

Goal \rightarrow of Statistics is to infer f

$Y \approx f(X) + \epsilon \rightarrow$ error term / bias.

$X = (x_1, x_2, x_3, \dots, x_p)$

\hookrightarrow features

dependent variable.

\rightarrow Linear Models good for inference, but often may not be the best predictors.

\rightarrow On contrast, some highly non-linear approaches will predict extremely well but are hard to interpret

Parametric Way \rightarrow 1) Select a Model

e.g. \rightarrow Linear Regression

2) Train the Model to estimate the parameters.

Non Parametric \rightarrow 1) Do not assume a f^u .

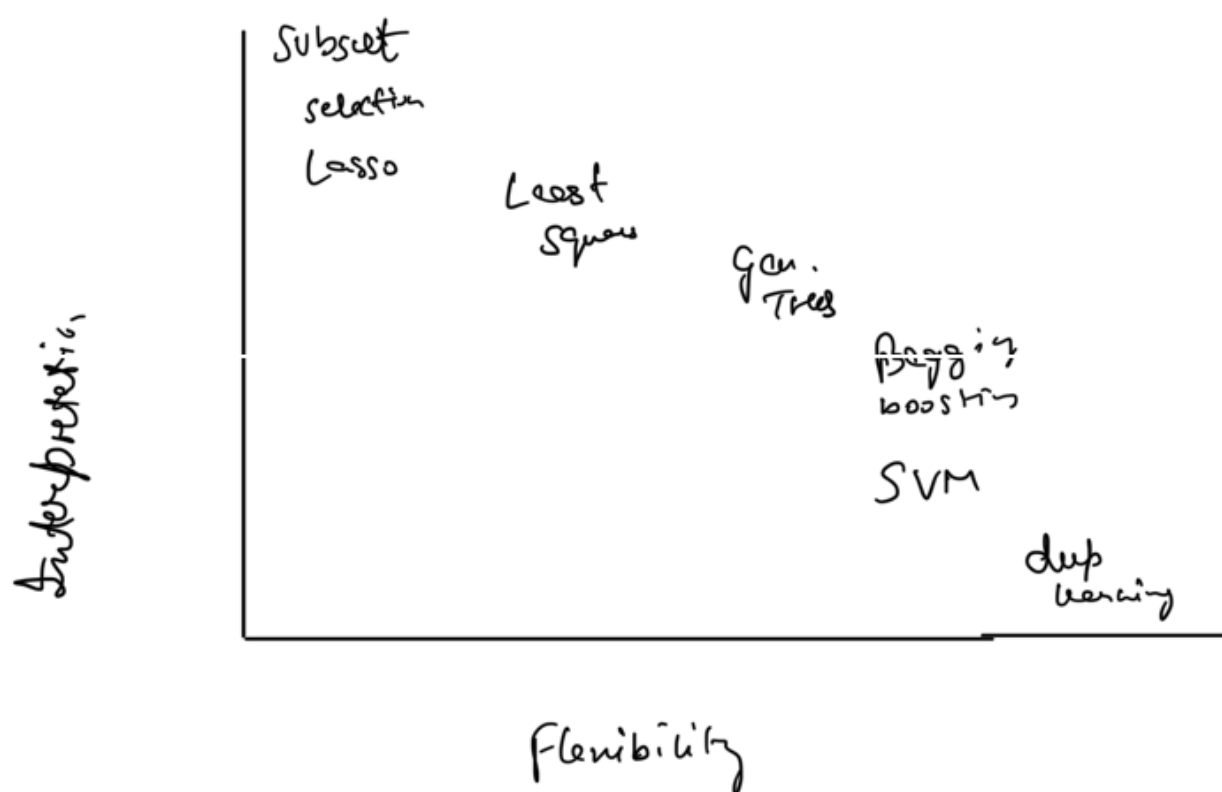
... observations

eg \rightarrow Thin-plate
splines

2) But lot more
are needed to estimate.

but have to be cautious of overfitting.

