

Announcements

Wednesday, February 29, 2012
1:25 PM

Feb 29

- Please submit the self-evaluation form by Friday (Mar 2).
- The form is posted on the course website and will be distributed in class (today).
- You can take a look at your midterm later today.

Mar 2

- Please submit the self-evaluation form by the end of today.
(5:20 PM)
- Note part I.7 is @ the copy center and website.

Mar 9

- HW 5 is posted : Due Mar 16 (Fri)

Mar 14

- Note part II.1 is now available.

Mar 16

- Submit HW 5 by 5 PM today
- Note Part II.2 is now available.

Mar 23

-Wed.

Mar 23

Wed.

- HW 6 (Mar 28)
- Note Part II.3 and Part II.4 are available.
- Start preparing the A4 sheet for the final exam.

2-sided
handwritten

Mar 28

- HW 7 and its solution will be posted later today.
- Note Part II.5 is available.
- There will be eval. form posted on the web site.

Mar 30

- Check the course website regularly (or at least before the final exam day)

9 Multiple RVs

Wednesday, February 29, 2012
1:25 PM

Previously, we talked about mean (average)
variance
standard deviation

of one set of data.

Now, we want to extend the calculation to > 1
sets of data.

Correlation
Covariance
correlation coefficient

One big example...

Suppose we have 10 students in a class...

To summarize these data, use a table

student id	Midterm score	Final score
1	10	40
2	10	40
3	10	50
4	20	40
5	20	40
6	20	50
7	20	50
8	20	50
9	20	50
10	20	50

$x \backslash y$	40	50
10	2	1
20	2	5

Randomly select one of the students

X = his/her midterm score

Y = his/her final score

what is the probability that $X = 10$ and $Y = 40$?

$$P[X = 10 \text{ and } Y = 40] = \frac{2}{10}$$

|||

$$P_{X,Y}(10, 40)$$

Defn $P_{X,Y}(x,y) = P[X=x \text{ and } Y=y]$

joint pmf

$$P_{X,Y}(10, 20) = 0$$

$P_{X,Y}(x,y)$

$x \backslash y$	40	50

$$P_{X,Y}(10,20) = 0$$

$\alpha \setminus Y$	40	50
10	1/5	1/10
20	1/5	1/2

What if my interest is only on the midterm score?

For example, randomly select a student

what is the probability that
his/her midterm score = 10?

$$\begin{aligned} P_X(10) = P[X=10] &= \frac{2+1}{10} = \frac{3}{10} \\ &= \frac{2}{10} + \frac{1}{10} = \frac{3}{10} \quad \leftarrow \text{same} \\ &= P_{X,Y}(10,40) + P_{X,Y}(10,50) \\ &= \sum_Y P_{X,Y}(10,Y) \end{aligned}$$

$$P_X(20) = P[X=20] = \sum_Y P_{X,Y}(20,Y) = \frac{2}{10} + \frac{5}{10} = \frac{7}{10}$$

Formula:
$$\left. \begin{aligned} P_X(\alpha) &= \sum_Y P_{X,Y}(\alpha, Y) \\ P_Y(y) &= \sum_{\alpha} P_{X,Y}(\alpha, y) \end{aligned} \right\} \text{marginal pmf's.}$$

$$P_X(\alpha) = \begin{cases} 3/10, & \alpha=10 \\ 7/10, & \alpha=20 \\ 0, & \text{otherwise} \end{cases} \quad P_Y(y) = \begin{cases} 2/5, & y=40 \\ 3/5, & y=50 \\ 0, & \text{otherwise.} \end{cases}$$

$$E[X] = \sum_{\alpha} \alpha P_X(\alpha) = \frac{2}{10} \times 10 + \frac{7}{10} \times 20 = 17 \quad E[Y] = \sum_Y Y P_Y(Y) = \frac{2}{5} \times 40 + \frac{3}{5} \times 50 = 46$$

Given that $X=10$, what is the probability that $Y=50$?

$$= \frac{1}{3}$$

$$P[Y=50 | X=10] = \frac{P(A \cap B)}{P(B)} = \frac{P[X=10 \text{ and } Y=50]}{P[X=10]}$$

$$P_{Y|X}(50|10) = \frac{P_{X,Y}(10,50)}{P_X(10)} = \frac{1/10}{3/10} = \frac{1}{3}$$

