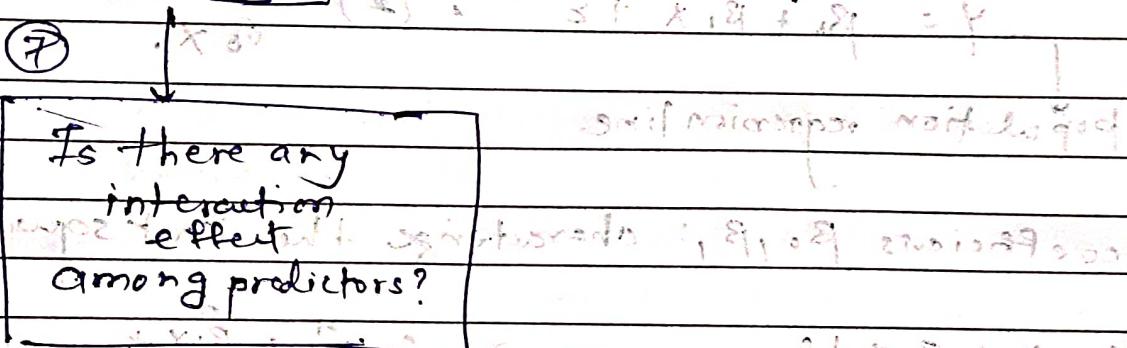
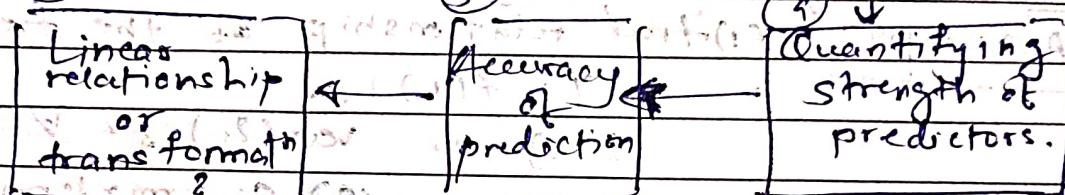
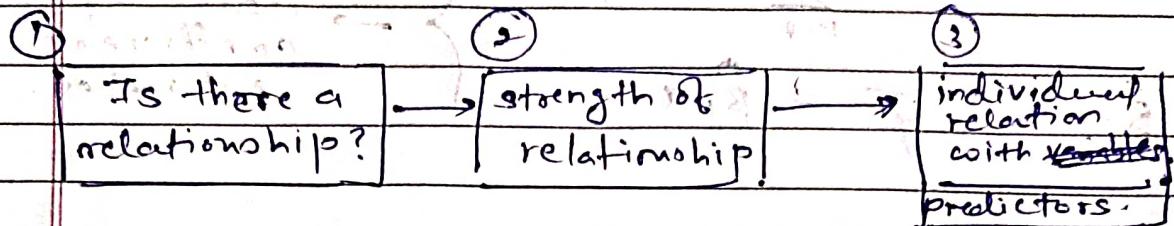


## Chapter 3: Linear Regression.

- Tool for predicting quantitative response.
- very simple approach.



- ### 3.1 Simple Linear Regression.
- predicting a quantitative response on the basis of a single predictor variable  $x$ .

$y = \beta_0 + \beta_1 x$ . ( $\beta_0, \beta_1$ : model coefficients or parameters.)

prediction:  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$  → average increase in  $y$  associated with one unit increase in  $x$ .

$$\text{Let } \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \\ e_i = y_i - \hat{y}_i \text{ (with residual)}$$

Residual Sum of Squares:

$$RSS = \sum_{i=1}^n e_i^2 \rightarrow \text{minimizing}$$

$$\sum_{i=1}^n e_i^2 = (\hat{y}_i - y_i)^2 \rightarrow \text{minimizing w.r.t } \hat{\beta}_0, \hat{\beta}_1$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

least square coefficient estimates.

