

Regression

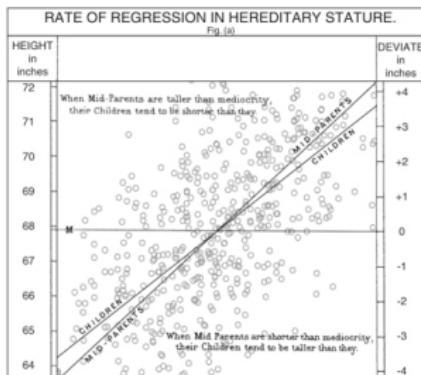
Schwartz

July 22, 2017

The Sophomore Slump

...or *sophomore jinx* or *sophomore jitters* refers to an instance in which a second, or sophomore, effort fails to live up to the standards of the first effort. It is commonly used to refer to the apathy of students (second year of high school, college or university), the performance of athletes (second season of play), singers/bands (second album), television shows (second seasons) and films (sequels/prequels). In the United Kingdom the *sophomore slump* is more commonly referred to as *second year blues*, particularly when describing university students. And in Australia it is known as *second year syndrome*, and is particularly common when referring to professional athletes who have a mediocre second season following a stellar debut. The phenomenon of a sophomore slump can be explained psychologically, where earlier success has a reducing effect on the subsequent effort, but it can also be explained statistically, as an effect of the regression towards the mean.

The concept of *regression* comes from genetics and was popularized by Sir Francis Galton's late 19th century publication of "Regression towards mediocrity in hereditary stature." Galton observed that extreme characteristics (e.g., height) in parents are not completely passed on to offspring, but rather the characteristics in the offspring "regress" towards a mediocre point. By measuring the heights of hundreds of people Galton was able to quantify this "regression" and in so doing invented linear regression analysis, thus laying the groundwork for much of modern statistical modeling. The term *regression* stuck.



Objectives

- ▶ Linear Model Regression
 - ▶ Terminology
 - ▶ Model Fitting (Least Squares)
 - ▶ Diagnostics (Evaluation and Critiquing)
- ▶ Multiple (not Multivariate) Linear Regression
 - ▶ Assumptions
 - ▶ Normal Distribution Theory
 - ▶ Model Selection
 - ▶ Coefficient Testing
- ▶ Alternatives to linear forms

Linear Models and Regression Terminology

- $Y_i = \beta_0 + x_i\beta_1 + \epsilon_i, \epsilon_i \stackrel{i.i.d.}{\sim} Normal(0, \sigma^2)$

Linear Models and Regression Terminology

Outcome / Response / Label / Dependent/Endogenous Var.

- $Y_i = \beta_0 + x_i\beta_1 + \epsilon_i, \epsilon_i \stackrel{i.i.d.}{\sim} Normal(0, \sigma^2)$

Linear Models and Regression Terminology

Outcome / Response / Label / Dependent/Endogenous Var.

► $Y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \epsilon_i \stackrel{i.i.d.}{\sim} Normal(0, \sigma^2)$

Feature / Covariate / Independent/Exogenous Var.

Linear Models and Regression Terminology

Outcome / Response / Label / Dependent/Endogenous Var.

► $Y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \epsilon_i \stackrel{i.i.d.}{\sim} Normal(0, \sigma^2)$

Feature / Covariate / Independent/Exogenous Var.

I don't like to call these *Predictors*...

Linear Models and Regression Terminology

- $Y_i = \beta_0 + x_i\beta_1 + \epsilon_i, \epsilon_i \stackrel{i.i.d.}{\sim} Normal(0, \sigma^2)$
Intercept

Linear Models and Regression Terminology

Coefficient

$$\blacktriangleright Y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \epsilon_i \stackrel{i.i.d.}{\sim} Normal(0, \sigma^2)$$

Intercept

Linear Models and Regression Terminology

Coefficient

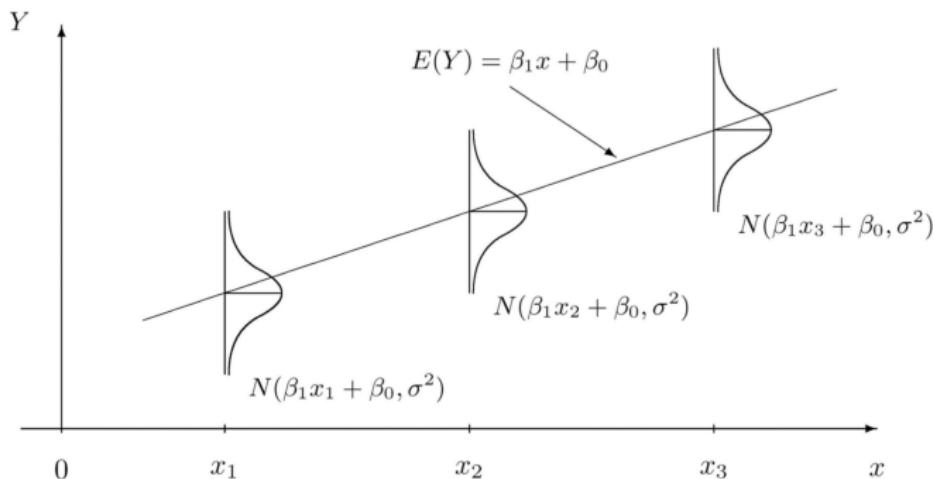
- $Y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \epsilon_i \stackrel{i.i.d.}{\sim} Normal(0, \sigma^2)$
Intercept Error/Noise

Linear Models and Regression Terminology

Coefficient

$$\blacktriangleright Y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \epsilon_i \stackrel{i.i.d.}{\sim} \text{Normal}(0, \sigma^2)$$

Intercept Error/Noise

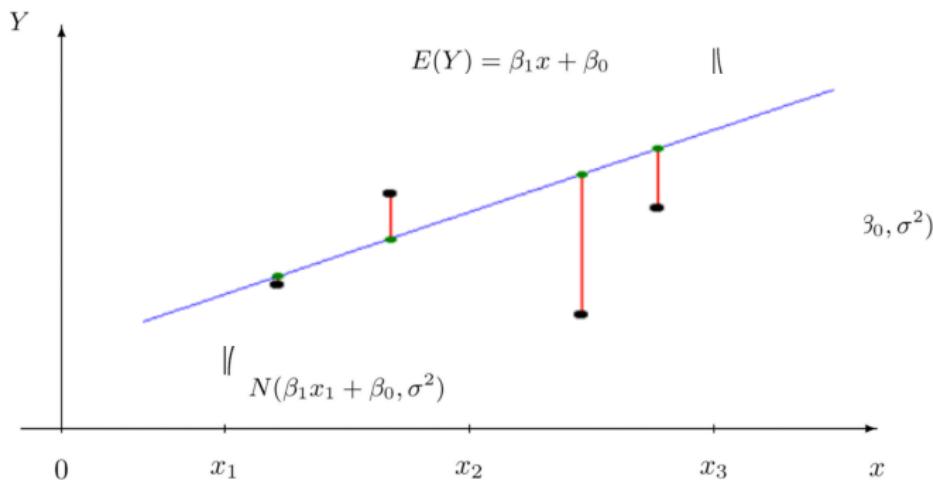


Linear Models and Regression Terminology

Coefficient

$$\blacktriangleright Y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \epsilon_i \stackrel{i.i.d.}{\sim} \text{Normal}(0, \sigma^2)$$

Intercept Error/Noise



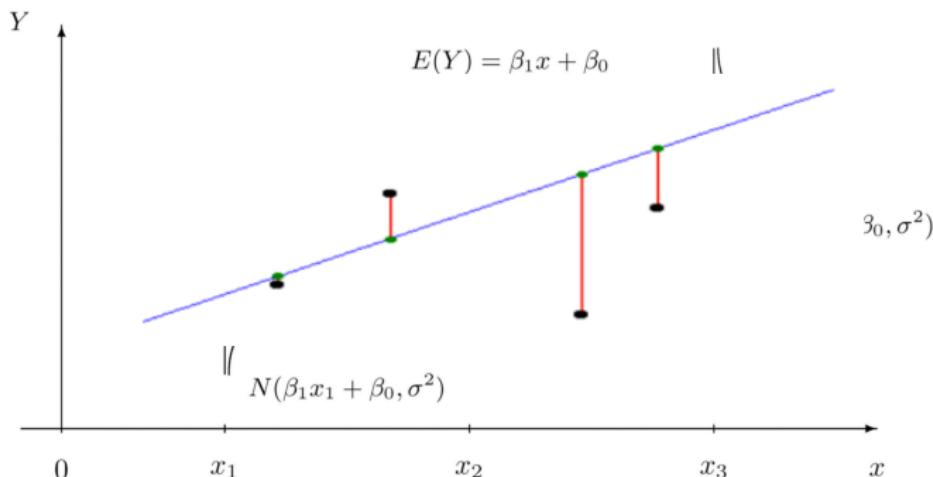
$$\blacktriangleright Y_i = \hat{\beta}_0 + x_i \hat{\beta}_1 + \hat{\epsilon}_i, \quad \hat{\sigma}^2 = \frac{\sum_{i=1}^n \hat{\epsilon}_i^2}{n-p-1} \quad (p = \# \text{of coefficients})$$

Linear Models and Regression Terminology

Coefficient

$$\blacktriangleright Y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \epsilon_i \stackrel{i.i.d.}{\sim} \text{Normal}(0, \sigma^2)$$

