

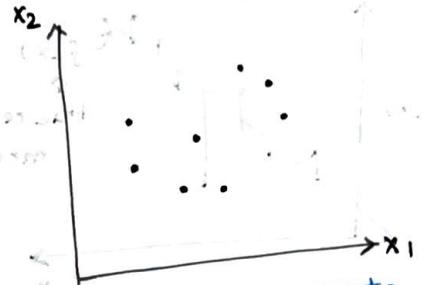
For a customer segmentation problem, Euclidean distance helps us to group and differentiate light buyers from heavy buyers. While the cosine similarity helps to group customer with respect to the similarities in terms of their buying habits and tastes.

Distance matrix: which customers are more similar? Diagonal elements → old customers

Similarity matrix: Correlation matrix (cosine similarity) Diagonal elements are all 1 hence all customers are similar

$x_1 \rightarrow$ no. of units of product A that has been bought

$x_2 \rightarrow$ no. of units of products B that has been purchased! no. of units of products C that has been purchased!



Single linkage: distance between 2 clusters is minimum of the all possible distances between 2 points.

Each dot represents a customer, no. of units

Complete linkage: maximum distance between points

This framework won't allow us to handle new customers

Average: Mean of pairwise distance between points

divide into 2 different clusters

Dendrogram → used to visualize the number of clusters. Observations are along x-axis and distance along y-axis.

Cutting the dendrogram →

Within cluster → homogeneous collection of observations.

Between clusters → heterogeneous collection of clusters

The main objective of cluster analysis is to separate clusters find groups within other data such that similar elements are present

inside a group and very different elements are present in different groups. Finding groups within a data helps us to benefit

economically and increases action effectiveness.

Targeting potential customers: conversion rate is high

Time and economically beneficial.

Customer with maximum conversion rate can be targeted.

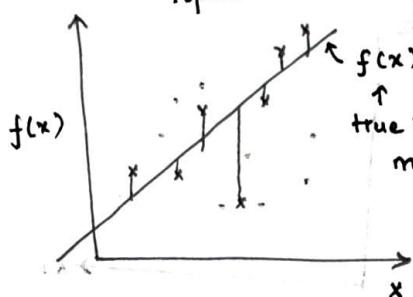
31/10/2022

SOURCES OF ERROR: (Most X: Conceptual topic)

ERRORS

Irrducible Error

Population \rightarrow cannot be minimized.



Reducible Error

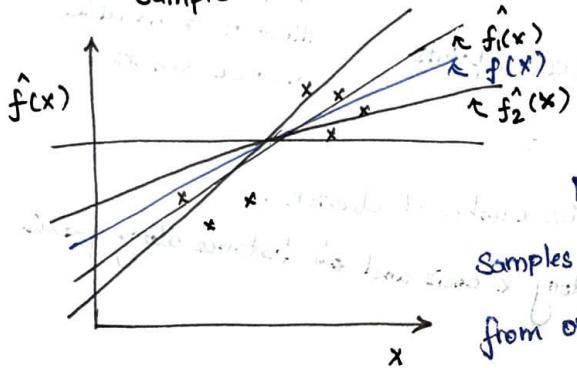
Error

$$\sigma^2 = \text{Var}(\text{residuals}) = \text{irreducible error}$$

\rightarrow obtained by fitting the population data on linear regression model on population data.

\rightarrow focus on errors that are reducible.

Sample:



$$\hat{f}(x) = \text{LR model on the sample } s_1$$

$$\hat{f}_1(x) = \text{LR model fitted on the sample } s_1$$

For different samples, we have different

Samples since y-intercept and slope changes from one to another sample.

\rightarrow The sample mean is a function of the sample we are choosing.

\rightarrow When the sample size is smaller the variability of sample means (the

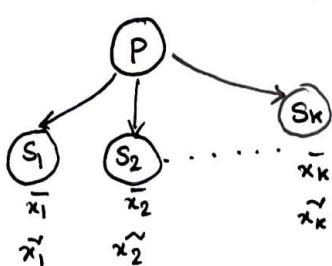
sample mean vary from one sample to another).

