CSCI 3022-002 Intro to Data Science

Regression \Box

4=mx+b

Consider the graph below. Can you draw a straight line that passes through each and every

point? What about a parabola? 20

(a can be/not be o)

Mullen: SLR Fall 2020

Mathematical Interpolation

These problems are known as *interpolation*. We actually don't care about them! As statisticians we ask a different problem: how can we *best* draw a line through data values if the process that's generated them is subject to **error**?

In that sense, what we do in data science is we take data that's highly variable - or noisy - and try to describe it with a simpler process. "Simpler" here means a couple of things: relying on less parameters/estimations, and nicer mathematically. The latter is often referred to as **smooth**, and it should come as no surprise that data scientists are **smooth** to the core.

Regression is the way that we smooth out data: take something that's noisy and has *lots* of observations and try to describe it with a single line!

Mullen: SLR Fall 2020 2 / 35

Announcements and Reminders

- ► Short HW for next week. (posted w?)
- ► NB day Friday.

Mullen: SLR Fall 2020

Where we at?

We're "done" with CI's and hypothesis testing (as standalones...)! If:

- 1. We're estimating means:
 - 1.1 Underlying normal: use t) if both: use t. (≈ 2)
 - 1.2 Large sample: use z
 - 1.3 Small sample, no assumptions: bootstrap.
- 2. We're not estimating means:
 - 2.1 Underlying normal: we can use χ^2 to make Cls/tests on variances for one sample or F to make Cls/tests on comparisons of variances. (Not done, we'll talk about these later!)
 - 2.2 Otherwise, we bootstrap!

Mullen: SLR Fall 2020

Bootstrapping Recap

Bootstrapping is an attempt to bring your data set from rags to riches.

When we aren't making assumptions like "underlying normal," we mean that **shape matters**. Bootstrapping is a way to understand how much the shape of a (maybe small) sample matters!

- 1. Take subsamples with replacement! of your original sample
- 2. Compute descriptive statistics from your subsample: anything from means, variances, maxima, to .72nd quantiles, and use those to get a feel for how those parameters behave from the underlying population!
- 3. **Example**: what's your best guess for the mean rainfall of Boulder in the last 10 years? Variance?

Bootstrapping

This process: simulating a data set, calculating a desired *sample statistic* from it, and then creating a *distribution* of that sample statistic is called a *non-parametric* bootstrap since it doesn't make distributional assumptions.

Definition: parametric statistics assume that sample data comes from a population that follows a probability distribution on a fixed set of parameters.

Examples:

- 1. μ and σ are the parameters of the Normal distribution.
- 2. λ is the parameters of the Poisson and Exponential distributions.
- 3. p is the parameter of the geometric and Bernoulli distributions.

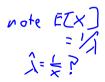
Mullen: SLR Fall 2020

Parametric Bootstrapping

Sometimes we really want to know about various statistics on e.g. the Poisson or Exponential without solving some challenging integral or sum or whatever else equations.

Definition: parametric bootstraps estimate a CI for a desired property in two steps.

- 1. Estimate the parameters of the known distribution from your sample. estimate
- 2. Draw bootstrap resamples from the distribution, assuming the estimated parameter $\sum_{i=1}^{n} \frac{1}{n} \frac{1}$
- 3. Compute a CI for the desired property from your resamples.



Mullen: SLR Fall 2020

Parametric Bootstrapping

Example: If we want to estimate the median of a sample that we *assume* is Poisson, we might:

- 1. Assume the data is $\operatorname{Pois}(\lambda)$. Estimate the parameter, e.g. $\lambda \approx \bar{X}$.

 2. Simulate a bootstrapped sample from $\operatorname{Pois}(\bar{X})$.
- 3. Create a CI for the median from that pool of bootstrapped samples.

