

Math Camp

Justin Grimmer

Associate Professor
Department of Political Science
Stanford University

September 13th, 2016

Multivariate Optimization

Optimizing multivariate functions

- Parameters $\beta = (\beta_1, \beta_2, \dots, \beta_n)$ such that $f(\beta|\mathbf{X}, \mathbf{Y})$ is maximized
- Policy $\mathbf{x} \in \mathfrak{R}^n$ that maximizes $U(\mathbf{x})$
- Weights $\pi = (\pi_1, \pi_2, \dots, \pi_K)$ such that a weighted average of forecasts $\mathbf{f} = (f_1, f_2, \dots, f_k)$ have minimum loss

$$\min_{\pi} = -\left(\sum_{j=1}^K \pi_j f_j - y\right)^2$$

Today we'll describe analytic and computational approaches to optimization

- Analytic recipe for optimization
- Computational optimization
 - Multivariate Newton-Raphson
 - BFGS
 - Approximate Optimization: k-means

Multivariate Optimization

Definition

Let $\mathbf{x} \in \mathbb{R}^n$ and let $\delta > 0$. Define a *neighborhood* of \mathbf{x} , $B(\mathbf{x}, \delta)$, as the set of points such that,

$$B(\mathbf{x}, \delta) = \{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{x} - \mathbf{y}\| < \delta\}$$

Definition

Suppose $f : X \rightarrow \mathbb{R}$ with $X \subset \mathbb{R}^n$. A vector $\mathbf{x}^* \in X$ is a *global maximum* if, for all other $\mathbf{x} \in X$

$$f(\mathbf{x}^*) > f(\mathbf{x})$$

A vector \mathbf{x}^{local} is a *local maximum* if there is a neighborhood around \mathbf{x}^{local} , $Q \subset X$ such that, for all $\mathbf{x} \in Q$,

$$f(\mathbf{x}^{local}) > f(\mathbf{x})$$

Multivariate Optimization

Definition

A set $X \subset \mathbb{R}^n$ is **compact** if it is closed and bounded

Theorem

Multivariate Extreme Value Theorem Suppose $f : X \rightarrow \mathbb{R}$ be continuous and $X \subset \mathbb{R}^n$ and X compact. Then f takes on its **maximum** and **minimum** values on X .

We're going to come up with the multivariate equivalent of the **first order** and **second order** conditions now

Gradient

Definition

Suppose $f : X \rightarrow \mathbb{R}^n$ with $X \subset \mathbb{R}^1$ is a differentiable function. Define the gradient vector of f at \mathbf{x}_0 , $\nabla f(\mathbf{x}_0)$ as,

$$\nabla f(\mathbf{x}_0) = \left(\frac{\partial f(\mathbf{x}_0)}{\partial x_1}, \frac{\partial f(\mathbf{x}_0)}{\partial x_2}, \frac{\partial f(\mathbf{x}_0)}{\partial x_3}, \dots, \frac{\partial f(\mathbf{x}_0)}{\partial x_n} \right)$$

Gradient First Order Condition

Theorem

Suppose $f : X \rightarrow \mathbb{R}^1$, $X \subset \mathbb{R}^n$. Suppose $\mathbf{a} \in X$ is a *local* extremum. Then,

$$\begin{aligned}\nabla f(\mathbf{a}) &= \mathbf{0} \\ &= (0, 0, \dots, 0)\end{aligned}$$

- Proof (intuition): same as one dimensional case (left-hand, right hand), just do it dimension by dimension
- **Critical Values:**
 - 1) Maximum
 - 2) Minimum
 - 3) **Saddle point**
- **Second Derivative Test!**

Second Order Conditions: Hessian

Definition

Suppose $f : X \rightarrow \mathbb{R}^1$, $X \subset \mathbb{R}^n$, with f a twice differentiable function. We will define the **Hessian** matrix as the matrix of second derivatives at $\mathbf{x}^* \in X$,

$$\mathbf{H}(f)(\mathbf{x}^*) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1}(\mathbf{x}^*) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(\mathbf{x}^*) & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n}(\mathbf{x}^*) \\ \frac{\partial^2 f}{\partial x_2 \partial x_1}(\mathbf{x}^*) & \frac{\partial^2 f}{\partial x_2 \partial x_2}(\mathbf{x}^*) & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n}(\mathbf{x}^*) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1}(\mathbf{x}^*) & \frac{\partial^2 f}{\partial x_n \partial x_2}(\mathbf{x}^*) & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n}(\mathbf{x}^*) \end{pmatrix}$$

General test \rightsquigarrow Two Dimensional Test \rightsquigarrow Example

Hessians

Definition

Consider $n \times n$ matrix \mathbf{A} . If, for all $\mathbf{x} \in \mathbb{R}^n$ where $\mathbf{x} \neq \mathbf{0}$:

$$\mathbf{x}' \mathbf{A} \mathbf{x} > 0 \quad \mathbf{A} \text{ is } \textit{positive definite}$$

$$\mathbf{x}' \mathbf{A} \mathbf{x} < 0 \quad \mathbf{A} \text{ is } \textit{negative definite}$$

If $\mathbf{x}' \mathbf{A} \mathbf{x} > 0$ for some \mathbf{x} and $\mathbf{x}' \mathbf{A} \mathbf{x} < 0$ for other \mathbf{x} , then we say \mathbf{A} is *indefinite*

Approximating functions and second order conditions

Theorem

Taylor's Theorem Suppose $f : \mathbb{R} \rightarrow \mathbb{R}$, $f(x)$ is infinitely differentiable function. Then, the Taylor expansion of $f(x)$ around a is given by

$$f(x) = f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f'''(a)}{3!}(x-a)^3 + \dots$$

$$f(x) = \sum_{n=0}^{\infty} \frac{f^n(a)}{n!}(x-a)^n$$

Example Function

Suppose $a = 0$ and $f(x) = e^x$. Then,

$$f'(x) = e^x$$

$$f''(x) = e^x$$

$$\vdots \quad \vdots \quad \vdots$$

$$f^n(x) = e^x$$

This implies

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} \dots + \frac{x^n}{n!} + \dots$$

Multivariate Taylor's Theorem

Theorem

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a three-times continuously differentiable function, then around $\mathbf{a} \in \mathbb{R}^n$,

$$f(\mathbf{x}) = f(\mathbf{a}) + \nabla f(\mathbf{a})(\mathbf{x} - \mathbf{a}) + \frac{1}{2}(\mathbf{x} - \mathbf{a})' \mathbf{H}(f)(\mathbf{a})(\mathbf{x} - \mathbf{a}) + R(\mathbf{a}, \mathbf{x})$$

where $\frac{R(\mathbf{x}, \mathbf{a})}{\|\mathbf{x} - \mathbf{a}\|^2} \rightarrow 0$ as $\mathbf{x} \rightarrow \mathbf{a}$