① A salesman who quotes 15Rs, 20Rs, 25Rs, 30Rs for a marker.

② Another salesman quotes Rs 18, Rs. 19, Rs. 20, Rs. 21 for a marker.

Actual marker price = Rs. 20.

② is more reliable because $SD_2 < SD_1$.



$$sd(\bar{x}_1, \bar{x}_2, \dots, \bar{x}_k)$$

$$sd(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_k)$$

$$sd(s_{d1}, s_{d2}, \dots, s_{dk})$$

STANDARD ERROR

\uparrow
sd of the sample statistics.

Standard error determines or measures the reliability of the estimate.

~~It is a range~~

on basis of estimation of μ and s^2

95% confidence interval of μ .

Sample mean = 10, the first estimate of μ is 10.

$$SE(\bar{x}) = 2$$

$$10 \pm 1.96 \times 2$$

estimate of μ is between 6 and 14.

$$[6, 14]$$

We are 95% sure that pop. mean will be between 6 and 14.

We are 95% sure that pop. mean cannot be true.
because μ is in $(6, 14)$ and there's 5% chance lying outside this range. We're here we say there

$$\Rightarrow P(\mu \text{ is in } (6, 14)) = 0.95$$

here μ has 5% chance lying outside this range. So there is 5% chance that μ will fall outside

here saying μ is moving. (to 6 or 14) and there is 5% chance that μ is within this range.

here saying μ is moving. (to 6 or 14) and there is 5% chance that μ is within this range.

$$P(\text{interval contains } \mu) = 0.95$$

here the interval can move while μ remains fixed.

most fixed. (the same with prediction) If the true regression line differs a lot from the fitted line, then there's a problem.

If the average of all the regression lines is away from the true model is

How far is the average of all the models is away from the true model?

called the bias. (the difference between the true regression line and the fitted line.)

But in real life, we don't have the true regression line.

Two sources of reducible error:

(1) Bias

(2) Variance

Let $f_i(x)$ represent the true regression model on the i th sample and $f(x)$ be the average of all the fitted models (i.e. the average of all possible fitted models for the input x).

Let $f_i(x)$ be the fitted regression model on the i th sample and $f(x)$ be the average of all the fitted values from all possible fitted models for the input x .

of all the fitted values from all possible fitted models for the input x .

the input x .

The bias of the model at \tilde{x} is defined as:

$$\text{Bias}(\tilde{x}) = f(\tilde{x}) - \bar{f}(\tilde{x})$$

Bias measures how far is the average of all possible fitted models from the true model.

models is away from the true model.

→ If this is large then we got a problem with the predicted values.

Variance → how much are the individually fitted models varying across the average of the fitted models.

Doesn't deal with true model at all.

The variance at \tilde{x} is defined as:

$$\text{Var}(\tilde{x}) = E[(\hat{f}_i(\tilde{x}) - \bar{f}(\tilde{x}))^2]$$

→ Variance measures how much does the outputs of each models

(fitted over all possible samples of the same size) varies from the average of the models.

The mean squared error at \tilde{x} is:

$$\text{MSE}(\tilde{x}) = \text{Irreducible}(x) + \text{bias}^2(\tilde{x}) + \text{variance}(\tilde{x})$$

Bias is dependent on one particular value of x and is not the same for all x .

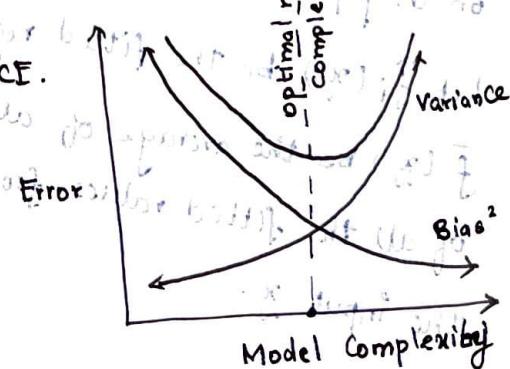
The model should not look different for different samples of same size.

