

## Model

True model :  $y = f(x) + \epsilon$

- Unknown to us
- Generates data that we observe
- Expanded form :

$$y^{(i)} = w_0 + w_1 x_1^{(i)} + w_2 x_2^{(i)} + \dots + w_p x_p^{(i)} + \epsilon^{(i)}$$

where  $i$  denotes the observation index

$p$  = # features

$N$  = # observations

$w$  = coefficients (parameters to be estimated)

Model that we wish to estimate (based on observed data) :

$$\hat{y}^{(i)} = \hat{w}_0 + \hat{w}_1 x_1^{(i)} + \hat{w}_2 x_2^{(i)} + \dots + \hat{w}_p x_p^{(i)}$$

where  $\hat{w}$  represents estimates of coefficients

$\hat{y}$  represent the response that we get after estimating the model.

Note: The  $\hat{}$  is dropped in further sections to be concise.

Matrix form :

$$\mathbf{X} \mathbf{w} = \mathbf{y}$$

$$\begin{bmatrix} | & x_1^{(1)} & x_2^{(1)} & \dots & x_p^{(1)} \\ | & x_1^{(2)} & x_2^{(2)} & \dots & x_p^{(2)} \\ | & \vdots & \vdots & \ddots & \vdots \\ | & x_1^{(N)} & x_2^{(N)} & \dots & x_p^{(N)} \end{bmatrix} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_p \end{bmatrix} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix}$$

$(N \times p)$        $(p \times 1)$        $(N \times 1)$

OR

$$w^T \mathbf{X} = \mathbf{y}$$

$$\begin{bmatrix} w_0 & w_1 & \dots & w_p \end{bmatrix} \begin{bmatrix} | & x_1^{(1)} & x_2^{(1)} & \dots & x_p^{(1)} \\ | & x_1^{(2)} & x_2^{(2)} & \dots & x_p^{(2)} \\ | & \vdots & \vdots & \ddots & \vdots \\ | & x_1^{(N)} & x_2^{(N)} & \dots & x_p^{(N)} \end{bmatrix} = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(N)} \end{bmatrix}$$

$(1 \times p)$        $(p \times N)$        $(1 \times N)$

### Interpreting the Coefficients

Numeric predictor : One unit change in  $x_i$  corresponds to  $w_i$  units of change in  $y$ .

Categorical predictor : Selecting  $x_i$  category (i.e.  $x_i=1$ ) corresponds to  $w_i$  units of change in  $y$  as compared to the default category (absorbed in the intercept)

## Fitting the Model

True model:  $Y = Xw + \epsilon$  (unknown)

We wish to estimate  $\hat{w}$  based on observed data that will give us  $\hat{Y} = X\hat{w}$  that minimizes RSS.

$$RSS = \sum_{i=1}^N (y^{(i)} - \hat{y}^{(i)})^2$$

(residual sum  
of squares)

Minimizing RSS is equivalent to minimizing mean squared error (MSE) or root mean squared error (RMSE)

i) Closed-form solution (analytical solution that minimizes RSS)

$$\begin{aligned} Xw &= y && \text{For simple linear regression case} \\ \Rightarrow X^T X w &= X^T y && \text{with only one predictor, there is} \\ \Rightarrow w &= (X^T X)^{-1} X^T y && \text{formula to directly compute } w_0 \text{ and } w_1 \end{aligned}$$

ii) Gradient descent (minimizes the cost function MSE ( $\approx$  RSS))

$$w_{\text{new}} = w_{\text{old}} - \gamma \cdot \nabla_w C(w)$$

where,

$$C(w) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \hat{y}^{(i)})^2 \quad [\text{cost function}]$$

$$\text{and } \nabla_w C(w) = -\frac{2}{N} \sum X^{(i)\top} (y^{(i)} - X^{(i)}w) \quad [\text{gradient of cost fn}]$$

## Statistics of Linear Regression

- Accuracy of coefficients

- i) Estimate coefficients: using closed form solution.
- ii) Compute standard errors: of each coefficient using (SE)  
a formula (see ISLR 3.1.2)  
that depends on N.
- iii) Compute 95% CI: for each coefficient
$$w_i = [\hat{w}_i - 2 \cdot SE(\hat{w}_i), \hat{w}_i + 2 \cdot SE(\hat{w}_i)]$$
- iv) Compute p-value: two hypothesis are tested