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

3.1.2 : Assessing the accuracy of the coefficient estimates.

Assumption : (1) true relationship of the form :

$$Y = f(x) + \varepsilon \quad E(\varepsilon) = 0, \quad \text{var}(\varepsilon) = \sigma^2 \quad \forall i$$

$$Y = \beta_0 + \beta_1 x + \varepsilon, \quad (2) \quad \varepsilon: \text{error term independent of } x.$$

population regression line

coefficients  $\beta_0, \beta_1$  characterize the least squares line.

$$Y = \beta_0 + \beta_1 x + \varepsilon \rightarrow \hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

population regression line. estimates from least squares

$\hat{\beta}_0, \hat{\beta}_1$  unbiased for  $\beta_0, \beta_1$  respectively.  $\rightarrow$  close to true population regression line.

$\pm$  standard error of  $\hat{\beta}_1$ : the avg amount that ~~this~~ it will differ from  $\beta_1$ .

$$SE(\hat{\beta}_0) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)$$

$$SE(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\text{where } \sigma^2 = \text{Var}(\varepsilon)$$

Residual standard error:

$$RSE = \sqrt{\frac{RSS}{n-2}}$$

\* Meaning of 95% confidence interval:   
 range of values s.t. with 95% probability, the range will contain the true unknown value of the parameter.

(using assumption  $\epsilon \sim N(0, \sigma^2)$ )

\* 95% CI for  $\beta_1$ :  $[\hat{\beta}_1 - 2SE(\hat{\beta}_1), \hat{\beta}_1 + 2SE(\hat{\beta}_1)]$

95% CI for  $\beta_0$ :  $[\hat{\beta}_0 - 2SE(\hat{\beta}_0), \hat{\beta}_0 + 2SE(\hat{\beta}_0)]$

$H_0$ : There is no relationship (i.e.  $\beta_1 = 0$ ) between  $x$  &  $y$ .

$H_1$ : There is some relationship ( $\beta_1 \neq 0$ ) between  $x$  &  $y$ .

Under  $H_0$ :  $\hat{\beta}_1 = 0$    
  $t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)} \sim t_{n-2}$

3.1.3: Assessing the accuracy of the model.

To quantify: The extent to which the model fits the data.

(1) Residual Standard Error.

- estimate of the standard deviation of  $\epsilon$ .

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

$$RSE = \sqrt{\frac{RSS}{n-2}} \rightarrow \text{measures the lack of fit. (absolute measure).}$$

(2)  $R^2$  statistic,  $\in [0, 1]$ .

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

For simple linear regression

$$R^2 = \rho^2 \text{ for } \{x_i\}$$

$$TSS := \text{total sum of squares} = \sum_{i=1}^n (y_i - \bar{y})^2$$

$R^2$ : proportion of variability in  $y$  that can be explained using  $x$ .

## 8.2 : Multiple linear regression

doubt

shark attacks & ice cream sales example: Is temperature a latent variable? what if it is some unknown problem.

edit do we have problems such as multicollinearity?

assumption of functional form:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$$

estimates:

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_p x_p$$

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{minimizing w.r.t } \hat{\beta}_j \rightarrow \text{estimates } \hat{\beta}_j$$

Important questions:

(1) Is there a relationship between response & predictors?

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

$H_1$ : at least one  $\beta_j$  is non zero.

$$F = \frac{(TSS - RSS)/p}{RSS/(n-p-1)}$$

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

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

$$F = \frac{RSS}{n-p-1} = \sigma^2$$

$$\text{Under } H_0: E\left[\frac{TSS - RSS}{p}\right] = \sigma^2$$

→ value of F statistics depends on n and p.

\* To test whether particular of the coefficients ( $\beta$ ) are zero.

$$H_0: \beta_{p-q+1} = \beta_{p-q+2} = \dots = \beta_p = 0$$

$$\rightarrow F = \frac{(RSS_0 - RSS)/q}{RSS/(n-p-q)}$$

RSS: RSS for whole model

RSS<sub>0</sub>: RSS for model which excludes q coefficients.