\Rightarrow are more flexible on the brighter side.



restricted models are good for inference

Unsupervised

Learning

\rightarrow Categorising data without
a response / output variable.

Clustering problems.

Regression v/s classification problems

variables \rightarrow Quantitative
 \rightarrow Qualitative (categorical)

Linear regression suitable for quantitative.

Assessing Model Accuracy \rightarrow

* very important to select which method is suitable for the given problem.

For Regression \rightarrow

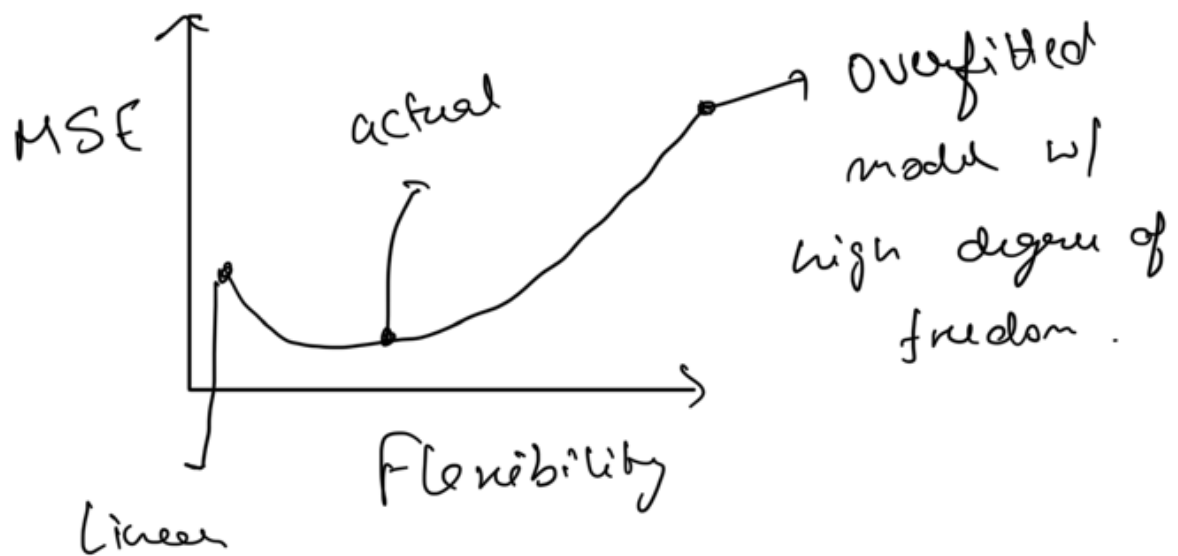
We use MSE (Mean squared error) to quantify the accuracy

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

MSE calculated on \ggg MSE on training data.
..... Test data

Unseen

Some overfitted curves may give very less training MSE, but when tested with unseen data will give larger MSE.



Overfitting → is the case when a linear model would have given less test MSE than high degree of freedom.

The Bias - Variance Trade off

The U-shaped curve of MSE is

the result of this trade off

MSE on test data

$$E(y_0 - \hat{f}(x_0))^2 = \underbrace{\text{Var}(\hat{f}(x_0))}_{\text{variance of } \hat{f}(x_0)} + \underbrace{[\text{Bias}(\hat{f}(x_0))]^2}_{\text{Squared of Bias of } \hat{f}(x_0)} + \underbrace{\text{Var}(\epsilon)}_{\text{variance of the error term}}$$

3 fundamental properties

Our aim to find something that reduces the variance and the bias.

{ both terms $\text{var}(\hat{y})$ & $(\text{Bias}(\hat{y}))^2$ are non -ve so the MSE can never be less than irreducible error }

for a more flexible f , the variance will be high as for different testing

data, the output is highly varied.

Bias \rightarrow error introduced by approximating a real life problem (with high complexity) by a simple model like linear regression.

