# LINEAR REGRESSION: VALIDATION

## Connor K. Brubaker

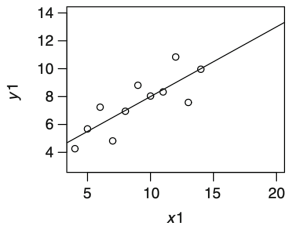### Department of Statistics
### Texas A&M University

# MODEL VALIDATION

Inferences and predictions made from a fitted model only make sense if the assumptions of that model are fulfilled.
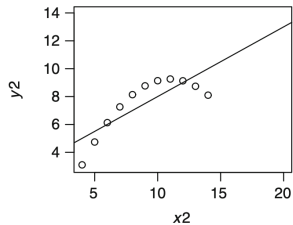
Model validation is the process of evaluating if the assumptions are met.
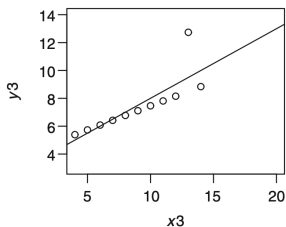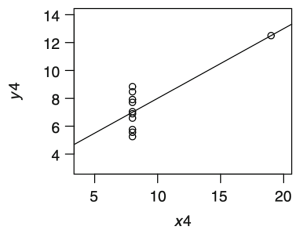
# ANSCOMBE'S DATA SETS

# ANSCOMBE'S DATA SETS

Each of these four data sets results in the same fitted regression line,

$$\hat{Y} = 3 + 0.5X$$

but only one (Data Set 1) satisfies the assumptions of the model–only one of these results in a valid model!

**Don't just look at numerical output of a model. Always check it visually!**

# ASSUMPTIONS OF THE LINEAR MODEL

The simple linear model makes four assumptions:

1. $X$ and $Y$ are linearly related,
2. the errors $\varepsilon_1, \ldots, \varepsilon_n$ are independent of each other,
3. the errors $\varepsilon_1, \ldots, \varepsilon_n$ have a common variance $\sigma^2$ (homoscedasticity), and
4. the errors $\varepsilon_1, \ldots, \varepsilon_n$ are normally distributed with a mean of 0 and variance $\sigma^2$.

We will look at ways of visually evaluating these.

# LINEARITY AND CONSTANT VARIANCE

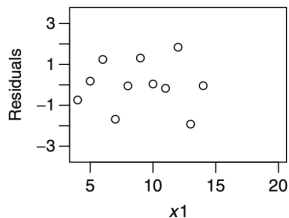Linearity and homoscedasticity is assessed by looking at scatter plots of the fitted residual

$$\hat{\varepsilon}_i = Y_i - \hat{Y}_i$$
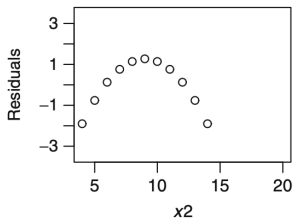
against $X_i$ (called a residual plot).

A residual plot should have no discernable pattern.
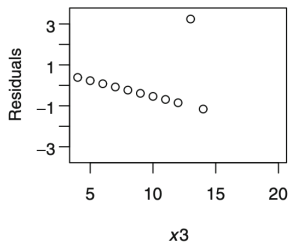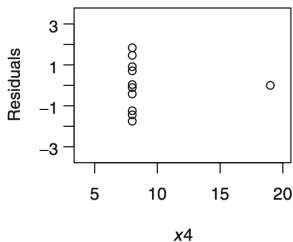
# RESIDUAL PLOTS FOR ANSCOMBE'S DATA SETS
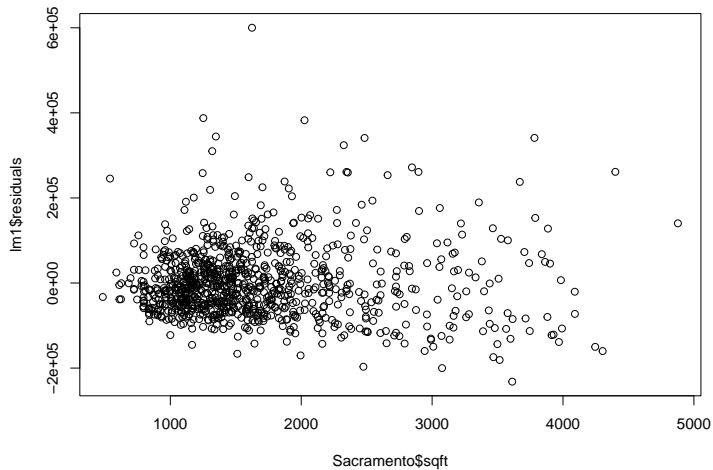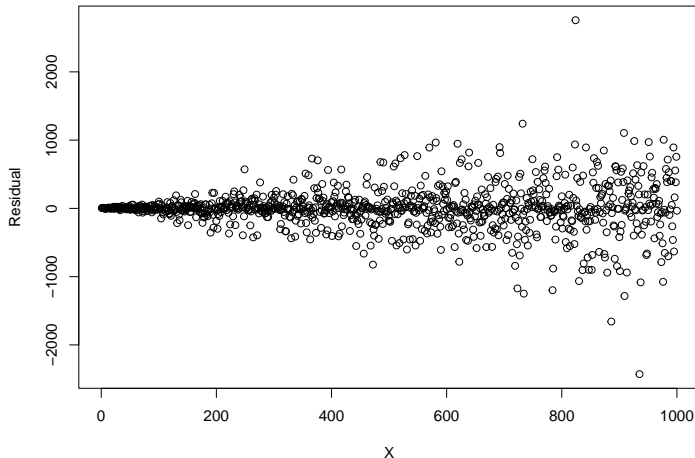
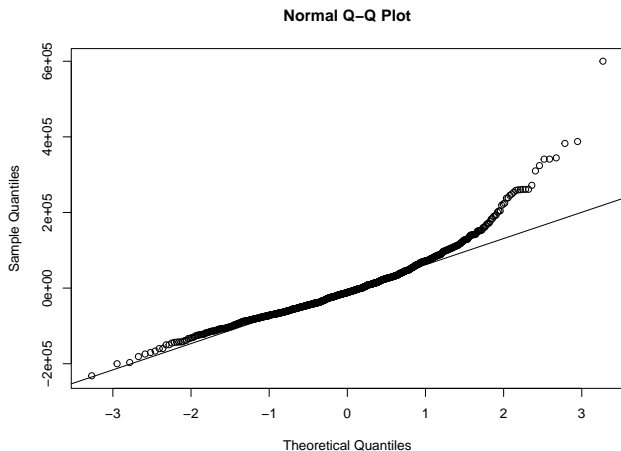# RESIDUAL PLOT FOR HOUSING DATA