Intercept Error/Noise



Fitted/Predicted value \hat{Y}_i

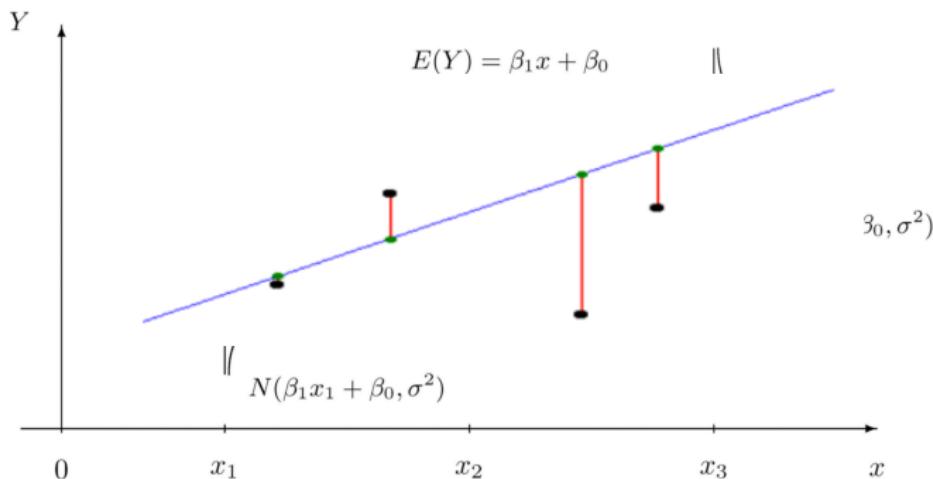
$$\blacktriangleright \hat{Y}_i = \hat{\beta}_0 + x_i \hat{\beta}_1 + \hat{\epsilon}_i, \quad \hat{\sigma}^2 = \frac{\sum_{i=1}^n \hat{\epsilon}_i^2}{n-p-1} \quad (p = \# \text{of coefficients})$$

Linear Models and Regression Terminology

Coefficient

$$\blacktriangleright Y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \epsilon_i \stackrel{i.i.d.}{\sim} \text{Normal}(0, \sigma^2)$$

Intercept Error/Noise



Fitted/Predicted value \hat{Y}_i

$$\blacktriangleright Y_i = \hat{\beta}_0 + x_i \hat{\beta}_1 + \hat{\epsilon}_i, \quad \hat{\sigma}^2 = \frac{\sum_{i=1}^n \hat{\epsilon}_i^2}{n-p-1} \quad (p = \# \text{of coefficients})$$

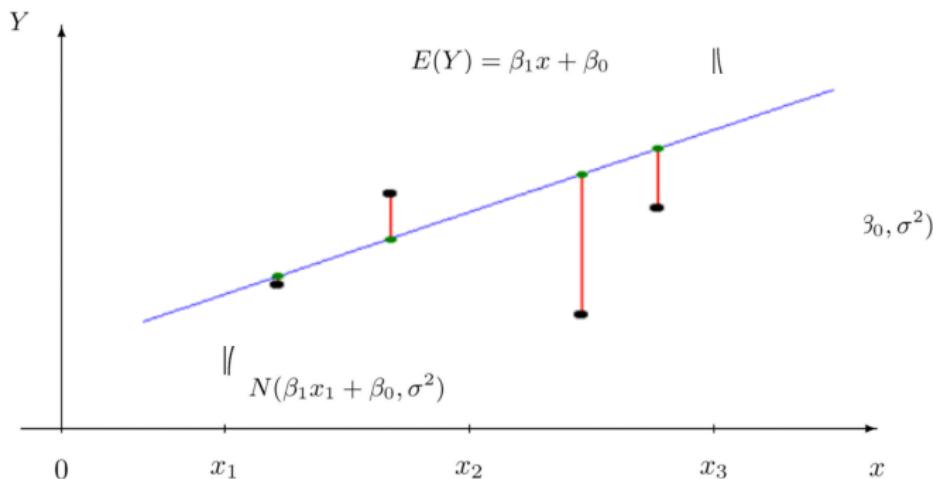
Residual

Linear Models and Regression Terminology

Coefficient

$$\blacktriangleright Y_i = \beta_0 + x_i \beta_1 + \epsilon_i, \epsilon_i \stackrel{i.i.d.}{\sim} \text{Normal}(0, \sigma^2)$$

Intercept Error/Noise



Fitted/Predicted value \hat{Y}_i Residual Variance

$$\blacktriangleright Y_i = \hat{\beta}_0 + x_i \hat{\beta}_1 + \hat{\epsilon}_i, \quad \hat{\sigma}^2 = \frac{\sum_{i=1}^n \hat{\epsilon}_i^2}{n-p-1} \quad (p = \# \text{of coefficients})$$

Residual

Quiz: what are these things and their parts?

$$Y_i = \beta_0 + x_i \beta_1 + \epsilon_i$$

$$Y_i = \hat{\beta}_0 + x_i \hat{\beta}_1 + \hat{\epsilon}_i, \quad \hat{\sigma}^2 = \frac{\sum_{i=1}^n \hat{\epsilon}_i}{n - p - 1}$$

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

Least Squares Fit

- $Y_i = \hat{\beta}_0 + x_i \hat{\beta}_1 + \hat{\epsilon}_i$

Least Squares Fit

► $Y_i = \hat{\beta}_0 + x_i \hat{\beta}_1 + \hat{\epsilon}_i$

$$[\hat{\beta}_0, \hat{\beta}_1] = \underset{[\beta_0, \beta_1]}{\operatorname{argmin}} \sum_{i=1}^n \hat{\epsilon}_i^2 = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (Y_i - \mathbf{x}_i^T \beta)^2$$

where $\mathbf{x}_i^T = [1, x_i]$ and $\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$

Least Squares Fit

► $Y_i = \hat{\beta}_0 + x_i \hat{\beta}_1 + \hat{\epsilon}_i$

$$[\hat{\beta}_0, \hat{\beta}_1] = \underset{[\beta_0, \beta_1]}{\operatorname{argmin}} \sum_{i=1}^n \hat{\epsilon}_i^2 = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (Y_i - \mathbf{x}_i^T \beta)^2$$

where $\mathbf{x}_i^T = [1, x_i]$ and $\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$

$$\sum_{i=1}^n (Y_i - \mathbf{x}_i^T \beta)^2 = (\mathbf{Y} - \mathbf{x}\beta)^T (\mathbf{Y} - \mathbf{x}\beta)$$

where $\mathbf{Y}^T = [Y_1, Y_2, \dots, Y_n]$ and $\mathbf{x}^T = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \end{bmatrix}$

Least Squares Fit

► $Y_i = \hat{\beta}_0 + x_i \hat{\beta}_1 + \hat{\epsilon}_i$

$$[\hat{\beta}_0, \hat{\beta}_1] = \underset{[\beta_0, \beta_1]}{\operatorname{argmin}} \sum_{i=1}^n \hat{\epsilon}_i^2 = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (Y_i - \mathbf{x}_i^T \beta)^2$$

where $\mathbf{x}_i^T = [1, x_i]$ and $\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$

$$\sum_{i=1}^n (Y_i - \mathbf{x}_i^T \beta)^2 = (\mathbf{Y} - \mathbf{x}\beta)^T (\mathbf{Y} - \mathbf{x}\beta)$$

where $\mathbf{Y}^T = [Y_1, Y_2, \dots, Y_n]$ and $\mathbf{x}^T = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \end{bmatrix}$

$$\nabla_{\beta} \beta^T (\mathbf{x}^T \mathbf{x}) \beta - 2 \mathbf{Y}^T \mathbf{x} \beta + \mathbf{Y}^T \mathbf{Y}$$

$$= 2(\mathbf{x}^T \mathbf{x}) \beta - 2 \mathbf{Y}^T \mathbf{x} \quad (\text{set to } \mathbf{0} \text{ to minimize})$$

$$\implies \hat{\beta} = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{Y} \implies \text{fitted values } \hat{\mathbf{Y}} = \mathbf{x} \hat{\beta}$$

$$\hat{\mathbf{Y}} = \mathbf{x} (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{Y}$$

Least Squares Fit *bonus*

1. Maximum likelihood estimation (MLE) \iff to least squares!

$$\begin{aligned} & \underset{\beta}{\operatorname{argmax}} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(Y_i - \mathbf{x}_i^T \beta)^2} \\ &= \underset{\beta}{\operatorname{argmax}} (2\pi\sigma^2)^{-\frac{n}{2}} e^{-\frac{1}{2\sigma^2}(\mathbf{Y} - \mathbf{x}\beta)^T(\mathbf{Y} - \mathbf{x}\beta)} \\ &= \underset{\beta}{\operatorname{argmax}} -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2}(\mathbf{Y} - \mathbf{x}\beta)^T(\mathbf{Y} - \mathbf{x}\beta) \\ &= \underset{\beta}{\operatorname{argmin}} (\mathbf{Y} - \mathbf{x}\beta)^T(\mathbf{Y} - \mathbf{x}\beta) \quad [\text{same as least squares!!}] \end{aligned}$$

2. In simple linear regression the $\underset{\beta}{\operatorname{argmin}} (\mathbf{Y} - \mathbf{x}\beta)^T(\mathbf{Y} - \mathbf{x}\beta)$ is

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (Y_i - \bar{Y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{R_{xY} S_Y}{S_x}$$