- F statistics test is useful if p is relatively small.

- If  $p > n$ : forward selection or high dimensional models.

(2) Deciding on important variables.

- Classical approaches to variable selection.

- For p features:  $2^p$  models to be considered.



Forward Selection

start with null model (only intercept)

fit 'p' simple linear regressions.

add the predictor with least RSS.

repeat until no predictor

add the predictor with least RSS

until some stopping criterion is satisfied.

→ can be always

used

→ might include variables early that later become redundant.

Backward Selection

remove var with largest p-value

fit for  $p-1$  predictors

remove least significant.

repeat until some stopping criterion is reached.

if  $p > n$  cannot be used.

→ until all p values

below certain threshold.

- combination of forward & backward.

- start with null-model

- go on adding variables, one by one

which provide best fit.

- if p value becomes larger, remove that variable

- until all p values

below certain threshold.

→ mixed selection.

doubt what exactly is the difference between confidence interval & prediction interval?



(3) Model fit: How well does the model fit the data?

$$R^2 = \text{Cor}(\mathbf{y}, \hat{\mathbf{y}})^2$$

- MLR maximizes  $\{\text{Cor}(\mathbf{y}, \hat{\mathbf{y}})\}^2$  among all possible linear models.

$$RSE = \sqrt{\frac{RSS}{n-p-1}} \quad (p=1 : \text{simple linear regression})$$

(4) Prediction: Given a set of predictor values, what response value should we predict and how accurate is our prediction?

Three uncertainties involved:

(1)  $\mathbf{y} = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$  <sup>only estimates</sup> for the true population regression plane  
 $f(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$  <sup>inaccuracy of est.</sup>

- inaccuracy in estimates - related to reducible error
- compute SE for  $\hat{y}$ .

(2) model bias.

(3) Using prediction intervals to find how much  $y$  varies from  $\hat{y}$ .

they are wider than CI

- prediction interval  $\rightarrow$  reducible error: error in estimate of  $f(x)$

$\rightarrow$  Intrinsic error: how much an individual point will differ from the population regression plane.

3.3. Other considerations in the regression model

3.3.1 Qualitative Predictors

(i) predictors with only two levels.

ex: Credit data set:

$\rightarrow$  A factor with  $x_i = \begin{cases} 1, & \text{if } i\text{th person owns a house} \\ 0, & \text{if } i\text{th person does not own a house} \end{cases}$

$$y_i = \beta_0 + \beta_1 x_{i1} + \epsilon_i = \begin{cases} \beta_0 + \beta_1 + \epsilon_i & \text{if } i\text{-th person owns a house} \\ \beta_0 + \epsilon_i & \text{if } i\text{-th person does not own a house.} \end{cases}$$

interpretation: these coefficients are the average credit card balance for non-owners.

$\beta_0$ : avg. credit card balance who do not own a house (0)

$\beta_0 + \beta_1$ : avg. credit card balance among those who own their houses.

$\beta_1$ : avg. diff. in credit card b. balance between owners & non owners.

→ interpretations change accordingly with definition of dummy variable. (coded as -1 or 1, interpretation will change, but values will remain same).

(2) Qualitative predictors with more than two levels.

$x \rightarrow$  qualitative variable  $\xrightarrow{\text{m levels}} \{ \text{m levels} \}$

We consider  $m-1$  dummy variables:

$x_{i1} = \begin{cases} 1 & \text{if } i\text{-th obs. } \in 1^{\text{st}} \text{ level} \\ 0 & \text{else} \end{cases}$

$x_{i2} = \begin{cases} 1 & \text{if } i\text{-th obs. } \in 2^{\text{nd}} \text{ level} \\ 0 & \text{else} \end{cases}$

$x_{im-1} = \begin{cases} 1 & \text{if } i\text{-th obs. } \in (m-1)^{\text{th}} \text{ level} \\ 0 & \text{else} \end{cases}$

The level with no dummy variable is known as baseline.