Formula:
$$P_{Y|X}(y|\alpha) \equiv P[Y=y | X=\alpha] = \frac{P_{X,Y}(\alpha, y)}{P_X(\alpha)}$$

conditional pmf

$$\downarrow P_{X|Y}(x|y) = \frac{P_{X,Y}(x,y)}{P_Y(y)}$$

To calculate the total score we use

$$t = \frac{x}{20} \times 40 + \frac{y}{50} \times 40 + 20$$

Annotations:
 - $\frac{x}{20} \times 40$: max midterm score
 - $\frac{y}{50} \times 40$: max final score
 - 20 : HW/quiz/participation
 - $\frac{x}{20} \times 40 + \frac{y}{50} \times 40$: midterm = 40% of total score
 - $\frac{y}{50} \times 40 + 20$: final = 40% of total score

$x \setminus y$	40	50
10	72	80
20	92	100

Let T be the total score of a randomly selected student.

$$P_T(t) = \begin{cases} 1/5 & t=72, \\ 1/10 & t=80, \\ 1/5 & t=92, \\ 1/2 & t=100, \\ 0, & \text{otherwise} \end{cases}$$

$$E_T = \sum_t t P_T(t) = 90.8$$

Alternatively,

$$\begin{aligned} E_T &= E\left[\frac{x}{20} \times 40 + \frac{y}{50} \times 40 + 20\right] = E\left[2X + \frac{4}{5}Y + 20\right] \\ &= 2E_X + \frac{4}{5}E_Y + 20 \\ &= 2(17) + \frac{4}{5}(46) + 20 = 90.8 \end{aligned}$$

In general, if $T = g(X, Y)$

$$E_T = E[g(X, Y)] = \sum_x \sum_y g(x, y) P_{X,Y}(x, y)$$

9 Least Square

Friday, March 02, 2012
11:33 AM

Earlier in our class

Basic statistics ...

Suppose we have n observations : $x_1, x_2, x_3, x_4, \dots, x_n$

Find "b" to represent this data set.

To determine how good "b" is,
we will calculate

$$\sum_{i=1}^n (x_i - b)^2$$

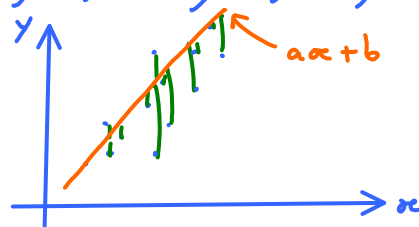
and minimize it.

→ the best "b" = $\frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$
Calculus

$$\frac{d}{db} = 0$$

Suppose that we have n pairs of observations

$$(x_1, y_1), (x_2, y_2), (x_3, y_3), \dots, (x_n, y_n)$$



Find the "best" straight line that "fits" the observations.

Again, want to minimize square error

$$\sum_{i=1}^n (y_i - (ax_i + b))^2$$

$$\text{EX } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

→

$\bar{x} = \bar{x}$

calculus
 $\frac{\partial}{\partial a} = 0$
 $\frac{\partial}{\partial b} = 0$

$$a = \frac{\overline{xy} - \bar{x}\bar{y}}{\overline{x^2} - (\bar{x})^2} \quad \mathbb{E}Y \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

$$b = \bar{y} - a\bar{x} \quad \mathbb{E}[X^2] \quad \overline{x^2} = \frac{1}{n} \sum_{i=1}^n x_i^2$$

$$\underbrace{\mathbb{E}[XY]}_{\text{correlation}} \quad \overline{xy} = \frac{1}{n} \sum_{i=1}^n x_i y_i$$

$$a = \frac{\mathbb{E}[XY] - \mathbb{E}X \mathbb{E}Y}{\mathbb{E}[X^2] - (\mathbb{E}X)^2} = \frac{\text{Cov}(X, Y)}{\Delta_X^2} = \rho \frac{\Delta_Y}{\Delta_X}$$

$$\rho = \frac{\text{Cov}(X, Y)}{\Delta_X \Delta_Y} \Rightarrow \text{Cov}(X, Y) = \rho \Delta_X \Delta_Y$$

$$y \approx ax + b$$

$$Y \approx aX + b = aX + (\mathbb{E}Y - a\mathbb{E}X)$$

$$= \mathbb{E}Y + a(X - \mathbb{E}X)$$

$$= \mathbb{E}Y + \rho \frac{\Delta_Y}{\Delta_X} (X - \mathbb{E}X)$$

Here is how things get crazy... in statistics...

$$(X_1, Y_1), (X_2, Y_2), (X_3, Y_3), \dots, (X_n, Y_n)$$

Post-midterm Review

Friday, March 09, 2012
10:37 AM

Two RVs X, Y

Correlation coefficient

$$\rho_{X,Y} = \frac{\text{Cov}[X,Y]}{\sigma_X \sigma_Y}$$

$$-1 \leq \rho_{X,Y} \leq 1$$

$\rho_{X,Y}$ captures linear (affine) relationship btw X and Y .

