

CS/ECE 148 –

Data Science Fundamentals

Linear Regression

UCLA Computer Science

Lecture Outline

- Linear models
- Estimate of the regression coefficients
 - Brute Force
 - Exact method
 - Gradient Descent
- **Confidence intervals for the predictors estimates**
- Bootstrap
- Evaluating Significance of Predictors
- How well we know the model \hat{f}

Interpretation of Predictors

Question: What do you think a predictor coefficient means?

$$Sales = 7.5 + 0.04 TV$$

What does 7.5 mean and what does 0.04 mean?

If we increase the TV by \$1000, what would you expect the increase in sales to be?

What if?

$$Sales = 7.5 + 1.01 TV$$

The interpretation of the predictors depends on the values but decisions depend on how much we trust these values.

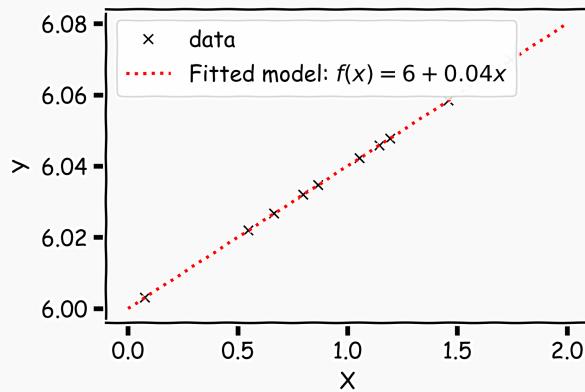
Confidence intervals for the predictors estimates

We interpret the ε term in our observation

$$y = f(x) + \varepsilon$$

to be noise introduced by random variations in natural systems or imprecisions of our scientific instruments.

If we knew the exact form of $f(x)$, for example, $f(x) = \beta_0 + \beta_1 x$, and there was no ε , then estimating the $\hat{\beta}$'s would have been exact (so is 1.01 worth it?).



Confidence intervals for the predictors estimates (cont)

However, three things happen, which result in mistrust of the values of $\hat{\beta}'$ s :

- ε is always there
- we do not know the exact form of $f(x)$
- limited sample size

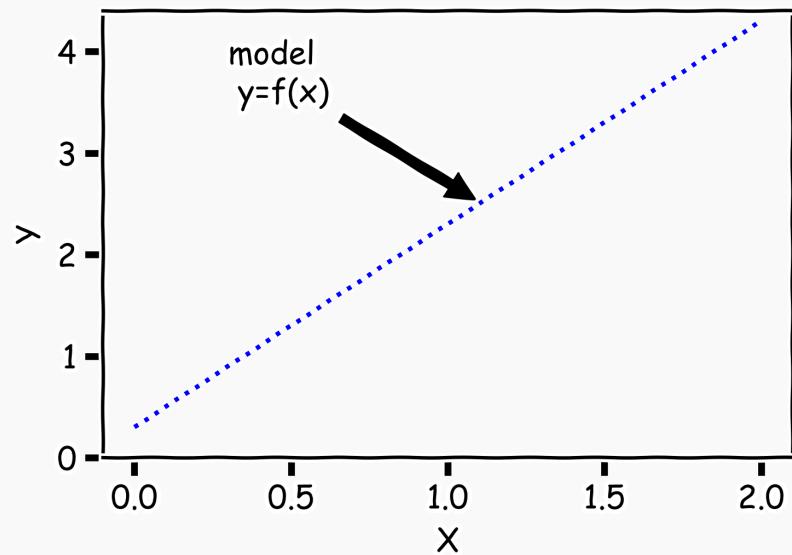
We will first address ε

We call ε the measurement error or ***irreducible error***. Since even predictions made with the actual function f will not match observed values of y .

Because of ε , every time we measure the response Y for a fix value of X , we will obtain a different observation, and hence a different estimate of $\hat{\beta}'$ s.

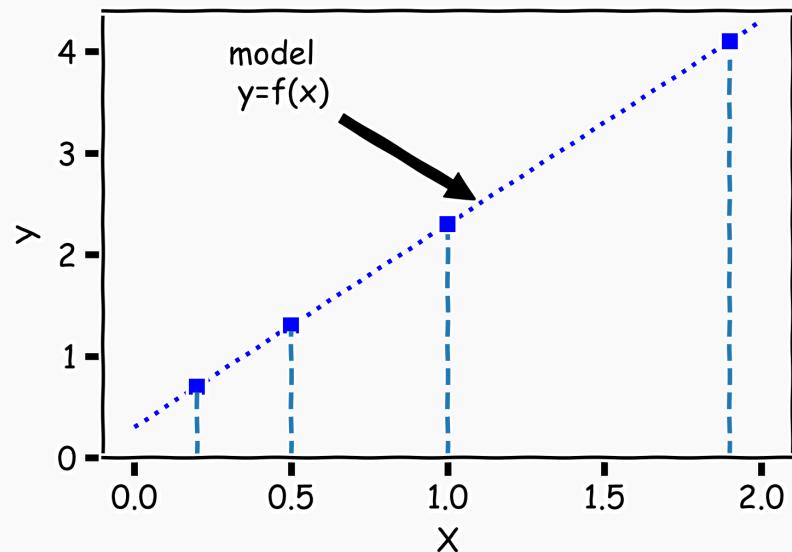
Confidence intervals for the predictors estimates (cont)

Start with a model



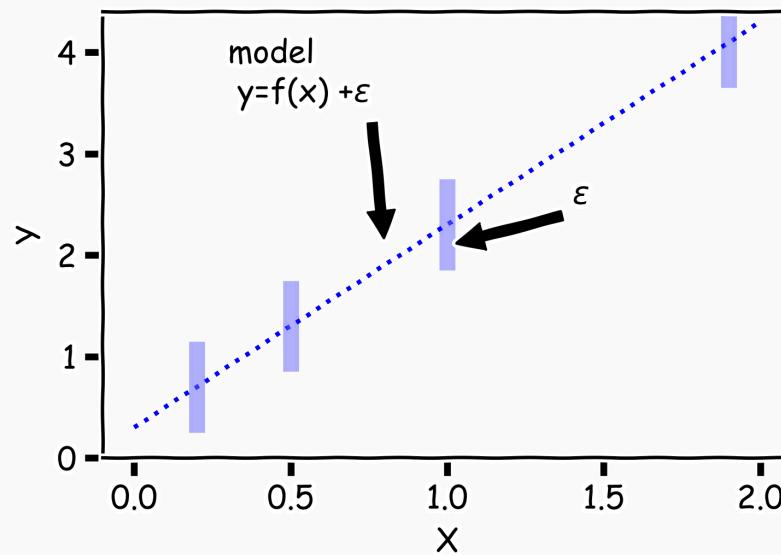
Confidence intervals for the predictors estimates (cont)

For some values of X , $Y = f(X)$



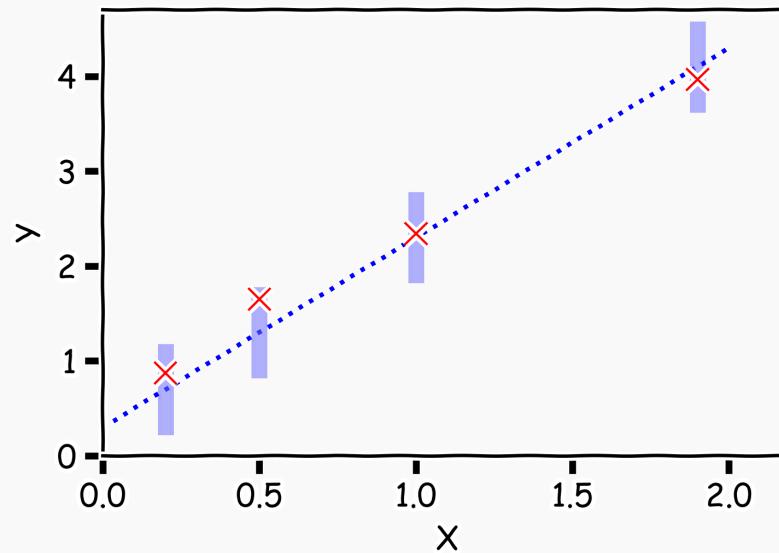
Confidence intervals for the predictors estimates (cont)

But due to error, every time we measure the response Y for a fixed value of X we will obtain a different observation.



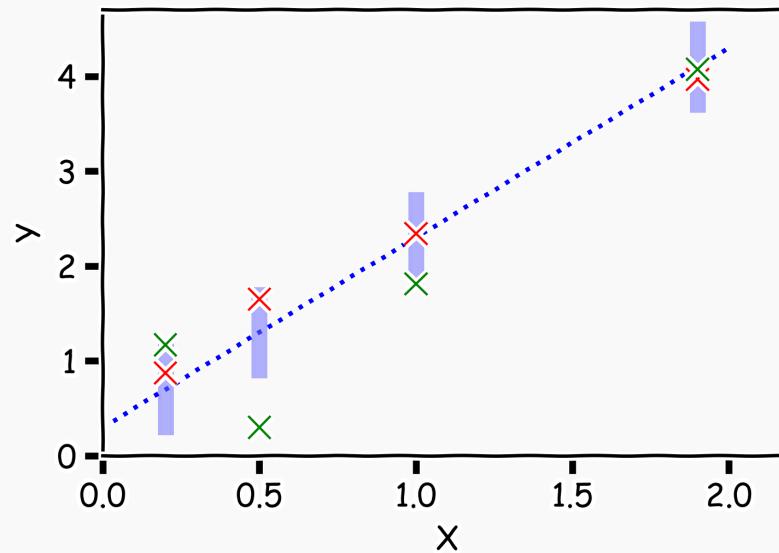
Confidence intervals for the predictors estimates (cont)

One set of observations, “one realization” we obtain one set of Ys (red crosses).



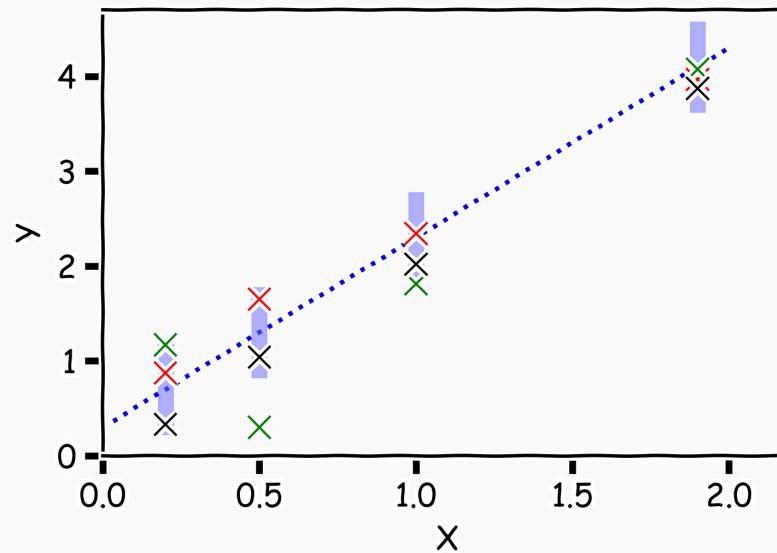
Confidence intervals for the predictors estimates (cont)

Another set of observations, “another realization” we obtain another set of Ys (green crosses).



