# Lecture 7: Linear Regression (continued)

## Reading: Chapter 3

### STATS 202: Data mining and analysis

Jonathan Taylor, 10/8
Slide credits: Sergio Bacallado

# Potential issues in linear regression

1. Interactions between predictors

2. Non-linear relationships

3. Correlation of error terms

4. Non-constant variance of error (heteroskedasticity).

5. Outliers

6. High leverage points

7. Collinearity

# Correlation of error terms

We assumed that the errors for each sample are independent:

$$y_i = f(x_i) + \varepsilon_i \quad ; \quad \varepsilon_i \sim \mathcal{N}(0, \sigma) \text{ i.i.d.}$$

What if this breaks down?

The main effect is that this invalidates any assertions about Standard Errors, confidence intervals, and hypothesis tests:

**Example**: Suppose that by accident, we double the data (we use each sample twice). Then, the standard errors would be artificially smaller by a factor of $\sqrt{2}$.
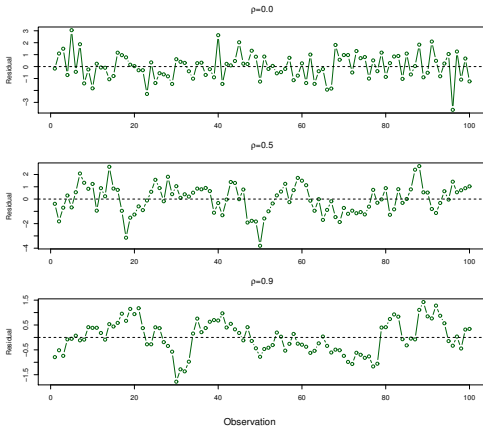
# Correlation of error terms

When could this happen in real life:

- **Time series:** Each sample corresponds to a different point in time. The errors for samples that are close in time are correlated.

- **Spatial data:** Each sample corresponds to a different location in space.

- **Predicting height from weight at birth:** Suppose some of the subjects in the study are in the same family, their shared environment could make them deviate from $f(x)$ in similar ways.
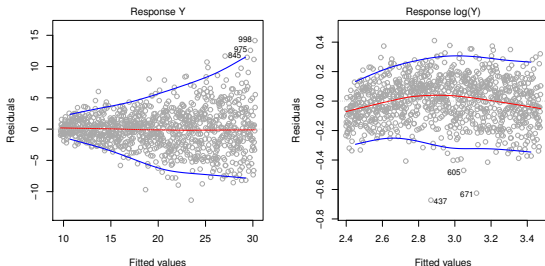
# Correlation of error terms

Simulations of time series with increasing correlations between $\varepsilon_i$.

# Non-constant variance of error (heteroskedasticity)

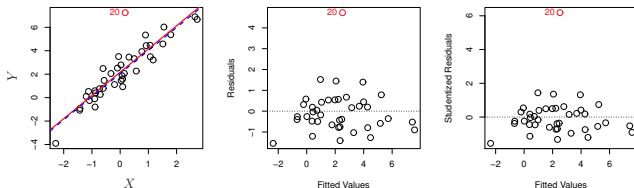The variance of the error depends on the input.

To diagnose this, we can plot residuals vs. fitted values:



**Solution:** If the trend in variance is relatively simple, we can transform the response using a logarithm, for example.
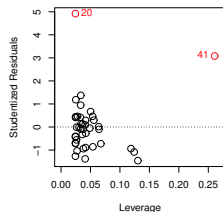
# Outliers
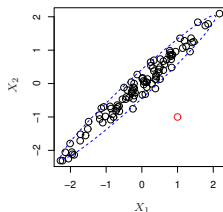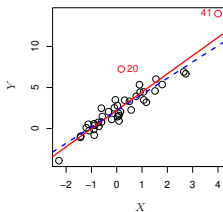
Outliers are points with very high errors.



While they may not affect the fit, they might affect our assessment of model quality.

Possible solutions:

- If we believe an outlier is due to an error in data collection, we can remove it.

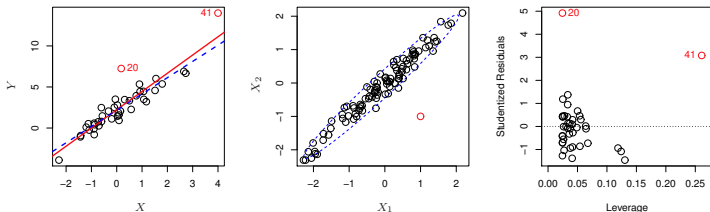- An outlier might be evidence of a missing predictor, or the need to specify a more complex model.

# High leverage points

Some samples with extreme inputs have an outsized effect on $\hat{\beta}$.

