

Stats 401 Lab 8

Sanjana

10/25/2018

Outline

- ▶ Quick Reminder: If you are thinking about withdrawing from the course, the deadline is **November 9th!**
- ▶ Review of expectation, variance and covariance
- ▶ Midterm Q3 solved
- ▶ Introduction to confidence intervals

Review: Univariate Random Variable

- ▶ Recall: A random variable X is a value whose outcome is determined by a random process.

Can be thought of as a random number with probabilities assigned to its outcomes.

- ▶ Each random variable X has a mean (μ_x) and variance (σ_x^2) associated with it. Then

- ▶ $\mu_x = E[X]$.



$$\begin{aligned}\sigma_x^2 = \text{Var}(\mathbf{X}) &= E[(X - E[X])^2] \\ &= E[X^2 + (E[X])^2 - 2XE[X]] \\ &= E[X^2] + E[(E[X])^2] - E[2XE[X]] \\ &= E[X^2] + (E[X])^2 - 2E[X]E[X] \\ &= E[X^2] + (E[X])^2 - 2(E[X])^2 \\ &= E[\mathbf{X}^2] - (E[\mathbf{X}])^2\end{aligned}$$

Review: Univariate Random Variable

- ▶ Linear combinations of X : consider $Y = aX + b$
 - ▶ $E[aX + b] = aE[X] + b$
 - ▶ $\text{Var}(aX + b) = a^2 \text{Var}(X)$
- ▶ Normal Variable: let $X \sim \mathcal{N}(\mu_x, \sigma_x^2)$,
Let $Y = aX + b$
Then $E[Y] = a\mu_x + b$ and $\text{Var}(Y) = a^2\sigma_x^2$
 $Y \sim \mathcal{N}(a\mu_x + b, a^2\sigma_x^2)$

Review: Bivariate Random Variables

- ▶ A bivariate random variable (X, Y) is a vector of length 2 whose values are each random variables.
- ▶ *Correlation* is a measure of the linear association between two random variables
 - ▶ Correlation is always between -1 and 1 (inclusive)
 - ▶ Correlation is symmetric: $\text{Cor}(X, Y) = \text{Cor}(Y, X)$
- ▶ *Covariance* is the unscaled version of correlation

$$\begin{aligned}\text{Cov}(\mathbf{X}, \mathbf{Y}) &= \mathbf{E}[(\mathbf{X} - \mathbf{E}(\mathbf{X}))(\mathbf{Y} - \mathbf{E}(\mathbf{Y}))] \\ &= \mathbf{E}[(XY - X\mathbf{E}(Y) - \mathbf{E}(X)Y + \mathbf{E}(X)\mathbf{E}(Y))] \\ &= \mathbf{E}[XY] - \mathbf{E}[X\mathbf{E}(Y)] - \mathbf{E}[\mathbf{E}(X)Y] + \mathbf{E}[\mathbf{E}(X)\mathbf{E}(Y)] \\ &= \mathbf{E}[XY] - \mathbf{E}(X)\mathbf{E}(Y) - \mathbf{E}(X)\mathbf{E}(Y) + \mathbf{E}(X)\mathbf{E}(Y) \\ &= \mathbf{E}[\mathbf{XY}] - \mathbf{E}(\mathbf{X})\mathbf{E}(\mathbf{Y})\end{aligned}$$

$$\text{Corr}(\mathbf{X}, \mathbf{Y}) = \frac{\text{Cov}(\mathbf{X}, \mathbf{Y})}{\sqrt{\text{Var}(\mathbf{X})}\sqrt{\text{Var}(\mathbf{Y})}}$$

Review: Bivariate Random Variables

▶ $E[(X, Y)] = (E(X), E(Y))$

▶ $\text{Var}(X, Y) = \begin{bmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) \\ \text{Cov}(X_2, X_1) & \text{Var}(X_2) \end{bmatrix}$

▶ **Remember:** $\text{Var}(X) = \text{Cov}(X, X)$

▶ Linear Combinations of multivariate $\mathbf{X} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix}$: consider $\mathbb{A}\mathbf{X}$

▶ $E(\mathbb{A}\mathbf{X}) = \mathbb{A}E(\mathbf{X})$

▶ $\text{Var}(\mathbb{A}\mathbf{X}) = \mathbb{A}\text{Var}(\mathbf{X})\mathbb{A}^T$

Review: Bivariate Normal Variables

If (X, Y) is bivariate normal where $X \sim \mathcal{N}(\mu_X, \sigma_X^2)$ and $Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$ then $aX + bY$ is also normal.