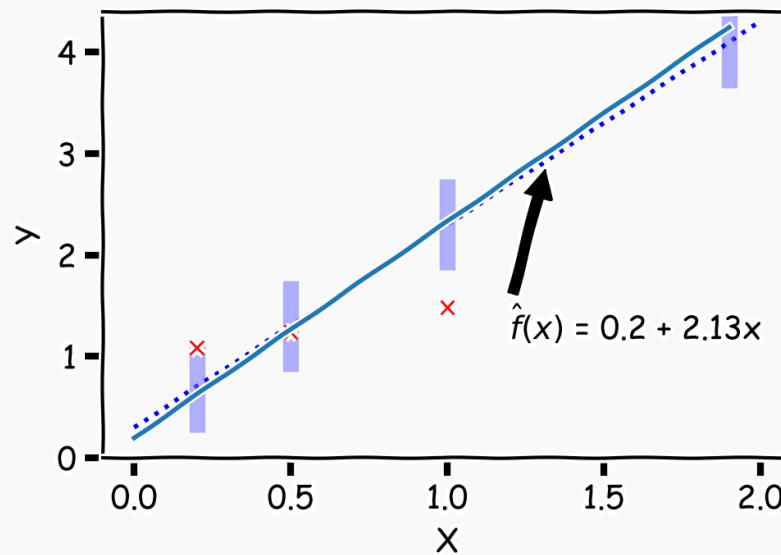
Confidence intervals for the predictors estimates (cont)

Another set of observations, “another realization” we obtain another set of Y s (black crosses).



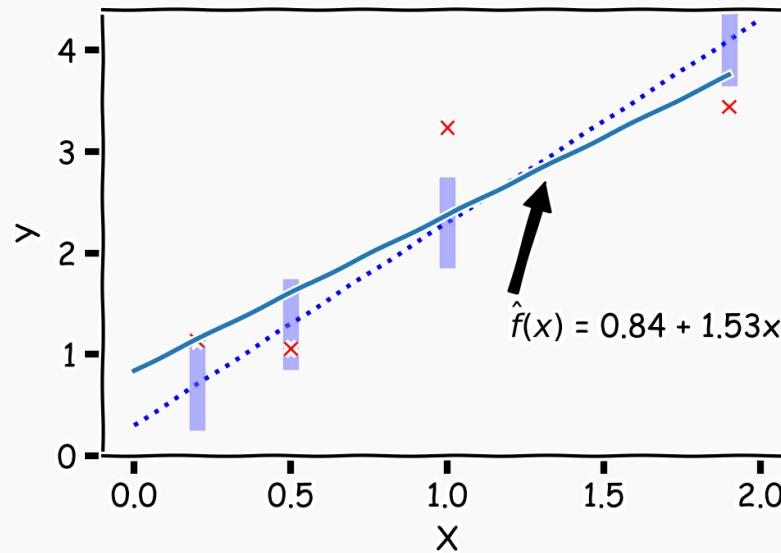
Confidence intervals for the predictors estimates (cont)

For each one of those “realizations”, we could fit a model and estimate $\hat{\beta}_0$ and $\hat{\beta}_1$.



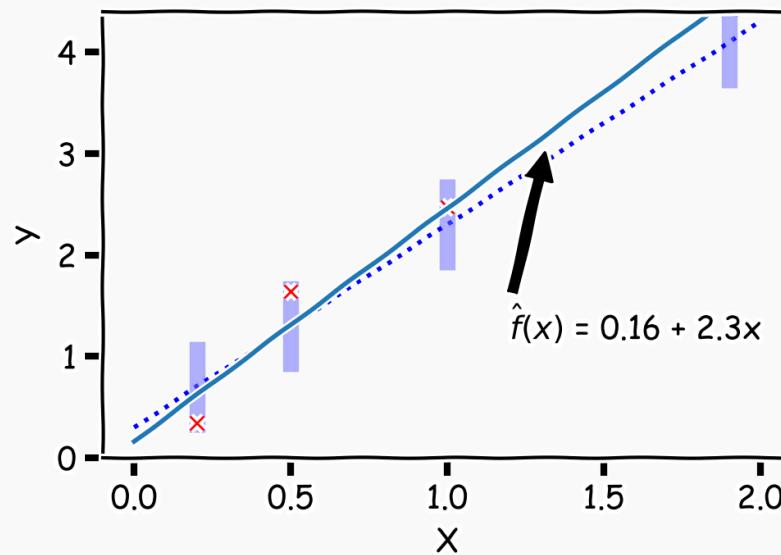
Confidence intervals for the predictors estimates (cont)

For each one of those “realizations”, we could fit a model and estimate, $\hat{\beta}_0$ and $\hat{\beta}_1$.



Confidence intervals for the predictors estimates (cont)

For each one of those “realizations”, we could fit a model and estimate, $\hat{\beta}_0$ and $\hat{\beta}_1$.

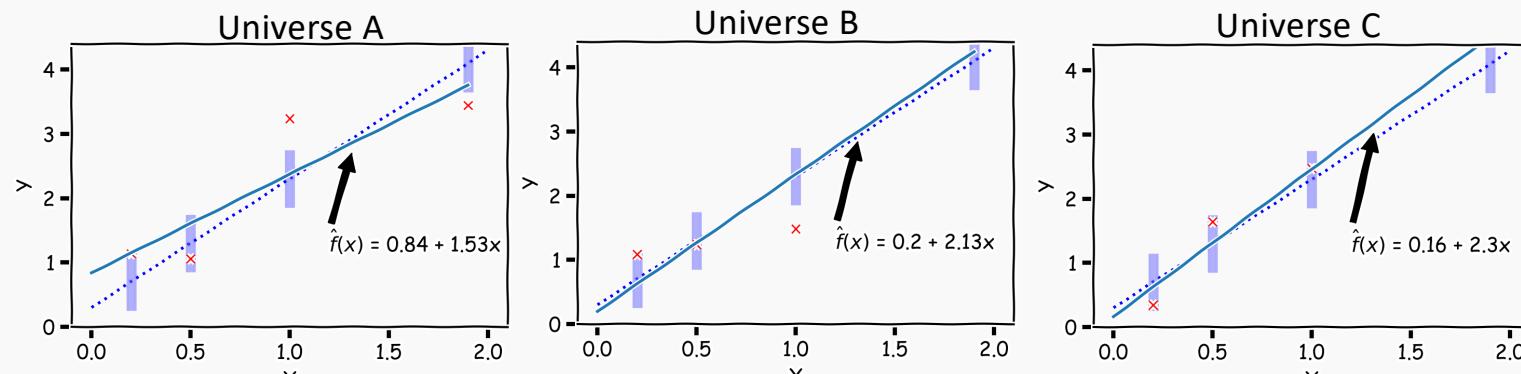


Confidence intervals for the predictors estimates (cont)

So if we just have one set of measurements of $\{X, Y\}$, our estimates of $\hat{\beta}_0$ and $\hat{\beta}_1$ are just for this particular realization.

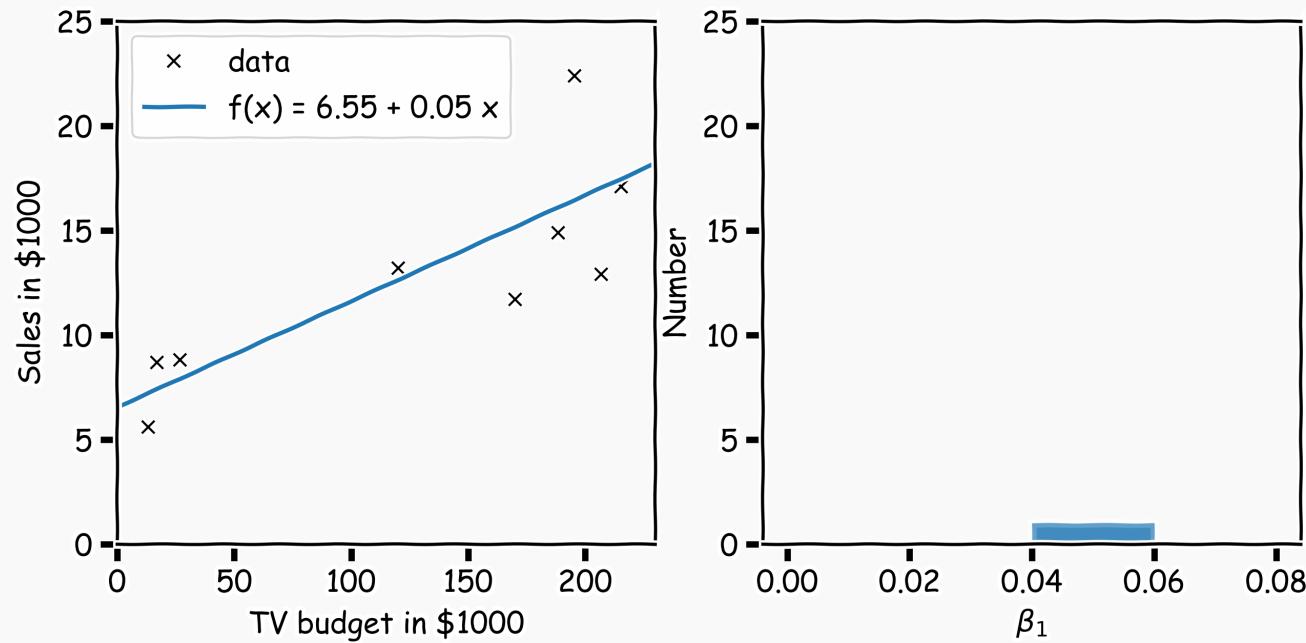
Question: If this is just one realization of the reality how do we know the truth? How do we deal with this conundrum?

Imagine (magic realism) we have parallel universes and we repeat this experiment on each of the other universes.



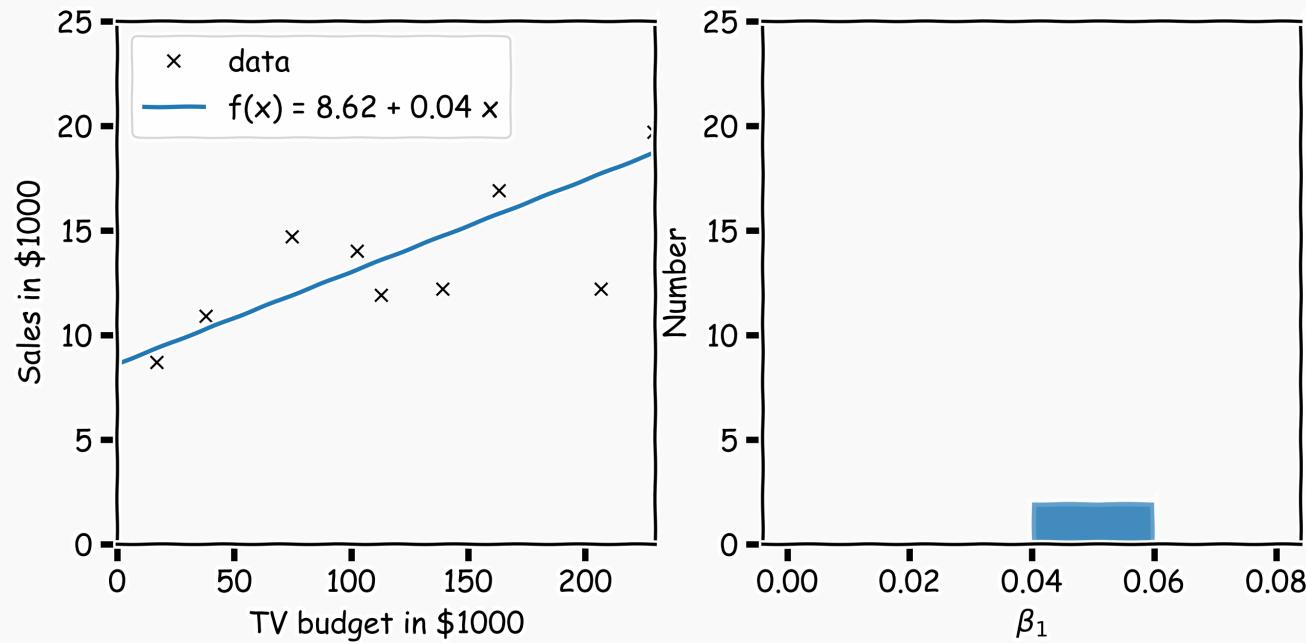
Confidence intervals for the predictors estimates (cont)

