

# Assessing accuracy of linear regression

Chris Cornwell

Oct 9, 2025

# Outline

Confidence intervals with linear regression

Measuring how well LSR line fits

Multiple variables

Polynomial fitting

# Outline

Confidence intervals with linear regression

Measuring how well LSR line fits

Multiple variables

Polynomial fitting

## Difference between parameters from “population” vs. from data

- Modeling relationship between independent variables and  $y$  with a linear model (with noise in the  $y$ -coordinate direction). In other words, the modeled relationship is that

$$y \approx b + \mathbf{x}^T \mathbf{w}.$$

for some parameters  $\mathbf{w}$ ,  $b$ .

---

<sup>1</sup>There is some number  $M$  so that  $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$ .

## Difference between parameters from “population” vs. from data

- Modeling relationship between independent variables and  $y$  with a linear model (with noise in the  $y$ -coordinate direction). In other words, the modeled relationship is that

$$y \approx b + \mathbf{x}^T \mathbf{w}.$$

for some parameters  $\mathbf{w}$ ,  $b$ .

- Alternatively, we may have some function<sup>1</sup> of the variables  $f(\mathbf{x})$  and a modeled relationship

$$y \approx b + f(\mathbf{x})^T \mathbf{w}.$$

---

<sup>1</sup>There is some number  $M$  so that  $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$ .

## Difference between parameters from “population” vs. from data

- Modeling relationship between independent variables and  $y$  with a linear model (with noise in the  $y$ -coordinate direction). In other words, the modeled relationship is that

$$y \approx b + \mathbf{x}^T \mathbf{w}.$$

for some parameters  $\mathbf{w}$ ,  $b$ .

- Alternatively, we may have some function<sup>1</sup> of the variables  $f(\mathbf{x})$  and a modeled relationship

$$y \approx b + f(\mathbf{x})^T \mathbf{w}.$$

However, we find values for  $\mathbf{w}$  and  $b$  by using observed data, *sampled* from a “population.” Among the entire population, there is a best fit linear model having some parameters. However, from the observed data  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_p, y_p)$  if our procedure determines best fit parameters  $b^*$  and  $\mathbf{w}^*$ , these are not (necessarily) the parameters for the population linear model.

<sup>1</sup>There is some number  $M$  so that  $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$ .

## Example

**Simulate noisy linear data:** make 30 points, from a line with slope  $-1.6$  and  $y$ -intercept  $0.8$ ; put noise into the  $y$ -coordinate with standard deviation  $\sigma = 0.5$ .

## Example

**Simulate noisy linear data:** make 30 points, from a line with slope  $-1.6$  and  $y$ -intercept  $0.8$ ; put noise into the  $y$ -coordinate with standard deviation  $\sigma = 0.5$ .

```
1 | x = np.random.uniform(0, 2, size=30)
2 |
3 | def simulate_data(x, std):
4 |     return -1.6*x + 0.8 + np.random.normal(0, std, size=len(x))
5 | y = simulate_data(x, 0.5)
```



## Example

**Simulate noisy linear data:** make 30 points, from a line with slope  $-1.6$  and y-intercept  $0.8$ ; put noise into the y-coordinate with standard deviation  $\sigma = 0.5$ .

```
1 | x = np.random.uniform(0, 2, size=30)
2 |
3 | def simulate_data(x, std):
4 |     return -1.6*x + 0.8 + np.random.normal(0, std, size=len(x))
5 | y = simulate_data(x, 0.5)
```

With this simulated data set, compute the slope  $w^*$  and intercept  $b^*$  from linear regression; then put  $w^*$ , and  $b^*$ , into a list of slopes, intercepts respectively. Iterate this 1000 times  $\rightarrow$  a list of 1000 slopes, another with 1000 intercepts.

## Example

**Simulate noisy linear data:** make 30 points, from a line with slope  $-1.6$  and y-intercept  $0.8$ ; put noise into the y-coordinate with standard deviation  $\sigma = 0.5$ .

```
1 | x = np.random.uniform(0, 2, size=30)
2 |
3 | def simulate_data(x, std):
4 |     return -1.6*x + 0.8 + np.random.normal(0, std, size=len(x))
5 | y = simulate_data(x, 0.5)
```

With this simulated data set, compute the slope  $w^*$  and intercept  $b^*$  from linear regression; then put  $w^*$ , and  $b^*$ , into a list of slopes, intercepts respectively. Iterate this 1000 times  $\rightarrow$  a list of 1000 slopes, another with 1000 intercepts.

What is the mean of these slopes and intercepts that were found?

## Sample statistic, relation to population statistic

This is fundamental to statistics.

- Say that a sample of 2000 people are selected from around the country and their height is measured. Mean of these 2000 heights: sample mean.

# Sample statistic, relation to population statistic

This is fundamental to statistics.

- ▶ Say that a sample of 2000 people are selected from around the country and their height is measured. Mean of these 2000 heights: sample mean.
- ▶ Sample mean differs from the true mean height of the entire population of the country. (Not by much, perhaps.)
  - ▶ Weak Law of Large Numbers: if  $s$  random samples of 2000 people taken, and each sample mean calculated then, as  $s \rightarrow \infty$ , the average of those  $s$  sample means limits (in probability) to the population mean.

