 1$ , where  $R = [\mathbf{U}_0 \ \mathbf{U}_1 \ \mathbf{U}_2]$  is a rotation matrix and  $D = \text{Diag}(1/e_0^2, 1/e_1^2, 1/e_2^2)$ .

The mean of the ellipsoid points is the average of those points in the sense of integration over the ellipsoid,

$$\begin{aligned}
\boldsymbol{\mu} &= \frac{1}{4\pi} \int_0^{2\pi} \int_0^\pi \mathbf{P}(\theta, \phi) \sin(\phi) d\phi d\theta \\
&= \frac{1}{4\pi} \int_0^{2\pi} \int_0^\pi (\mathbf{C} + e_0 \cos(\theta) \sin(\phi) \mathbf{U}_0 + e_1 \sin(\theta) \sin(\phi) \mathbf{U}_1 + e_2 \cos(\phi) \mathbf{U}_2) \sin(\phi) d\phi d\theta \\
&= \mathbf{C} + \frac{e_0}{4\pi} \left( \int_0^{2\pi} \int_0^\pi \cos(\theta) \sin^2(\phi) d\phi d\theta \right) \mathbf{U}_0 + \frac{e_1}{4\pi} \left( \int_0^{2\pi} \int_0^\pi \sin(\theta) \sin^2(\phi) d\phi d\theta \right) \mathbf{U}_1 \\
&\quad + \frac{e_2}{4\pi} \left( \int_0^{2\pi} \int_0^\pi \cos(\phi) \sin(\phi) d\phi d\theta \right) \mathbf{U}_2 \\
&= \mathbf{C}
\end{aligned} \tag{91}$$

It is natural to choose the initial center  $\mathbf{C}$  to be the mean  $\boldsymbol{\mu}$  of the sample points.

The covariance matrix for the ellipsoid is

$$\begin{aligned}
\mathcal{C} &= \frac{1}{4\pi} \int_0^{2\pi} \int_0^\pi (\mathbf{P}(\theta) - \mathbf{C})(\mathbf{P}(\theta) - \mathbf{C})^\top \sin(\phi) d\phi d\theta \\
&= \frac{1}{4\pi} \int_0^{2\pi} \int_0^\pi (e_0 \cos(\theta) \sin(\phi) \mathbf{U}_0 + e_1 \sin(\theta) \sin(\phi) \mathbf{U}_1 + e_2 \cos(\phi) \mathbf{U}_2) \\
&\quad (e_0 \cos(\theta) \sin(\phi) \mathbf{U}_0 + e_1 \sin(\theta) \sin(\phi) \mathbf{U}_1 + e_2 \cos(\phi) \mathbf{U}_2)^\top \sin(\phi) d\phi d\theta \\
&= \frac{e_0^2}{4\pi} \left( \int_0^{2\pi} \int_0^\pi \cos^2(\theta) \sin^3(\phi) d\phi d\theta \right) \mathbf{U}_0 \mathbf{U}_0^\top \\
&\quad + \frac{e_1^2}{4\pi} \left( \int_0^{2\pi} \int_0^\pi \sin^2(\theta) \sin^3(\phi) d\phi d\theta \right) \mathbf{U}_1 \mathbf{U}_1^\top \\
&\quad + \frac{e_2^2}{4\pi} \left( \int_0^{2\pi} \int_0^\pi \cos^2(\phi) \sin(\phi) d\phi d\theta \right) \mathbf{U}_2 \mathbf{U}_2^\top \\
&\quad + \frac{e_0 e_1}{4\pi} \left( \int_0^{2\pi} \int_0^\pi \cos(\theta) \sin(\theta) \sin^3(\phi) d\phi d\theta \right) (\mathbf{U}_0 \mathbf{U}_1^\top + \mathbf{U}_1 \mathbf{U}_0^\top) \\
&\quad + \frac{e_0 e_2}{4\pi} \left( \int_0^{2\pi} \int_0^\pi \cos(\theta) \cos(\phi) \sin^2(\phi) d\phi d\theta \right) (\mathbf{U}_0 \mathbf{U}_2^\top + \mathbf{U}_2 \mathbf{U}_0^\top) \\
&\quad + \frac{e_1 e_2}{4\pi} \left( \int_0^{2\pi} \int_0^\pi \sin(\theta) \cos(\phi) \sin^2(\phi) d\phi d\theta \right) (\mathbf{U}_1 \mathbf{U}_2^\top + \mathbf{U}_2 \mathbf{U}_1^\top) \\
&= \frac{e_0^2}{3} \mathbf{U}_0 \mathbf{U}_0^\top + \frac{e_1^2}{3} \mathbf{U}_1 \mathbf{U}_1^\top + \frac{e_2^2}{3} \mathbf{U}_2 \mathbf{U}_2^\top \\
&= R \text{Diag} \left( \frac{e_0^2}{3}, \frac{e_1^2}{3}, \frac{e_2^2}{3} \right) R^\top
\end{aligned} \tag{92}$$

Computing  $\mathcal{C}$  to be the covariance matrix for the sample points, factor  $\mathcal{C} = R \text{Diag}(\lambda_0, \lambda_1, \lambda_2) R^\top$  using an eigendecomposition. The estimates for the axis extents are then  $e_j = \sqrt{3\lambda_j}$  for  $j = 0, 1, 2$ .

In summary, compute

$$\boldsymbol{\mu} = \frac{1}{m} \sum_{i=1}^m \mathbf{X}_i, \quad \mathcal{C} = \frac{1}{m} \sum_{i=1}^m (\mathbf{X}_i - \boldsymbol{\mu})(\mathbf{X}_i - \boldsymbol{\mu})^\top \tag{93}$$

Choose the initial center  $\mathbf{C} = \boldsymbol{\mu}$ . Factor  $\mathcal{C} = R \text{Diag}(\lambda_0, \lambda_1, \lambda_2) R^\top$  where  $R$  is a rotation matrix. Choose the initial orientation to be  $R$ . Choose the initial diagonal matrix  $D = \text{Diag}(1/(3\lambda_0), 1/(3\lambda_1), 1/(3\lambda_2))$ .

### 6.2.3 Initial Choice for the Parameters of a Hyperellipsoid

Let  $S^n$  be the hypersphere of center  $\mathbf{0}$  and radius 1 that lives in  $\mathbb{R}^{n+1}$ ; for example,  $S^1$  is the unit circle in  $\mathbb{R}^2$  centered at the origin and  $S^2$  is the unit sphere in  $\mathbb{R}^3$  centered at the origin. Generally, the hypersphere is a smooth  $n$ -dimensional manifold that lives in  $(n+1)$ -dimensional space. It is implicitly defined by  $x_0^2 + \dots + x_n^2 = 1$  and can be parameterized by

$$\begin{aligned}
x_0 &= \cos \phi_0 \\
x_1 &= \sin \phi_0 \cos \phi_1 \\
x_2 &= \sin \phi_0 \sin \phi_1 \cos \phi_2 \\
&\vdots \\
x_{n-1} &= \sin \phi_0 \cdots \sin \phi_{n-2} \cos \phi_{n-1} \\
x_n &= \sin \phi_0 \cdots \sin \phi_{n-2} \sin \phi_{n-1}
\end{aligned} \tag{94}$$

where  $\phi_0 \in [0, \pi]$  for  $0 \leq i \leq n-2$  and  $\phi_{n-2} \in [0, 2\pi)$ . The parameterization for the circle  $S^1$  is  $x_0 = \cos \phi_0$  and  $x_1 = \sin \phi_0$ . The parameterization for the sphere  $S^2$  is  $x_0 = \cos \phi_0$ ,  $x_2 = \sin \phi_0 \cos \phi_1$  and

$x_3 = \sin \phi_0 \sin \phi_1$ . The parameterization for the hypersphere  $S^3$  is  $x_0 = \cos \phi_0$ ,  $x_1 = \sin \phi_0 \cos \phi_1$ ,  $x_2 = \sin \phi_0 \sin \phi_1 \cos \phi_2$  and  $x_3 = \sin \phi_0 \sin \phi_1 \sin \phi_2$ . The angles can be listed as components of a vector  $\boldsymbol{\phi} = (\phi_0, \dots, \phi_{n-1})$ .

The infinitesimal measure of surface area on  $S^n$  is

$$dA_n = \sin^{n-1}(\phi_0) \sin^{n-2}(\phi_1) \cdots \sin(\phi_{n-2}) d\phi_0 \cdots d\phi_{n-1} \quad (95)$$

For example, the infinitesimal arc length on the circle  $S^1$  (the specialization of surface area to a curve) is  $dA_1 = d\phi_0$ . The infinitesimal surface area of the sphere  $S^2$  is  $dA_2 = \sin(\phi_0) d\phi_0 d\phi_1$ . The infinitesimal surface area of the hypersphere  $S^3$  is  $dA_3 = \sin^2(\phi_0) \sin(\phi_1) d\phi_0 d\phi_1 d\phi_2$ . The total surface area of  $S^n$  is

$$A_n = \int_{S^n} dA_n = \frac{2\pi^{(n+1)/2}}{\Gamma((n+1)/2)} \quad (96)$$

where  $\Gamma(z)$  is the gamma function.

The hyperellipsoid is parameterized by

$$\mathbf{P}(\boldsymbol{\phi}) = \mathbf{C} + \sum_{j=0}^{n-1} e_j x_j(\boldsymbol{\phi}) \mathbf{U}_j \quad (97)$$

The algebraic equation that implicitly defines the ellipsoid is

$$(\mathbf{P}(\boldsymbol{\phi}) - \mathbf{C})^\top R D R^\top (\mathbf{P}(\boldsymbol{\phi}) - \mathbf{C}) = 1 \quad (98)$$

where  $R = [\mathbf{U}_0 \cdots \mathbf{U}_{n-1}]$  is a rotation matrix and  $D = \text{Diag}(1/e_0^2, \dots, 1/e_{n-1}^2)$ .

The mean of the hyperellipsoid points is the average of those points in the sense of integration over the hyperellipsoid,

$$\begin{aligned} \boldsymbol{\mu} &= \frac{1}{A_n} \int_{S^n} \mathbf{P}(\boldsymbol{\phi}) dA_n \\ &= \frac{1}{A_n} \int_{S^n} \left( \mathbf{C} + \sum_{j=0}^{n-1} e_j x_j(\boldsymbol{\phi}) \mathbf{U}_j \right) dA_n \\ &= \mathbf{C} + \sum_{j=0}^{n-1} \left( \int_{S^n} x_j(\boldsymbol{\phi}) dA_n \right) \mathbf{U}_j \\ &= \mathbf{C} \end{aligned} \quad (99)$$

All the integrals evaluate to zero. It is natural to choose the initial center  $\mathbf{C}$  to be the mean  $\boldsymbol{\mu}$  of the sample points.

The covariance matrix for the hyperellipsoid is

$$\begin{aligned} \mathcal{C} &= \frac{1}{A_n} \int_{S^n} (\mathbf{P}(\boldsymbol{\phi}) - \mathbf{C})(\mathbf{P}(\boldsymbol{\phi}) - \mathbf{C})^\top dA_n \\ &= \frac{1}{A_n} \int_{S^n} \left( \sum_{j=0}^{n-1} e_j x_j(\boldsymbol{\phi}) \mathbf{U}_j \right) \left( \sum_{k=0}^{n-1} e_k x_k(\boldsymbol{\phi}) \mathbf{U}_k \right)^\top dA_n \\ &= \sum_{j=0}^{n-1} \sum_{k=0}^{n-1} \frac{e_j e_k}{A_n} \left( \int_{S^n} x_j(\boldsymbol{\phi}) x_k(\boldsymbol{\phi}) dA_n \right) \mathbf{U}_j \mathbf{U}_k^\top \\ &= \sum_{j=0}^{n-1} \frac{e_j^2}{n} \mathbf{U}_j \mathbf{U}_j^\top \\ &= R \text{Diag} \left( \frac{e_0^2}{n}, \dots, \frac{e_{n-1}^2}{n} \right) R^\top \end{aligned} \quad (100)$$

The integrals evaluate to  $\int_{S^n} x_j(\boldsymbol{\phi}) x_k(\boldsymbol{\phi}) dA_n = 0$  for  $j \neq k$  and  $\int_{S^n} x_j^2(\boldsymbol{\phi}) dA_n = A_n/n$ . Computing  $\mathcal{C}$  to be the covariance matrix for the sample points, factor  $\mathcal{C} = R \text{Diag}(\lambda_0, \dots, \lambda_{n-1}) R^\top$  using an eigendecomposition. The estimates for the axis extents are then  $e_j = \sqrt{n\lambda_j}$  for  $0 \leq j < n$ .

In summary, compute

$$\boldsymbol{\mu} = \frac{1}{m} \sum_{i=1}^m \mathbf{X}_i, \quad \mathbf{C} = \frac{1}{m} \sum_{i=1}^m (\mathbf{X}_i - \boldsymbol{\mu})(\mathbf{X}_i - \boldsymbol{\mu})^\top \quad (101)$$

Choose the initial center  $\mathbf{C} = \boldsymbol{\mu}$ . Factor  $\mathbf{C} = R \text{Diag}(\lambda_0, \dots, \lambda_{n-1}) R^\top$  where  $R$  is a rotation matrix. Choose the initial orientation to be  $R$ . Choose the initial diagonal matrix  $D = \text{Diag}(1/(n\lambda_0), \dots, 1/(n\lambda_{n-1}))$ .

