



Lecture 22: Multiple Linear Regression



If she loves you more each and every day,
by linear regression she hated you before you met.

Announcements and reminders

- HW 5 due next Friday



If she loves you more each and every day,
by linear regression she hated you before you met.

Previously on CSCI 3022...

Given data $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, fit a simple linear regression of the form

$$y_i = \alpha + \beta x_i + \epsilon_i \text{ where } \epsilon_i \sim N(0, \sigma^2)$$

Estimates of the intercept and slope parameters are given by minimizing

$$SSE = \sum_{i=1}^n (y_i - (\alpha + \beta x_i))^2$$

$$R^2 = 1 - \frac{SSE}{SST}$$

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2$$

The least-squares estimates of the parameters are:

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}$$

and

$$\hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Previously on CSCI 3022...

We can perform inference on slope parameter β to determine if relationship is significant

$$\hat{\sigma}^2 = \frac{SSE}{n-2} \quad SE(\hat{\beta}) = \frac{\hat{\sigma}}{\sqrt{\sum_i (x_i - \bar{x})^2}} \quad \text{CI: } \hat{\beta} \pm t_{\alpha/2, n-2} \times SE(\hat{\beta})$$

We can use the Coefficient of Determination to evaluate the goodness-of-fit of SLR model

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2 \quad SSE = \sum_{i=1}^n (y_i - (\alpha + \beta x_i))^2 \quad R^2 = 1 - \frac{SSE}{SST}$$

If R^2 is close to 1, then the model fits the data relatively well.

$$R^2 \approx 0, \quad SSE \approx SST$$

Regression with multiple features

Turns out, in most practical applications, there are multiple features/predictors that potentially have an effect on the response

Example: S'pose that Y represents the sale price of a house. What are some reasonable features associated with the sale price?



Regression with multiple features

Turns out, in most practical applications, there are multiple features/predictors that potentially have an effect on the response

Example: S'pose that Y represents the sale price of a house. What are some reasonable features associated with the sale price?

- x_1 -- interior size of the house
 - x_2 -- size of the lot
 - x_3 -- number of bedrooms
 - x_4 -- number of bathrooms
 - x_5 -- age of the house
- ... and others?



Regression with multiple features

$$SLR : \begin{cases} H_0 : \beta = 0 \\ H_1 : \beta \neq 0 \end{cases}$$

Questions we would like to answer:

- Is at least one of the features useful in predicting the response?
- Do all of the feature help to explain the response? Or can we reduce to just a few?
- How well does the model fit the data? How well does just a subset of features do?
- Given a set of predictor values, what response should we predict, and how accurate is our prediction?

We'll look at these questions today and into next week

But first, let's do a little exploration of a multiple feature data set and remind ourselves of SLR

Advertising budget example

Get in groups, get out your laptops, and let's open the **nb 22** in-class notebook

Example: Data is provided about the sales of a particular product in 200 different markets, along with the advertising budgets for each market for three different media types: TV, radio and newspaper.

The sales response is given in thousands of units, and each of the advertising budget features are given in thousands of dollars.

Let's begin by fitting three individual SLR models with the advertising budget as the feature and the sales as the response.



Multiple linear regression

We've seen from the Advertising example that SLR analysis has indicated that there is a significant relationship between each of the media types (TV/radio/newspaper) and the sales of the product.

But individual SLR models only show the effect of each media type in a vacuum.

To get a clearer picture of what's going on, we want to consider the effect of all three advertising types on sales simultaneously

→ This is where Multiple Linear Regression (MLR) comes in!