Bias for decision tree, model complexity is independent on the depth of the decision tree.

MODEL COMPLEXITY - BIAS - VARIANCE.

Feature selection is quite hard to determine the point @ which the

Total Error



model learns significantly less → underfitting
Model learns sufficiently large → overfitting.
High variance will make the model overfit the data. While high bias makes the model underfit the data. Both are bad.

A KNN model with $K=2$, $K=20 \rightarrow$ smoother model.
↓ agitated model.

Book: ISLR.

Lesser the value of $K \rightarrow$ overfitting
Higher the value of $K \rightarrow$ underfitting
Always work with samples of larger size.

Model complexity?
for KNN.

Exp 1: Let $P =$ cars data with 100 observations and the variables MPG, weight, horsepower and displacement. Generate 100 samples of size 50 each from P . Over each sample, fit a decision tree model to predict MPG of depth = 2 and a decision tree model of depth = 8. For a fixed input, get the predicted values for both the kinds of models. Compute the variance of the predicted values separately for ① DT (depth = 2) ② DT (depth = 8). Report your observations and comment.

Exp 2: ✓ if set to large value, model complexity is reduced
The minimum no. of observations at leaf node controls the complexity of a decision tree. For each sample fit a decision tree model with minimum samples - leaf = 5 vs. minimum samples - leaf = 20. For a fixed input compute the variance of the predicted values separately. Compute the variance of a 2NN model. (if time permits)

Exp 3: Using similar experiment as above, compare the variance of a 2NN model and a 10NN model. Print out your CV's.

Books

Lambs and Lambs

ISLR

- Take notes on AUC-ROC curve given to senior batch
- If necessary talk with Anchal Khanna.
- Notes on data leakage.
- Final Assignment on Graphs, WhatsApp Assignment on House Price prediction
- Decision Tree Assignment on House Price prediction

VIF → doesn't measure multicollinearity.

Does VIF measure multicollinearity? ← AIM magazine.

Multicollinearity → Multi-correlation if high, we drop the column.

VIF:

Predictor	Target	Predictor
① x_1	$\sim \beta_0 + \beta_2 x_2 + \beta_3 x_3$	$R_1^2 = 0.90 \rightarrow 90\% \text{ of var of } x_1 \text{ is explained by } x_2 \text{ and } x_3$
② x_2	$\sim \beta_0 + \beta_1 x_1 + \beta_3 x_3$	$R_2^2 = \underline{\quad}$
③ x_3	$\sim \beta_0 + \beta_1 x_1 + \beta_2 x_2$	$R_3^2 = \underline{\quad}$

Find individual R^2 assuming one variable as target and rest as predictors.

$VIF = 1$, the variable is uncorrelated with all other variables.

$VIF = \frac{1}{1-R^2} \Rightarrow$

- $VIF_1 = \frac{1}{1-0.9} = 10 \quad 0 < R^2 < 1 \rightarrow$
- $VIF_2 = \frac{1}{1-0.8} = 5 \quad 1 < VIF < \infty$
- $VIF_3 = \frac{1}{1-0.7} = 3.33 \quad \text{furthermore, it is not a good benchmark}$

The particular column (but not the first) drop the particular column at once.

should be set).

Using VIF to reduce multicollinearity:

Step 1: Fit a model with all the predictors.

Step 2: Calculate the VIF of all the predictors in the model.

Step 3: Let θ_1 be the VIF provided it exceeds the max. VIF.

Note: Only drop the highest VIF predictor, the VIF of all predictors will significantly fall.

Step 4: Refit the model by dropping the predictor (with all the other predictors).

Step 5: Repeat step 2-4 until all the VIF's of the selected variables are less than 1.