In our magical realisms, we can now sample multiple times



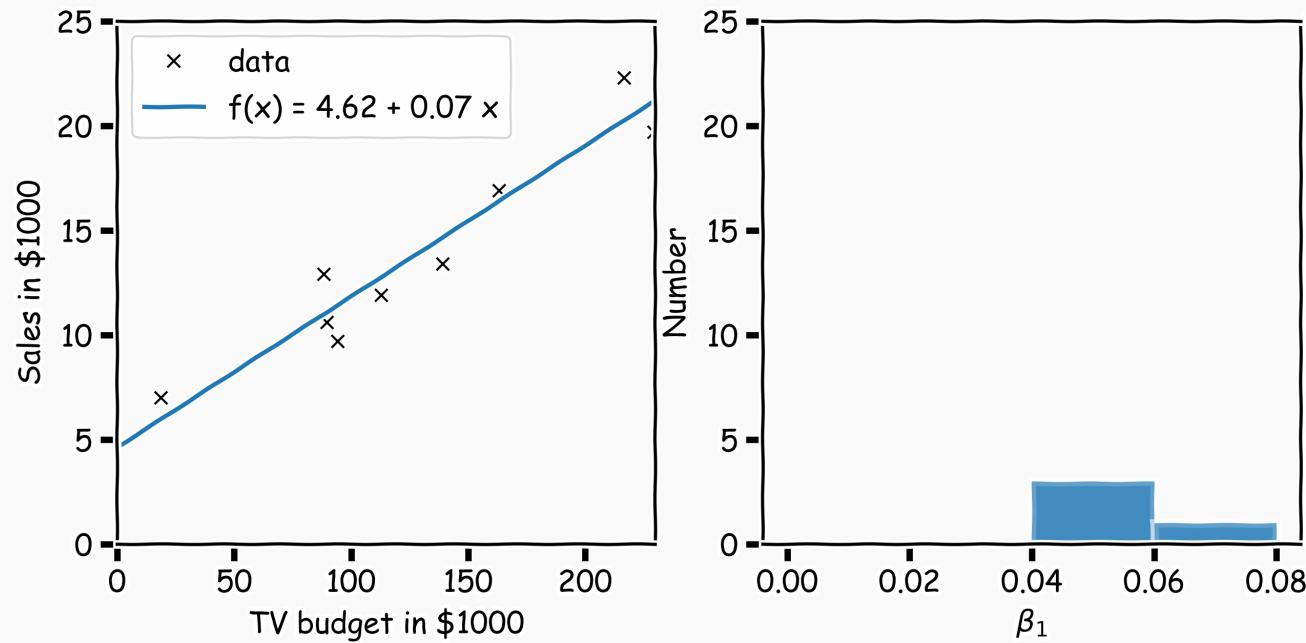
Confidence intervals for the predictors estimates (cont)

Another sample



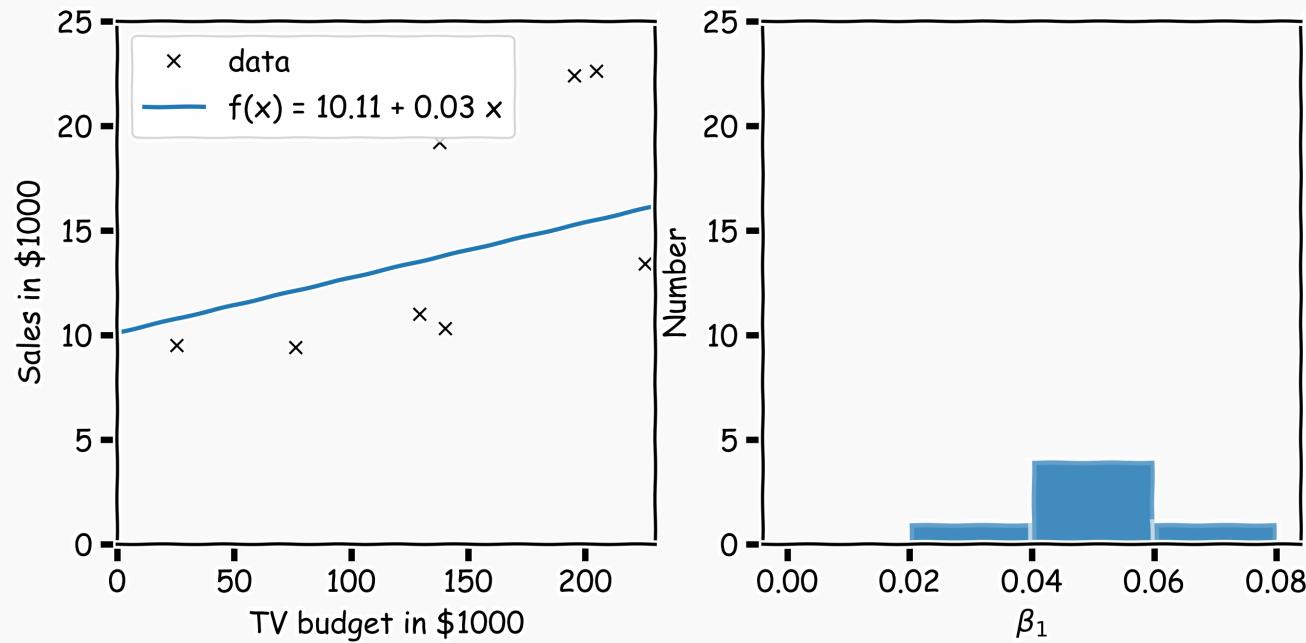
Confidence intervals for the predictors estimates (cont)

Another sample



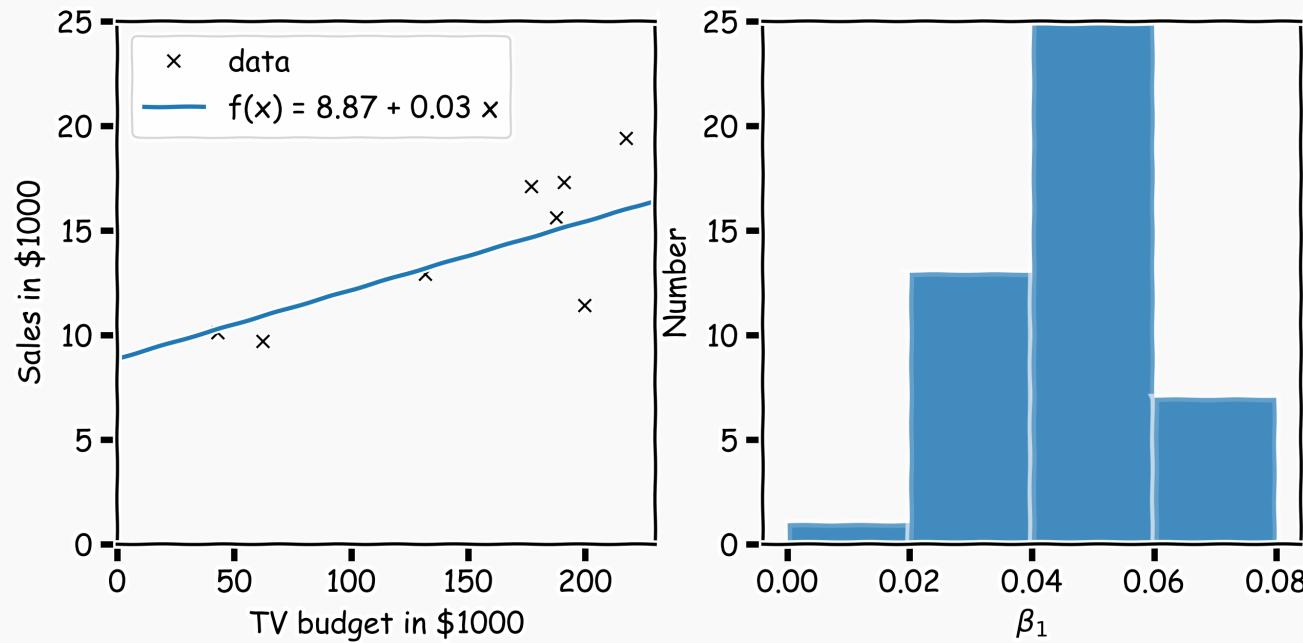
Confidence intervals for the predictors estimates (cont)

And another sample



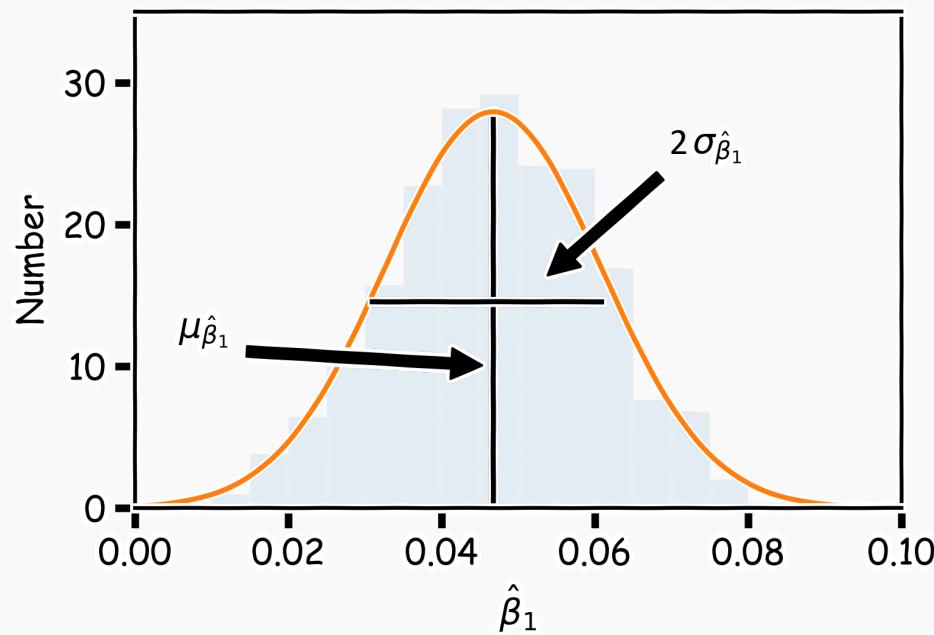
Confidence intervals for the predictors estimates (cont)

Repeat this for 100 times



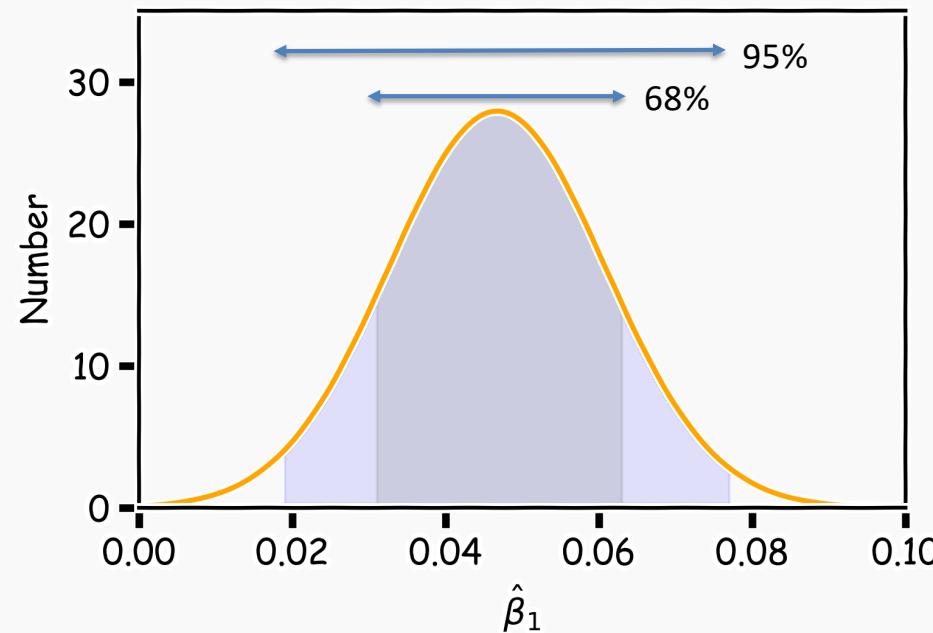
Confidence intervals for the predictors estimates (cont)

We can now estimate the mean and standard deviation of all the estimates $\hat{\beta}_1$. Square root of the variances of $\hat{\beta}_0$ and $\hat{\beta}_1$ are also called their **standard errors**, $SE(\hat{\beta}_0), SE(\hat{\beta}_1)$.



