

Stat 134 Sec 38

Last time

Sec 6.4 Correlation

$$\text{Corr}(x, y) = \frac{\text{Cov}(x, y)}{\text{SD}(x)\text{SD}(y)} = E(x^*y^*)$$

$$-1 \leq \text{Corr}(x, y) \leq 1$$

$$\text{Corr}(ax+b, cy+d) = \pm \text{Corr}(x, y) \quad a, b, c, d \in \mathbb{R}, \quad a, c \neq 0.$$

$N_1 + \dots + N_k = c$ a sum of exchangeable RVs

$$\text{Corr}(N_1, N_2) = -\frac{1}{k-1}$$

x in
std units

$+ \text{if } ac > 0$

Today

- ① Review student comments from last time.
- ② Sec 6.5 the bivariate normal distribution

① Review of student comments

Stat 134

Friday April 26 2019

1. An urn contains 90 marbles, of which there are 20 green, 25 black, and 45 red marbles. Alice randomly picks 10 marbles without replacement, and Bob randomly picks another 10 marbles without replacement. Let X_1 be the number of green marbles that Alice has and X_2 the number of green marbles that Bob has.

To find $\text{Corr}(X_1, X_2)$ is

$$X_1 + X_2 + \dots + X_9 = 20$$

a true identity that is useful? Explain.

a yes

b no

c not enough info to decide

Discuss with your neighbor for 1 minute how you did this.

Not true.

Imagine example of deck of 52 cards.

What is the chance there are 4 hearts in cards 14-26?

$$\text{answ } \frac{\binom{13}{4} \binom{39}{13}}{\binom{52}{13}}$$

b

The x_i are not identically distributed as the distribution of x_i depends on the values of x_j $j < i$

a

The X_n s are identical, each being hypergeometric(90, 10, 20).



a

Assign the marbles in lots of 10, this makes X_1 through X_9 identically distributed (but not remotely independent), so we can use the trick in the guided example to find the correlation

Pull marbles
 $M_1 - M_{20}$ out
without looking at other marbles.
 $X_2 \sim \text{Hyper}(90, 10, 20)$

a

After 9 turns, 90 marbles have been picked and those 90 marbles must contain all of the 20 green marbles since the marbles are picked without replacement

X_1, \dots, X_9 i.d.
and equal to 20
so

$$\text{Cov}(X_1, X_2) = -\frac{1}{q-1} = -\frac{1}{8}$$

② sec 6.5 Bivariate normal.

- A generalization of joint distribution of 2 iid normal to 2 correlated normals.
- An introduction to linear regression. We can learn how to predict final exam score from midterm score.

let X, Z iid $N(0, 1)$, $-1 \leq \rho \leq 1$

$$Y = \rho X + \sqrt{1-\rho^2} Z$$

Y is normal being a linear combination of independent normals.

$$E(Y) = \rho E(X) + \sqrt{1-\rho^2} E(Z) = 0$$

$$\text{Var}(Y) = \rho^2 \text{Var}(X) + (1-\rho^2) \text{Var}(Z) = 1$$

$$\Rightarrow Y \sim N(0, 1)$$

Find $\text{Corr}(X, Y)$

$$= \frac{\text{Cov}(X, Y)}{\text{SD}(X)\text{SD}(Y)}$$

$$\begin{aligned} \text{Cov}(X, Y) &= \text{Cov}(X, \rho X + \sqrt{1-\rho^2} Z) \\ &= \underset{n}{\text{Cov}}(X, \rho X) + \underset{0}{\text{Cov}}(X, \sqrt{1-\rho^2} Z) \\ &= \rho \text{Var}(X) \end{aligned}$$

$$\text{Corr}(X, Y) = \frac{\rho \text{Var}(X)}{\text{SD}(X)\text{SD}(Y)} = \boxed{\rho}$$