Multiple linear regression

$$\text{sales} = \beta_0 + \beta_1 \text{TV} + \beta_2 \text{Radio} + \beta_3 \text{News}$$

Definition: In MLR, the data is assumed to come from a model of the form:

$p=3$ features

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \epsilon$$

SLP: $Y = \alpha + \beta x$

$$Y = \beta_0 + \beta_1 x_1$$

So for each of the n data points $(x_{i1}, x_{i2}, \dots, x_{ip}, y_i)$, for $i = 1, 2, \dots, n$, we assume:

p features

data point

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i$$

$\beta_j = \text{coefficient for the } j^{\text{th}} \text{ feature}$
 $x_{ij} = j^{\text{th}} \text{ feature from the } i^{\text{th}}$

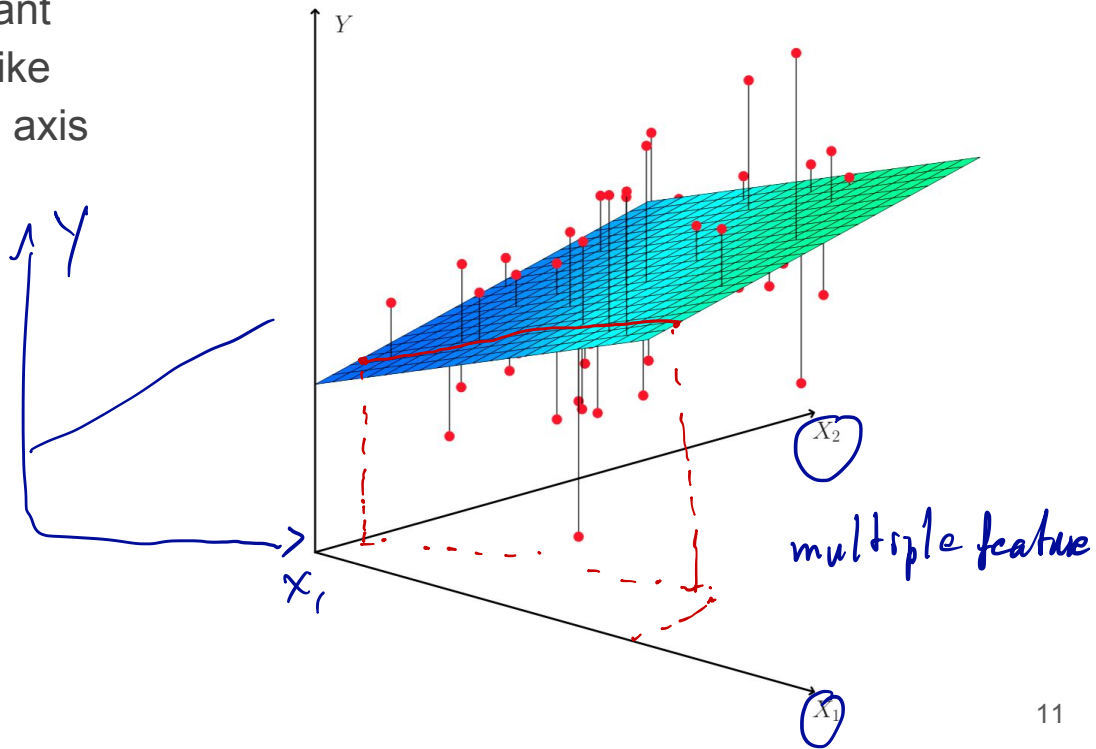
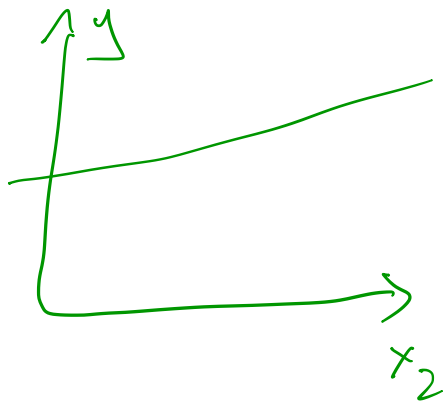
We make similar assumptions as in the case of SLR:

