



STA2GO-lecture 4 July 14, 2025

Bias - Variance trade-off

Criteria for goodness of estimators

- 1) Principle of unbiasness $E[\hat{\theta}] = \theta$
- 2) Consistency (convergence in prob.) $\hat{\theta} \xrightarrow{P} \theta, \lim_{n \rightarrow \infty} \text{Var}(\hat{\theta}) = 0$
- 3) Principle of Minimum Variance: an estimator is preferable with low variance
- 4) Known Sampling dist.: for most distributions, we have approximate distributions

Def. $\hat{\theta}$ is unbiased estimator of θ if $E[\hat{\theta}] = \theta$.
 Otherwise we call it biased.

The bias of $\hat{\theta}$ is $B(\hat{\theta}) = E[\hat{\theta}] - \theta$

In reality Sampling space can't go to $n \rightarrow \infty$,

Def. Mean Square Error

$$MSE(\hat{\theta}) = E[(\hat{\theta} - \theta)^2] \xrightarrow{\substack{\text{used in machine} \\ \text{learning}}}$$

$$\text{Bias-Variance} \quad \text{Trade-off (decomposition)}$$

$$\text{MSB}(\hat{\theta}) = \text{Var}(\hat{\theta}) + \underbrace{B(\hat{\theta})}_{\geq 0} + \underbrace{V(\hat{\theta})}_{\geq 0}$$

$$\begin{aligned} \text{Prove: } \text{MSB}(\hat{\theta}) &= E[(\hat{\theta} - \theta)^2] \\ &= E[(\hat{\theta} - E[\hat{\theta}]) + (E[\hat{\theta}] - \theta)]^2 \\ &= E[(\theta - E[\theta])^2 + (E[\hat{\theta}] - \theta) + 2E[(E[\hat{\theta}] - \theta)(\hat{\theta} - E[\hat{\theta}])]_0] \\ &= \text{Var}(\hat{\theta}) + B(\hat{\theta}) + 2B(\theta) E[(\theta - E[\theta])]_0 \\ &= \text{Var}(\hat{\theta}) + B(\hat{\theta}) + 2B(\theta) \underbrace{E[\hat{\theta}] - E[E[\hat{\theta}]]}_0 \end{aligned}$$

Can always convert biased est. into an unbiased one

$$E[\hat{\theta}] = a\theta + b, \quad a, b > 0$$

What is $B(\hat{\theta})$? $B(\hat{\theta}) = E[\hat{\theta}] - \theta = a\theta + b - \theta = (a-1)\theta + b$

$$\Rightarrow B(\hat{\theta}) \neq 0 \Rightarrow \hat{\theta} \text{ is biased}$$

- Find a function of $\hat{\theta}$ which is unbiased estimator of θ

$$E[\hat{\theta}] = a\theta + b \Rightarrow E[\hat{\theta}] - b = a\theta \Rightarrow \frac{1}{a}E[\hat{\theta} - b] = \theta$$

$$\hat{\theta} = \frac{\hat{\theta} - b}{a} \quad \text{unbiased for } \Rightarrow E\left[\frac{\hat{\theta} - b}{a}\right] = \theta$$

Suppose we have two independent estimators

$$E[\hat{\theta}_1] = \theta \quad E[\hat{\theta}_2] = \theta$$

$$\text{Var}(\hat{\theta}_1) = \sigma_1^2 \quad \text{Var}(\hat{\theta}_2) = \sigma_2^2$$

linear combination $\hat{\theta} = w\hat{\theta}_1 + (1-w)\hat{\theta}_2, \quad 0 < w < 1$

$$w_1x_1 + w_2x_2$$

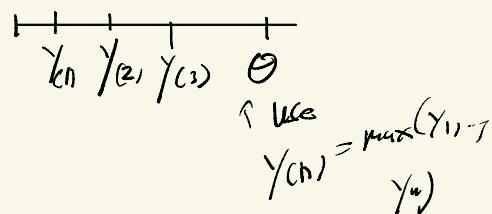
Convex combination

$$w_1x_1 + w_2x_2 \Rightarrow w_2 = 1 - w_1$$

$$0 < w_1, w_2 < 1$$

$$Y_1, Y_2, \dots, Y_n \stackrel{iid}{\sim} \text{Unif}(0, \Theta)$$

Convex combination of estimators gives lower var. unbiased result.