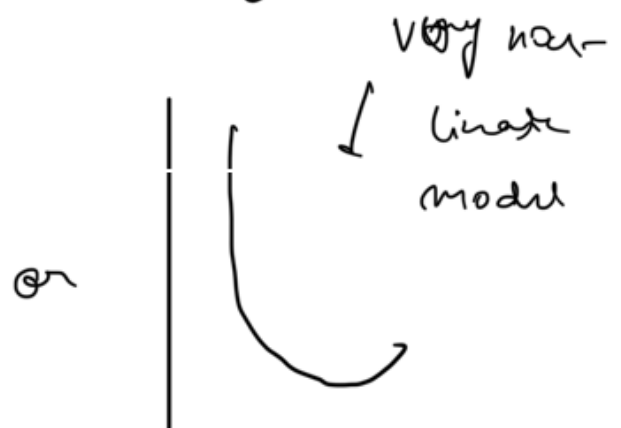
Simple Linear Regression = generally higher Bias.

if the approximated model is closer to the actual underlying model then the Bias is low

More Flexible Model

Variance \uparrow

Bias \downarrow



mean

we try and find the
Minimum of this curve to have
low variance and Bias.

for qualitative data
(categorical variable)

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(y_i \neq \hat{y})$$

training error rate

Bayes Classifier

$$P(Y = j / X = x_0)$$

for a problem w/
2 possible outcomes

$$P(Y = 1 / X = x_0) > 0.5$$

$$\text{Bayes's error rate} = 1 - E(\max_j P(Y = j / X))$$

↳ analogous to the irreducible error.

1

↳ this is basically the overlap.

K - nearest neighbours -

In reality we don't know the conditional relation $Pr(Y=j / X=x_0)$
So Bayes Classifier is the Gold Standard.

given a +ve integer K we see

K nearest neighbours and assign
the class which most of the
neighbours are a part of

$$Pr(Y=j / X=x_0) = \frac{1}{K} \sum_{i \in N_0} I(y_i = j)$$

↳ class is assigned with the
largest probability.

less K = more flexible
{ high variance }
{ low bias }

high k = Less flexible (tends to linearity)

Choosing the correct level of flexibility is very imp.

Chapter - 3 Linear Regression

Statement/
Problem - 1 \rightarrow expenditure on advertisement
v/s sales.

3 modes $\left\{ \begin{array}{l} \text{news} \\ \text{TV} \\ \text{radio} \end{array} \right\}$ set of 3 which
is the best use

Interaction / Synergy effect

Spending 50K on TV and Radio
instead of 100K on radio or TV.

We can solve this by Linear Regression

1. $n \rightarrow$ single predictor variable X

Lin

$$Y = \beta_0 + \beta_1 X$$

$$\text{eg} \rightarrow \text{sales} = \beta_0 + \beta_1 \times \text{TV}$$

$\downarrow \qquad \qquad \downarrow$
intercept slope

Residual $\rightarrow e_i = y_i - \hat{y}_i$

RSS \rightarrow Residual sum of squares

$$\text{RSS} = e_1^2 + e_2^2 + \dots + e_n^2$$

$$\text{RSS} = (y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1)^2 + (y_2 - \hat{\beta}_0 - \hat{\beta}_1 x_2)^2 + \dots + (y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n)^2$$

taking derivatives w.r.t. β_1 & β_0

$$\beta_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

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

$$\beta_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$r = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

In general the sample mean \bar{y} is a good estimator of the population mean μ .

over a large set of samples and taking mean of the mean it will converge to μ . Law of large numbers.

$$\text{var}(\hat{\mu}) = \frac{\sigma^2}{n}$$

as n increases the deviation from true mean decreases

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

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

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

estimate of σ is RSE

Residual standard error

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

Confidence interval

CI of 95% for β_1

$$= \hat{\beta}_1 \pm 2 \cdot SE(\hat{\beta}_1)$$

$$\Rightarrow \beta_1 \in [\hat{\beta}_1 - 2SE(\hat{\beta}_1), \hat{\beta}_1 + 2SE(\hat{\beta}_1)]$$

_____ X _____

Hypothesis Testing

most common: {null hypothesis}

$H_0 \rightarrow$ There is no relationship b/t
X and Y

{alternative hypothesis}

H_a : There is some relationship b/t
X and Y.

$$H_0: \beta_1 = 0$$

$$H_a: \beta_1 \neq a$$

So β_1 must be sufficiently distant
from 0.

how do we define sufficiently distant

We compute the t -statistic

$$t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)}$$

this measures how many S.D.s away
is our result from mean.