- Each ϵ_i is independent
- $\epsilon_i \sim N(0, \sigma^2)$

Multiple linear regression

Note that our model is no longer a simple line. Instead, it is a **linear surface**

- If you held all of the variables constant except for one of them, would look like a line as viewed from that variable's axis



Multiple linear regression

The interpretation of the model parameters is similar to that of SLR

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i$$

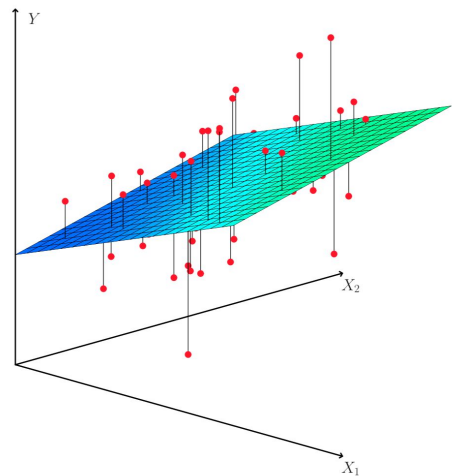
Parameter β_k is the expected change in the response associated with a unit change in the value of feature x_k **while all of the other features are held fixed**

\$ los k

Example: House sale prices:

$$Y = \frac{15 + 50x_1 + 25x_2 + 0.1x_3}{x_1 = \text{house s.q. ft. (1000 sq. ft.)}$$

$x_2 = \# \text{ bedrooms}$
 $x_3 = \# \text{ pizza that come w/ house}$



Estimating the MLR parameters

Just as in the case of SLR, we have no hope of discovering the true model parameters.

→ Need to **estimate** them from the data. Our estimated model will be:

As before, we will find the estimated parameters by minimizing the **sum of squared errors**:

The SSE is again interpreted as the measure of how much variation is left in the data that cannot be explained by the model.

Estimating the MLR parameters

Just as in the case of SLR, we have no hope of discovering the true model parameters.

→ Need to **estimate** them from the data. Our estimated model will be:

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_p x_p$$

As before, we will find the estimated parameters by minimizing the **sum of squared errors**:

$$SSE = \sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_p x_p))^2$$

1, our model will predict

The SSE is again interpreted as the measure of how much variation is left in the data that cannot be explained by the model.

coeff of data: $R^2 = 1 - \frac{SSE}{SST}$

Estimating the MLR parameters

Just as in the case of SLR, we have no hope of discovering the true model parameters.

→ Need to **estimate** them from the data. Our estimated model will be:

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_p x_p$$

As before, we will find the estimated parameters by minimizing the **sum of squared errors**:

$$SSE = \sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_1 \overset{x_{i1}}{\cancel{x_1}} + \hat{\beta}_2 \overset{x_{i2}}{\cancel{x_2}} + \dots + \hat{\beta}_p \overset{x_{ip}}{\cancel{x_p}}))^2$$

Note: Without linear algebra, it is difficult to write down an explicit expression for the parameter estimates. For now, we'll see how to obtain them in Python. Later, we can return to this problem using **stochastic gradient descent**!

Advertising budget example

Back to our groups! Let's see how we can find a MLR model for the advertising data.

Advertising budget example

Back to our groups! Let's see how we can find a MLR model for the advertising data.

... $\text{sale} \quad (TV=1, \text{Radio}=2, \text{News}=0) = 2.94 + 0.046 \cdot 1 + 0.189 \cdot 2 - 0.001 \cdot 0$