Intuition for Quadratic Form

Suppose \mathbf{x}^* is some critical value,

$$f(\mathbf{x}) = f(\mathbf{x}^*) + \nabla f(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) + \left(\mathbf{x} - \frac{1}{2}\mathbf{x}^*\right)\mathbf{H}(f)(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) + R(\mathbf{x}^*, \mathbf{x})$$

$$f(\mathbf{x}) - f(\mathbf{x}^*) = 0(\mathbf{x} - \mathbf{x}^*) + \left(\mathbf{x} - \frac{1}{2}\mathbf{x}^*\right)\mathbf{H}(f)(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) + R(\mathbf{x}^*, \mathbf{x})$$

For \mathbf{x} near \mathbf{x}^* , $R(\mathbf{x}^*, \mathbf{x}) \approx 0$

$\mathbf{H}(f)(\mathbf{x}^*)$ positive definite $\rightarrow f(\mathbf{x}) > f(\mathbf{x}^*) \rightarrow$ local minimum

$\mathbf{H}(f)(\mathbf{x}^*)$ negative definite $\rightarrow f(\mathbf{x}) < f(\mathbf{x}^*) \rightarrow$ local maximum

Theorem

Second Derivative Test

- If $\mathbf{H}(f)(\mathbf{a})$ is *positive definite* then \mathbf{a} is a local minimum
- If $\mathbf{H}(f)(\mathbf{a})$ is *negative definite* then \mathbf{a} is a local maximum
- If $\mathbf{H}(f)(\mathbf{a})$ is *indefinite* then \mathbf{a} is a saddle point

Second Derivative Test

Many ways to assess definiteness \rightsquigarrow use determinant

Theorem

Two Dimensional, Second Derivative Test. Suppose $f : X \rightarrow \mathbb{R}$ with $X \subset \mathbb{R}^2$ and f twice differentiable. Write the *Hessian* of f at a critical value \mathbf{a} ,

$$\mathbf{H}(f)(\mathbf{a}) = \begin{pmatrix} A & B \\ B & C \end{pmatrix}$$

Then, we can conduct the second derivative test as:

- $AC - B^2 > 0$ and $A > 0$ \rightsquigarrow *positive definite* \rightsquigarrow \mathbf{a} is a local minimum
- $AC - B^2 > 0$ and $A < 0$ \rightsquigarrow *negative definite* \rightsquigarrow \mathbf{a} is a local maximum
- $AC - B^2 < 0$ \rightsquigarrow *indefinite* \rightsquigarrow saddle point
- $AC - B^2 = 0$ *inconclusive*

Multivariate Recipe

- 1) Calculate **gradient**
- 2) Set equal to zero, solve system of equations
- 3) Calculate **Hessian**
- 4) Assess **Hessian** at critical values
- 5) **Boundary values?** (if relevant)

Example 1: A Simple Optimization Problem

Suppose $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ with

$$f(x_1, x_2) = 3(x_1 + 2)^2 + 4(x_2 + 4)^2$$

Calculate gradient

$$\nabla f(\mathbf{x}) = (6x_1 + 12, 8x_2 + 32)$$

$$\mathbf{0} = (6x_1^* + 12, 8x_2^* + 32)$$

We now solve the system of equations to yield $x_1^* = -2$ and $x_2^* = -4$

Example 1: A Simple Optimization Problem

$$\mathbf{H}(f)(\mathbf{x}^*) = \begin{pmatrix} 6 & 0 \\ 0 & 8 \end{pmatrix}$$

$\det(\mathbf{H}(f)(\mathbf{x}^*)) = 48$ and $6 > 0$ so $\mathbf{H}(f)(\mathbf{x}^*)$ is positive definite. **local minimum**

Example 2: Two Dimensional Ideal Points

Suppose legislators are considering legislation $\mathbf{x} \in \mathbb{R}^2$. And suppose legislator i has utility function $U_i : \mathbb{R}^2 \rightarrow \mathbb{R}$,

$$U(\mathbf{x})_i = -(x_1 - \mu_1)^2 - (x_2 - \mu_2)^2$$

What is legislator i 's **optimal** policy?

$$\nabla f(\mathbf{x}) = (-2(x_1 - \mu_1), -2(x_2 - \mu_2))$$

$$\nabla f(\mathbf{x}) = \mathbf{0}$$

$$-2(x_1^* - \mu_1) = 0$$

$$-2(x_2^* - \mu_2) = 0$$

Solving yields $x_1^* = \mu_1$ and $x_2^* = \mu_2$.

Example 2: Two Dimensional Ideal Points

$$U(\mathbf{x})_i = -(x_1 - \mu_1)^2 - (x_2 - \mu_2)^2$$

Call $\boldsymbol{\mu} = (\mu_1, \mu_2)$

The Hessian at the critical value is

$$\begin{aligned} \mathbf{H}(f)(\boldsymbol{\mu}) &= \begin{pmatrix} \frac{\partial^2 U_i}{\partial x_1 \partial x_1}(\boldsymbol{\mu}) & \frac{\partial^2 U_i}{\partial x_1 \partial x_2}(\boldsymbol{\mu}) \\ \frac{\partial^2 U_i}{\partial x_2 \partial x_1}(\boldsymbol{\mu}) & \frac{\partial^2 U_i}{\partial x_2 \partial x_2}(\boldsymbol{\mu}) \end{pmatrix} \\ &= \begin{pmatrix} -2 & 0 \\ 0 & -2 \end{pmatrix} \end{aligned}$$

So, $-2 * -2 - 0 = 4 > 0$ and $-2 < 0 \rightsquigarrow$ **negative definite, maximum**
 $\boldsymbol{\mu} = (\mu_1, \mu_2)$ are legislator i 's two dimensional ideal point.

