

PS5841

Data Science in Finance & Insurance

Method of Least Squares

Yubo Wang

Spring 2022

Squared Error Loss for Prediction

- Find $f(\mathbf{X})$ for predicting Y given values of \mathbf{X} , where
 - Random input vector $\mathbf{X} \in \mathcal{R}^p$
 - Random output variable $Y \in \mathcal{R}$

- Squared error loss

$$Loss = [Y - f(\mathbf{X})]^2$$

- Expected (squared) prediction error

$$\begin{aligned} EPE(f) &= E([Y - f(\mathbf{X})]^2) \\ &= E_{\mathbf{X}} E_{Y|\mathbf{X}}([Y - f(\mathbf{X})]^2 | \mathbf{X}) \end{aligned}$$

Regression Function

- Expected (squared) prediction error

$$\begin{aligned} EPE(f) &= E([Y - f(\mathbf{X})]^2) \\ &= E_{\mathbf{X}} E_{Y|\mathbf{X}}([Y - f(\mathbf{X})]^2 | \mathbf{X}) \end{aligned}$$

- Minimize EPE pointwise

$$f(\mathbf{x}) = \underset{c}{\operatorname{argmin}} E_{Y|\mathbf{X}}([Y - c]^2 | \mathbf{X} = \mathbf{x})$$

- Regression function

$$f(\mathbf{x}) = E(Y | \mathbf{X} = \mathbf{x})$$

Scenario 1

- When an analytical solution for loss minimization is available
 - e.g. Linear Regression Function

Basic/Simple Linear Regression Model

- Training set of size n : $\{(x_i, y_i)\}$
 - y_i is the observed value of Y_i
 - The Y_i s are independent

- Regression function

$$E(Y|X = x) = \beta_0 + \beta_1 x$$

Fitted by the Method of Least Squares

- Assume $Var(Y_i) = \sigma^2 \forall i$ for now
- Minimize loss $\sum_i (y_i - \beta_0 - \beta_1 x_i)^2$
- Estimated parameters

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \sum_i \omega_i y_i$$

$$\omega_i = \frac{x_i - \bar{x}}{(n-1)s_X^2} \rightarrow \sum_i \omega_i = 0$$

$$s_X^2 = \frac{1}{n-1} \sum_i (x_i - \bar{x})^2$$

- Fitted model $\hat{y}(x) = \hat{\beta}_0 + \hat{\beta}_1 x$

Useful Results

- $\bar{\hat{y}} = \bar{y}$
- $\sum_i \hat{e}_i = 0, \hat{e}_i = \hat{y}_i - y_i$
- $\sum_i x_i \hat{e}_i = 0$
- $TSS = ESS + RSS$
$$\sum_i (y_i - \bar{y})^2 = \sum_i (\hat{y}_i - y_i)^2 + \sum_i (\hat{y}_i - \bar{y})^2$$
- Estimating σ^2 : $s^2 = \frac{1}{n-2} \sum_i \hat{e}_i^2$
- Coefficient of Determination: $R^2 = \frac{ESS}{TSS}$
- ANOVA Table: keep track of variability

Multiple Linear Regression Model

- Training set of size n : $\{(\mathbf{x}_i, y_i)\}$
 - y_i is the observed value of Y_i
 - The Y_i s are independent
 - $\text{Var}(\mathbf{Y}) = \mathbf{V}$
- Regression function

$$E(Y|\mathbf{X} = \mathbf{x}) = \mathbf{x}^T \boldsymbol{\beta}$$

Fitted by the Method of Least Squares (1)

- Minimize loss

$$\sum_i \frac{1}{\sigma_i^2} (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2$$
$$= (\mathbf{y} - \boldsymbol{\mu})^T V^{-1} (\mathbf{y} - \boldsymbol{\mu})$$

- Matrix of 2nd derivatives positive definite
- Solutions of the normal equations vs local minima at parameter space boundaries

Fitted by the Method of Least Squares (2)

- Estimated parameters (X_0 is the intercept if used)

$$\hat{\boldsymbol{\beta}} = (X^T V^{-1} X)^{-1} X^T V^{-1} \mathbf{y}$$

$$X = \begin{pmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{pmatrix} = (X_0 \quad X_1 \quad \cdots \quad X_k)$$