# AN EXAMPLE OF HETEROSCEDASTICITY

A residual plot like this indicates that the homoscedasticity assumption has been violated.

# NORMALITY ASSUMPTION

The normality of errors assumption can be checked visually with a normal Q-Q plot like we've seen before. Here is the Q-Q plot of the residuals from the housing data.
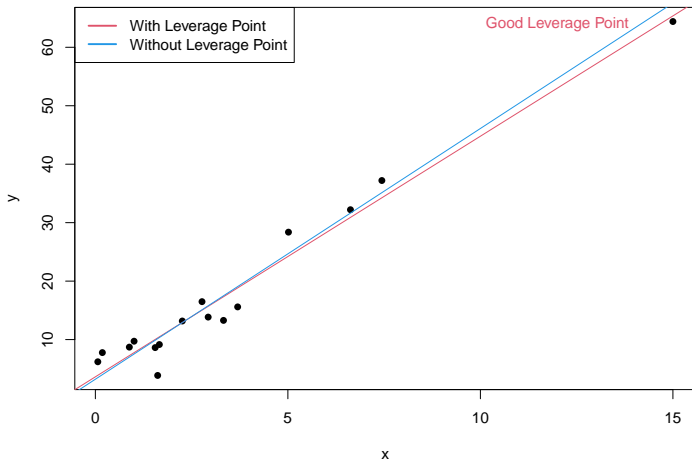


**Normal Q–Q Plot**

# LEVERAGE POINTS

▶ Data points which have considerable influence on the fitted model are called **leverage points**.

▶ They have $X$ values that lie away from the rest of the data points.

▶ Leverage points can be either "good" or "bad".

▶ The leverage of a point can be examined by removing it from the data, refitting the model, and evaluating the effect.
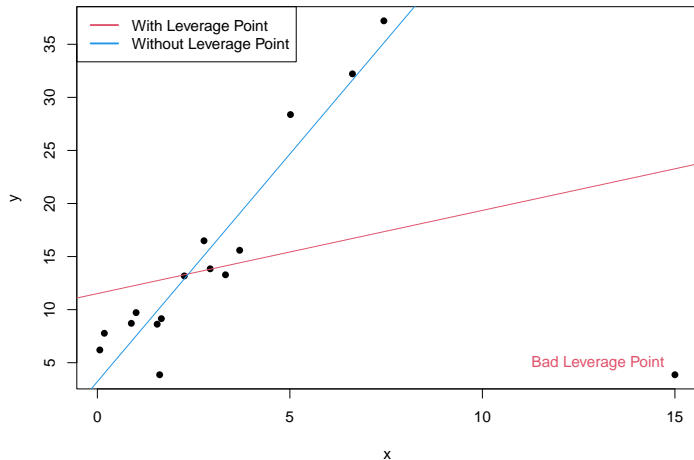
# GOOD LEVERAGE POINTS

A leverage point that follows the linear pattern of the data is a good leverage point.

# GOOD LEVERAGE POINTS

A leverage point that is also an outlier is a bad leverage point.

# DEALING WITH BAD LEVERAGE POINTS

- ▶ Bad leverage points need investigation. They may be candidates for removal from the data due to wrong data entry, bad experimental conditions, etc.
- ▶ If valid, the existence of bad leverage points may suggest a different model is needed (e.g., non-linear model)
- ▶ Even if a leverage point is "good", they do affect standard errors and the value of $R^2$, so they always warrant investigation to check they are valid.