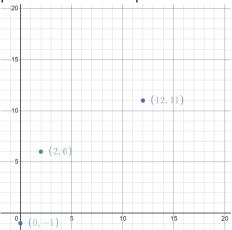
Why make *more* assumptions, like assuming the distribution of the random variable at all? The advantage of the parametric bootstrap is that it can be shown to do a better job in particular scenarios. (if we're 1:4+)

The downside? The parametric bootstrap does a very poor job if the population does not have the same population as you assumed. This is called *model misspecification*, and is a risk **any** time we assume things have **any** underlying distribution, including in hypothesis testing!

Mullen: SLR Fall 2020 8 / 35

Opening

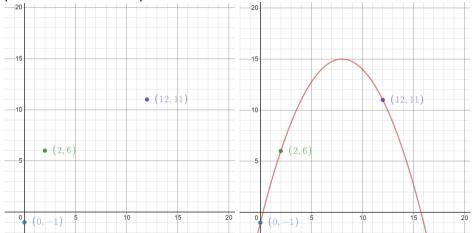
Consider the graph below. Can you draw a straight line that passes through each and every point? What about a parabola?



:: SLR Fall 2020 9 / 35

Opening

Consider the graph below. Can you draw a straight line that passes through each and every point? What about a parabola?



Fall 2020 9/35

Mathematical Interpolation

These problems are known as *interpolation*. We actually don't care about them! As statisticians we ask a different problem: how can we *best* draw a line through data values if the process that's generated them is subject to **error**?

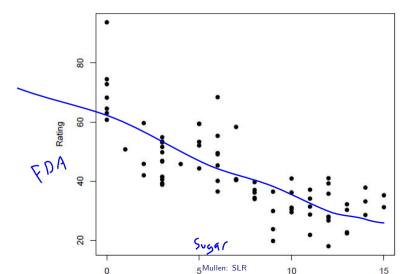
In that sense, what we do in data science is we take data that's highly variable - or noisy - and try to describe it with a simpler process. "Simpler" here means a couple of things: relying on less parameters/estimations, and nicer mathematically. The latter is often referred to as **smooth**, and it should come as no surprise that data scientists are **smooth** to the core.

Regression is the way that we smooth out data: take something that's noisy and has *lots* of observations and try to describe it with a single line!

Mullen: SLR Fall 2020

Simple Linear Regression

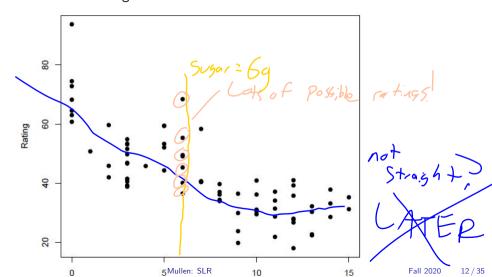
We see bivariate data: a set of points (X_i, Y_i) .



Fall 2020 11 / 35

Simple Linear Regression

...and we draw the "best" line through them.



SLR

Simple Linear Regression Model

Simple Linear Regression Model

The simplest mathematical relationship between two variables x and y is a linear relationship:

Mullen: SLR Fall 2020 13 / 35

The simplest mathematical relationship between two variables x and y is a linear relationship:

$$y = mx + b$$

The objective of this section is about equivalent linear probabilistic models.

Mullen: SLR Fall 2020

The simplest mathematical relationship between two variables x and y is a linear relationship:

$$y = mx + b$$

The objective of this section is about equivalent linear probabilistic models.

If two random variables are probabilistically related, then for a fixed value of x, there is uncertainty in the value of the second variable.

The simplest mathematical relationship between two variables x and y is a linear relationship:

$$y = mx + b$$

The objective of this section is about equivalent linear probabilistic models.

If two random variables are probabilistically related, then for a fixed value of x, there is uncertainty in the value of the second variable.

So we assume where ε is a random variable and:

ariable and:
$$MODE$$

$$y = \beta_0 + \beta_1 x + \varepsilon$$

$$y = b + m \times + error / Probab: |; + y|$$

Mullen: SLR Fall 2020

Definition: Simple Linear Regression (SLR)

With 3 assumptions on ε :

Mullen: SLR Fall 2020 14 / 35

depends on X

Definition: Simple Linear Regression (SLR)

The Simple Linear Regression model is a model of the form

1.

