

# Summary of Week 2

## Monday

- Expectation / Variance
  - ↳ allow us to summarize important aspects of distribution
- Binomial distribution
  - ↳ example of iid sum model of election/survey

## Wednesday

- Histograms in python
- Properties of binomial and large  $N$  limit

## Friday

- LLN, CLT
  - ↳ allows us to easily describe dist of iid sum
- Normal dist and density
  - ↳ density is mathematical identification of histogram

## This week

- linear regression Model basics
  - ↳ covariance, least squares
- Working w/ tabular data in python

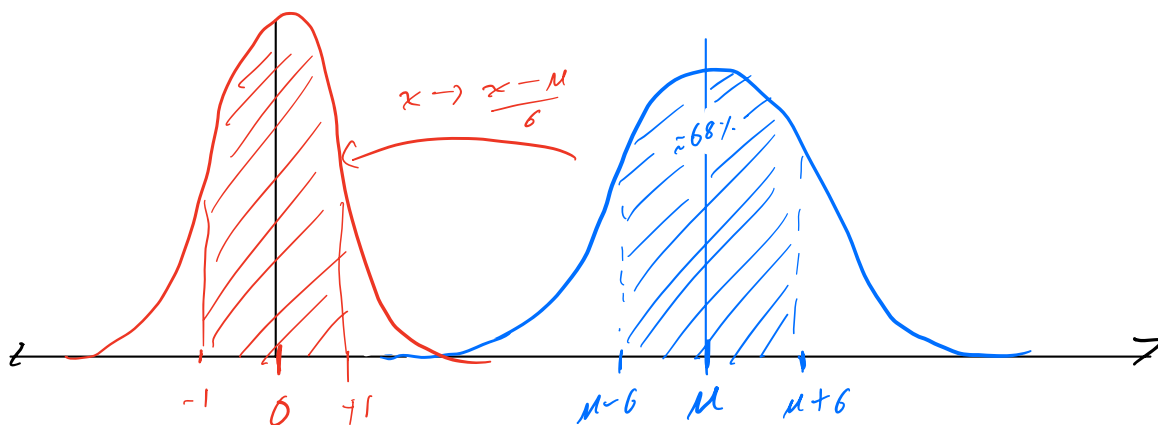
---

## Properties of Normal random variables (ch 4.6)

$$\underline{X} \sim \text{Normal}(\mu, \sigma^2) \text{ or } \underline{X} \sim N(\mu, \sigma^2)$$

is a normal r.v. and has density

$$f_{\underline{X}}(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} = \frac{1}{\sigma} \phi\left(\frac{x-\mu}{\sigma}\right)$$



## Combining Normal r.v.s (Theorem 4.6.1)

Let  $Z_1 \sim N(0,1)$ ,  $\mu, \sigma$  constants. Then let

$$X = \sigma Z_1 + \mu \Rightarrow E[X] = \sigma E[Z_1] + \mu = \mu$$

$$\text{Var}(X) = \sigma^2 \text{Var}(Z_1) = \sigma^2$$

1)  $X$  is normal  $\Rightarrow X \sim N(\mu, \sigma^2)$

Let  $X_i \sim \text{Normal}(\mu_i, \sigma_i^2)$  be independent

2) Then  $Y = \sum_{i=1}^N X_i \sim N\left(\sum_{i=1}^N \mu_i, \sum_{i=1}^N \sigma_i^2\right)$

- might talk about Theorem 4.6.2 later

- will talk about chi-squared in week 4

---

Note on CLT: CLT says  $Z_1 = \frac{S - n\mu}{\sqrt{n}\sigma} \rightarrow N(0,1)$

*not technically true*

meaning  $P(a < Z < b) \rightarrow \int_a^b \phi(x) dx$

$$S = \sqrt{n}\sigma Z_1 + n\mu, \quad S \not\rightarrow N(n\mu, \sqrt{n}\sigma)$$

Still, often think of CLT as saying  $S$  is approximately Normal

# Linear Regression Model (ch 10.)

related variables (Def 10.1.1) = not independent

## linear regression w/ one predictor (Ex 10.1.1)

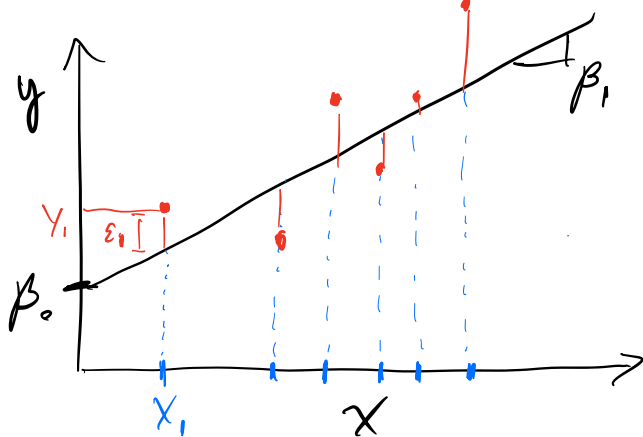
Let  $X$  be any r.v. and

$$Y | X \sim N(\beta_0 + \beta_1 X, \sigma^2)$$

Can also write

$$Y = \beta_0 + \beta_1 X + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

"error" or "residual" term



Very important model because  
simple linear relationship  
+ Normal error  
 $\Rightarrow$  we can easily  
estimate  $\beta_0, \beta_1, \sigma^2$

