

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

Aside: Adjusted  $R^2$  ( $R_{adj}^2$ ):

$$R_{adj}^2 = 1 - \frac{SSE/(n-p-1)}{SST/(n-1)} = 1 - \frac{(n-1)}{(n-p-1)}(1-R^2)$$

- \*  $R_{adj}^2 < R^2$ , but as  $n \rightarrow \infty$ ,  $R_{adj}^2 \rightarrow R^2$
- \* This adjusted version penalizes you for adding unimportant explanatory variables into the model. Whereas  $R^2$  will always increase when additional explanatory variables are added to a model, the adjusted  $R^2$  may increase or decrease.

If the reduction in SSE brought about by the addition of extra explanatory variables is not big enough to offset the penalty, then  $R_{adj}^2$  will go down. But, if the reduction in SSE is big enough,  $R_{adj}^2$  will go up.

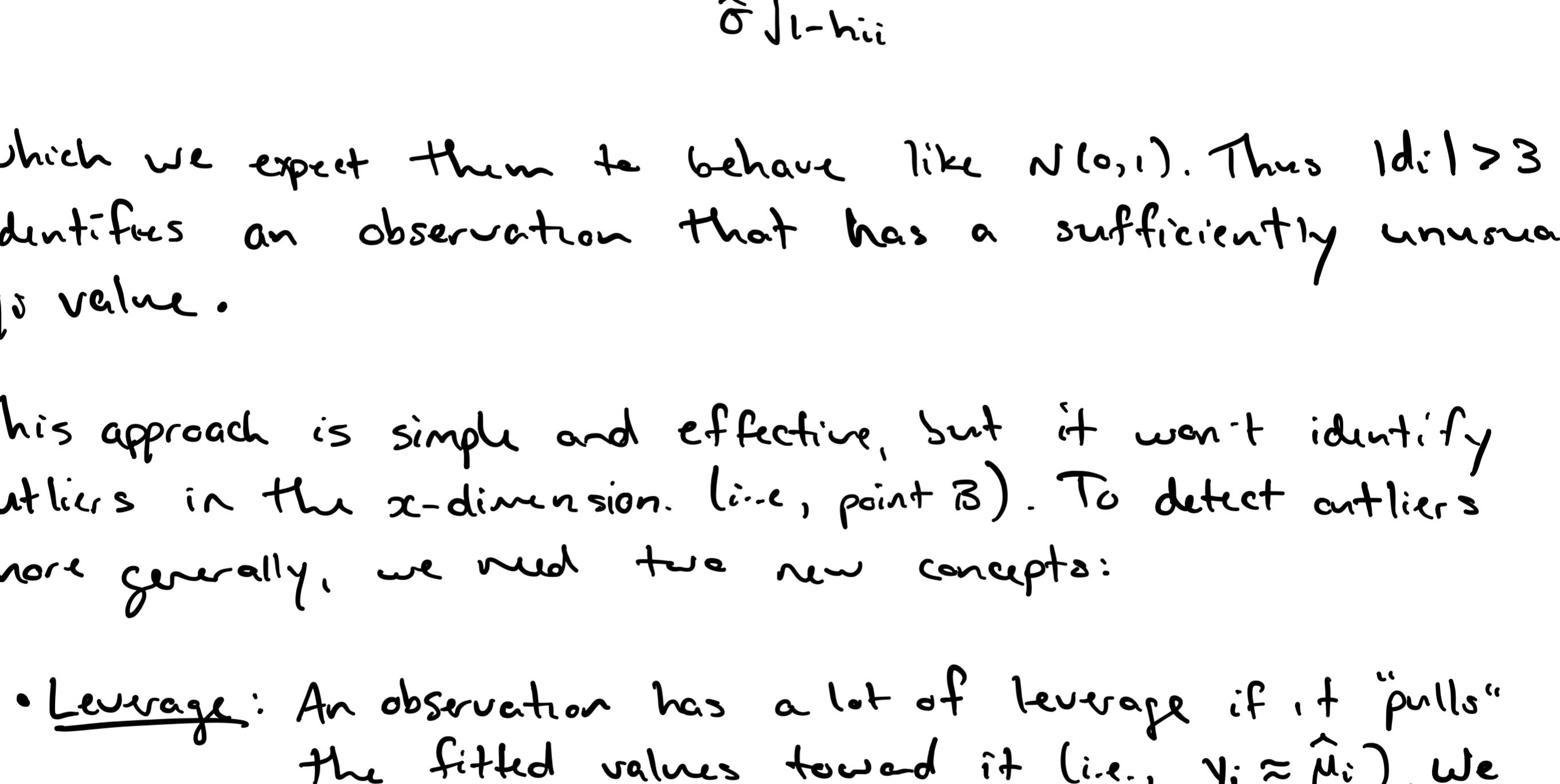
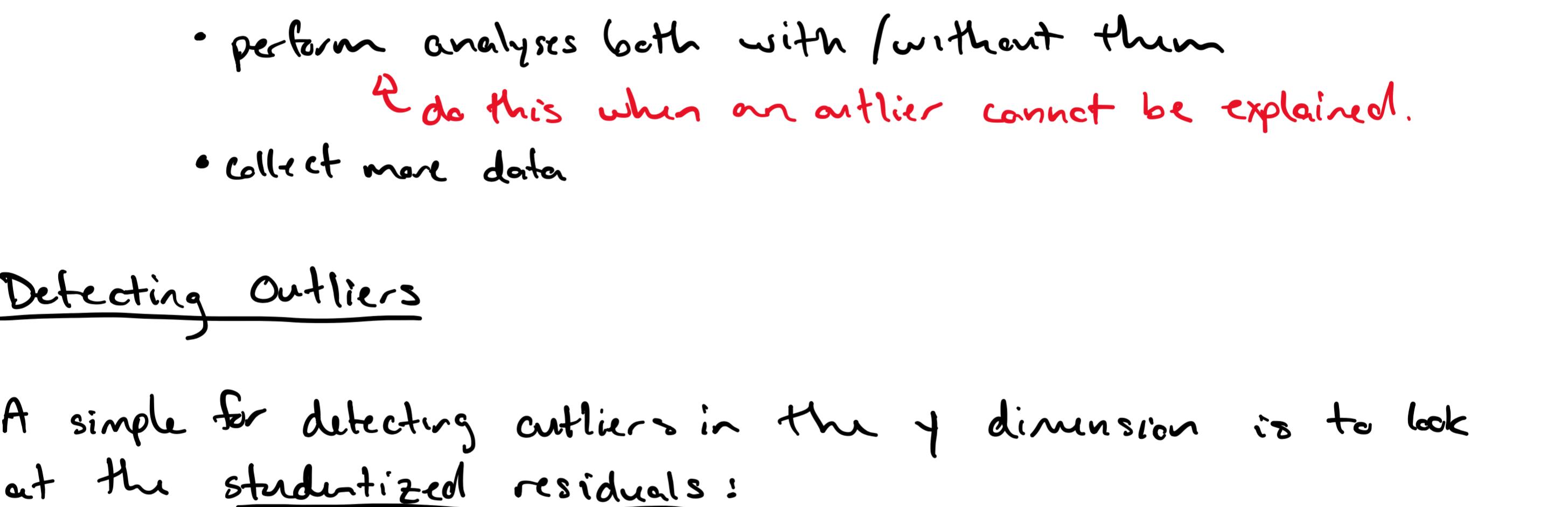
- \* Note that  $R_{adj}^2$  doesn't have a useful interpretation, but it does a better job of describing goodness-of-fit than the ordinary  $R^2$ .