- If $V = \text{Var}(\mathbf{Y}) = \sigma^2 I$

$$\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y} = \sum_i \boldsymbol{\omega}_i y_i$$

$$\boldsymbol{\omega}_i = (X^T X)^{-1} (1, x_{i1}, \dots, x_{ik})^T$$

- Fitted model $\hat{y}(\mathbf{x}) = \mathbf{x}^T \hat{\boldsymbol{\beta}}$ or $\hat{\mathbf{y}}(X) = X \hat{\boldsymbol{\beta}}$

Useful Results (1)

When $E(\mathbf{Y}) = \mathbf{X}\boldsymbol{\beta}$, $\text{Var}(\mathbf{Y}) = \sigma^2 \mathbf{I}$, Y_i s are independent RV

- $E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$
- $\text{Var}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$
- LSE is BLUE (Gauss-Markov)
- $s^2 = \hat{\sigma}^2 = \frac{1}{n-p} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$

Useful Results (2)

- Hat Matrix $H = X(X^T X)^{-1} X^T$
 $\hat{\mathbf{y}} = X\hat{\boldsymbol{\beta}} = X(X^T X)^{-1} X^T \mathbf{y} = H(\mathbf{y})$
 $H^T = H, \quad HH = H$
 $HX = X, \quad HX_j = X_j$
- $\hat{\mathbf{e}}^T X_j = 0, \hat{\mathbf{e}} = (I - H)\mathbf{y}$
- $\mathbf{y}^T \hat{\mathbf{y}} = \hat{\mathbf{y}}^T \hat{\mathbf{y}}$
- $TSS = ESS + RSS$
$$TSS = \mathbf{y}^T \mathbf{y} - n(\bar{y})^2$$
$$ESS = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \hat{\mathbf{y}}$$
$$RSS = \hat{\mathbf{y}}^T \hat{\mathbf{y}} - n(\bar{y})^2$$
- Coefficient of Determination: $R^2 = \frac{RSS}{TSS} = (r_{\mathbf{y}, \hat{\mathbf{y}}})^2$
- ANOVA Table: keep track of variability

Normal Linear Model: Inference (1)

- $\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y} \sim \mathcal{N}(\boldsymbol{\beta}, (X^T X)^{-1} \sigma^2)$
- $\hat{\sigma}^2 = \frac{ESS}{n-p}$
- $H_0: \hat{\beta}_j = 0$ vs $H_1: \hat{\beta}_j \neq 0 \rightarrow t = \frac{\hat{\beta}_j}{\hat{\sigma} \sqrt{(X^T X)^{-1}_{jj}}} \sim t_{n-p}$
- 100(1 - α)% confidence interval
$$\hat{\beta}_j \pm t_{n-p, 1-\alpha/2} \hat{\sigma} \sqrt{(X^T X)^{-1}_{jj}}$$
- 100(1 - α)% prediction interval for y at \mathbf{x}
$$\mathbf{x}^T \hat{\boldsymbol{\beta}} \pm t_{n-p, 1-\alpha/2} \hat{\sigma} (1 + \mathbf{x}^T (X^T X)^{-1} \mathbf{x})^{1/2}$$

Normal Linear Model: Inference (2)

$$H_0: \hat{\boldsymbol{\beta}} = (\hat{\beta}_1, \dots, \hat{\beta}_{p_0})^T$$

$$H_1: \hat{\boldsymbol{\beta}} = (\hat{\beta}_1, \dots, \hat{\beta}_{p_1})^T$$

$$p_0 < p_1 < n$$

Equivalently, the general linear hypothesis

$$H_0: C\hat{\boldsymbol{\beta}} = \mathbf{0}$$

$$C = (0_{p_0}, I_{p_1-p_0}), \quad \hat{\boldsymbol{\beta}} = (\hat{\beta}_1, \dots, \hat{\beta}_{p_1})^T$$

$$F = \frac{ESS_0 - ESS_1}{p_1 - p_0} \div \frac{ESS_1}{n - p_1} \sim \mathcal{F}_{p_1-p_0, n-p_1}$$

Power Transforms

- To make the assumption of normality (if desired) more plausible

Strictly Positive Data

- Box-Cox family of transforms ($y > 0$)

$$y^{(\lambda)} = \begin{cases} \frac{y^\lambda - 1}{\lambda}, & \lambda \neq 0 \\ \ln y, & \lambda = 0 \end{cases}$$