Example 3: Maximum Likelihood Estimation, Normal Distribution

Example 3: Maximum Likelihood Estimation, Normal Distribution

Suppose that we draw an independent and identically distributed random sample of n observations from a normal distribution,

Example 3: Maximum Likelihood Estimation, Normal Distribution

Suppose that we draw an independent and identically distributed random sample of n observations from a normal distribution,

$$Y_i \sim \text{Normal}(\mu, \sigma^2)$$

Example 3: Maximum Likelihood Estimation, Normal Distribution

Suppose that we draw an independent and identically distributed random sample of n observations from a normal distribution,

$$Y_i \sim \text{Normal}(\mu, \sigma^2)$$
$$\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)$$

Example 3: Maximum Likelihood Estimation, Normal Distribution

Suppose that we draw an independent and identically distributed random sample of n observations from a normal distribution,

$$Y_i \sim \text{Normal}(\mu, \sigma^2)$$
$$\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)$$

Our task:

Example 3: Maximum Likelihood Estimation, Normal Distribution

Suppose that we draw an independent and identically distributed random sample of n observations from a normal distribution,

$$Y_i \sim \text{Normal}(\mu, \sigma^2)$$
$$\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)$$

Our task:

- Obtain likelihood (summary estimator)

Example 3: Maximum Likelihood Estimation, Normal Distribution

Suppose that we draw an independent and identically distributed random sample of n observations from a normal distribution,

$$Y_i \sim \text{Normal}(\mu, \sigma^2)$$
$$\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)$$

Our task:

- Obtain likelihood (summary estimator)
- Derive maximum likelihood estimators for μ and σ^2

Example 3: Maximum Likelihood Estimation, Normal Distribution

Suppose that we draw an independent and identically distributed random sample of n observations from a normal distribution,

$$Y_i \sim \text{Normal}(\mu, \sigma^2)$$
$$\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)$$

Our task:

- Obtain likelihood (summary estimator)
- Derive maximum likelihood estimators for μ and σ^2
- Characterize sampling distribution

Example 3: Maximum Likelihood Estimation, Normal Distribution

Example 3: Maximum Likelihood Estimation, Normal Distribution

$$L(\mu, \sigma^2 | \mathbf{Y}) \propto \prod_{i=1}^n f(Y_i | \mu, \sigma^2)$$

Example 3: Maximum Likelihood Estimation, Normal Distribution

$$\begin{aligned} L(\mu, \sigma^2 | \mathbf{Y}) &\propto \prod_{i=1}^n f(Y_i | \mu, \sigma^2) \\ &\propto \prod_{i=1}^n \frac{\exp\left[-\frac{(Y_i - \mu)^2}{2\sigma^2}\right]}{\sqrt{2\pi\sigma^2}} \end{aligned}$$

Example 3: Maximum Likelihood Estimation, Normal Distribution

$$\begin{aligned}L(\mu, \sigma^2 | \mathbf{Y}) &\propto \prod_{i=1}^n f(Y_i | \mu, \sigma^2) \\ &\propto \prod_{i=1}^n \frac{\exp[-\frac{(Y_i - \mu)^2}{2\sigma^2}]}{\sqrt{2\pi\sigma^2}} \\ &\propto \frac{\exp[-\sum_{i=1}^n \frac{(Y_i - \mu)^2}{2\sigma^2}]}{(2\pi)^{n/2} \sigma^{2n/2}}\end{aligned}$$

Example 3: Maximum Likelihood Estimation, Normal Distribution

$$\begin{aligned}L(\mu, \sigma^2 | \mathbf{Y}) &\propto \prod_{i=1}^n f(Y_i | \mu, \sigma^2) \\ &\propto \prod_{i=1}^n \frac{\exp[-\frac{(Y_i - \mu)^2}{2\sigma^2}]}{\sqrt{2\pi\sigma^2}} \\ &\propto \frac{\exp[-\sum_{i=1}^n \frac{(Y_i - \mu)^2}{2\sigma^2}]}{(2\pi)^{n/2} \sigma^{2n/2}}\end{aligned}$$

Taking the logarithm, we have

Example 3: Maximum Likelihood Estimation, Normal Distribution

$$\begin{aligned}L(\mu, \sigma^2 | \mathbf{Y}) &\propto \prod_{i=1}^n f(Y_i | \mu, \sigma^2) \\ &\propto \prod_{i=1}^n \frac{\exp[-\frac{(Y_i - \mu)^2}{2\sigma^2}]}{\sqrt{2\pi\sigma^2}} \\ &\propto \frac{\exp[-\sum_{i=1}^n \frac{(Y_i - \mu)^2}{2\sigma^2}]}{(2\pi)^{n/2} \sigma^{2n/2}}\end{aligned}$$

Taking the logarithm, we have

$$l(\mu, \sigma^2 | \mathbf{Y}) = -\sum_{i=1}^n \frac{(Y_i - \mu)^2}{2\sigma^2} - \frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) + c$$

Example 3: Maximum Likelihood Estimation, Normal Distribution

$$\begin{aligned}L(\mu, \sigma^2 | \mathbf{Y}) &\propto \prod_{i=1}^n f(Y_i | \mu, \sigma^2) \\ &\propto \prod_{i=1}^n \frac{\exp[-\frac{(Y_i - \mu)^2}{2\sigma^2}]}{\sqrt{2\pi\sigma^2}} \\ &\propto \frac{\exp[-\sum_{i=1}^n \frac{(Y_i - \mu)^2}{2\sigma^2}]}{(2\pi)^{n/2} \sigma^{2n/2}}\end{aligned}$$

Taking the logarithm, we have

$$\begin{aligned}l(\mu, \sigma^2 | \mathbf{Y}) &= -\sum_{i=1}^n \frac{(Y_i - \mu)^2}{2\sigma^2} - \frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) + c \\ &= -\sum_{i=1}^n \frac{(Y_i - \mu)^2}{2\sigma^2} - \frac{n}{2} \log(\sigma^2) + c'\end{aligned}$$

Example 3: Log-Likelihood Plot

- In **R**, drew 10,000 realizations from

Example 3: Log-Likelihood Plot

- In **R**, drew 10,000 realizations from

$$Y_i \sim \text{Normal}(0.25, 100)$$

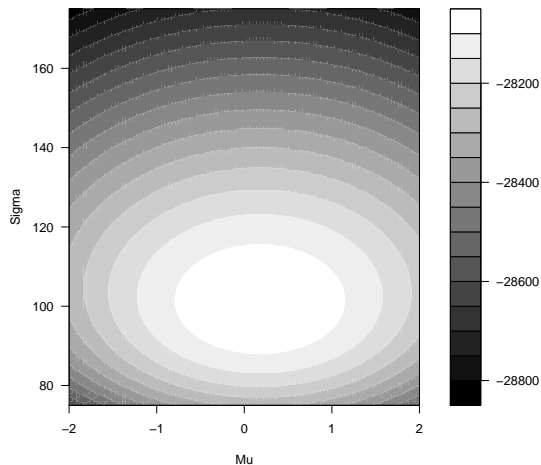
Example 3: Log-Likelihood Plot

- In **R**, drew 10,000 realizations from

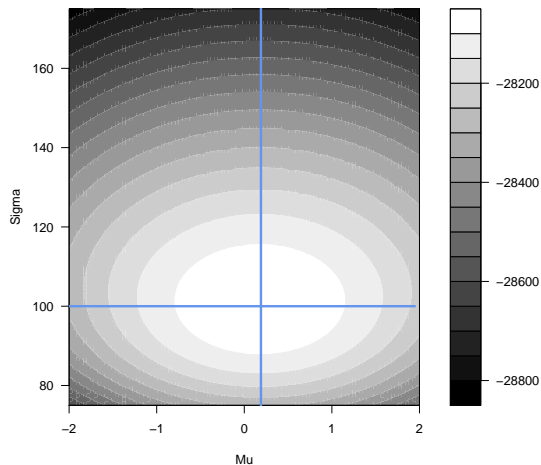
$$Y_i \sim \text{Normal}(0.25, 100)$$

- Used realized values y_i evaluate $l(\mu, \sigma^2 | \mathbf{y})$

Example 3: Log-Likelihood Plot



Example 3: Log-Likelihood Plot



Example 3: Maximum Likelihood Estimation, Normal Distribution

Let's find $\hat{\mu}$ and $\hat{\sigma}^2$ that maximizes log-likelihood.