### The Effect of Individual Observations

[Q]: Do all observations look like they come from the same model?

We call an observation  $(y_i, x_{i1}, x_{i2}, \dots, x_{ip})$  an outlier if it differs substantially from the other observations. Since several variables are involved we must distinguish outliers in the  $y$  (response) dimension from outliers in the  $x$  (explanatory variable) dimension.

To visualize this distinction assume we have one  $y$  and one  $x$ :



Outliers can arise for many reasons (i.e., data recording issues or unexpected natural phenomenon), but no matter the reason, we want to be able to identify them. We may then decide to:

- delete them
- correct them
- perform analyses both with / without them
- do this when an outlier cannot be explained.
- collect more data

### Detecting Outliers

A simple for detecting outliers in the  $y$  dimension is to look at the studentized residuals:

$$d_i = \frac{e_i}{\hat{\sigma} \sqrt{1-h_{ii}}}$$

which we expect them to behave like  $N(0,1)$ . Thus  $|d_i| > 3$  identifies an observation that has a sufficiently unusual  $y_i$  value.

This approach is simple and effective, but it won't identify outliers in the  $x$ -dimension (like, point B). To detect outliers more generally, we need two new concepts:

- Leverage: An observation has a lot of leverage if it "pulls" the fitted values toward it (i.e.,  $\hat{y}_i \approx \bar{y}_i$ ). We will see that the extent of the "pull" is determined by the  $x$ 's.
- Influence: An observation has a major influence if the fitted regression model is significantly altered when the observation is omitted. (i.e. point C)

Recall that the fitted values are:

$$\hat{y} = X \hat{\beta} = X(X^T X)^{-1} X^T \bar{y} = H \bar{y}$$

For a given observation  $i$  the fitted value is:

$$\hat{y}_i = h_{ii} \bar{y} + \sum_{j \neq i} h_{ij} y_j$$

in element of  $\hat{y}$

( $i$ th row of  $H$ )  $\times \bar{y}$

The coefficient  $h_{ii}$  indicates how heavily  $y_i$  contributes to the fitted value  $\hat{y}_i$ . If  $h_{ii}$  is large (compared to the other  $h_{ij}$ 's) then  $y_i$  dominates  $\hat{y}_i$ .

It can be shown (see Assy) that  $0 \leq h_{ii} \leq 1$  and note that if  $h_{ii} \approx 1$  then  $\text{Var}[e_i] = \sigma^2(1-h_{ii}) \approx 0$ . Thus, if  $h_{ii} \approx 1$  then  $y_i \approx \hat{y}_i$  and we'd say that observation  $i$  has high leverage.

We will define  $h_{ii}$  to be the leverage of observation  $i$ . We will say that observation  $i$  has high leverage if:

- $h_{ii} \approx 1$
- $h_{ii} > 2\bar{h}$  and / or

$$\bar{h} = \frac{1}{n} \sum_{i=1}^n h_{ii}$$

the leverage is much bigger than the other points' leverages.

\* Since  $H$  depends only on the  $x$ 's, the concept of leverage helps find outliers in the  $x$  dimension.