Defⁿ (Standard Bivariate normal distribution)

let $X, Z \sim \text{iid } N(0, 1)$, $-1 \leq \rho \leq 1$

$$Y = \rho X + \sqrt{1-\rho^2} Z$$

We call the joint distribution (X, Y) the

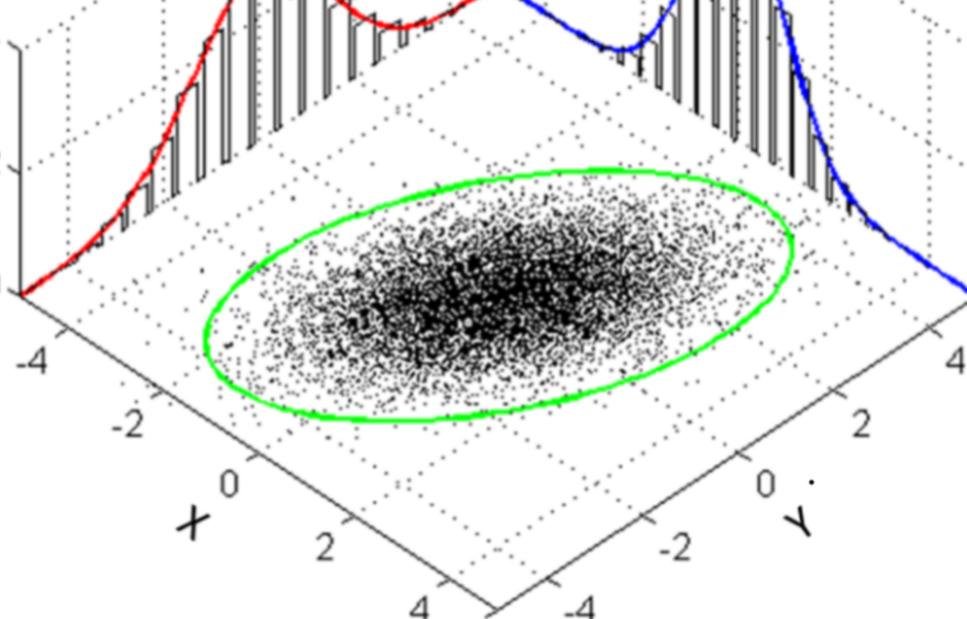
Standard bivariate normal with $\text{corr}(X, Y) = \rho$

Picture

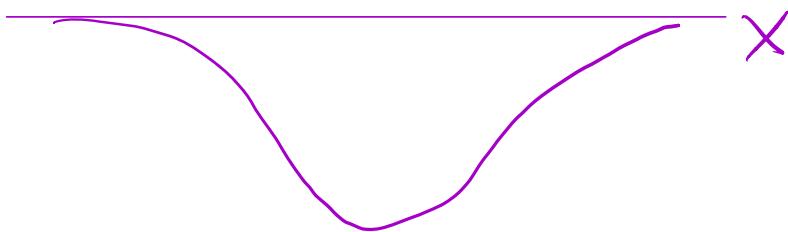
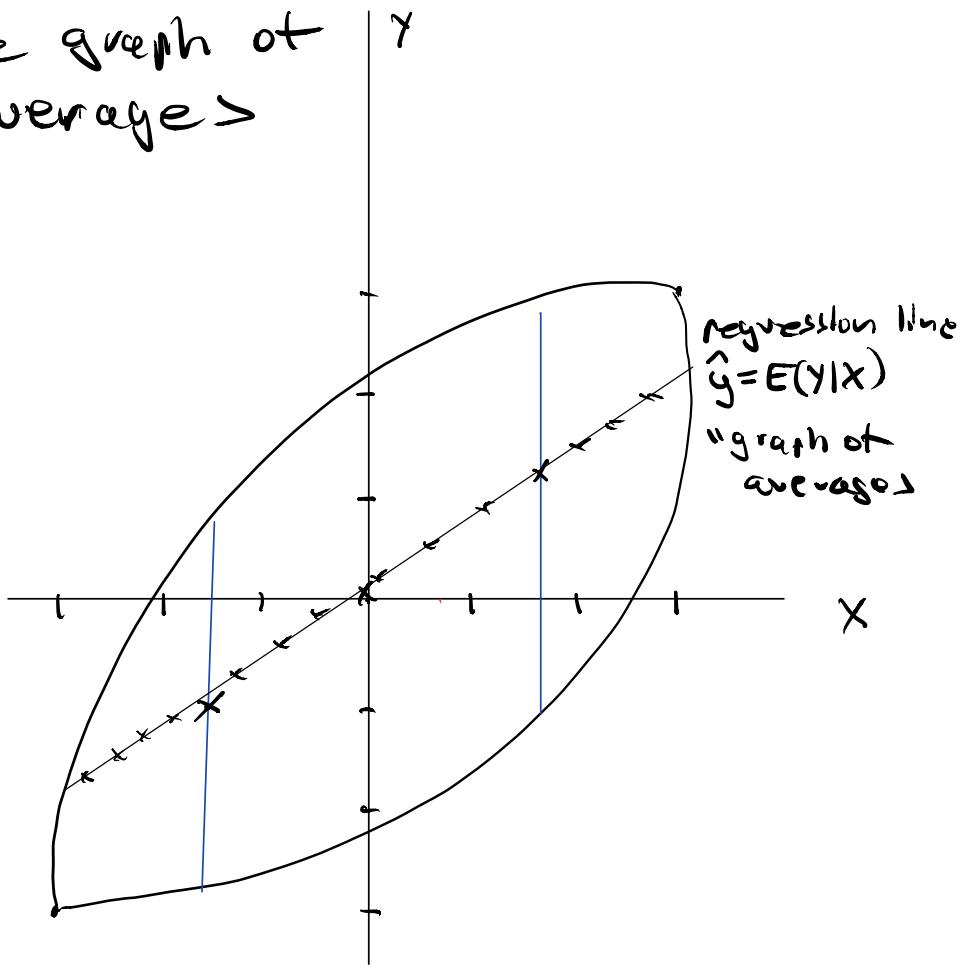
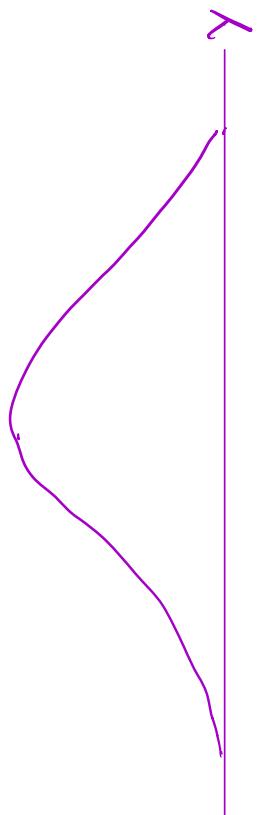
$$\rho > 0$$