Example 3: Maximum Likelihood Estimation, Normal Distribution

Let's find $\hat{\mu}$ and $\hat{\sigma}^2$ that maximizes log-likelihood.

$$l(\mu, \sigma^2 | \mathbf{Y}) = - \sum_{i=1}^n \frac{(Y_i - \mu)^2}{2\sigma^2} - \frac{n}{2} \log(\sigma^2) + c'$$

Example 3: Maximum Likelihood Estimation, Normal Distribution

Let's find $\hat{\mu}$ and $\hat{\sigma}^2$ that maximizes log-likelihood.

$$l(\mu, \sigma^2 | \mathbf{Y}) = -\sum_{i=1}^n \frac{(Y_i - \mu)^2}{2\sigma^2} - \frac{n}{2} \log(\sigma^2) + c'$$
$$\frac{\partial l(\mu, \sigma^2 | \mathbf{Y})}{\partial \mu} = \sum_{i=1}^n \frac{2(Y_i - \mu)}{2\sigma^2}$$

Example 3: Maximum Likelihood Estimation, Normal Distribution

Let's find $\hat{\mu}$ and $\hat{\sigma}^2$ that maximizes log-likelihood.

$$l(\mu, \sigma^2 | \mathbf{Y}) = -\sum_{i=1}^n \frac{(Y_i - \mu)^2}{2\sigma^2} - \frac{n}{2} \log(\sigma^2) + c'$$

$$\frac{\partial l(\mu, \sigma^2) | \mathbf{Y}}{\partial \mu} = \sum_{i=1}^n \frac{2(Y_i - \mu)}{2\sigma^2}$$

$$\frac{\partial l(\mu, \sigma^2) | \mathbf{Y}}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (Y_i - \mu)^2$$

Example 3: Maximum Likelihood Estimation, Normal Distribution

$$0 = -\sum_{i=1}^n \frac{2(Y_i - \hat{\mu})}{2\hat{\sigma}^2}$$
$$0 = -\frac{n}{2\hat{\sigma}^2} + \frac{1}{2\hat{\sigma}^4} \sum_{i=1}^n (Y_i - \mu^*)^2$$

Example 3: Maximum Likelihood Estimation, Normal Distribution

$$0 = -\sum_{i=1}^n \frac{2(Y_i - \hat{\mu})}{2\hat{\sigma}^2}$$
$$0 = -\frac{n}{2\hat{\sigma}^2} + \frac{1}{2\hat{\sigma}^4} \sum_{i=1}^n (Y_i - \mu^*)^2$$

Solving for $\hat{\mu}$ and $\hat{\sigma}^2$ yields,

Example 3: Maximum Likelihood Estimation, Normal Distribution

$$0 = -\sum_{i=1}^n \frac{2(Y_i - \hat{\mu})}{2\hat{\sigma}^2}$$
$$0 = -\frac{n}{2\hat{\sigma}^2} + \frac{1}{2\hat{\sigma}^4} \sum_{i=1}^n (Y_i - \mu^*)^2$$

Solving for $\hat{\mu}$ and $\hat{\sigma}^2$ yields,

$$\hat{\mu} = \frac{\sum_{i=1}^n Y_i}{n}$$

Example 3: Maximum Likelihood Estimation, Normal Distribution

$$0 = -\sum_{i=1}^n \frac{2(Y_i - \hat{\mu})}{2\hat{\sigma}^2}$$
$$0 = -\frac{n}{2\hat{\sigma}^2} + \frac{1}{2\hat{\sigma}^4} \sum_{i=1}^n (Y_i - \mu^*)^2$$

Solving for $\hat{\mu}$ and $\hat{\sigma}^2$ yields,

$$\hat{\mu} = \frac{\sum_{i=1}^n Y_i}{n}$$
$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

Example 3: Maximum Likelihood Estimation, Normal Distribution

Example 3: Maximum Likelihood Estimation, Normal Distribution

$$\mathbf{H}(f)(\hat{\mu}, \hat{\sigma}^2) = \begin{pmatrix} \frac{\partial^2 l(\mu, \sigma^2 | \mathbf{Y})}{\partial \mu^2} & \frac{\partial^2 l(\mu, \sigma^2 | \mathbf{Y})}{\partial \sigma^2 \partial \mu} \\ \frac{\partial^2 l(\mu, \sigma^2 | \mathbf{Y})}{\partial \sigma^2 \partial \mu} & \frac{\partial^2 l(\mu, \sigma^2 | \mathbf{Y})}{\partial^2 \sigma^2} \end{pmatrix}$$

Example 3: Maximum Likelihood Estimation, Normal Distribution

$$\mathbf{H}(f)(\hat{\mu}, \hat{\sigma}^2) = \begin{pmatrix} \frac{\partial^2 l(\mu, \sigma^2 | \mathbf{Y})}{\partial \mu^2} & \frac{\partial^2 l(\mu, \sigma^2 | \mathbf{Y})}{\partial \sigma^2 \partial \mu} \\ \frac{\partial^2 l(\mu, \sigma^2 | \mathbf{Y})}{\partial \sigma^2 \partial \mu} & \frac{\partial^2 l(\mu, \sigma^2 | \mathbf{Y})}{\partial^2 \sigma^2} \end{pmatrix}$$

Taking derivatives and evaluating at MLE's yields,

Example 3: Maximum Likelihood Estimation, Normal Distribution

$$\mathbf{H}(f)(\hat{\mu}, \hat{\sigma}^2) = \begin{pmatrix} \frac{\partial^2 l(\mu, \sigma^2 | \mathbf{Y})}{\partial \mu^2} & \frac{\partial^2 l(\mu, \sigma^2 | \mathbf{Y})}{\partial \sigma^2 \partial \mu} \\ \frac{\partial^2 l(\mu, \sigma^2 | \mathbf{Y})}{\partial \sigma^2 \partial \mu} & \frac{\partial^2 l(\mu, \sigma^2 | \mathbf{Y})}{\partial^2 \sigma^2} \end{pmatrix}$$

