

LECTURE 5

BOGDAN ICHIM

Assessing the Accuracy of the Model

(via Dispersion Analysis)

Consider the sample mean $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$. We want to *understand* $y_i - \bar{y}$.

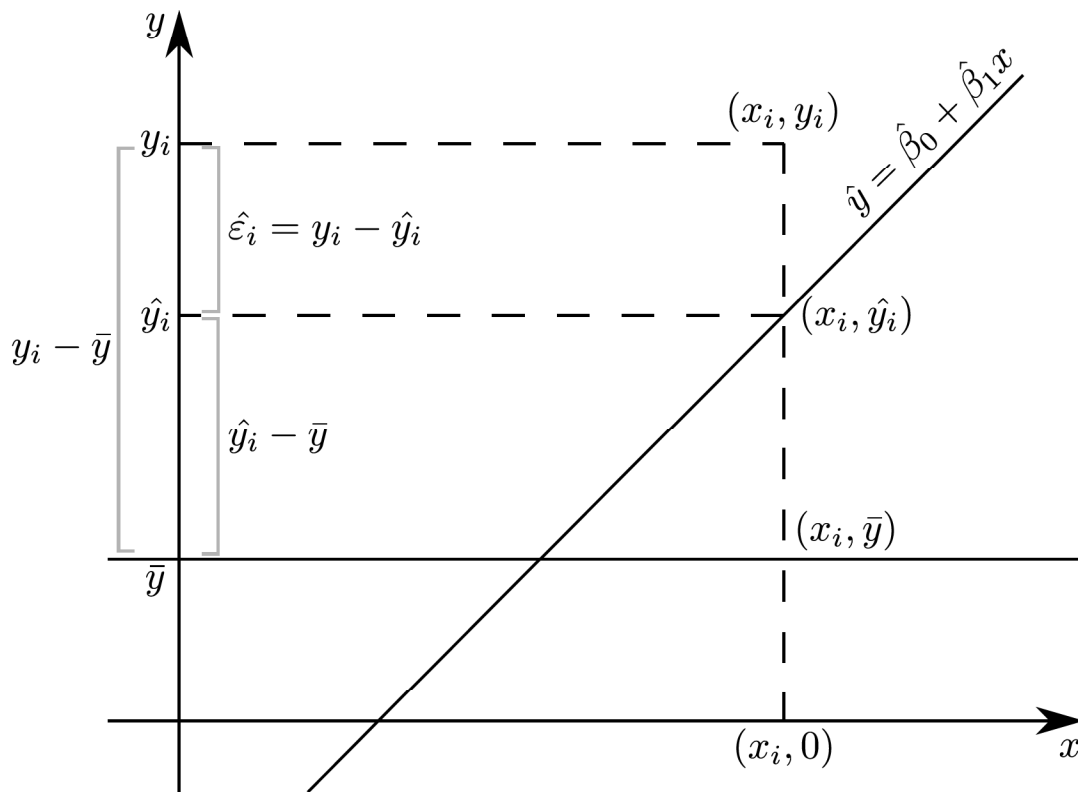


FIGURE 1

Recall from **Lecture 5**:

$$\begin{cases} \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \\ \hat{\epsilon}_i = y_i - \hat{y}_i \end{cases} \implies y_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + \hat{\epsilon}_i \quad (1)$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (2)$$

Remark 1. We have that

$$y_i - \bar{y} = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y}) \quad (3)$$

We compute

$$\begin{aligned}(1) + (2) &\implies y_i = \bar{y} + \hat{\beta}_1 x_i - \hat{\beta}_1 \bar{x} + \hat{\epsilon}_i \\ &\implies y_i - \bar{y} = \hat{\beta}_1 (x_i - \bar{x}) + \hat{\epsilon}_i \quad \Big| \quad \sum_{i=1}^n \\ &\implies 0 = \hat{\beta}_1 \cdot 0 + \sum_{i=1}^n \hat{\epsilon}_i \\ &\implies \boxed{\sum_{i=1}^n \hat{\epsilon}_i = 0.}\end{aligned}$$

Then we have

$$\begin{aligned}&\implies \sum_{i=1}^n \hat{\epsilon}_i = 0 \\ &\implies \sum_{i=1}^n (y_i - \hat{y}_i) = 0 \\ &\implies \sum_{i=1}^n y_i - \sum_{i=1}^n \hat{y}_i = 0 \\ &\iff \sum_{i=1}^n y_i = \sum_{i=1}^n \hat{y}_i \\ &\iff \boxed{\bar{y} = \bar{\hat{y}}.}\end{aligned}$$

