

Some useful quantities

$$SSE = RSS = \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \quad R^2 = 1 - \frac{SSE}{S_{yy}} = 1 - \frac{SSR}{TSS}$$

$$S_{yy} = TSS = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2$$

$$S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$SSR = SSE - S_{yy} \Rightarrow$$

$$S_{yy} = \underbrace{SSE}_{\text{irreducible error}} + \underbrace{SSR}_{\text{reducible error}}$$

Training and validation

In ML we don't really care about the past (statistically).

We split out data in two (or even three) sets. E.g.: $\text{data} = \text{train} \mid \text{validation}$

Now we need comparison of means between samples (training set and validation set).

E_{MSE} : expectation for the mean of SSE

n.b $MSE = \frac{1}{n} SSE$

$E_{MSE}(\bar{Y} - \hat{\bar{Y}})$: Expectation for difference of mean between \bar{Y} and our approximation.

$$E_{MSE}(\bar{Y} - \hat{\bar{Y}})^2 = \underbrace{\text{Var}(\bar{Y}) + \text{Bias}(\hat{\bar{Y}})^2}_{\text{reducible error}} + \underbrace{\text{Var}(\epsilon)}_{\text{irreducible error}}$$

RSE: $\sqrt{\frac{1}{n-2} SSE}$, unbiased estimator for unknown data!

RSE is run with the model we trained on new data (not even the validation).

ie

train, validation, test

during training

(R^2 , MSE)

↑
RSE for quality-checks

Feature engineering

$$\alpha(\lambda u + \mu v) = \lambda \alpha(u) + \mu \alpha(v)$$

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

- i, Test significance of predictors, remove those that fail the test.
- ii, Add new features (possibly from other data sources)
- iii, Break the rules! Add non-linearity.

$$\underline{f(x_1, x_2)} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2$$

$$f(x_1, x_2, \overset{\uparrow}{x_3}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

Linearisation!

$\hat{f}(x) \Rightarrow$ numbers (real) regression

$\hat{f}(x) =$ labels (0, 1, 2) classification

 ↓ ↓ ↓
 "red" "green" "purple"

labels are qualitative, categorical data