Taking derivatives and evaluating at MLE's yields,

$$\mathbf{H}(f)(\hat{\mu}, \hat{\sigma}^2) = \begin{pmatrix} \frac{-n}{\hat{\sigma}^2} & 0 \\ 0 & \frac{-n}{(\hat{\sigma}^2)^2} \end{pmatrix}$$

Example 3: Maximum Likelihood Estimation, Normal Distribution

$$\mathbf{H}(f)(\hat{\mu}, \hat{\sigma}^2) = \begin{pmatrix} \frac{\partial^2 l(\mu, \sigma^2 | \mathbf{Y})}{\partial \mu^2} & \frac{\partial^2 l(\mu, \sigma^2 | \mathbf{Y})}{\partial \sigma^2 \partial \mu} \\ \frac{\partial^2 l(\mu, \sigma^2 | \mathbf{Y})}{\partial \sigma^2 \partial \mu} & \frac{\partial^2 l(\mu, \sigma^2 | \mathbf{Y})}{\partial^2 \sigma^2} \end{pmatrix}$$

Taking derivatives and evaluating at MLE's yields,

$$\mathbf{H}(f)(\hat{\mu}, \hat{\sigma}^2) = \begin{pmatrix} \frac{-n}{\hat{\sigma}^2} & 0 \\ 0 & \frac{-n}{(\hat{\sigma}^2)^2} \end{pmatrix}$$

$$\det(\mathbf{H}(f)(\hat{\mu}, \hat{\sigma}^2)) = n^2 / \hat{\sigma}^5 \text{ and } -n / \hat{\sigma}^2 < 0 \rightsquigarrow \text{maximum}$$

Computational Optimization

Analytic solutions: often hard.

Computational Optimization

Analytic solutions: often hard.

Computational solutions: simplify. Trade offs

Computational Optimization

Analytic solutions: often hard.

Computational solutions: simplify. Trade offs

- Newton-Raphson: expensive

Computational Optimization

Analytic solutions: often hard.

Computational solutions: simplify. Trade offs

- Newton-Raphson: expensive
- BFGS: less expensive

Computational Optimization

Analytic solutions: often hard.

Computational solutions: simplify. Trade offs

- Newton-Raphson: expensive
- BFGS: less expensive
- EM-like optimization: solve intractable problems, parallelizable

Multivariate Newton Raphson

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Suppose we have guess \mathbf{x}_t .

Multivariate Newton Raphson

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Suppose we have guess \mathbf{x}_t . Then our update is:

Multivariate Newton Raphson

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Suppose we have guess \mathbf{x}_t . Then our update is:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \mathbf{H}(f)(\mathbf{x}_t)^{-1} \nabla f(\mathbf{x}_t)$$

Multivariate Newton Raphson

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Suppose we have guess \mathbf{x}_t . Then our update is:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \mathbf{H}(f)(\mathbf{x}_t)^{-1} \nabla f(\mathbf{x}_t)$$

Derivation (intuition):

Multivariate Newton Raphson

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Suppose we have guess \mathbf{x}_t . Then our update is:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \mathbf{H}(f)(\mathbf{x}_t)^{-1} \nabla f(\mathbf{x}_t)$$

Derivation (intuition): Approximate function with **tangent plane**.

Multivariate Newton Raphson

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Suppose we have guess \mathbf{x}_t . Then our update is:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \mathbf{H}(f)(\mathbf{x}_t)^{-1} \nabla f(\mathbf{x}_t)$$

Derivation (intuition): Approximate function with **tangent plane**. Find value of x_{t+1} that makes the plane equal to zero. Update again.

Multivariate Newton Raphson

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Suppose we have guess \mathbf{x}_t . Then our update is:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \mathbf{H}(f)(\mathbf{x}_t)^{-1} \nabla f(\mathbf{x}_t)$$

Derivation (intuition): Approximate function with **tangent plane**. Find value of x_{t+1} that makes the plane equal to zero. Update again.

R Code

Multivariate Newton Raphson

Multivariate Newton Raphson

- Expensive to calculate (requires inverting Hessian)

Multivariate Newton Raphson

- Expensive to calculate (requires inverting Hessian)
- Very sensitive to starting points

Multivariate Newton Raphson

- Expensive to calculate (requires inverting Hessian)
- Very sensitive to starting points
- Ideally: method that exploits Newton-like structure, but is cheaper and more robust

Multivariate Newton Raphson

- Expensive to calculate (requires inverting Hessian)
- Very sensitive to starting points
- Ideally: method that exploits Newton-like structure, but is cheaper and more robust

BFGS: **Quasi-Newton** method

Multivariate Newton Raphson

- Expensive to calculate (requires inverting Hessian)
- Very sensitive to starting points
- Ideally: method that exploits Newton-like structure, but is cheaper and more robust

BFGS: **Quasi-Newton** method

R code

Optimization that is Both Discrete and Continuous

K-means: most commonly used clustering algorithm.

Optimization that is Both Discrete and Continuous

K-means: most commonly used clustering algorithm.

Story: Data are grouped in K clusters and each cluster has a **center** or mean.

Optimization that is Both Discrete and Continuous

K-means: most commonly used clustering algorithm.

Story: Data are grouped in K clusters and each cluster has a **center** or mean.

→ Two **types** of parameters to estimate

Optimization that is Both Discrete and Continuous

K-means: most commonly used clustering algorithm.

Story: Data are grouped in K clusters and each cluster has a **center** or mean.

→ Two **types** of parameters to estimate

1) For each cluster j , ($j = 1, \dots, K$)

r_{ij} = Indicator, Document i assigned to cluster j

$\mathbf{r}_j = (r_{1j}, r_{2j}, \dots, r_{Nj})$

$\mathbf{r} = (\mathbf{r}'_1, \mathbf{r}'_2, \dots, \mathbf{r}'_K)$ ($N \times K$ matrix)

Optimization that is Both Discrete and Continuous

K-means: most commonly used clustering algorithm.

Story: Data are grouped in K clusters and each cluster has a **center** or mean.

→ Two **types** of parameters to estimate

1) For each cluster j , ($j = 1, \dots, K$)

r_{ij} = Indicator, Document i assigned to cluster j

$$\mathbf{r}_j = (r_{1j}, r_{2j}, \dots, r_{Nj})$$

$$\mathbf{r} = (\mathbf{r}'_1, \mathbf{r}'_2, \dots, \mathbf{r}'_K) \quad (N \times K \text{ matrix})$$

2) For each cluster j

μ_j a **cluster center** for cluster j .

$$\mu_j = (\mu_{1j}, \mu_{2j}, \dots, \mu_{Mj})$$

Optimization that is Both Discrete and Continuous

K-means: most commonly used clustering algorithm.

