

Agenda

1. Linear model w/ Gaussian/normal error terms
2. OLS as maximum likelihood
3. Inference under normality (of error terms)

Last time: normal eq's

$$X^T X \hat{\beta} = X^T y$$

$$\text{OLS estimator } \hat{\beta} = (X^T X)^{-1} X^T y$$

1. Linear model w/ Gaussian/normal error terms

Assume (X, y) satisfies

$$\boxed{y|X} \sim N(X\beta_*, \sigma^2 I_n) \text{ for some } \beta_* \in \mathbb{R}^p$$

random vector in \mathbb{R}^n

$X = \begin{pmatrix} -x_1^T- \\ \vdots \\ -x_n^T- \end{pmatrix}$ → vector of covariates for sample 1 each row has p components

$y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$

- Ex:
1. generate X from some dist. on covariates
 2. compute: $m \leftarrow X\beta_*$ for some vector $\beta_* \in \mathbb{R}^p$
 3. generate y from $N(m, \sigma^2 I_n)$
concretely: $y_i \stackrel{\text{ind}}{\sim} N(m_i, \sigma^2)$

Equivalent statement of Assumption:

1. (linearity) $E[y|X] = X\beta_*$

$$\begin{cases} \sigma^2 I_n: \\ \sigma^2 \begin{bmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{bmatrix} \\ \cdot \begin{bmatrix} \sigma^2 & & 0 \\ & \ddots & \\ 0 & & \sigma^2 \end{bmatrix} \end{cases}$$

2 (Gaussianity/normality) Define error terms: $\epsilon \leftarrow y - \underbrace{E[y|X]}_{\text{by linearity}} = y - X\beta$

$\epsilon | X \sim N(0, \sigma^2 1_n)$

Aside (max. likelihood):

Data: observation of weight of chalk $x \in \mathbb{R}$

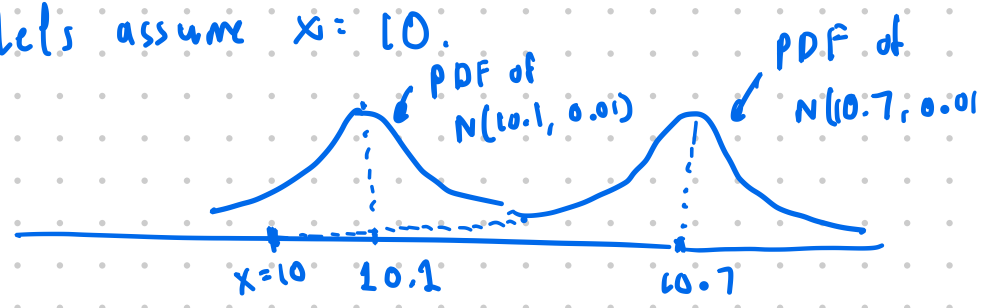
Assume $x \sim N(\underbrace{\mu}_{\text{unknown}}, 0.01)$

PDF of x is $p(x; \mu) = \frac{1}{\sqrt{2\pi} \cdot 0.1} \exp\left(-\frac{(x-\mu)^2}{2 \cdot 0.01}\right)$

likelihood: $L(\mu) \equiv p(x; \mu)$

maximum likelihood: $\hat{\mu}_{MLE} \leftarrow \arg\max_{\mu} L(\mu)$

Let's assume $x = 10$.



OLS as Max Likelihood under normality

likelihood of (X, y)

$L(\beta, \sigma) \equiv p(\underbrace{y|X}_{\text{random vector } y|X}; \underbrace{\beta, \sigma}_{\text{parameters of dist of } y|X}) = \prod_{i=1}^n p(y_i | X_i; \beta, \sigma)$

$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2} (y_i - x_i^T \beta)^2\right\}$

$\sum_{i=1}^n \epsilon_i = \epsilon_1 + \dots + \epsilon_n$

$\prod_{i=1}^n \epsilon_i = \epsilon_1 \cdot \epsilon_2 \cdot \dots \cdot \epsilon_n$

PDF of a $N(x^T \beta, \sigma^2)$ random var

$$(\hat{\beta}, \hat{\sigma}) = \arg \max_{\beta, \sigma} L(\beta, \sigma)$$

$$= \arg \max_{\beta, \sigma} \prod_{i=1}^n p(y_i | X_i; \beta, \sigma)$$

$$= \arg \max_{\beta, \sigma} \log \left(\prod_{i=1}^n p(y_i | X_i; \beta, \sigma) \right)$$

$$= \arg \max_{\beta, \sigma} \sum_{i=1}^n \log p(y_i | X_i; \beta, \sigma)$$

$$= \arg \max_{\beta, \sigma} \sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) - \frac{1}{2\sigma^2} (y_i - x_i^T \beta)^2$$

plug in

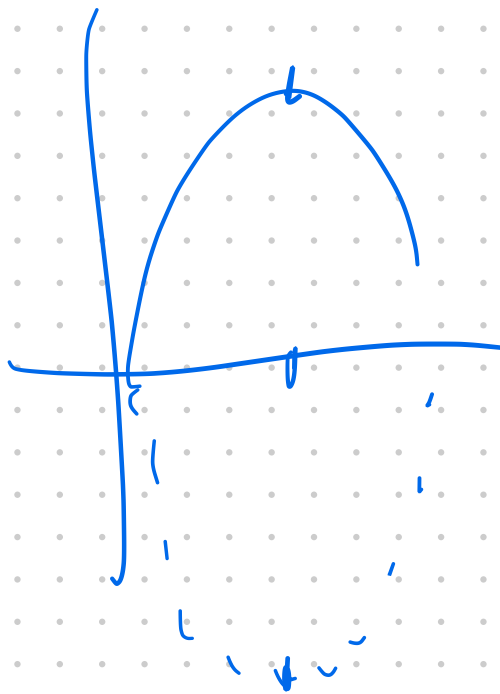
Focus on solving for $\hat{\beta}$

The argmax in β is

$$\arg \max_{\beta} \sum_{i=1}^n -\frac{1}{2\sigma^2} (y_i - x_i^T \beta)^2$$

$$= \arg \max_{\beta} \sum_{i=1}^n -\frac{1}{2} (y_i - x_i^T \beta)^2$$

$$= \arg \min_{\beta} \underbrace{\sum_{i=1}^n \frac{1}{2} (y_i - x_i^T \beta)^2}_{\text{SSR}(\beta)}$$



Dist. of $\hat{\beta}$ under normality:

$$\hat{\beta} = (X^T X)^{-1} X^T y \quad (\text{from normal eq})$$

$$= (X^T X)^{-1} X^T (X \beta_0 + \epsilon) \quad (\text{from linearity})$$

$$= \cancel{(X^T X)^{-1}} \cancel{X^T X} \beta_0 + (X^T X)^{-1} X^T \epsilon$$

$$= \beta_0 + (X^T X)^{-1} X^T \epsilon$$

$$\hat{\beta} - \beta_0 = (X^T X)^{-1} X^T \epsilon$$

If $z \sim N(\mu, \Sigma)$ then

$$A z \sim N(A\mu, A \Sigma A^T)$$

$$\sim N\left(\cancel{(X^T X)^{-1} X^T} \cdot 0, \sigma^2 \cancel{(X^T X)^{-1} X^T} 1_n X \cancel{(X^T X)^{-1}}\right)$$

$$= N\left(0, \sigma^2 (X^T X)^{-1}\right)$$