Step 3: Divide into train-test-validation-split

21 (x) (or) use cross validation

Step 4: Run EDA tools (create a package of your own): boxplot, histogram

Step 4: Run EDA tools to create a package boxplot, histogram

Create a package that will help you generate numeric variables.

Create a package for bivariate (numerical predictors target and for univariate) as well as for categorical target.

categorical predictors → side by side box plots. for univariate, for bivariate as well as for multivariate analysis

→ A predictor that's highly collinear with target should be

Kept: $\frac{2}{3} = \frac{1}{3}$ + 70% (1) \rightarrow it is shrinkage? 11%

→ Transformation like \log (to increase linearity, shrinkage), $\frac{1}{x}$ (to study reciprocal relationship)

Square root transformation (ii) a longitudinal view

starts to look like Undecidability w/ goals ↑ & input log ↑ + new vars

We take log on both variables.

Step 5: Fit a linear regression model.

baseline model.
based on all the x_i (features) \leftarrow baseline model.

A model which performs worse than random
because of how we're manipulating

then R^2 can be negative, unless otherwise we require $R^2 \geq 0$

then we can now fit the intercept or failed to fit the intercept and using its data model

OLS regression $\rightarrow R^2$ can never be negative! (assuming no data points outside the regression line)

~~has fitted on training data~~

R^2 → how well model has fitted on training data.

R^2 → how well model fits data → measure of goodness of fit, only to access fit of model.

3. It must not be measured on test data.

R^2 → should not be measured on test data.
MAPE → same as with MAE, RMSE and MAE not R^2 on test data.

To measure error metric go with MAE, RMSE and MAE not R².
To not drop variables,

To measure error metric go to
high overfitting occurs. If you drop var

→ If R^2 is high, overfitting occurs. If you →
overfitting reduced, validation performance will increase

Adjusted R^2 :

No matter what you add a feature, R^2 will increase. So we add some penalty, (i.e.) penalized $R^2 \approx$ adjusted R^2 .

→ Adding relevant predictors will lead to increase R^2 . Adding features may or may not increase the adjusted R^2 .

→ Adjusted R^2 can be used for feature selection.

Check for multicollinearity → to drop columns (VIF)

adjusted R^2 vs. R^2

Adding features keeps R^2 constant

Con-adjusted R^2 can be used to measure error on test data?

Engineering Features:

→ Generating new features from existing features. To improve model accuracy and decrease in linear state.

Forward Selection: which features can finally go into our model.

↳ Use adjusted R^2 to figure out which variable can go into the data.

① $y \approx x_1$, Adjusted R^2 vs. R^2 based on the first predictor.

↳ $y \approx x_1$, Adjusted R^2 vs. R^2 based on the first predictor that goes into model.

↳ $y \approx x_2$, Adjusted R^2 vs. R^2 based on the second predictor that goes into model.

② $y \approx x_3 + x_1$, Adjusted R^2 vs. R^2 based on the third predictor that goes into model.

↳ $y \approx x_3 + x_2$, Adjusted R^2 vs. R^2 based on the fourth predictor that goes into model.

↳ $y \approx x_3 + x_4$, Adjusted R^2 vs. R^2 based on the fifth predictor that goes into model.

③ $y \approx x_3 + x_4 + x_1$, Adjusted R^2 vs. R^2 based on the sixth predictor that goes into model.

↳ $y \approx x_3 + x_4 + x_2$, Adjusted R^2 vs. R^2 based on the seventh predictor that goes into model.

Then we stop. We consider only x_3 and x_4 in my final model.

Why have you selected this model? Objective is reducing the error for features.

Whichever model gives best results.

→ For every instance of test data, KNN goes through each and every observation on test data. For e.g. test data (100) training data (100000).

100 * 100000. So KNN is lazy learner while then for loop runs for k nearest neighbors.

Decision tree is greedy learner.

→ Interpretation → Linear regression vs. DT (but bushy tree is not easily interpretable).

