## IDS 702 Influential points

#### Outliers

- rest of the data for a given variable
- (and corrected if possible)
- it



An outlier is a data point whose value does not follow the general trend of the

Outliers should be investigated as potential data errors/implausible values

• If the value does not result from a data error/implausible value, do not remove





# Influential points

- $R^2$ , RSE)

Individual observations can have a large impact on the model (estimates, SE,

• Sometimes the points are obvious from EDA, but other times they are not

#### Note about the model in matrix notation

#### Leverage

- Points with extreme predictor/covariate/feature values are called high leverage points
- variables for the *i*th observation are from those of other observations
- of the hat matrix:  $\mathbf{H} = \mathbf{X}(\mathbf{X}^{T}\mathbf{X})^{-1}\mathbf{X}^{T}$

• 
$$0 \le h_{ii} \le 1$$
 and  $\sum_{i=1}^{n} h_{ii} = p + 1$ 

• The leverage score measures how far away the values of the independent

• The leverage score for an observation is defined as the *i*th diagonal element

# High leverage: what to do?

- Make sure they do not result from data entry errors
- Make sure you look at the impact of those points on the estimates: just regression!
- influential
  - This depends on the value of y

because a point is high leverage does not mean it will have a large effect on

When a point has a large effect on the regression, we say the observation is

### Cook's distance

- Quantifies the influence of the *i*th observation
- Measures how far, on average, predicted values will move if a given observation is removed from the dataset
- $\hat{y}_{j(i)}$  is the predicted value after excluding the *i*th observation

$$D_{i} = \sum_{j=1}^{n} \frac{(\hat{y}_{j} - \hat{y}_{j(i)})^{2}}{s_{e}^{2}(p+1)}$$

## Large Cook's distance: what to do?

- General consensus is that  $D_i > 1$  indicates an observation is an influential value, but we generally pay attention to  $D_i > 0.5$
- For each observation with high Cook's distance, fit the model with and without that observation, and compare the results

# **Diagnostic plots in R**

plot(model, which=5) gives residuals vs leverage plot



Fitted values Im(InfctRsk ~ Stay + Age + Beds)



Leverage Im(InfctRsk ~ Stay + Age + Beds)



Fitted values



Leverage Im(InfctRsk ~ Stay + Age + Beds)

#### Summary: What to do with outliers/influential observations?

- Make sure the observation does not arise from a data entry error
  - If it does, it can be changed or excluded
- May want to report results with and without influential observation(s)
  - With a large sample size, you may not see a difference
- Be wary of numeric cutoffs!