# Sample statistic, relation to population statistic

This is fundamental to statistics.

- ▶ Say that a sample of 2000 people are selected from around the country and their height is measured. Mean of these 2000 heights: sample mean.
- ▶ Sample mean differs from the true mean height of the entire population of the country. (Not by much, perhaps.)
  - ▶ Weak Law of Large Numbers: if  $s$  random samples of 2000 people taken, and each sample mean calculated then, as  $s \rightarrow \infty$ , the average of those  $s$  sample means limits (in probability) to the population mean.
- ▶ Analogous thing happens with data from linear relationship with noise – think of parameters  $w^*$  and  $b^*$  as sample statistics (like the sample mean).

## Confidence intervals

How close do we suspect  $w^*$  and  $b^*$  to be to the *true* (population) slope and intercept?

---

<sup>2</sup>These formulae are just for single-variable linear regression.

## Confidence intervals

How close do we suspect  $w^*$  and  $b^*$  to be to the *true* (population) slope and intercept?

**Standard Error (SE):** Suppose that for our error term  $\varepsilon$ , we have  $\text{Var}(\varepsilon) = \sigma^2$ . Sample size:  $P$ .

---

<sup>2</sup>These formulae are just for single-variable linear regression.

## Confidence intervals

How close do we suspect  $w^*$  and  $b^*$  to be to the *true* (population) slope and intercept?

**Standard Error (SE):** Suppose that for our error term  $\varepsilon$ , we have  $\text{Var}(\varepsilon) = \sigma^2$ . Sample size:  $P$ .

Using  $\bar{x}$  for the average of  $x_1, \dots, x_P$ ,

$$(SE(w^*))^2 = \frac{\sigma^2}{\sum_{i=1}^P (x_i - \bar{x})^2};$$

$$(SE(b^*))^2 = \sigma^2 \left( \frac{1}{P} + \frac{\bar{x}^2}{\sum_{i=1}^P (x_i - \bar{x})^2} \right).$$

---

<sup>2</sup>These formulae are just for single-variable linear regression.



## Confidence intervals

How close do we suspect  $w^*$  and  $b^*$  to be to the *true* (population) slope and intercept?

**Standard Error (SE):** Suppose that for our error term  $\varepsilon$ , we have  $\text{Var}(\varepsilon) = \sigma^2$ . Sample size:  $P$ .

Using  $\bar{x}$  for the average of  $x_1, \dots, x_P$ ,

$$(SE(w^*))^2 = \frac{\sigma^2}{\sum_{i=1}^P (x_i - \bar{x})^2};$$
$$(SE(b^*))^2 = \sigma^2 \left( \frac{1}{P} + \frac{\bar{x}^2}{\sum_{i=1}^P (x_i - \bar{x})^2} \right).$$

*Roughly*, the Standard Error is the amount, on average, that  $w^*$  (resp.  $b^*$ ) differs from true slope  $w$  (resp. true intercept  $b$ ).<sup>2</sup>

---

<sup>2</sup>These formulae are just for single-variable linear regression.

## Confidence intervals

How close do we suspect  $w^*$  and  $b^*$  to be to the *true* (population) slope and intercept?

**Standard Error (SE):** Suppose that for our error term  $\varepsilon$ , we have  $\text{Var}(\varepsilon) = \sigma^2$ . Sample size:  $P$ .

Using  $\bar{x}$  for the average of  $x_1, \dots, x_P$ ,

$$(SE(w^*))^2 = \frac{\sigma^2}{\sum_{i=1}^P (x_i - \bar{x})^2};$$
$$(SE(b^*))^2 = \sigma^2 \left( \frac{1}{P} + \frac{\bar{x}^2}{\sum_{i=1}^P (x_i - \bar{x})^2} \right).$$

*Roughly*, the Standard Error is the amount, on average, that  $w^*$  (resp.  $b^*$ ) differs from true slope  $w$  (resp. true intercept  $b$ ).<sup>2</sup>

$\sigma$  is unknown, but can estimate it with **residual standard error**:

$$\hat{\sigma}^2 = RSE^2 = \frac{\sum_{i=1}^P (y_i - \hat{y}_i)^2}{P - 2}.$$

---

<sup>2</sup>These formulae are just for single-variable linear regression.

## Confidence intervals, cont'd

How close do we suspect  $w^*$  and  $b^*$  to be to the *true* (population) slope and intercept?

Formulae:

$$(SE(w^*))^2 = \frac{\sigma^2}{\sum_{i=1}^P (x_i - \bar{x})^2};$$

$$(SE(b^*))^2 = \sigma^2 \left( \frac{1}{P} + \frac{\bar{x}^2}{\sum_{i=1}^P (x_i - \bar{x})^2} \right).$$

Estimate:

$$\sigma^2 \approx RSE^2 = \frac{\sum_{i=1}^P (y_i - \hat{y}_i)^2}{P - 2}.$$

---

<sup>3</sup>95% of the time, these intervals contain population  $w$ ,  $b$ .

## Confidence intervals, cont'd

How close do we suspect  $w^*$  and  $b^*$  to be to the *true* (population) slope and intercept?