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

With 3 assumptions on ε :

Mullen: SLR Fall 2020 14 / 35

Definition: Simple Linear Regression (SLR)

The Simple Linear Regression model is a model of the form

1.

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

With 3 assumptions on ε :

2.

Independence of errors

$$\begin{aligned} & (\text{ovariance} := \text{O}) \text{ independence} \\ & [\text{Cov}[\varepsilon_i, \varepsilon_j] = 0 \quad \forall i, j \quad \text{il } \text{knowledge of} \\ & \text{error of observation in} \\ & \text{telly us nothing about the} \\ & \text{error of observation j} \end{aligned}$$
 Result: Data Probably not $X = t$ ine

Mullen: SLR Fall 2020 14 / 35

Definition: Simple Linear Regression (SLR)

The Simple Linear Regression model is a model of the form

1.

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

With 3 assumptions on ε :

 $\mathsf{Cov}[\varepsilon_i, \varepsilon_j] = 0 \quad \forall i, j$

Independence of errors

Homoskedasticity of errors

$$Var(\varepsilon_i) = \sigma^2$$
 $\forall i$ their variance.

Fall 2020 14 / 35

Definition: Simple Linear Regression (SLR)

The Simple Linear Regression model is a model of the form

1.) LINE IS CORRECT
$$y_i = eta_0 + eta_1 x_i + arepsilon_i$$

With 3 assumptions on ε :

With 3 assumptions on
$$\varepsilon$$

$$\mathsf{Cov}[\varepsilon_i, \varepsilon_j] = 0 \qquad \forall i, j$$

Independence of errors

$$Var(\varepsilon_i) = \sigma^2 \qquad \forall i$$

Homoskedasticity of errors

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

(0,X) Rine is correct mode)

Important Terminology:

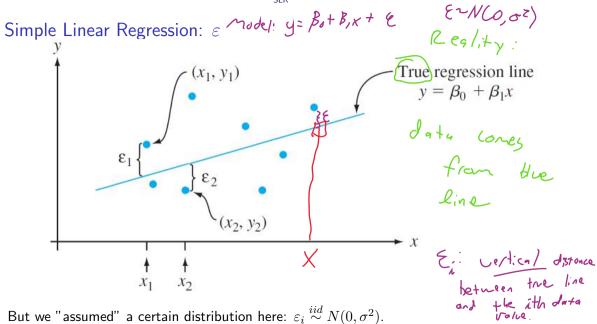
 \triangleright x: the independent variable, predictor, or explanatory variable (usually known). x is not random.

 \triangleright Y: The dependent variable or response variable. For fixed x, Y is random.

 \triangleright ε : The random deviation or random error term. For fixed x, ε is random.

What exactly does ε do?

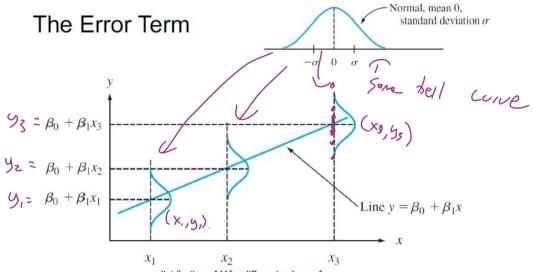
Mullen: SLR Fall 2020



SLR

But we "assumed" a certain distribution here: $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$. Mullen: SLR Fall 2020 16 / 35

Simple Linear Regression: ε



distribution of Y for different values of x

Mullen: SLR

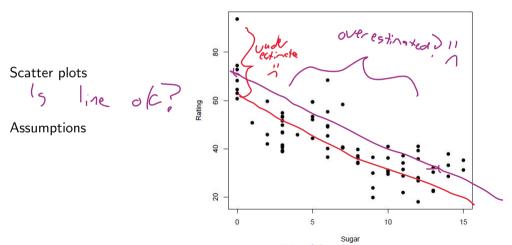
In practice, using a linear model largely comes down to diagnosing and fixing deviations from these assumptions. What does that look like?