Recall that

$$\begin{aligned}\frac{\partial \text{RSS}}{\partial \hat{\beta}_1} = 0 &\iff \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i \\ &\implies \sum_{i=1}^n x_i (y_i - \underbrace{(\hat{\beta}_0 + \hat{\beta}_1 x_i)}_{\hat{y}_i}) = 0 \\ &\implies \sum_{i=1}^n x_i \underbrace{(y_i - \hat{y}_i)}_{\hat{\epsilon}_i} = 0 \\ &\implies \boxed{\sum_{i=1}^n x_i \hat{\epsilon}_i = 0}\end{aligned}$$

We also have that

$$\begin{aligned}
 \hat{y}_i &= \hat{\beta}_0 + \hat{\beta}_1 x_i \quad \Big| \cdot \hat{\epsilon}_i \\
 \implies \hat{y}_i \hat{\epsilon}_i &= \hat{\beta}_0 \hat{\epsilon}_i + \hat{\beta}_1 x_i \hat{\epsilon}_i \quad \Big| \sum_{i=1}^n \\
 \implies \sum_{i=1}^n \hat{y}_i \hat{\epsilon}_i &= \hat{\beta}_0 \underbrace{\sum_{i=1}^n \hat{\epsilon}_i}_{=0} + \hat{\beta}_1 \underbrace{\sum_{i=1}^n x_i \hat{\epsilon}_i}_{=0} \\
 \implies \boxed{\sum_{i=1}^n \hat{y}_i \hat{\epsilon}_i = 0.}
 \end{aligned}$$

Then

$$\sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) = \sum_{i=1}^n \hat{\epsilon}_i(\hat{y}_i - \bar{y}) = \underbrace{\sum_{i=1}^n \hat{\epsilon}_i \hat{y}_i}_{=0} - \bar{y} \underbrace{\sum_{i=1}^n \hat{\epsilon}_i}_{=0} = 0.$$

Now we ready to make the final computation. We square both sides of formula (3).

$$\begin{aligned}
 (3)^2 \implies (y_i - \bar{y})^2 &= (y_i - \hat{y}_i)^2 + (\hat{y}_i - \bar{y})^2 + 2(y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) \quad \Big| \sum_{i=1}^n \\
 \implies \sum_{i=1}^n (y_i - \bar{y})^2 &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2} + 2 \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y})}_{=0} \\
 \implies \underbrace{\sum_{i=1}^n (y_i - \bar{y})^2}_{\text{variation in the } y \text{ data}} &= \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{\text{variation in the } \hat{y} \text{ data}} + \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i - \overbrace{(y - \bar{y})}^{\hat{\epsilon}_i})^2}_{\text{variation in the } \hat{\epsilon} \text{ data}}
 \end{aligned}$$

We introduce several **notations**

$$\begin{aligned}
 \sum_{i=1}^n (y_i - \bar{y})^2 &= \text{TSS} = \text{total sum of squares;} \\
 &= \text{RSS}_0 = \text{residual sum of squares for a model with 0 predictors;} \\
 &= \text{SST} = \text{sum of squares total;} \\
 \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 &= \text{SSR} = \text{sum of squares explained by regression;} \\
 &= \text{ESS} = \text{explained sum of squares;} \\
 \sum_{i=1}^n (y_i - \hat{y}_i)^2 &= \text{RSS} = \text{residual sum of squares;} \\
 &= \text{RSS}_k = \text{residual sum of squares for a model with } k \text{ predictors;} \\
 &= \text{SSE} = \text{sum of squares of errors.}
 \end{aligned}$$

Above we have computed the **Regression Identity (Linear Model Identity)**

$$\boxed{\text{TSS} = \text{SSR} + \text{RSS}.}$$

The Coefficient of Determination

Note that RSS, MSE and RSE are absolute measurements of the performance of the regression model. But since they are measured in the units of Y , it is not always clear what is a good RSS, MSE or RSE.

Definition 2. The **coefficient of determination** R^2 (pronounced “R-squared”) provides an alternative measure for the model’s performance, which is independent of the scale of Y . It is defined by the following formula

$$R^2 = 1 - \frac{\text{RSS}}{\text{TSS}} = \frac{\text{SSR}}{\text{TSS}}.$$

In the case of the linear model $R^2 \in [0, 1]$ and may be interpreted as the proportion of the total variation in the response variable Y which is explained by the model.