Assumptions of Linear Regression

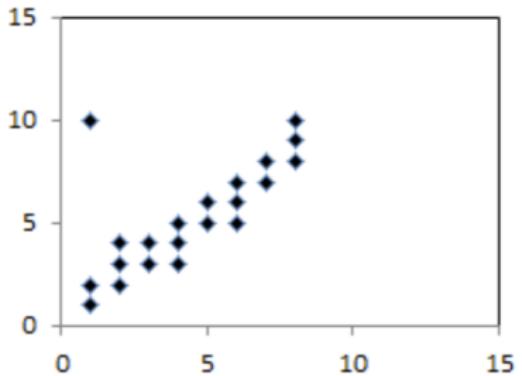
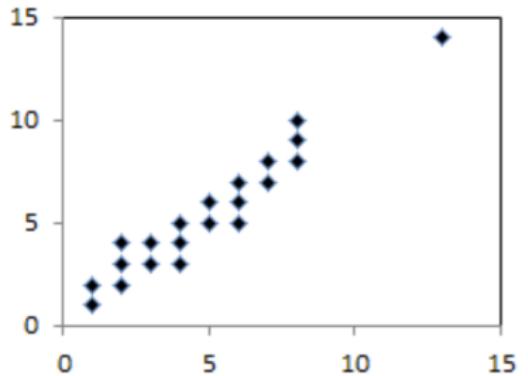
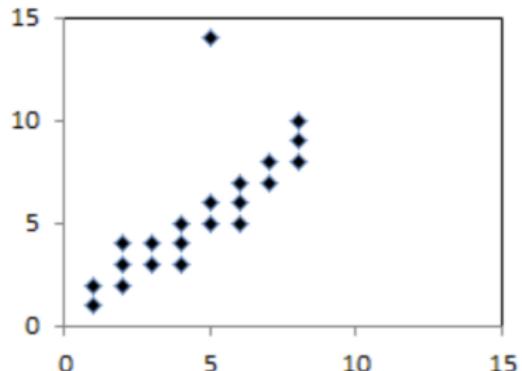
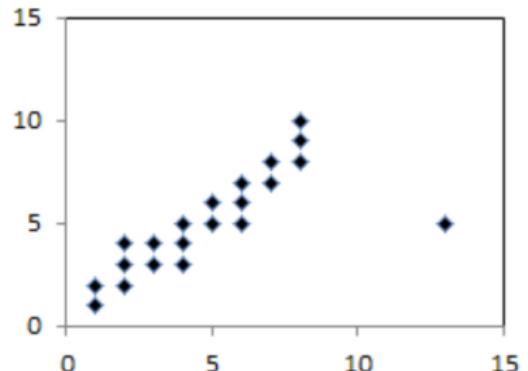
1. Sample data is representative of the population.
2. True relationship between X and y is linear.
3. Feature matrix X has full rank (rows and columns are linearly independent).
4. Residuals are independent.
5. Residuals are normally distributed.
6. Variance of the residuals is constant (homoscedastic).

Note: Linear regression does *not* assume anything about the distributions of x and y , it only makes assumptions about the distribution of the residuals, and this is all that's needed for the statistical tests to be valid.

Colinearity of the Feature Matrix

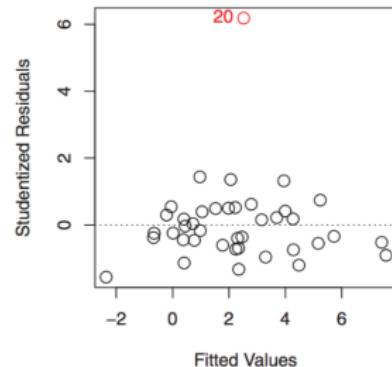
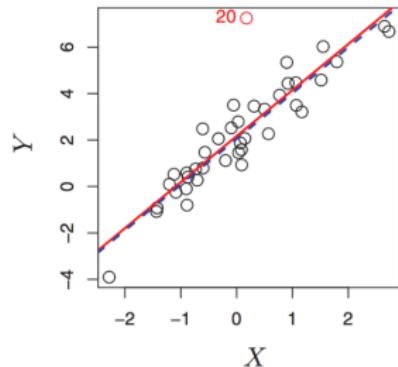
- Colinearity occurs in a dataset when two or more features (columns of X) are highly correlated. These features provide redundant information.
- *Perfect* colinearity violates the full-rank assumption, and the feature matrix becomes singular or degenerate.
- Signs of colinearity include:
 - Opposing signs for the coefficients of the effected variables, when it's expected that both would have the same sign.
 - Standard errors of the regression coefficients of the effected variables tend to be large.
 - Large changes to the regression coefficients when a feature is added or deleted (unstable solution).
 - Rule of thumb: a variance inflation factor (VIF) > 5 indicates a multicollinearity problem.
 - $VIF = 1 / (1 - R_j^2)$
 - R_j^2 is the coefficient of determination of a regression of feature j on all the other features.
- Remedies to colinearity include:
 - Regularization (Ridge and Lasso -- tomorrow)
 - Principal component analysis (PCA – next week)
 - Engineering a feature that combines the affected features (ahem first case study)
 - Simply dropping one of the features (lazy, but viable, option)

What makes these data points “unusual”?



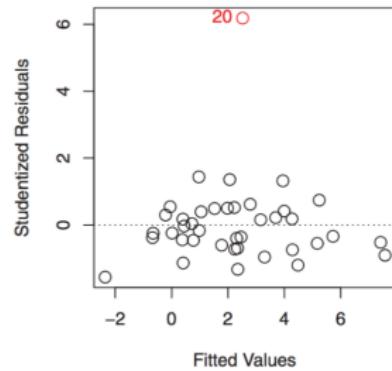
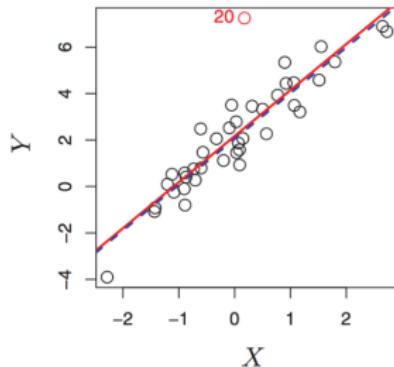
Regression Diagnostics

Outliers

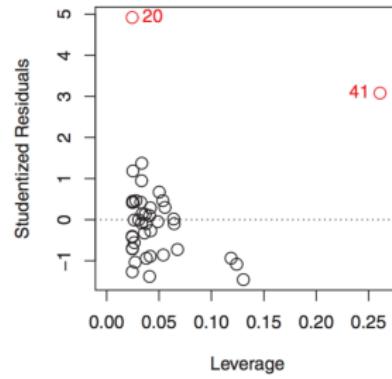
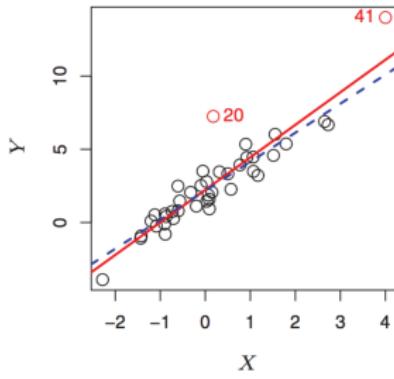


Regression Diagnostics

Outliers

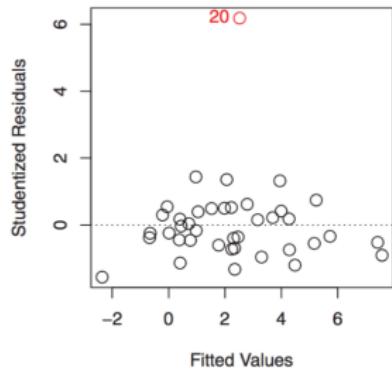
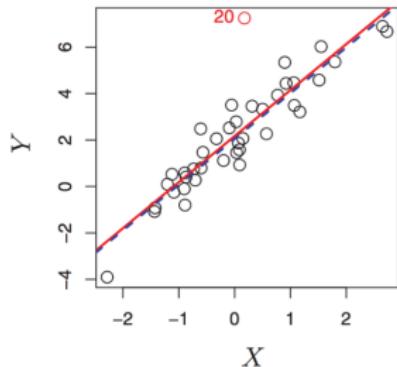


High Leverage Points

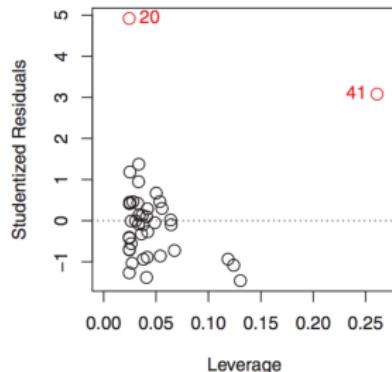
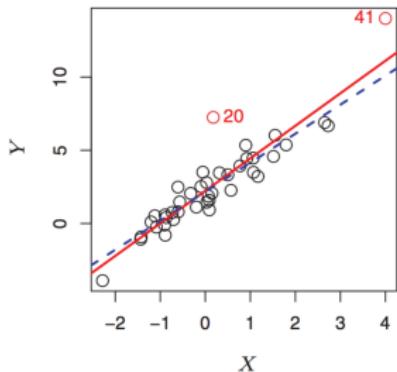


Regression Diagnostics

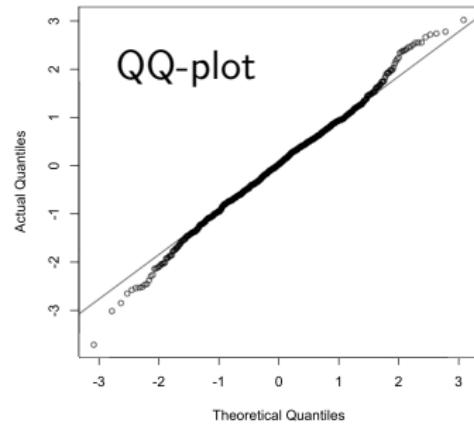
Outliers impact residual variance estimates



High Leverage Points impact prediction estimates

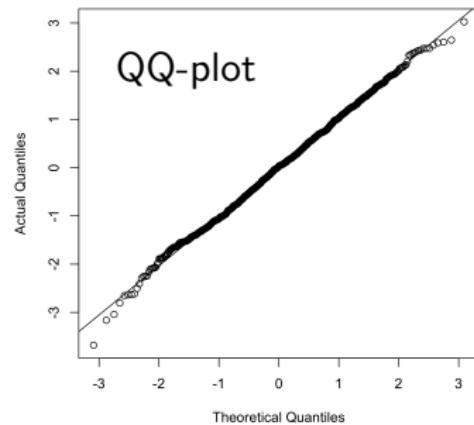


Regression Diagnostics (with residuals)



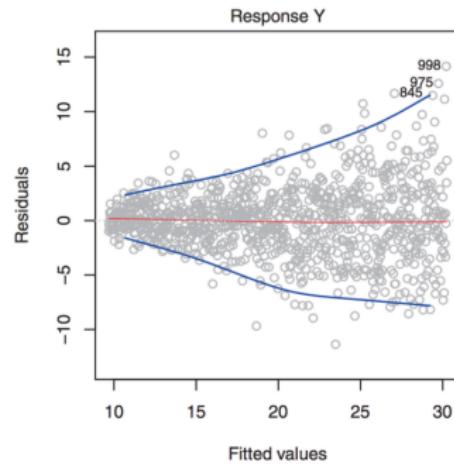
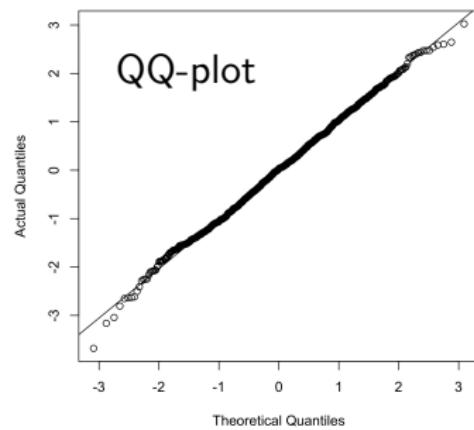
- ▶ What's wrong with the residual distribution?