### 6.3 Fitting Using Derivative-Based Methods

Three methods for minimization via iteration and that use derivative information for the error function are gradient descent, the Gauss–Newton method and the Levenberg–Marquardt method. Gradient descent is a conservative method that can require a large number of iterations until convergence. The Gauss–Newton method is an aggressive method, one that hopes to converge in a small number of iterations. The Levenberg–Marquardt method is adaptive and uses step size that vary between those of gradient descent and those of Gauss–Newton.

Each of the methods has a well-known general framework. All that is required for using the framework is to implement the error function and its derivatives. The error function is equation (63). The derivative functions are equations (64), (65) and (70).

### 6.4 Fitting by Reduction to Polynomial Equations

We are now ready to formulate equations for computing where the gradient of  $E(\mathbf{C}, R, D)$  is the zero vector. Setting equation (64) to zero and using the invertibility of  $RDR^\top$ , we obtain

$$\sum_{i=1}^m \left( \boldsymbol{\Delta}_i^\top RDR^\top \boldsymbol{\Delta}_i - 1 \right) \boldsymbol{\Delta}_i = \mathbf{0} \quad (102)$$

The  $n(n+1)/2$  equations (71) form a linear system of  $n(n+1)/2$  unknowns whose matrix of coefficients is invertible because of the orthonormality of the columns of  $R$ . The vector of unknowns is therefore the zero vector, summarized by

$$\mathbf{U}_j^\top \left( \sum_{i=1}^m \left( \boldsymbol{\Delta}_i^\top RDR^\top \boldsymbol{\Delta}_i - 1 \right) \boldsymbol{\Delta}_i \boldsymbol{\Delta}_i^\top \right) \mathbf{U}_k = 0 \quad (103)$$

for  $0 \leq j \leq k < n$ . The matrices  $\boldsymbol{\Delta}_i \boldsymbol{\Delta}_i^\top$  are symmetric, so equation (103) is forcing all the entries of the weighted sum of these matrices to be zero; therefore,

$$\sum_{i=1}^m \left( \boldsymbol{\Delta}_i^\top RDR^\top \boldsymbol{\Delta}_i - 1 \right) \boldsymbol{\Delta}_i \boldsymbol{\Delta}_i^\top = 0 \quad (104)$$

The solutions to the nonlinear equations (102) and (104) must be solved for  $\mathbf{C}$ ,  $R$  and  $D$  that minimize  $E(\mathbf{C}, R, D)$ . Define  $P = [p_{ij}] = RDR^\top$ , which is a positive definite matrix. As a symmetric matrix, it has  $n(n+1)/2$  independent entries (the diagonal and upper triangular elements). Recall that  $\boldsymbol{\Delta}_i = \mathbf{X}_i - \mathbf{C}$ . The aforementioned two equations are linear in the entries of  $P$ . Equation (102) is cubic in the entries of  $\mathbf{C}$  and equation (104) is quartic in the entries of  $\mathbf{C}$ . Therefore, we have a system of polynomial equations that can be reduced to a single polynomial equation using elimination theory and resultants.

Define  $P = RDR^T$ , a positive definite matrix. Equation (104) can be viewed as a system of  $n(n+1)/2$  linear equations in the  $n(n+1)/2$  unique entries of  $P$ . Only the upper-triangular entries of the various matrices need to be considered. Define  $\mathbf{p}$  to be the upper-triangular elements of  $P$  listed in row-major order except that the off-diagonal terms are multiplied by 2. The system can be written as  $A\mathbf{p} = \mathbf{b}$ , where  $A$  and  $\mathbf{b}$  are dependent only on  $\mathbf{C}$  and the sample points. Generally, we expect  $A$  to be invertible. Symbolically, invert  $A$  to obtain  $A^{\text{adj}}/\det(A)$ , where  $A^{\text{adj}}$  is the adjugate matrix of  $A$  (the transpose of matrix of cofactors), and compute  $\mathbf{p} = A^{\text{adj}}\mathbf{b}/\det(A)$ . Substitute this into equation (102) and multiply through to obtain a system of polynomial equations that depend only on  $\mathbf{C}$ . These equations can be reduced to a single polynomial equation in one of the components of  $\mathbf{C}$ . The roots of that polynomial are computed, and other components of  $\mathbf{C}$  are evaluated based on the intermediate steps of the reduction. For all roots, keep track of the one that minimizes  $E$ .

### 6.4.1 Reduction for an Ellipse

Let the samples be  $(x_i, y_i)$  and let the center be  $\mathbf{C} = (c_0, c_1)$ . Define  $s_{jk} = \sum_{i=1}^m (x_i - c_0)^j (y_i - c_1)^k$ . The matrix  $A$  and vectors  $\mathbf{b}$  and  $\mathbf{p}$  mentioned in the general discussion are the following for dimension 2. Notice that the off-diagonal entry of  $P = RDR^T$  is scaled by 2,

$$A = \begin{bmatrix} s_{40} & s_{31} & s_{22} \\ s_{31} & s_{22} & s_{13} \\ s_{22} & s_{13} & s_{04} \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} s_{20} \\ s_{11} \\ s_{02} \end{bmatrix}, \quad \mathbf{p} = \begin{bmatrix} p_{00} \\ 2p_{01} \\ p_{11} \end{bmatrix} \quad (105)$$

Define

$$\mathbf{r}_0 = \begin{bmatrix} s_{30} \\ s_{21} \\ s_{12} \end{bmatrix}, \quad \mathbf{r}_1 = \begin{bmatrix} s_{21} \\ s_{12} \\ s_{03} \end{bmatrix} \quad (106)$$

Solve  $A\mathbf{p} = \mathbf{b}$  for  $\mathbf{p}$  and substitute into equation (102) to obtain

$$\mathbf{r}_0^T A^{\text{adj}} \mathbf{b} - s_{10} \det(A) = 0, \quad \mathbf{r}_1^T A^{\text{adj}} \mathbf{b} - s_{01} \det(A) = 0 \quad (107)$$

The degree of the  $\mathbf{r}_i$  is 3, the degree of the elements of  $A^{\text{adj}}$  is 8, the degree of  $\mathbf{b}$  is 2 and the degree of  $\det(A)$  is 12. Therefore, each of these equations is a polynomial of degree 13 in the 2 components of  $\mathbf{C}$ . The reduced polynomial of one component will be of symbolic degree  $13^2 = 169$ . Naturally, you need access to a robust polynomial root finder. The computation of the coefficients can be performed using arbitrary precision arithmetic (for exact rational arithmetic). After that, computing the roots will require modern-day methods for obtaining bounding intervals for roots, after which a switch to floating-point arithmetic can be made to bisect the roots.

As an alternative to elimination theory, it might be possible to solve numerically the two polynomial equations by other means. For example, a GPGPU approach could be used that samples  $(c_0, c_1)$  in a rectangle containing the initial center guess and looks for that sample leading to the minimum sum of squares of the polynomials.

### 6.4.2 Reduction for an Ellipsoid

Let the samples be  $(x_i, y_i, z_i)$  and let the center be  $\mathbf{C} = (c_0, c_1, c_2)$ . Define  $s_{jkl} = \sum_{i=1}^m (x_i - c_0)^j (y_i - c_1)^k (z_i - c_2)^\ell$ . The matrix  $A$  and vectors  $\mathbf{b}$  and  $\mathbf{p}$  mentioned in the general discussion are the following for dimension 3. Notice that the off-diagonal entries of  $P = RDR^\top$  are scaled by 2,

$$A = \begin{bmatrix} s_{400} & s_{310} & s_{301} & s_{220} & s_{211} & s_{202} \\ s_{310} & s_{220} & s_{211} & s_{130} & s_{121} & s_{112} \\ s_{301} & s_{211} & s_{202} & s_{121} & s_{112} & s_{103} \\ s_{220} & s_{130} & s_{121} & s_{040} & s_{031} & s_{022} \\ s_{211} & s_{121} & s_{112} & s_{031} & s_{022} & s_{013} \\ s_{202} & s_{112} & s_{103} & s_{022} & s_{013} & s_{004} \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} s_{200} \\ s_{110} \\ s_{101} \\ s_{020} \\ s_{011} \\ s_{002} \end{bmatrix}, \quad \mathbf{p} = \begin{bmatrix} p_{00} \\ 2p_{01} \\ 2p_{02} \\ p_{11} \\ 2p_{12} \\ p_{22} \end{bmatrix} \quad (108)$$

Define

$$\mathbf{r}_0 = \begin{bmatrix} s_{300} \\ s_{210} \\ s_{201} \\ s_{120} \\ s_{111} \\ s_{102} \end{bmatrix}, \quad \mathbf{r}_1 = \begin{bmatrix} s_{210} \\ s_{120} \\ s_{111} \\ s_{030} \\ s_{021} \\ s_{012} \end{bmatrix}, \quad \mathbf{r}_2 = \begin{bmatrix} s_{201} \\ s_{111} \\ s_{102} \\ s_{021} \\ s_{012} \\ s_{003} \end{bmatrix} \quad (109)$$

Solve  $A\mathbf{p} = \mathbf{b}$  for  $\mathbf{p}$  and substitute into equation (102) to obtain

$$\mathbf{r}_0^\top A^{\text{adj}} \mathbf{b} - s_{100} \det(A) = 0, \quad \mathbf{r}_1^\top A^{\text{adj}} \mathbf{b} - s_{010} \det(A) = 0, \quad \mathbf{r}_2^\top A^{\text{adj}} \mathbf{b} - s_{001} \det(A) = 0 \quad (110)$$

The degree of the  $\mathbf{r}_i$  is 3, the degree of the elements of  $A^{\text{adj}}$  is 20, the degree of  $\mathbf{b}$  is 2 and the degree of  $\det(A)$  is 24. Therefore, each of these equations is a polynomial of degree 25 in the 3 components of  $\mathbf{C}$ . The reduced polynomial of one component will be of symbolic degree  $25^3 = 15625$ . The large degree makes this not an attractive approach.

As an alternative to elimination theory, it might be possible to solve numerically the three polynomial equations by other means. For example, a GPGPU approach could be used that samples  $(c_0, c_1, c_2)$  in a rectangular solid containing the initial center guess and looks for that sample leading to the minimum sum of squares of the polynomials.

## 7 Fitting a Cylinder to 3D Points

This document describes an algorithm for fitting a set of 3D points with a cylinder. The assumption is that the underlying data is modeled by a cylinder and that errors have caused the points not to be exactly on the cylinder. You could very well try to fit a random set of points, but the algorithm is not guaranteed to produce a meaningful solution.

## 7.1 Representation of a Cylinder

An infinite cylinder is specified by an axis containing a point  $\mathbf{C}$  and having unit-length direction  $\mathbf{W}$ . The radius of the cylinder is  $r > 0$ . Two more unit-length vectors  $\mathbf{U}$  and  $\mathbf{V}$  may be defined so that  $\{\mathbf{U}, \mathbf{V}, \mathbf{W}\}$  is a right-handed orthonormal set; that is, all vectors are unit-length, mutually perpendicular, and with  $\mathbf{U} \times \mathbf{V} = \mathbf{W}$ ,  $\mathbf{V} \times \mathbf{W} = \mathbf{U}$ , and  $\mathbf{W} \times \mathbf{U} = \mathbf{V}$ . Any point  $\mathbf{X}$  may be written uniquely as

$$\mathbf{X} = \mathbf{C} + y_0\mathbf{U} + y_1\mathbf{V} + y_2\mathbf{W} = \mathbf{C} + R\mathbf{Y} \quad (111)$$

where  $R$  is a rotation matrix whose columns are  $\mathbf{U}$ ,  $\mathbf{V}$ , and  $\mathbf{W}$  and where  $\mathbf{Y}$  is a column vector whose rows are  $y_0$ ,  $y_1$ , and  $y_2$ . To be on the cylinder, we need

$$\begin{aligned} r^2 &= y_0^2 + y_1^2 \\ &= (\mathbf{U} \cdot (\mathbf{X} - \mathbf{C}))^2 + (\mathbf{V} \cdot (\mathbf{X} - \mathbf{C}))^2 \\ &= (\mathbf{X} - \mathbf{C})^\top (\mathbf{U}\mathbf{U}^\top + \mathbf{V}\mathbf{V}^\top) (\mathbf{X} - \mathbf{C}) \\ &= (\mathbf{X} - \mathbf{C})^\top (I - \mathbf{W}\mathbf{W}^\top) (\mathbf{X} - \mathbf{C}) \end{aligned} \quad (112)$$

where  $I$  is the identity matrix. Because the unit-length vectors form an orthonormal set, it is necessary that  $I = \mathbf{U}\mathbf{U}^\top + \mathbf{V}\mathbf{V}^\top + \mathbf{W}\mathbf{W}^\top$ . A finite cylinder is obtained by bounding the points in the axis direction,

$$|y_2| = |\mathbf{W} \cdot (\mathbf{X} - \mathbf{C})| \leq h/2 \quad (113)$$

where  $h > 0$  is the height of the cylinder.

## 7.2 The Least-Squares Error Function

Let  $\{\mathbf{X}_i\}_{i=1}^n$  be the input point set. An error function for a cylinder fit based on Equation (112) is

$$E(r^2, \mathbf{C}, \mathbf{W}) = \sum_{i=1}^n [F(\mathbf{X}_i; r^2, \mathbf{C}, \mathbf{W})]^2 = \sum_{i=1}^n \left[ (\mathbf{X}_i - \mathbf{C})^\top (I - \mathbf{W}\mathbf{W}^\top) (\mathbf{X}_i - \mathbf{C}) - r^2 \right]^2 \quad (114)$$

where the cylinder axis is a line containing point  $\mathbf{C}$  and having unit-length direction  $\mathbf{W}$  and the cylinder radius is  $r$ . Thus, the error function involves 6 parameters: 1 for the squared radius  $r^2$ , 3 for the point  $\mathbf{C}$ , and 2 for the unit-length direction  $\mathbf{W}$ . These parameters form the 6-tuple  $\mathbf{q}$  in the generic discussion presented previously.