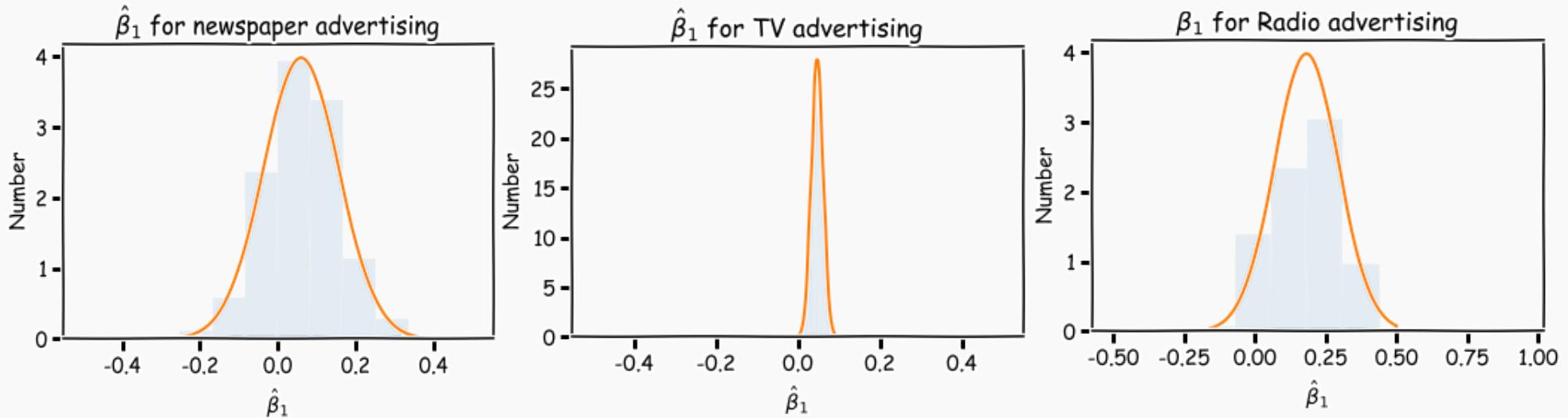
Confidence intervals for the predictors estimates (cont)

Finally we can calculate the confidence intervals, which are the ranges of values such that the **true** value of β_1 is contained in this interval with n percent probability.



And also we can answer the question, 'how significant are the predictors?' Here we show the same analysis for all three predictors.

Question: Which ones are important?



Before we answer this question, we need to answer another question.

Lecture Outline

- Linear models
- Estimate of the regression coefficients
 - Brute Force
 - Exact method
 - Gradient Descent
- Confidence intervals for the predictors estimates
- **Bootstrap**
- Evaluating Significance of Predictors
- How well we know the model \hat{f}

Bootstrap

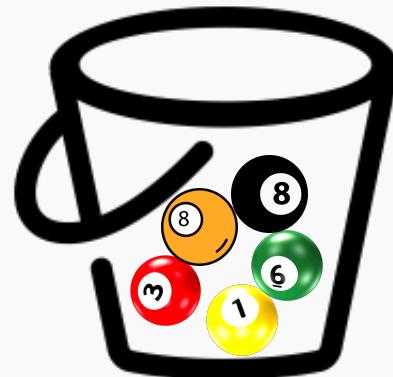
In the lack of active imagination, parallel universes and the likes, we need an alternative way of producing fake data set that resemble the parallel universes.

.

Bootstrapping is the practice of sampling from the observed data (X, Y) in estimating statistical properties.

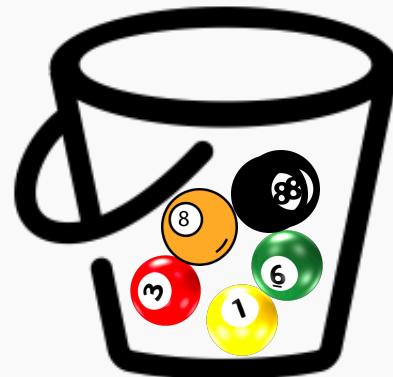
Bootstrap

Imagine we have 5 billiard balls in a bucket.

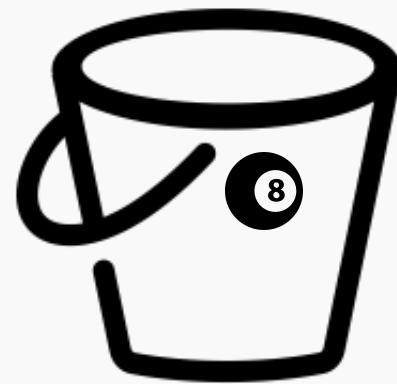
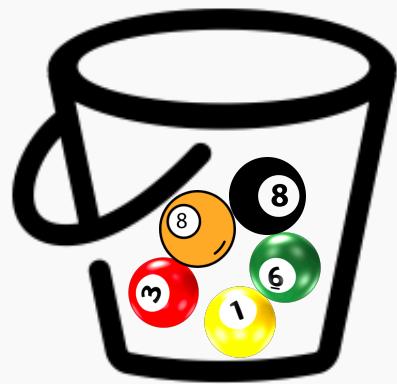


Bootstrap

We first pick randomly a ball and replicate it. This is called **sampling with replacement**. We move the replicated ball to another bucket.

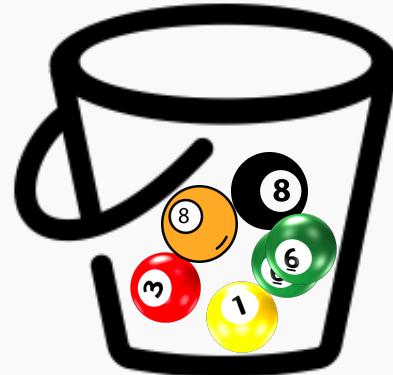


Bootstrap

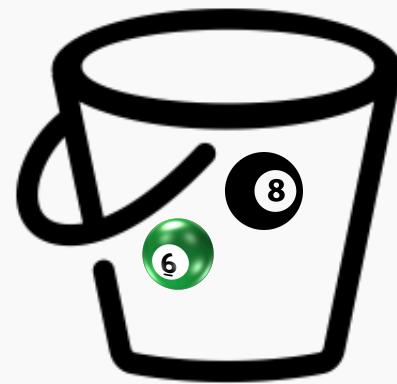
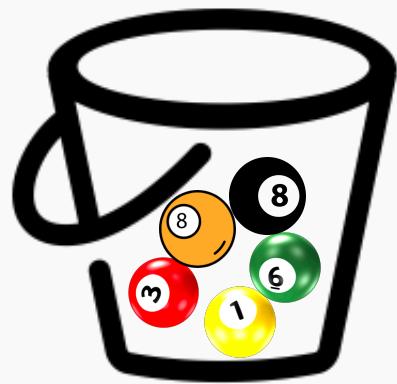


Bootstrap

We then randomly pick another ball and again we replicate it.
As before, we move the replicated ball to the other bucket.

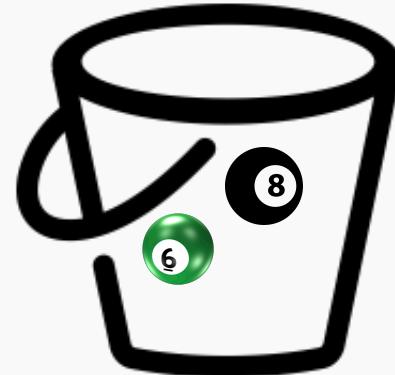
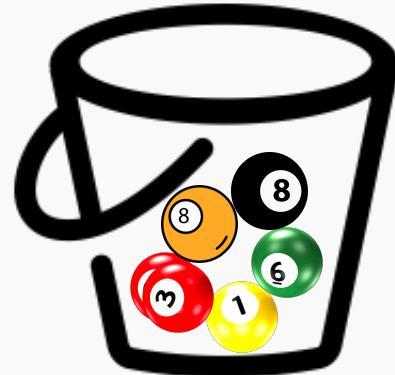


Bootstrap



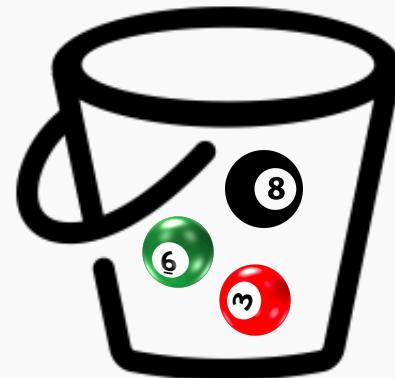
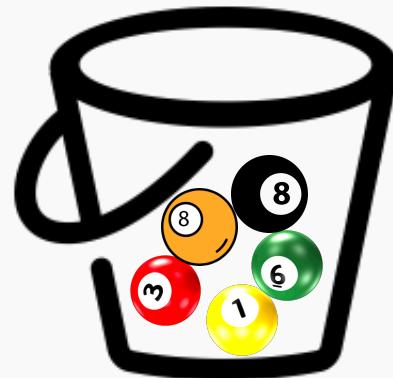
Bootstrap

We repeat this process.