- 1. Linearity: Is a straight line actually appropriate? $5 = B_0 + B_1 \times 0 \neq 0$
- 2. Independence: Are there patterns or structures in my data not covered by a straight line? Clumps of overluder estimates.
- 3. **Homoskedasticity:** Do the sizes of my errors vary depending on where I am in the data set? Larger x, larger y?
- 4. **Normality:** Do the errors appear to be coming from a normal distribution?

Mullen: SLR Fall 2020

Simple Linear Regression: Assumptions

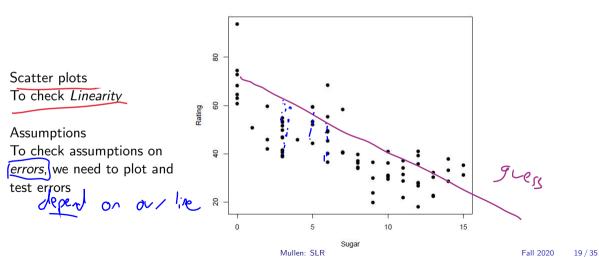
How do we know when the SLR model is appropriate?



Mullen: SLR Fall 2020

Simple Linear Regression: Assumptions

How do we know when the SLR model is appropriate?



Simple Linear Regression: Interpretations

Interpreting parameters:

y = Bo+ Bix terror \triangleright β_0 : the intercept of the true regression line. The average value of Y when x is zero. Usually this is called the "baseline average".

 \triangleright β_1 : the slope of the true regression line. The average change in Y associated with a 1-unit increase in the value of x.

Fall 2020

Estimating SLR Parameters

Goal:

Given sample data, which consists of n observed pairs, $(x_1, y_1), (x_2, y_2), \dots (x_n, y_n)$, construct an estimated "line of best fit":

This line can then be used to make predictions or provide explanations for unobserved phenomena.

How do we construct this line?

Mullen: SLR Fall 2020

Estimating SLR Parameters

Goal:

Given sample data, which consists of n observed pairs, $(x_1, y_1), (x_2, y_2), \dots (x_n, y_n)$, construct an estimated "line of best fit":

$$y = \hat{\beta}_0 + \hat{\beta}_1 x \quad \text{our estimated slope}$$

This line can then be used to make predictions or provide explanations for unobserved phenomena.

How do we construct this line? Find estimates of the β values with nice statistical properties, and call them $\hat{\beta}$!

Mullen: SLR Fall 2020 21 / 35

SLR Parameters

What are the properties of our ith data point, Y_i ?

Consider $E[Y_i]$:

E[y:] = E[
$$B_0 + B_1 X_1 + E_2$$
] = $B_0 + B_1 X_2 + E[E]$

Consider $Var[Y_i]$:

What if we got a new data point, like $X_{new} = X_i + 1$?

Fall 2020

SLR Parameters

What are the properties of our ith data point, Y_i ?

Consider $E[Y_i]$:

$$E[Y_i] = E[\beta_0 + \beta_1 X_i + \varepsilon] = \beta_0 + \beta_1 X_i + E[\varepsilon] = \beta_0 + \beta_1 X_i$$

Consider $Var[Y_i]$:

$$Var[Y_i] = Var[\beta_0 + \beta_1 X_i + \varepsilon] = Var[\varepsilon] = \sigma^2$$

What if we got a new data point, like $X_{new} = X_i + 1$?

$$E[Y_{new}] = E[\beta_0 + \beta_1 X_{new} + \varepsilon] = \beta_0 + \beta_1 X_{new} + E[\varepsilon] = \beta_0 + \beta_1 (X_i + 1) = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} + \frac{\beta_0 + \beta_1 X_i}{\beta_0 + \beta_1 X_i} = \frac{\beta_0 +$$

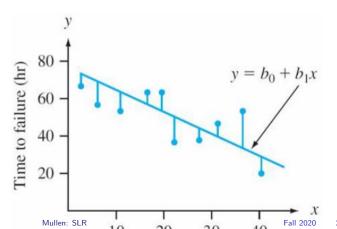
So we interpret β_1 as a regular slope: it's the expected change in Y per 1 unit change in X.