For numerical robustness, it is advisable to subtract the sample mean  $\mathbf{A} = (\sum_{i=1}^n \mathbf{X}_i)/n$  from the samples,  $\mathbf{X}_i \leftarrow \mathbf{X}_i - \mathbf{A}$ . This preconditioning is assumed in the mathematical derivations to follow, in which case  $\sum_{i=1}^n \mathbf{X}_i = \mathbf{0}$ .

In the following discussion, define

$$P = I - \mathbf{W}\mathbf{W}^\top, \quad r_i^2 = (\mathbf{C} - \mathbf{X}_i)^\top P (\mathbf{C} - \mathbf{X}_i) \quad (115)$$

The matrix  $P$  represents a projection onto a plane with normal  $\mathbf{W}$ , so  $P^2 = P$  and depends only on the direction  $\mathbf{W}$ . The term  $r_i^2$  depends on the center  $\mathbf{C}$  and the direction  $\mathbf{W}$ . The error function is written concisely as  $E = \sum_{i=1}^n (r_i^2 - r^2)^2$ .



### 7.3 An Equation for the Radius

The partial derivative of the error function with respect to the squared radius is  $\partial E/\partial r^2 = -2\sum_{i=1}^n(r_i^2 - r^2)$ . Setting this to zero, we have the constraint

$$0 = \sum_{i=1}^n(r_i^2 - r^2) \quad (116)$$

which leads to

$$r^2 = \frac{1}{n} \sum_{i=1}^n r_i^2 \quad (117)$$

Thus, the squared radius is the average of the squared distances of the projections of  $\mathbf{X}_i - \mathbf{C}$  onto a plane containing  $\mathbf{C}$  and having normal  $\mathbf{W}$ . The right-hand side depends on the parameters  $\mathbf{C}$  and  $\mathbf{W}$ .

Observe that

$$\begin{aligned} r_i^2 - r^2 &= r_i^2 - \frac{1}{n} \sum_{j=1}^n r_j^2 \\ &= (\mathbf{C} - \mathbf{X}_i)^\top P(\mathbf{C} - \mathbf{X}_i) - \frac{1}{n} \sum_{j=1}^n (\mathbf{C} - \mathbf{X}_j)^\top P(\mathbf{C} - \mathbf{X}_j) \\ &= \mathbf{C}^\top P\mathbf{C} - 2\mathbf{X}_i^\top P\mathbf{C} + \mathbf{X}_i^\top P\mathbf{X}_i - \frac{1}{n} \sum_{j=1}^n \left( \mathbf{C}^\top P\mathbf{C} - 2\mathbf{X}_j^\top P\mathbf{C} + \mathbf{X}_j^\top P\mathbf{X}_j \right) \\ &= \mathbf{C}^\top P\mathbf{C} - 2\mathbf{X}_i^\top P\mathbf{C} + \mathbf{X}_i^\top P\mathbf{X}_i - \mathbf{C}^\top P\mathbf{C} + \frac{2}{n} \left( \sum_{j=1}^n \mathbf{X}_j^\top \right) P\mathbf{C} - \frac{1}{n} \sum_{j=1}^n \mathbf{X}_j^\top P\mathbf{X}_j \\ &= \left( \frac{1}{n} \sum_{j=1}^n \mathbf{X}_j^\top - \mathbf{X}_i^\top \right) 2P\mathbf{C} + \left( \mathbf{X}_i^\top P\mathbf{X}_i - \frac{1}{n} \sum_{j=1}^n \mathbf{X}_j^\top P\mathbf{X}_j \right) \\ &= -\mathbf{X}_i^\top 2P\mathbf{C} + \left( \mathbf{X}_i^\top P\mathbf{X}_i - \frac{1}{n} \sum_{j=1}^n \mathbf{X}_j^\top P\mathbf{X}_j \right) \end{aligned} \quad (118)$$

The last equality is based on the precondition  $\sum_{j=1}^n \mathbf{X}_j = \mathbf{0}$ .

### 7.4 An Equation for the Center

The partial derivative with respect to the center is  $\partial E/\partial \mathbf{C} = -4\sum_{i=1}^n(r_i^2 - r^2)P(\mathbf{X}_i - \mathbf{C})$ . Setting this to zero, we have the constraint

$$\begin{aligned} 0 &= \sum_{i=1}^n(r_i^2 - r^2)P(\mathbf{X}_i - \mathbf{C}) \\ &= \sum_{i=1}^n(r_i^2 - r^2)P\mathbf{X}_i - \left[ \sum_{i=1}^n(r_i^2 - r^2) \right] P\mathbf{C} \\ &= \sum_{i=1}^n(r_i^2 - r^2)P\mathbf{X}_i \end{aligned} \quad (119)$$

where the last equality is a consequence of Equation (116). Multiply equation (118) by  $P\mathbf{X}_i$ , sum over  $i$ , and use equation (119) to obtain

$$\begin{aligned} 0 &= -2P \left( \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top \right) P\mathbf{C} + \sum_{i=1}^n \left( \mathbf{X}_i^\top P\mathbf{X}_i \right) P\mathbf{X}_i - \left( \frac{1}{n} \sum_{j=1}^n \mathbf{X}_j^\top P\mathbf{X}_j \right) P \sum_{i=1}^n \mathbf{X}_i \\ &= -2P \left( \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top \right) P\mathbf{C} + \sum_{i=1}^n \left( \mathbf{X}_i^\top P\mathbf{X}_i \right) P\mathbf{X}_i \end{aligned} \quad (120)$$

where the last equality is based on the precondition  $\sum_{i=1}^n \mathbf{X}_i = \mathbf{0}$ . We wish to solve this equation for  $\mathbf{C}$ , but observe that  $\mathbf{C} + t\mathbf{W}$  are valid centers for all  $t$ . It is sufficient to compute a center that has no component

in the  $\mathbf{W}$ -direction; that is, we may construct a point for which  $\mathbf{C} = P\mathbf{C}$ . It suffices to solve Equation (120) for  $P\mathbf{C}$  written as the linear system  $A(P\mathbf{C}) = \mathbf{B}/2$  where

$$A = P \left( \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top \right) P, \quad \mathbf{B} = \frac{1}{n} \sum_{i=1}^n \left( \mathbf{X}_i^\top P \mathbf{X}_i \right) P \mathbf{X}_i \quad (121)$$

The projection matrix is symmetric,  $P = P^\top$ , a condition that leads to the right-hand side of the equation defining  $A$ . We have used  $P = P^2$  to introduce an additional  $P$  factor,  $\mathbf{X}_i^\top P \mathbf{X}_i = \mathbf{X}_i^\top P^2 \mathbf{X}_i = \mathbf{X}_i^\top P^\top P \mathbf{X}_i$ , which leads to the right-hand side of the equation defining  $\mathbf{B}$ .

The matrix  $A$  is singular because the projection matrix  $P$  is singular, so we cannot directly invert  $A$  to solve the equation. The linear system involves terms that live only in the plane perpendicular to  $\mathbf{W}$ , so in fact the linear system reduces to two equations in two unknowns in the projection space and is solvable as long as the coefficient matrix is invertible.

Choose  $\mathbf{U}$  and  $\mathbf{V}$  so that  $\{\mathbf{U}, \mathbf{V}, \mathbf{W}\}$  is a right-handed orthonormal set; then  $P\mathbf{X}_i = \mu_i \mathbf{U} + \nu_i \mathbf{V}$  and  $P\mathbf{C} = k_0 \mathbf{U} + k_1 \mathbf{V}$ , where  $\mu_i = \mathbf{U} \cdot \mathbf{X}_i$ ,  $\nu_i = \mathbf{V} \cdot \mathbf{X}_i$ ,  $k_0 = \mathbf{U} \cdot P\mathbf{C}$ , and  $k_1 = \mathbf{V} \cdot P\mathbf{C}$ . The matrix  $A$  becomes

$$A = \left( \frac{1}{n} \sum_{i=1}^n \mu_i^2 \right) \mathbf{U} \mathbf{U}^\top + \left( \frac{1}{n} \sum_{i=1}^n \mu_i \nu_i \right) (\mathbf{U} \mathbf{V}^\top + \mathbf{V} \mathbf{U}^\top) + \left( \frac{1}{n} \sum_{i=1}^n \nu_i^2 \right) \mathbf{V} \mathbf{V}^\top \quad (122)$$

and the vector  $\mathbf{B}$  becomes

$$\mathbf{B} = \frac{1}{n} \sum_{i=1}^n (\mu_i^2 + \nu_i^2) (\mu_i \mathbf{U} + \nu_i \mathbf{V}) \quad (123)$$

The vector  $A(P\mathbf{C})$  becomes

$$A(P\mathbf{C}) = \left( \frac{k_0}{n} \sum_{i=1}^n \mu_i^2 + \frac{k_1}{n} \sum_{i=1}^n \mu_i \nu_i \right) \mathbf{U} + \left( \frac{k_0}{n} \sum_{i=1}^n \mu_i \nu_i + \frac{k_1}{n} \sum_{i=1}^n \nu_i^2 \right) \mathbf{V} \quad (124)$$

Equating this to  $\mathbf{B}/2$  and grouping the coefficients for  $\mathbf{U}$  and  $\mathbf{V}$  leads to the linear system

$$\begin{bmatrix} \frac{1}{n} \sum_{i=1}^n \mu_i^2 & \frac{1}{n} \sum_{i=1}^n \mu_i \nu_i \\ \frac{1}{n} \sum_{i=1}^n \mu_i \nu_i & \frac{1}{n} \sum_{i=1}^n \nu_i^2 \end{bmatrix} \begin{bmatrix} k_0 \\ k_1 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n (\mu_i^2 + \nu_i^2) \mu_i \\ \frac{1}{n} \sum_{i=1}^n (\mu_i^2 + \nu_i^2) \nu_i \end{bmatrix} \quad (125)$$

The coefficient matrix is the covariance matrix of the projection of the samples onto the plane perpendicular to  $\mathbf{W}$ . Intuitively, this matrix is invertible as long as the projections do not lie on a line. If the matrix is singular (or nearly singular numerically), the original samples lie on a plane (or nearly lie on a plane numerically). They are not fitted well by a cylinder or, if you prefer, they are fitted by a cylinder with infinite radius.

The matrix system of Equation (125) has solution

$$\begin{bmatrix} k_0 \\ k_1 \end{bmatrix} = \frac{1}{2 \left( \frac{1}{n} \sum_{i=1}^n \mu_i^2 \frac{1}{n} \sum_{i=1}^n \nu_i^2 - \left( \frac{1}{n} \sum_{i=1}^n \mu_i \nu_i \right)^2 \right)} \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n \nu_i^2 & -\frac{1}{n} \sum_{i=1}^n \mu_i \nu_i \\ -\frac{1}{n} \sum_{i=1}^n \mu_i \nu_i & \frac{1}{n} \sum_{i=1}^n \mu_i^2 \end{bmatrix} \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n (\mu_i^2 + \nu_i^2) \mu_i \\ \frac{1}{n} \sum_{i=1}^n (\mu_i^2 + \nu_i^2) \nu_i \end{bmatrix} \quad (126)$$

which produces the cylinder center  $P\mathbf{C} = k_0 \mathbf{U} + k_1 \mathbf{V}$ ; use this instead of  $\mathbf{C}$  in Equation (112).

Although the solution appears to depend on the choice of  $\mathbf{U}$  and  $\mathbf{V}$ , it does not. Let  $\mathbf{W} = (w_0, w_1, w_2)$  and define the skew symmetric matrix

$$S = \begin{bmatrix} 0 & -w_2 & w_1 \\ w_2 & 0 & -w_0 \\ -w_1 & w_0 & 0 \end{bmatrix} \quad (127)$$

By definition of skew symmetry,  $S^\top = -S$ . This matrix represents the cross product operation:  $S\xi = \mathbf{W} \times \xi$  for any vector  $\xi$ . Because  $\{\mathbf{U}, \mathbf{V}, \mathbf{W}\}$  is a right-handed orthonormal set, it follows that  $S\mathbf{U} = \mathbf{V}$  and  $S\mathbf{V} = -\mathbf{U}$ . It can be shown also that  $S = \mathbf{V}\mathbf{U}^\top - \mathbf{U}\mathbf{V}^\top$ . Define matrix  $\hat{A}$  by

$$\hat{A} = \left( \frac{1}{n} \sum_{i=1}^n \nu_i^2 \right) \mathbf{U}\mathbf{U}^\top - \left( \frac{1}{n} \sum_{i=1}^n \mu_i \nu_i \right) (\mathbf{U}\mathbf{V}^\top + \mathbf{V}\mathbf{U}^\top) + \left( \frac{1}{n} \sum_{i=1}^n \mu_i^2 \right) \mathbf{V}\mathbf{V}^\top = SAS^\top \quad (128)$$

Effectively, this generates a  $2 \times 2$  matrix that is the adjugate of the  $2 \times 2$  matrix representing  $A$ . It has the property

$$\hat{A}A = \delta P, \quad \delta = \sum_{i=1}^n \mu_i^2 \sum_{i=1}^n \nu_i^2 - \left( \sum_{i=1}^n \mu_i \nu_i \right)^2 \quad (129)$$

The trace of a matrix is the sum of the diagonal entries. Observe that  $\text{trace}(P) = 2$  because  $|\mathbf{W}|^2 = 1$ . Taking the trace of  $\hat{A}A = \delta P$ , we obtain  $2\delta = \text{trace}(\hat{A}A)$ . The cylinder center is obtained by multiplying  $A(PC) = \mathbf{B}/2$  by  $\hat{A}$ , using the definitions in equation (121) and using equation (129),

$$PC = \frac{\hat{A}}{\text{trace}(\hat{A}A)} \left( \frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i^\top P \mathbf{X}_i) P \mathbf{X}_i \right) = \frac{\hat{A}}{\text{trace}(\hat{A}A)} \left( \frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i^\top P \mathbf{X}_i) \mathbf{X}_i \right) \quad (130)$$

where the last equality uses  $\hat{A}P = \hat{A}$  because  $S^\top P = S^\top$ . Equation (130) is independent of  $\mathbf{U}$  and  $\mathbf{V}$  but dependent on  $\mathbf{W}$ .

