# STAT 88: Lecture 33

**Contents**

Section 10.4: Normal Distribution
Section 11.1: Bias and Variance

Warm up:

*population dist'n.*
$$E(X_i) = \mu, \ Var(X_i) = \sigma^2$$

(a) If $X_1, \ldots, X_n \overset{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$, what distribution is $\bar{X}$?

(b) If $X_1, \ldots, X_n \overset{\text{iid}}{\sim} \mathcal{N}(\mu_X, \sigma_X^2)$ and $Y_1, \ldots, Y_m \overset{\text{iid}}{\sim} \mathcal{N}(\mu_Y, \sigma_Y^2)$ and two samples are independent, what distribution is $\bar{X} - \bar{Y}$?

*pop. distrib.  $E(X_i) = p, \ Var(X_i) = p(1-p)$*

(c) If $X_1, \ldots, X_n \overset{\text{iid}}{\sim}$ Bernoulli$(p)$, (approximately) what distribution is $\bar{X}$?

(d) If $X_1, \ldots, X_n \overset{\text{iid}}{\sim}$ Bernoulli$(p_X)$ and $Y_1, \ldots, Y_m \overset{\text{iid}}{\sim}$ Bernoulli$(p_Y)$ and two samples are independent, (approximately) what distribution is $\bar{X} - \bar{Y}$?

(a) $E(\bar{X}) = \mu, \ Var(\bar{X}) = \frac{\sigma^2}{n}$.

$\qquad \bar{X} = \frac{1}{n}(X_1 + \cdots + X_n) \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$ (Exact distribution)

(b) $\bar{X} \sim \mathcal{N}\left(\mu_X, \frac{\sigma_X^2}{n}\right). \quad \bar{Y} \sim \mathcal{N}\left(\mu_Y, \frac{\sigma_Y^2}{m}\right)$

$\qquad\qquad\qquad \underbrace{\phantom{XXXX}}_{\text{indep.}}$

$\qquad \Rightarrow \bar{X} - \bar{Y} \sim \mathcal{N}\left(\mu_X - \mu_Y, \frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}\right)$

(c) $E(\bar{X}) = p, \ Var(\bar{X}) = \frac{p(1-p)}{n}$.

$\qquad \bar{X} \sim \mathcal{N}\left(p, \frac{p(1-p)}{n}\right)$ (by CLT)

(d) $\bar{X} \sim \mathcal{N}\left(p_X, \frac{p_X(1-p_X)}{n}\right), \quad \bar{Y} \sim \mathcal{N}\left(p_Y, \frac{p_Y(1-p_Y)}{m}\right)$

$\qquad \Rightarrow \bar{X} - \bar{Y} \sim \mathcal{N}\left(p_X - p_Y, \frac{p_X(1-p_X)}{n} + \frac{p_Y(1-p_Y)}{m}\right)$
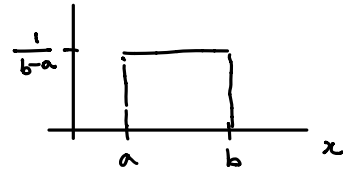
**Last time**

Continuous probability distributions:

Uniform

Let $X \sim \text{Unif}(a, b)$. Then the density is

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{if } a < x < b \\ 0 & \text{otherwise} \end{cases}$$
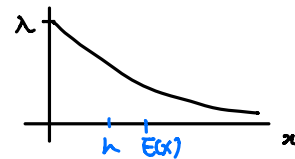
We have $E(X) = (a + b)/2$ and $\text{SD}(X) = (b - a)/\sqrt{12}$.

Exponential

Let $X \sim \text{Exp}(\lambda)$. Then the density is

E.g. X= time until failure
of a mechanical device

$$f(x) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

We have $E(X) = \frac{1}{\lambda}$ and $\text{SD}(X) = \frac{1}{\lambda}$. The Half Life is $h = \log 2/\lambda$.

Median

Normal

Let $X \sim \mathcal{N}(\mu, \sigma^2)$. Then the density is

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \text{ for } -\infty < x < \infty.$$

We have $E(X) = \mu$ and $\text{SD}(X) = \sigma$. If $X \sim \mathcal{N}(\mu_X, \sigma_X^2)$ and $Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$ and $X$ and $Y$ are independent, then

$$X + Y \sim \mathcal{N}(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2).$$

This result extends to linear combinations of independent normal random variables.