> Mullen: SLR Fall 2020 22 / 35

Estimating SLR Parameters

One way to define "best fit" line is motivated by the principle of least squares, which can be traced back to the German mathematician Gauss (1777–1855):

A line provides the **best fit** to the data if the sum of the squared vertical distances (deviations) from the observed points to that line is as small as it can be.



Estimating SLR Parameters

The sum of squared vertical deviations from the data points to the line $y = \beta_0 + \beta_1 x$ is then



Mullen: SLR Fall 2020 24 / 35

Estimating SLR Parameters "best" $\lim_{n \to \infty} \hat{\beta} + \hat{\beta} \chi$

The sum of squared vertical deviations from the data points to the line $y=\beta_0+\beta_1 x$ is then $\mathcal{C}(x): \qquad \mathcal{C} = (\hat{\mathcal{C}}_0 + \hat{\mathcal{C}}_1 \times \hat{\mathcal{C}}_2)$

Fall 2020 24 / 35

The sum of squared vertical deviations from the data points to the line $y = \beta_0 + \beta_1 x$ is then

$$\sum_{i=1}^{n} \left(\overbrace{Y_i}^{\mathsf{Data}} - \overbrace{\beta_0 - \beta_1 X_i}^{\mathsf{Line}} \right)^2$$

The point estimates of β_0 and β_1 , denoted $\hat{\beta}_0$; $\hat{\beta}_1$ are called the *least squares estimates*. They are those values that minimize SSE or sum of squared errors.

> Fall 2020 24 / 35

The sum of squared vertical deviations from the data points to the line $y = \beta_0 + \beta_1 x$ is then

$$\sum_{i=1}^n \left(\underbrace{Y_i}_{\text{Data}} - \underbrace{\beta_0 - \beta_1 X_i}_{\text{Data}} \right)^2$$

The point estimates of β_0 and β_1 , denoted $\hat{\beta}_0$; $\hat{\beta}_1$ are called the *least squares estimates*. They are those values that minimize SSE or sum of squared errors.

$$\hat{\beta}_0, \hat{\beta}_1 = \underset{\beta_0, \beta_1}{\operatorname{argmin}} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$$

Mullen: SLR 24 / 35

Estimating SLR Parameters: Pen and Paper

Goal: find the minimizers of the function $f(\underline{\beta_0, \beta_1}) = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$. Sounds like a Calculus problem!

Mullen: SLR Fall 2020 25 / 35

Estimating SLR Parameters: Pen and Paper

Goal: find the minimizers of the function $f(\beta_0, \beta_1) = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$. Sounds like a Calculus problem!

$$\frac{df}{d\beta_0} = \frac{d}{d\beta_0} \sum_{i=1}^{n} (Y_i - \beta_0)$$

$$\frac{df}{d\beta_0} = \frac{d}{d\beta_0} \sum_{i=1}^{n} (Y_i - \beta_0 - \beta_1 X_i)^2 = \bigcirc$$

$$\frac{df}{d\beta_1} = \frac{d}{d\beta_1} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$$

25 / 35

Mullen: SLR Fall 2020

Estimating SLR Parameters: Pen and Paper

Goal: find the minimizers of the function $f(\beta_0, \beta_1) = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$. Sounds like a Calculus problem!

$$\frac{df}{d\beta_0} = \frac{d}{d\beta_0} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$$

$$\frac{df}{d\beta_1} = \frac{d}{d\beta_1} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$$

For finding the joint maximum/minimum of multiple inputs, we end up with a system of equations: set both equal to zero and find the values that make both equal to zero.

