

Introduction to Data Science With Probability and Statistics

Lecture 23: Regression Continued & MLR

CSCI 3022 - Summer 2020

Sourav Chakraborty

Dept. of Computer Science

University of Colorado Boulder

Simple Linear Regression (SLR)

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

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

Estimates of the intercept and slope parameters are given by minimizing

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

The least-squares estimates of the parameters are:

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x} \quad \text{and} \quad \hat{\beta} = \frac{\bar{x} \bar{y} - \overline{xy}}{(\bar{x})^2 - \bar{x}^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Simple Linear Regression (SLR)

$$\begin{aligned}\beta &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \\&= \frac{\sum_{i=1}^n (x_i y_i - x_i \bar{y} - \bar{x} y_i + \bar{x} \bar{y})}{\sum_{i=1}^n (x_i^2 - 2x_i \bar{x} + (\bar{x})^2)} \\&= \frac{n \bar{x} \bar{y} - n \bar{x} \bar{y} - n \bar{x} \bar{y} + n \bar{x} \bar{y}}{n \bar{x}^2 - 2n(\bar{x})^2 + n (\bar{x})^2} \\&= \frac{n \bar{x} \bar{y} - n \bar{x} \bar{y}}{n \bar{x}^2 - n(\bar{x})^2} \\&= \frac{\bar{x} \bar{y} - \bar{x} \bar{y}}{\bar{x}^2 - (\bar{x})^2}\end{aligned}$$



Estimating the variance

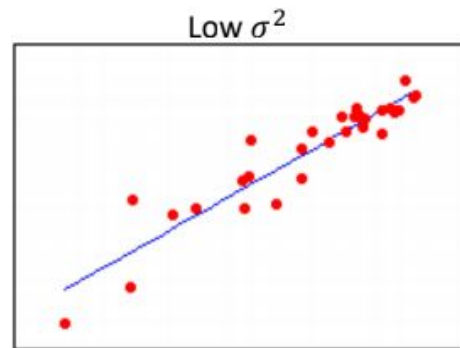
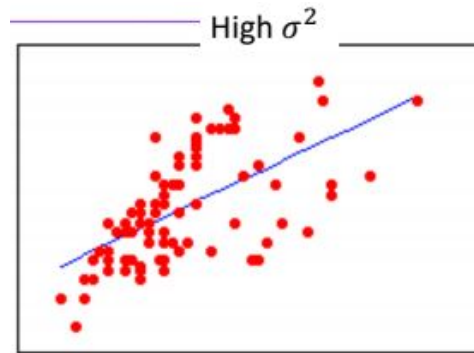
The parameter σ^2 determines the spread of the data about the true regression line.

An estimate of σ^2 will be used in computing confidence intervals and doing hypothesis testing on the estimated regression parameters.

Recall that the sum of squared errors is given by:

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Our estimate of the variance $\hat{\sigma}^2$ is given by: $\hat{\sigma}^2 = \frac{SSE}{n-2}$



The Coefficient of Determination

The coefficient of determination, R^2 , quantifies how well the model explains the data.

R^2 is a value between 0 and 1

The **sum of squared errors (SSE)** can be thought of as a measure of how much variation in Y is left unexplained by the model.

The **regression sum of squares** is given by: $SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$

It gives a sense of how much variation in Y is explained by our model.

A quantitative measure of the total amount of variation in observed Y values is given by the **total sum of squares**: $SST = \sum_{i=1}^n (y_i - \bar{y})^2$