$$\mathbb{E}[\hat{\theta}] = w\mathbb{E}[\hat{\theta}_1] + (1-w)\mathbb{E}[\hat{\theta}_2]$$

to estimate Θ

$$+ \Theta - w\Theta = \Theta \Rightarrow$$

$$\mathbb{E}[\hat{\theta}] = \Theta$$

$\Rightarrow \hat{\theta}$ is unbiased estimator

$\bar{Y}_{(1)}$

\bar{Y}

Find the optimal value of w s.t. variance of $\hat{\theta}$ is minimized

$$\begin{aligned} \text{var}(\hat{\theta}) &= \text{var}(w\hat{\theta}_1) + \text{var}((1-w)\hat{\theta}_2) \\ &= w^2 \text{var}(\hat{\theta}_1) + (1-w)^2 \text{var}(\hat{\theta}_2) \\ &= w^2 \sigma_1^2 + (1-w)^2 \sigma_2^2 \end{aligned}$$

$$g(w) = w^2 \sigma_1^2 + (1-w)^2 \sigma_2^2$$

$$\frac{d}{dw} g(w) = 2w\sigma_1^2 - (1-w)\sigma_2^2 = 0$$

$$\Rightarrow 2w\sigma_1^2 - 2(1-w)\sigma_2^2 = 0$$

$$\Rightarrow w\sigma_1^2 + \sigma_2^2 + w\sigma_2^2 = 0 \Rightarrow$$

$$w^* = \frac{\sigma_2^2}{\sigma_1^2 + \sigma_2^2}$$

$$\frac{d^2 g(w)}{dw^2} = \Gamma_1^2 + \Gamma_2^2 > 0 \Rightarrow w^* \text{ is global minimum}$$

ex: $\hat{\theta} = 0.8\hat{\theta}_1 + 0.2\hat{\theta}_2$

Note: All machine learning algorithms is an optimization problem

One more example:
 $y_1, y_2, \dots, y_n \stackrel{iid}{\sim} \text{Exp}(\theta)$ $f(x) = \frac{1}{\theta} e^{-x/\theta}, x > 0$
 $F(x) = 1 - e^{-x/\theta}$

Show that $\hat{\theta} = \bar{y}_{(1)}$ is unbiased for θ and find
 Use dist. function technique to find sampling dist of
 $\hat{\theta}$

$$\begin{aligned} F_{Y_{(1)}}(y) &= P(Y_{(1)} < y) = P(\min(Y_1, \dots, Y_n) \leq y) \\ &= 1 - P[\min(Y_1, \dots, Y_n) \geq y] \\ &= 1 - P[Y_1 > y, Y_2 > y, \dots, Y_n > y] \end{aligned}$$

$$= 1 - P[Y_1 > y] P[Y_2 > y] \dots P[Y_n > y]$$

$$= 1 - (P[Y_1 > y])^n = 1 - (1 - F_y(y))^n$$

$$\frac{d}{dy} F_{Y_{(1)}}(y) = n(1 - F(y))^{n-1} f(y)$$

Using dist: $f_{Y_{(1)}}(y) = n \left(1 - (1 - e^{-\frac{y}{\theta}})\right)^{n-1} \cdot \frac{1}{\theta} e^{-\frac{y}{\theta}}$

$$= n e^{-\frac{(n-1)y}{\theta}} \cdot \frac{1}{\theta} e^{-\frac{y}{\theta}}$$

$$= \frac{n}{\theta} e^{-\frac{ny}{\theta}}$$

$$f_{Y_{(1)}}(y) = \frac{n}{\theta} e^{-\left(\frac{n}{\theta}\right)y} \Rightarrow Y_1 \sim \text{Exp}\left(\frac{\theta}{n}\right)$$

$$E[Y_{(1)}] = \frac{\theta}{n} \Rightarrow E[nY_{(1)}] = n \left(\frac{\theta}{n}\right) = \theta$$

$$\Rightarrow \hat{\theta} = \frac{nY_{(1)}}{n} \text{ is unbiased}$$

$$\text{MSE}(\hat{\theta}) = \text{Var}\left(\frac{y_{(1)}}{n}\right) + \left[\cancel{B\left(\frac{y_{(1)}}{n}\right)^T}\right]^2$$