Story: Data are grouped in K clusters and each cluster has a **center** or mean.

→ Two **types** of parameters to estimate

1) For each cluster j , ($j = 1, \dots, K$)

r_{ij} = Indicator, Document i assigned to cluster j

$\mathbf{r}_j = (r_{1j}, r_{2j}, \dots, r_{Nj})$

$\mathbf{r} = (\mathbf{r}'_1, \mathbf{r}'_2, \dots, \mathbf{r}'_K)$ ($N \times K$ matrix)

2) For each cluster j

μ_j a **cluster center** for cluster j .

$\mu_j = (\mu_{1j}, \mu_{2j}, \dots, \mu_{Mj})$

Notation. Representation of document i :

$$\mathbf{y}_i = (y_{i1}, y_{i2}, \dots, y_{iM})$$

Specifying the Method

- 1) Assume Euclidean distance between objects.
- 2) **Objective function**

$$f(\mathbf{r}, \boldsymbol{\mu}, \mathbf{y}) = \sum_{i=1}^N \sum_{j=1}^K r_{ij} \left(\sum_{m=1}^M (y_{im} - \mu_{km})^2 \right)$$

Specifying the Method

- 1) Assume Euclidean distance between objects.
- 2) **Objective function**

$$f(\mathbf{r}, \boldsymbol{\mu}, \mathbf{y}) = \sum_{i=1}^N \sum_{j=1}^K r_{ij} \left(\sum_{m=1}^M (y_{im} - \mu_{km})^2 \right)$$

Goal:

Specifying the Method

- 1) Assume Euclidean distance between objects.
- 2) **Objective function**

$$f(\mathbf{r}, \boldsymbol{\mu}, \mathbf{y}) = \sum_{i=1}^N \sum_{j=1}^K r_{ij} \left(\sum_{m=1}^M (y_{im} - \mu_{km})^2 \right)$$

Goal:

Choose \mathbf{r}^* and $\boldsymbol{\mu}^*$ to minimize $f(\mathbf{r}, \boldsymbol{\mu}, \mathbf{y})$

Specifying the Method

- 1) Assume Euclidean distance between objects.
- 2) **Objective function**

$$f(\mathbf{r}, \boldsymbol{\mu}, \mathbf{y}) = \sum_{i=1}^N \sum_{j=1}^K r_{ij} \left(\sum_{m=1}^M (y_{im} - \mu_{km})^2 \right)$$

Goal:

Choose \mathbf{r}^* and $\boldsymbol{\mu}^*$ to minimize $f(\mathbf{r}, \boldsymbol{\mu}, \mathbf{y})$

Two observations:

- If $K = N$ $f(\mathbf{r}^*, \boldsymbol{\mu}^*, \mathbf{y}) = 0$ (Minimum)
 - Each observation in own cluster
 - $\boldsymbol{\mu}_i = \mathbf{y}_i$
- If $K = 1$, $f(\mathbf{r}^*, \boldsymbol{\mu}^*, \mathbf{y}) = N \times \sigma^2$
 - Each observation in one cluster
 - Center: average of documents

Specifying the Method

- 1) Assume Euclidean distance between objects
- 2) Objective function
- 3) Algorithm for optimization

Iterative algorithm, Each Iteration t

- Conditional on μ^{t-1} (from previous iteration), choose r^t
- Conditional on r^t , choose μ^t

Repeat until convergence, measured as change in f .

$$\text{Change} = f(\mu^t, r^t, \mathbf{y}) - f(\mu^{t-1}, r^{t-1}, \mathbf{y})$$

Specifying the Method

$$f(\mathbf{r}, \boldsymbol{\mu}, \mathbf{y}) = \sum_{i=1}^N \sum_{j=1}^K r_{ij} \left(\sum_{m=1}^M (y_{im} - \mu_{km})^2 \right)$$

Algorithm for estimation:

Begin: initialize $\boldsymbol{\mu}_1^{t-1}, \boldsymbol{\mu}_2^{t-1}, \dots, \boldsymbol{\mu}_K^{t-1}$

Choose \mathbf{r}^t

$$r_{ij}^t = \begin{cases} 1 & \text{if } j = \arg \min_k \sum_{m=1}^M (y_{im} - \mu_{km})^2 \\ 0 & \text{otherwise,} \end{cases}$$

In words: Assign each document \mathbf{y}_i to the closest center $\boldsymbol{\mu}_k$

$$f(\mathbf{r}, \boldsymbol{\mu}, \mathbf{y}) = \sum_{i=1}^N \sum_{j=1}^K r_{ij} \left(\sum_{m=1}^M (y_{im} - \mu_{km})^2 \right)$$

Conditional on \mathbf{r}^t , choose $\boldsymbol{\mu}^t$

Let's focus on $\boldsymbol{\mu}_k$

$$f(\mathbf{r}, \boldsymbol{\mu}_k, \mathbf{y})_k = \sum_{i=1}^N r_{ik} \left(\sum_{m=1}^M (y_{im} - \mu_{km})^2 \right)$$

Focus on just μ_{km}

$$f(\mathbf{r}, \mu_{km}, \mathbf{y})_{km} = \sum_{i=1}^N r_{ik} (y_{im} - \mu_{km})^2$$

Quadratic: take derivative, set equal to zero (second derivative test works)

$$\frac{\partial f(\mathbf{r}, \mu_{km}, \mathbf{y})_{km}}{\partial \mu_{km}} = -2 \sum_{i=1}^N r_{ik} (y_{im} - \mu_{km})$$

$$2 \sum_{i=1}^N r_{ik} (y_{im} - \mu_{km}^t) = 0$$

$$\sum_{i=1}^N r_{ik} y_{im} - \mu_{km}^t \sum_{i=1}^N r_{ik} = 0$$

$$\frac{\sum_{i=1}^N r_{ik} y_{im}}{\sum_{i=1}^N r_{ik}} = \mu_{km}^t$$

$$\mu_k^t = \frac{\sum_{i=1}^N r_{ik} \mathbf{y}_i}{\sum_{i=1}^N r_{ik}}$$

In words:

- μ_k^t is the average of documents assigned to the k^{th} cluster

Algorithm, In Words

- Conditional on center estimates, assign documents to closest cluster centers
- Conditional on document assignments, cluster centers are averages of documents assigned to the cluster