Small p value = There is no
association

hence we conclude $\beta_1 \neq 0$
and reject null hypothesis

Assessing the Accuracy of the model

- 1) Residual standard error
- 2) R^2 statistic

and \dots

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

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

R^2 statistic

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

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

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

if it is close to 0 that means that
our prediction model is garbage

R^2 is the measure of linear relationship b/w X & Y

$$\text{Cor}(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{SD(X)} \sqrt{SD(Y)}}$$

Multiple Linear Regression

in the case of Advertising.csv, we need a relation which accounts for all 3 factors.

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

$$q \rightarrow \text{sales} = \beta_0 + \beta_1 \times \text{TV} + \beta_2 \times \text{news} + \beta_3 \text{radio}$$

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

$$= \sum_{i=1}^N \left(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_p x_{ip} \right)^2$$

doing multiple regression here reveals that newspaper advertisement

is garbage for brand

Some important questions.

- 1) Is at least one of the p predictors a good estimator.
- 2) Do all predictors help predict y or only a subset is useful

Null hypothesis

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

H_a : at least one β_j is $\neq 0$

hypothesis test is provided by computing the F -statistic

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

$$TSS = \sum (y_i - \bar{y})^2 ; \quad RSS = \sum (y_i - \hat{y}_i)^2$$

F is close to 1 is H_0 is true

Q) how large does the F -statistic need to be before we reject it.

it follows F -distribution so we can find the p -value

Null

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

{ q of the coefficients are zero }

$$F = \frac{(RSS_0 - RSS)/q}{\left(\frac{RSS}{n-p-1} \right)}$$

variable selection

Selecting a subset of p that actually matters.

there are 2^p subsets and we cannot consider every subset unless p is very small.

Methods

- 1) Forward Selection: begin with a null model
→ then add the variable with lowest RSS
→ then continue till a stopping rule.
- 2) Backward Selection → start with p features then reduce
- 3) Mixed.

Regression with qualitative variables

only have 2 levels

$$x_i = \begin{cases} 1 & \text{if } y=0 \\ 0 & \text{else} \end{cases}$$

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i = \begin{cases} \beta_0 + \beta_1 + \epsilon_i & 1 \\ \beta_0 + \epsilon_i & 0 \end{cases}$$

β_0 = avg. value when $x=0$

β_1 = avg. diff when $x=1$ + $x=0$

for 3 eg East, West, South

$x_1 \rightarrow$ South

$x_2 \rightarrow$ West

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

$$y = \begin{cases} \beta_0 + \beta_1 + \epsilon & \rightarrow \text{South} \\ \beta_0 + \beta_2 + \epsilon & \rightarrow \text{West} \\ \beta_0 + \epsilon & \rightarrow \text{East} \end{cases}$$

Regression assumes linearity and additive nature

{ basically superimposing k_i one thing doesn't affect the output of others }

but there is synergy effect in marketing in stats it is known as interactive effect

for the case of advertisement and capturing the relation b/w TV and radio

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_2 X_1 + \varepsilon$$

$$= \beta_0 + (\beta_1 + \beta_3 X_2) X_1 + \beta_2 X_2 + \varepsilon$$

$$= \beta_0 + \beta_1' X_1 + \beta_2 X_2 + \varepsilon$$

If interaction b/w two variables
is statistically significant

then I should keep both of
the main effects as well.