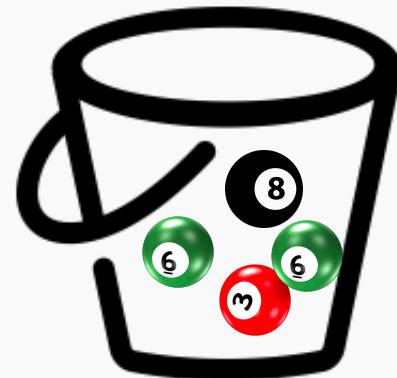
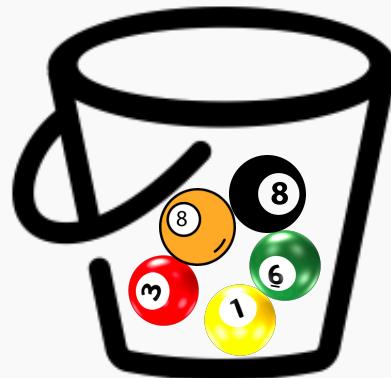
Bootstrap

Again



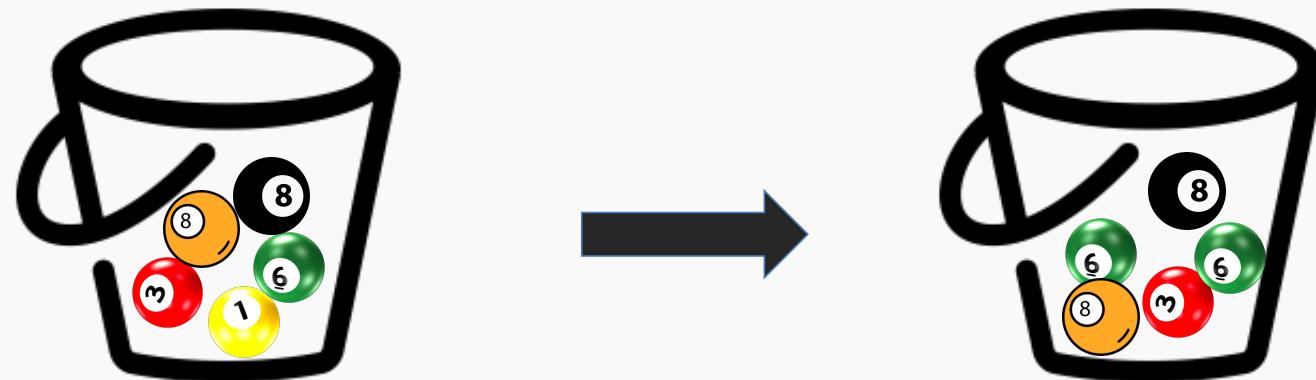
Bootstrap

And again



Bootstrap

Until the “other” bucket has **the same number of balls** as the original one.



This new bucket represents a new parallel universe

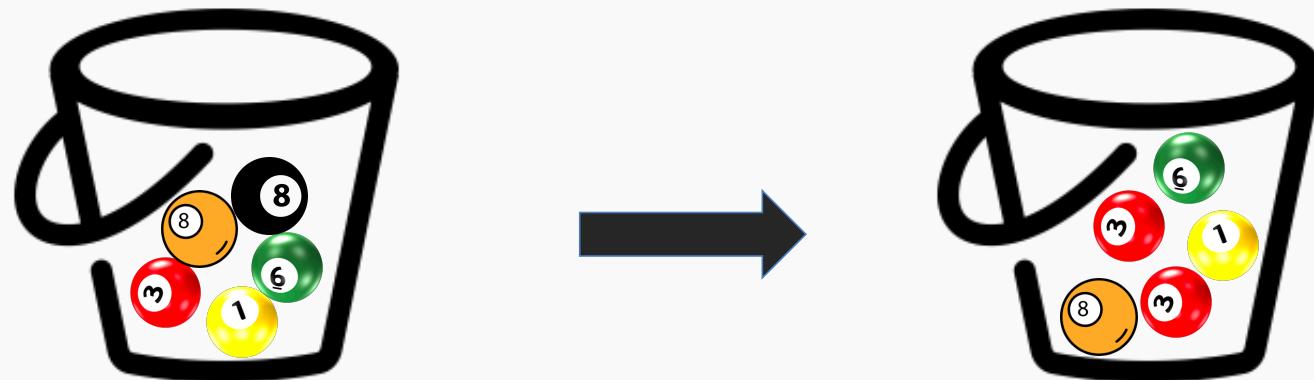
Bootstrap

We repeat the same process and acquire another sample.



Bootstrap

We repeat the same process and acquire another sample.



These new buckets represents the parallel universes

Bootstrapping for Estimating Sampling Error

Definition

Bootstrapping is the practice of estimating properties of an estimator by measuring those properties by, for example, sampling from the observed data.

For example, we can compute $\hat{\beta}_0$ and $\hat{\beta}_1$ multiple times by randomly sampling from our data set. We then use the variance of our multiple estimates to approximate the true variance of $\hat{\beta}_0$ and $\hat{\beta}_1$.

Confidence intervals for the predictors estimates: **Standard Errors**

We can empirically estimate the **standard errors**, $SE(\hat{\beta}_0), SE(\hat{\beta}_1)$ of β_0 and β_1 through bootstrapping.

If for each bootstrapped sample the estimated betas are: $\hat{\beta}_{0,i}, \hat{\beta}_{1,i}$, then

$$SE(\hat{\beta}_0) = \sqrt{\text{var}(\hat{\beta}_0)}$$

$$SE(\hat{\beta}_1) = \sqrt{\text{var}(\hat{\beta}_1)}$$

Confidence intervals for the predictors estimates: **Standard Errors**

Alternatively:

If we know the variance σ_ϵ^2 of the noise ϵ , we can compute $SE(\hat{\beta}_0), SE(\hat{\beta}_1)$ analytically using the formulae below (no need to bootstrap):

$$SE(\hat{\beta}_0) = \sigma_\epsilon \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum_i (x_i - \bar{x})^2}}$$
$$SE(\hat{\beta}_1) = \frac{\sigma_\epsilon}{\sqrt{\sum_i (x_i - \bar{x})^2}}$$

Standard Errors

More data: $n \uparrow$ and $\sum_i (x_i - \bar{x})^2 \uparrow \Rightarrow SE \downarrow$

Larger coverage: $var(x)$ or $\sum_i (x_i - \bar{x})^2 \uparrow \Rightarrow SE \downarrow$

Better data: $\sigma^2 \downarrow \Rightarrow SE \downarrow$

$$SE(\hat{\beta}_0) = \sigma \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum_i (x_i - \bar{x})^2}}$$
$$SE(\hat{\beta}_1) = \frac{\sigma}{\sqrt{\sum_i (x_i - \bar{x})^2}}$$

In practice, we do not know the theoretical value of σ since we do not know the exact distribution of the noise ϵ .

Standard Errors

However, if we make the following assumptions,

- the errors $\epsilon_i = y_i - \hat{y}_i$ and $\epsilon_j = y_j - \hat{y}_j$ are uncorrelated, for $i \neq j$,
- each ϵ_i has a mean 0 and variance σ_ϵ^2 ,

then, we can empirically estimate σ^2 , from the data and our regression line:

$$\sigma_\epsilon \approx \sqrt{\frac{n \cdot \text{MSE}}{n - 2}} = \sqrt{\frac{\sum_i (y_i - \hat{y}_i)^2}{n - 2}}$$

Remember:

$$y_i = f(x_i) + \epsilon_i \Rightarrow \epsilon_i = y_i - f(x_i)$$

Standard Errors

More data: $n \uparrow$ and $\sum_i(x_i - \bar{x})^2 \uparrow \Rightarrow SE \downarrow$

Larger coverage: $var(x)$ or $\sum_i(x_i - \bar{x})^2 \uparrow \Rightarrow SE \downarrow$

Better data: $\sigma^2 \downarrow \Rightarrow SE \downarrow$

$$SE(\hat{\beta}_0) = \sigma \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum_i (x_i - \bar{x})^2}}$$
$$SE(\hat{\beta}_1) = \frac{\sigma}{\sqrt{\sum_i (x_i - \bar{x})^2}}$$

Better model: $(\hat{f} - y_i) \downarrow \Rightarrow \sigma \downarrow \Rightarrow SE \downarrow$

$$\sigma \approx \sqrt{\sum \frac{(\hat{f}(x) - y_i)^2}{n - 2}}$$

Question: What happens to the $\widehat{\beta}_0$, $\widehat{\beta}_1$ under these scenarios?

Standard Errors

The following results are for the coefficients for TV advertising:

Method	$SE(\hat{\beta}_1)$
Analytic Formula	0.0061
Bootstrap	0.0061

The coefficients for TV advertising but restricting the coverage of x are:

Method	$SE(\hat{\beta}_1)$
Analytic Formula	0.0068
Bootstrap	0.0068

The coefficients for TV advertising but with added **extra** noise:

Method	$SE(\hat{\beta}_1)$
Analytic Formula	0.0028
Bootstrap	0.0023

Lecture Outline

- Linear models
- Estimate of the regression coefficients
 - Brute Force
 - Exact method
 - Gradient Descent
- Confidence intervals for the predictors estimates
- Bootstrap
- **Evaluating Significance of Predictors**
 - **Hypothesis Testing**
- How well we know the model \hat{f}

Interpretation of Predictors

Question: What do you think a predictor coefficient means?

$$Sales = 7.5 + 0.04 TV$$

What does 7.5 mean and what does 0.04 mean?

If we increase the TV by \$1000, what would you expect the increase in sales to be?

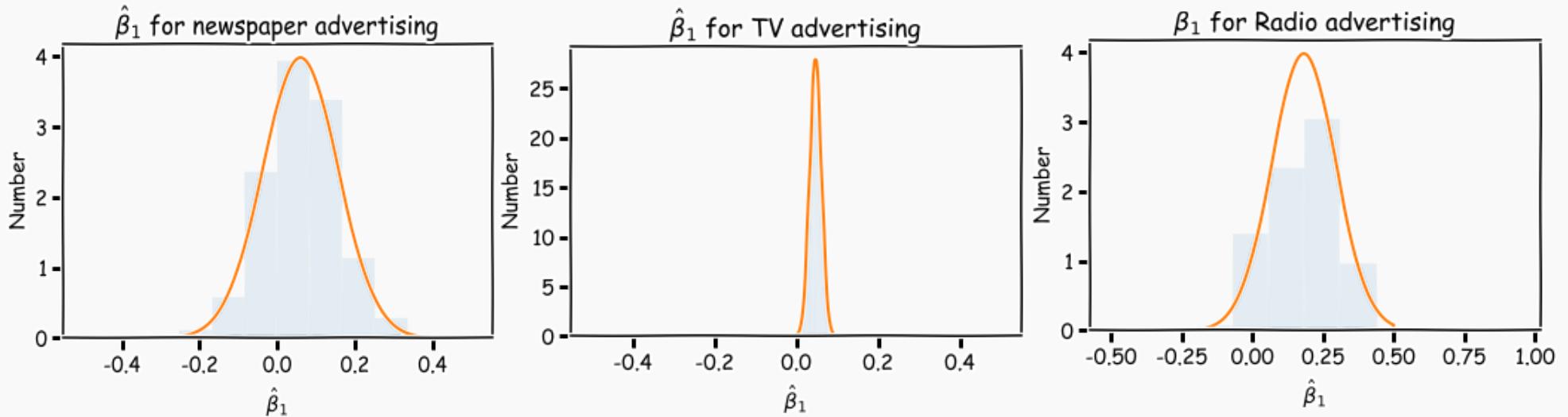
What if?

$$Sales = 7.5 + 1.01 TV$$

The interpretation of the predictors depends on the values but decisions depend on how much we trust these values.

And also we can answer the question, 'how significant are the predictors?' Here we show the same analysis for all three predictors.

Question: Which ones are important?



Now we know how to generate these distributions we are ready to answer
'how significant are the predictors?'