$$p(Y)$$

$$p(X)$$



Draw the graph of
average >



$$y = \rho x + \sqrt{1-\rho^2} z$$

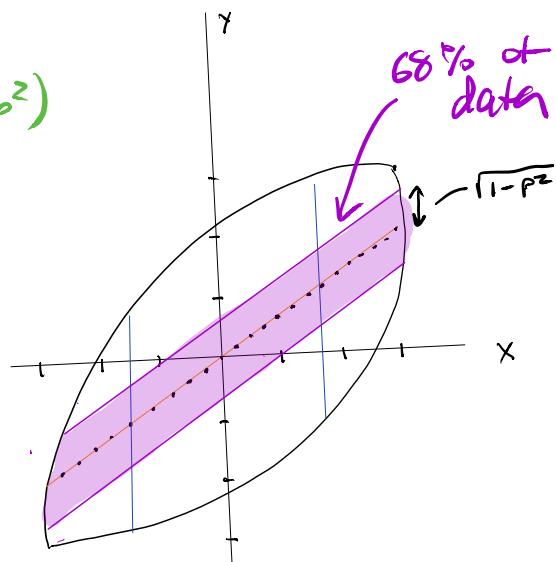
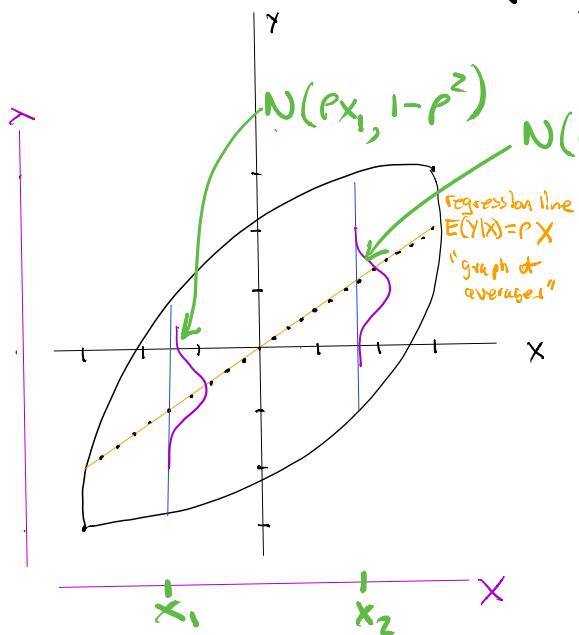
Find the "regression line" $E(y|x)$

$$\begin{aligned} E(y|x) &= E(\rho x + \sqrt{1-\rho^2} z | x) \\ &= E(\rho x | x) + E(\sqrt{1-\rho^2} z | x) \\ &= \rho x + \sqrt{1-\rho^2} E(z) \\ &= \rho x \end{aligned}$$