$$- E[aX + bY] = a\mu_X + b\mu_Y$$

-

$$\begin{aligned}\mathbf{Var}(a\mathbf{X} + b\mathbf{Y}) &= \mathbf{Cov}(a\mathbf{X} + b\mathbf{Y}, a\mathbf{X} + b\mathbf{Y}) \\ &= \mathbf{Cov}(a\mathbf{X} + b\mathbf{Y}, a\mathbf{X}) + \mathbf{Cov}(a\mathbf{X} + b\mathbf{Y}, b\mathbf{Y}) \\ &= a\mathbf{Cov}(a\mathbf{X} + b\mathbf{Y}, \mathbf{X}) + b\mathbf{Cov}(a\mathbf{X} + b\mathbf{Y}, \mathbf{Y}) \\ &= a\mathbf{Cov}(a\mathbf{X}, \mathbf{X}) + a\mathbf{Cov}(b\mathbf{Y}, \mathbf{X}) + b\mathbf{Cov}(a\mathbf{X}, \mathbf{Y}) + b\mathbf{Cov}(b\mathbf{Y}, \mathbf{Y}) \\ &= a^2\mathbf{Cov}(\mathbf{X}, \mathbf{X}) + ab\mathbf{Cov}(\mathbf{Y}, \mathbf{X}) + ba\mathbf{Cov}(\mathbf{X}, \mathbf{Y}) + b^2\mathbf{Cov}(\mathbf{Y}, \mathbf{Y}) \\ &= a^2\mathbf{Var}(\mathbf{X}) + ab\mathbf{Cov}(\mathbf{X}, \mathbf{Y}) + ab\mathbf{Cov}(\mathbf{X}, \mathbf{Y}) + b^2\mathbf{Var}(\mathbf{Y}) \\ &= \mathbf{a}^2\mathbf{Var}(\mathbf{X}) + \mathbf{2abCov}(\mathbf{X}, \mathbf{Y}) + \mathbf{b}^2\mathbf{Var}(\mathbf{Y}) \\ &= a^2\sigma_X^2 + 2ab\mathbf{Cov}(\mathbf{X}, \mathbf{Y}) + b^2\sigma_Y^2\end{aligned}$$

So, if (X, Y) is bivariate normal as above, then

$$aX + bY \sim \mathcal{N}(a\mu_X + b\mu_Y, a^2\sigma_X^2 + b^2\sigma_Y^2 + 2ab\mathbf{Cov}(\mathbf{X}, \mathbf{Y}))$$

In-lab Activity 1: Midterm Q3

Qtn) X and Y are bivariate random variables with respective means $\mu_X = \mu_Y = 0$, standard deviations $\sigma_X = 1$ and $\sigma_Y = 2$ and correlation $\text{cor}(X, Y) = 0.5$. Find the distributions of $X + Y$ and $X - Y$.

Soln) $X + Y$ and $X - Y$ are both normally distributed since they are linear combinations of normal random variables. So, we need to find the following

- ▶ $E(X + Y)$ and $E(X - Y)$
- ▶ $\text{Var}(X + Y)$ and $\text{Var}(X - Y)$

We are given that

- $E(X) = 0$, $\text{Var}(X) = \sigma_X^2 = 1$
- $E(Y) = 0$, $\text{Var}(Y) = \sigma_Y^2 = 2^2 = 4$
- $\text{Cor}(X, Y) = 0.5$

So, $\text{Cov}(X, Y) = \text{Cor}(X, Y)\sqrt{\sigma_X^2\sigma_Y^2} = \text{Cor}(X, Y)\sigma_X\sigma_Y = 0.5(1)(2) = 1$

In-lab Activity 1: Midterm Q3

- ▶ $E[X + Y] = E[X] + E[Y] = 0 + 0 = 0$ and
 $E[X - Y] = E[X] - E[Y] = 0 - 0 = 0$
- ▶ $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y) = 1 + 4 + 2(1) = 7$
- ▶ $\text{Var}(X - Y) = \text{Var}(X) + \text{Var}(-Y) + 2\text{Cov}(X, -Y) = 1 + 4 - 2(1) = 3$

So, $X + Y \sim \mathcal{N}(0, 7)$ and $X - Y \sim \mathcal{N}(0, 3)$

In-lab Activity 2

Let (W, X, Y) be a multivariate normal vector such that

$$E(X) = 2E(Y) = 2 \text{ and } E(W) = 0.$$

$$\text{Var}(X) = \text{Var}(Y) = \text{Var}(W) = 2.$$

$$\text{Cor}(X, Y) = -0.5, \text{Cor}(Y, W) = -0.5, \text{Cor}(X, W) = 0.$$

Find the distribution of $W + X - Y$.