$$= n^2 \text{Var}(y_{(1)}) = n^2 \left(\frac{\theta}{n}\right)^2 = \theta^2$$

$\therefore \text{MSE}(\hat{\theta}) = \theta^2 \rightarrow \text{Variance is Constant}$

$$\mathbb{E}[\bar{y}] = \theta \quad \text{while } \text{Sel} \quad \text{MSE}(\hat{\theta}) = \text{Var}(\bar{y})$$

$$= \frac{\text{Var}(y_{(1)})}{n} = \frac{\theta^2}{n}$$

$$\therefore \text{MSE}(\hat{\theta}) = \frac{\theta^2}{n}$$

After break:

$$Y_1, Y_2, \dots, Y_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$$

$$S_1^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

$$S_2^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

- a) Find bias of each estimator
- b) Find MSE of each estimator
- c) Compare them

$$\frac{(n-1)s_1^2}{\sigma^2} \sim \chi^2_{(n-1)} \Rightarrow E\left[\frac{(n-1)s_1^2}{\sigma^2}\right] = n-1$$

$$\Rightarrow E[s_1^2] = \sigma^2$$

$$B(s_1^2) = E[s_1^2] - \sigma^2 = \sigma^2 - \sigma^2 = 0 \Rightarrow B(s_1^2) = 0$$

$$s_2^2 = \left(\frac{n-1}{n}\right) s_1^2 \Rightarrow E[s_2^2] = \left(\frac{n-1}{n}\right) E[s_1^2]$$

$$= \left(\frac{n-1}{n}\right) \sigma^2$$

$$B(s_2^2) = E[s_2^2] - \sigma^2 = \left(\frac{n-1}{n}\right) \sigma^2 - \sigma^2 = \frac{1}{n} \sigma^2$$

$$\Rightarrow B(s_2^2) = -\frac{1}{n} \sigma^2$$

1st one is better, according to principle of unbiasedness

s_1^2 is unbiased, whereas s_2^2 is biased.

However, the bias of s_2^2 vanishes as $n \rightarrow \infty$

$$\frac{(n-1)s_1^2}{\sigma^2} \sim \chi^2_{(n-1)} \Rightarrow \text{Var}\left(\frac{(n-1)s_1^2}{\sigma^2}\right) = 2(n-1)$$

$$\begin{aligned} \text{MSE}(s_1^2) &= \left[\text{Var}(s_1^2) + \left(\mathbb{E}(s_1^2) - s_1^2 \right)^2 \right] \\ &= \left(\frac{n-1}{\sigma^2} \right) \text{Var}(s_1^2) = 2(n-1) \\ &\Rightarrow \text{Var}(s_1^2) = \left(\frac{2}{n-1} \right) \sigma^4 \end{aligned}$$

$$S_2^2 = \left(\frac{n-1}{n}\right) S_1^2 \Rightarrow \text{var}(S_2^2) = \left(\frac{n-1}{n}\right)^2 \text{var}(S_1^2) =$$

$$\left(\frac{n-1}{n}\right)^2 \cdot \left(\frac{2}{n-1}\right) \sigma^4 = \left(\frac{2(n-1)}{n^2}\right) \sigma^4$$

$$MSE(S_2^2) = \left(\frac{2(n-1)}{n^2}\right) \sigma^4 + \left(\frac{1}{n}\right)^2 \sigma^4 = \left(\frac{2n-1}{n^2}\right) \sigma^4$$

$$\therefore MSE(S_1^2) = \left(\frac{2}{n-1}\right) \sigma^4$$

$$MSE(S_2^2) = \left(\frac{2n-1}{n^2}\right) \sigma^4$$

We have $B(S_1^2) > 0$ and $B(S_2^2) \neq 0$, so
the first estimator is preferable according to principle of
unbiasedness

$n > 1 \Rightarrow \left(\frac{2}{n-1}\right) S_1^4 > \left(\frac{2n-1}{n^2}\right) S_2^4 \Rightarrow \text{MSE}(S_1^2) > \text{MSE}(S_2^2)$

So, S_2^2 is preferable according to principle of minimax

Variance whichever one you use depends on sample size n .

$$n \rightarrow \infty \Rightarrow \left(\frac{n}{n-1}\right) \approx \left(\frac{2n-1}{n^2}\right) \Rightarrow \text{MSE}(S_1^2) \approx \text{MSE}(S_2^2)$$