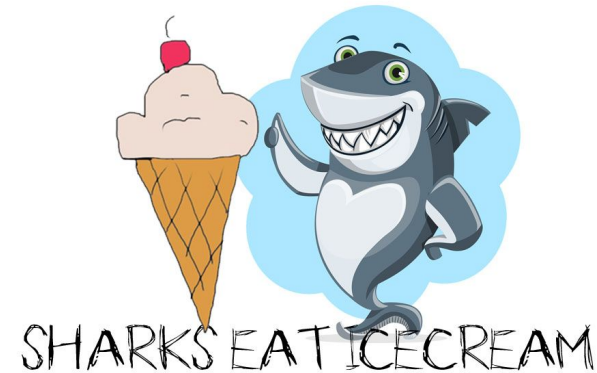
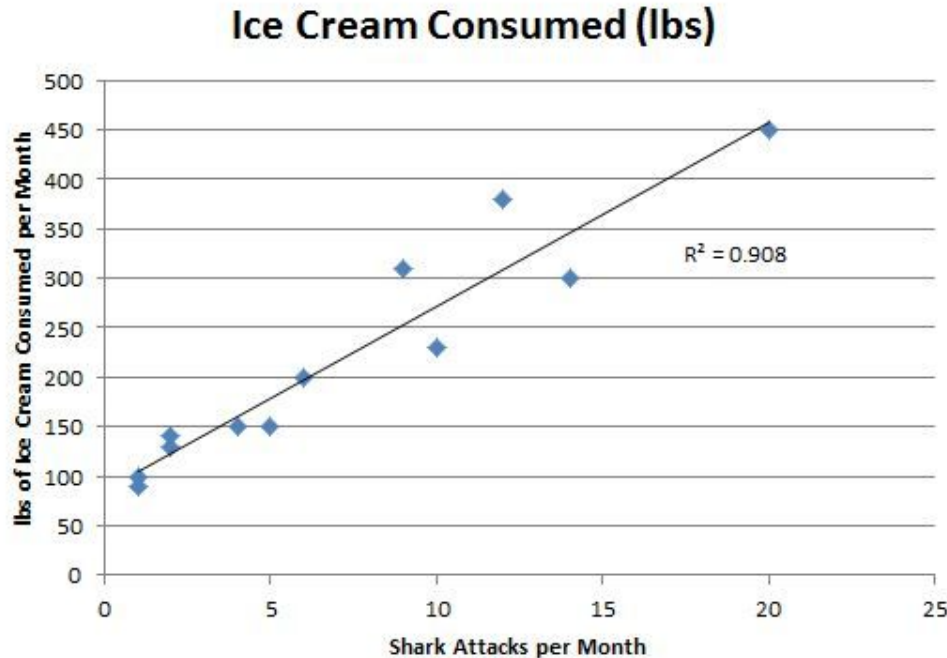
So. We see that the MLR model for the advertising data is:

$$\text{sales} = \underbrace{2.94}_{\text{Baseline}} + \underbrace{0.046}_{\text{cost how much we spend on TV}} \times \text{TV} + 0.189 \times \text{radio} - 0.001 \times \text{news}$$

Question: Why did our SLR models indicate a positive relationship between **newspaper advertising** and **sales**, but our MLR model did not?

A correlation parable

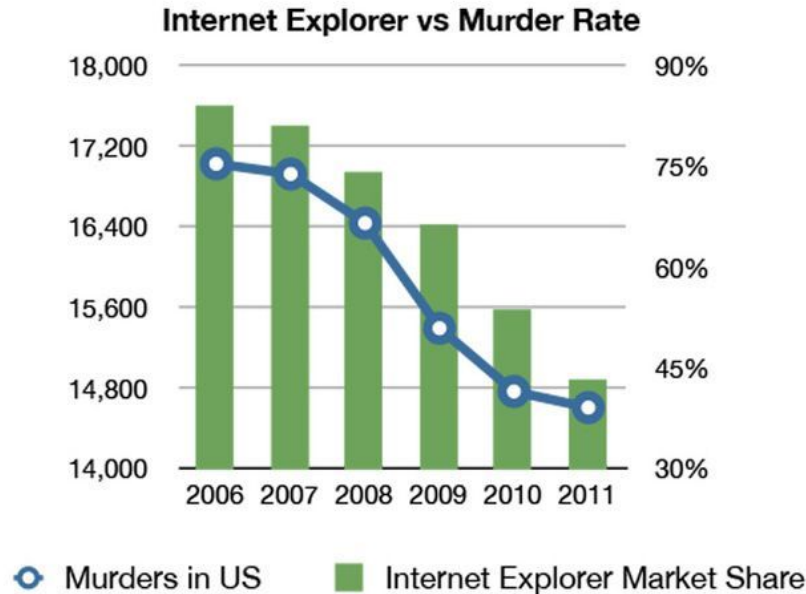
Example: A SLR analysis of shark attacks vs ice cream sales at a Southern California beach indicates that there is a strong relationship between the two.