Formulae:

$$(SE(w^*))^2 = \frac{\sigma^2}{\sum_{i=1}^P (x_i - \bar{x})^2};$$

$$(SE(b^*))^2 = \sigma^2 \left( \frac{1}{P} + \frac{\bar{x}^2}{\sum_{i=1}^P (x_i - \bar{x})^2} \right).$$

Estimate:

$$\sigma^2 \approx RSE^2 = \frac{\sum_{i=1}^P (y_i - \hat{y}_i)^2}{P - 2}.$$

Can get (roughly) 95% confidence interval<sup>3</sup> with  $\pm 2SE$ :

$$(w^* - 2SE(w^*), w^* + 2SE(w^*))$$

and

$$(b^* - 2SE(b^*), b^* + 2SE(b^*)).$$

---

<sup>3</sup>95% of the time, these intervals contain population  $w$ ,  $b$ .

# Outline

Confidence intervals with linear regression

**Measuring how well LSR line fits**

Multiple variables

Polynomial fitting

# Mean Squared Error

How to measure how well the data fits to regression line?

## Mean Squared Error

How to measure how well the data fits to regression line?

In linear regression, we found “predicted”  $\hat{y}_i = b^* + \mathbf{x}_i^T \mathbf{w}^*$ ,  $1 \leq i \leq P$  so that the points  $(\mathbf{x}_1, \hat{y}_1), \dots, (\mathbf{x}_P, \hat{y}_P)$  fit a line. Could use the Mean Squared Error as our measure.

$$\text{MSE} = \frac{1}{P} \sum_{i=1}^P (\hat{y}_i - y_i)^2.$$

## Mean Squared Error

How to measure how well the data fits to regression line?

In linear regression, we found “predicted”  $\hat{y}_i = b^* + \mathbf{x}_i^T \mathbf{w}^*$ ,  $1 \leq i \leq P$  so that the points  $(\mathbf{x}_1, \hat{y}_1), \dots, (\mathbf{x}_P, \hat{y}_P)$  fit a line. Could use the Mean Squared Error as our measure.

$$\text{MSE} = \frac{1}{P} \sum_{i=1}^P (\hat{y}_i - y_i)^2.$$

- For the same sample size, the larger the MSE is the farther  $y_i$  is from  $\hat{y}_i$ , on average.



## Mean Squared Error

How to measure how well the data fits to regression line?

In linear regression, we found “predicted”  $\hat{y}_i = b^* + \mathbf{x}_i^T \mathbf{w}^*$ ,  $1 \leq i \leq P$  so that the points  $(\mathbf{x}_1, \hat{y}_1), \dots, (\mathbf{x}_P, \hat{y}_P)$  fit a line. Could use the Mean Squared Error as our measure.

$$\text{MSE} = \frac{1}{P} \sum_{i=1}^P (\hat{y}_i - y_i)^2.$$

- For the same sample size, the larger the MSE is the farther  $y_i$  is from  $\hat{y}_i$ , on average.

Closely related to RSE (residual standard error). Recall,

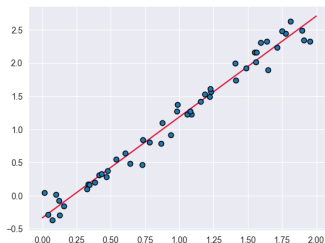
$$\text{RSE} = \sqrt{\frac{1}{P-2} \sum_{i=1}^P (y_i - \hat{y}_i)^2}.$$

$$\text{So } \text{MSE} = \frac{P-2}{P} \text{RSE}^2.$$

## Mean Squared Error, example

Recall, 'Example1.csv' data. Its best fit line is

$$y = 1.520275x - 0.33458.$$

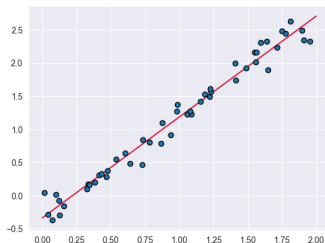


## Mean Squared Error, example

Recall, 'Example1.csv' data. Its best fit line is

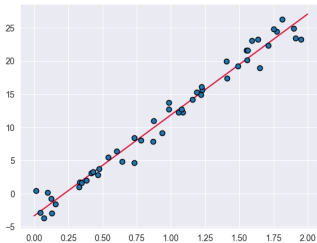
$$y = 1.520275x - 0.33458.$$

The MSE for this data and its predictions is  $\approx 0.0197$ .  
Does that mean that the linear model is a “good fit”?



## Mean Squared Error, scaling

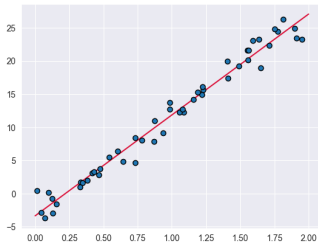
What about the following data and its best fit line? Here, the MSE is 1.9746.



Is it still a good fit?

## Mean Squared Error, scaling

What about the following data and its best fit line? Here, the MSE is 1.9746.



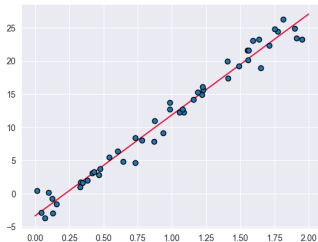
Is it still a good fit?