- Even before fitting model, in linear regression we assume that there exists a linear relation between the target and predictor.
- DecisionTree Regressor automatically figures out things and splitting (like decision tree).

If we're really clear about x and y , it's parametric model like linear regression performs better than non-parametric model.

KNN vs Linear regression → ISLR

non-parametric (parametric)

non-parametric (decision tree)

When you will use linear regression or decision tree?

How will you choose between LR or decision tree?

In higher dimensional space, the plane is a flat surface

and looks like straight line in n-d plane whereas in 2D it can be a curve.

Every time we try to fit the model we'll be making an assumption,

$$y = a + bx + \epsilon$$

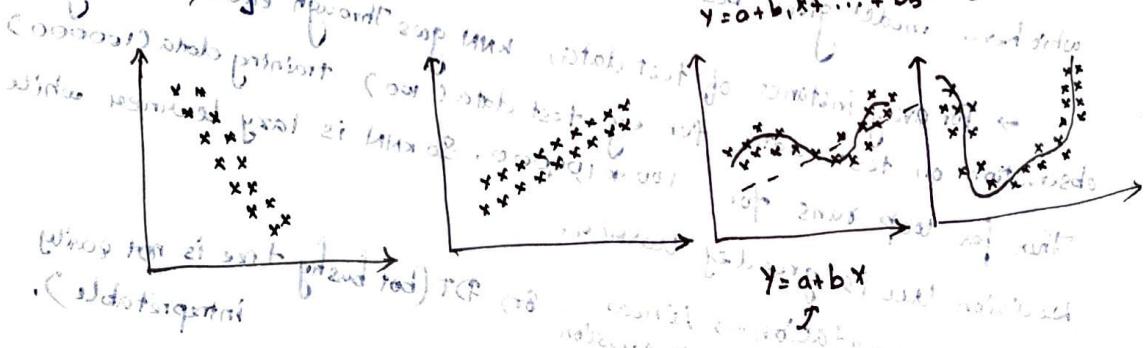
We'll use assumption to fit linear regression

The 1st assumption we're making is y is linear

relation with features.

Assumption 1 cannot be checked.

Assumption 2 is "linearly related with coefficients"



- The assumption you made is wrong, linear regression doesn't fail.
- Use linear regression only after EDA after you get clear picture between target and feature.

Decision tree: splitting happens on feature space based on homogeneity of target variable.

In LR, we need to prove hypothesis which explains the target with predictor in best possible way.

The target with predictor instead of we searching patterns.

We make decision tree to find patterns instead of manual searching.

Why not use decision tree all the time?

A LR model is non-parametric model and it's better than a linear regression on rel. b/w x and y well, a linear regression

high predictive power. If I know x and y well, a linear regression outperforms the decision tree.

Parametric model outperforms non-parametric model.

To achieve linearity, we can use logarithmic transformation, higher order transformation, mode with parameters.

DT → no assumptions between parameters.
 RMSE of DT → RMSE of LR
 Never compare the y with RMSE of DT
 log y, we must take antilog to find right RMSE.

$$y = a + bx + \epsilon$$

RMSE (LR) = 1000
 RMSE (DT) = 4.5

$$\log(y) = a + b \log(x)$$

homoscedacity

Always suspect the too good things.
 Increasing the power of the polynomial; you're adding more and more features and leads to overfitting.