In-lab Activity 2: Solution

$X - Y + W$ is a normal variable, since linear combinations of normal variables are normal.

$$\blacktriangleright E(X - Y + W) = E(X) + E(Y) + E(W) = 2 - 1 + 0 = 1$$

$$\blacktriangleright \text{Var} \begin{pmatrix} X \\ Y \\ W \end{pmatrix} = \begin{bmatrix} \text{Var}(X) & \text{Cov}(X, Y) & \text{Cov}(X, W) \\ \text{Cov}(Y, X) & \text{Var}(Y) & \text{Cov}(Y, W) \\ \text{Cov}(W, X) & \text{Cov}(W, Y) & \text{Var}(W) \end{bmatrix}$$

$$\blacktriangleright \text{Cov}(X, Y) = \text{Cor}(X, Y) \sqrt{\text{Var}(X)\text{Var}(Y)} = -0.5\sqrt{22} = -1$$

$$\blacktriangleright \text{Cov}(X, W) = \text{Cor}(X, W) \sqrt{\text{Var}(X)\text{Var}(W)} = 0\sqrt{22} = 0$$

$$\blacktriangleright \text{Cov}(Y, W) = \text{Cor}(Y, W) \sqrt{\text{Var}(Y)\text{Var}(W)} = -0.5\sqrt{22} = -1$$

$$\text{So, } \text{Var} \begin{pmatrix} X \\ Y \\ W \end{pmatrix} = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}$$

In-lab Activity 2: Solution

$$\text{Then, } X - Y + W = [1 \quad -1 \quad 1] \begin{bmatrix} X \\ Y \\ W \end{bmatrix} = \mathbb{A}\mathbf{X}, \text{ where } \mathbf{X} = \begin{bmatrix} X \\ Y \\ W \end{bmatrix}$$

$$\begin{aligned} \text{Var}(X - Y + W) &= \text{Var}(\mathbb{A}\mathbf{X}) = \mathbb{A}\text{Var}(\mathbf{X})\mathbb{A}^T \\ &= [1 \quad -1 \quad 1] \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} \\ &= [1 \quad -1 \quad 1] \begin{bmatrix} 2(1) + (-1)(-1) + 0(1) \\ (-1)1 + 2(-1) + (-1)1 \\ 0(1) + (-1)(-1) + 2(1) \end{bmatrix} \\ &= [1 \quad -1 \quad 1] \begin{bmatrix} 3 \\ -4 \\ 3 \end{bmatrix} \\ &= 1(3) + (-1)(-4) + 1(3) \\ &= 10 \end{aligned}$$

So, $\mathbf{X} - \mathbf{Y} + \mathbf{W} \sim \mathcal{N}(\mathbf{1}, \mathbf{10})$

Confidence Intervals

- ▶ We are we interested in studying confidence intervals?
 - ▶ CIs essentially perform a two-sided hypothesis test and provide you with a estimate the true population value
- ▶ There are several natural uses for confidence intervals in regression:
 - ▶ estimating population coefficients (β)
 - ▶ comparing means of different populations
 - ▶ predicting future values (prediction interval)
 - ▶ predicting mean future values (confidence interval)

Confidence Intervals: formulae

- ▶ Recall from Stats250 that a $100(1 - \alpha)\%$ confidence interval for a value is given by
 - ▶ $x \pm z_{\frac{\alpha}{2}} \text{s.e.}(x)$ (population s.d. is known) or
 - ▶ $x \pm t_{(\frac{\alpha}{2}, n-2)} \text{s.e.}(x)$ (population s.d. is unknown)
- ▶ *Approximate Interval for Linear Model* An approximate $100(1 - \alpha)$ CI for β_1 is

$$\mathbf{b}_1 \pm z_{\frac{\alpha}{2}} \text{SE}(\mathbf{b}_1)$$

In-lab Activity 3: Constructing CI in R

Construct a 90% CI for the association between GPA and ACT scores

```
# read-in dataset
gpa <- read.table("gpa.txt",header=T)
# fit model and print coefficients summary
fit <- lm(GPA~ACT, data=gpa)
fit_coef_smry <- summary(fit)$coefficients; fit_coef_smry
```

```
##              Estimate Std. Error t value    Pr(>|t|)
## (Intercept) 1.55870218 0.138016669 11.29358 2.694524e-27
## ACT          0.05780049 0.005549788 10.41490 1.015955e-23
```

```
beta <- fit_coef_smry["ACT","Estimate"]
SE <- fit_coef_smry["ACT","Std. Error"]
z <- qnorm(1-0.1/2) #for a 90% CI using normal approximation
cat("CI = [", beta-z*SE, ", ", beta+z*SE, "]" )
```

```
## CI = [ 0.04867191 , 0.06692908 ]
```

Lab Ticket

- ▶ Let X and Y be two random variables such that $\text{Var}(X) = 4$ and $\text{Var}(Y) = 9$. Find the range for the $\text{Cov}(X, Y)$
- ▶ Let X, Y be bivariate normal such that they are independent. Further, $X \sim \mathcal{N}(1, 4)$ and $Y \sim \mathcal{N}(2, 9)$. Find the distribution of the following:
 - ▶ $2X + 1$
 - ▶ $X - 2Y$
- ▶ Would a 90% confidence interval be broader or a 95%? Justify your answer.