Regression Diagnostics (with residuals)



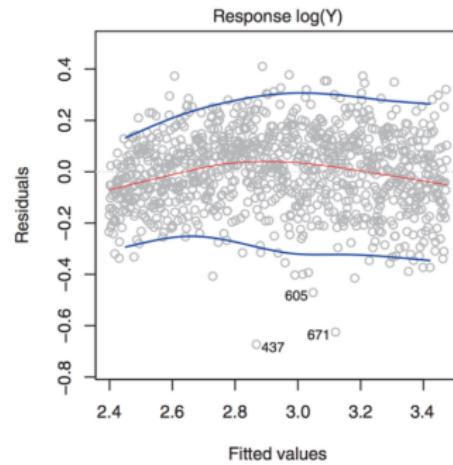
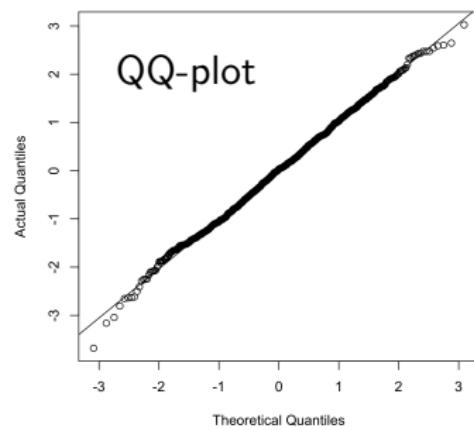
- ▶ What's wrong with the residual distribution?

Regression Diagnostics (with residuals)



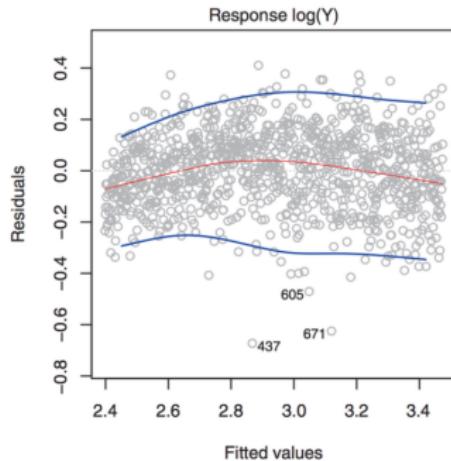
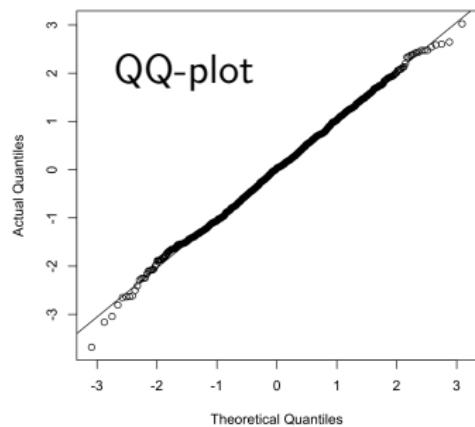
- ▶ What's wrong with the residual distribution?
- ▶ What's wrong with the residual variance?

Regression Diagnostics (with residuals)

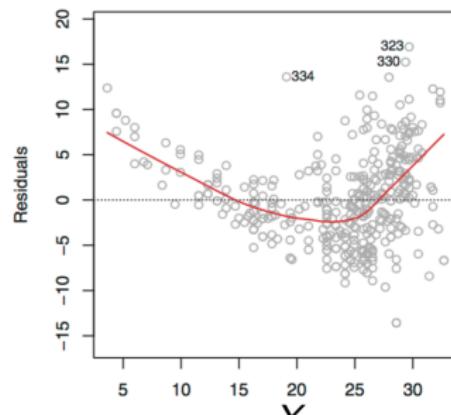


- ▶ What's wrong with the residual distribution?
- ▶ What's wrong with the residual variance?

Regression Diagnostics (with residuals)



- ▶ What's wrong with the residual distribution?
- ▶ What's wrong with the residual variance?
- ▶ What's wrong with this feature/outcome relationship?



Leverage

The *hat* matrix H “*puts the hat on*” \mathbf{Y} projecting \mathbf{Y} onto the (least squares) closest vector to \mathbf{Y} in the column space of \mathbf{x} , $\hat{\mathbf{Y}} \in \mathcal{R}(\mathbf{x})$

$$H = \mathbf{x}(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T$$

$$\begin{aligned}\hat{\mathbf{Y}} &= \mathbf{x}(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{Y} \\ &= H\mathbf{Y}\end{aligned}$$

Leverage

The *hat* matrix H “*puts the hat on*” \mathbf{Y} projecting \mathbf{Y} onto the (least squares) closest vector to \mathbf{Y} in the column space of \mathbf{x} , $\hat{\mathbf{Y}} \in \mathcal{R}(\mathbf{x})$

$$H = \mathbf{x}(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T$$

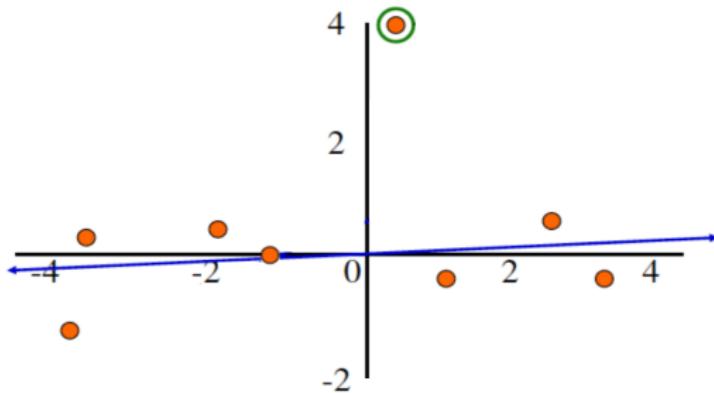
$$\begin{aligned}\hat{\mathbf{Y}} &= \mathbf{x}(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{Y} \\ &= H\mathbf{Y}\end{aligned}$$

- ▶ Diagonal element $H_{ii} \in [0, 1]$, and $\sum_{i=1}^n H_{ii} = \text{rank}(\mathbf{x})$
 H_{ii} is called the *leverage* of observation i

Leverage

The *hat matrix* H “*puts the hat on*” \mathbf{Y} projecting \mathbf{Y} onto the (least squares) closest vector to \mathbf{Y} in the column space of \mathbf{x} , $\hat{\mathbf{Y}} \in \mathcal{R}(\mathbf{x})$

$$\begin{aligned} H &= \mathbf{x}(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \\ \hat{\mathbf{Y}} &= \mathbf{x}(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{Y} \\ &= H\mathbf{Y} \end{aligned}$$

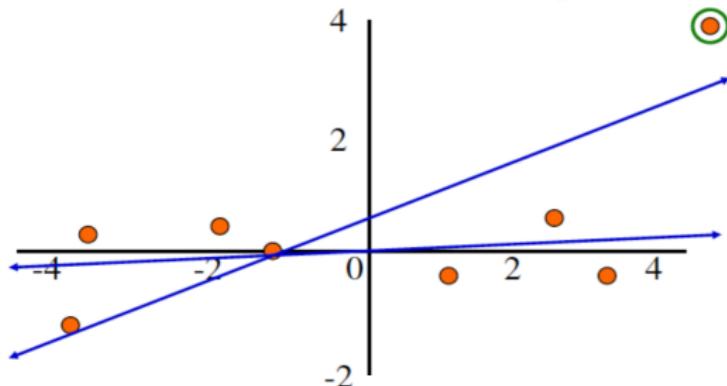


- ▶ Diagonal element $H_{ii} \in [0, 1]$, and $\sum_{i=1}^n H_{ii} = \text{rank}(\mathbf{x})$
 H_{ii} is called the *leverage* of observation i
- ▶ H_{ii} shows much \hat{Y}_i depends on Y_i
which depends on the “extremeness” of x_i

Leverage

The *hat matrix* H “*puts the hat on*” \mathbf{Y} projecting \mathbf{Y} onto the (least squares) closest vector to \mathbf{Y} in the column space of \mathbf{x} , $\hat{\mathbf{Y}} \in \mathcal{R}(\mathbf{x})$

$$\begin{aligned} H &= \mathbf{x}(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \\ \hat{\mathbf{Y}} &= \mathbf{x}(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{Y} \\ &= H\mathbf{Y} \end{aligned}$$