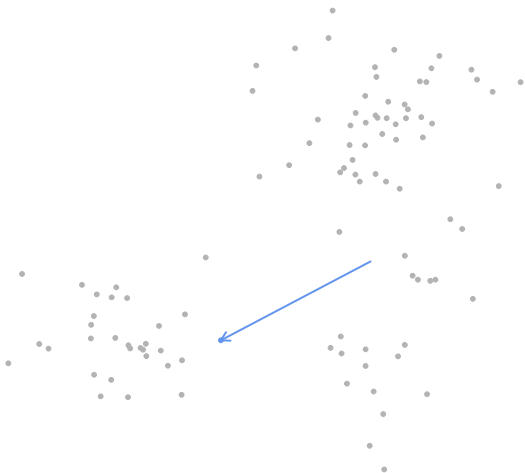
Expectation-Maximization (EM) [connection guarantees convergence]

- Estimation of $r \rightsquigarrow$ Expectation step (data augmentation)
- Estimation of $\mu_k \rightsquigarrow$ Maximization Step

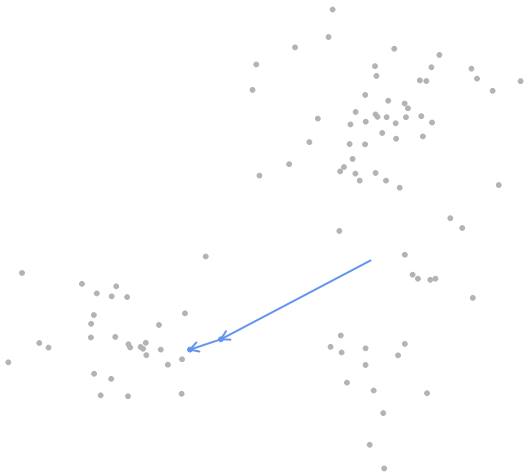
Visual Example



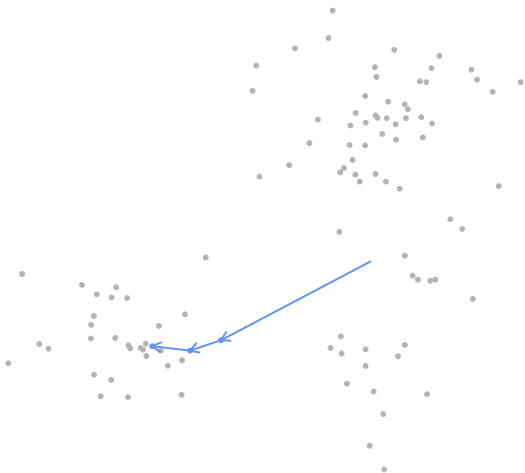
Visual Example



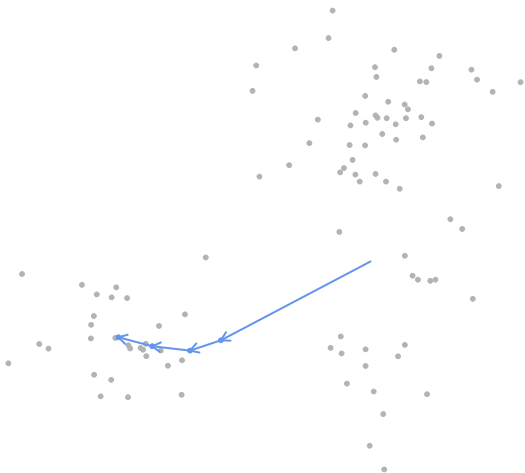
Visual Example



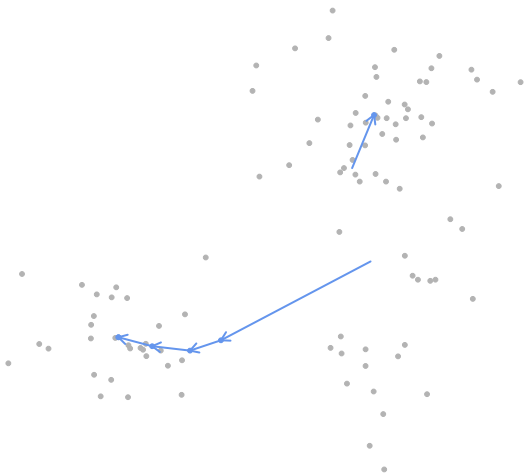
Visual Example



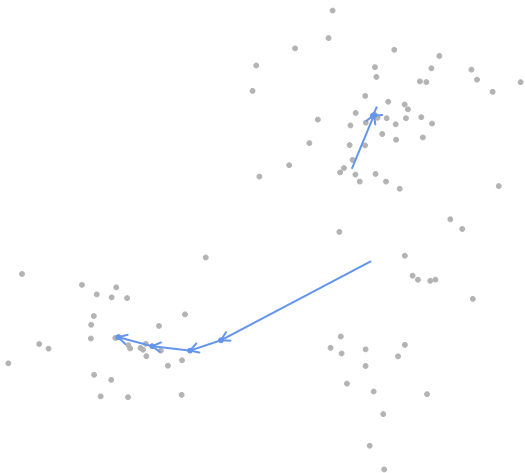
Visual Example



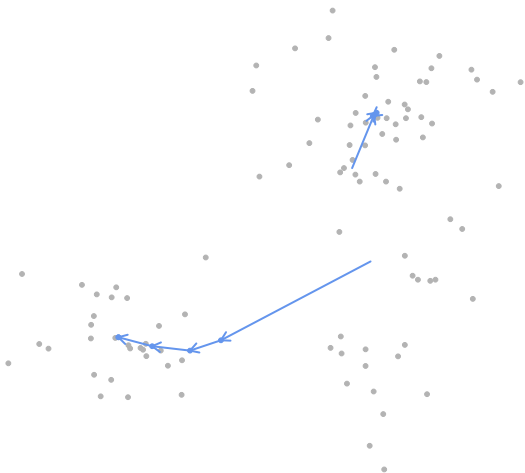
Visual Example



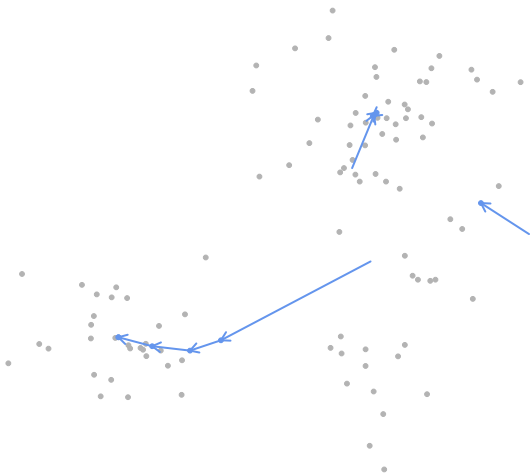
Visual Example



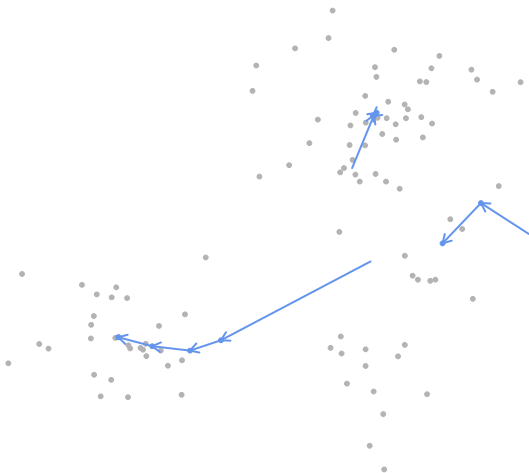
Visual Example



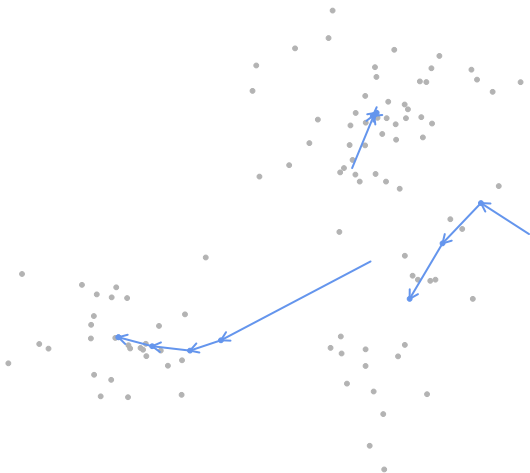
Visual Example



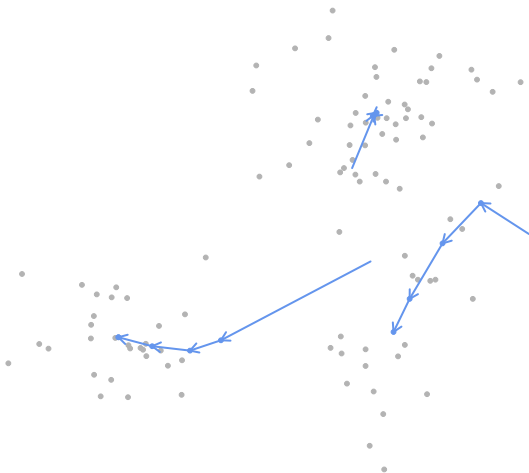
Visual Example



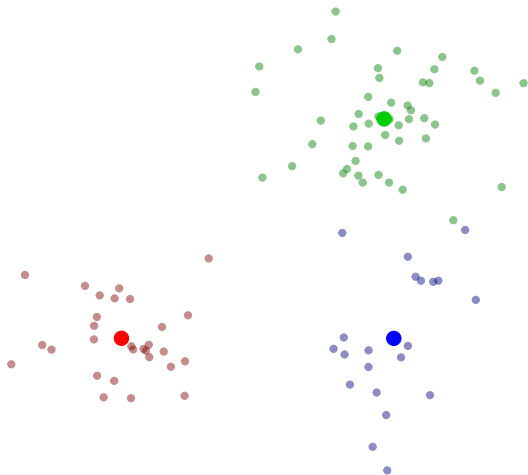
Visual Example



Visual Example



Visual Example



Many Optimization Procedures!!!

Many Optimization Procedures!!!

Nelder Mead:

Many Optimization Procedures!!!

Nelder Mead:

- Evaluate points on a simplex (triangle)

Many Optimization Procedures!!!

Nelder Mead:

- Evaluate points on a simplex (triangle)
- Either Reflect, Expand, or Contract (based on values)

Many Optimization Procedures!!!

Nelder Mead:

- Evaluate points on a simplex (triangle)
- Either Reflect, Expand, or Contract (based on values)
- Converges to local extrema

Many Optimization Procedures!!!

Nelder Mead:

- Evaluate points on a simplex (triangle)
- Either Reflect, Expand, or Contract (based on values)
- Converges to local extrema

Stochastic Optimization:

Many Optimization Procedures!!!

Nelder Mead:

- Evaluate points on a simplex (triangle)
- Either Reflect, Expand, or Contract (based on values)
- Converges to local extrema

Stochastic Optimization:

- Sample a subset of data, perform optimization

Many Optimization Procedures!!!

Nelder Mead:

- Evaluate points on a simplex (triangle)
- Either Reflect, Expand, or Contract (based on values)
- Converges to local extrema

Stochastic Optimization:

- Sample a subset of data, perform optimization
- Sample a new subset, perform optimization, combine with previous sample

Many Optimization Procedures!!!

Nelder Mead:

- Evaluate points on a simplex (triangle)
