

# Linear Regression Models

## P8111

### Lecture 04

Jeff Goldsmith  
January 28, 2016



THE DEPARTMENT OF  
**BIostatISTICS**



Columbia University  
**MAILMAN SCHOOL  
OF PUBLIC HEALTH**

# Today's lecture

- Simple Linear Regression
- Least Squares Estimation



# Regression modeling

$$\underline{E(y|x)} = ???$$

- Want to use predictors to learn about the outcome distribution, particularly conditional expected value.
- Formulate the problem parametrically

$$E(y | x) = f(x; \beta) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

- (Note that other useful quantities, like covariance and correlation, tell you about the joint distribution of  $y$  and  $x$ )

## Covariance and Correlation

$$\text{COV}(X, Y) = E((X - \mu_X)(Y - \mu_Y))$$

$$\hat{\text{COV}}(X, Y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n}$$

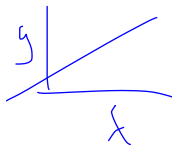
$$\text{COR}(X, Y) = \frac{\text{COV}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$$

# Simple linear regression

- Linear models are a special case of all regression models; simple linear regression is the simplest place to start
- Only one predictor:

$$\underline{E(y | x)} = f(x; \beta) = \underline{\beta_0 + \beta_1 x_1}$$

*(Handwritten blue annotations: a blue '1' above the plus sign, a blue arrow pointing down to the plus sign, and a blue underline under the entire right-hand side of the equation.)*

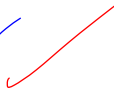


- Useful to note that  $x_0 = 1$  (implicit definition)
- Somehow, estimate  $\beta_0, \beta_1$  using observed data.

## Coefficient interpretation

$$E(y|x) = \beta_0 + \beta_1 x$$

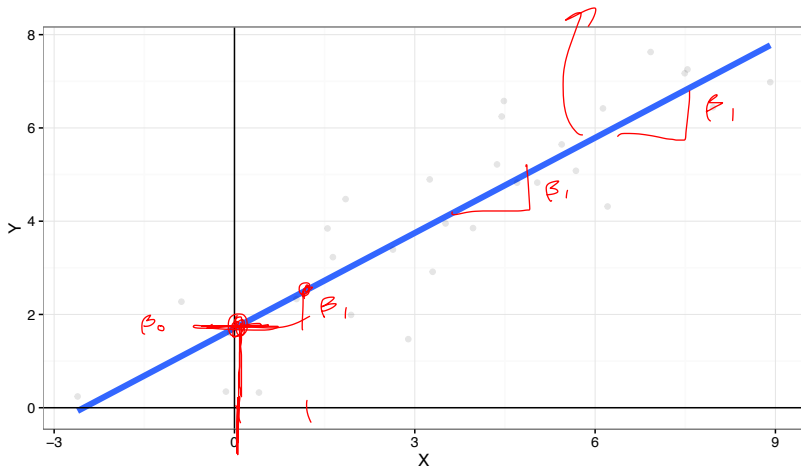
$$\underline{E(y|x=0) = \beta_0}$$



$$\begin{aligned}\beta_1 &= (\beta_0 + \beta_1 \textcolor{red}{17}) - \beta_0 - \beta_1 \textcolor{red}{16} \\ &= E(y|x=\textcolor{red}{1}) - E(y|x=\textcolor{red}{0})\end{aligned}$$

# Coefficient interpretation

$$E(y|x) = \beta_0 + \beta_1 x$$



# Look at the data

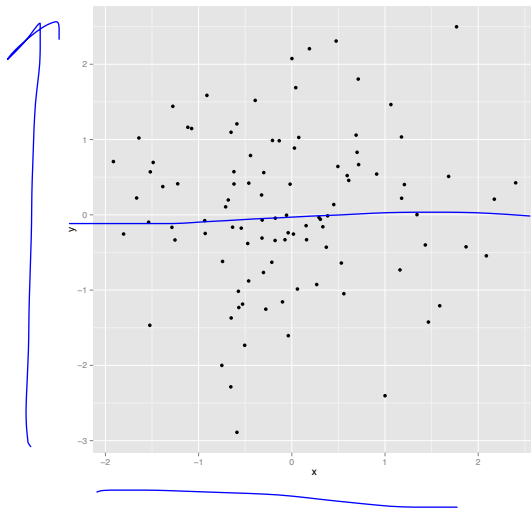
## • Before Modeling

- Plot the data (using ggplot ...)
- Do the data look like the assumed model? ✓
- ✓ ■ Should you be concerned about outliers? ✓
- ✓ ■ Define what you expect to see before fitting any model.