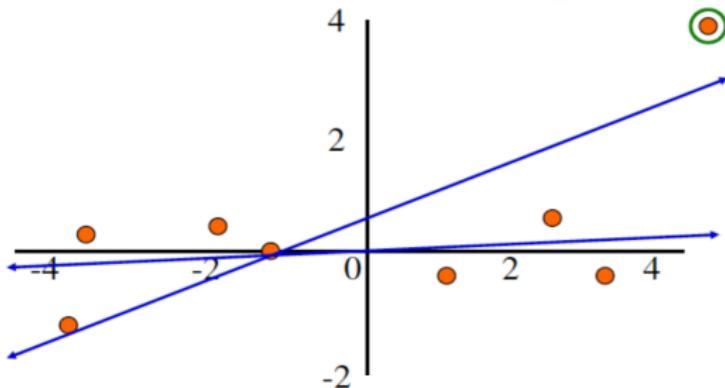


- ▶ Diagonal element $H_{ii} \in [0, 1]$, and $\sum_{i=1}^n H_{ii} = \text{rank}(\mathbf{x})$
 H_{ii} is called the *leverage* of observation i
- ▶ H_{ii} shows much \hat{Y}_i depends on Y_i
which depends on the “extremeness” of x_i

Leverage

The *hat matrix* H “*puts the hat on*” \mathbf{Y} projecting \mathbf{Y} onto the (least squares) closest vector to \mathbf{Y} in the column space of \mathbf{x} , $\hat{\mathbf{Y}} \in \mathcal{R}(\mathbf{x})$

$$\begin{aligned} H &= \mathbf{x}(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \\ \hat{\mathbf{Y}} &= \mathbf{x}(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{Y} \\ &= H\mathbf{Y} \end{aligned}$$

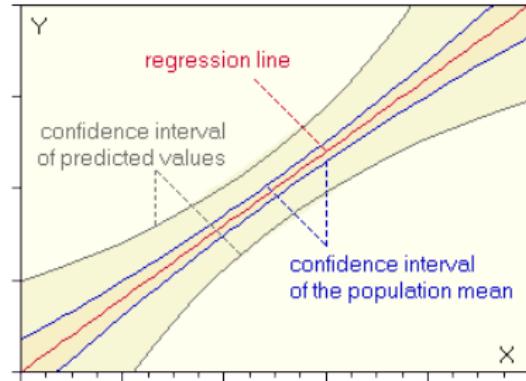


- ▶ Diagonal element $H_{ii} \in [0, 1]$, and $\sum_{i=1}^n H_{ii} = \text{rank}(\mathbf{x})$
 H_{ii} is called the *leverage* of observation i
- ▶ H_{ii} shows much \hat{Y}_i depends on Y_i
which depends on the “extremeness” of x_i
- ▶ Relative comparison of H_{ii} 's id.'s “high leverage observations”

Influential Data Points

Studentized Residuals
have a t-distribution...

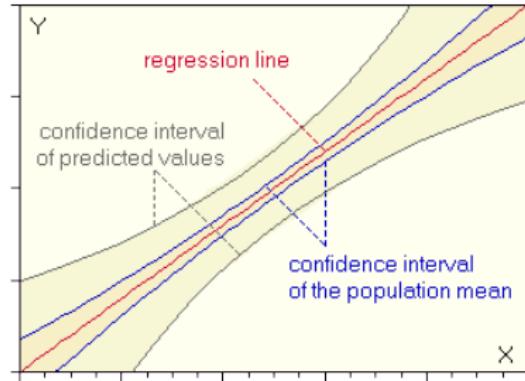
$$r_i = \frac{\hat{\epsilon}_i}{\hat{\sigma} \sqrt{1 - h_{ii}}}$$



Influential Data Points

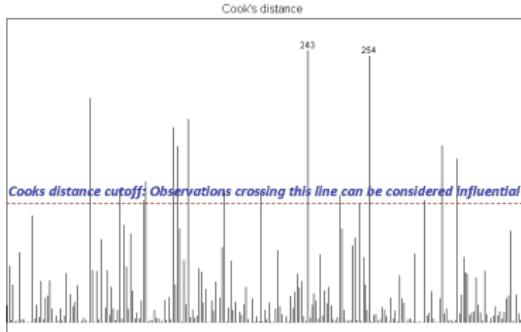
Studentized Residuals
have a t-distribution...

$$r_i = \frac{\hat{\epsilon}_i}{\hat{\sigma} \sqrt{1 - h_{ii}}}$$



Cook's Distance is

$$\begin{aligned} D_i &= \frac{\sum_{j=1}^n (\hat{y}_j - \hat{y}_{j(-i)})^2}{\hat{\sigma}^2 p} \\ &= \frac{\hat{\epsilon}_i}{\hat{\sigma}^2 p} \frac{h_{ii}}{(1 - h_{ii})^2} \end{aligned}$$



Influential data point i may have $D_i > \{3 \times \bar{D}, 1, 4/n, F_{p,n-p}^{1-\alpha}\}$

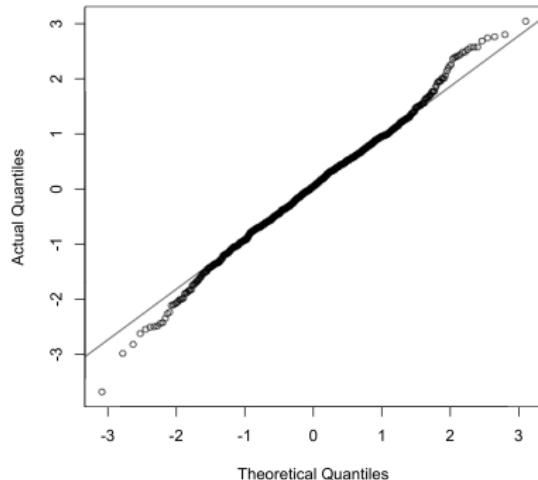
Assumptions, violations, and remedial measures

$$Y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \epsilon_i, \quad \epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

Assumptions, violations, and remedial measures

$$Y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \epsilon_i, \quad \epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

- ▶ Normality



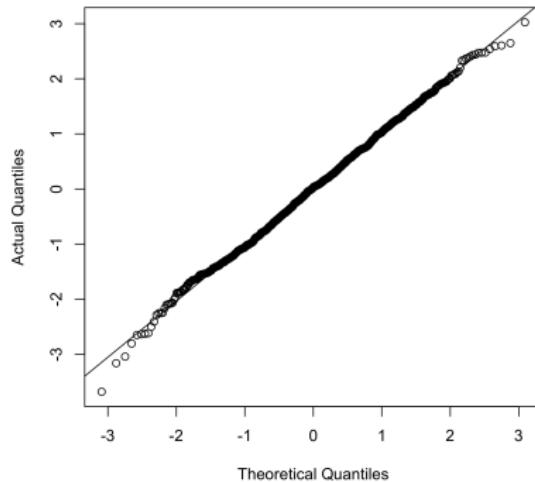
Q-Q Plot

Hypothesis testing depends on
distributional assumptions

Assumptions, violations, and remedial measures

$$Y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \epsilon_i, \quad \epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

- ▶ Normality



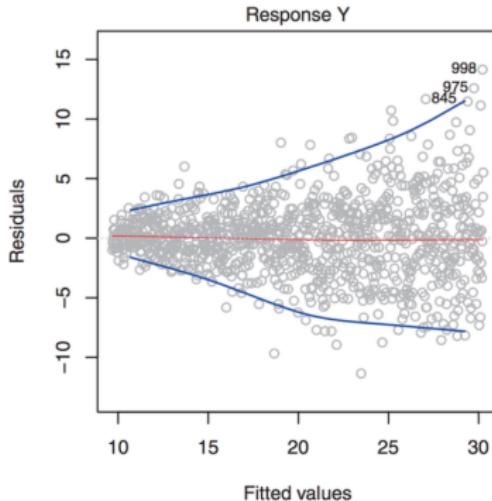
Q-Q Plot

Hypothesis testing depends on
distributional assumptions

Assumptions, violations, and remedial measures

$$Y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \epsilon_i, \quad \epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

- ▶ Normality
- ▶ Homoskedasticity



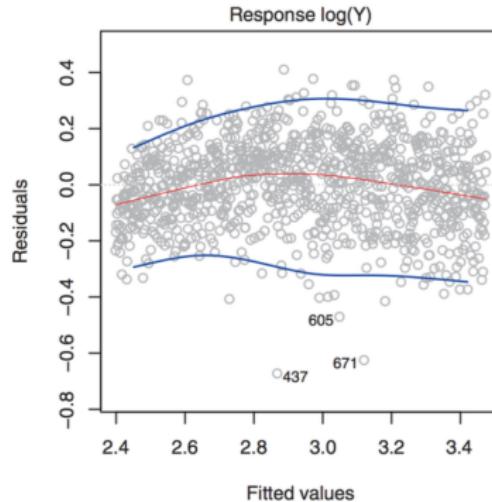
Residuals versus Fitted Values

Box-Cox transformations $\frac{Y^\lambda - 1}{\lambda}$ can help

Assumptions, violations, and remedial measures

$$Y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \epsilon_i, \quad \epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

- ▶ Normality
- ▶ Homoskedasticity



Residuals versus Fitted Values

Box-Cox transformations $\frac{Y^\lambda - 1}{\lambda}$ can help

Assumptions, violations, and remedial measures

$$Y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \epsilon_i, \quad \epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

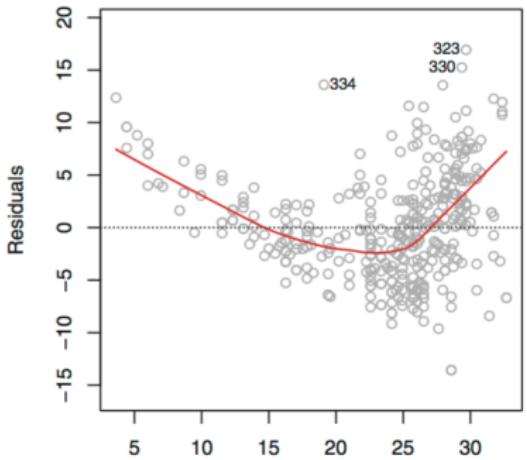
- ▶ Normality
- ▶ Homoskedasticity
- ▶ Independence

$$\text{Cov}[\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}] \approx \begin{bmatrix} \sigma^2 & 0 & \cdots & 0 \\ 0 & \sigma^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma^2 \end{bmatrix}$$

Assumptions, violations, and remedial measures

$$Y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \epsilon_i, \quad \epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

- ▶ Normality
- ▶ Homoskedasticity
- ▶ Independence
- ▶ Linear form



Residuals versus Feature Values

“All models are wrong, some are useful”
– George Box

Assumptions, violations, and remedial measures

$$Y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \epsilon_i, \quad \epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

- ▶ Normality
- ▶ Homoskedasticity
- ▶ Independence
- ▶ Linear form
- ▶ Fixed x 's

$$\mathbf{Y} \sim \text{MVN}(\mathbf{x}\boldsymbol{\beta}, \sigma^2 I)$$

Quiz: assumptions?

$$Y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \epsilon_i, \quad \epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

- 1.
- 2.
- 3.
- 4.
- 5.