Mullen: SLR Fall 2020 25 / 35

Estimating SLR Parameters: Results

For a model of the form $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$; $\varepsilon \sim N(0, \sigma^2)$

1.
$$\hat{\beta_0} =$$

2.
$$\hat{\beta}_1 =$$

What happens if $\beta_0 \approx 0$? If $\beta_1 \approx 0$?

Mullen: SLR Fall 2020

Estimating SLR Parameters: Results

For a model of the form $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$; $\varepsilon \sim N(0, \sigma^2)$

1.
$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

2.
$$\hat{\beta}_1 = \frac{Cov[X,Y]}{\sum_{i=1}^n (X_i - \bar{X})^2} = \underbrace{\left[\frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}\right]} \int \int d\mathbf{p} d\mathbf{r}$$

What happens if $\beta_0 \approx 0$? If $\beta_1 \approx 0$?

Mullen: SLR Fall 2020

Estimating SLR Parameters: Results

For a model of the form $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$; $\varepsilon \sim N(0, \sigma^2)$

$$1. \ \hat{\beta_0} = \bar{Y} - \hat{\beta_1} \bar{X}$$

2.
$$\hat{\beta}_1 = \frac{Cov[X,Y]}{\sum_{i=1}^n (X_i - \bar{X})^2} = \underbrace{\left(\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}_{\sum_{i=1}^n (X_i - \bar{X})^2}\right)}$$

What happens if
$$\beta_0 \approx 0$$
? If $\beta_1 \approx 0$?

 $2. \ \hat{\beta_1} = \frac{Cov[X,Y]}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - Y)}{\sum_{i=1}^n (X_i - \bar{X})^2} \text{ Numerator lacks like} \\ \text{Cov}(\chi, \chi) \cdot + ; f \times \chi \text{ note} \\ \text{Tope the Volume} \\ \text{One result: the regression line goes through } (0, \beta_0). \ \text{It also goes through } (\bar{X}, \bar{Y})! \\ \text{aportions} \\ \text{a$

Mullen: SLR

Definitions:

1. The *fitted (or predicted) values* __ are obtained by plugging in __ to the equation of the estimated regression line:

2. The *residuals* are the differences between the observed and fitted y values:

Residuals are estimates of the true error. Why?

Mullen: SLR Fall 2020 27 / 35

Definitions:

1. The fitted (or predicted) values $\underline{\hat{Y}_i}$ are obtained by plugging in $\underline{\hat{X}_i}$ to the equation of the estimated regression line:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

2. The *residuals* are the differences between the observed and fitted y values:

Residuals are estimates of the true error. Why?

Mullen: SLR Fall 2020

Definitions:

1. The *fitted (or predicted) values* __ are obtained by plugging in __ to the equation of the estimated regression line:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

2. The *residuals* are the differences between the observed and fitted y values:

$$\hat{\varepsilon_i} = r_i = \hat{e_I} = Y_i - \hat{Y_i} = Y_i - \hat{\beta_0} + be\hat{t}a_1X_i$$

Residuals are estimates of the true error. Why?

Mullen: SLR Fall 2020

Definitions:

1. The *fitted (or predicted) values* __ are obtained by plugging in __ to the equation of the estimated regression line:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

2. The *residuals* are the differences between the observed and fitted y values:

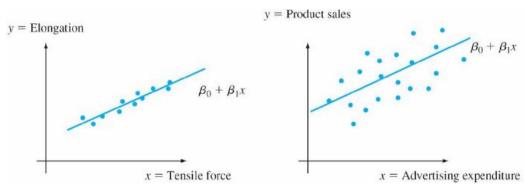
$$\hat{\varepsilon_i} = r_i = \hat{e_I} = Y_i - \hat{Y_i} = Y_i - \hat{\beta_0} + be\hat{t}a_1X_i$$

Residuals are estimates of the true error. Why?