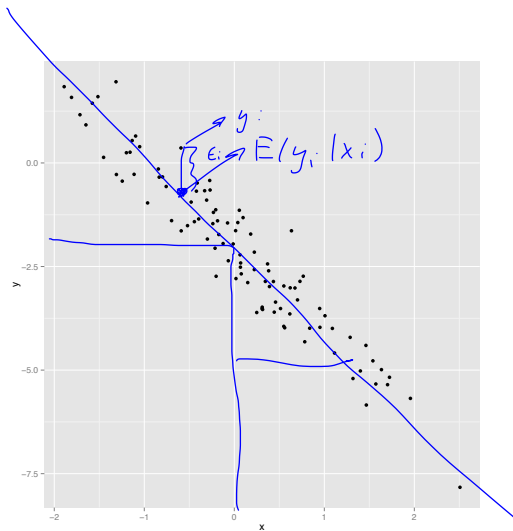
# Look at the data

$$E(y|x) = \beta_0 + \beta_1 x$$

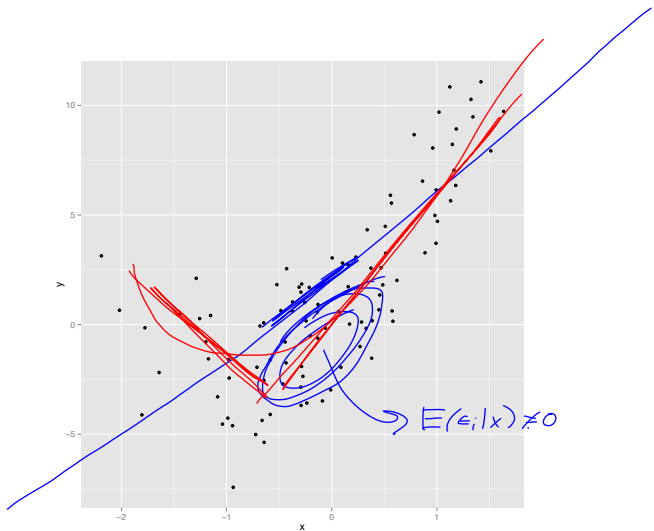


$$\begin{aligned} \beta_0 &= 0 \\ \hat{\beta}_1 &\approx 0 \end{aligned}$$

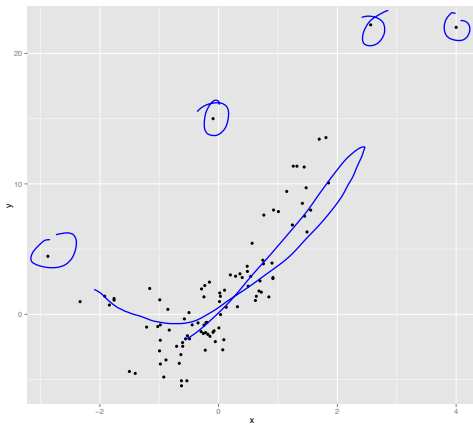
# Look at the data



# Look at the data



# Look at the data



# Least squares estimation

$$E(y|x)$$

- Observe data  $(y_i, x_i)$  for subjects  $1, \dots, n$ . Want to estimate  $\beta_0, \beta_1$  in the model

$$y_i = E(y_i | x_i) + \epsilon_i; \epsilon_i \sim N(0, \sigma^2)$$

$y_i = \beta_0 + \beta_1 x_i + \epsilon_i; \epsilon_i \stackrel{iid}{\sim} (0, \sigma^2)$

- Note the assumptions on the variance:

- $E(\epsilon | x) = E(\epsilon) = 0$

- Constant variance

- Independence

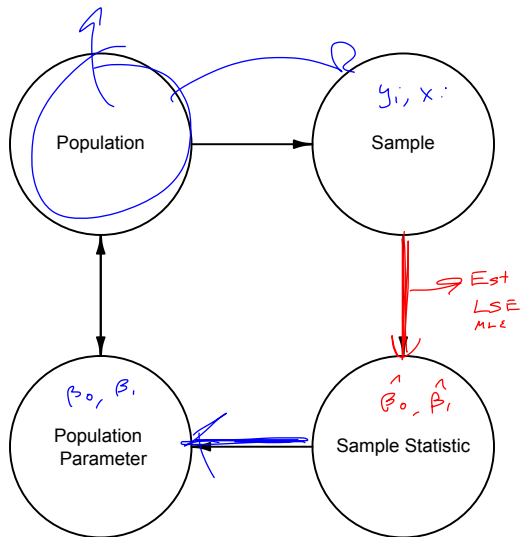
- [Normally distributed is not needed for least squares, but is needed for inference]

nice / useful

needed for MLE



# Circle of Life



# Least squares estimation

- Recall that for a single sample  $y_i, i \in 1, \dots, n$ , the sample mean  $\hat{\mu}_y$  minimizes the sum of squared deviations.