## 10.4. The Normal Distribution

**Confidence Interval for the Difference Between Means**   Suppose you have two independent samples as follows:

- $X_1, X_2, \ldots, X_n$ are i.i.d. with mean $\mu_X$ and SD $\sigma_X$.

- $Y_1, Y_2, \ldots, Y_m$ are i.i.d. with mean $\mu_Y$ and SD $\sigma_Y$.

You want to estimate the difference $\mu_X - \mu_Y$. Then $\bar{X} - \bar{Y}$ is an unbiased estimator for $\mu_X - \mu_Y$.

By CLT, we know

- $\bar{X}$ is approximately $\mathcal{N}(\mu_X, \frac{\sigma_X^2}{n})$.

- $\bar{Y}$ is approximately $\mathcal{N}(\mu_Y, \frac{\sigma_Y^2}{m})$.

$W \sim \mathcal{N}(\mu_W, \sigma_W^2)$

$P(W - 2 \cdot \sigma_W \le \mu_W \le W + 2 \cdot \sigma_W) = 95\%$

Then $\bar{X} - \bar{Y}$ is approximately $\mathcal{N}(\mu_X - \mu_Y, \frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m})$ and an approximate 95% CI for $\mu_X - \mu_Y$ is given by

$$\bar{X} - \bar{Y} \pm 2\sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}.$$

$W = \bar{X} - \bar{Y}$.
$\mu_W = \mu_X - \mu_Y$
$\sigma_W = \frac{\hat{\sigma_X^2}}{n} + \frac{\hat{\sigma_Y^2}}{m}$

Example: Suppose you have drawn samples of people independently from two cities, and suppose you have collected the following data:

$\mu_X$
- The incomes of the 400 sampled people in City X have an average of 70,000 dollars and an SD of 40,000 dollars. $n=400, \quad \bar{X} = 70000, \quad \sigma_X \approx 40000$

$\mu_Y$
- The incomes of the 600 sampled people in City Y have an average of 80,000 dollars and an SD of 50,000 dollars. $m=600, \quad \bar{Y} = 80000, \quad \sigma_Y = 50000$

Find a 95% CI for the difference between the mean incomes in the two citie

$$(70000 - 80000) \pm 2\sqrt{\frac{40000^2}{400} + \frac{50000^2}{600}} = (-15,916, \ -4,824)$$

**Test for the Equality of Two Means (A/B Test)** We wish to determine if two independent populations have the same mean, i.e. $\mu_X - \mu_Y = 0$.

Example: Suppose you have drawn samples of people independently from two cities, and suppose you have collected the following data:

- The incomes of the 400 sampled people in City X have an average of 70,000 dollars and an SD of 40,000 dollars.

- The incomes of the 600 sampled people in City Y have an average of 80,000 dollars and an SD of 50,000 dollars.

$H_0 : \mu_X = \mu_Y$, the mean income in City X is the same as the mean income in City Y.

$H_A : \mu_Y > \mu_X.$    $(\mu_Y - \mu_X > 0)$

$= 80,000 - 70,000$

Test statistic: $\bar{Y} - \bar{X}$. Our observed value is 10,000. We reject $H_0$ if $\bar{Y} - \bar{X}$ is large.

Under $H_0$, we have

$$\bar{Y} - \bar{X} \sim \mathcal{N}\left(\underline{0}, \underbrace{\frac{\sigma_X^2}{400}}_{= \frac{40,000^2}{}} + \underbrace{\frac{\sigma_Y^2}{600}}_{= 50,000^2}\right). \qquad \sim \mathcal{N}\left(\mu_Y - \mu_X, \frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}\right)$$

$$\sim \mathcal{N}(0, 2858^2) \qquad \underset{0 \text{ under } H_0}{\overset{\shortparallel}{}}$$

$p$-value:

$$P(\bar{Y} - \bar{X} \geq 10,000)$$

$$= P\left(\frac{\bar{Y} - \bar{X} - 0}{2858} \geq \frac{10,000 - 0}{2858}\right)$$

$$= P(Z \geq 3.5)$$

$$= 1 - \Phi(3.5)$$

**Confidence Interval for the Difference Between Proportions**   This is a special case of the above where now populations are 0's and 1's.