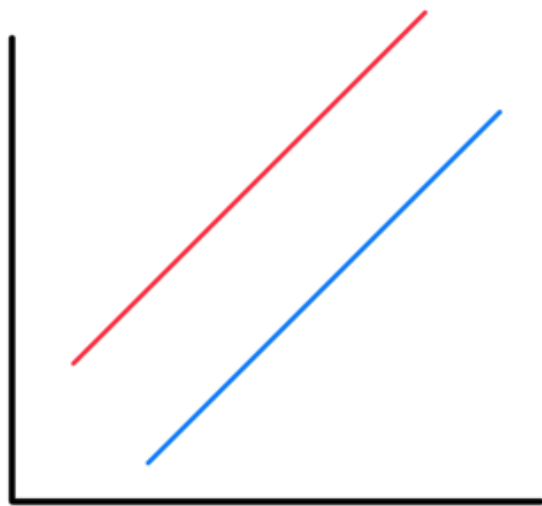
This interaction variable is specially
very useful in the case of
quantitative and qualitative
variables.

Ex \rightarrow In an case of credit
With $X_1 \rightarrow$ Income (quant)
 $X_2 \rightarrow$ Student (qualitative)

$$Y = \beta_0 + \beta_1 \times \text{income} + \begin{cases} \beta_2 + \beta_3 \times \text{income} & 1 \\ 0 & 0 \end{cases}$$

1 0 1 0 0 1

$$= \begin{cases} (\beta_0 + \beta_2) + (\beta_1 + \beta_3) \text{ income} \\ \beta_0 + \beta_1 \times \text{income} \end{cases}$$



W/o interaction



With interaction

Non Linear Relationships

extending linear model to accommodate
non-linear relationships by
using polynomial Regression

eg \rightarrow in Auto. CSV,

$$\text{mpg} = \beta_0 + \beta_1 * \text{hp} + \beta_2 * \text{hp}^2 + \varepsilon_j$$

to capture the quadratic
nature of

mpg v/s hp

To check for correlated error.

We plot residuals as a function
of time and there should be no
discernible pattern.

If error terms are positively correlated
then we may see tracking in the
residuals, i.e. adjacent residuals
may have similar values

3) Non constant Variance of Error terms that $\text{var}(\epsilon_i) \propto x^2$.

Identified by heteroscedasticity or the presence of a funnel shape in residual plot.

Solⁿ \rightarrow transform it to a constant f^* like $\log(y)$ or \sqrt{y}

4) Problem of outlier

Outlier can be seen in the residual plot easily.

\Rightarrow Instead of plotting residuals, we can plot the studentized residuals computed by dividing each residual e_i by its estimated std. error.

$> 3 \rightarrow$ then outlier.

5) High leverage observations

obs that have high residual

& leverage \rightarrow change the predictors a lot.

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

this is a bigger problem to identify in multiple regression

6) Collinearity

If i make two predictors that are collinear then I test my null hypothesis then I

might get an incorrect answer
So I should always plot the
Scatter matrix

A better way to assess multi-collinearity
is Variance Inflation Factor (VIF)

$$VIF \rightarrow \frac{\text{var of } \hat{\beta}_j \text{ w/full model}}{\text{var of } \hat{\beta}_j \text{ fitted on its own}}$$

Smallest value of VIF is 1,

which means there is no collinearity

$VIF > 5$ or 10 is problematic

$$VIF(\hat{\beta}_j) = \frac{1}{1 - R^2_{x_j/x_{-j}}}$$

KNN with small k will have high variance & less bias as it is flexible - whereas with high k it tends to regression which is not very flexible.

large k = smoother fit

In general a non parametric approach will predict better than ~~parametric~~ approach unless the parameters are close to true form.

When $p = 1$ or 2

KNN outperforms LR.

with $p = 3$ mixed results

but $p \geq 4$ $\{p = \text{no. of predictors}\}$

Linear regression outperforms KNN

as dimensions increase they
pose a challenge for KNN

curse of dimensionality

$\{ \text{there is no near neighbour} \}$

eg $\rightarrow n = 50, p = 20$