## 7.5 An Equation for the Direction

Let the direction be parameterized as  $\mathbf{W}(\mathbf{s}) = (w_0(\mathbf{s}), w_1(\mathbf{s}), w_2(\mathbf{s}))$ , where  $\mathbf{s}$  is a 2-dimensional parameter. For example, spherical coordinates is such a parameterization:  $\mathbf{W} = (\cos s_0 \sin s_1, \sin s_0 \sin s_1, \cos s_1)$  for  $s_0 \in [0, 2\pi)$  and  $s_1 \in [0, \pi/2]$ , where  $w_2(s_0, s_1) \geq 0$ . The partial derivatives of  $E$  are

$$\frac{\partial E}{\partial s_k} = 2 \sum_{i=1}^n (r_i^2 - r^2) (\mathbf{C} - \mathbf{X}_i)^\top \frac{\partial P}{\partial s_k} (\mathbf{C} - \mathbf{X}_i) \quad (131)$$

Solving  $\partial E / \partial s_k = 0$  in closed form is not tractable. It is possible to generate a system of polynomial equations in the components of  $\mathbf{W}$ , use elimination theory to obtain a polynomial in one variable, and then find its roots. This approach is generally tedious and not robust numerically.

Alternatively, we can skip root finding for  $\partial E / \partial s_k = 0$ , instead substituting Equations (118) and (130) directly into the error function  $E/n = \frac{1}{n} \sum_{i=1}^n (r_i^2 - r^2)^2$  to obtain a nonnegative function,

$$G(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \left[ \mathbf{X}_i^\top P \mathbf{X}_i - \frac{1}{n} \sum_{j=1}^n \mathbf{X}_j^\top P \mathbf{X}_j - 2\mathbf{X}_i^\top \frac{\hat{A}}{\text{trace}(\hat{A}A)} \left( \frac{1}{n} \sum_{j=1}^n (\mathbf{X}_j^\top P \mathbf{X}_j) \mathbf{X}_j \right) \right]^2 \quad (132)$$

A numerical algorithm for locating the minimum of  $G$  can be used. Or, as is shown in the sample code, the domain for  $(s_0, s_1)$  may be partitioned into samples at which  $G$  is evaluated. The sample producing the minimum  $G$ -value determines a reasonable direction  $\mathbf{W}$ . The center  $\mathbf{C}$  and squared radius  $r^2$  are inherent in the evaluation of  $G$ , so in the end we have a fitted cylinder. The evaluations of  $G$  are expensive for a large number of samples: equation (132) contains summations in a term that is then squared, followed by another summation.

Some algebraic manipulation lead to encapsulating the summations, allowing us to precompute the summations and represent  $G(\mathbf{W})$  as a rational polynomial of the components of  $\mathbf{W}$ . This approach increases the performance on the CPU. It also allows an efficient implementation for massively parallel performance on the GPU—one GPU thread per direction vector that is sampled from a hemisphere.

The projection matrix  $P$  is determined by its upper-triangular elements  $\mathbf{p} = (p_{00}, p_{01}, p_{02}, p_{11}, p_{12}, p_{22})$ . We can write the following, where  $\mathbf{p}$  is represented as a  $6 \times 1$  vector,

$$\mathbf{X}_i^\top P \mathbf{X}_i - \frac{1}{n} \sum_{j=1}^n \mathbf{X}_j^\top P \mathbf{X}_j = \mathbf{p} \cdot (\boldsymbol{\xi}_i - \boldsymbol{\mu}) = \mathbf{p} \cdot \boldsymbol{\delta}_i \quad (133)$$

The  $6 \times 1$  vectors  $\boldsymbol{\xi}_i$ ,  $\boldsymbol{\mu}$  and  $\boldsymbol{\delta}_i$  are defined next (written as 6-tuples). As a 3-tuple, let  $\mathbf{X}_i = (x_i, y_i, z_i)$ ,

$$\boldsymbol{\xi}_i = (x_i^2, 2x_i y_i, 2x_i z_i, y_i^2, 2y_i z_i, z_i^2), \quad \boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \boldsymbol{\xi}_i, \quad \boldsymbol{\delta}_i = \boldsymbol{\xi}_i - \boldsymbol{\mu} \quad (134)$$

We can also write

$$\frac{1}{n} \sum_{j=1}^n (\mathbf{X}_j^\top P \mathbf{X}_j) \mathbf{X}_j = \left( \frac{1}{n} \sum_{j=1}^n \mathbf{X}_j \boldsymbol{\xi}_j^\top \right) \mathbf{p} = \left( \frac{1}{n} \sum_{j=1}^n \mathbf{X}_j \boldsymbol{\delta}_j^\top \right) \mathbf{p} \quad (135)$$

The last equality is true because  $\sum_{j=1}^n \mathbf{X}_j = \mathbf{0}$  implies  $\sum_{j=1}^n \mathbf{X}_j \boldsymbol{\mu}^\top = 0$ . Define

$$Q = \hat{A} / \text{trace}(\hat{A}A), \quad F_0 = \frac{1}{n} \sum_{j=1}^n \mathbf{X}_j \mathbf{X}_j^\top, \quad F_1 = \frac{1}{n} \sum_{j=1}^n \mathbf{X}_j \boldsymbol{\delta}_j^\top, \quad F_2 = \frac{1}{n} \sum_{j=1}^n \boldsymbol{\delta}_j \boldsymbol{\delta}_j^\top \quad (136)$$

where  $Q$  and  $F_0$  are  $3 \times 3$ ,  $F_1$  is  $3 \times 6$  and  $F_2$  is  $6 \times 6$ ; then

$$\begin{aligned} G(\mathbf{W}) &= \frac{1}{n} \sum_{i=1}^n \left[ \mathbf{p}^\top \boldsymbol{\delta}_i - 2 \mathbf{X}_i^\top Q F_1 \mathbf{p} \right]^2 \\ &= \frac{1}{n} \sum_{i=1}^n \left[ (\mathbf{p}^\top \boldsymbol{\delta}_i)^2 - 4 (\mathbf{p}^\top \boldsymbol{\delta}_i) (\mathbf{X}_i^\top Q F_1 \mathbf{p}) + 4 (\mathbf{X}_i^\top Q F_1 \mathbf{p})^2 \right] \\ &= \mathbf{p}^\top F_2 \mathbf{p} - 4 \mathbf{p}^\top F_1^\top Q F_1 \mathbf{p} + 4 \mathbf{p}^\top F_1^\top Q^\top F_0 Q F_1 \mathbf{p} \end{aligned} \quad (137)$$

The precomputations for the input samples  $\{\mathbf{Y}_i\}_{i=1}^n$  is the following and must be done so in the order specified. These steps are independent of direction vectors  $\mathbf{W}$ .

1. Compute  $\mathbf{A} = \frac{1}{n} (\sum_{i=1}^n \mathbf{Y}_i)$ .
2. Compute  $\mathbf{X}_i = \mathbf{Y}_i - \mathbf{A}$  for all  $i$ .
3. Compute  $\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n (x_i^2, 2x_i y_i, 2x_i z_i, y_i^2, 2y_i z_i, z_i^2)$ ,  $\boldsymbol{\delta}_i = (x_i^2, 2x_i y_i, 2x_i z_i, y_i^2, 2y_i z_i, z_i^2) - \boldsymbol{\mu}$  for all  $i$ .

4. Compute  $F_0 = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top$ ,  $F_1 = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \boldsymbol{\delta}_i^\top$ ,  $F_2 = \frac{1}{n} \sum_{i=1}^n \boldsymbol{\delta}_i \boldsymbol{\delta}_i^\top$ .

For each specified direction  $\mathbf{W}$ , do the following steps.

1. Compute the projection matrix  $P = I - \mathbf{W}\mathbf{W}^\top$  and the skew-symmetric matrix  $S$ .
2. Compute  $A = PF_0P$ ,  $\hat{A} = SAS^\top$ ,  $\hat{A}A$ ,  $\text{trace}(\hat{A}A)$ .
3. Compute  $Q = \hat{A} / \text{trace}(\hat{A}A)$ .
4. Store the upper-triangular entries of  $P$  in  $\mathbf{p}$ .
5. Compute  $\boldsymbol{\alpha} = F_1\mathbf{p}$  and  $\boldsymbol{\beta} = Q\boldsymbol{\alpha}$ .
6. Compute  $G = \mathbf{p}^\top F_2 \mathbf{p} - 4\boldsymbol{\alpha}^\top \boldsymbol{\beta} + 4\boldsymbol{\beta}^\top F_0 \boldsymbol{\beta}$ .
7. The corresponding center is  $PC = \boldsymbol{\beta}$ .
8. The corresponding squared radius is  $r^2 = \frac{1}{n} \sum_{i=1}^n r_i^2$  where  $r_i^2 = (\mathbf{P}\mathbf{C} - \mathbf{P}\mathbf{X}_i)^\top (\mathbf{P}\mathbf{C} - \mathbf{P}\mathbf{X}_i)$ . This factors to  $r^2 = \mathbf{p} \cdot \boldsymbol{\mu} + \boldsymbol{\beta}^\top \boldsymbol{\beta}$ .

The sample application that used equation (132) directly was really slow. On an Intel<sup>®</sup> Core<sup>™</sup> i7-6700 CPU at 3.40 GHz, the single-threaded version for 10765 points required 129 seconds and the multithreaded version using 8 hyperthreads required 22 seconds. The evaluation of  $G$  using the precomputed summations is much faster. The single-threaded version required 85 milliseconds and the multithreaded version using 8 hyperthreads required 22 milliseconds.

## 7.6 Fitting for a Specified Direction

Although one may apply root-finding or minimization techniques to estimate the global minimum of  $E$ , as shown previously, in practice it is possible first to obtain a good estimate for the direction  $\mathbf{W}$ . Using this direction, we may solve for  $\mathbf{C} = \mathbf{P}\mathbf{C}$  in Equation (130) and then  $r^2$  in Equation (117).

For example, suppose that the  $\mathbf{X}_i$  are distributed approximately on a section of a cylinder so that the least-squares line that fits the data provides a good estimate for the direction  $\mathbf{W}$ . This vector is a unit-length eigenvector associated with the largest eigenvalue of the covariance matrix  $A = \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top$ . We may use a numerical eigensolver to obtain  $\mathbf{W}$ , and then solve the aforementioned equations for the cylinder center and squared radius.

The distribution can be such that the estimated direction  $\mathbf{W}$  from the covariance matrix is not good, as is shown in the experiments of the next section.

## 7.7 Pseudocode and Experiments

The simplest algorithm to implement involves the minimization of the function  $G$  in Equation (132). An implementation is shown in Listing 14.