- Estimate λ via maximum likelihood
- In practice,
 - Typically, $\lambda = 1, 0.5, 0, -1$
 - No need to -1 and $\div \lambda$ for operations unaffected by location and scale shifts (e.g. some regressions)

General Data

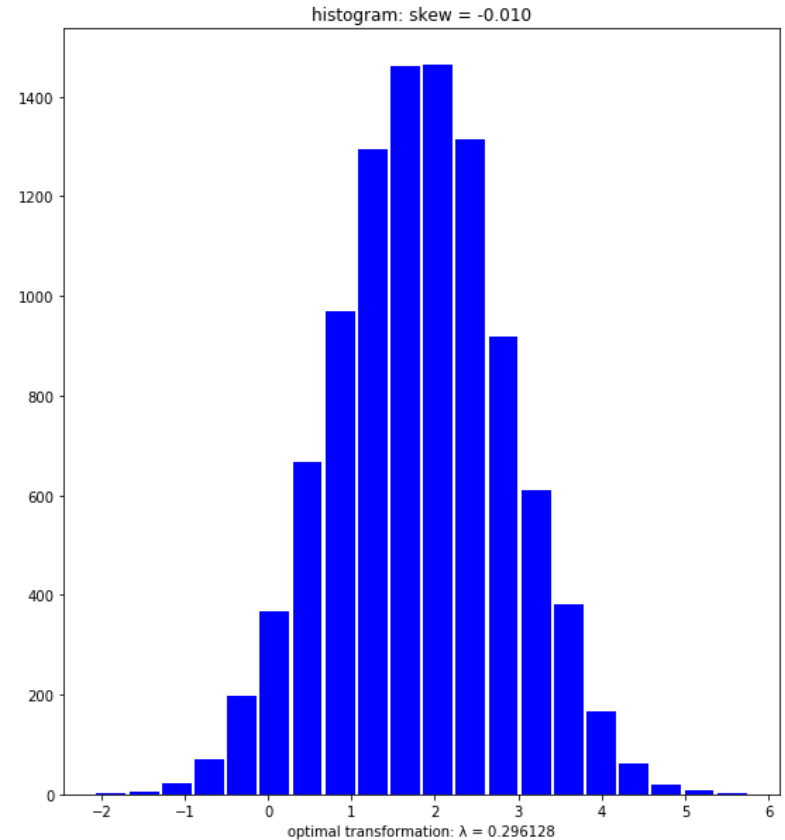
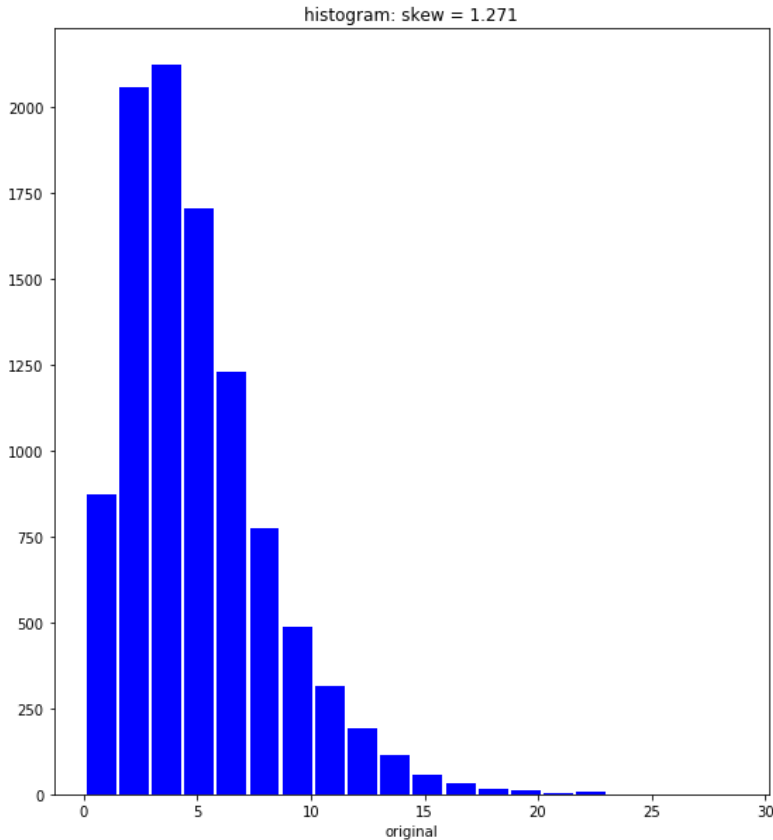
- Box-Cox family of transforms ($y > -\alpha$)