Find the conditional variance $\text{Var}(y|x)$

$$\begin{aligned} \text{Var}(y|x) &= \text{Var}(\rho x + \sqrt{1-\rho^2} z | x) \\ &= \text{Var}(\rho x | x) + \text{Var}(\sqrt{1-\rho^2} z | x) \quad \text{since } x, z \text{ indep.} \\ &= \rho^2 \text{Var}(x|x) + (1-\rho^2) \text{Var}(z|x) \\ &\quad \text{since } \text{Var}(z) = 1 \\ &= 1 - \rho^2 \end{aligned}$$

Hence $y|x \sim N(\rho x, 1-\rho^2)$



Defⁿ (Bivariate Normal Distribution)

Random variables U and V have bivariate normal distribution with parameters $\mu_U, \mu_V, \sigma_U^2, \sigma_V^2, \rho$ iff the standardized variables $X = \frac{U - \mu_U}{\sigma_U}$, $Y = \frac{V - \mu_V}{\sigma_V}$ have std. bivariate normal distribution with corr ρ .

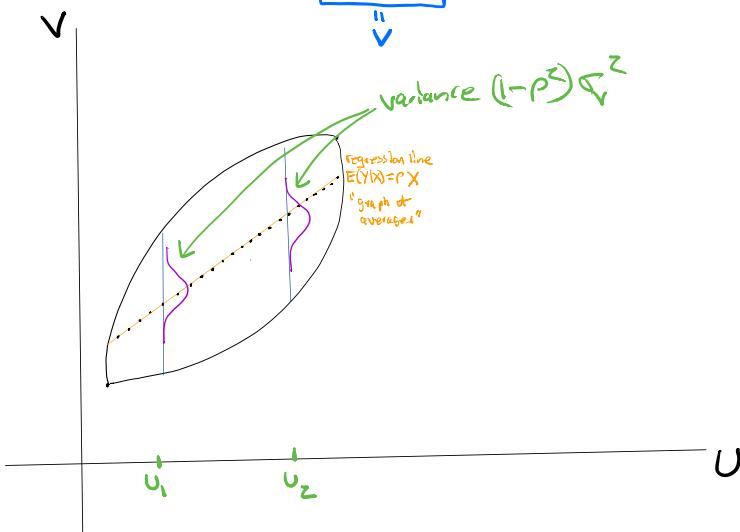
Then $\rho = \text{corr}(X, Y) = \text{corr}(U, V)$.

regression line of bivariate normal distribution

$$\begin{aligned} & \frac{V - \mu_V}{\sigma_V} = \rho \frac{U - \mu_U}{\sigma_U} \quad \text{regression line in S.V.} \\ \Leftrightarrow & V - \mu_V = \frac{\sigma_V}{\sigma_U} \rho (U - \mu_U) \\ \Leftrightarrow & \hat{V} = \left(\frac{\sigma_V}{\sigma_U} \rho \right) U + \mu_V - \frac{\sigma_V}{\sigma_U} \rho \mu_U \quad \text{regression line.} \end{aligned}$$

Facts

$$\begin{aligned} \hat{V} &= E(V|U) = \frac{\sigma_V}{\sigma_U} \rho U + \mu_V - \frac{\sigma_V}{\sigma_U} \rho \mu_U \quad \text{regression line} \\ \text{and } \text{var}(V|U) &= \text{var}(\underbrace{\sigma_V y + \mu_V}_{\hat{V}}|U) = \sigma_V^2 \text{var}(y|U) = (1 - \rho^2) \sigma_V^2 \end{aligned}$$



$$\begin{array}{ll}
 \text{Test 1 is } & M_U = 60 \\
 & \sigma_U = 20 \\
 \text{Test 2 is } & M_V = 60 \\
 & \sigma_V = 20
 \end{array}
 \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \rho = .6$$

a) Find the regression line

$$\begin{aligned}
 \hat{V} &= E(V|U) = \frac{\sigma_V}{\sigma_U} \rho U + M_V - \frac{\sigma_V}{\sigma_U} \rho M_U \quad \text{regression line} \\
 &= \frac{20}{20} (.6) U + 60 - \frac{20}{20} (.6)(60) \\
 \boxed{\hat{V}} &= .6U + 24
 \end{aligned}$$

b) If you get a 70 on Test 1 what score do you predict to get on Test 2?

$$\hat{V} = .6(70) + 24 = \boxed{66}$$

Notice you did relatively worse. Your test 1 score was $\frac{10}{20}$ sd above average but your test 2 score was only $\frac{6}{20}$ sd above average.