The data here is from 'Example1.csv' again, except that the y-coordinates have been multiplied by 10. Its regression line is

$$y = 15.20275x - 3.3458.$$

## Mean Squared Error, scaling

What about the following data and its best fit line? Here, the MSE is 1.9746.



Is it still a good fit?

The data here is from 'Example1.csv' again, except that the y-coordinates have been multiplied by 10. Its regression line is

$$y = 15.20275x - 3.3458.$$

MSE is still a good measure to think about, but its size depends on scale of y-coordinates (equivalently, depends on units y is measured in).

## $R^2$ : Proportion of “variance explained”

Get a measure that is unchanged by scaling: first, set **total sum of squares** (TSS) to

$$\text{TSS} = \sum_{i=1}^P (y_i - \bar{y})^2,$$

where  $\bar{y} = \frac{1}{P} \sum_{i=1}^P y_i$ .

## $R^2$ : Proportion of “variance explained”

Get a measure that is unchanged by scaling: first, set **total sum of squares** (TSS) to

$$\text{TSS} = \sum_{i=1}^P (y_i - \bar{y})^2,$$

where  $\bar{y} = \frac{1}{P} \sum_{i=1}^P y_i$ .

Then,

$$R^2 = \frac{\text{TSS} - \text{PMSE}}{\text{TSS}} = 1 - \frac{\sum_{i=1}^P (y_i - \hat{y}_i)^2}{\sum_{i=1}^P (y_i - \bar{y})^2}.$$



## $R^2$ : Proportion of “variance explained”

Get a measure that is unchanged by scaling: first, set **total sum of squares** (TSS) to

$$\text{TSS} = \sum_{i=1}^P (y_i - \bar{y})^2,$$

where  $\bar{y} = \frac{1}{P} \sum_{i=1}^P y_i$ .

Then,

$$R^2 = \frac{\text{TSS} - \text{PMSE}}{\text{TSS}} = 1 - \frac{\sum_{i=1}^P (y_i - \hat{y}_i)^2}{\sum_{i=1}^P (y_i - \bar{y})^2}.$$

- $R^2$  does *not* depend on the scale of the  $y$ -coordinates.

## $R^2$ : Proportion of “variance explained”

Get a measure that is unchanged by scaling: first, set **total sum of squares** (TSS) to

$$\text{TSS} = \sum_{i=1}^P (y_i - \bar{y})^2,$$

where  $\bar{y} = \frac{1}{P} \sum_{i=1}^P y_i$ .

Then,

$$R^2 = \frac{\text{TSS} - \text{PMSE}}{\text{TSS}} = 1 - \frac{\sum_{i=1}^P (y_i - \hat{y}_i)^2}{\sum_{i=1}^P (y_i - \bar{y})^2}.$$

- ▶  $R^2$  does *not* depend on the scale of the  $y$ -coordinates.
- ▶ Any data set, have  $0 \leq R^2 \leq 1$  (provided  $R^2$  is defined; i.e., we do not have  $y_1, y_2, \dots, y_n$  all the same).
  - ▶ Can you prove this?

## $R^2$ : Proportion of “variance explained”

Get a measure that is unchanged by scaling: first, set **total sum of squares** (TSS) to

$$\text{TSS} = \sum_{i=1}^P (y_i - \bar{y})^2,$$

where  $\bar{y} = \frac{1}{P} \sum_{i=1}^P y_i$ .

Then,

$$R^2 = \frac{\text{TSS} - \text{PMSE}}{\text{TSS}} = 1 - \frac{\sum_{i=1}^P (y_i - \hat{y}_i)^2}{\sum_{i=1}^P (y_i - \bar{y})^2}.$$

- ▶  $R^2$  does *not* depend on the scale of the  $y$ -coordinates.
- ▶ Any data set, have  $0 \leq R^2 \leq 1$  (provided  $R^2$  is defined; i.e., we do not have  $y_1, y_2, \dots, y_n$  all the same).
  - ▶ Can you prove this?
- ▶ Checking that  $R^2$  is “close” to 1 is often done to indicate that a linear model is a very good one.

# Outline

Confidence intervals with linear regression

Measuring how well LSR line fits

**Multiple variables**

Polynomial fitting

# Linear regression with multiple independent variables

We'll use the 'Advertising.csv' data set, found on the class GitHub site, in DataSets.

# Linear regression with multiple independent variables

We'll use the 'Advertising.csv' data set, found on the class GitHub site, in DataSets.

- Before, looked at the Sales ( $y$ ) as a function of TV (advertising budget,  $x$ ). In data set, budgets for other media: Radio and Newspaper.

# Linear regression with multiple independent variables

We'll use the 'Advertising.csv' data set, found on the class GitHub site, in DataSets.

- ▶ Before, looked at the Sales ( $y$ ) as a function of TV (advertising budget,  $x$ ). In data set, budgets for other media: Radio and Newspaper.
- ▶ Fitting Sales to each one with single-variable regression (one for TV, one for Radio, one for Newspaper) is inadequate.

# Linear regression with multiple independent variables

We'll use the 'Advertising.csv' data set, found on the class GitHub site, in DataSets.