**Listing 14.** Pseudocode for preprocessing the sample points. The output  $X[]$  is the array of sample points translated by the average. The other outputs are  $\mu$  for  $\boldsymbol{\mu}$ ,  $F0$  for  $F_0$ ,  $F1$  for  $F_1$  and  $F2$  for  $F_2$ . The pseudocode is also given for evaluating the function  $G(\mathbf{W})$  and generating the corresponding  $PC$  and  $r^2$ .

```

void Preprocess(int n, Vector3 points[n], Vector3 X[n],
Vector3& average, Vector6& mu, Matrix3x3& F0, Matrix3x6& F1, Matrix6x6&F2)
{
    average = { 0, 0, 0 };
    for (int i = 0; i < n; ++i)
    {
        average += points[i];
    }
    average /= n;
    for (int i = 0; i < n; ++i)
    {
        X[i] = points[i] - average;
    }

    Vector6 zero = { 0, 0, 0, 0, 0, 0 };
    Vector6 products[n];
    MakeZero(products);
    mu = zero;
    for (int i = 0; i < n; ++i)
    {
        products[i][0] = X[i][0] * X[i][0];
        products[i][1] = X[i][0] * X[i][1];
        products[i][2] = X[i][0] * X[i][2];
        products[i][3] = X[i][1] * X[i][1];
        products[i][4] = X[i][1] * X[i][2];
        products[i][5] = X[i][2] * X[i][2];
        mu[0] += products[i][0];
        mu[1] += 2 * products[i][1];
        mu[2] += 2 * products[i][2];
        mu[3] += products[i][3];
        mu[4] += 2 * products[i][4];
        mu[5] += products[i][5];
    }
    mu /= n;

    MakeZero(F0);
    MakeZero(F1);
    MakeZero(F2);
    for (int i = 0; i < n; ++i)
    {
        Vector6 delta;
        delta[0] = products[i][0] - mu[0];
        delta[1] = 2 * products[i][1] - mu[1];
        delta[2] = 2 * products[i][2] - mu[2];
        delta[3] = products[i][3] - mu[3];
        delta[4] = 2 * products[i][4] - mu[4];
        delta[5] = products[i][5] - mu[5];
        F0(0, 0) += products[i][0];
        F0(0, 1) += products[i][1];
        F0(0, 2) += products[i][2];
        F0(1, 1) += products[i][3];
        F0(1, 2) += products[i][4];
        F0(2, 2) += products[i][5];
        F1 += OuterProduct(X[i], delta);
        F2 += OuterProduct(delta, delta);
    }
    F0 /= n;
    F0(1, 0) = F0(0, 1);
    F0(2, 0) = F0(0, 2);
    F0(2, 1) = F0(1, 2);
    F1 /= n;
    F2 /= n;
}

Real G(int n, Vector3 X[n], Vector3 mu, Matrix3x3 F0, Matrix3x6 F1, Matrix6x6 F2, Vector3 W,
Vector3& PC, Real& rSqr)
{

```

```

Matrix3x3 P = Matrix3x3::Identity() - OuterProduct(W, W); // P = I - W * W^T
// S = {{0, -w2, w1}, {w2, 0, -w0}, {-w1, w0, 0}}, inner braces are rows
Matrix3x3 S(0, -W[2], W[1], W[2], 0, -W[0], -W[1], W[0], 0);
Matrix3x3 A = P * F0 * P;
Matrix3x3 hatA = -(S * A * S);
Matrix3x3 hatAA = hatA * A;
Real trace = Trace(hatAA);
Matrix3x3 Q = hatA / trace;
Vector6 p = { P(0, 0), P(0, 1), P(0, 2), P(1, 1), P(1, 2), P(2, 2) };
Vector3 alpha = F1 * p;
Vector3 beta = Q * alpha;
Real error = (Dot(p, F2 * p) - 4 * Dot(alpha, beta) + 4 * Dot(beta, F0 * beta)) / n;
PC = beta;
rsqr = Dot(p, mu) + Dot(beta, beta);
return error;
}

```

---

The fitting is performed by searching a large number of directions  $\mathbf{W}$ , as shown in Listing 15.

---

**Listing 15.** Fitting a cylinder to a set of points.

```

// The X[] are the points to be fit. The outputs rSqr, C, and W are the
// cylinder parameters. The function return value is the error function
// evaluated at the cylinder parameters.
Real FitCylinder (int n, Vector3 points[n], Real& rSqr, Vector3& C, Vector3& W)
{
    Vector3 X[n];
    Vector3 mu;
    Matrix3x3 F0;
    Matrix3x6 F1;
    Matrix6x6 F2;
    Preprocess(n, points, X, average, mu, F0, F1, F2);

    // Choose imax and jmax as desired for the level of granularity you
    // want for sampling W vectors on the hemisphere.
    Real minError = infinity;
    W = Vector3::Zero();
    C = Vector3::Zero();
    rSqr = 0;
    for (int j = 0; j <= jmax; ++j)
    {
        Real phi = halfPi * j / jmax; // in [0, pi/2]
        Real csphi = cos(phi), snphi = sin(phi);
        for (int i = 0; i < imax; ++i)
        {
            Real theta = twoPi * i / imax; // in [0, 2*pi)
            Real csttheta = cos(theta), sntheta = sin(theta);
            Vector3 currentW(csttheta * snphi, sntheta * snphi, csphi);
            Vector3 currentC;
            Real currentRSqr;
            Real error = G(n, X, mu, F0, F1, F2, currentW, currentC, currentRSqr);
            if (error < minError)
            {
                minError = error;
                W = currentW;
                C = currentC;
                rSqr = currentRSqr;
            }
        }
    }

    // Translate the center to the original coordinate system.
    C += average;
}

```

```
    }  
    return minError;  
}
```

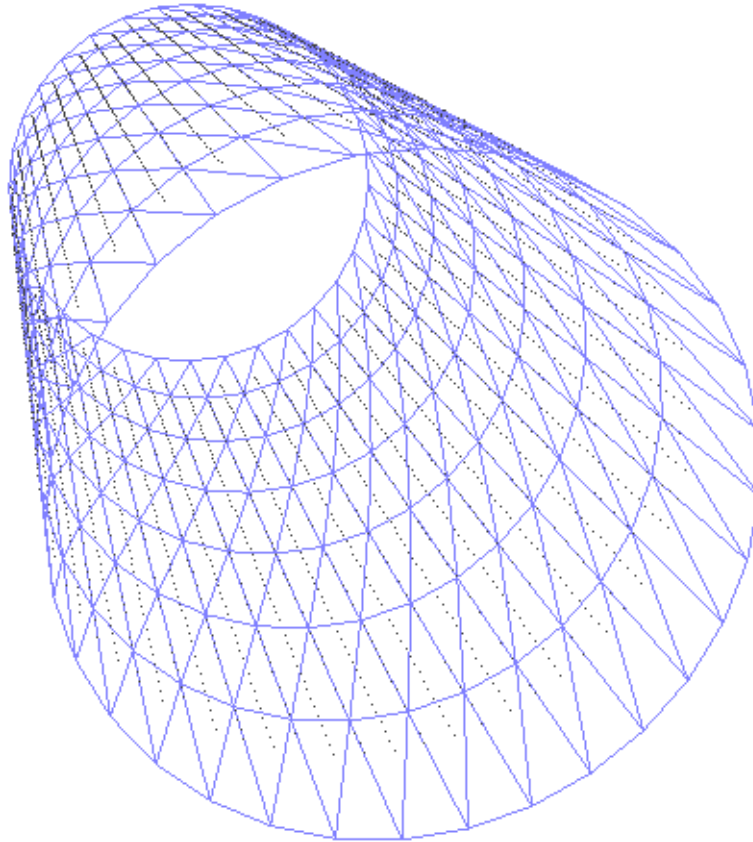
---

The following experiments show how well the cylinder fitting works.

A regular lattice of  $64 \times 65$  samples were chosen for a cylinder centered at the origin, with axis direction  $(0, 0, 1)$ , with radius 1, and height 4. The samples are  $(\cos \theta_j, \sin \theta_j, z_{ij})$ , where  $\theta_j = 2\pi j/64$  and  $z_{ij} = -2 + 4i/64$  for  $0 \leq i \leq 64$  and  $0 \leq j < 64$ . The fitted center, radius, and axis direction are the actual ones (modulo numerical round-off errors). Figure 1 shows a rendering of the points (in black) and a wire frame of the fitted cylinder (in blue).

---

**Figure 1.** Fitting of samples on a cylinder ring that is cut perpendicular to the cylinder axis.



---

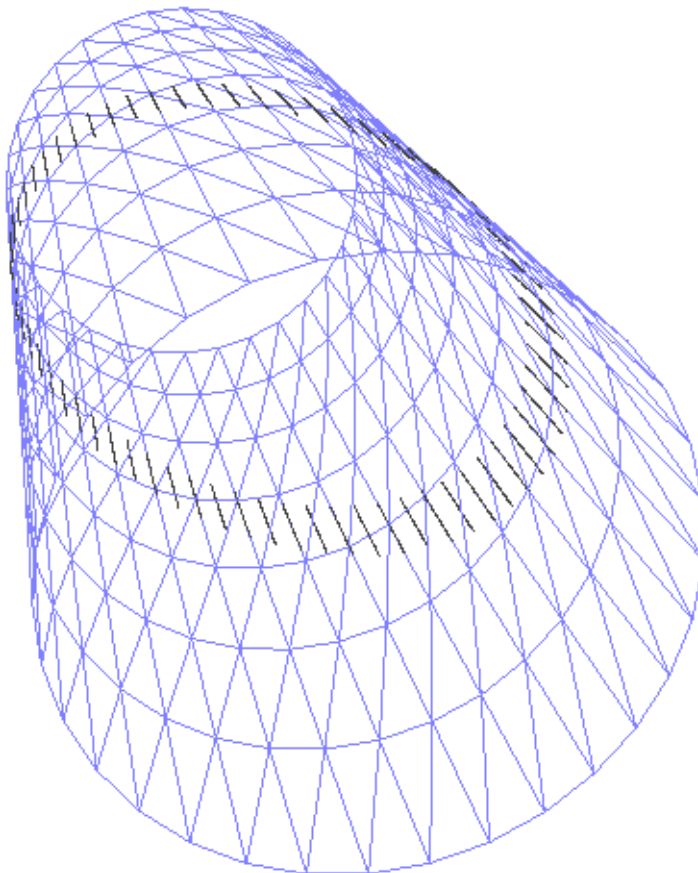
A lattice of samples was chosen for the same cylinder but the samples are skewed to lie on a cut that is not perpendicular to the cylinder axis. The samples are  $(\cos \theta_j, \sin \theta_j, z_{ij})$ , where  $\theta_j = 2\pi j/64$  for  $0 \leq j < 64$



and  $z_{ij} = -b + \cos \theta_j + 2bi/64$  for  $b = 1/4$  and  $0 \leq i \leq 64$ . The fitted center, radius, and axis direction are the actual ones (modulo numerical round-off errors). Figure 2 shows a rendering of the points (in black) and a wire frame of the fitted cylinder (in blue).

---

**Figure 2.** Fitting of samples on a cylinder ring that is cut skewed relative to the cylinder axis.

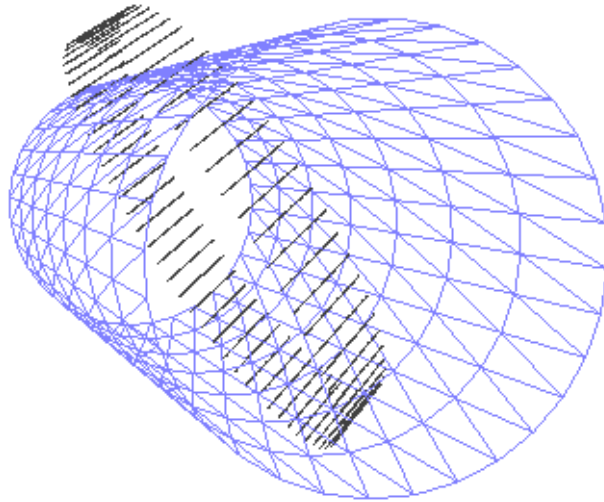


---

In this example, if you were to compute the covariance matrix of the samples and choose the cylinder axis direction to be the eigenvector in the direction of maximum variance, that direction is skewed relative to the original cylinder axis. The fitted parameters are (approximately)  $\mathbf{W} = (0.699, 0, 0.715)$ ,  $PC = (0, 0, 0)$ , and  $r^2 = 0.511$ . Figure 3 shows a rendering of the points (in black) and a wire frame of the fitted cylinder (in blue).

---

**Figure 3.** Fitting of samples on a cylinder ring that is cut skewed relative to the cylinder axis. The direction  $\mathbf{W}$  was chosen to be an eigenvector corresponding to the maximum eigenvalue of the covariance matrix of the samples.



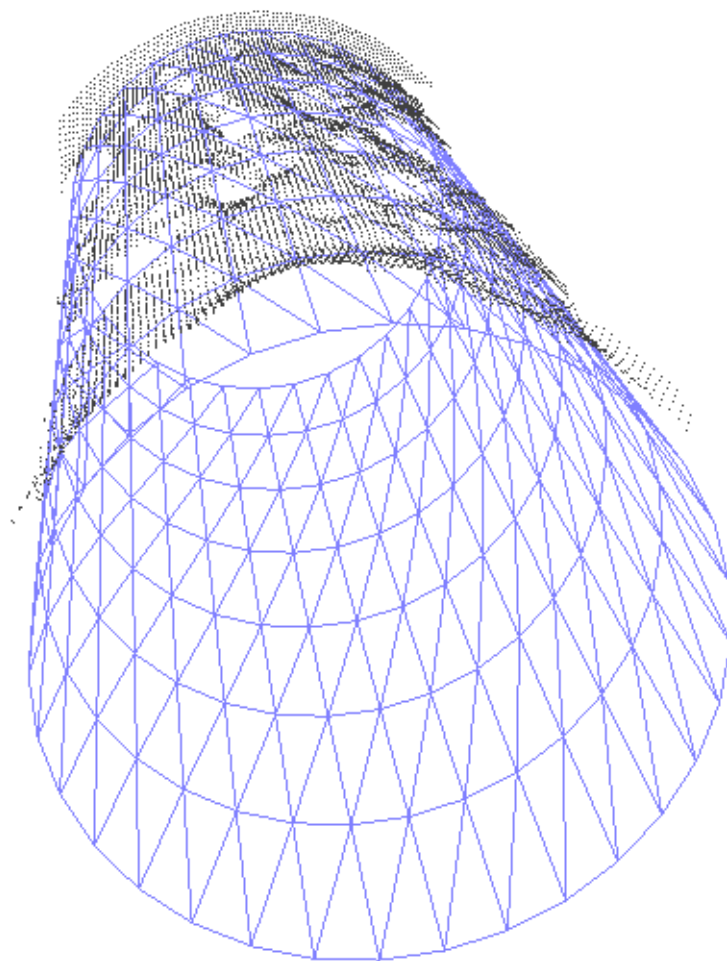
---

You can see that the fitted cylinder is not a good approximation to the points.

Figure 4 shows a point cloud generated from a DIC file (data set courtesy of Carolina Lavecchia) and the fitted cylinder. The data set has nearly 11000 points.