Continuous RV X, Y, Z

~~pdf~~ $P[X=5] = P[X=3] = P[X=1.3] = P[X=\alpha] = 0$

use pdf $f_X(\alpha)$

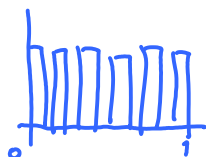
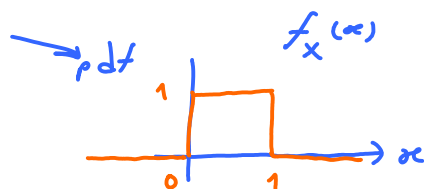
$$P[3 \leq X < 5] = \int_3^5 f_X(\alpha) d\alpha$$

$$P[X^2 > 1] = \int_{-\infty}^{-1} f_X(\alpha) d\alpha + \int_1^{\infty} f_X(\alpha) d\alpha$$

Important example of continuous random variables:

$x = \text{rand}()$

↓
histogram



Important properties of pdf $\begin{cases} \geq 0 \\ \int = 1 \end{cases}$

Another important example: Gaussian/Normal RV

$$f_X(x) = \frac{1}{\sqrt{2\pi}\Delta} e^{-\frac{1}{2}\left(\frac{x-m}{\Delta}\right)^2}$$

$$\mathbb{E}X = \int_{-\infty}^{\infty} x \frac{1}{\sqrt{2\pi}\Delta} e^{-\frac{1}{2}\left(\frac{x-m}{\Delta}\right)^2} dx = m$$

$$\mathbb{E}[X^2] = \int_{-\infty}^{\infty} x^2 \frac{1}{\sqrt{2\pi}\Delta} e^{-\frac{1}{2}\left(\frac{x-m}{\Delta}\right)^2} dx = m^2 + \Delta^2$$

calculus

$$\text{Var} X = \mathbb{E}[X^2] - (\mathbb{E}X)^2 = m^2 + \Delta^2 - m^2 = \Delta^2$$

$$\Delta_X = \sqrt{\text{Var} X} = \Delta$$

General Gaussian/Normal RV : $X \sim \mathcal{N}(m, \Delta^2)$

$$\text{CDF: } F_X(x) = P[X \leq x] \quad \begin{array}{l} \downarrow \\ m=0 \\ \Delta=1 \end{array}$$

Standard Gaussian/Normal RV : $Z \sim \mathcal{N}(0, 1)$

$$\text{CDF: } F_Z(z) = P[Z \leq z] = \Phi(z)$$

↑
can evaluate this using Table

$$P[a \leq Z \leq b] = \Phi(b) - \Phi(a)$$

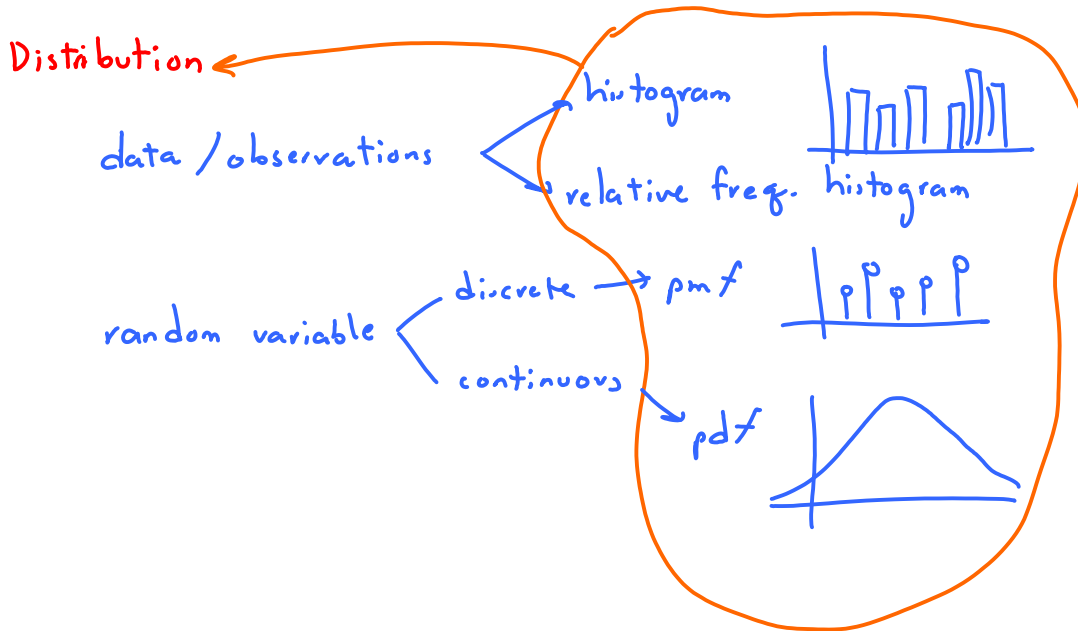
This can be extended to find the probability involving general Gaussian RV $X \sim \mathcal{N}(m, \Delta^2)$.