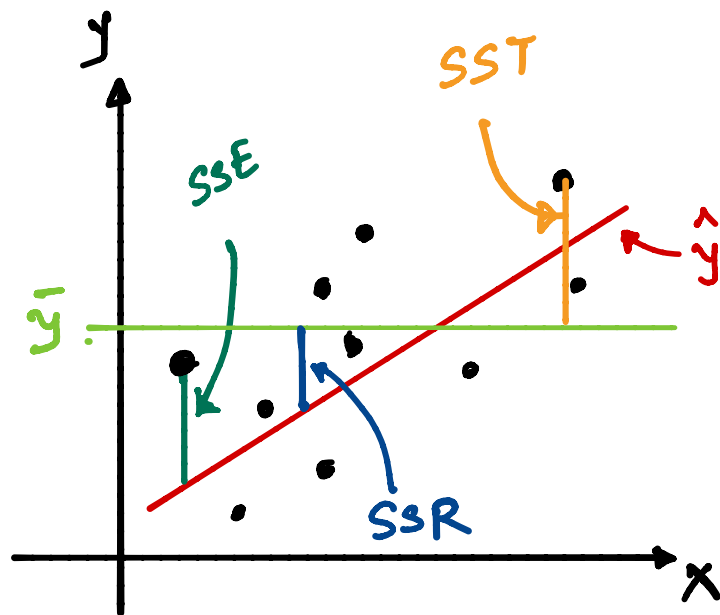
SST is what we would get for SSE if we used the mean of the data as our model.

The Coefficient of Determination

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

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

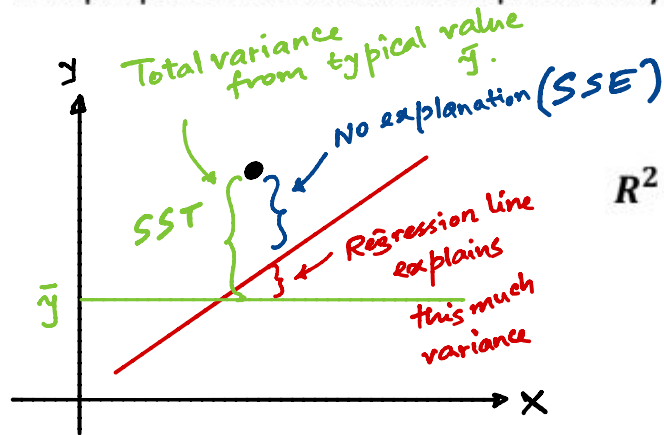


The Coefficient of Determination

The sum of squared deviations about the least-squares line is smaller than the sum of squared deviations about any other line.

$$SSE \leq SST$$

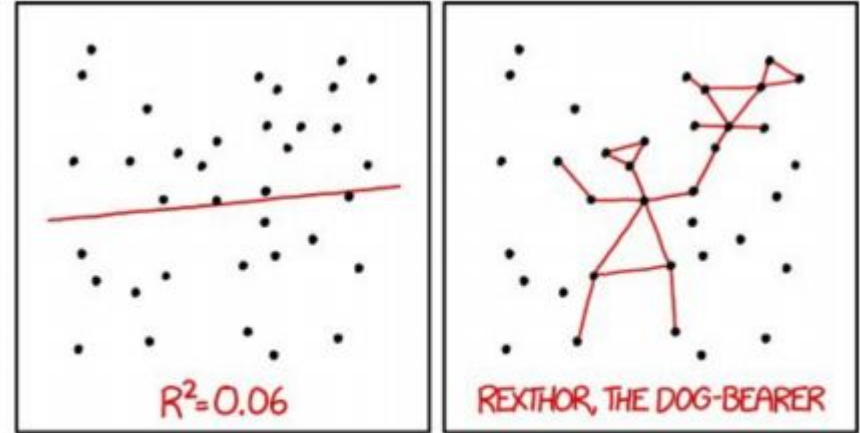
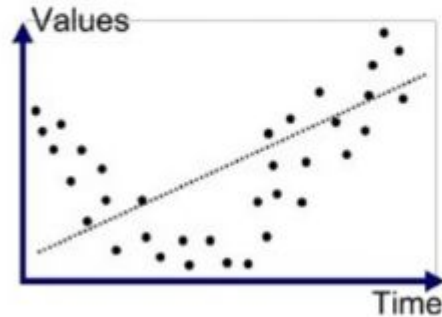
The ratio SSE/SST is the proportion of total variation in the data (SST) that cannot be explained by the SLR model (SSE). So we define the coefficient of determination R^2 to be the proportion that can be explained by the model.