$$y^{(\lambda)} = \begin{cases} \frac{(y + \alpha)^\lambda - 1}{\lambda}, & \lambda \neq 0 \\ \ln(y + \alpha), & \lambda = 0 \end{cases}$$

- Yeo-Johnson family of transforms

$$y^{(\lambda)} = \begin{cases} \frac{(y + 1)^\lambda - 1}{\lambda}, & \lambda \neq 0, y \geq 0 \\ \ln(y + 1), & \lambda = 0, y \geq 0 \\ -\frac{(-y + 1)^{2-\lambda} - 1}{2 - \lambda}, & \lambda \neq 2, y < 0 \\ -\ln(-y + 1), & \lambda = 0, y < 0 \end{cases}$$

Box-Cox Transform Example



Categorical Features

- Often modeled by binary (dummy) variables
- For a feature (factor) with J levels
 - Need J binary variables if there is no intercept
 - Need (J-1) binary variables if there is intercept
 - Baseline is the level without a dummy variable
- Example: 1 factor with 3 levels (A,B,C)
 - with baseline A (A as the reference level)
$$y = \beta_0 + \beta_1 x_B + \beta_2 x_C + \beta_3 x_3 + \beta_4 x_4 + \epsilon$$
 - How are $\hat{\beta}_j$'s interpreted?

Example

- 1 factor with 3 levels (A,B,C)
 - with baseline A (A as the reference level)

$$y = \beta_0 + \beta_1 x_B + \beta_2 x_C + \beta_3 x_3 + \beta_4 x_4 + \epsilon$$

- If level A is present, all else being equal

$$\hat{y}_A = \hat{\beta}_0 + \hat{\beta}_3 x_3 + \hat{\beta}_4 x_4$$

- If level B is present, all else being equal

$$\hat{y}_B = \hat{\beta}_0 + \hat{\beta}_1 + \hat{\beta}_3 x_3 + \hat{\beta}_4 x_4$$

- Interpretation of $\hat{\beta}_1$

$$\hat{y}_B - \hat{y}_A = \hat{\beta}_1$$

Application

- Examples
 - Equity beta
 - Market model or CAPM
 - Demand for life insurance
 - Features are family characteristics that influence the amount of insurance purchased

Scenario 2

- When the loss needs to be numerically minimized
 - e.g. Non-Linear Regression Function
- (Stochastic) Gradient Descent
- Newton Raphson

Directional Derivative

- Directional derivative of $f: \mathcal{R}^p \rightarrow \mathcal{R}$ at \mathbf{x} in the direction of a unit vector \mathbf{u} , the rate of change of f at \mathbf{x} in the direction of \mathbf{u}

$$\mathbf{D}_{\mathbf{u}}f(\mathbf{x}) = \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{u}) - f(\mathbf{x})}{h} = \nabla f(\mathbf{x}) \cdot \mathbf{u}$$

- Gradient operator (del operator)

$$\nabla = \left(\frac{\partial}{\partial x_1}, \dots, \frac{\partial}{\partial x_p} \right)^T$$

- Gradient vector

$$\nabla(f) = \nabla f = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_p} \right)^T = \frac{\partial f}{\partial \mathbf{x}}$$

Maximum Rate of Change

- When $\nabla f(\mathbf{x}) \neq \mathbf{0}$
 - the maximum rate of increase of f is $|\nabla f(\mathbf{x})|$ and is in the direction of $\nabla f(\mathbf{x})$
 - the maximum rate of decrease of f is $|\nabla f(\mathbf{x})|$ and is in the direction of $-\nabla f(\mathbf{x})$

Example (1)