# High leverage points

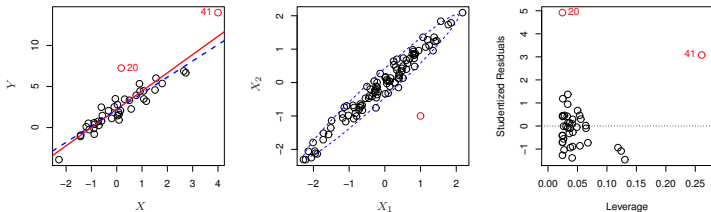Some samples with extreme inputs have an outsized effect on $\hat{\beta}$.



This can be measured with the **leverage statistic** or **self influence**:

$$h_{ii} = \frac{\partial \hat{y}_i}{\partial y_i} = (\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T)_{i,i} \in [1/n, 1].$$

# High leverage points

Some samples with extreme inputs have an outsized effect on $\hat{\beta}$.
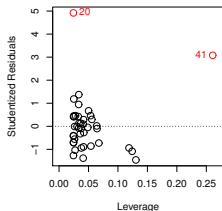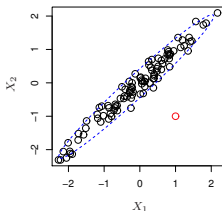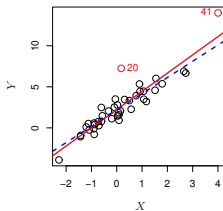


This can be measured with the **leverage statistic** or **self influence**:

$$h_{ii} = \frac{\partial \hat{y}_i}{\partial y_i} = (\underbrace{\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T}_{\text{Hat matrix}})_{i,i} \in [1/n, 1].$$

# Studentized residuals

- The residual $\hat{\epsilon}_i = y_i - \hat{y}_i$ is an estimate for the noise $\epsilon_i$.
- The standard error of $\hat{\epsilon}_i$ is $\sigma\sqrt{1 - h_{ii}}$.
- A **studentized residual** is $\hat{\epsilon}_i$ divided by its standard error.
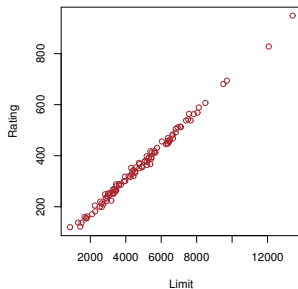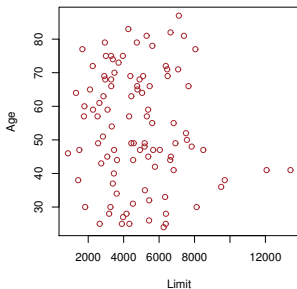- It follows a Student-t distribution with $n - p - 2$ degrees of freedom.

# Collinearity

Two predictors are collinear if one explains the other well:

$$\texttt{limit} = a \times \texttt{rating} + b$$
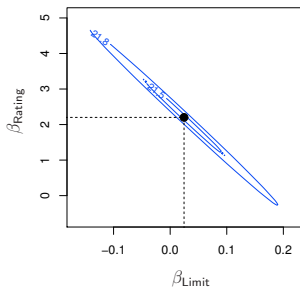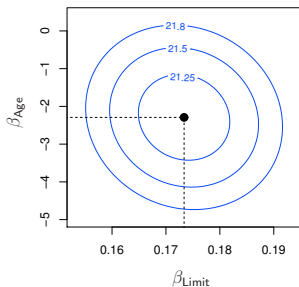
i.e. they contain the same information

# Collinearity

**Problem:** The coefficients become *unidentifiable*. Consider the extreme case of using two identical predictors `limit`:

$$\texttt{balance} = \beta_0 + \beta_1 \times \texttt{limit} + \beta_2 \times \texttt{limit}$$
$$= \beta_0 + (\beta_1 + 100) \times \texttt{limit} + (\beta_2 - 100) \times \texttt{limit}$$

The fit $(\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2)$ is just as good as $(\hat{\beta}_0, \hat{\beta}_1 + 100, \hat{\beta}_2 - 100)$.

# Collinearity

If 2 variables are collinear, we can easily diagnose this using their correlation.

A group of $q$ variables is **multilinear** if these variables "contain less information" than $q$ independent variables. Pairwise correlations may not reveal multilinear variables.

The Variance Inflation Factor (VIF) measures how *necessary* a variable is, or how predictable it is given the other variables:
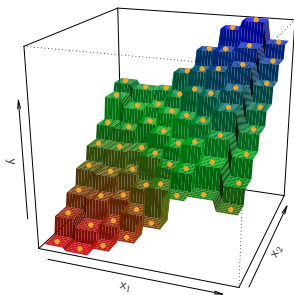
$$VIF(\hat{\beta}_j) = \frac{1}{1 - R^2_{X_j|X_{-j}}},$$

where $R^2_{X_j|X_{-j}}$ is the $R^2$ statistic for Multiple Linear regression of the predictor $X_j$ onto the remaining predictors.