Half time

Assessing Model Fit (more Machine Learning-ish)

Residual Sum of Squares

$$RSS = \sum(Y_i - \hat{Y}_i)^2 = \sum \hat{\epsilon}_i^2$$

Total Sum of Squares

$$\begin{aligned} TSS &= \sum(Y_i - \bar{Y})^2 \\ &= \sum(\hat{Y}_i - \bar{Y})^2 + RSS \end{aligned}$$

Residual Standard Deviation

$$\begin{aligned} \hat{\sigma} &= \sqrt{\frac{1}{n-p-1} RSS} \\ &= \sqrt{\frac{\sum(Y_i - \hat{Y}_i)^2}{n-p-1}} \end{aligned}$$

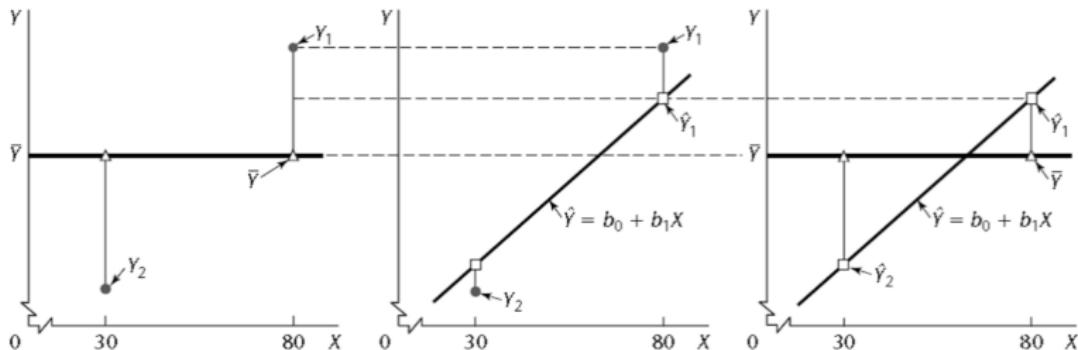
Proportion of Variance Explained

$$\begin{aligned} R^2 &= \frac{TSS - RSS}{TSS} \\ &= 1 - \frac{RSS}{TSS} \end{aligned}$$

F-statistic

$$F = \frac{(TSS - RSS)/p}{RSS/(n-p-1)} = \frac{\sum(\hat{Y}_i - \bar{Y})^2/p}{\sum(Y_i - \hat{Y})^2/(n-p-1)}$$

Decomposition of Total Variation



$$\begin{aligned} TSS &= \sum(Y_i - \bar{Y})^2 = \sum(Y_i - \hat{Y}_i + \hat{Y}_i - \bar{Y})^2 \\ &= \sum(Y_i - \hat{Y}_i)^2 + 2\sum(Y_i - \hat{Y}_i)(\hat{Y}_i - \bar{Y}) + \sum(\hat{Y}_i - \bar{Y})^2 \\ &= \sum(Y_i - \hat{Y}_i)^2 + 2\sum\hat{\epsilon}_i(\hat{Y}_i - \bar{Y}) + \sum(\hat{Y}_i - \bar{Y})^2 \\ &\quad \sum\hat{\epsilon}_i = 0 \uparrow \uparrow \sum\hat{\epsilon}_i \hat{Y}_i = 0 \\ &= \sum(Y_i - \hat{Y}_i)^2 + \sum(\hat{Y}_i - \bar{Y})^2 = RSS + \sum(\hat{Y}_i - \bar{Y})^2 \end{aligned}$$

MODEL SELECTION!!!

- ▶ R^2 (model fit) is insufficient – more features means larger R^2

MODEL SELECTION!!!

- ▶ R^2 (model fit) is insufficient – more features means larger R^2
- ▶ Spuriously improving model fit to data is called *overfitting*

MODEL SELECTION!!!

- ▶ R^2 (model fit) is insufficient – more features means larger R^2
 - ▶ Spuriously improving model fit to data is called *overfitting*
 - ▶ We want model fits to generalize to *population* phenomenon
-

MODEL SELECTION!!!

- ▶ R^2 (model fit) is insufficient – more features means larger R^2
 - ▶ Spuriously improving model fit to data is called *overfitting*
 - ▶ We want model fits to generalize to *population* phenomenon
-
- ▶ Classical Statistics Approaches:
Model Selection Criterion (choose smallest)

$$\text{Mallow's } C_p = \frac{1}{n}(RSS + 2p\hat{\sigma}^2)$$

$$AIC = -2 \log L + 2p$$

$$BIC = -2 \log L + p \log n$$

$$\text{Adjusted } R^2 = 1 - \frac{RSS/(n-p-1)}{TSS/(n-1)}$$

$$D_M = -2 \log f(Y|\hat{\theta}^{M_p}) + 2 \log f(Y|Y)$$

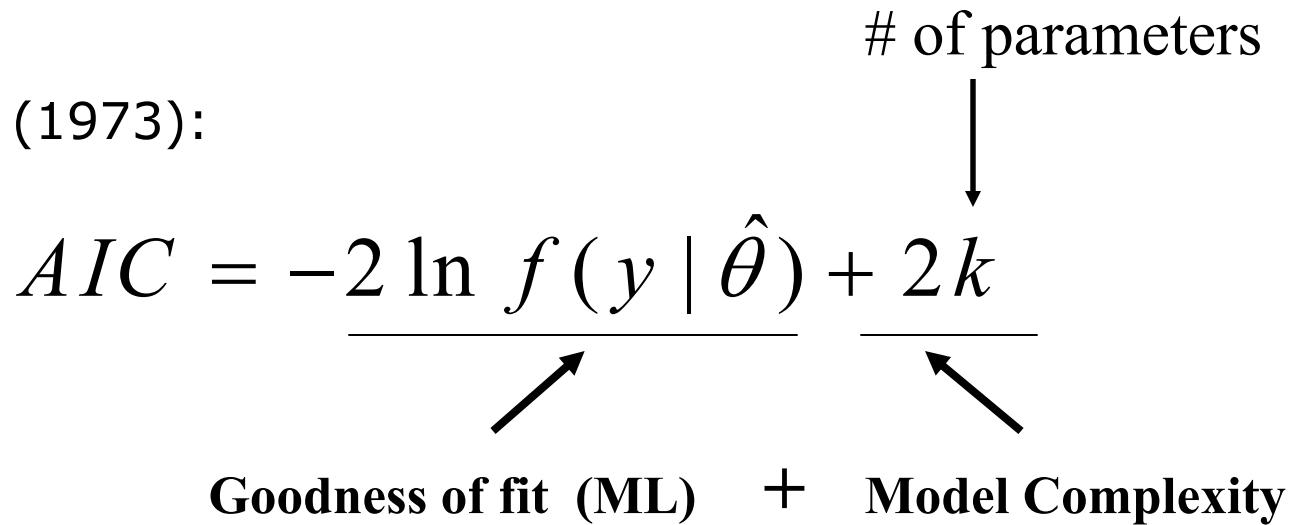
$$D_M \stackrel{\text{approx.}}{\sim} \chi^2_{n-p-1}$$

Akaike Information Criterion (AIC) as a Method of Model Selection

Akaike (1973):

$$AIC = \frac{-2 \ln f(y | \hat{\theta})}{\text{Goodness of fit (ML)}} + \frac{2k}{\text{Model Complexity}}$$

of parameters
↓
 $\hat{\theta}$
k



The model that minimizes AIC should be preferred

Bayesian Information Criterion (BIC)

Schwarz (1978):

$$BIC = \underbrace{-2 \ln f(y | \hat{\theta})}_{\text{Goodness of fit (ML)}} + \underbrace{k \ln n}_{\text{Model Complexity}}$$

Assessing Parameter Uncertainty (definitely Statistics)

For $\hat{\beta} = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{Y}$, since (under H_0)

$$f(\hat{\beta} | \beta, \sigma^2, \mathbf{x}) = MVN\left(\beta, \sigma^2(\mathbf{x}^T \mathbf{x})^{-1}\right)$$

we have that

Assessing Parameter Uncertainty (definitely Statistics)

For $\hat{\beta} = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{Y}$, since (under H_0)

$$f(\hat{\beta} | \beta, \sigma^2, \mathbf{x}) = MVN\left(\beta, \sigma^2(\mathbf{x}^T \mathbf{x})^{-1}\right)$$

we have that

$$\frac{\hat{\beta}_i - \beta_i}{\text{SD}(\hat{\beta}_i)} \sim N(0, 1)$$

Assessing Parameter Uncertainty (definitely Statistics)

For $\hat{\beta} = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{Y}$, since (under H_0)

$$f(\hat{\beta} | \beta, \sigma^2, \mathbf{x}) = MVN\left(\beta, \sigma^2(\mathbf{x}^T \mathbf{x})^{-1}\right)$$

we have that

$$\frac{\hat{\beta}_i - \beta_i}{\text{SD}(\hat{\beta}_i)} \sim N(0, 1)$$

and if we estimate

$$\hat{\sigma}^2 = \frac{(\mathbf{Y} - \mathbf{x}\hat{\beta})^T (\mathbf{Y} - \mathbf{x}\hat{\beta})}{n - p - 1} = \frac{\sum_{i=1}^n \hat{\epsilon}_i^2}{n - p - 1}$$

(where p is the number of coefficients) then we have that

Assessing Parameter Uncertainty (definitely Statistics)

For $\hat{\beta} = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{Y}$, since (under H_0)

$$f(\hat{\beta} | \beta, \sigma^2, \mathbf{x}) = MVN\left(\beta, \sigma^2(\mathbf{x}^T \mathbf{x})^{-1}\right)$$

we have that

$$\frac{\hat{\beta}_i - \beta_i}{\text{SD}(\hat{\beta}_i)} \sim N(0, 1)$$

and if we estimate

$$\hat{\sigma}^2 = \frac{(\mathbf{Y} - \mathbf{x}\hat{\beta})^T (\mathbf{Y} - \mathbf{x}\hat{\beta})}{n - p - 1} = \frac{\sum_{i=1}^n \hat{\epsilon}_i^2}{n - p - 1}$$

(where p is the number of coefficients) then we have that

$$\frac{\hat{\beta}_i - \beta_i}{\widehat{\text{SD}}(\hat{\beta}_i)} \sim t_{n-p-1}$$

And this works for any number of feature variables...