- $f(\mathbf{x}) = x_1 + x_2^2, \nabla f(\mathbf{x}) = \begin{pmatrix} 1 \\ 2x_2 \end{pmatrix}$
- The unit vector in the direction of the gradient at $\mathbf{x}_0 = (1,1)^T$ is $\mathbf{u} = \left(\frac{1}{\sqrt{5}}, \frac{2}{\sqrt{5}}\right)^T$ since

$$\nabla f(\mathbf{x}_0) = (1,2)^T, \quad |\nabla f(\mathbf{x}_0)| = \sqrt{5}$$

- Moving d units from $\mathbf{x}_0 = (1,1)^T$ in the direction of the gradient vector will land at $\mathbf{x}_1 = \mathbf{x}_0 + \mathbf{u}d = \left(1 + \frac{d}{\sqrt{5}}, 1 + \frac{2d}{\sqrt{5}}\right)^T$, with

$$f(\mathbf{x}_1) = f(\mathbf{x}_0) + \sqrt{5}d + \frac{4d^2}{5}$$

rate of
increase

predicted
increase

“error”
on large
steps

Example (2)

Same as
predicted
rate

gradient ascent						
	$f[1+(1/\sqrt{5})d,$					
d	$1+(2/\sqrt{5})d]$	$f(1,1)$	chg (f)	$\sqrt{5}*d$	error	
1	5.036067977	2	3.036068	2.23606798	0.8	
0.1	2.231606798	2	0.231607	0.2236068	0.008	
0.01	2.02244068	2	0.022441	0.02236068	8E-05	
non-optimal direction						
	$f[1+(1/\sqrt{2})d,$					
d	$1+(1/\sqrt{2})d]$	$f(1,1)$	chg (f)	See! It's true!		
1	4.621320344	2	2.62132	<	3.036068	
0.1	2.217132034	2	0.217132	<	0.231607	
0.01	2.021263203	2	0.021263	<	0.022441	

Gradient Descent

- To minimize the loss function

$$R(\boldsymbol{\beta}) = \sum_{i=1}^n R_i(\boldsymbol{\beta}) = \sum_{i=1}^n (y_i - \hat{y}_i(\boldsymbol{\beta}))^2$$

with learning rate $\eta > 0$, updating involves all observations

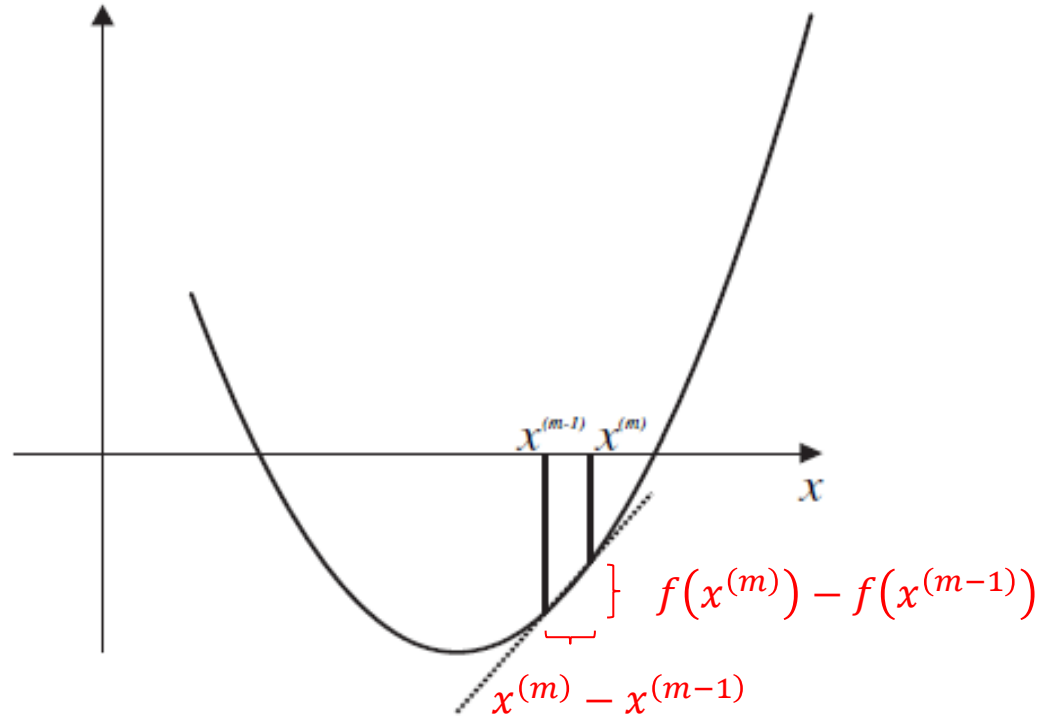
$$\begin{aligned}\boldsymbol{\beta}^{(r+1)} &= \boldsymbol{\beta}^{(r)} - \eta \nabla R(\boldsymbol{\beta}^{(r)}) \\ &= \boldsymbol{\beta}^{(r)} - \eta \sum_{i=1}^n \nabla R_i(\boldsymbol{\beta}^{(r)})\end{aligned}$$