Regularisation? What is it?
 All categorical variables are considered one (all under one dummy).
 All dummy variables should be looked as same variable (should not be looked as different columns, values should be looked as vectors).

$$MPG = b_0 + b_1 O_{UX} + b_2 O_{US} + \epsilon$$

MPG	Origin	O _{US}	O _{UK}
US	(1, 0)	0	1
UK	(0, 1)	1	0
JAP	(0, 0)	0	0

When origin changes from Japan to UK, the MPG ↑ by b_1 units.

first & residual noise of random variable

$$\textcircled{1} \quad y_{ij} = b_0 + b_i x_{ij} + \epsilon_{ij}, \quad i = 1(1)n$$

y is linearly related with coefficient b_i and error ϵ_{ij} .
 errors should not be auto-correlated.
 mean

component iid $\epsilon_{ij} \sim N(\mu, \sigma^2)$, $\epsilon_{ij} \sim I(1)n$

ϵ_{ij} is independent \rightarrow covariance between pairs $\text{cov}(\epsilon_{ij}, \epsilon_{kj}) = 0, \forall i, j$
 ϵ_{ij} is identically distributed \rightarrow Normal distribution with same mean and variance.

ϵ_{ij} is distributed \rightarrow Errors are normally distributed \rightarrow error variance is homoscedasticity

most of professors say that the covariance between them is zero.

For any pairs of errors, the covariance will not be zero \rightarrow no matter what error I bring, values won't be identical \rightarrow because all errors follow normal distribution

able to distinguish between mean and variance
 able to distinguish between distribution with same mean and variance
 able to distinguish between normal distribution with same mean and variance

follow normal distribution with same mean and variance
 errors are constant \rightarrow Homoscedasticity \rightarrow Variance of errors are not constant \rightarrow Heteroscedasticity \rightarrow error variance are not constant

error variance are not constant \rightarrow error variance are not constant \rightarrow error variance are not constant

$\hat{\beta} = (X^T X)^{-1} X^T y$

$\hat{\beta}$ can be represented as

\rightarrow If data has one column which can be represented as a linear combination of other 2, then we have a dependent column

Considering as matrix $\begin{bmatrix} 2 & 6 & 0 \\ 5 & 15 & 0 \\ 8 & 24 & 0 \end{bmatrix} \leftarrow \text{singular matrix} \quad \Delta = 0$

so we can't use $(X^T X)^{-1}$ if we use gradient descent at above point

$\hat{\beta} = (X^T X)^{-1} X^T y \leftarrow$ inverse is not possible since the feature matrix is not invertible.

\rightarrow So we need to remove strong multicollinearity so that we get pretty good results.

Complete case analysis: A case is in observation or record.

Handling missing values:

A case doesn't have any missing values in it.

Complete Case Analysis: Restrict the analysis to only those cases which are complete. Delete all cases which contain atleast one missing values in it. Only if it legitimate.

Caution: This may lead to a massive loss of information.

If we're randomly sample, the sample we pick is not random if we're removing observations resulting in more bias.

We're removing observations.

Never drop observations.

For sequential data, we need to look into the frequency like time series, we can use forward and backward fill.

Check for missing values with some arbitrary

+ Do not replace missing values with some arbitrary values.

In case of timeseries, the timestamps should be of same interval.

If distribution is very skewed we go with median imputation.

If distribution makes lot of sense because it's the expected value of the random variable.

Mean imputation makes lot of sense if the variable has a lot of missing values, replacing missing values with mean.

If there's a lot of missing values, replacing missing values with mean.

Missing data handling
if $b^l = b^u$ then $b^l - b^u = 0$

② Zombie Analysis

if $b^l < b^u$

Decision variables: $b^l < b^u$ or $b^l = b^u$

Two categories possible: $b^l < b^u$ or $b^l = b^u$

$b^l = b^u \Rightarrow b^l = 100 \Rightarrow b^l = 100 \text{ days}$

$b^l < b^u \Rightarrow b^l = 100 \Rightarrow b^l = 100 \text{ days}$

$b^l < b^u \Rightarrow b^l = 100 \Rightarrow b^l = 100 \text{ days}$

③ $b^l < b^u \Rightarrow b^l = 100 \Rightarrow b^l = 100 \text{ days}$

Two categories possible:

Count of $b^l < b^u$ in dataset

Count of $b^l = b^u$ in dataset

Count of $b^l > b^u$ in dataset

Count of $b^l < b^u$ in dataset