$$\text{Link: } \frac{X - \mathbb{E}X}{\Delta_X} = \frac{X - m}{\Delta} \sim \mathcal{N}(0, 1)$$

$$\Rightarrow P[a \leq X \leq b] = P\left[\frac{a-m}{\Delta} \leq \frac{X-m}{\Delta} \leq \frac{b-m}{\Delta}\right]$$

↑
 $Z \sim \mathcal{N}(0, 1)$

$$= \Phi\left(\frac{b-m}{\delta}\right) - \Phi\left(\frac{a-m}{\delta}\right)$$



Sampling Distribution

For us, we focus on sample mean (\bar{X})



estimate population mean (μ)

Assumption: For this part of the class, we will assume that the population is large.

(infinite)

We will model the members of the population by random variables all with the same distribution (pmf, pdf).

Now, under the assumption above, we consider the distribution of \bar{X} .

Theorem 1: If the population is governed by

Normal/Gaussian distribution, $\mathcal{N}(\mu, \sigma^2)$

then $\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$

Theorem 2: For any population, if n is large (CLT) ($n \geq 30$), then

$$\bar{X} \approx \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

Goal: Estimate parameter of the population

↓
we focus on estimating the mean μ

└─ point estimate: \bar{X}
└─ interval estimate

↓
Find confidence interval (CI).

Section 13 Confidence interval (CI) (CL)

↳ interval estimate + confidence level
[L, u] 100(1- α)%

Quality of CI

1) CL: 100(1- α)%

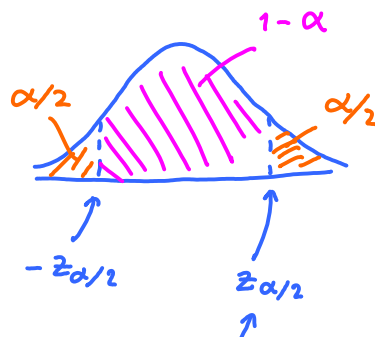
2) Precision (inversely related to the width of CI)

Construction of CI

↳ 13.1) σ is known

$$L = \bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$$u = \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$



get this value by finding the value of z in the Φ table such that $\Phi(z) = 1 - \frac{\alpha}{2}$.

13.2) Δ is unknown

$$L = \bar{x} - t_{\alpha/2, n-1} \frac{\hat{\Delta}}{\sqrt{n}}$$

$$u = \bar{x} + t_{\alpha/2, n-1} \frac{\hat{\Delta}}{\sqrt{n}}$$

14. Tests of Hypotheses

Inferential Statistics ; reaching conclusions based on sample information.

To do this, we use hypothesis testing

↓
a statement about population parameter(s)

for us, we focus on the population mean (μ)