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

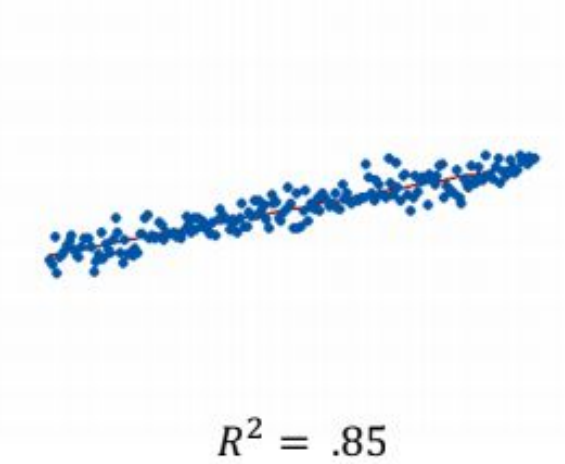
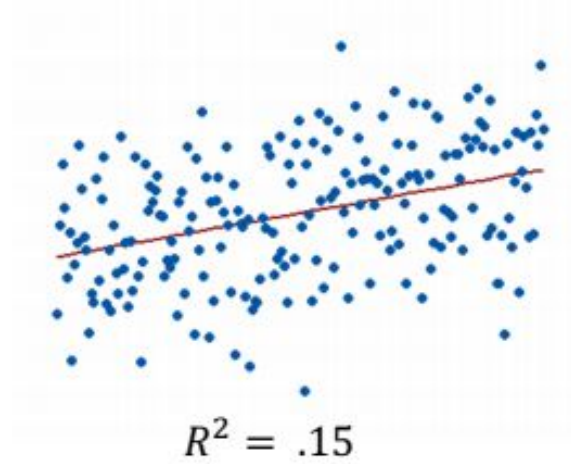
The Coefficient of Determination

R^2 is the proportion of total variation in the data that is explained by the model. It does not tell you whether you have the correct model.



I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER TO GUESS THE DIRECTION OF THE CORRELATION FROM THE SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

The Coefficient of Determination



Regression with Multiple Features

In most practical applications, there are multiple features/predictors that potentially have an effect on the response.

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



Regression with Multiple Features

In most practical applications, there are multiple features/predictors that potentially have an effect on the response.

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

y : sale price

x_1 : interior size

x_2 : size of the lot

x_3 : # Bedrooms.

x_4 : distance from
airport etc.



Regression with Multiple Features

Questions we would like to answer:

- Is at least one of the features useful in predicting the response?
- Do all of the features 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?



Multiple Linear Regression


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

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

For each of the n data points $(x_{i1}, x_{i2}, \dots, x_{ip}, y_i)$, for $i = 1, 2, 3, \dots, n$, we assume:

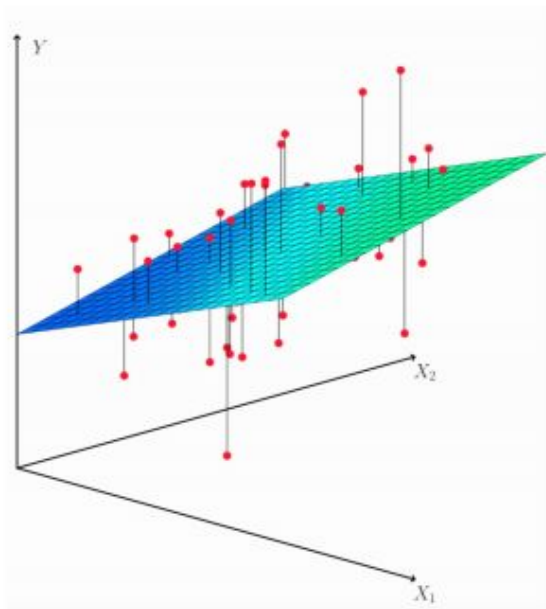
$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i$$

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

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, it 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 the 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.

Example: Consider this model for House sale prices:

$$Y = 15 + 50x_1 + 25x_2 + 0.1x_3$$

where $x_1 = \text{house sq. feet}$, $x_2 = \# \text{ bedrooms}$,
 $x_3 = \# \text{ new appliances (appliances)}$



Every increase in house sq. feet by 1 sq. feet,
Sale price increases by \$50.