- ▶ Before, looked at the Sales ( $y$ ) as a function of TV (advertising budget,  $x$ ). In data set, budgets for other media: Radio and Newspaper.
- ▶ Fitting Sales to each one with single-variable regression (one for TV, one for Radio, one for Newspaper) is inadequate.
  - ▶ Ignores that all are contributing together to Sales.
  - ▶ Doesn't give predictive ability that matches data.



## Working with multiple independent variables

Rather than separate single-variable linear regressions, use multivariable linear regression.

---

<sup>4</sup>When a “real world” data set with  $P \geq N + 1$ , almost surely.

## Working with multiple independent variables

Rather than separate single-variable linear regressions, use multivariable linear regression.

With  $\mathbf{x} = (x_1, x_2, \dots, x_N)$  as the variables, use the model

$$y \approx b + \mathbf{x}^T \mathbf{w} = b + w_1 x_1 + w_2 x_2 + \dots + w_N x_N.$$

---

<sup>4</sup>When a “real world” data set with  $P \geq N + 1$ , almost surely.

## Working with multiple independent variables

Rather than separate single-variable linear regressions, use multivariable linear regression.

With  $\mathbf{x} = (x_1, x_2, \dots, x_N)$  as the variables, use the model

$$y \approx b + \mathbf{x}^T \mathbf{w} = b + w_1 x_1 + w_2 x_2 + \dots + w_N x_N.$$

With Advertising data set: independent variables are  $x_1 = \text{TV}$ ,  $x_2 = \text{Radio}$ ,  $x_3 = \text{Newspaper}$ .

---

<sup>4</sup>When a “real world” data set with  $P \geq N + 1$ , almost surely.

## Working with multiple independent variables

Rather than separate single-variable linear regressions, use multivariable linear regression.

With  $\mathbf{x} = (x_1, x_2, \dots, x_N)$  as the variables, use the model

$$y \approx b + \mathbf{x}^T \mathbf{w} = b + w_1 x_1 + w_2 x_2 + \dots + w_N x_N.$$

With Advertising data set: independent variables are  $x_1 = \text{TV}$ ,  $x_2 = \text{Radio}$ ,  $x_3 = \text{Newspaper}$ .

As we discussed, set  $A$  to be  $P \times (N + 1)$  matrix with a column of ones, and a column for each independent variable. That is,

$$A = \begin{bmatrix} \vec{1} & \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_N \end{bmatrix}.$$

$(b^*, w_1^*, \dots, w_N^*)$  is the solution to normal equations:  $(A^T A)^{-1} (A^T \mathbf{y})$ .

*Note:* the matrix  $A^T A$  is invertible when  $A$  has rank  $N + 1$  (when  $\vec{1}, \mathbf{x}_1, \dots, \mathbf{x}_N$  are linearly independent).<sup>4</sup>

---

<sup>4</sup>When a “real world” data set with  $P \geq N + 1$ , almost surely.

## Working with multiple independent variables

Rather than separate single-variable linear regressions, use multivariable linear regression.

With  $\mathbf{x} = (x_1, x_2, \dots, x_N)$  as the variables, use the model

$$y \approx b + \mathbf{x}^T \mathbf{w} = b + w_1 x_1 + w_2 x_2 + \dots + w_N x_N.$$

With Advertising data set: independent variables are  $x_1 = \text{TV}$ ,  $x_2 = \text{Radio}$ ,  $x_3 = \text{Newspaper}$ .

As we discussed, set  $A$  to be  $P \times (N + 1)$  matrix with a column of ones, and a column for each independent variable. That is,

$$A = [\vec{1}, \quad \mathbf{x}_1, \quad \mathbf{x}_2, \quad \dots, \quad \mathbf{x}_N].$$

$(b^*, w_1^*, \dots, w_N^*)$  is the solution to normal equations:  $(A^T A)^{-1} (A^T \mathbf{y})$ .

Note: the matrix  $A^T A$  is invertible when  $A$  has rank  $N + 1$  (when  $\vec{1}, \mathbf{x}_1, \dots, \mathbf{x}_N$  are linearly independent).<sup>4</sup>

- Larger  $N \rightarrow$  more likely there are numerical issues computing inverse of  $A^T A$ .

<sup>4</sup>When a “real world” data set with  $P \geq N + 1$ , almost surely.

## Advertising example

$x_0$  for TV budget;    $x_1$  for Radio budget;    $x_2$  for Newspaper budget.

## Advertising example

$x_0$  for TV budget;  $x_1$  for Radio budget;  $x_2$  for Newspaper budget.

Writing  $y$  for Sales, multiple linear regression model for Advertising data is approximately

$$\hat{y} = 0.0458x_0 + 0.1885x_1 - 0.001x_2 + 2.9389.$$

## Advertising example

$x_0$  for TV budget;  $x_1$  for Radio budget;  $x_2$  for Newspaper budget.

Writing  $y$  for Sales, multiple linear regression model for Advertising data is approximately

$$\hat{y} = 0.0458x_0 + 0.1885x_1 - 0.001x_2 + 2.9389.$$

Interpretation: given fixed budget for radio and newspaper ads, increasing TV ad budget by \$1000 will increase sales by around 46 units (in each market, on average).



