

Section 3D. Coefficient Uncertainty

Statistics for Data Science

Victor M. Preciado, PhD MIT EECS
Dept of Electrical & Systems Engineering
University of Pennsylvania
preciado@seas.upenn.edu

Coefficients Accuracy

- Assume that we observe data generated by an additive linear model of the form

$$\mathbf{x}_i \sim f_X \text{ and } \mathbf{y} = M_X \beta + \varepsilon, \text{ where } \varepsilon_i \sim f_\varepsilon$$

- We can estimate the linear coefficients as $\hat{\beta} = (M_X^\top M_X)^{-1} M_X^\top \mathbf{y}$. Note that, since the datapoints (\mathbf{x}_i, y_i) are random, *the coefficients $\hat{\beta}_i$ are random variables themselves!*
- **Mean analysis:** One can prove that $\mathbb{E} [\hat{\beta}_i] = \beta_i$ for all i (unbiased estimator)
- **Covariance analysis:** Assuming that we are given a data matrix M_X , the covariance matrix of $\hat{\beta}$ satisfies (without proof)

$$\text{Cov} [\hat{\beta} | M_X] = \sigma^2 (M_X^\top M_X)^{-1}$$

Coefficients Accuracy: Univariate Case

- ▶ For the particular case $p = 1$, the linear regression takes the form $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$
- ▶ The diagonal elements of the covariance matrix are the variances of the estimated coefficients. These variances are:

$$\begin{aligned} \text{SD}(\hat{\beta}_1)^2 &= \text{Var}[\hat{\beta}_1 | \{x_i\}_{i=1}^N] = \frac{\sigma^2}{\sum_{i=1}^N (x_i - \bar{x})^2} \\ \text{SD}(\hat{\beta}_0)^2 &= \text{Var}[\hat{\beta}_0 | \{x_i\}_{i=1}^N] = \sigma^2 \left(\frac{1}{N} + \frac{\bar{x}^2}{\sum_{i=1}^N (x_i - \bar{x})^2} \right) \end{aligned}$$

with $\sigma^2 = \text{Var}(\varepsilon)$, $\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$ and $\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$. In this course, we use $\text{SD}(\cdot)$ to denote the Standard Deviation of a r.v.

Coefficients Accuracy: Numerical Example

Numerical validation: Let us consider an unrealistic, but illustrative situation

- ▶ Consider 1,000 independent realizations of a training dataset, $\mathcal{D}_{\text{Tr}}^k$ for $k = 1, 2, \dots, 1000$. Notice that, in practice, we will only have access to a single realization of a training dataset!
- ▶ Each dataset contains $N = 100$ sample points $\mathcal{D}_{\text{Tr}}^k = \{(\mathbf{x}_1^k, y_1^k), \dots, (\mathbf{x}_{100}^k, y_{100}^k)\}$
- ▶ For each $\mathcal{D}_{\text{Tr}}^k$, we compute the corresponding estimates $\hat{\beta}_0^k$ and $\hat{\beta}_1^k$ for $k = 1, 2, \dots, 1000$

Coefficients Accuracy: Numerical Example

In this example, $\beta_1 = \mathbb{E} [\hat{\beta}_1] \neq 0.5$ and $\text{Var} [\hat{\beta}_1 | \{x_i\}_{i=1}^N] \neq 0.05^2$.

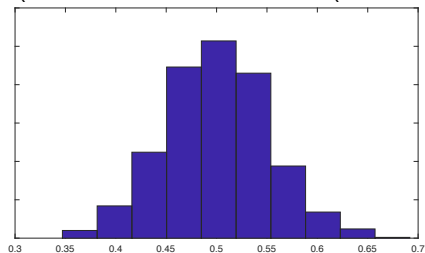
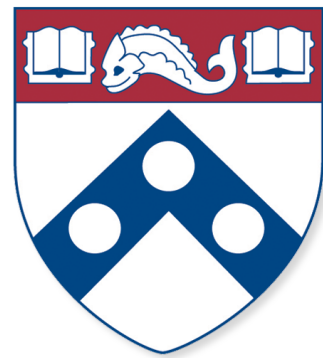


Figure: Histogram of the values of $\hat{\beta}_1^k$. We have used a linear model $Y = 1 + 0.5X + \varepsilon$, where $X \sim \mathcal{N}(0, 1)$ and $\text{Var}(\varepsilon) = \sigma^2 = 1$.



Penn Engineering

Copyright 2020 University of Pennsylvania
No reproduction or distribution without permission.