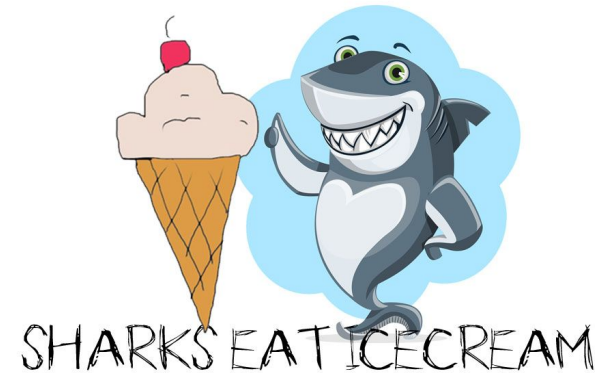
A correlation parable

Example: A SLR analysis of shark attacks vs ice cream sales at a Southern California beach indicates that there is a strong relationship between the two.

Another example: Internet Explorer use vs Murders



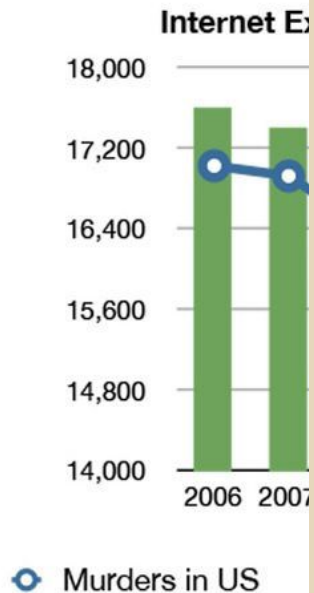
Question: Do you think that these relationships are real?



A correlation para

Example: A SLR anal
indicates that there is a

Another example: Int



What are we?



BROWSERS!



BROWSERS!



BROWSERS!



What do we want?

MORE
SPEED!



And when do we
want it?

RIGHT
NOW!!!



BROWSERS!



in California beach

think that these
hips are real?



A correlation parable

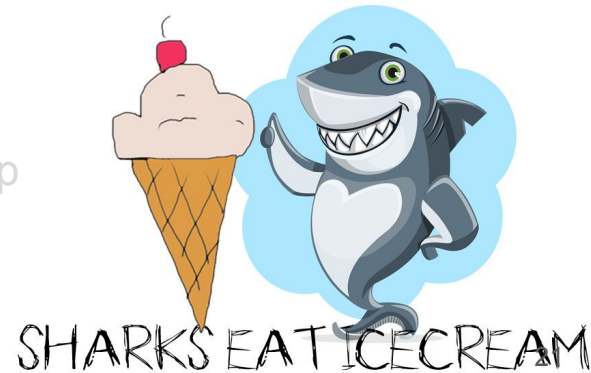
Example: A SLR analysis of shark attacks vs ice cream sales at a Southern California beach indicates that there is a strong relationship between the two.

Another example: Internet Explorer use vs Murders

Question: Do you think that these relationships are real?

Answer: Probably not. **Temperature** explains both shark attacks and ice cream consumption.