We don't have the true values of β_0 , β_1 , so when we estimate them we get variance and error in our estimates.

Mullen: SLR Fall 2020 27 / 35

The parameter σ^2 determines the amount of spread about the true regression line. Two separate examples:



Mullen: SLR Fall 2020

An estimate of σ^2 will be used in confidence interval formulas and hypothesis testing procedures presented in the next days. Recall that the residual sum of squares or sum of squared errors (SSE) is:

SSE =

So, our estimate of the variance of the model is like a measure for an average of this summand:

Mullen: SLR Fall 2020

An estimate of σ^2 will be used in confidence interval formulas and hypothesis testing procedures presented in the next days. Recall that the residual sum of squares or sum of squared errors (SSE) is:

$$SSE = \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^{n} (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2$$

So, our estimate of the variance of the model is like a measure for an average of this summand:

Mullen: SLR Fall 2020

An estimate of σ^2 will be used in confidence interval formulas and hypothesis testing procedures presented in the next days. Recall that the residual sum of squares or sum of squared errors (SSE) is:

$$SSE = \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^{n} (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2$$

So, our estimate of the variance of the model is like a measure for an average of this summand:

$$\hat{\sigma^2} = \frac{SSE}{n-2}$$

Wait, what? Why the n-2??

Fall 2020

An estimate of σ^2 will be used in confidence interval formulas and hypothesis testing procedures presented in the next days. Recall that the residual sum of squares or sum of squared errors (SSE) is:

$$SSE = \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^{n} (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2$$

So, our estimate of the variance of the model is like a measure for an average of this summand:

$$\hat{\sigma^2} = \frac{SSE}{n-2}$$

Wait, what? Why the n-2?? These are again degrees of freedom.

> Mullen: SLR Fall 2020

Degrees of Freedom Intuition

Suppose you have 3 (random) points on the XY plane.

1. Can you draw a line through them?

2. Can you draw a parabola through them?

3. Can you draw a cubic function through them?

4. Can you draw a quartic function through them?

Mullen: SLR Fall 2020

Degrees of Freedom Intuition

Suppose you have 3 (random) points on the XY plane.

- 1. Can you draw a line through them? It's very unlikely. In fact, for truly random (normal) points, this result has probability zero!
- Can you draw a parabola through them?
 It's very unlikely. In fact, for truly random (normal) points, this result has probability zero! Yes, but there's only one such parabola.
- 3. Can you draw a cubic function through them? Yes. Not only that, you could choose any one of a,b,c,d in the $ax^3+bx^2+cx+d=0$ and then solve for the others. You have **one degree of freedom**.
- 4. Can you draw a quartic function through them? Yes. Not only that, you could choose any two of a,b,c,d,e in the $ax^4+bx^3+cx^2+dx+e=0$ and then solve for the others. You have **two degrees of freedom**.

Mullen: SLR Fall 2020 30 / 35

Degrees of Freedom

The takeaway?

One property of mathematical estimation: the more you estimate, the more you risk overfitting. In this model we've estimated **2** means $(\hat{\beta}_0, \hat{\beta}_1)$ before we got to σ , which "costs" us two degrees of freedom.

The more we estimate, the less options - degrees of freedom - we get for the remaining terms.

Mullen: SLR Fall 2020 31 / 35

Some properties of our estimate:

1. The divisor n-2 in is the number of degrees of freedom (df) associated with SSE and $\hat{\sigma}^2$.

2. This is because to obtain $\hat{\sigma}^2$, two parameters must first be estimated, which results in a loss of 2 df.

3. Replacing each y_i in the formula for $\hat{\sigma}^2$ by the r.v. Y_i gives a random variable.

4. It can be shown that the r.v. $\hat{\sigma}^2$ is an unbiased estimator for σ^2 .

Mullen: SLR Fall 2020

On Optimization

In any data science technique, there are two important considerations:

- 1. What are we optimizing? What are solving for and why?
- 2. How do we solve for that?
 - 2.1 Subsequent data science classes e.g. "Advanced Data Science," "Machine Learning," etc. involve a *lot* of algorithmic considerations: memory allocation, flop counts, etc.
 - 2.2 Do we have to approximate, or can we solve for an exact solution?