- $X_1, X_2, \ldots, X_n$ are i.i.d. from Bernoulli($p_X$);

- $Y_1, Y_2, \ldots, Y_m$ are i.i.d. from Bernoulli($p_Y$).   95% CI: $\bar{X} - \bar{Y} \pm 2\sqrt{\frac{P_X(1-P_X)}{n} + \frac{P_Y(1-P_Y)}{m}}$

$\nearrow$

$\bar{X} - \bar{Y}$ is an unbiased estimator for $p_X - p_Y$:

$$\bar{X} - \bar{Y} \sim \mathcal{N}\left(p_X - p_Y, \frac{p_X(1-p_X)}{n} + \frac{p_Y(1-p_Y)}{m}\right).$$   by CLT

Example: Suppose we have independent samples from two cities, where sample sizes are $n = 400$ and $m = 600$ for City X and City Y, and:

$\bar{X} = 0.37$
- 37% of the City X sample are undecided about who they want as President;

- 28% of the City Y sample are undecided about who they want as President.

$\bar{Y} = 0.28$

Find a 95% CI for $p_X - p_Y$.

$$\bar{X} - \bar{Y} \pm 2\sqrt{\frac{P_X(1-P_X)}{n} + \frac{P_Y(1-P_Y)}{m}}$$

$$\approx (0.37 - 0.28) \pm 2\sqrt{\frac{0.37(1-0.37)}{400} + \frac{0.28(1-0.28)}{600}}$$

$$= (0.029, 0.15)$$

**Test for the Equality of Two Proportions**    Our hypotheses are:

- $H_0 : p_X = p_Y = p$; here $p$ is just a name we are giving to the common value of $p_X$ and $p_Y$.

- $H_A : p_X > p_Y$.    $p_X - p_Y > 0$

Test statistic: $\bar{X} - \bar{Y}$. Under $H_0$,

$$\bar{X} - \bar{Y} \sim \mathcal{N}\left(0, \frac{p(1-p)}{n} + \frac{p(1-p)}{m}\right).$$

Example: Suppose we have independent samples from two cities, where sample sizes are $n = 400$ and $m = 600$ for City X and City Y, and:

- 37% of the City X sample are undecided about who they want as President;

- 28% of the City Y sample are undecided about who they want as President.

$$\frac{\hat{p}_X(1-\hat{p}_X)}{n} + \frac{\hat{p}_Y(1-\hat{p}_Y)}{m}$$

Test $H_0 : p_X = p_Y$ vs $H_A : p_X > p_Y$ at level 5%.

Test statistic: $\bar{X} - \bar{Y}$.

Under $H_0$, $\bar{X} - \bar{Y} \sim N\left(0, \frac{p(1-p)}{n} + \frac{p(1-p)}{m}\right)$

What is $p$?   City X and City Y have common $p$.

$\Rightarrow$ estimate $p$ by $\hat{p} = \frac{400}{1000}\cdot\bar{X} + \frac{600}{1000}\cdot\bar{Y} = 0.316$

$\left(\hat{p} = \frac{n}{n+m}\cdot\bar{X} + \frac{m}{n+m}\cdot\bar{Y}\right)$

$\Rightarrow \bar{X} - \bar{Y} \sim N\left(0, \frac{0.316(1-0.316)}{400} + \frac{0.316(1-0.316)}{600}\right)$