For  $m=3$ , regression equation becomes  $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i$ ,  $\in I^{\text{st}}$  level

$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i$   $\left\{ \begin{array}{l} \beta_0 + \beta_1 + \epsilon_i; \epsilon_i \in II^{\text{nd}} \text{ level} \\ \beta_0 + \epsilon_i; \epsilon_i \in III^{\text{rd}} \text{ level.} \end{array} \right.$

Interpretation: Estimating coefficient  $\beta_1$  if we want to know the effect of variable  $x_1$  on  $y$ .

### 3.3.2 Extensions of the linear model



- Two restrictive assumptions that are violated often:

(i) additivity:

Association bet<sup>n</sup> ~~and~~  $x_j$  &  $y$  does not depend on the values of other predictors.

(ii) linearity:

functional form of population regression line is linear  
- change in  $y$  associated with a one-unit change in  $x_j$  is constant, regardless of the value of  $x_j$ .

(i) Removing additive assumption:

Considering interaction terms in a relationship between  $y$  and  $x_1, x_2, \dots, x_k$ .

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \epsilon$$

$\downarrow$  interaction term

$$\Rightarrow Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \epsilon$$

$$= \beta_0 + (\beta_1 + \beta_3 x_2) x_1 + \beta_2 x_2 + \epsilon$$

$$= \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

$$\beta_1 = \beta_1 + \beta_3 x_2 \Rightarrow \beta_1 \text{ is function of } x_2.$$

$\Rightarrow$  change in value of  $x_2$  will change the association between  $x_1$  and  $y$ .

How to know which variables are interacting?

using F-test

- The hierarchical principle (from last slide)

If we include an interaction in a model, we should also include the main effects, even if the p-values associated with their coefficients are not significant.

- For nonlinear relationships: polynomial regression.

### 3.3.3 : Potential problems.

(1) non-linearity of the response-predictor relationships

- residual plots:

SLR is similar to MLR

$e_i$

$e_i$

$x_i$

$y_i$

> ideally, residual plot should not show any pattern, if present, indicates severe problem with some aspect of linear model. → take transformations such as  $\log x$ ,  $\sqrt{x}$ ,  $x^2$ .

(2) Correlation of error terms

- If error terms are correlated: CI will be narrower

- funnel shape:  $\sqrt{y}$ .

U-shape:  $\log y$  } (2) heteroscedasticity.

- weighted least squares,  $\rightarrow$  OLS is not good

(3) heteroscedasticity.

(4) Outliers.

- a point for which  $y_i$  is far from the value predicted by the model.

What is a studentized residual & why we need it?

studentized residual =  $\frac{e_i}{\sqrt{SE(e_i)}}$

studentized residual =  $\frac{e_i}{\sqrt{SE(e_i)}}$

(5) High leverage points.

- points with unusual value of  $x_i$

- difficult to identify in multiple linear regression.

$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\frac{1}{n} \leq h_i \leq 1, \text{ average } = \frac{p+1}{n}$$

→ so if a given observation has an leverage statistic that greatly exceeds  $\frac{p+1}{n}$  → corresponding point might have high leverage ↗

### (e) Collinearity:

- two or more variables closely related to one another.
- ~~multiple~~ for multiple estimates  $\rightarrow$  RSS is minimum.
- CI for  $\beta$  widens  $\rightarrow$  wider confidence interval
- t statistics significantly decreases  $\rightarrow$  p-value ↑  
 $\rightarrow$  i.e. fail to reject Null hypothesis ↗
- $\rightarrow$  marks importance of some variable.

\* How do we decide whether CI is 'too wide' or not?

\* VIF: variance inflation factor.  $\rightarrow$  measure?

\*  $VIF = 1 \rightarrow$  complete absence of collinearity.