Remark 3. (1) $R^2 = 1$ indicates a perfect fit. We have

$$\begin{aligned} R^2 = 1 &\iff \\ \text{RSS} = 0 &\iff \\ \sum_{i=1}^n (y_i - \hat{y}_i)^2 = 0 &\iff \\ y_i = \hat{y}_i \quad \forall i = \overline{1, n}. & \end{aligned}$$

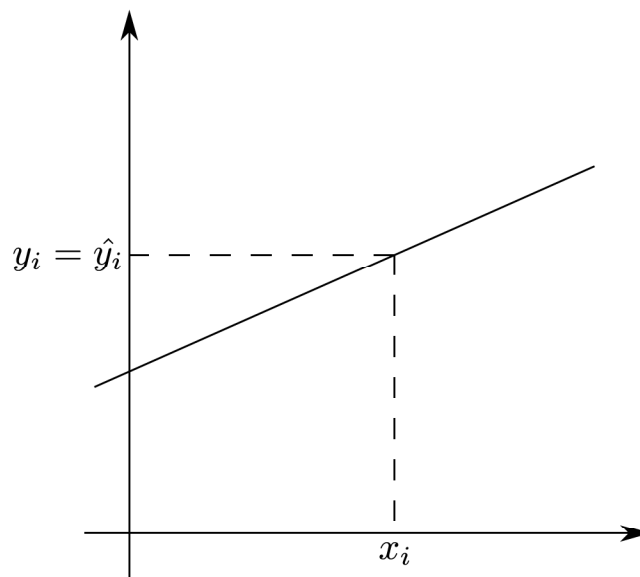


FIGURE 2

(2) If $R^2 \approx 0$ then the regression did not explain much of the variability of the response.

$$\begin{aligned}
 R^2 = 0 &\iff \\
 SSR = 0 &\iff \\
 \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = 0 &\iff \\
 \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = 0 &\iff \\
 \hat{y}_i = \bar{y} \quad \forall i = \overline{1, n}.
 \end{aligned}$$

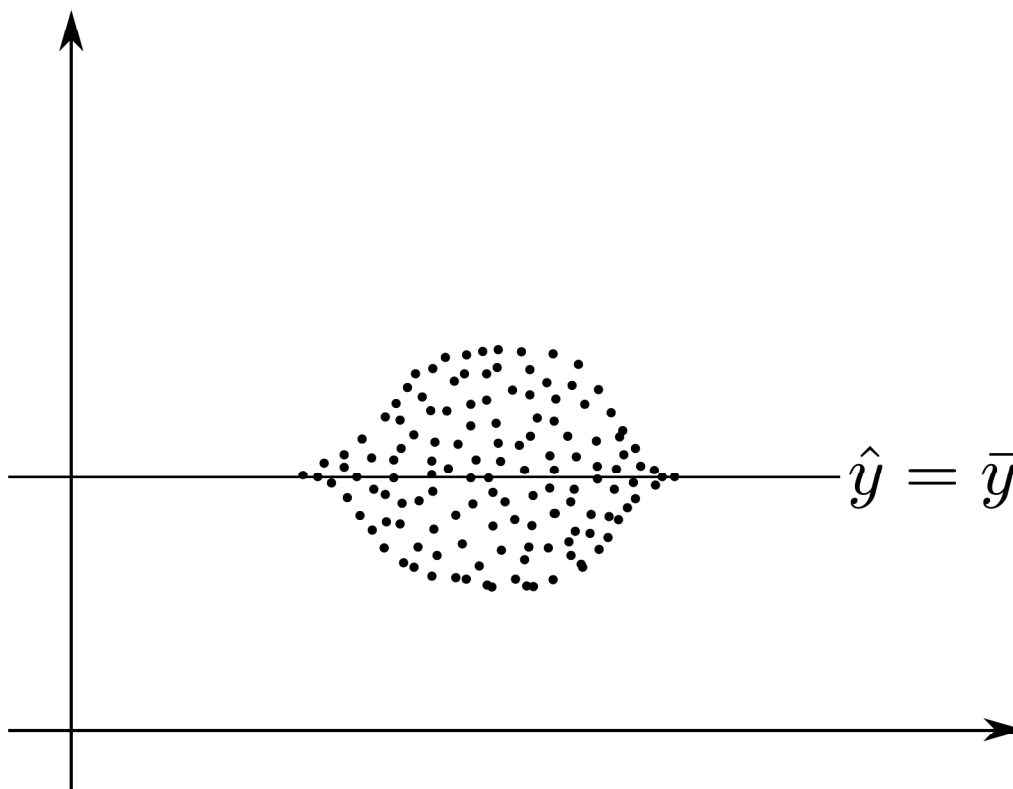


FIGURE 3

- (3) R^2 will always increase when more variables are added to the model, no matter how useless those variables are for prediction.
- (4) The formulas presented above may vary, depending on the implementation. For example in R if the intercept is removed, then the formula

$$R_0^2 = 1 - \frac{\text{RSS}}{\sum_{i=1}^n y_i^2}$$

is used.

This formula may be seen as anachronistic since makes no sense to use it for other machine learning models.