- If we did an MLR analysis with shark attacks as the response and both temperature and ice cream sales as the features, our model would show the strong relationship between temperature and shark attacks, and an insignificant relationship between ice cream sales and shark attacks.
- If we **adjust** or **control** for temperature, then the relationship between ice cream sales and shark attacks disappears

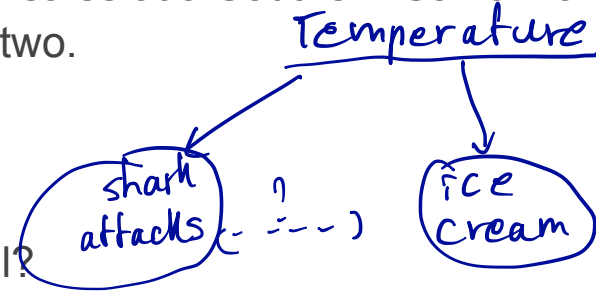


A correlation parable

Example: A SLR analysis of shark attacks vs ice cream sales at a Southern California beach indicates that there is a strong relationship between the two.

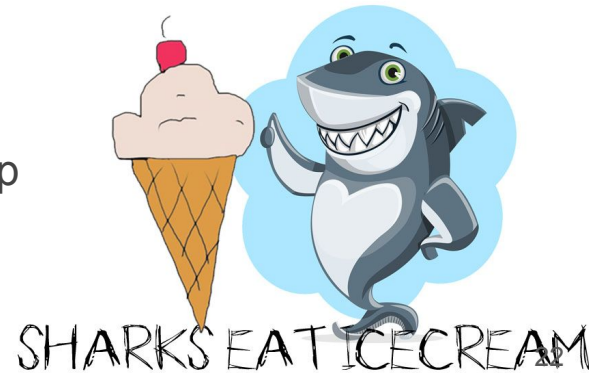
Another example: Internet Explorer use vs Murders

Question: Do you think that these relationships are real?



Answer: Probably not. **Temperature** explains both shark attacks and ice cream consumption.

- If we did an MLR analysis with shark attacks as the response and both temperature and ice cream sales as the features, our model would show the strong relationship between temperature and shark attacks, and an insignificant relationship between ice cream sales and shark attacks.
- If we **adjust** or **control** for temperature, then the relationship between ice cream sales and shark attacks disappears



Advertising budget example

Question: Based on our absurd shark attack example, can you explain why newspaper spending became less significant in our MLR of product sales?

$$\text{sales} = 2.94 + 0.046 \times \text{TV} + 0.189 \times \text{radio} - 0.001 \times \text{news}$$

controlled for radio & TV, relation b/w news & sales disappears
→ radio & TV advertising explain the response in sales
→ news doesn't give any

Covariance and correlation of features

$$\begin{aligned}\text{Var}(x) &= E[(x - E[x])^2] \\ &= E[(x - E[x])(x - E[x])]\end{aligned}$$

One way to discover these relationships among features is to do a correlation analysis.

We want to know: if the value of one features changes, how will this affect the other features?

Definition: Let X and Y be random variables. The **covariance** between X and Y is given by

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

Definition: The **correlation coefficient** $\rho(X, Y)$ is a measure between -1 and 1, and given by

$\rho(X, Y)$
1 latex

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$$

$\rho \approx 1 \rightarrow$ strongly positively correlated

$\rho \approx -1 \rightarrow$ strongly negatively correlated

$\rho \approx 0 \rightarrow$ uncorrelated.