\*  $[VIF > 5 \text{ or } VIF > 10] \rightarrow$  problematic amount of collinearity.

$$VIF(\hat{\beta}_j) = \frac{1}{1 - R_{x_j|x_{-j}}^2}$$

where  $R_{x_j|x_{-j}}^2$  is  $R^2$  from a regression of  $x_j$  onto all of the other predictors.

If  $R_{x_j|x_{-j}}^2 \approx 1 \Rightarrow$  collinearity present, large VIF.

Solutions:

(1) Drop one of the problematic variable.

(2) Combine the collinear variables together into a single predictor.

\* How to decide if we have to add interaction terms or not?  
 Which ones to add, consider?

### 3.5 Comparison of Linear Regression with K nearest neighbors

- linear regression is example of parametric.
- if specified functional form is far from truth, prediction accuracy will be very less.

Do all non parametric methods don't have a loss function optimization problem inside them? (like KNN).

- non parametric methods : K Nearest Neighbors. (KNN regression)

Algorithm:

- (1)  $x_0$ , prediction point  $\in \text{space}$  (representing  $A$ )
- (2) identify  $k$  training observations that are closest to  $x_0$  (represented by  $N_k$ )

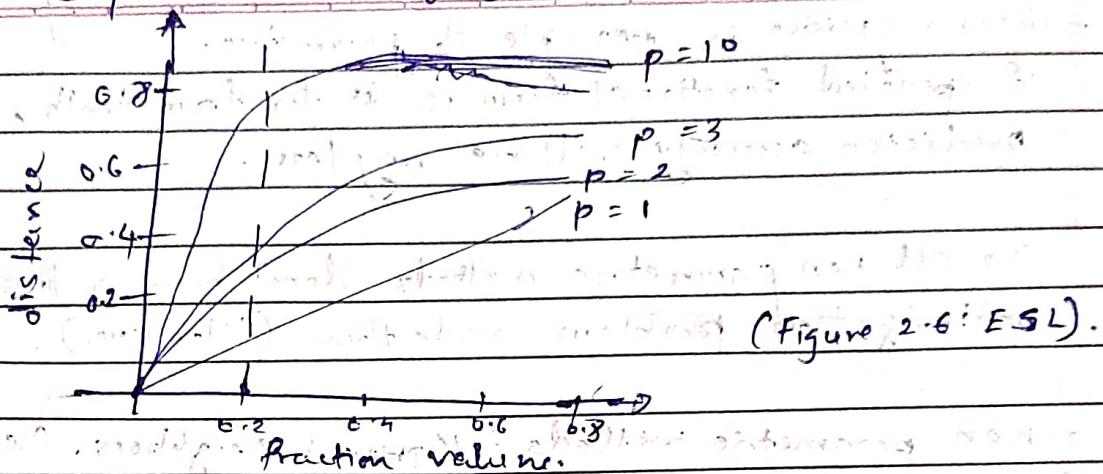
$$(3) \hat{f}(x_0) = \frac{1}{k} \sum_{i \in N_k} y_i$$

- optimum value of  $k$  depends on bias-variance trade off.  
Small  $k$  :- low bias, high variance.  
Large  $k$  :- smoother & less variable.  
 $\hookrightarrow$  smoothing may cause bias by masking some of the structure in  $f(x)$ .
- The parametric approach will outperform the non-parametric approach if the parametric form that has been selected is close to the true form of  $f$ .

How to know for which values  $K$  it is not overfitting?  
(for  $k=1$ , it is definitely over fit).

- performance of KNN decreases as no of features increase.
- Curse of dimensionality.  $k$  observations that are nearest to a given test observation  $x_0$  may be very far away from  $x_0$  in  $p$  dim space when  $p$  is large.  
 $\rightarrow$  poor prediction of  $f(x_0)$  in KNN.

e.g (ESL) In 10 dimensions we need to cover 80% of the range of each coordinate to capture 16% of cluster



- As a general rule : parametric methods will tend to out perform nonparametric approaches when there is small no. of observations per predictor.