Hypothesis Testing

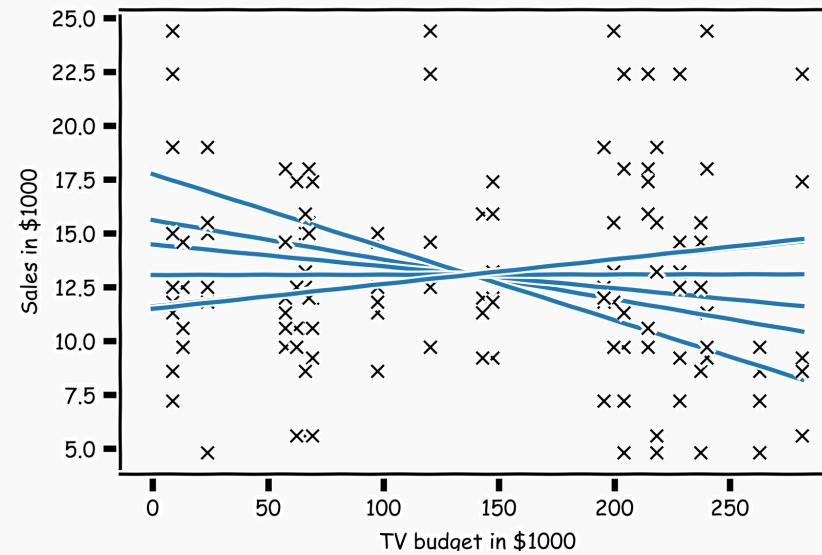
Hypothesis testing is a formal process through which we evaluate the validity of a statistical hypothesis by considering evidence **for** or **against** the hypothesis gathered by **random sampling of the data.**

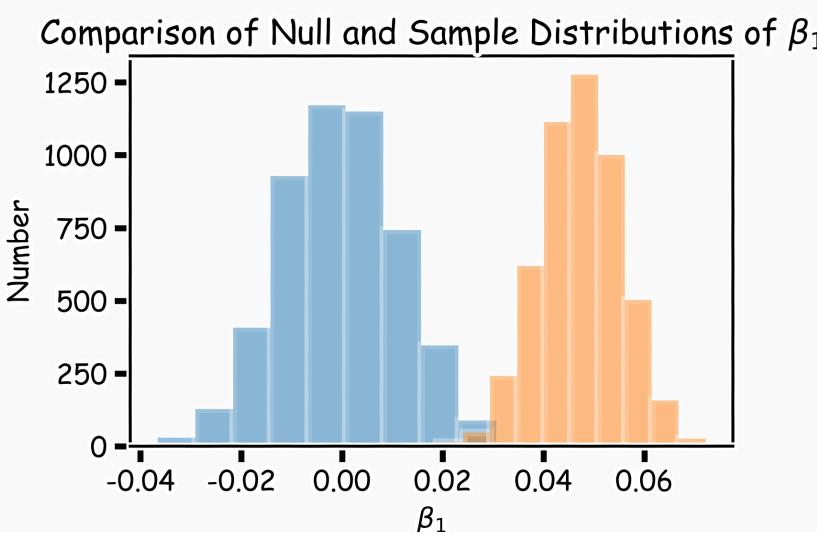
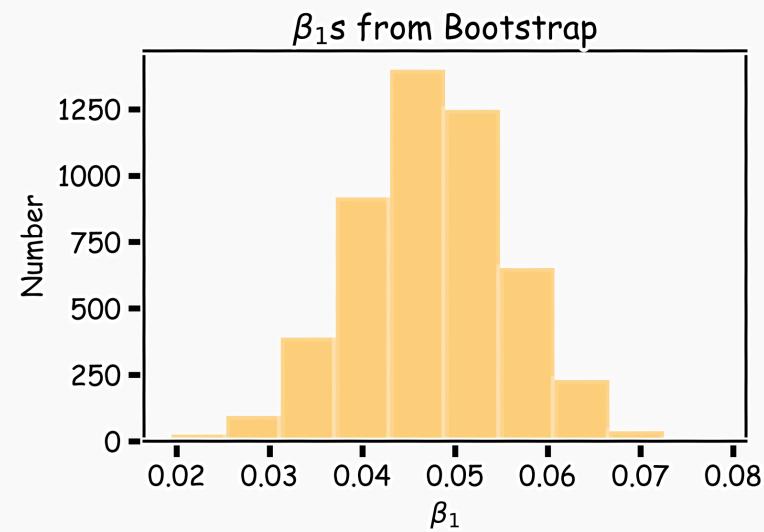
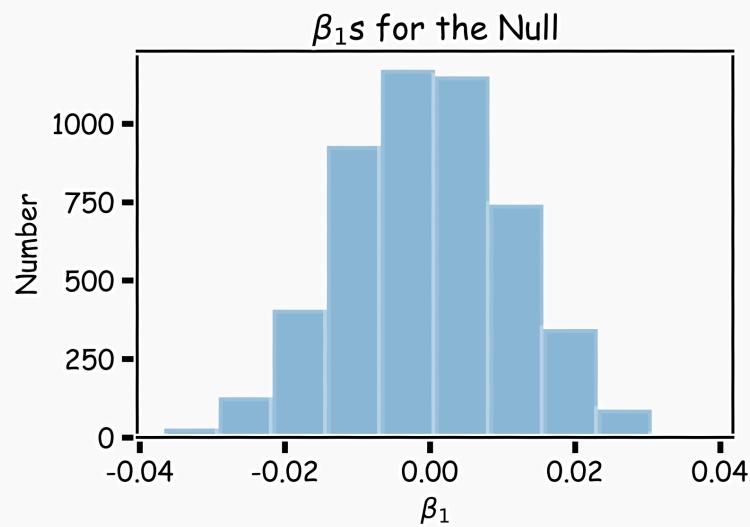
TV	sales
2364	22.1
2925	10.4
2028	9.3
1915	18.5
2829	12.9
2964	7.2
2054	11.8
2022	13.2
2009	4.8
1928	10.6
2022	8.6
1926	17.4
2322	9.2
1929	9.7
2021	19.0
2924	22.4
2324	12.5
2958	24.4

TV	sales
50.0	22.1
184.9	10.4
11.7	9.3
219.8	18.5
13.1	12.9
248.8	7.2
76.4	11.8
197.6	13.2
195.4	4.8
75.5	10.6
238.2	8.6
222.4	17.4
171.3	9.2
184.9	9.7
193.2	19.0
131.7	22.4
116.0	12.5
166.8	24.4

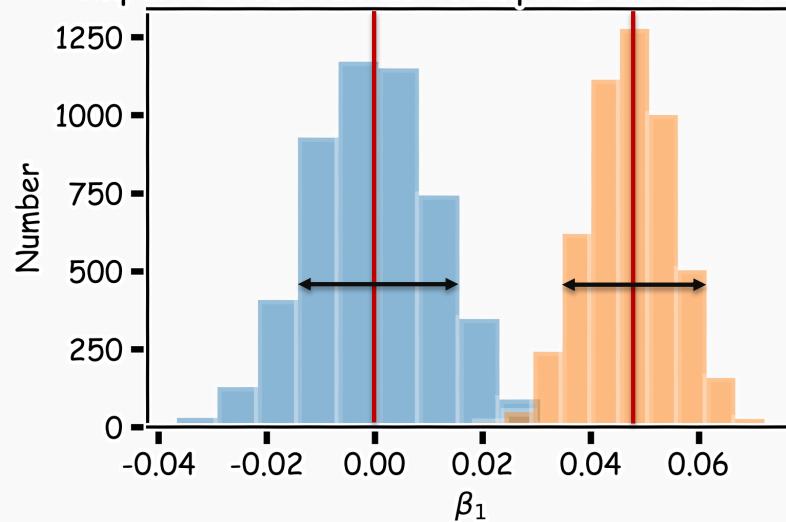
Random sampling of the data

Shuffle the values of the predictor variable





Comparison of Null and Sample Distributions of β_1



$$\mu_{Null} = 0$$

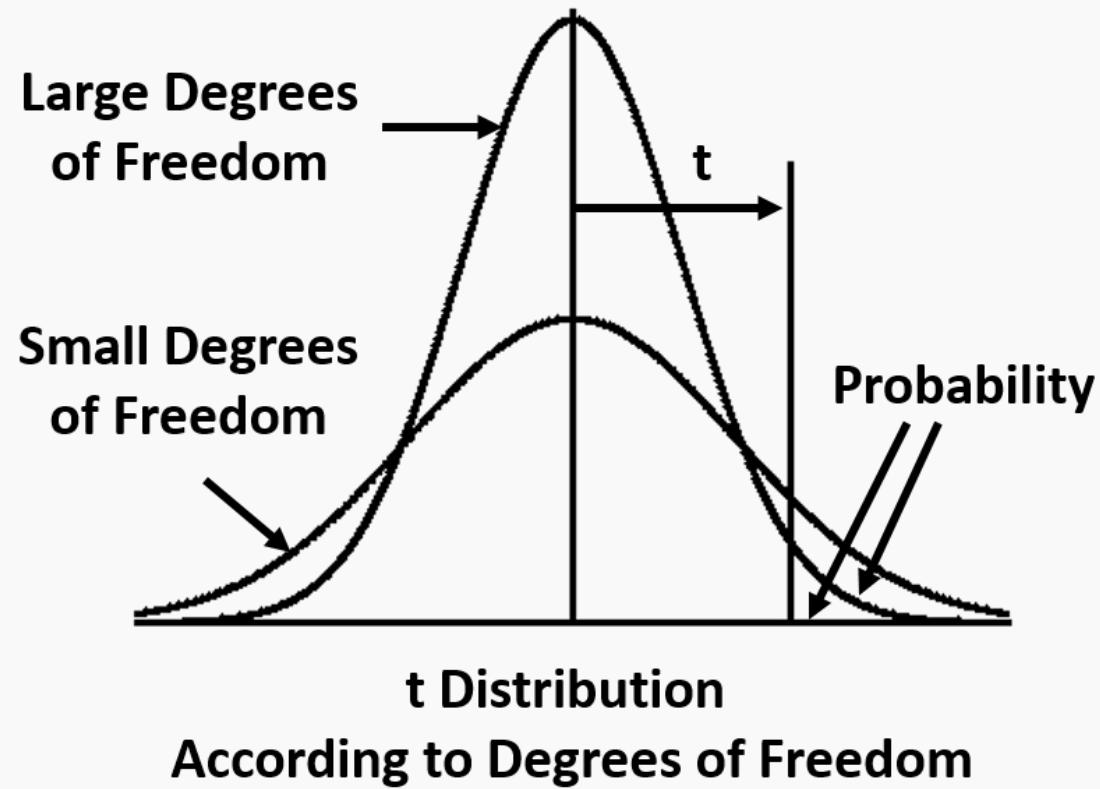
$$\mu_{\hat{\beta}} = \mu_{boot}$$

$$\sigma_{\hat{\beta}} = SE(\hat{\beta}) = \sigma_{boot}$$

$$\sigma_{Null} \approx \sigma_{\hat{\beta}}$$

Translate this to the significance. Let's look at the distance of the estimated value of the coefficient in units of $SE(\hat{\beta}_1) = \sigma_{\hat{\beta}_1}$.

$$D = \frac{\mu_{\hat{\beta}} - \mu_{Null}}{\sigma_{\hat{\beta}}}$$



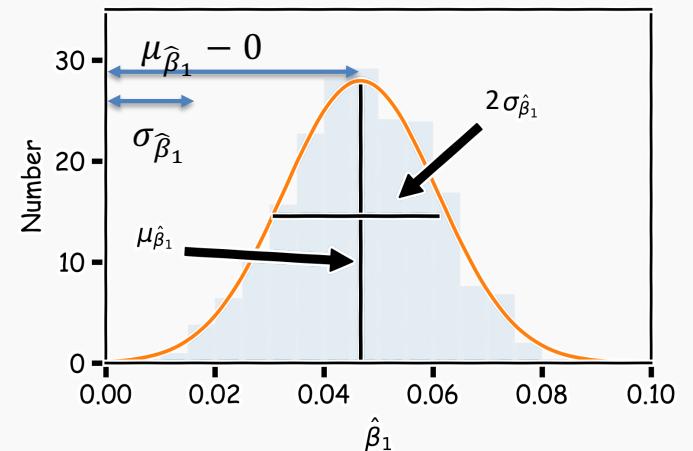
Importance of predictors

In practice, we do not need the distribution for Null.