$$RSS(\mu_y) = \sum (y_i - \mu_y)^2 \quad \checkmark \quad \xrightarrow{\quad} \quad E(y)$$

$$E(y|x) = \beta_0 = \mu_y$$

$$\frac{\partial RSS(\mu_y)}{\partial \mu_y} = -2 \sum (y_i - \mu_y) = 0$$

$$\sum y_i - 2\mu_y = 0$$

$$n\mu_y = \sum y_i$$

$$\hat{\mu}_y = \frac{\sum y_i}{n}$$

# Least squares estimation

- Find  $\hat{\beta}_0$ .

$$E(y|x) = \beta_0 + \beta_1 x$$

$$RSS(\beta_0, \beta_1) = \sum (y_i - \beta_0 - \beta_1 x_i)^2$$

$$\frac{\partial RSS(\beta_0)}{\partial \beta_0} = -2 \sum (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\sum y_i - n\beta_0 - \beta_1 \sum x_i = 0$$

$$n\beta_0 = \sum y_i - \beta_1 \sum x_i$$

$$\hat{\beta}_0 = \bar{y} - \beta_1 \bar{x}$$

↑



# Least squares estimation


- Now find  $\hat{\beta}_1$ .

$$\begin{aligned} RSS(\beta_1) &= \sum (y_i - \overbrace{\bar{y} + \beta_1 \bar{x}}^{\hat{\beta}_0} - \beta_1 x_i)^2 \\ &= \sum (\underbrace{y_i - \bar{y}} - \beta_1 (\underbrace{x_i - \bar{x}})) \end{aligned}$$

$$\begin{aligned} \frac{\partial RSS(\beta_1)}{\partial \beta_1} &= -2 \sum (y_i - \bar{y}) - \beta_1 \sum (x_i - \bar{x}) = 0 \\ &= \sum (y_i - \bar{y})(x_i - \bar{x}) - \beta_1 \sum (x_i - \bar{x})^2 = 0 \end{aligned}$$

$$\hat{\beta}_1 = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

# Note about correlation


$$\rho = \frac{\text{cov}(x, y)}{\sqrt{\text{var}(x)\text{var}(y)}}; \quad \beta_1 = \frac{\text{cov}(x, y)}{\text{var}(x)}$$

# R does exactly what we now expect

*> data =*

*> linmod = lm(y~x, data = data)*  
*> summary(linmod)*

Call:  
lm(formula = y ~ x, data = data)

Residuals:

Min	1Q	Median	3Q	Max
-1.5202	-0.5050	-0.2297	0.5753	1.8534

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	2.08743	0.22958	9.092	7.53e-10 ***
x	0.61396	0.05415	11.338	5.61e-12 ***

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.8084 on 28 degrees of freedom

Multiple R-squared: 0.8211, Adjusted R-squared: 0.8148

F-statistic: 128.6 on 1 and 28 DF, p-value: 5.612e-12

# R does exactly what we now expect

library(broom)

```
> tidy(linmod)
  term      estimate std.error statistic    p.value
1 (Intercept) 2.0874344 0.22958105  9.092364 7.529711e-10
2 x          0.6139621 0.05415004 11.338166 5.611585e-12
> glance(linmod)
  r.squared adj.r.squared    sigma statistic    p.value df logLik ...
1 0.821148    0.8147604 0.8084399   128.554 5.611585e-12  2 -35.1538 ...
>
> beta1 = with(data, sum((x - mean(x)) * (y - mean(y))) / sum((x - mean(x))^2))
> beta0 = with(data, mean(y) - beta1 * mean(x))
> c(beta0, beta1)
[1] 2.0874344 0.6139621
```

*Handwritten annotations:*

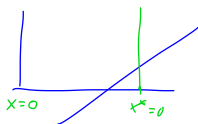
- Blue box around the `tidy(linmod)` output table.
- Red circle around the `beta0` and `beta1` calculations and their result.
- Blue arrows pointing from `data` to `data$x` and `data$y` in the `beta1` calculation.
- Blue arrow pointing from `attach(data)` to the `data` argument in the `beta1` calculation.