p-value of model

$$H_0: w_0 = w_1 = w_2 = \dots = w_p = 0$$

$$H_A: \text{At least one } w_i \neq 0$$

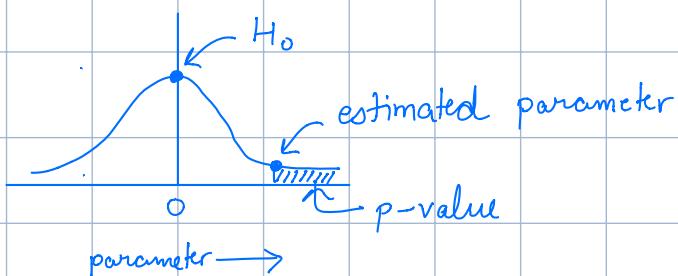
p-value of individual coefficients

$$H_0: w_i = 0$$

$$H_A: w_i \neq 0$$

- F-statistic is computed  
to test if at least one predictor is associated with the response. Needs to be  $> 1$  given N is large enough

- To test the above hypothesis, a model is fit with and without each predictor (both include all other predictors)



$H_0$  is rejected if p-value is small because it suggests it is unlikely for the parameter to take on this value given  $H_0$  is true.

## • Accuracy of the model

- Residual standard error (RSE) is the average deviation of RSS:

$$RSE = \sqrt{\frac{RSS}{N-p-1}} \quad (\text{similar to RMSE})$$

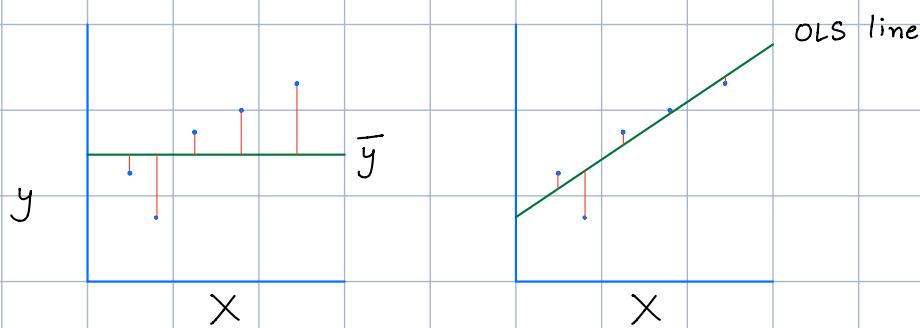
RSE estimates the average  $\epsilon$  if true model were known.

Disadvantage: RSE's scale makes it difficult to assess.

## - $R^2$

Range of  $R^2$  is  $[0, 1]$  so it is easier to assess.

Tells the proportion of variance explained by the model.



$$TSS = \text{Var}(y)$$

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

= sum of red lines  
in above plot

RSS = Variance left after fitting model

$$= \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$$

= sum of red lines in above plot

variance explained by model  
(or removed)

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

## Assumptions

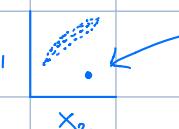
- $X$  are iid
- $y$  are iid
- $e \stackrel{iid}{\approx} N(0, \sigma^2)$

mean 0                                  constant variance (no heteroskedasticity)

Validate linear model by plotting  $(y - \hat{y})$  against  $\hat{y}$  (residual)

and if the plot has patterns, make sure that :

- $y$  is linearly dependent of  $X$  (if not, try  $\sqrt{x_i}$ ,  $\log x_i$  or  $x_i^k$ )
- residuals are uncorrelated
- residuals has no heteroscedasticity (if not, try  $\sqrt{y}$  or  $\log y$ )
- $y$  has no outliers (outliers are visible in residual plot or studentized residual plot where  $|studentized\ residuals| > 3$ )
- $X$  has no high leverage points (outliers for a predictor)  
(points can look normal for a predictor but can be high leverage if we plot multiple predictors)

$x_1$  |   
       .  
 $x_2$

Compute leverage statistic to identify such points.

- No multicollinearity (compute  $VIF = \frac{1}{1 - R^2_{x_j|x_{-j}}}$ ;  $VIF > 5 \Rightarrow$  multicollinearity)
- Residuals are normally distributed

## Parametric vs Non-parametric methods

Linear regression	KNN
<ul style="list-style-type: none"><li>- Better when the relationship b/w X and Y is truly linear</li></ul>	<ul style="list-style-type: none"><li>- Slightly worse than linear regression when relationship b/w X and Y is truly linear.</li></ul>
<ul style="list-style-type: none"><li>- When relationship is not linear, model is highly biased.</li></ul>	<ul style="list-style-type: none"><li>- Not biased</li></ul>
<ul style="list-style-type: none"><li>- Better for interpretability</li></ul>	<ul style="list-style-type: none"><li>- Need to have <math>&gt; \sim 100</math> observations per predictor to be better than linear regression.</li></ul>