Hypothesis Testing for Feature Selection

$$F = \frac{(TSS - RSS)/p}{RSS/(n - p - 1)} = \frac{\sum(\hat{Y}_i - \bar{Y})^2/p}{\sum(Y_i - \hat{Y})/(n - p - 1)}$$

$F \sim F_{p,n-p-1}$ (tests if any coefficient is *non-zero*)

$\frac{\hat{\beta}_i - \beta_i}{\widehat{SD}(\hat{\beta}_i)} \sim t_{n-p-1}$ (tests if a *specific* coefficient is non-zero*)

*in the presence of all the others (this is a “last-in” test)

Hypothesis Testing for Feature Selection

$$F = \frac{(TSS - RSS)/p}{RSS/(n - p - 1)} = \frac{\sum(\hat{Y}_i - \bar{Y})^2/p}{\sum(Y_i - \hat{Y})/(n - p - 1)}$$

$F \sim F_{p,n-p-1}$ (tests if any coefficient is *non-zero*)

$\frac{\hat{\beta}_i - \beta_i}{\widehat{SD}(\hat{\beta}_i)} \sim t_{n-p-1}$ (tests if a *specific* coefficient is non-zero*)

*in the presence of all the others (this is a “last-in” test)

OLS Regression Results

Dep. Variable:	y	R-squared:	0.933
Model:	OLS	Adj. R-squared:	0.928
Method:	Least Squares	F-statistic:	211.8
Date:	Mon, 03 Nov 2014	Prob (F-statistic):	6.30e-27 ←
Time:	14:45:06	Log-Likelihood:	-34.438
No. Observations:	50	AIC:	76.88
Df Residuals:	46	BIC:	84.52
Df Model:	3		
Covariance Type:	nonrobust		
coef	std err	t	P> t
x1	0.4687	0.026	17.751
x2	0.4836	0.104	4.659
x3	-0.0174	0.002	-7.507
const	5.2058	0.171	30.405

Hypothesis Testing for Feature Selection

$$F = \frac{(TSS - RSS)/p}{RSS/(n - p - 1)} = \frac{\sum(\hat{Y}_i - \bar{Y})^2/p}{\sum(Y_i - \hat{Y})/(n - p - 1)}$$

$F \sim F_{p,n-p-1}$ (tests if any coefficient is *non-zero*)

$\frac{\hat{\beta}_i - \beta_i}{\widehat{SD}(\hat{\beta}_i)} \sim t_{n-p-1}$ (tests if a *specific* coefficient is non-zero*)

*in the presence of all the others (this is a “last-in” test)

- ▶ Forward Selection

OLS Regression Results					
Dep. Variable:	y	R-squared:	0.933		
Model:	OLS	Adj. R-squared:	0.928		
Method:	Least Squares	F-statistic:	211.8		
Date:	Mon, 03 Nov 2014	Prob (F-statistic):	6.30e-27 ←		
Time:	14:45:06	Log-Likelihood:	-34.438		
No. Observations:	50	AIC:	76.88		
Df Residuals:	46	BIC:	84.52		
Df Model:	3				
Covariance Type:	nonrobust				
	coef	std err	t	P> t	[95.0% Conf. Int.]
x1	0.4687	0.026	17.751	0.000	0.416 0.522
x2	0.4836	0.104	4.659	0.000	0.275 0.693
x3	-0.0174	0.002	-7.507	0.000	-0.022 -0.013
const	5.2058	0.171	30.405	0.000	4.861 5.550

Hypothesis Testing for Feature Selection

$$F = \frac{(TSS - RSS)/p}{RSS/(n - p - 1)} = \frac{\sum(\hat{Y}_i - \bar{Y})^2/p}{\sum(Y_i - \hat{Y})/(n - p - 1)}$$

$F \sim F_{p,n-p-1}$ (tests if any coefficient is *non-zero*)

$\frac{\hat{\beta}_i - \beta_i}{\widehat{SD}(\hat{\beta}_i)} \sim t_{n-p-1}$ (tests if a *specific* coefficient is non-zero*)

*in the presence of all the others (this is a “last-in” test)

- ▶ Forward Selection

OLS Regression Results						
Dep. Variable:	y	R-squared:	0.933	Model:	OLS	Adj. R-squared:
Method:	Least Squares	F-statistic:	211.8	Date:	Mon, 03 Nov 2014	Prob (F-statistic):
Time:	14:45:06	Log-Likelihood:	-34.438	No. Observations:	50	AIC:
Df Residuals:	46	BIC:	76.88	Df Model:	3	BIC:
Covariance Type:	nonrobust		84.52			
=====						
	coef	std err	t	P> t	[95.0% Conf. Int.]	
x1	0.4687	0.026	17.751	0.000	0.416	0.522
x2	0.4836	0.104	4.659	0.000	0.275	0.693
x3	-0.0174	0.002	-7.507	0.000	-0.022	-0.013
const	5.2058	0.171	30.405	0.000	4.861	5.550
=====						

- ▶ Backward Selection

Hypothesis Testing for Feature Selection

$$F = \frac{(TSS - RSS)/p}{RSS/(n - p - 1)} = \frac{\sum(\hat{Y}_i - \bar{Y})^2/p}{\sum(Y_i - \hat{Y})/(n - p - 1)}$$

$F \sim F_{p,n-p-1}$ (tests if any coefficient is *non-zero*)

$\frac{\hat{\beta}_i - \beta_i}{\widehat{SD}(\hat{\beta}_i)} \sim t_{n-p-1}$ (tests if a *specific* coefficient is non-zero*)

*in the presence of all the others (this is a “last-in” test)

- ▶ Forward Selection

OLS Regression Results							
Dep. Variable:	y	R-squared:	0.933	Model:	OLS	Adj. R-squared:	
Method:	Least Squares	F-statistic:	211.8	Date:	Mon, 03 Nov 2014	Prob (F-statistic):	
Time:	14:45:06	Log-Likelihood:	-34.438	No. Observations:	50	AIC:	76.88
Df Residuals:	46	BIC:	84.52	Df Model:	3	Covariance Type:	nonrobust

	coef	std err	t	P> t	[95.0% Conf. Int.]	
x1	0.4687	0.026	17.751	0.000	0.416	0.522
x2	0.4836	0.104	4.659	0.000	0.275	0.693
x3	-0.0174	0.002	-7.507	0.000	-0.022	-0.013
const	5.2058	0.171	30.405	0.000	4.861	5.550

- ▶ Backward Selection

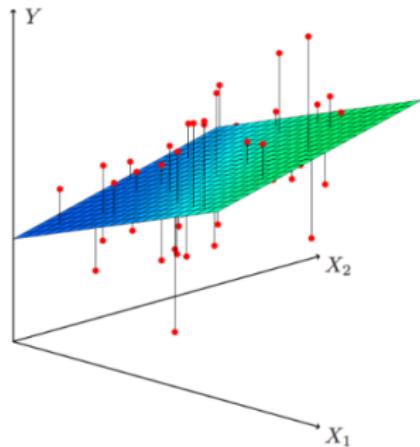
- ▶ Both

Linear Models

- ▶ Linear model... that sounds too simple...

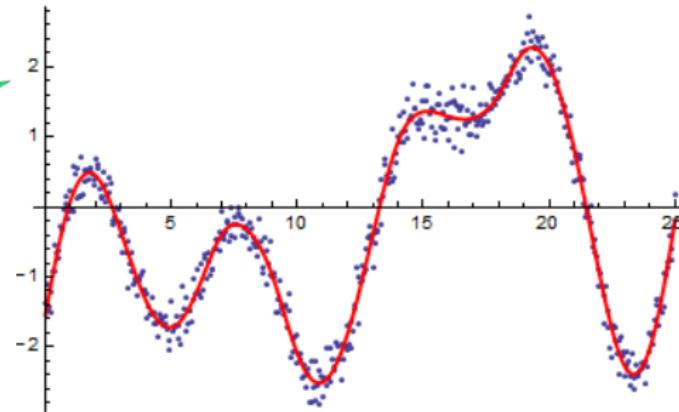
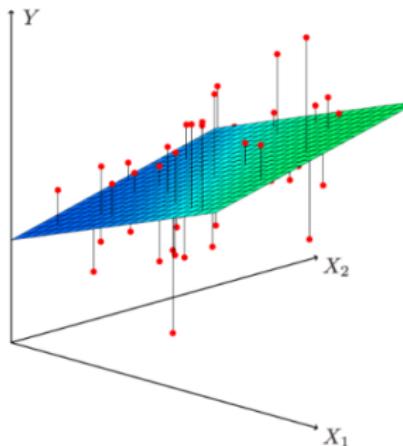
Linear Models

- ▶ Linear model... that sounds too simple...