- Either Reflect, Expand, or Contract (based on values)
- Converges to local extrema

Stochastic Optimization:

- Sample a subset of data, perform optimization
- Sample a new subset, perform optimization, combine with previous sample
- Converges on local extrema (given regulatory conditions)

Many Optimization Procedures!!!

Nelder Mead:

- Evaluate points on a simplex (triangle)
- Either Reflect, Expand, or Contract (based on values)
- Converges to local extrema

Stochastic Optimization:

- Sample a subset of data, perform optimization
- Sample a new subset, perform optimization, combine with previous sample
- Converges on local extrema (given regulatory conditions)

Genetic Optimization:

Many Optimization Procedures!!!

Nelder Mead:

- Evaluate points on a simplex (triangle)
- Either Reflect, Expand, or Contract (based on values)
- Converges to local extrema

Stochastic Optimization:

- Sample a subset of data, perform optimization
- Sample a new subset, perform optimization, combine with previous sample
- Converges on local extrema (given regulatory conditions)

Genetic Optimization:

- Evaluate fitness of solutions

Many Optimization Procedures!!!

Nelder Mead:

- Evaluate points on a simplex (triangle)
- Either Reflect, Expand, or Contract (based on values)
- Converges to local extrema

Stochastic Optimization:

- Sample a subset of data, perform optimization
- Sample a new subset, perform optimization, combine with previous sample
- Converges on local extrema (given regulatory conditions)

Genetic Optimization:

- Evaluate fitness of solutions
- Randomly select most fit, then combine

Many Optimization Procedures!!!

Nelder Mead:

- Evaluate points on a simplex (triangle)
- Either Reflect, Expand, or Contract (based on values)
- Converges to local extrema

Stochastic Optimization:

- Sample a subset of data, perform optimization
- Sample a new subset, perform optimization, combine with previous sample
- Converges on local extrema (given regulatory conditions)

Genetic Optimization:

- Evaluate fitness of solutions
- Randomly select most fit, then combine
- Can converge to **global** maximum, but might require extensive run time

Where We Are Going

- Done with math component
- Start probability tomorrow