$\nwarrow \quad \rho > 0$ $\searrow \quad \rho < 0$

Covariance and correlation of features

We can estimate these relationships from the data using formulas analogous to the sample variance

Definition: The sample covariance is given by

Definition: The sample correlation coefficient is given by

Covariance and correlation of features

We can estimate these relationships from the data using formulas analogous to the sample variance

Definition: The sample covariance is given by

$$S_{XY}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

✓ look exactly like the variance s_x^2
in sample x

Definition: The sample correlation coefficient is given by

$$\hat{\rho}(X, Y) = \frac{S_{XY}^2}{\sqrt{S_X^2 S_Y^2}}$$

↳ $s_x^2 = \frac{1}{n_x-1} \sum_{i=1}^n (x_i - \bar{x})^2$

Advertising budget example

Let's compute the pairwise correlation coefficients for the TV, radio and newspaper spending features in the advertising data

```
In [9]: dfAd[["tv", "radio", "news"]].corr()
```

```
Out[9]:
```

	tv	radio	news
tv	1.000000	0.054809	0.056648
radio	0.054809	1.000000	0.354104
news	0.056648	0.354104	1.000000

$p(x, y)$

$p(\text{radio}, \text{news})$

Question: What do you notice?

radio & news are correlated!
tv & radio explained the variance by on sales
: news had nothing else to add.

Polynomial regression

Now we will look at how we can use MLR to explain **nonlinear** relationships between single-feature data and the response

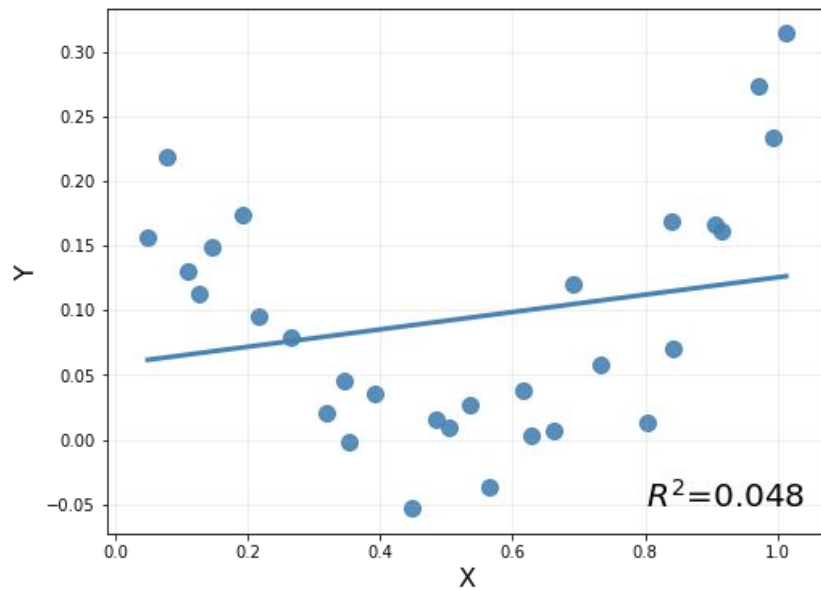
look alot more like a parabola!