---

**Figure 4.** A view of a fitted point cloud.

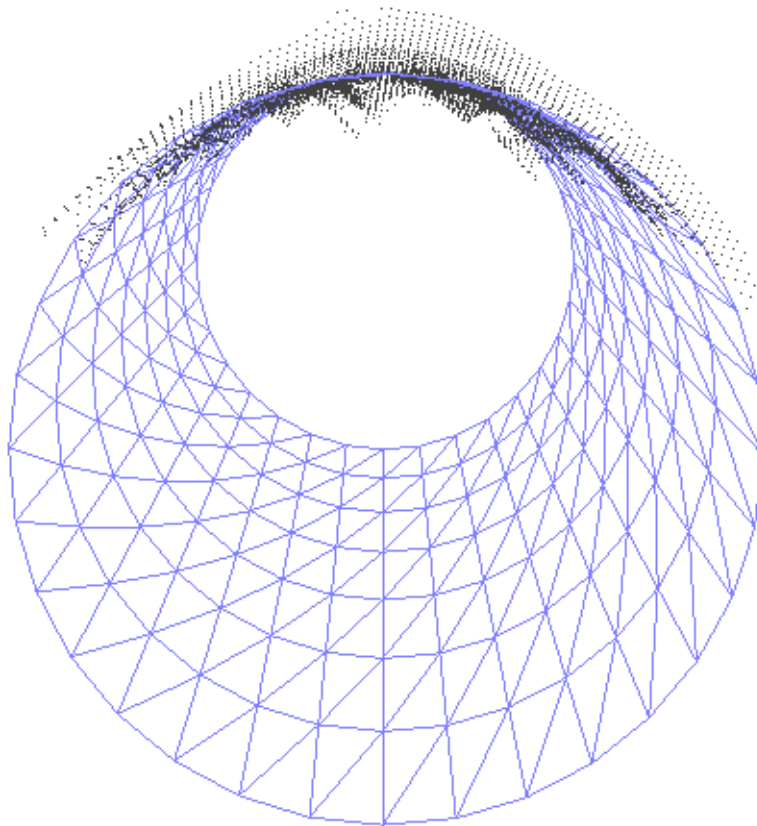


---

Figure 5 shows a different view of the point cloud and cylinder.

---

**Figure 5.** Another view of a fitted point cloud, looking along the top of the cylinder to get an idea of how well the cylinder fits the data.




---

## 8 Fitting a Cone to 3D Points

The cone vertex is  $\mathbf{V}$ , the unit-length axis direction is  $\mathbf{U}$  and the cone angle is  $\theta \in (0, \pi/2)$ . The cone is defined algebraically by those points  $\mathbf{X}$  for which

$$\mathbf{U} \cdot \frac{\mathbf{X} - \mathbf{V}}{|\mathbf{X} - \mathbf{V}|} = \cos(\theta) \quad (138)$$

This can be written as a quadratic equation

$$(\mathbf{X} - \mathbf{V})^T (\cos(\theta)^2 I - \mathbf{U}\mathbf{U}^T) (\mathbf{X} - \mathbf{V}) = 0 \quad (139)$$

with the implicit constraint that  $\mathbf{U} \cdot (\mathbf{X} - \mathbf{V}) > 0$ ; that is,  $\mathbf{X}$  is on the “positive” cone. Define  $\mathbf{W} = \mathbf{U} / \cos(\theta)$ , so  $|\mathbf{W}| > 1$  and

$$F(\mathbf{X}; \mathbf{V}, \mathbf{W}) = (\mathbf{X} - \mathbf{V})^T (I - \mathbf{W}\mathbf{W}^T) (\mathbf{X} - \mathbf{V}) = 0 \quad (140)$$

The nonlinear least-squares fitting of points  $\{\mathbf{X}_i\}_{i=0}^{n-1}$  computes  $\mathbf{V}$  and  $\mathbf{W}$  to minimize the error function

$$E(\mathbf{V}, \mathbf{W}) = \sum_{i=0}^{n-1} F(\mathbf{X}_i; \mathbf{V}, \mathbf{W})^2 \quad (141)$$

This can be solved using standard iterative minimizers such as the Gauss–Newton method and Levenberg–Marquardt method. The partial derivatives of the  $F$  terms are needed for  $\partial E/\partial \mathbf{V}$  and  $\partial E/\partial \mathbf{W}$  in order to compute the Jacobian matrix,

$$\frac{\partial F}{\partial \mathbf{V}} = 2(\Delta - (\mathbf{W} \cdot \Delta) \mathbf{W}), \quad \frac{\partial F}{\partial \mathbf{W}} = -2(\mathbf{W} \cdot \Delta) \Delta \quad (142)$$

where  $\Delta = \mathbf{X} - \mathbf{V}$ . Implementations are found in [GteApprCone3.h](#).

## 8.1 Initial Choice for the Parameters of the Error Function

Without any application-specific knowledge about the sample points, the simplest initial choice for the center, axis directions and axis extents is based on assuming the sample points are dense on a frustum of a cone and then using integrals of various quantities over that frustum. The cone frustum surface is parameterized by

$$\mathbf{P}(h, \phi) = \mathbf{V} + h\mathbf{U} + h \tan \theta (\cos \phi \mathbf{W}_0 + \sin \phi \mathbf{W}_1), \quad h \in [h_0, h_1], \quad \phi \in [0, 2\pi) \quad (143)$$

where  $0 \leq h_0 < h_1$ . The set  $\{\mathbf{U}, \mathbf{W}_0, \mathbf{W}_1\}$  is a right-handed orthonormal basis for  $\mathbb{R}^3$ ; that is, the vectors are unit length, mutually perpendicular, and  $\mathbf{U} = \mathbf{W}_0 \times \mathbf{W}_1$ . We want to use the sample points  $\mathbf{X}_i$  to determine an initial choice for the cone axis direction  $\mathbf{U}$ , the cone angle  $\theta$  and the height bounds  $h_0$  and  $h_1$  of the cone frustum. The integrals involve ratio expressions that are defined by

$$\rho_n = \frac{\int_{h_0}^{h_1} h^{n-1} dh}{\int_{h_0}^{h_1} h dh} = \frac{\frac{1}{n}(h_1^n - h_0^n)}{\frac{1}{2}(h_1^2 - h_0^2)} \quad (144)$$

Define the rotation matrix  $R = [\mathbf{U} \ \mathbf{W}_0 \ \mathbf{W}_1]$  whose columns are the specified basis vectors and define  $\mathbf{Y}(h, \phi)$  to be the  $3 \times 1$  vector which in 3-tuple form is  $(h, h \tan \theta \cos \phi, h \tan \theta \sin \phi)$ . The parameterization of the surface is concisely  $\mathbf{P}(h, \phi) = \mathbf{V} + R\mathbf{Y}(h, \phi)$ . Several integrals are formulated next and involve the following integrals. For notational convenience, I will use  $\mathbf{Y}$  with the understanding that it depends on  $h$  and  $\phi$ .

The element of surface area for the cone frustum surface is

$$\begin{aligned} dA &= \left| \frac{\partial \mathbf{P}}{\partial h} \times \frac{\partial \mathbf{P}}{\partial \phi} \right| dh d\phi \\ &= |(\mathbf{U} + \tan \theta (\cos \phi \mathbf{W}_0 + \sin \phi \mathbf{W}_1)) \times (h \tan \theta (-\sin \phi \mathbf{W}_0 + \cos \phi \mathbf{W}_1))| dh d\phi \\ &= (h \tan \theta) |-\sin \phi \mathbf{U} \times \mathbf{W}_0 + \cos \phi \mathbf{U} \times \mathbf{W}_1 + \tan \theta \mathbf{W}_0 \times \mathbf{W}_1| dh d\phi \\ &= h \tan \theta \sqrt{\sin^2 \phi + \cos^2 \phi + \tan^2 \theta} dh d\phi \\ &= h \tan \theta \sqrt{1 + \tan^2 \theta} dh d\phi \\ &= h \tan \theta \sqrt{1/\cos^2 \theta} dh d\phi \\ &= (h \tan \theta / \cos \theta) dh d\phi \end{aligned} \quad (145)$$

and the surface area of the cone frustum is

$$A = \int_{h_0}^{h_1} \int_0^{2\pi} dA = \int_{h_0}^{h_1} \int_0^{2\pi} (h \tan \theta / \cos \theta) dh d\phi = \pi \frac{\tan \theta}{\cos \theta} (h_1^2 - h_0^2) \quad (146)$$

The average of  $\mathbf{Y}$  over the cone frustum surface is

$$\begin{aligned} \bar{\mathbf{Y}} &= \frac{1}{A} \int_{h_0}^{h_1} \int_0^{2\pi} \mathbf{Y} dA \\ &= \frac{1}{\pi(h_1^2 - h_0^2)} \int_{h_0}^{h_1} \int_0^{2\pi} \mathbf{Y} h dh d\phi \\ &= \frac{1}{\pi(h_1^2 - h_0^2)} \int_{h_0}^{h_1} \int_0^{2\pi} \begin{bmatrix} 1 \\ \tan \theta \cos \phi \\ \tan \theta \sin \phi \end{bmatrix} h^2 dh d\phi \\ &= \rho_3 \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \end{aligned} \quad (147)$$

The average of  $\mathbf{Y}^\top \mathbf{Y}$  over the cone frustum surface is

$$\begin{aligned} \overline{\mathbf{Y}^\top \mathbf{Y}} &= \frac{1}{A} \int_{h_0}^{h_1} \int_0^{2\pi} \mathbf{Y}^\top \mathbf{Y} dA \\ &= \frac{1}{\pi(h_1^2 - h_0^2)} \int_{h_0}^{h_1} \int_0^{2\pi} \mathbf{Y}^\top \mathbf{Y} h dh d\phi \\ &= \frac{1}{\pi(h_1^2 - h_0^2)} \int_{h_0}^{h_1} \int_0^{2\pi} (1 + \tan^2 \theta) h^3 dh d\phi \\ &= \rho_4 \sec^2 \theta \end{aligned} \quad (148)$$

The average of  $\mathbf{Y} \mathbf{Y}^\top$  over the cone frustum surface is

$$\begin{aligned} \overline{\mathbf{Y} \mathbf{Y}^\top} &= \frac{1}{A} \int_{h_0}^{h_1} \int_0^{2\pi} \mathbf{Y} \mathbf{Y}^\top dA \\ &= \frac{1}{\pi(h_1^2 - h_0^2)} \int_{h_0}^{h_1} \int_0^{2\pi} \mathbf{Y} \mathbf{Y}^\top h dh d\phi \\ &= \frac{1}{\pi(h_1^2 - h_0^2)} \int_{h_0}^{h_1} \int_0^{2\pi} \begin{bmatrix} 1 & \tan \theta \cos \phi & \tan \theta \sin \phi \\ \tan \theta \cos \phi & \tan^2 \theta \cos^2 \phi & \tan^2 \theta \sin \phi \cos \phi \\ \tan \theta \sin \phi & \tan^2 \theta \sin \phi \cos \phi & \tan^2 \theta \sin^2 \phi \end{bmatrix} h^3 dh d\phi \\ &= \rho_4 \text{Diag} \left( 1, \frac{1}{2} \tan^2 \theta, \frac{1}{2} \tan^2 \theta \right) \end{aligned} \quad (149)$$

The average of  $\mathbf{Y}\mathbf{Y}^\top\mathbf{Y}$  over the cone frustum surface is

$$\begin{aligned}
\overline{\mathbf{Y}\mathbf{Y}^\top\mathbf{Y}} &= \frac{1}{A} \int_{h_0}^{h_1} \int_0^{2\pi} \mathbf{Y}\mathbf{Y}^\top\mathbf{Y} dA \\
&= \frac{1}{\pi(h_1^2-h_0^2)} \int_{h_0}^{h_1} \int_0^{2\pi} \mathbf{Y}\mathbf{Y}^\top\mathbf{Y} h dh d\phi \\
&= \frac{1}{\pi(h_1^2-h_0^2)} \int_{h_0}^{h_1} \int_0^{2\pi} (1 + \tan^2 \theta) \begin{bmatrix} 1 \\ \tan \theta \cos \phi \\ \tan \theta \sin \phi \end{bmatrix} h^4 dh d\phi \\
&= \rho_5 \sec^2 \theta \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}
\end{aligned} \tag{150}$$

In the following,  $\mathbf{P}$  is the surface parameterization with the understanding that it depends on  $h$  and  $\phi$ . The average of the points over the cone frustum surface is

$$\begin{aligned}
\mathbf{C} &= \frac{1}{A} \int_{h_0}^{h_1} \int_0^{2\pi} \mathbf{P} dA \\
&= \frac{1}{\pi(h_1^2-h_0^2)} \int_{h_0}^{h_1} \int_0^{2\pi} (h\mathbf{V} + h^2\mathbf{U} + h^2 \tan \theta (\cos \phi \mathbf{W}_0 + \sin \phi \mathbf{W}_1)) dh d\phi \\
&= \mathbf{V} + \rho_3 \mathbf{U}
\end{aligned} \tag{151}$$