Adjusted Coefficient of Determination

Definition 4. The **adjusted coefficient of determination** $\overline{R^2}$ (pronounced “adjusted R-squared”) is defined by the following formula

$$\overline{R^2} = 1 - \frac{\frac{\text{RSS}}{n - k - 1}}{\frac{\text{TSS}}{n - 1}},$$

where k = the number of explanatory variables = the number of predictors.

Remark 5. (1) $\overline{R^2} \leq R^2$.

(2) It is possible that $\overline{R^2} \leq 0$.

(3) $\overline{R^2}$ may be used for choosing one model among several models with similar performance. In principle, one should choose the one with the highest $\overline{R^2}$.

(4) Assume

$$\left. \begin{array}{l} \text{RSS}_{k_1} = \text{RSS}_{k_2} \\ k_1 < k_2 \end{array} \right\} \xrightarrow{\text{TSS}_{k_1} = \text{TSS}_{k_2} \text{ (always)}} \overline{R^2}_{k_1} > \overline{R^2}_{k_2}.$$

Occam’s razor: One should prefer the model with k_1 predictors since it is the simpler model.

Elements of Hypothesis Testing

We test the **Null Hypothesis**

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0$$

against the **Alternative Hypothesis**

$$H_a : \text{at least one } \beta_j \neq 0$$

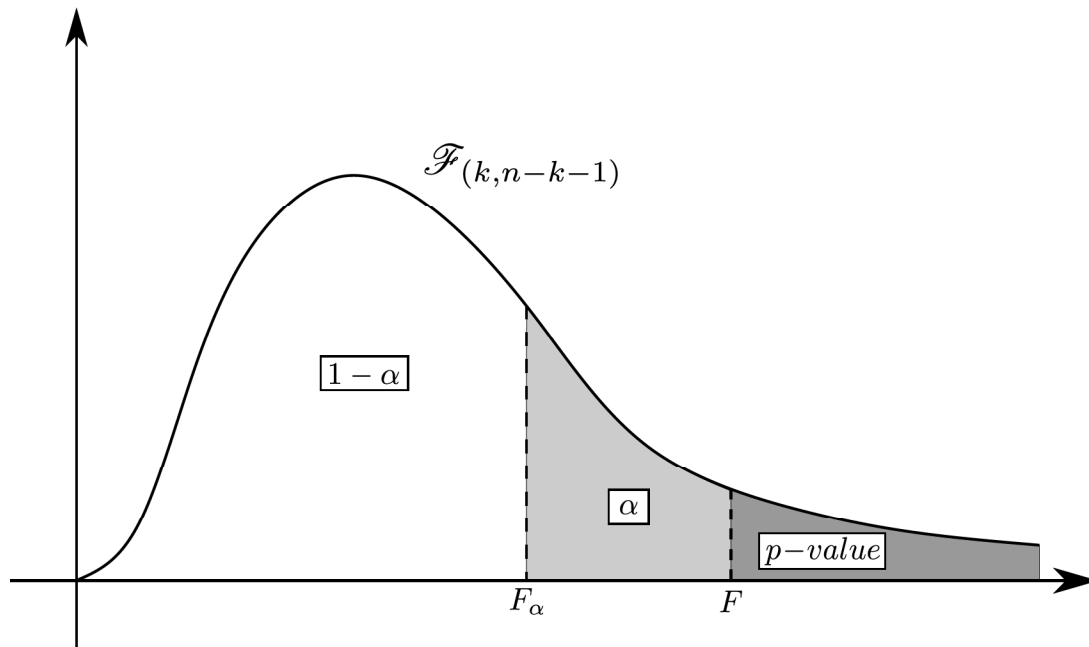
in order to answer the following.

Question 6. Is at least one of the predictors useful in predicting the response?

This hypothesis test is performed by computing the **F-statistic**

$$F = \frac{\frac{\text{RSS}_0 - \text{RSS}_k}{k}}{\frac{\text{RSS}_k}{n - k - 1}}.$$

When H_0 is true and the errors ϵ_i have a normal distribution, the F-statistic follows an F-distribution (Fisher-Snedecor distribution).

FIGURE 4. Fisher-Snedecor distribution $\mathcal{F}(k, n - k - 1)$

We consider the following.

$$\alpha = \text{significance level} = \int_{F_\alpha}^{\infty} \mathcal{F}(k, n - k - 1)$$

$$1 - \alpha = \text{confidence level} = \int_0^{F_\alpha} \mathcal{F}(k, n - k - 1)$$

$$p\text{-value} = \text{significance } F = \int_F^{\infty} \mathcal{F}(k, n - k - 1)$$

Decisions based on p-values:

- $p\text{-value} \leq \alpha \iff F \geq F_\alpha \implies$ reject H_0 and then H_a is true;
- $p\text{-value} > \alpha \iff F < F_\alpha \implies H_0$ is true.