$$y = \beta_0 + \beta_1 x + \beta_2 x^2$$

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

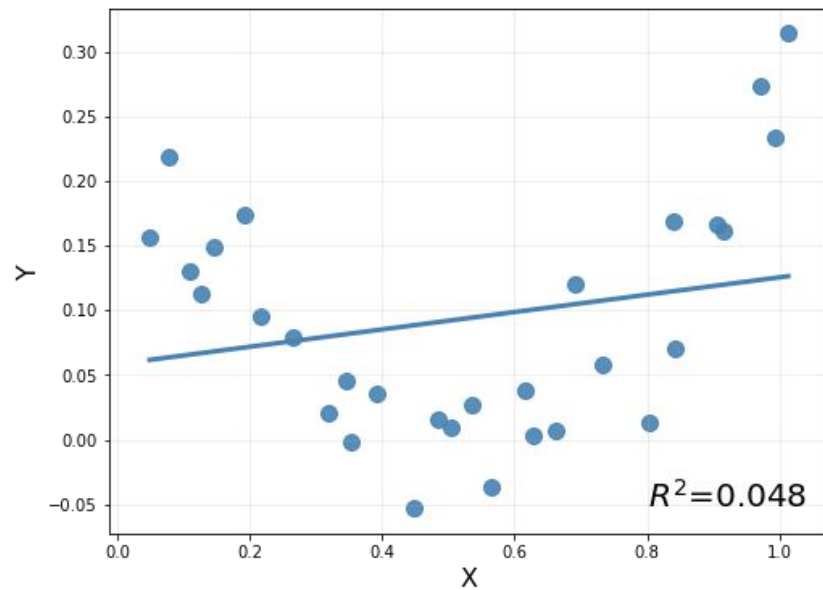
x	x^2	y
1	1	⋮
2	4	⋮
5	25	⋮
⋮	⋮	⋮

$$x_2 \approx x^2$$



Polynomial regression

For single-feature data, we can fit a polynomial regression model by casting it as a multiple linear regression, where the additional features are **powers** of the original single feature, x



Residual plots, in polynomial regression

$$\varepsilon \sim N(0, \sigma^2)$$

Recall that the assumed nature of our true model is: $Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon$

if true model $\Rightarrow Y = \beta_0 + \beta_1 x + \varepsilon$

& we assume: $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Residuals: $r = Y - \hat{Y} = \varepsilon \sim N(0, \sigma^2)$

if true model: $Y = \beta_0 + \beta_1 x + \beta_2 x^2 + \varepsilon$

& we assume: $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x$

residual: $r = y - \hat{y} = \beta_2 x^2 + \varepsilon \sim N(\beta_2 x^2, \sigma^2)$

residuals centred around the missing feature

Residual plots, in polynomial regression

Recall that the assumed nature of our true model is:

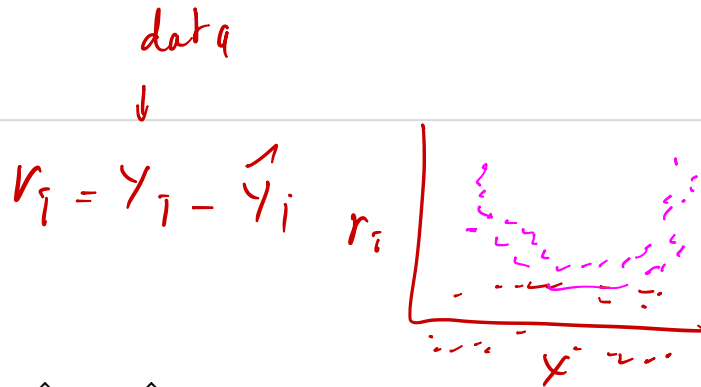
$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_p x^p + \epsilon$$

If true model is $y = \beta_0 + \beta_1 x + \epsilon$, and our model is $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x + \epsilon$

$$\text{then } r = y - \hat{y} \sim N(0, \sigma^2)$$

If true model is $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \epsilon$, and our model is $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x + \epsilon$

$$\text{then } r = y - \hat{y} \sim N(\beta_2 x^2, \sigma^2)$$



Residual plots, in polynomial regression

Recall that the assumed nature of our true model is:

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_p x^p + \epsilon$$

If true model is $y = \beta_0 + \beta_1 x + \epsilon$, and our model is $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x + \epsilon$

$$\text{then } r = y - \hat{y} \sim N(0, \sigma^2)$$

If true model is $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \epsilon$, and our model is $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x + \epsilon$

$$\text{then } r = y - \hat{y} \sim N(\beta_2 x^2, \sigma^2)$$

In general: If you plot the residuals $r = y - \hat{y}$, they should be **normally distributed around the missing feature**. So add that to your model!

Okay! Let's get (back) to work!

Get in groups, get out laptops, and open **nb 22** notebook back up

Let's...

- Dig into comparisons of **multiple linear regression** versus **simple linear regression**!
- Get some practice working with **polynomial regression**!
- Get some practice guessing **polynomial features** using **residual plots**!
- Figure out why **Internet Explorer use causes murders!** (not really)