Define the difference  $\mathbf{\Delta} = \mathbf{V} - \mathbf{C} = -\rho_3 \mathbf{U}$ , in which case  $\mathbf{P} - \mathbf{C} = \mathbf{\Delta} + R\mathbf{Y}$ . Define  $\mathbf{Z} = \mathbf{P} - \mathbf{C}$  and define  $\mathbf{i}$  be the  $3 \times 1$  vector which as a 3-tuple is  $(1, 0, 0)$ . Observe that  $R\mathbf{i} = \mathbf{U}$ . The average of  $\mathbf{Z}\mathbf{Z}^\top$  over the cone frustum surface is

$$\begin{aligned}
\overline{\mathbf{Z}\mathbf{Z}^\top} &= \frac{1}{A} \int_{h_0}^{h_1} \int_0^{2\pi} \mathbf{Z}\mathbf{Z}^\top dA \\
&= \frac{1}{A} \int_{h_0}^{h_1} \int_0^{2\pi} \left( \mathbf{\Delta}\mathbf{\Delta}^\top + \mathbf{\Delta}\mathbf{Y}^\top R^\top + R\mathbf{Y}\mathbf{\Delta}^\top + R\mathbf{Y}\mathbf{Y}^\top R^\top \right) dA \\
&= \mathbf{\Delta}\mathbf{\Delta}^\top + \mathbf{\Delta}\overline{\mathbf{Y}^\top} R^\top + R\overline{\mathbf{Y}}\mathbf{\Delta}^\top + R\overline{\mathbf{Y}\mathbf{Y}^\top} R^\top \\
&= (-\rho_3 \mathbf{U})(-\rho_3 \mathbf{U})^\top + (-\rho_3 \mathbf{U})(\rho_3 \mathbf{i})^\top R^\top + R(\rho_3 \mathbf{i})(-\rho_3 \mathbf{U})^\top \\
&\quad + R(\rho_4 \text{Diag}(1, \frac{1}{2} \tan^2 \theta, \frac{1}{2} \tan^2 \theta)) R^\top \\
&= (\rho_4 - \rho_3^2) \mathbf{U}\mathbf{U}^\top + \frac{1}{2} \rho_4 \tan^2 \theta (\mathbf{W}_0 \mathbf{W}_0^\top + \mathbf{W}_1 \mathbf{W}_1^\top) \\
&= R \text{Diag}(\rho_4 - \rho_3^2, \frac{1}{2} \rho_4 \tan^2 \theta, \frac{1}{2} \rho_4 \tan^2 \theta) R^\top
\end{aligned} \tag{152}$$

The matrix  $\overline{\mathbf{Z}\mathbf{Z}^\top}$  is the covariance matrix of the cone frustum surface points and its eigendecomposition is

given by the right-hand side of equation (152). The average of  $\overline{\mathbf{Z}\mathbf{Z}^\top\mathbf{Z}}$  over the cone frustum surface is

$$\begin{aligned}
\overline{\mathbf{Z}\mathbf{Z}^\top\mathbf{Z}} &= \frac{1}{A} \int_{h_0}^{h_1} \int_0^{2\pi} \mathbf{Z}\mathbf{Z}^\top\mathbf{Z} dA \\
&= \frac{1}{A} \int_{h_0}^{h_1} \int_0^{2\pi} \left( (\Delta^\top\Delta)\Delta + (2\Delta\Delta^\top R)\mathbf{Y} + (\mathbf{Y}^\top\mathbf{Y})\Delta(\Delta^\top\Delta)R\mathbf{Y} + \right. \\
&\quad \left. + (2R\mathbf{Y}\mathbf{Y}^\top R^\top)\Delta + R\mathbf{Y}\mathbf{Y}^\top\mathbf{Y} \right) \\
&= (\Delta^\top\Delta)\Delta + (2\Delta\Delta^\top R)\overline{\mathbf{Y}} + \overline{\mathbf{Y}^\top\mathbf{Y}}\Delta + (\Delta^\top\Delta)R\overline{\mathbf{Y}} + 2\overline{R\mathbf{Y}\mathbf{Y}^\top R^\top}\Delta + \overline{R\mathbf{Y}\mathbf{Y}^\top\mathbf{Y}} \\
&= (-\rho_3^3 + 2\rho_3^3 - \rho_3\rho_4 \sec^2 \theta + \rho_3^3 - 2\rho_3\rho_4 + \rho_5 \sec^2 \theta)\mathbf{U} \\
&= (2\rho_3(\rho_3^2 - \rho_4) + (\rho_5 - \rho_3\rho_4) \sec^2 \theta)\mathbf{U}
\end{aligned} \tag{153}$$

Equations (152) and (153) can be manipulated to obtain 3 equations in the 3 unknowns  $h_0$ ,  $h_1$  and  $\tan \theta$ . In this sense, if we have good estimates for  $\mathbf{C}$ ,  $\overline{\mathbf{Z}\mathbf{Z}^\top}$  and  $\overline{\mathbf{Z}\mathbf{Z}^\top\mathbf{Z}}$ , we can reconstruct the cone frustum from a dense set of samples on or near the frustum.

### 8.1.1 Simple Attempt to Reconstruct Height Extremes

The cone parameterization in terms of the centroid  $\mathbf{C}$  is

$$\mathbf{P}(h, \phi) = \mathbf{C} + (h - p_3)\mathbf{U} + h \tan \theta (\cos \phi \mathbf{W}_0 + \sin \phi \mathbf{W}_1) \tag{154}$$

The minimum and maximum of the projections onto  $\mathbf{U}$  relative to  $\mathbf{C}$  are  $\hat{h}_0 = h_0 - p_3$  and  $\hat{h}_1 = h_1 - p_3$ , respectively. Recall that  $p_3 = ((h_1^3 - h_0^3)/3)/((h_1^2 - h_0^2)/2)$ , so we have two equations in  $h_0$  and  $h_1$  that we can solve for in terms of  $\hat{h}_0$  and  $\hat{h}_1$ ,

$$h_0 = \frac{\hat{h}_0^2 + \hat{h}_0\hat{h}_1 - 2\hat{h}_1^2}{3(\hat{h}_0 + \hat{h}_1)}, \quad h_1 = \frac{\hat{h}_0^1 + \hat{h}_1\hat{h}_0 - 2\hat{h}_0^2}{3(\hat{h}_1 + \hat{h}_0)} \tag{155}$$

Because  $\mathbf{C}$  is the centroid, it is necessary that  $\hat{h}_0 < 0 < \hat{h}_1$ . Also notice that  $h_1 - h_0 = \hat{h}_1 - \hat{h}_0$ , and the right-hand side is a robust estimate of  $h_1 - h_0$ .

Although this is a simple mathematical problem given infinitely many points on a cone frustum surface, in practice it has problems. Listing 16 contains code to generate a rectangular mesh of points on a cone frustum that leads to  $\hat{h}_1$  being nearly equal to  $-\hat{h}_0$ , which causes the denominator in the  $h_0$  and  $h_1$  reconstruction equations to be nearly zero. The reconstructed  $h_0$  and  $h_1$  are nowhere near what they should be.

---

**Listing 16.** The reconstructed  $h_0$  and  $h_1$  for a dense sample of points on the cone frustum are grossly incorrect. The cone axis direction is estimated using equation (153) rather than as an eigenvector from equation (152); however, both estimates are nearly identical and nearly equal to the  $(3, 2, 1)/\sqrt{14}$ .

```

Vector3<double> V = { 3.0, 2.0, 1.0 };
Vector3<double> U = { 1.0, 2.0, 3.0 };
Vector3<double> basis [3];
basis [0] = U;
ComputeOrthogonalComplement (1, basis);
U = basis [0];
Vector3<double> W0 = basis [1];

```



```

Vector3<double> W1 = basis [2];
double h0 = 1.0;
double h1 = 2.0;
double theta = GTE_C_PI / 4.0;
double tantheta = std::tan(theta);

size_t const numH = 512, numR = 512;
std::vector<Vector3<double>> X(numH * numR);
for (size_t ih = 0, i = 0; ih < numH; ++ih)
{
    double h = h0 + (h1 - h0) * (double)ih / (double)(numH - 1);
    double r = h * tantheta;
    for (size_t ir = 0; ir < numR; ++ir, ++i)
    {
        double phi = GTE_C_TWO_PI * (double)ir / (double)numR;
        double csphi = std::cos(phi);
        double snphi = std::sin(phi);
        X[i] = V + h * U + r * (csphi * W0 + snphi * W1);
    }
}

Vector3<double> C{ 0.0, 0.0, 0.0 }; // the centroid
for (size_t i = 0; i < X.size(); ++i)
{
    C += X[i];
}
C /= (double)X.size();

Vector3<double> Ufit{ 0.0, 0.0, 0.0 }; // the third-order term
for (size_t i = 0; i < X.size(); ++i)
{
    Vector3<double> diff = X[i] - C;
    Ufit += Dot(diff, diff) * diff;
}
Ufit /= (double)X.size();
Normalize(Ufit);

double H1cen = std::numeric_limits<double>::max();
double h1hat = -h0hat;
for (size_t i = 0; i < X.size(); ++i)
{
    Vector3<double> diff = X[i] - C;
    double h = Dot(Ufit, diff);
    h0hat = std::min(h, h0hat);
    h1hat = std::max(h, h1hat);
}
double hsum = h1hat + h0hat;
double h0reconstruct = (h0hat * (h0hat + h1hat) - 2.0 * h1hat * h1hat) / (3.0 * hsum);
double h1reconstruct = (h1hat * (h0hat + h1hat) - 2.0 * h0hat * h0hat) / (3.0 * hsum);
// h0hat = -0.50000000000238787
// h1hat = 0.50000000000458455
// hsum = -8.7140517258189930e-14
// h0reconstruct = -75871822289.547775 (original was 1.0)
// h1reconstruct = -75871822288.547775 (original was 2.0)

```

Other equations can be derived to attempt reconstruction of  $h_0$  and  $h_1$ . For example, we can compute averages of powers of the height relative to the centroid,

$$\frac{1}{A} \int_{h_0}^{h_1} \int_0^{2\pi} (\mathbf{U} \cdot (\mathbf{P} - \mathbf{C}))^2 dA = \frac{1}{\pi(h_1^2 - h_0^2)} \int_{h_0}^{h_1} \int_0^{2\pi} (h - p_3)^2 h dh d\phi = p_4 - p_3^2 \quad (156)$$

Given a set of sample points, the equation to reconstruct the height extremes is

$$p_4 - p_3^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{U} \cdot (\mathbf{X}_i - \mathbf{C}))^2 = b \quad (157)$$

where the last equality defines  $b$ . We know the estimate  $h_1 - h_0 = \hat{h}_1 - \hat{h}_0 = a$  where the last equality defines  $a$ . We need to solve  $h_1 - h_0 = a$  and  $p_4 - p_3^2 = b$ . The solution is

$$h_0 = \frac{+(36ab - 3a^3) \pm \sqrt{3} a^2 \sqrt{a^2 - 12b}}{6(a^2 - 12b)}, \quad h_1 = \frac{-(36ab - 3a^3) \pm \sqrt{3} a^2 \sqrt{a^2 - 12b}}{6(a^2 - 12b)} \quad (158)$$

In Listing 16, after reconstruction of  $h_0$  and  $h_1$ , add the code

```
double a = h1hat - h0hat; // 1.000000000069724
double b = 0.0;
for (size_t i = 0; i < X.size(); ++i)
{
    Vector3<double> diff = X[i] - C;
    double h = Dot(Ufit, diff);
    b += h * h;
}
b /= (double)X.size(); // 0.083659491193733351
double discr = a * a - 12.0 * b; // -0.0039138943108554258
```

The discriminant  $a^2 - 12b^2$  is negative, so there are no real-valued solutions for  $h_0$  and  $h_1$ .

Attempts to estimate  $h_0$  and  $h_1$  from the average integral of  $(\mathbf{U} \cdot (\mathbf{P} - \mathbf{C}))^3$  also failed when using the point samples to estimate the integrals and then solve the resulting equations.

### 8.1.2 Attempt to Reconstruct the Cone Axis Direction

The eigendecomposition in equation (152) provides a symbolic expression involving the cone axis direction  $\mathbf{U}$ . In a numerical implementation, we need to identify which eigenvector output by the eigensolver corresponds to  $\mathbf{U}$ . The covariance matrix of the point samples is used to estimate  $\overline{\mathbf{Z}\mathbf{Z}^T}$ . We cannot expect two distinct eigenvalues  $\lambda_1 = p_4 - p_3^2$  and  $\lambda_2 = (p_4 \tan^2 \theta)/4$ , the first with multiplicity 1 and the second with multiplicity 2. For if we obtained such a result numerically, selection of  $\mathbf{U}$  is trivial. Even theoretically we could have a problem when there is a single eigenvalue of multiplicity 3.

If the point samples are nearly on a cone frustum, let the numerically computed eigenvalues be sorted as  $\lambda_1 \leq \lambda_2 \leq \lambda_3$ . We could compute  $\lambda_2 - \lambda_1$  and  $\lambda_3 - \lambda_2$  and claim that the minimum difference is due to numerical rounding errors. The conclusion is that those two eigenvalues are theoretically a single eigenvalue. The other eigenvalue is the one associated with  $\mathbf{U}$ . Some experiments showed that this is not a reliable way to select  $\mathbf{U}$ .

Instead we can use equation (153) to select  $\mathbf{U}$  unambiguously. The integral  $\overline{\mathbf{Z}\mathbf{Z}^T\mathbf{Z}}$  is estimated by a summation involving the point samples as shown in Listing 16. The normalization of the summation is the estimate for  $\mathbf{U}$ , named  $\mathbf{Ufit}$  in the listing.

### 8.1.3 Attempt to Reconstruct the Cone Vertex

As mentioned previously, reconstructing  $h_0$  and  $h_1$  would allow us to compute  $p_3$  and then  $\mathbf{V} = \mathbf{C} - p_3\mathbf{U}$ , but the reconstruction attempts can fail numerically. A more reliable algorithm for estimating  $\mathbf{V}$  is shown here. The goal is to estimate  $p_3$  directly rather than estimating  $h_0$  and  $h_1$  and then computing  $p_3$  from them.