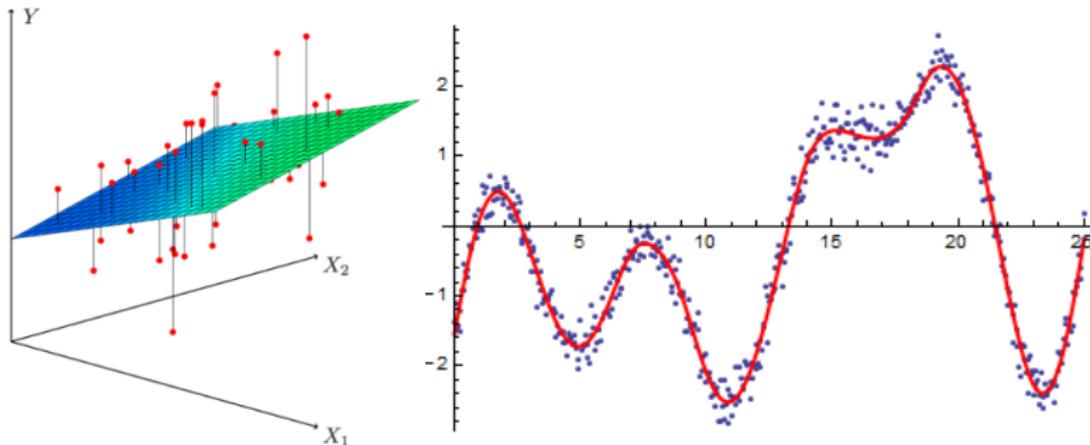
Linear Models

- ▶ Linear model... that sounds too simple...



Linear Models

- ▶ Linear model... that sounds too simple...



- ▶ “Linear” models are only linear in the coefficients

$$Y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \epsilon_i$$

- ▶ The x 's can be pretty wild...

Features that produce “non-linear” response surfaces?

Features that produce “non-linear” response surfaces?

- ▶ Higher order terms: $X_1^{\frac{1}{2}}, X_1^2, X_1^3$

Features that produce “non-linear” response surfaces?

- ▶ Higher order terms: $X_1^{\frac{1}{2}}, X_1^2, X_1^3$
- ▶ Qualitative variables: $X_1 \in \{0, 1\}, X_1 \in \{a, b, c, d\}$

Features that produce “non-linear” response surfaces?

- ▶ Higher order terms: $X_1^{\frac{1}{2}}, X_1^2, X_1^3$
- ▶ Qualitative variables: $X_1 \in \{0, 1\}, X_1 \in \{a, b, c, d\}$

$$\begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} \neq \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

Features that produce “non-linear” response surfaces?

- ▶ Higher order terms: $X_1^{\frac{1}{2}}, X_1^2, X_1^3$
- ▶ Qualitative variables: $X_1 \in \{0, 1\}, X_1 \in \{a, b, c, d\}$

$$\begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} \neq \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

- ▶ Interactions: $X_1 \cdot X_2 + X_1 + X_2$ (*interpretation?*)

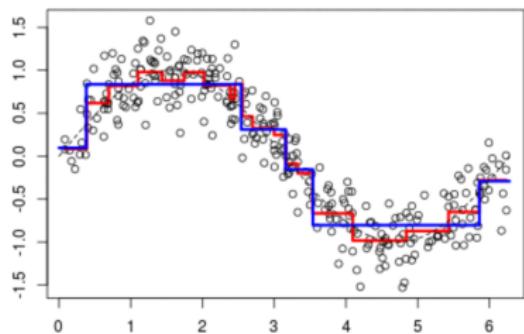
Features that produce “non-linear” response surfaces?

- ▶ Higher order terms: $X_1^{\frac{1}{2}}, X_1^2, X_1^3$
- ▶ Qualitative variables: $X_1 \in \{0, 1\}, X_1 \in \{a, b, c, d\}$

$$\begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} \neq \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

- ▶ Interactions: $X_1 \cdot X_2 + X_1 + X_2$ (*interpretation?*)
- ▶ Step functions

$$Y_i = \beta_j : \text{if } a_j \leq X_i < b_j$$

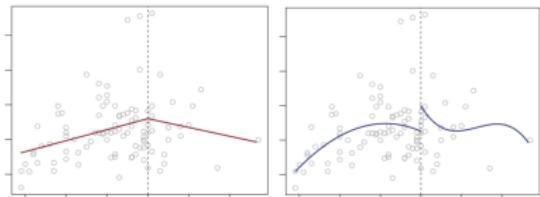


Features that produce “non-linear” response surfaces?

- ▶ Higher order terms: $X_1^{\frac{1}{2}}, X_1^2, X_1^3$
- ▶ Qualitative variables: $X_1 \in \{0, 1\}, X_1 \in \{a, b, c, d\}$

$$\begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} \neq \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

- ▶ Interactions: $X_1 \cdot X_2 + X_1 + X_2$ (*interpretation?*)
- ▶ Step functions
- ▶ Regression Splines



$$Y_i = \begin{cases} \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \epsilon_i : \text{if } X_i \leq c \\ \beta_0^* + \beta_1 X_i + \beta_2^* X_i^2 + \epsilon_i : \text{if } X_i > c \end{cases}$$

Features that produce “non-linear” response surfaces?

- ▶ Higher order terms: $X_1^{\frac{1}{2}}, X_1^2, X_1^3$
- ▶ Qualitative variables: $X_1 \in \{0, 1\}, X_1 \in \{a, b, c, d\}$

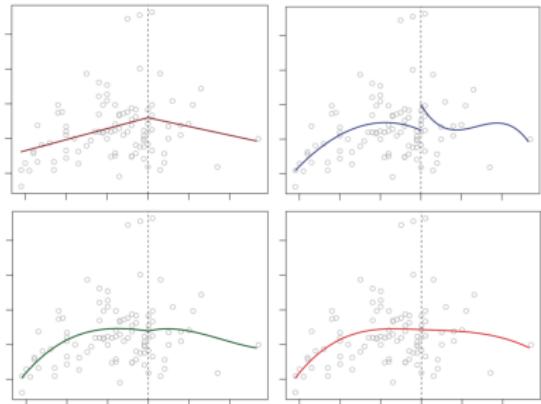
$$\begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} \neq \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

- ▶ Interactions: $X_1 \cdot X_2 + X_1 + X_2$ (*interpretation?*)
- ▶ Step functions
- ▶ Regression Splines

$$h(X_i, \xi) = \begin{cases} (x - \xi)^3 & : \text{if } X_i > \xi \\ 0 & : \text{if } X_i \leq \xi \end{cases}$$

basis functions & knots

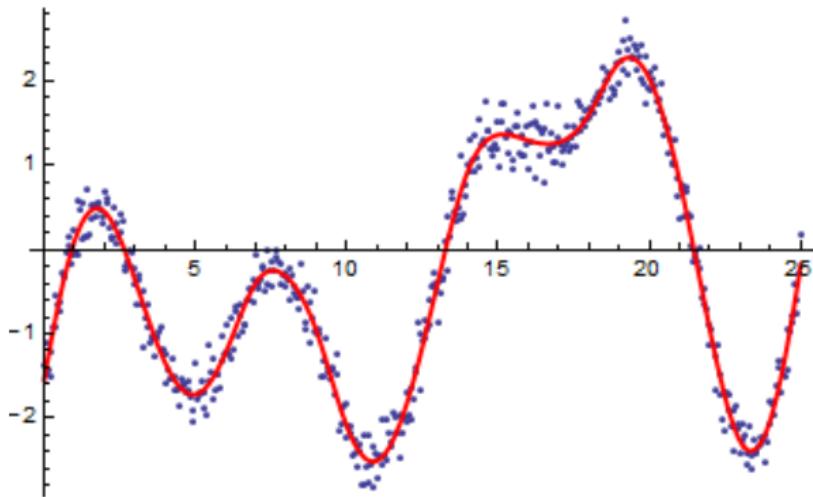
$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \beta_3 X_i^3 + \beta_{s1} h(X_i, \xi_1) + \dots + \epsilon_i$$



Linear models aren't really so “linear”

Other ways to get “non linear regressions”

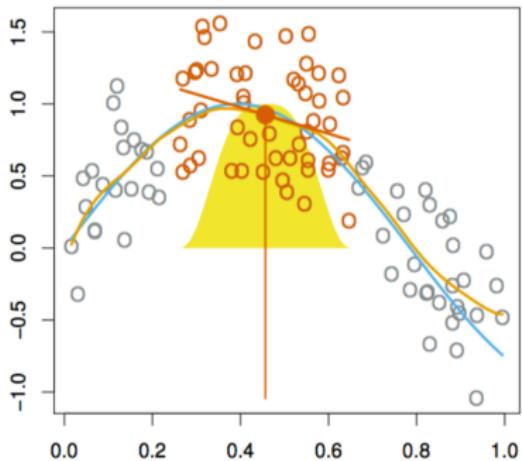
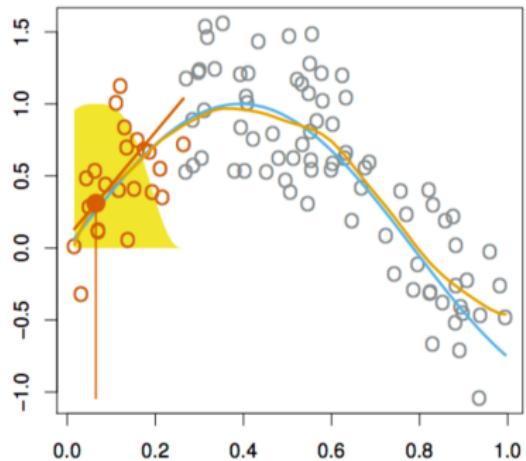
- ▶ Smoothing Splines



$$\min_g \sum_{i=1}^n (Y_i - g(x_i))^2 + \lambda \int g''(t)^2 dt$$

Other ways to get “non linear regressions”

- ▶ Local Regression (LOESS)



Other ways to get “non linear regressions”

- ▶ Generalized Additive Models

$$y_i = \beta_0 + f_1(x_{i1}) + f_2(x_{i2}) + \cdots + f_p(x_{ip}) + \epsilon_i.$$