So S_1^2 is preferable

Y ~ Poisson(θ). Show that $(-1)^y$ is unbiased estimator for $e^{-2\theta}$

$$\text{Compute } E[(-1)^y] = \sum_{i=1}^{\infty} (-1)^y \frac{\theta^y e^{-\theta}}{y!} = e^{-\theta} \sum_{i=0}^{\infty} \frac{(-\theta)^y}{y!} = e^{-\theta} \cdot e^{-\theta} = e^{-2\theta}$$

$(-1)^y$ is unbiased

Q1. Why isn't this a good estimator, $(-1)^y \rightarrow -1, 1, -1, 1, \dots$
b/c $\text{O}(e) \leq 1$

Q2. Are we violating principles of unbiasedness? Yes.
Principles are not theorems to hold true in
all cases, and seldom do they yield satisfactory
result in all cases.

Standard Error

The standard deviation of the sampling
dist. of $\hat{\theta}$ is called standard error and is
denoted by $S_{\hat{\theta}} = \sqrt{\text{Var}(\hat{\theta})}$

$$\text{Var}(\bar{Y}) = \frac{\sigma^2}{n} \Rightarrow \sigma_{\bar{Y}} = \frac{\sigma}{\sqrt{n}}$$

$$Y \sim \text{Binomial}(n, p) \Rightarrow \text{Var}(\hat{p}) = \text{Var}\left(\frac{Y}{n}\right)$$

$$= \frac{1}{n^2} \text{Var}(Y) = \frac{np(1-p)}{n^2} = \frac{p(1-p)}{n}$$

Error of Estimation:

$$E = |\hat{\theta} - \theta|$$

$$\therefore \sigma_E = \sqrt{\frac{p(1-p)}{n}}$$

$$Y_1, \dots, Y_n \stackrel{iid}{\sim} \text{Exp}(\theta)$$

Suggest an unbiased estimator for θ and find its standard error.

unbiased estimator

$$E[\bar{Y}] = \theta \quad \text{unbiased estimator}$$

$$\text{Var}(\bar{Y}) = \frac{\text{Var}(Y)}{n} = \frac{\theta^2}{n} \Rightarrow \sigma_{\bar{Y}} = \frac{\theta}{\sqrt{n}} \approx \frac{\hat{\theta}}{\sqrt{n}}$$

plug-in
principle

$$\Rightarrow \sigma_{\bar{Y}} = \frac{\hat{\theta}}{\sqrt{n}}$$

An engineer observes $n=10$ independent length-of-life measurements for an electronic component. The coverage of these observations is 1020 hours. If these observations come from exponential dist., with mean \bar{Y} then estimate θ and place a 2-standard-error bound on the error of estimation.

$$\text{2.6) } \theta = \bar{Y} = 1020 \Rightarrow 2\text{-standard-error} = (2)\sqrt{\frac{\bar{Y}}{n}}$$

$$\approx 2 \cdot \frac{\bar{Y}}{\sqrt{n}}$$

Last example:
The # of grain nucleation sites per unit volume
is modeled as poisson dist. with mean λ

50 samples are collected with average 20
sites per unit volume. Estimate
mean λ and place a 2-standard-error bound
on the error of estimation.

$$= 2 \left(\frac{1020}{\sqrt{50}} \right)$$

$$= 645.1$$

Let $y_i \in \mathbb{R}^n$ denoting # of grain nucleation sites
 $y_i \sim \text{Poisson}(\lambda)$, $i=1, \dots, n$ $E[y_i] = \lambda$
 $\text{Var}(y_i) = \lambda$

$$E[\bar{y}] = \frac{1}{50} \sum_{i=1}^{50} E[y_i] = \frac{1}{50} \sum_{i=1}^{50} \lambda = \lambda \quad \text{unbiased estimator}$$

$$\text{Var}(\bar{y}) = \frac{\text{Var}(y_1)}{n} = \frac{\lambda}{50} \Rightarrow S_{\bar{y}} = \sqrt{\frac{\lambda}{50}} \approx \sqrt{\frac{\lambda}{50}}$$

$$= \sqrt{\frac{8}{50}} = 0.6324$$

bound on error:

$$20 \pm 2 \times 0.6324 = 20 \pm 1.265 = (18.7, 21.2)$$