## Advertising example

$x_0$  for TV budget;  $x_1$  for Radio budget;  $x_2$  for Newspaper budget.  
Writing  $y$  for Sales, multiple linear regression model for Advertising data is approximately

$$\hat{y} = 0.0458x_0 + 0.1885x_1 - 0.001x_2 + 2.9389.$$

Interpretation: given fixed budget for radio and newspaper ads, increasing TV ad budget by \$1000 will increase sales by around 46 units (in each market, on average).

Contrast with result of three separate linear regressions, below.

Variable	TV	Radio	Newspaper
LSR line	$0.0475x_0 + 7.0326$	$0.2025x_1 + 9.3116$	$0.0547x_2 + 12.3514$
$R^2$	0.612	0.332	0.052

$R^2$

Previously:  $R^2$  for predicting Sales from TV significantly higher than from either Radio or Newspaper.

$R^2$ 

Previously:  $R^2$  for predicting Sales from TV significantly higher than from either Radio or Newspaper.

Can get  $R^2$  from regression model, either using 2 of the variables, or using all 3.

First, recall result from simple regression:

Independent var.	TV	Radio	Newspaper
$R^2$	0.612	0.332	0.052

$R^2$

Previously:  $R^2$  for predicting Sales from TV significantly higher than from either Radio or Newspaper.

Can get  $R^2$  from regression model, either using 2 of the variables, or using all 3.

First, recall result from simple regression:

Independent var.	TV	Radio	Newspaper
$R^2$	0.612	0.332	0.052

Now,  $R^2$  for all possible pairs of two:

Two vars.	TV, Radio	TV, Newspaper	Radio, Newspaper
$R^2$	0.89719	0.646	0.333

$R^2$ 

Previously:  $R^2$  for predicting Sales from TV significantly higher than from either Radio or Newspaper.

Can get  $R^2$  from regression model, either using 2 of the variables, or using all 3.

First, recall result from simple regression:

Independent var.	TV	Radio	Newspaper
$R^2$	0.612	0.332	0.052

Now,  $R^2$  for all possible pairs of two:

Two vars.	TV, Radio	TV, Newspaper	Radio, Newspaper
$R^2$	0.89719	0.646	0.333

The value of  $R^2$  with all three predictor (independent) variables is:

0.89721. What conclusion can we draw?

## How small to decide not significant?

---

<sup>5</sup>Recall, *SE* how far  $\hat{p}_i$  is from population coeff.  $p_i$ , on average.

## How small to decide not significant?

Hypothesis testing: choose a  $p$ -value threshold (often  $< 0.05$  or  $< 0.01$ ). The  $p$ -value corresponds to some  $t$ -statistic – use regression coefficient ( $\hat{p}_i$  for  $x_i$ ) and standard error.

- In example, if using simple linear regression on Newspaper, would get the variable is significant. However, using multiple regression with TV, Radio, and Newspaper, get very large  $p$ -value  $\rightarrow$  so, not significant.

---

<sup>5</sup>Recall, SE how far  $\hat{p}_i$  is from population coeff.  $p_i$ , on average.

## How small to decide not significant?

Hypothesis testing: choose a  $p$ -value threshold (often  $< 0.05$  or  $< 0.01$ ). The  $p$ -value corresponds to some  $t$ -statistic – use regression coefficient ( $\hat{p}_i$  for  $x_i$ ) and standard error.

- In example, if using simple linear regression on Newspaper, would get the variable is significant. However, using multiple regression with TV, Radio, and Newspaper, get very large  $p$ -value  $\rightarrow$  so, not significant.

Alternatively: if sample regression coeff. varies a lot (relative to size) compared to coeff.s of the other var's, that variable is not significant.

---

<sup>5</sup>Recall, SE how far  $\hat{p}_i$  is from population coeff.  $p_i$ , on average.



## How small to decide not significant?

Hypothesis testing: choose a  $p$ -value threshold (often  $< 0.05$  or  $< 0.01$ ). The  $p$ -value corresponds to some  $t$ -statistic – use regression coefficient ( $\hat{p}_i$  for  $x_i$ ) and standard error.

- In example, if using simple linear regression on Newspaper, would get the variable is significant. However, using multiple regression with TV, Radio, and Newspaper, get very large  $p$ -value  $\rightarrow$  so, not significant.

Alternatively: if sample regression coeff. varies a lot (relative to size) compared to coeff.s of the other var's, that variable is not significant.

- $p$ -value large when  $t$ -statistic is small, which is when  $SE$  is large *relative to size of  $\hat{p}_i$ .* <sup>5</sup>

---

<sup>5</sup>Recall,  $SE$  how far  $\hat{p}_i$  is from population coeff.  $p_i$ , on average.

## Intuitive estimate of significance

Checking whether fluctuation of regression coefficient for an independent variable, relative to coeff.'s size, is large.

---

<sup>6</sup>\* Some evidence in literature (Goodhue-Lewis, 2012) that not much precision is to be gained with more than 100 samples, for bootstrapping standard errors.

<sup>7</sup>This is an example of a bootstrapping procedure: the whole sample is used as a proxy for the population and the subsamples, or resamplings, are simulating samples from the population.

# Intuitive estimate of significance

Checking whether fluctuation of regression coefficient for an independent variable, relative to coeff.'s size, is large.