$X$  = predictor variable

(e.g. <sup>(quantitative)</sup> weight of dog or <sup>(qualitative)</sup> breed)

$Y$  = response variable

(e.g. lifespan)

# Meaning of parameters

Param	formula	description	units
$\beta_0$	$E[Y   X=0]$	avg. of $Y$ when $X=0$	units of $Y$
$\beta_1$	$E[Y   X=x+1] - E[Y   X=x]$ ( $x$ can be anything)	avg. change in $Y$ when we change $X$ by unit	$\frac{\text{units of } Y}{\text{units of } X}$
$\sigma^2$	$\text{Var}(Y   X=x)$ ( $x$ can be anything)	var of $Y$ for fixed $x$	units of $Y^2$

## Example (model of height)

$$X \sim \text{Bernoulli}(1/2) \quad (1 = \text{male}, 0 = \text{female})$$

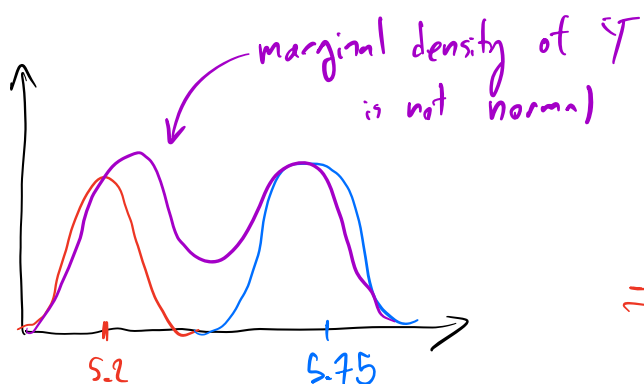
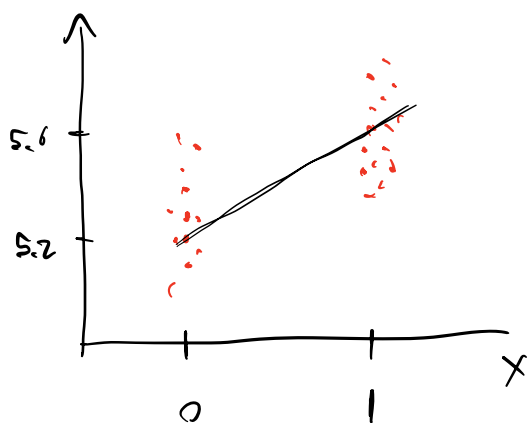
$$Y | X \sim N(5.2 + 0.55X, 0.25^2) \quad (\text{height in inches})$$

Note:  $E[Y | X=0] = 5.2 = \beta_0$

$E[Y | X=1] = 5.75 = \beta_1$

$\text{Var}(Y | X=1) = \text{Var}(Y | X=0) = 0.25^2$

$\Rightarrow$  in order to estimate  $\beta_0, \beta_1, \sigma$  we merely need to compute mean and variance within each group



Q: what is prob. male  $> 6.25$  ft?

$$6.25 \approx 5.75 + 0.25 \times 2$$

= mean of male height + 2 standard dev

$\Rightarrow P(Y > 6.1 | X=1) =$   $\approx 2\%$

Let  $X$  have mean and variance  $\mu_X$  and  $\sigma_X^2$

$$Y|X \sim N(\beta_0 + \beta_1 X, \sigma)$$

## Marginal mean and variance

$$\mathbb{E}[Y] = \mathbb{E}[\beta_0 + X\beta_1 + \varepsilon] = \beta_0 + \beta_1 \mu_X$$

$$\text{Var}(Y) = \beta_1^2 \sigma_X^2 + \sigma^2$$

more spread of  $X$  points  $\rightarrow$  more  $Y$  variation  
larger slope  $\rightarrow$  more  $Y$  variation

## Least Squares

Question: How do we estimate slope?