# Note on interpretation of $\beta_0$

Recall  $\beta_0 = E(y|x = 0)$

- This often makes no sense in context ✓
- “Centering”  $x$  can be useful:  $x^* = x - \bar{x}$
- Center by mean, median, minimum, etc
- Effect of centering on slope:

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

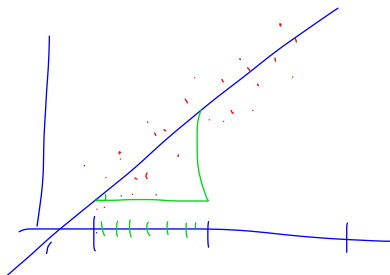


# Note on interpretation of $\beta_0, \beta_1$

- The interpretations are sensitive to the scale of the outcome and predictors (in reasonable ways)
- You can't get a better model fit by rescaling variables ✓

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad X^* = cX$$

$$\hat{\beta}_1^* = \frac{1}{c} \hat{\beta}_1$$



# R example

```
> data = mutate(data, x.cen = x - mean(x), x2 = x*2)
> linmod.cen = lm(y ~ x.cen, data = data)
> tidy(linmod.cen)
```

	term	estimate	std.error	statistic	p.value
1	(Intercept)	4.0811993	0.14760027	27.65035	7.172437e-22
2	x.cen	0.6139621	0.05415004	11.33817	5.611585e-12

x.scaled.by.2



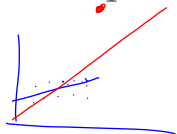
# R example

```
> linmod.x2 = lm(y ~ x2, data = data)
> tidy(linmod.x2)
```

	term	estimate	std.error	statistic	p.value
1	(Intercept)	<u>2.0874344</u>	0.22958105	<u>9.092364</u>	7.529711e-10
2	x2	<u>0.3069811</u>	0.02707502	<u>11.338166</u>	<u>5.611585e-12</u>



# Least squares notes and foreshadowing



$$RSS(\beta_0, \beta_1) = \sum (y_i - E(y_i|x_i))^2 \rightarrow \text{sample mean}$$

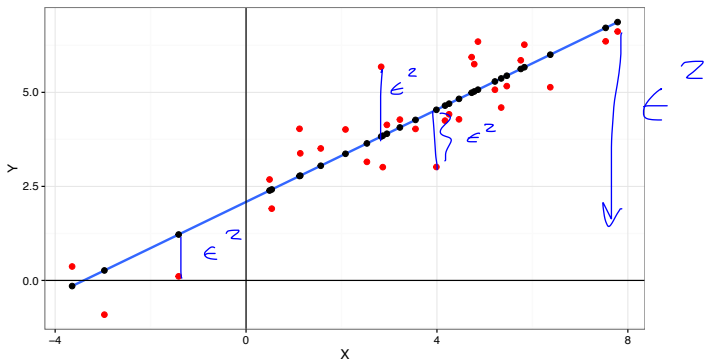
$$\sum |y_i - E(y_i|x_i)| \rightarrow \text{sample median}$$

- Didn't have to choose to minimize squares – could minimize absolute value, for instance.
- Least squares estimates turn out to be a “good idea” – unbiased, BLUE.
- Later we'll see about maximum likelihood as well.

# Geometric interpretation of least squares

Least squares minimizes the sum of squared vertical distances between observed and estimated  $y$ 's:

$$\min_{\beta_0, \beta_1} \sum_{i=1}^I (y_i - (\beta_0 + \beta_1 x_i))^2$$



# Least squares in regression generally

Broadly speaking, in regression we often are concerned with minimizing

$$E[\underbrace{f(x) + \epsilon - \hat{f}(x)}^{\text{Bias}}]^2$$

by choosing a “good”  $\hat{f}$ . For a given  $\hat{f}$  this decomposes into

$$\underbrace{E[f(x) - \hat{f}(x)]^2}_{\text{Bias}^2} + \text{Var}(\epsilon)$$

$\uparrow$                        $\uparrow$                        $\nwarrow$

- Some variance isn't explainable (we just don't know how much)
- Focus on getting the left component right
- Minimizing squared error for *unseen* data is the real goal

# Today's big ideas

- Simple linear regression – model and interpretation
- Least squares estimation

- 
- Suggested reading: Faraway Ch 1, 2.1; ISLR 3.1