Estimating the MLR parameters

Just as in the case of SLR, we likely will not discover the true model parameters, We need to estimate them from the data. Our estimated model will be:

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

$$\hat{y} = \hat{\alpha} + \hat{\beta} x$$

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

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

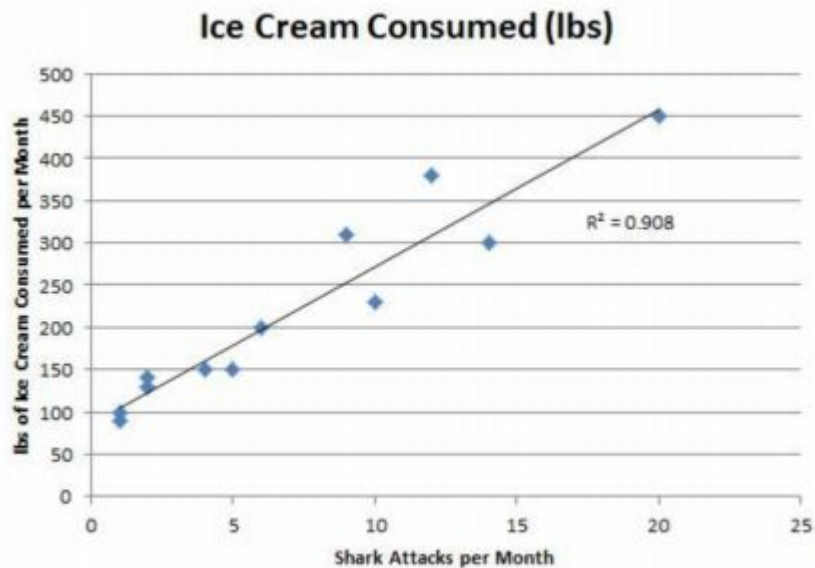
The SSE is interpreted as the measure of how much variation is left in the data that cannot be explained by the model.



Correlation?

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

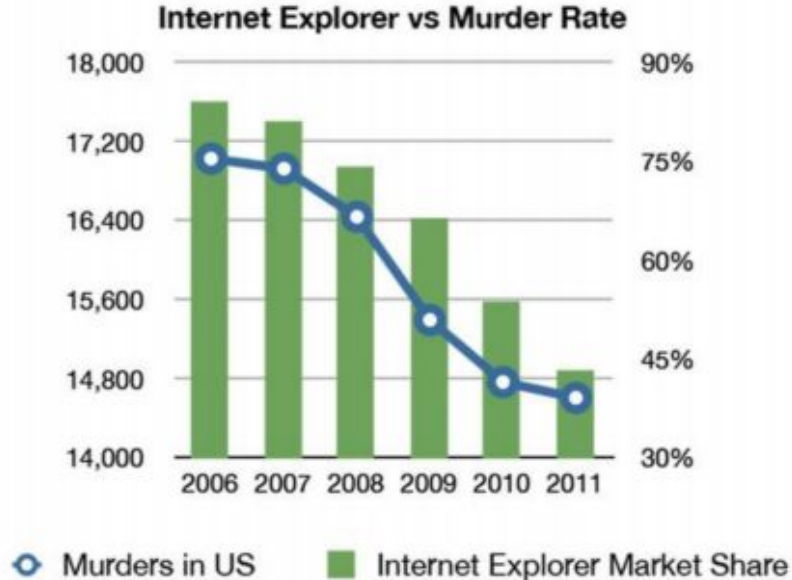
Do you think this relationship is real?



Correlation?

Example: Another study that examined internet explorer usage vs. murder rate is shown below.

Do you think this relationship is real?



Correlation?

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

This relationship is probably not real.

- 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.



Covariance and correlation of features

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

If the value of one feature changes, how will this affect the other features?

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

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

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

Covariance and correlation of features

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

The sample covariance is given by: $S_{XY}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$

The sample correlation coefficient is given by: $\hat{\rho}(X, Y) = \frac{S_{XY}^2}{\sqrt{S_X^2 S_Y^2}}$

$\rho \approx 1$: Strong +ve correlation
 $\rho \approx -1$: Strong -ve correlation
 $\rho \approx 0$: No correlation.

Covariance and correlation of features

Example: Looking ahead to Notebook 15. Suppose we have a data frame with data corresponding to TV, radio, and newspaper spending features as it pertains to advertising.

```
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