We wish to test

$$H_0: \mu = \mu_0$$

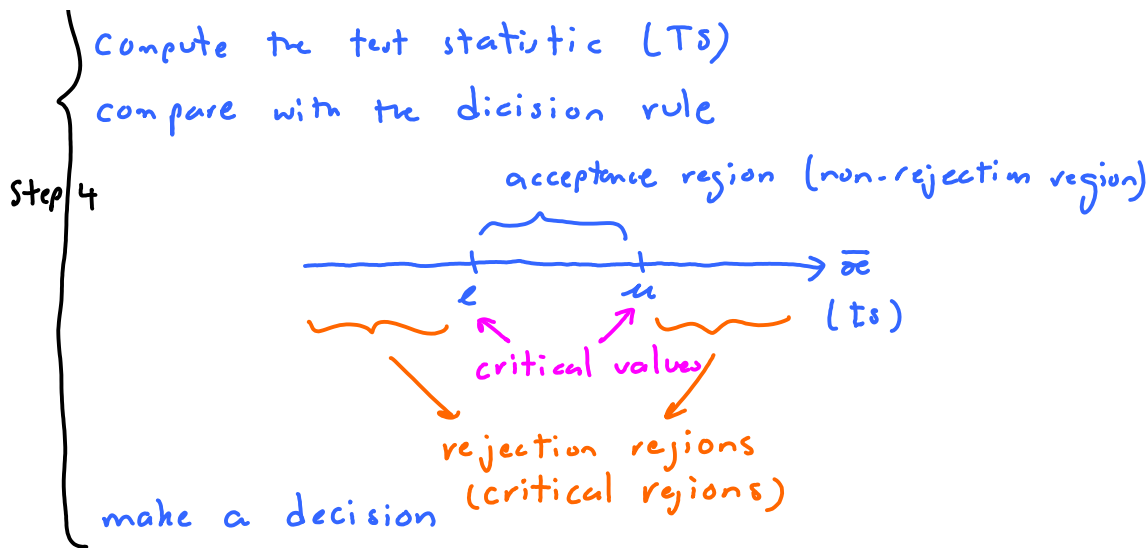
$$H_1: \mu \neq \mu_0$$

First step: form H_0, H_1

Two types of error

- False positive : Type I \Rightarrow probability = α
significant level
- False negative : Type II \Rightarrow probability = β

(Take a sample



If $\bar{x} \in [l, u]$, then conclude that ~~H_0 is true.~~
we fail to reject H_0

If $\bar{x} \notin [l, u]$, then conclude that H_1 is true.
we reject H_0

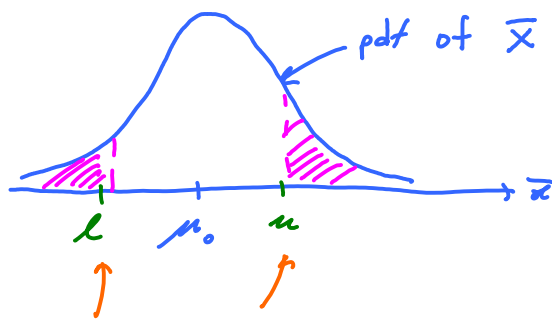
What is α ?

↳ probability of type I error.

= probability of rejecting H_0 when H_0 is true.

$$\mu = \mu_0$$

$$\bar{x} \sim \mathcal{N}(\mu_0, \frac{\sigma^2}{n})$$



suppose the critical values are l and u

Then $\alpha =$ the shaded area.

Observe: small $\alpha \rightarrow$ larger $u-l \rightarrow$ more demanding test

Second step: choose the value of α

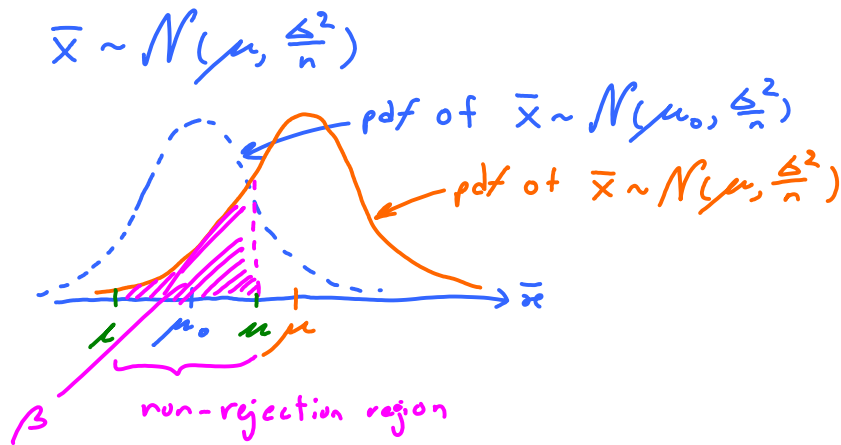
↳ usually,

$$\alpha = 0.05$$

Third step: Determine l, u

What about β ? (more difficult)

↑ probability of the false negative
= probability of failing to reject H_0 when
 H_1 is true.
↓
 $\mu \neq \mu_0$



Note: smaller $\alpha \rightarrow$ non-rejection region \rightarrow larger β
will be expanded