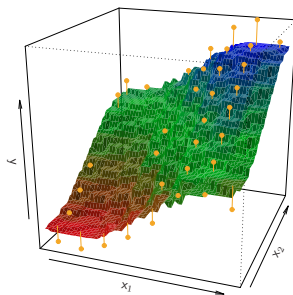# Comparing Linear Regression to $K$-nearest neighbors

**Linear regression:** prototypical parametric method.
**KNN regression:** prototypical nonparametric method.

$$\hat{f}(x) = \frac{1}{K} \sum_{i \in N_K(x)} y_i$$



$K = 1$          $K = 9$

# Comparing Linear Regression to $K$-nearest neighbors

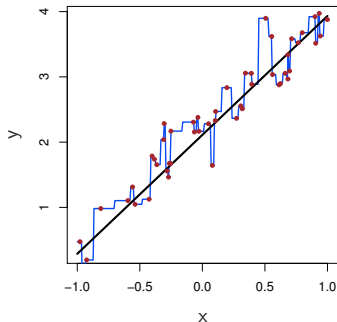**Linear regression:** prototypical parametric method.

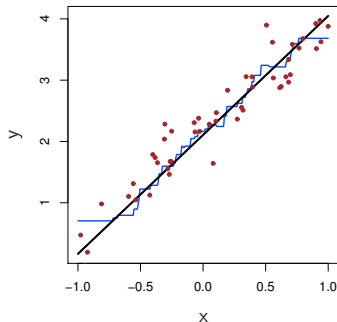**KNN regression:** prototypical nonparametric method.

Long story short:

- KNN is only better when the function $f$ is not linear.
- When $n$ is not much larger than $p$, even if $f$ is nonlinear, Linear Regression can outperform KNN. KNN has smaller bias, but this comes at a price of higher variance.
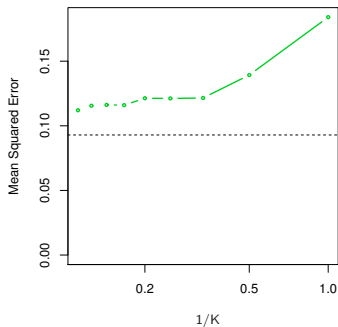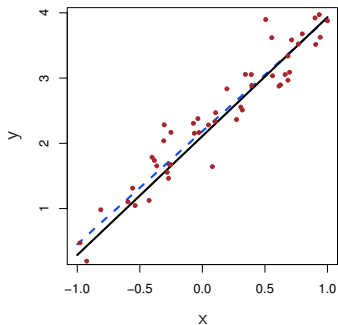
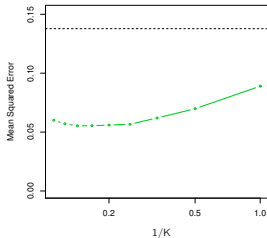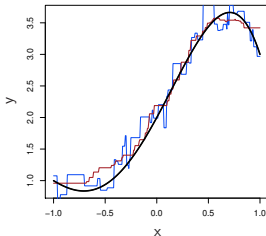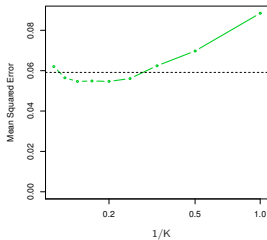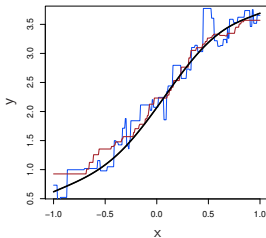# KNN estimates for a simulation from a linear model
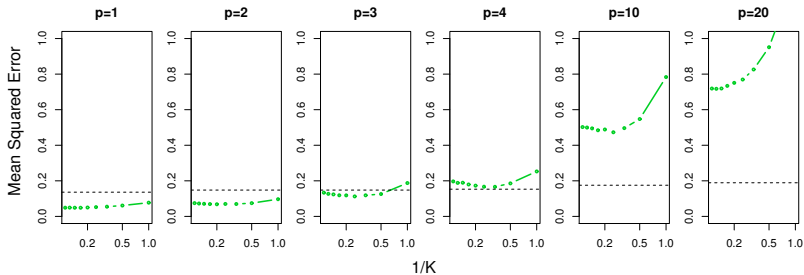
$K = 1$

$K = 9$

# Linear models dominate KNN

# Increasing deviations from linearity

# When there are more predictors than observations, Linear Regression dominates



When $p \gg n$, each sample has no nearest neighbors, this is known as the *curse of dimensionality*. The variance of KNN regression is very large.

# Next time: Classification

Supervised learning with a **qualitative or categorical** response.

Just as common, if not more common than regression:

- *Medical diagnosis:* Given the symptoms a patient shows, predict which of 3 conditions they are attributed to.

- *Online banking:* Determine whether a transaction is fraudulent or not, on the basis of the IP address, client's history, etc.

- *Web searching:* Based on a user's history, location, and the string of a web search, predict which link a person is likely to click.

- *Online advertising:* Predict whether a user will click on an ad or not.