In the following,  $\mathbf{U}$  is the estimate of the cone axis direction obtained from equation (153). We want to estimate a positive  $t$  for which  $\mathbf{V} = \mathbf{C} - t\mathbf{U}$ . At the same time we need to estimate  $s = \cos^2 \theta$  for the cone

angle. The algorithm will use the discrete least-squares error function for a specified  $\mathbf{U}$ . Define  $\mathbf{\Delta}_i = \mathbf{X}_i - \mathbf{C}$ ,  $a_i = \mathbf{U} \cdot \mathbf{\Delta}_i$  and  $b_i = \mathbf{\Delta}_i \cdot \mathbf{\Delta}_i$ . Define

$$F_i = (\mathbf{X}_i - \mathbf{V})^\top \left( sI - \mathbf{U}\mathbf{U}^\top \right) (\mathbf{X}_i - \mathbf{V}) = s(t^2 + 2a_it + b_i) - (t + a_i)^2 \quad (159)$$

The least-squares error function is  $E = (\sum_{i=1}^n F_i^2)/n$ . The derivatives of the  $F_i$  terms are

$$\frac{\partial F_i}{\partial s} = t^2 + 2a_it + b_i, \quad \frac{\partial F_i}{\partial t} = 2(s-1)(t + a_i) \quad (160)$$

Some algebra will show that  $\partial E/\partial t = 0$  and  $\partial E/\partial s = 0$  lead to

$$sp_3(t) - q_3(t) = 0, \quad sp_4(t) - q_4(t) = 0 \quad (161)$$

where  $p_d(t)$  and  $q_d(t)$  are polynomials of degree  $d$ . Specifically,

$$\begin{aligned} p_3(t) &= t^3 + e_0t^2 + e_1t + e_2 \\ q_3(t) &= t^3 + e_0t^2 + e_3t + e_4 \\ p_4(t) &= t^4 + f_0t^3 + f_1t^2 + f_2t + f_3 \\ q_4(t) &= t^4 + f_0t^3 + f_4t^2 + f_5t + f_6 \end{aligned} \quad (162)$$

where

$$\begin{aligned} e_0 &= \frac{1}{n} \sum_{i=1}^n 3a_i & f_0 &= \frac{1}{n} \sum_{i=1}^n 4a_i \\ e_1 &= \frac{1}{n} \sum_{i=1}^n (2a_i^2 + b_i) & f_1 &= \frac{1}{n} \sum_{i=1}^n (4a_i^2 + 2b_i) \\ e_2 &= \frac{1}{n} \sum_{i=1}^n a_i b_i & f_2 &= \frac{1}{n} \sum_{i=1}^n 4a_i b_i \\ e_3 &= \frac{1}{n} \sum_{i=1}^n 3a_i^2 & f_3 &= \frac{1}{n} \sum_{i=1}^n b_i^2 \\ e_4 &= \frac{1}{n} \sum_{i=1}^n a_i^3 & f_4 &= \frac{1}{n} \sum_{i=1}^n (5a_i^2 + b_i) \\ & & f_5 &= \frac{1}{n} \sum_{i=1}^n (2a_i^3 + 2a_i b_i) \\ & & f_6 &= \frac{1}{n} \sum_{i=1}^n a_i^2 b_i \end{aligned} \quad (163)$$

Define  $c_{rs} = \frac{1}{n} \sum_{i=1}^n a_i^r b_i^s$ . We can solve one  $st$ -equation for  $s = q_4(t)/p_4(t)$ , substitute into the other  $st$ -equation and then multiply by  $p_4(t)$  to obtain

$$0 = p_3(t)q_4(t) - p_4(t)q_3(t) = g_4t^4 + g_3t^3 + g_2t^2 + g_1t + g_0 \quad (164)$$

where

$$\begin{aligned} g_4 &= c_{30} - c_{11} + c_{10}(c_{01} - c_{20}) \\ g_3 &= c_{21} - c_{02} + c_{01}(c_{01} + c_{20}) + 2(c_{10}(c_{30} - c_{11}) - c_{20}^2) \\ g_2 &= 3(c_{11}(c_{01} - c_{20}) + c_{10}(c_{21} - c_{02})) \\ g_1 &= c_{01}c_{21} - 3c_{02}c_{20} + 2(c_{20}c_{21} - c_{11}(c_{30} - c_{11})) \\ g_0 &= c_{11}c_{21} - c_{02}c_{30} \end{aligned} \quad (165)$$

The least-squares error function is smooth and nonnegative, so there must be at least one  $(s, t)$  for which  $\nabla E(s, t) = (0, 0)$ . In theory, this means  $g(t) = 0$  must have at least one real-valued root. Naturally, the

quartic root solver must guard against floating-point round-off errors to guarantee this. The GTEngine quartic root solver can be executed with rational arithmetic to correctly classify the roots, after which rounding errors will not change the classification. For each root  $t \geq 0$ , compute  $s = q_3(t)/p_3(t)$ . If  $s \in (0, 1)$ , we have a valid angle because  $s = \cos^2 \theta$ . Evaluate  $E(s, t)$ . Of all such valid pairs  $(s, t)$ , choose the one whose  $E$ -value is minimum. This pair is used to initialize  $\mathbf{V}$  and  $\cos \theta$ .

Using floating-point arithmetic, it is unlikely but possible that there is no pair  $(s, t)$  generated by the algorithm of the previous paragraph. An implementation must guard against this. Although most likely bad estimates, choose the initial  $\mathbf{V}$  to be the average  $\mathbf{C}$  and choose the angle arbitrarily to be  $\pi/4$  which is the midpoint of the domain of the angles  $(0, \pi/2)$ .

Listing 16 can be modified by replacing all the code after the line `Normalize(Ufit)` with that shown in Listing 17. The new listing shows an implementation for the initial guesses for  $\mathbf{V}$  and  $\cos^2 \theta$ .

---

**Listing 17.** Code for the initial guesses for the cone vertex and cone angle. The numbers in the comments are those produced by the example in Listing 16.

```

double c10 = 0.0, c20 = 0.0, c30 = 0.0, c01 = 0.0, c02 = 0.0, c11 = 0.0, c21 = 0.0;
for (size_t i = 0; i < X.size(); ++i)
{
    Vector3<double> diff = X[i] - C;
    double ai = Dot(Ufit, diff);
    double bi = Dot(diff, diff);
    c10 += ai;
    c20 += ai * ai;
    c30 += ai * ai * ai;
    c01 += bi;
    c02 += bi * bi;
    c11 += ai * bi;
    c21 += ai * ai * bi;
}
c10 /= (double)X.size();
c20 /= (double)X.size();
c30 /= (double)X.size();
c01 /= (double)X.size();
c02 /= (double)X.size();
c11 /= (double)X.size();
c21 /= (double)X.size();

// The coefficients for polynomials p3(t) and q3(t).
double e0 = 3.0 * c10;
double e1 = 2.0 * c20 + c01;
double e2 = c11;
double e3 = 3.0 * c20;
double e4 = c30;

// The coefficients for polynomials p4(t) and q4(t).
double f0 = 4.0 * c10;
double f1 = 4.0 * c20 + 2.0 * c01;
double f2 = 4.0 * c11;
double f3 = c02;
double f4 = 5.0 * c20 + c01;
double f5 = 2.0 * c30 + 2.0 * c11;
double f6 = c21;

// The coefficients for the quartic polynomial g(t).
double g0 = c11 * c21 - c02 * c30; // 0.053566286603418618
double g1 = c01 * c21 - 3.0 * c02 * c20 + 2.0 * (c20 * c21 - c11 * (c30 - c11)); // -0.98354780514088114
double g2 = 3.0 * (c11 * (c01 - c20) + c10 * (c21 - c02)); // 1.7570948908739785
double g3 = c21 - c02 + c01 * (c01 + c20) + 2.0 * (c10 * (c30 - c11) - c20 * c20); // -0.37366801803251715
double g4 = c30 - c11 + c10 * (c01 - c20); // -0.25097847358125042

// Compute the roots of g(t) = 0.

```

```

std::map<double, int> rmMap;
RootsPolynomial<double>::SolveQuartic(g0, g1, g2, g3, g4, rmMap);
// root[0] = -3.6829478517723415
// root[1] = 0.061025506205900915
// root[2] = 0.63307745500821921
// root[3] = 1.499999999961315

std::vector<std::array<double, 3>> info;
double s, t;
for (auto const& element : rmMap)
{
    t = element.first;
    if (t > 0.0)
    {
        s = (e4 + t * (e3 + t * (e0 + t))) / (e2 + t * (e1 + t * (e0 + t)));
        if (0.0 < s && s < 1.0)
        {
            double error = 0.0;
            for (size_t i = 0; i < X.size(); ++i)
            {
                Vector3<double> diff = X[i] - C;
                double ai = Dot(Ufit, diff);
                double bi = Dot(diff, diff);
                double tpai = t + ai;
                double Fi = s * (bi + t * (2.0 * ai + t)) - tpai * tpai;
                error += Fi * Fi;
            }
            error /= (double)X.size();
            std::array<double, 3> item = { s, t, error };
            info.push_back(item);
        }
    }
}

double minError = std::numeric_limits<double>::max();
std::array<double, 3> minItem = { 0.0, 0.0, minError };
for (auto const& item : info)
{
    if (item[2] < minError)
    {
        minItem = item;
    }
}

Vector3<double> Vfit;
double cosAngleSqrFit;
if (minItem[2] < std::numeric_limits<double>::max())
{
    // minItem = { minS, minT, minError }
    Vfit = C - minItem[1] * Ufit;
    cosAngleSqrFit = minItem[0];

    // We do not need this for the minimizer, but let's reconstruct
    // the heights to see how close we get to the original ones
    // (h0 = 1, h1 = 2). Knowing the estimate h1 - h0 = 1, solve
    // minT = p3 = ((h1^3 - h0^3)/3)/((h1^2 - h0^2)/2)
    // for h0 and h1.
    double h0 = (6.0 + std::sqrt(276)) / 24.0; // 0.94221865524317294
    double h1 = 1.0 + h0; // 1.9422186552431731

    // The original cone has p3 = 14/9 = 1.55555 and the minimizing
    // t-value is 1.499999999961315.
}
else
{
    Vfit = C;
    cosAngleSqrFit = sqrt(0.5); // angleFit is pi/4
}

```

## 9 Fitting a Paraboloid to 3D Points of the Form $(x, y, f(x, y))$

Given a set of samples  $\{(x_i, y_i, z_i)\}_{i=1}^m$  and assuming that the true values lie on a paraboloid

$$z = f(x, y) = p_1x^2 + p_2xy + p_3y^2 + p_4x + p_5y + p_6 = \mathbf{P} \cdot \mathbf{Q}(x, y) \quad (166)$$

where  $\mathbf{P} = (p_1, p_2, p_3, p_4, p_5, p_6)$  and  $\mathbf{Q}(x, y) = (x^2, xy, y^2, x, y, 1)$ , select  $\mathbf{P}$  to minimize the sum of squared errors

$$E(\mathbf{P}) = \sum_{i=1}^m (\mathbf{P} \cdot \mathbf{Q}_i - z_i)^2 \quad (167)$$

where  $\mathbf{Q}_i = \mathbf{Q}(x_i, y_i)$ . The minimum occurs when the gradient of  $E$  is the zero vector,

$$\nabla E = 2 \sum_{i=1}^m (\mathbf{P} \cdot \mathbf{Q}_i - z_i) \mathbf{Q}_i = \mathbf{0} \quad (168)$$

Some algebra converts this to a system of 6 equations in 6 unknowns:

$$\left( \sum_{i=1}^m \mathbf{Q}_i \mathbf{Q}_i^T \right) \mathbf{P} = \sum_{i=1}^m z_i \mathbf{Q}_i \quad (169)$$

The product  $\mathbf{Q}_i \mathbf{Q}_i^T$  is a product of the  $6 \times 1$  matrix  $\mathbf{Q}_i$  with the  $1 \times 6$  matrix  $\mathbf{Q}_i^T$ , the result being a  $6 \times 6$  matrix.

Define the  $6 \times 6$  symmetric matrix  $A = \sum_{i=1}^m \mathbf{Q}_i \mathbf{Q}_i^T$  and the  $6 \times 1$  vector  $\mathbf{B} = \sum_{i=1}^m z_i \mathbf{Q}_i$ . The choice for  $\mathbf{P}$  is the solution to the linear system of equations  $A\mathbf{P} = \mathbf{B}$ . The entries of  $A$  and  $\mathbf{B}$  indicate summations over the appropriate product of variables. For example,  $s(x^3y) = \sum_{i=1}^m x_i^3 y_i$ :

$$\begin{bmatrix} s(x^4) & s(x^3y) & s(x^2y^2) & s(x^3) & s(x^2y) & s(x^2) \\ s(x^3y) & s(x^2y^2) & s(xy^3) & s(x^2y) & s(xy^2) & s(xy) \\ s(x^2y^2) & s(xy^3) & s(y^4) & s(xy^2) & s(y^3) & s(y^2) \\ s(x^3) & s(x^2y) & s(xy^2) & s(x^2) & s(xy) & s(x) \\ s(x^2y) & s(xy^2) & s(y^3) & s(xy) & s(y^2) & s(y) \\ s(x^2) & s(xy) & s(y^2) & s(x) & s(y) & s(1) \end{bmatrix} \begin{bmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \\ p_5 \\ p_6 \end{bmatrix} = \begin{bmatrix} s(zx^2) \\ s(zxy) \\ s(zy^2) \\ s(zx) \\ s(zy) \\ s(z) \end{bmatrix} \quad (170)$$

An implementation is [GteApprParaboloid3.h](#).