Define a test statistic, which we call t-test statistic

$$t = \frac{\hat{\beta}_1}{\sigma_{\hat{\beta}_1}}$$

Which measures the distance from zero in units of standard deviation.



We evaluate how often a particular value of t can occur by accident. We expect that t will have a *t-distribution with $n-2$ degrees of freedom*.

To compute the probability of observing any value equal to $|t|$ or larger, assuming $\hat{\beta}_1 = 0$ is easy. We call this probability the **p-value**.

a small p-value (<0.05) indicates that it is unlikely to observe such a substantial association between the predictor and the response due to chance.

Hypothesis Testing

Hypothesis testing is a formal process through which we evaluate the validity of a statistical hypothesis by considering evidence for or against the hypothesis gathered by **random sampling** of the data.

1. State the hypotheses, typically a **null hypothesis**, H_0 and an **alternative hypothesis**, H_1 , that is the negation of the former.
2. Choose a type of analysis, i.e. how to use sample data to evaluate the null hypothesis. Typically this involves choosing a single test statistic.
3. **Sample** data and compute the test statistic.
4. Use the value of the test statistic to either **reject** or **not reject** the null hypothesis.

Hypothesis testing

1. State Hypothesis:

Null hypothesis:

H_0 : There is no relation between X and Y

The alternative:

H_a : There is some relation between X and Y

2: Choose test statistics

To test the null hypothesis, we need to determine whether, our estimate for $\hat{\beta}_1$, is sufficiently far from zero that we can be confident that $\hat{\beta}_1$ is non-zero. We use the following test statistic:

$$t = \frac{\mu_{\hat{\beta}_1}}{\sigma_{\hat{\beta}_1}}$$

Hypothesis testing

3. Sample:

Using bootstrap we can estimate $\hat{\beta}'$ s, and therefore $\mu_{\hat{\beta}_1}$ and $\sigma_{\hat{\beta}_1}$.

4. Reject or not reject the hypothesis:

If there is really no relationship between X and Y , then we expect that will have a *t-distribution with n-2 degrees of freedom*.

To compute the probability of observing any value equal to $|t|$ or larger, assuming $\hat{\beta}_1 = 0$ is easy. We call this probability the p-value.

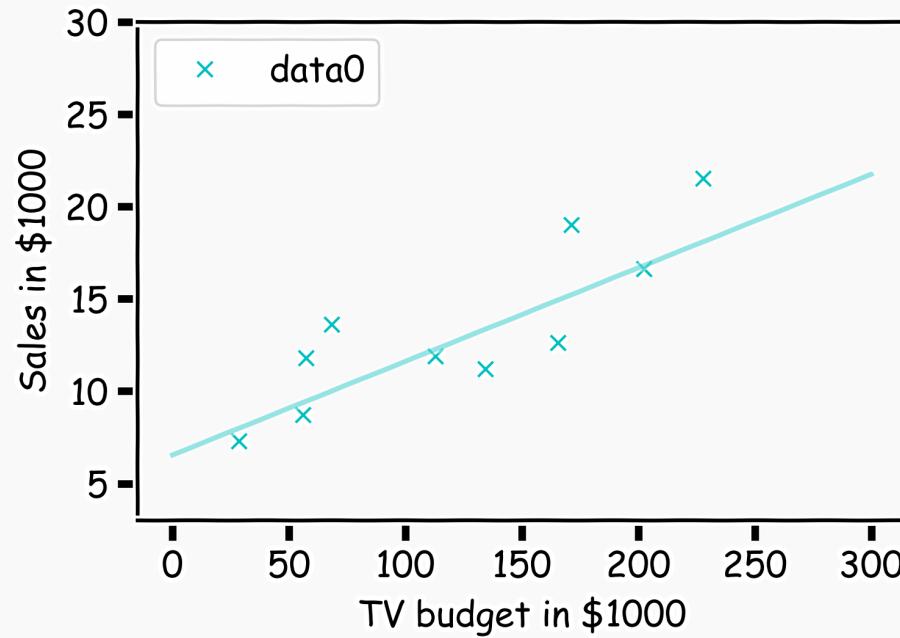
a small p-value indicates that it is unlikely to observe such a substantial association between the predictor and the response due to chance

Lecture Outline

- Linear models
- Estimate of the regression coefficients
 - Brute Force
 - Exact method
 - Gradient Descent
- Confidence intervals for the predictors estimates
- Bootstrap
- Evaluating Significance of Predictors
 - Hypothesis Testing
- **How well we know the model \hat{f}**

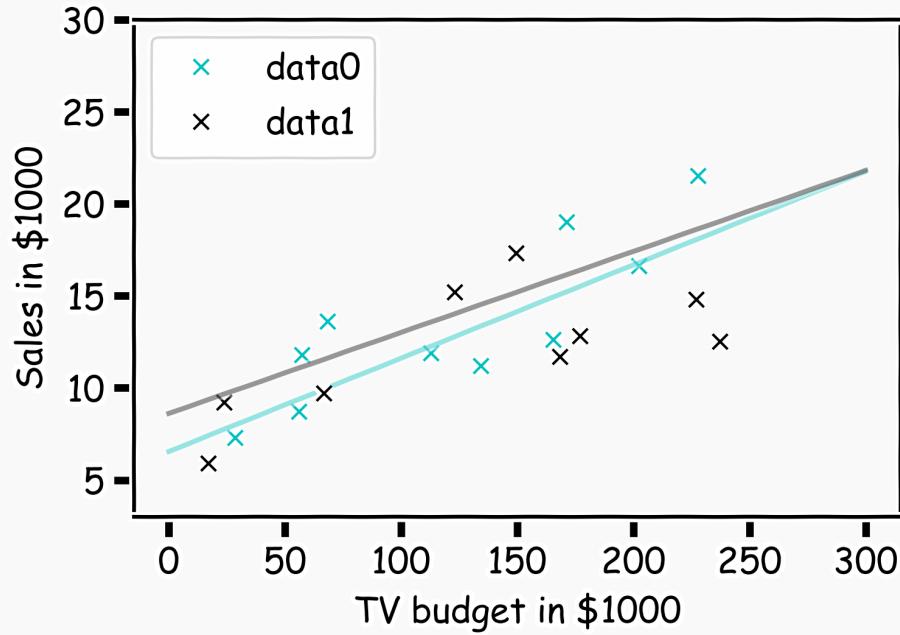
How well do we know \hat{f} ?

Our confidence in f is directly connected with the confidence in β s. So for each bootstrap sample, we have one β which we can use to determine the model.



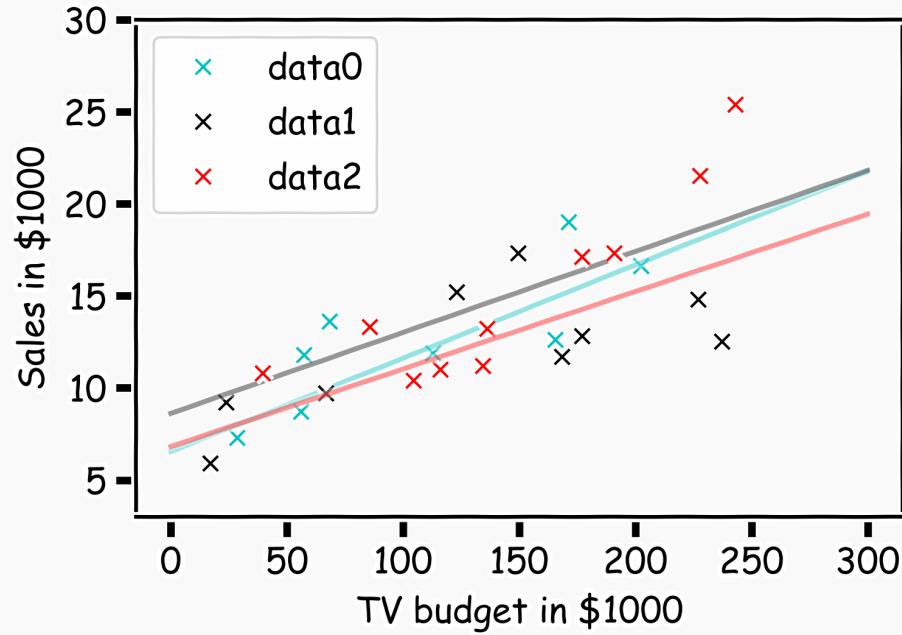
How well do we know \hat{f} ?

Here we show two different sets of models given the fitted coefficients.



How well do we know \hat{f} ?

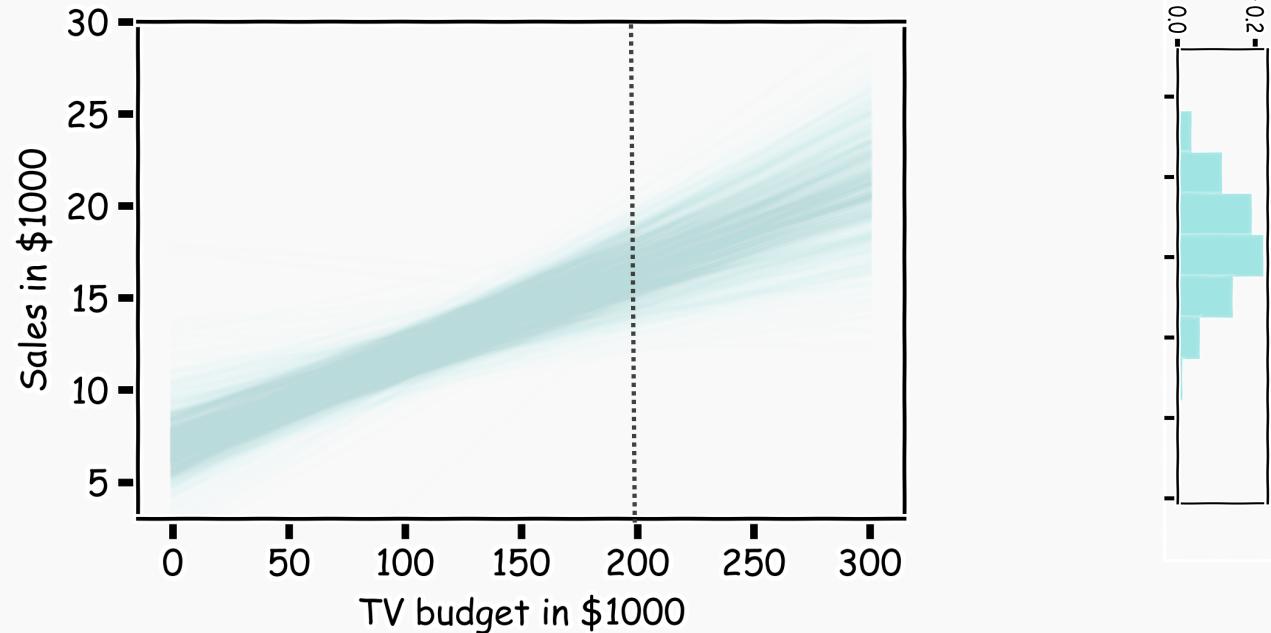
There is one such regression line for every bootstrapped sample.