1. Take around 100 random subsamples<sup>6</sup> of data (or, resamplings with replacement); compute  $\hat{p}_i$  for those. Standard deviation of them  $\approx SE$ .

---

<sup>6</sup>\*Some evidence in literature (Goodhue-Lewis, 2012) that not much precision is to be gained with more than 100 samples, for bootstrapping standard errors.

<sup>7</sup>This is an example of a bootstrapping procedure: the whole sample is used as a proxy for the population and the subsamples, or resamplings, are simulating samples from the population.

# Intuitive estimate of significance

Checking whether fluctuation of regression coefficient for an independent variable, relative to coeff.'s size, is large.

1. Take around 100 random subsamples<sup>6</sup> of data (or, resamplings with replacement); compute  $\hat{p}_i$  for those. Standard deviation of them  $\approx SE$ .
2. Use regression coeff. from whole data set,  $\approx p_i$ . If standard dev. found in 1., divided by this coeff., is larger than about 0.5, variable is not significant.
  - ▶ Since we are *estimating some things* here, don't use as a hard cutoff. Getting 0.48, versus 0.59, would perhaps both be *weakly* significant. However, if larger than 1.5, say, definitely not significant.

---

<sup>6</sup>\*Some evidence in literature (Goodhue-Lewis, 2012) that not much precision is to be gained with more than 100 samples, for bootstrapping standard errors.

<sup>7</sup>This is an example of a bootstrapping procedure: the whole sample is used as a proxy for the population and the subsamples, or resamplings, are simulating samples from the population.

# Outline

Confidence intervals with linear regression

Measuring how well LSR line fits

Multiple variables

Polynomial fitting

## Powers of $x$ in place of multiple variables

Often, a linear model does not seem like a good fit for our data. What about trying to fit the data to a polynomial?

---

<sup>8</sup>Could do mix, multivariate regression and powers.

## Powers of $x$ in place of multiple variables

Often, a linear model does not seem like a good fit for our data. What about trying to fit the data to a polynomial?

i.e., consider the model

$$y = p_0x^d + p_1x^{d-1} + \dots + p_{d-1}x + p_d + \varepsilon$$

for some degree  $d$ , and find the coefficients which give best fit polynomial.

---

<sup>8</sup>Could do mix, multivariate regression and powers.

## Powers of $x$ in place of multiple variables

Often, a linear model does not seem like a good fit for our data. What about trying to fit the data to a polynomial?

i.e., consider the model

$$y = p_0x^d + p_1x^{d-1} + \dots + p_{d-1}x + p_d + \varepsilon$$

for some degree  $d$ , and find the coefficients which give best fit polynomial.

For the procedure, use essentially the same idea for the matrix  $A$ , but using powers of single variable  $x$  instead of using different independent variables<sup>8</sup>. Given data with  $x$ -coordinates  $x_1, x_2, \dots, x_n$ , the matrix  $A$  is known as a **Vandermonde matrix**.

---

<sup>8</sup>Could do mix, multivariate regression and powers.



## Powers of $x$ in place of multiple variables

Often, a linear model does not seem like a good fit for our data. What about trying to fit the data to a polynomial?

i.e., consider the model

$$y = p_0x^d + p_1x^{d-1} + \dots + p_{d-1}x + p_d + \varepsilon$$

for some degree  $d$ , and find the coefficients which give best fit polynomial.

For the procedure, use essentially the same idea for the matrix  $A$ , but using powers of single variable  $x$  instead of using different independent variables<sup>8</sup>. Given data with  $x$ -coordinates  $x_1, x_2, \dots, x_n$ , the matrix  $A$  is known as a **Vandermonde matrix**.

$$A = \begin{bmatrix} x_1^d & \dots & x_1^2 & x_1 & 1 \\ x_2^d & \dots & x_2^2 & x_2 & 1 \\ \vdots & & \vdots & \vdots & \\ x_n^d & \dots & x_n^2 & x_n & 1 \end{bmatrix}$$

---

<sup>8</sup>Could do mix, multivariate regression and powers.

## Example

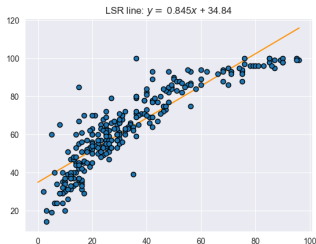
Taking the `'College.csv'` data set from the `DataSets` folder. Two of the columns are `'Top10perc'` and `'Top25perc'`. For the schools in the data set, these columns give the percentage of the entering class that were in the top 10% (resp. 25%) of their graduating high school class.<sup>9</sup>

---

<sup>9</sup>Removed rows that contained schools receiving fewer than 2500 applications.

## Example

Taking the 'College.csv' data set from the DataSets folder. Two of the columns are 'Top10perc' and 'Top25perc'. For the schools in the data set, these columns give the percentage of the entering class that were in the top 10% (resp. 25%) of their graduating high school class.<sup>9</sup> Here is the data set with a least squares line. The value of  $R^2$  is 0.791.

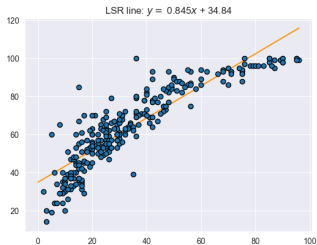


---

<sup>9</sup>Removed rows that contained schools receiving fewer than 2500 applications.