$\approx N(0, 0.03^2)$

$p\text{-val} = P(\bar{X} - \bar{Y} \geq 0.37 - 0.28)$

$= P(\bar{X} - \bar{Y} \geq 0.09)$

$= P\left(\frac{\bar{X} - \bar{Y} - 0}{0.03} \geq \frac{0.09 - 0}{0.03}\right)$

$= P(Z \geq 3)$

$= 1 - \Phi(3)$

Under $H_0$

$X_1, \cdots, X_n \sim \text{Ber}(p_X) = \text{Ber}(p)$

$Y_1, \cdots, Y_m \sim \text{Ber}(p_Y) = \text{Ber}(p)$

Under $H_0$

$\hat{p} = \frac{X_1 + \cdots + X_n + Y_1 + \cdots + Y_m}{n+m}$

$= \frac{X_1 + \cdots + X_n}{n+m} + \frac{Y_1 + \cdots + Y_m}{n+m}$

$= \frac{n}{n+m}\cdot\bar{X} + \frac{m}{n+m}\cdot\bar{Y}$
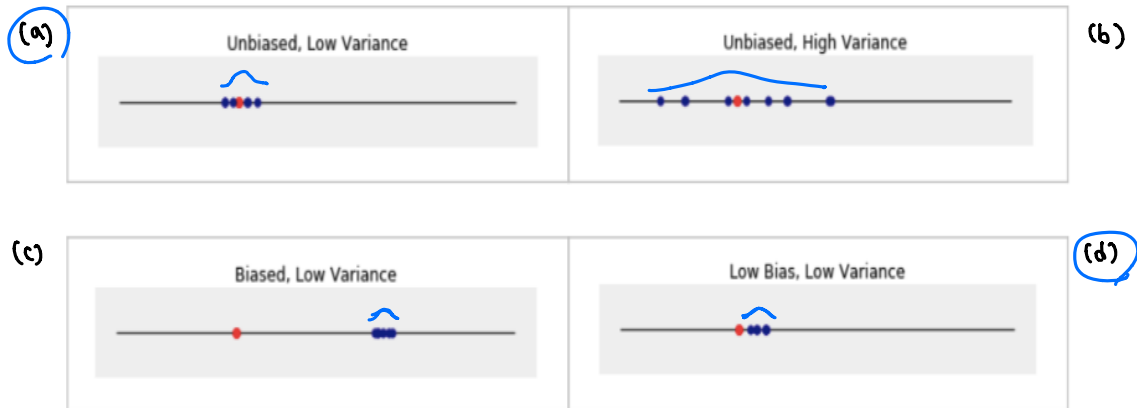
# 11.1. Bias and Variance

*E.g. $\mu$*

Suppose we are trying to estimate a constant numerical parameter, $\theta$, and our estimator is the statistic $T$. Below $\theta$ is red and $T$ is blue for different samples.

*E.g. $\bar{X}$*     *fixed*     *random*

What are the two best estimators?

(a) Unbiased, Low Variance    (b) Unbiased, High Variance

(c) Biased, Low Variance    (d) Low Bias, Low Variance

Lets make a quantitative analysis.

*E.g. $T = \bar{X}$, $\theta = \mu$*

**Mean Squared Error**   The error in our estimate is $T - \theta$. Then

$$\mathrm{MSE}_\theta(T) = E_\theta\left((T - \underline{\theta})^2\right).$$

We are using $\theta$ as a subscript to remind us that the expectation is calculated under the assumption that $\theta$ is the true value of the parameter.

Think of this as the average distance squared of $T$ from $\theta$. We want $\mathrm{MSE}_\theta(T)$ to be as small as possible.

**Decomposition of Error**

Deviation:
$$D_\theta(T) = T - E_\theta(T).$$

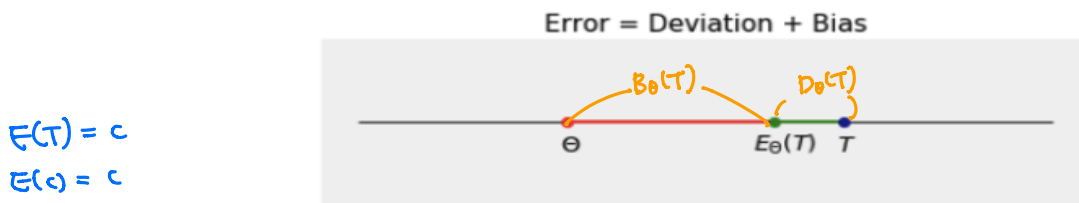( deviation of T from the mean )

Is it random or constant?

Bias:
$$B_\theta(T) = E_\theta(T) - \theta.$$

( Unbiased means $E_\theta(T) = \theta$ )

Is it random or constant?

We have a decomposition of the error as the sum of the deviation and the bias:
$$T - \theta = \underbrace{(T - E_\theta(T))}_{=D_\theta(T)} + \underbrace{(E_\theta(T) - \theta)}_{=B_\theta(T)}.$$



Error = Deviation + Bias

$E(T) = c$

$E(c) = c$

What is $E_\theta(D_\theta(T))$? What is $E_\theta(D_\theta^2(T))$?

$E(T - E(T))$

$E(T) - E(E(T))$
$\quad \quad \searrow E(T)$

$''$
$0$

$E((T - E(T))^2)$

$Var(T)$

**Bias-Variance Decomposition**

$$
\begin{aligned}
\text{MSE}_\theta(T) &= E_\theta\left((T - \theta)^2\right) \\
&= E_\theta\left((D_\theta(T) + B_\theta(T))^2\right) \\
&= E_\theta\left(D_\theta^2(T) + 2B_\theta(T)D_\theta(T) + B_\theta^2(T)\right) \\
&= E(D^2(T)) + 2E(B(T) \cdot D(T)) + E(B^2(T)) \\
&= Var(T) + 2B(T) \cdot E(D(T)) + B^2(T) \\
&= Var(T) + B^2(T)
\end{aligned}
$$

$''$
$0$