Stochastic Gradient Descent

- SGD is a stochastic approximation of GD.
- SGD uses randomly selected samples/subset from the training set for each iteration
- At extreme, updating would involve only a single (randomly selected) observation

$$\boldsymbol{\beta}^{(r+1)} = \boldsymbol{\beta}^{(r)} - \eta \nabla R_i(\boldsymbol{\beta}^{(r)})$$

Find a root of a function



Root: $f(x^{(m)}) = 0$ for some m .

$$f'(x^{(m-1)}) = \frac{f(x^{(m)}) - f(x^{(m-1)})}{x^{(m)} - x^{(m-1)}}$$

$$x^{(m)} = x^{(m-1)} - \frac{f(x^{(m-1)})}{f'(x^{(m-1)})}$$

Newton Raphson

- $f: \mathcal{R} \rightarrow \mathcal{R}$

$$x^{(r+1)} = x^{(r)} - \frac{f(x^{(r)})}{f'(x^{(r)})}$$

- $f: \mathcal{R}^p \rightarrow \mathcal{R}^q$

$$\mathbf{x}^{(r+1)} = \mathbf{x}^{(r)} - [J(\mathbf{x}^{(r)})]^{-1} f(\mathbf{x}^{(r)})$$

$$J = \frac{\partial f}{\partial \mathbf{x}}, \quad J_{ij} = \frac{\partial f_i}{\partial x_j}$$

Minimizing Squared Loss

- To minimize the loss function

$$R(\boldsymbol{\beta}) = \sum_{i=1}^n R_i(\boldsymbol{\beta}) = \sum_{i=1}^n (y_i - \hat{y}_i(\boldsymbol{\beta}))^2$$

- Updating

$$\boldsymbol{\beta}^{(r+1)} = \boldsymbol{\beta}^{(r)} - [J(\boldsymbol{\beta}^{(r)})]^{-1} \nabla R(\boldsymbol{\beta}^{(r)})$$

$$\boldsymbol{\beta}^{(r+1)} = \boldsymbol{\beta}^{(r)} - [H(\boldsymbol{\beta}^{(r)})]^{-1} \left(\frac{\partial R(\boldsymbol{\beta}^{(r)})}{\partial \mathbf{x}} \right)$$

$$J = \frac{\partial}{\partial \boldsymbol{\beta}} \left(\frac{\partial R}{\partial \boldsymbol{\beta}^T} \right), \quad J_{ij} = \frac{\partial}{\partial x_i} \left(\frac{\partial \sum_k R_k}{\partial x_j} \right)$$

$$H = \frac{\partial^2 R}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T}, \quad H_{ij} = \frac{\partial^2 \sum_k R_k}{\partial x_i \partial x_j}$$

Example

- To minimize the loss function

$$R(\boldsymbol{\beta}) = \sum_{i=1}^n R_i(\boldsymbol{\beta}) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \beta_2 x_{i2})^2$$

- $\nabla R(\boldsymbol{\beta}) = \begin{pmatrix} \sum_i e_i \\ \sum_i e_i x_{i1} \\ \sum_i e_i x_{i2} \end{pmatrix}$

- $J(\boldsymbol{\beta}) = \begin{pmatrix} X_0^T X_0 & X_0^T X_1 & X_0^T X_2 \\ X_1^T X_0 & X_1^T X_1 & X_1^T X_2 \\ X_2^T X_0 & X_2^T X_1 & X_2^T X_2 \end{pmatrix}$

Example

- To maximize a log-likelihood function (or minimize $-l^*$)

$$l(\boldsymbol{\beta}; \mathbf{y}) = \sum_{i=1}^n l_i(\boldsymbol{\beta}; y_i)$$

That was



to be continued