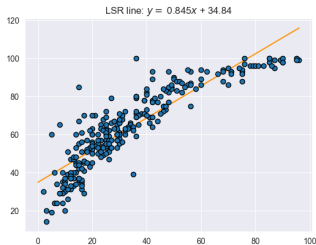
## Example

Here is the data set with a least squares line. The value of  $R^2$  is 0.791.

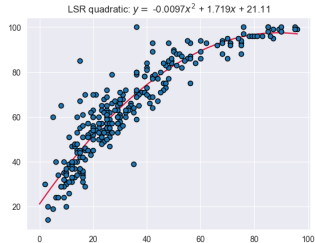


## Example

Here is the data set with a least squares line. The value of  $R^2$  is 0.791.



Next, the data set with a least squares quadratic polynomial fit. The  $R^2$  value is 0.854.



## Value of $R^2$ as polynomial degree increases

What will happen to the value of  $R^2$  if we increase the degree of the polynomial that we fit to the data?

---

<sup>10</sup>So,  $A_1$  has all the columns of  $A_0$ , and one additional column.

## Value of $R^2$ as polynomial degree increases

What will happen to the value of  $R^2$  if we increase the degree of the polynomial that we fit to the data?

- Note: Suppose that  $n > d$ . A Vandermonde matrix for  $x$ -values  $x_1, x_2, \dots, x_n$ , which has  $d + 1$  columns (so, highest power is  $x_i^d$ ), will have rank  $d + 1$  if and only if there are  $d + 1$  of the  $x_i$  that are distinct.

---

<sup>10</sup>So,  $A_1$  has all the columns of  $A_0$ , and one additional column.

## Value of $R^2$ as polynomial degree increases

What will happen to the value of  $R^2$  if we increase the degree of the polynomial that we fit to the data?

- Note: Suppose that  $n > d$ . A Vandermonde matrix for  $x$ -values  $x_1, x_2, \dots, x_n$ , which has  $d + 1$  columns (so, highest power is  $x_i^d$ ), will have rank  $d + 1$  if and only if there are  $d + 1$  of the  $x_i$  that are distinct.

*If  $x_1, x_2, \dots, x_{d+1}$  are pairwise distinct, say, then the determinant of the  $(d + 1) \times (d + 1)$  submatrix for their corresponding rows is*

$$\prod_{1 \leq i < j \leq d+1} (x_j - x_i).$$

---

<sup>10</sup>So,  $A_1$  has all the columns of  $A_0$ , and one additional column.



## Value of $R^2$ as polynomial degree increases

What will happen to the value of  $R^2$  if we increase the degree of the polynomial that we fit to the data?

- Note: Suppose that  $n > d$ . A Vandermonde matrix for  $x$ -values  $x_1, x_2, \dots, x_n$ , which has  $d + 1$  columns (so, highest power is  $x_i^d$ ), will have rank  $d + 1$  if and only if there are  $d + 1$  of the  $x_i$  that are distinct.

*If  $x_1, x_2, \dots, x_{d+1}$  are pairwise distinct, say, then the determinant of the  $(d + 1) \times (d + 1)$  submatrix for their corresponding rows is*

$$\prod_{1 \leq i < j \leq d+1} (x_j - x_i).$$

set  $A_0$ : the Vandermonde matrix used to fit polynomial of degree  $d$ ; set

$A_1$ : the one used for polynomial of degree  $d + 1$ .<sup>10</sup>

From Note, as long as enough of the  $x_i$  are distinct,

$$\text{rank}(A_1) = \text{rank}(A_0) + 1.$$

---

<sup>10</sup>So,  $A_1$  has all the columns of  $A_0$ , and one additional column.

## Value of $R^2$ as polynomial degree increases

What will happen to the value of  $R^2$  if we increase the degree of the polynomial that we fit to the data?

- Note: Suppose that  $n > d$ . A Vandermonde matrix for  $x$ -values  $x_1, x_2, \dots, x_n$ , which has  $d + 1$  columns (so, highest power is  $x_i^d$ ), will have rank  $d + 1$  if and only if there are  $d + 1$  of the  $x_i$  that are distinct.

*If  $x_1, x_2, \dots, x_{d+1}$  are pairwise distinct, say, then the determinant of the  $(d + 1) \times (d + 1)$  submatrix for their corresponding rows is*

$$\prod_{1 \leq i < j \leq d+1} (x_j - x_i).$$

set  $A_0$ : the Vandermonde matrix used to fit polynomial of degree  $d$ ; set

$A_1$ : the one used for polynomial of degree  $d + 1$ .<sup>10</sup>

From Note, as long as enough of the  $x_i$  are distinct,

$$\text{rank}(A_1) = \text{rank}(A_0) + 1.$$

Meaning:  $\text{Col}(A_0)$  is proper subspace of  $\text{Col}(A_1)$ . So, using  $A_1$  makes  $|y - \hat{y}|^2$  smaller. Since  $\sum (y - \bar{y})^2$  is unchanged, makes  $R^2$  closer to 1.

---

<sup>10</sup>So,  $A_1$  has all the columns of  $A_0$ , and one additional column.