How well do we know \hat{f} ?

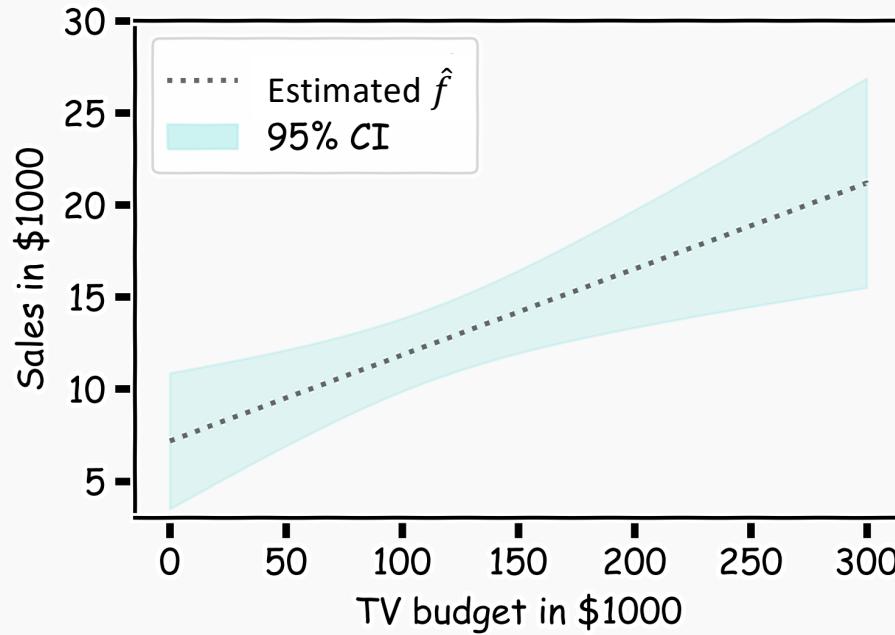
Below we show all regression lines for a thousand of such bootstrapped samples.

For a given x , we examine the distribution of \hat{f} , and determine the mean and standard deviation.

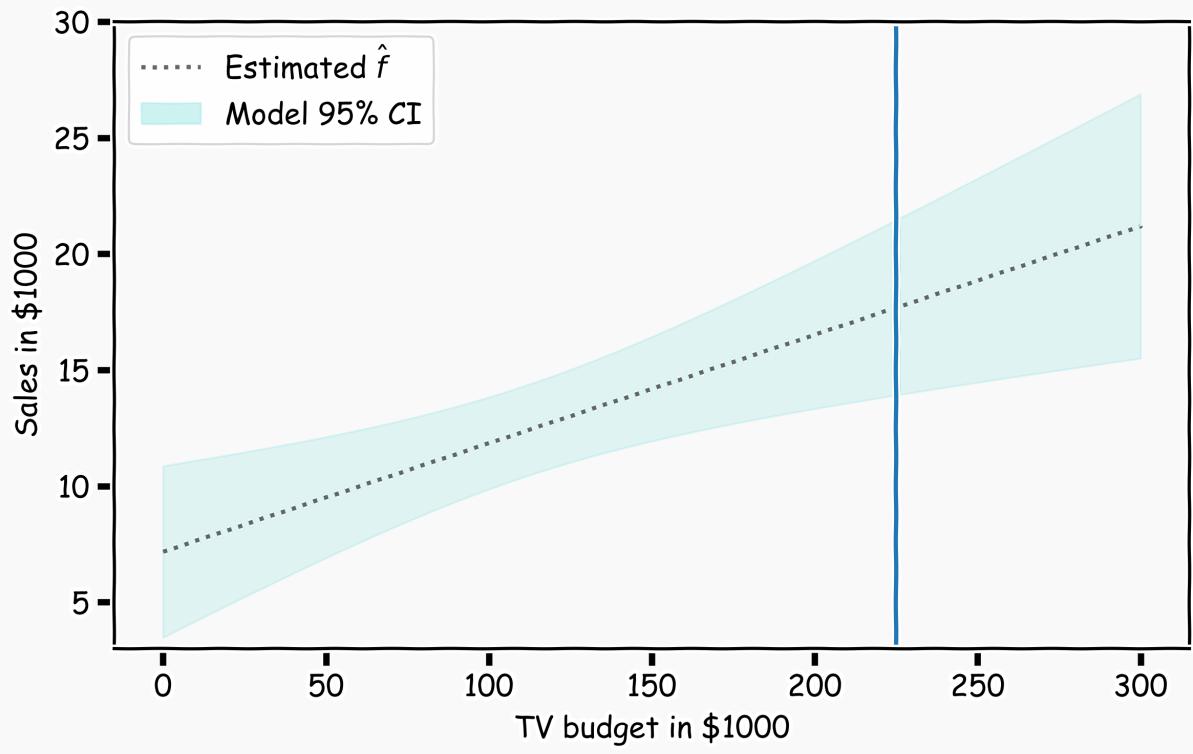


How well do we know \hat{f} ?

For every x , we calculate the mean of the models, \hat{f} (shown with dotted line) and the 95% CI of those models (shaded area).

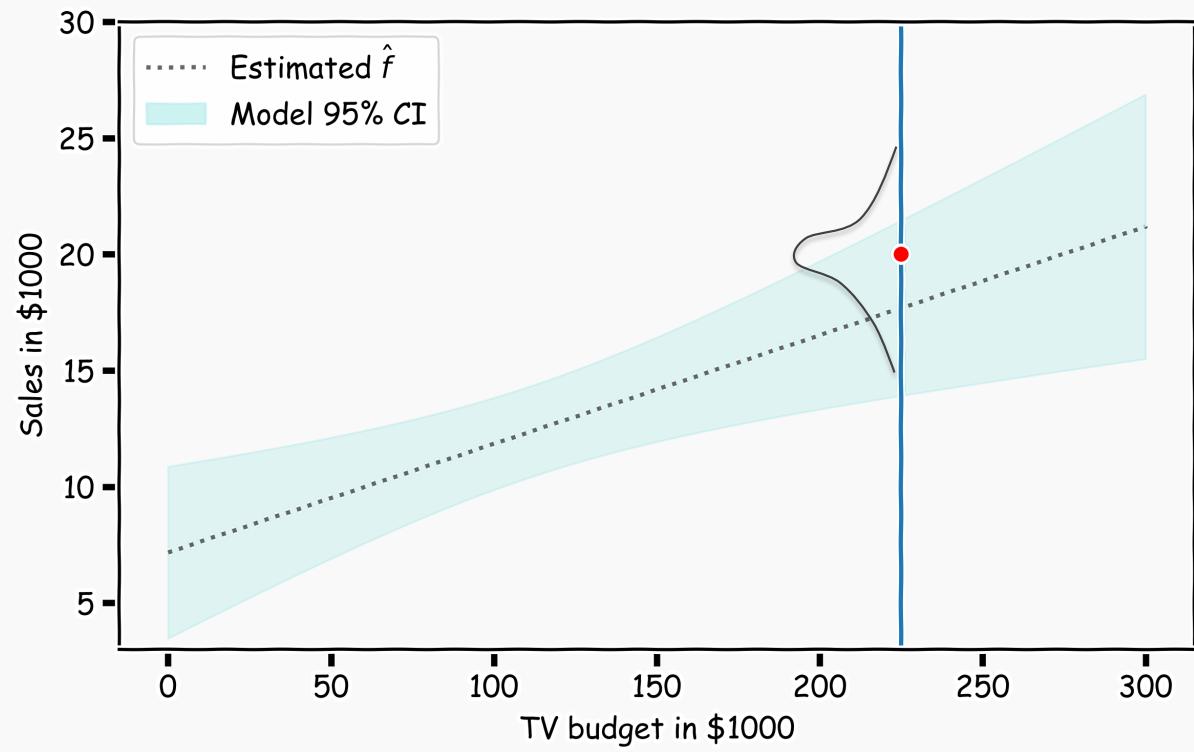


Confidence in predicting \hat{y}



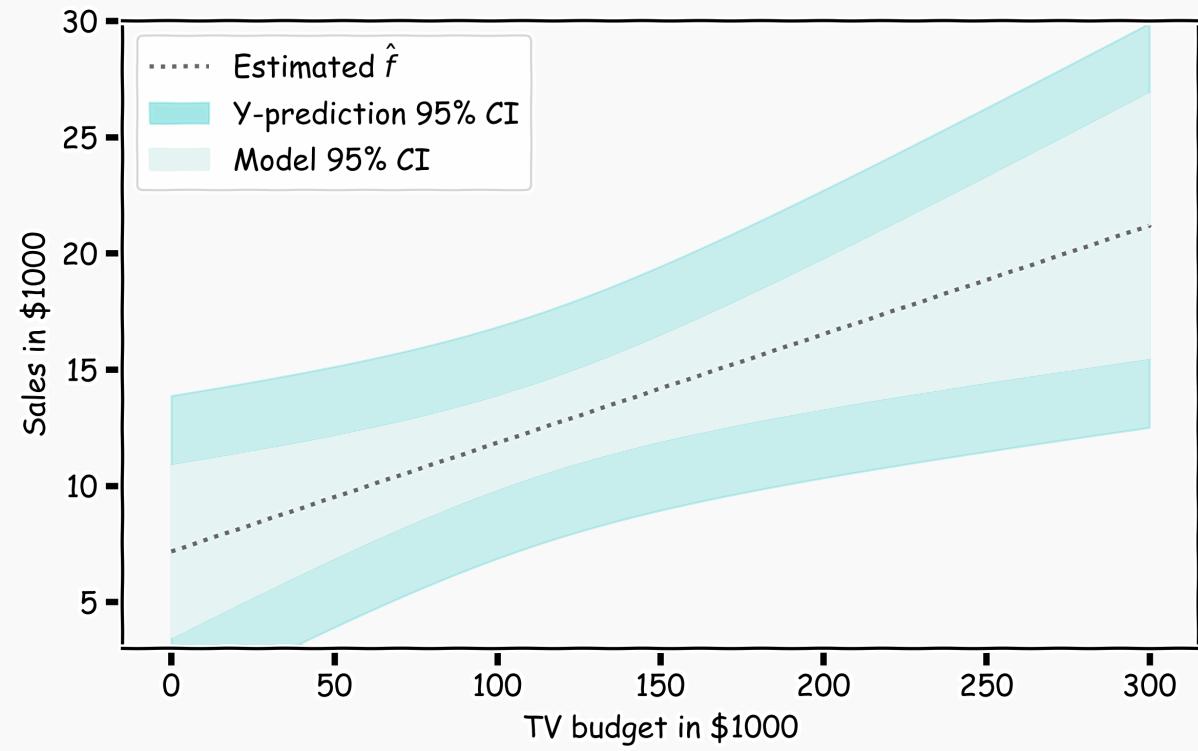
Confidence in predicting \hat{y}

- for a given x , we have a distribution of models $f(x)$
- for each of these $f(x)$, the prediction for $y \sim N(f, \sigma_\epsilon)$



Confidence in predicting \hat{y}

- for a given x , we have a distribution of models $f(x)$
- for each of these $f(x)$, the prediction for $y \sim N(f, \sigma_\epsilon)$
- The prediction confidence intervals are then



Summary so far

Model Fitness

How does the model perform predicting?

Comparison of Two Models

How do we choose from two different models?

Evaluating Significance of Predictors

Does the outcome depend on the predictors?

How well do we know \hat{f}

The confidence intervals of our \hat{f}

What's next?

Multiple predictors

Collinearity

Categorical variables

Polynomial regression

Interaction terms