In principle we could use sample avg.:

$$\beta_1 = \mathbb{E}[Y|X=x+1] - \mathbb{E}[Y|X=x]$$

$$\approx \overline{Y|X=x+1} - \overline{Y|X=x} \quad \text{for some } x$$

Won't work so well if we only have one value for each  $x$

Covariance (Def 3.3.3)

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)]$$

$$= \mathbb{E}[XY] - \mu_X \mu_Y \quad (\text{unit of } X \cdot Y)$$

$\uparrow$  Theorem 3.3.3

lets calculate  $\text{Cov}(X, Y)$  for linear regression model

need to compute  $E[XY]$ . Can either tower property (HW 3B)

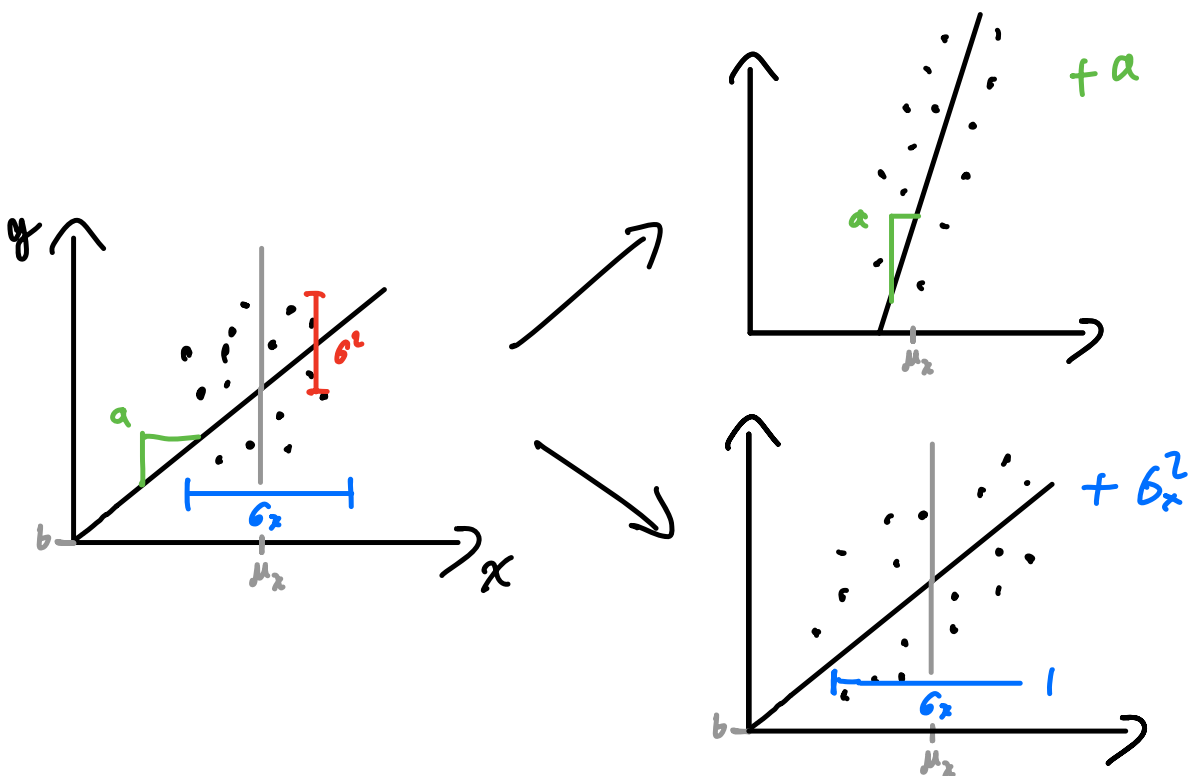
or

$$\begin{aligned}
 E[XY] &= E[X(\beta_0 + \beta_1 X + \varepsilon)] \\
 &= E[\beta_0 X] + \beta_1 E[X^2] + E[X\varepsilon] \\
 &= \beta_0 \mu_X + \beta_1 (\text{Var}(X) + \mu_X^2) + E[X]E[\varepsilon] \\
 &= \beta_0 \mu_X + \beta_1 \sigma_X^2 + \beta_1 \mu_X^2
 \end{aligned}$$

$$\begin{aligned}
 E[X]E[Y] &= \mu_X E[\beta_0 + \beta_1 X + \varepsilon] \\
 &= \mu_X \beta_0 + \mu_X^2 \beta_1
 \end{aligned}$$

$$\Rightarrow E[XY] - E[X]E[Y] = \sigma_X^2 \beta_1 \quad (\star)$$

check units



Formula for covariance gives us a way to express slope from samples:

Samples  $(X_1, Y_1), \dots, (X_N, Y_N)$

$$\text{Var } \sigma_X^2 = E[(X - E[X])^2] \\ \approx \overline{(X - \bar{X})^2} \approx \frac{1}{N} \sum_{i=1}^N (X_i - \bar{X})^2$$

technically should be done on the later

$$E[(X - \mu_X)(Y - \mu_Y)] \approx \frac{1}{N} \sum_{i=1}^N (X_i - \bar{X})(Y_i - \bar{Y})$$

$$\Rightarrow \hat{\beta}_1 = \frac{\sum_{i=1}^N (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^N (X_i - \bar{X})^2}$$

Next, we can estimate  $\beta_0$ :

$$\beta_0 = E[Y] - \beta_1 \mu_X$$

$$\approx \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

important!

these are formula in theorem 12.3.1 derived differently

We can also estimate  $\sigma$  by

$$\hat{\sigma} = \frac{(\bar{Y} - (\hat{\beta}_0 + \bar{X} \hat{\beta}_1))^2}{n-2}$$