Mullen: SLR Fall 2020 33 / 35

Estimating SLR Parameters: the MLE

An alternative method for estimating model parameters is to create a likelihood function that quantifies the goodness-of-fit between the model and the data, and choose the values of the parameters that maximizes it

Turns out, we've done this before! But we didn't call it Maximum Likelihood Estimation at the time.

Example: Suppose you have a biased coin. You flip it 6 times, and get 5 Heads and 1 Tails. Estimate the parameter p for the coin.

Mullen: SLR Fall 2020 34 / 35

Example: Suppose you have a biased coin. You flip it 6 times, and get 5 Heads and 1 Tails. Estimate the parameter p for the coin.

We want the function that gives the probability of outcomes as a function of p. In this case, that's:

$$l(p) = P(\mathsf{data}\;\mathsf{GIVEN}\;\mathsf{P}) = P(\mathsf{5}\;\mathsf{heads}\;\mathsf{and}\;\mathsf{one}\;\mathsf{tails}|p)$$

Mullen: SLR Fall 2020

Example: Suppose you have a biased coin. You flip it 6 times, and get 5 Heads and 1 Tails. Estimate the parameter p for the coin.

We want the function that gives the probability of outcomes as a function of p. In this case, that's:

$$l(p) = P(\text{data GIVEN P}) = P(\text{5 heads and one tails}|p)$$

If we know p, this is a binomial, and the function is $l(p) = \binom{6}{5} p^5 (1-p) = 6p^5 (1-p)$. We want to find the value of p that maximizes this!

Mullen: SLR Fall 2020

Example: Suppose you have a biased coin. You flip it 6 times, and get 5 Heads and 1 Tails. Estimate the parameter p for the coin.

We want the function that gives the probability of outcomes as a function of p. In this case, that's:

$$l(p) = P(\text{data GIVEN P}) = P(\text{5 heads and one tails}|p)$$

If we know p, this is a binomial, and the function is $l(p) = \binom{6}{5} p^5 (1-p) = 6p^5 (1-p)$. We want to find the value of p that maximizes this!

$$\frac{dl}{dp} = 30p^4(1-p) - 6p^5 = p^4(30 - 30p - 6) = p^4(24 - 30p)$$

Mullen: SLR Fall 2020

Example: Suppose you have a biased coin. You flip it 6 times, and get 5 Heads and 1 Tails. Estimate the parameter p for the coin.

We want the function that gives the probability of outcomes as a function of p. In this case, that's:

$$l(p) = P(\text{data GIVEN P}) = P(\text{5 heads and one tails}|p)$$

If we know p, this is a binomial, and the function is $l(p) = \binom{6}{5} p^5 (1-p) = 6p^5 (1-p)$. We want to find the value of p that maximizes this!

$$\frac{dl}{dp} = 30p^4(1-p) - 6p^5 = p^4(30 - 30p - 6) = p^4(24 - 30p)$$

which equals zero at p = 0 and p = 5/6.

Mullen: SLR Fall 2020

Why care?

For any problem with underlying probability distributions - a pmf or a pmf - we can typically write down a likelihood function, which often reduces our data science problem to a numerical maximization problem.

For other problems, we may instead solve a least-squares or cost minimization problem. In either case, there's some metric by which we're coming up with the best solution.

For simple linear regression, they provide the same values of $\hat{\beta}$! This isn't always true. For example, the MLE for σ^2 of a normal data set is $\frac{\sum (X_i - \bar{X})^2}{n}$, which is a different denominator than our usual s^2 .

> Mullen: SLR Fall 2020

Daily Recap

Today we learned

1. Regression!

Moving forward:

- nb day Friday

Next time in lecture:

- More Regression!

Mullen: SLR Fall 2020 37 / 35