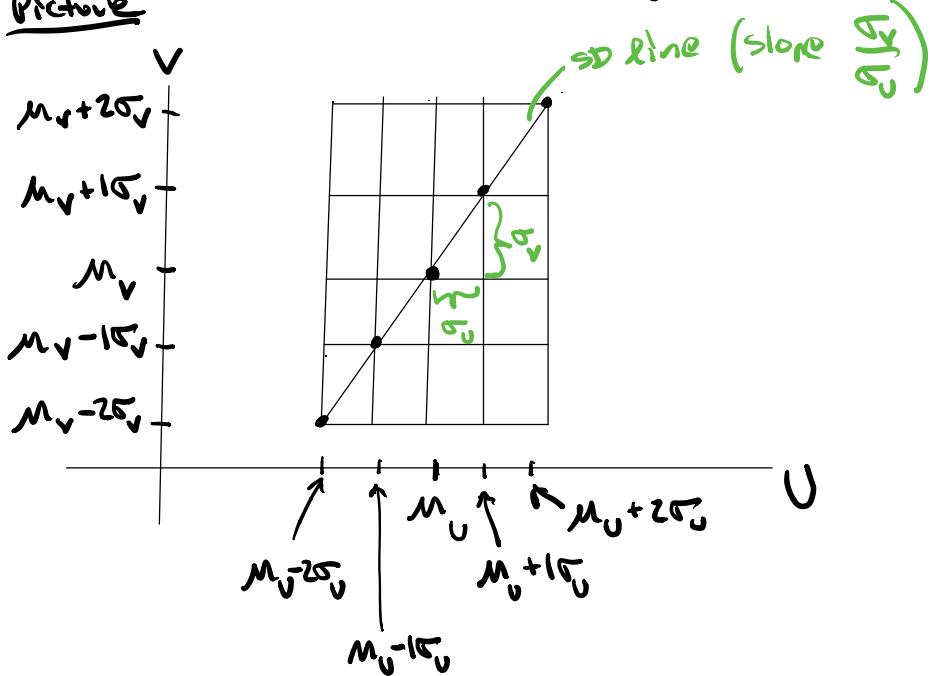
This is the "regression effect".

(3)

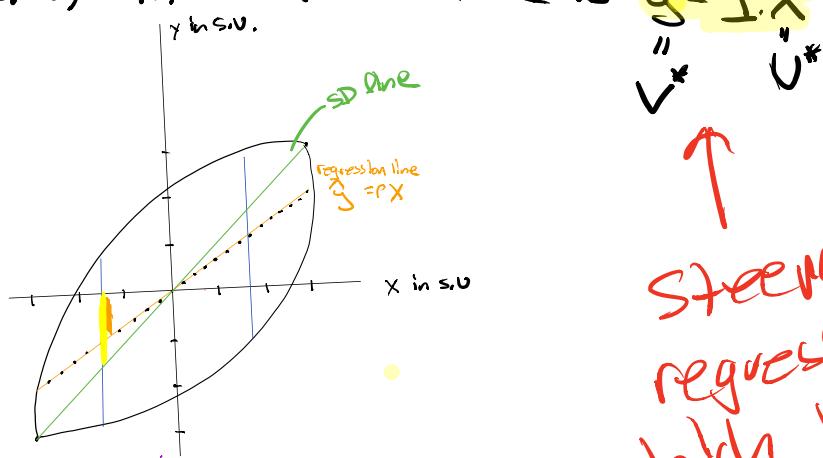
Regression line vs. SD line and regression effect

Def'n the SD line is $V - \mu_V = \frac{\sigma_V}{\sigma_U} (U - \mu_U)$.

Picture



For $U, V \sim s.u.$ the SD line is $y = 1 \cdot x$



↑
Steeper than
regression line
which has
slope p .

Regression effect

The regression line has slope $-1 \leq p \leq 1$ compared with the SD line which has slope 1. For a fixed x you predict $y = px$ which will be less than x if $x > 0$ and greater than x if $x < 0$. This means that if you do really well on a midterm (at least greater than 50th percentile — so $x > 0$ in s.u.) then you won't do as well on the final relative to the class (i.e. you won't go all the way up to the SD line). The opposite is true however if you do poorly on the midterm.

Picture

