

card closing \rightarrow should we try to retain the customer.

~~Challenges~~

Credit card \rightarrow rewards, lounges, interest free credit

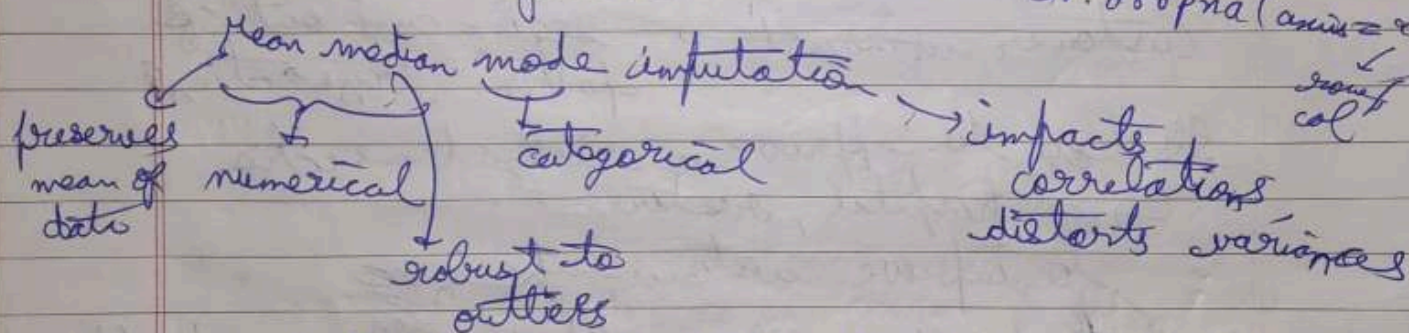
\downarrow
merchants dislike it due to \neq transaction fees

\downarrow
user wants so merchants end up in non-negotiable position

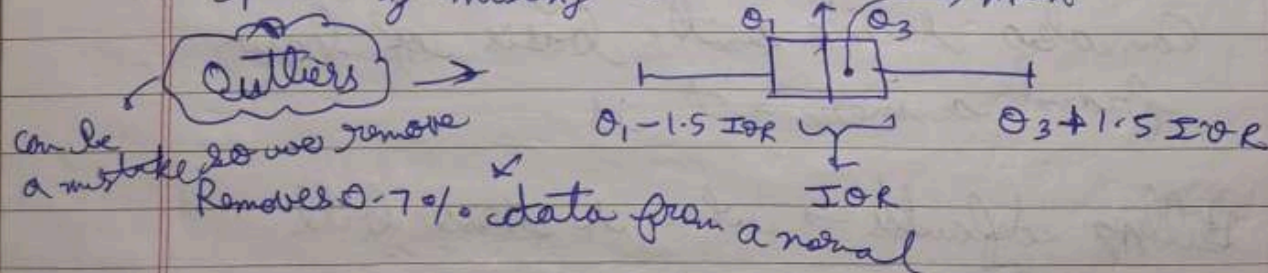
Machine Learning Revision

Missing values \rightarrow MCAR, not MCAR, Removing rows/cols
fillna()

df.dropna(axis=0,
subset=col)



predicting missing values median \rightarrow mean



$$Z\text{score} = \frac{x - \mu}{\sigma} \rightarrow > 3?$$

\downarrow
How many std dev away from mean we are?

Feature Scaling → suppresses outliers
faster optimization

Standard Scaling $x' = \frac{x - \mu}{\sigma}$ → makes $\mu = 0$
 $\sigma = 1$

Minmax scaling $x' = \frac{x - \min(x)}{\max(x) - \min(x)}$ → good for
question data
↓
makes data $[0, 1]$

Robust Scaling → $x' = \frac{x - Q_1(x)}{Q_3(x) - Q_1(x)}$

Good for removing outliers

encoding categorical features

OHE → One hot encoding

Ordinal Encoding

Count, frequency encoding

Target guided encoding

Num to Cat

Binning

Quantization

Prediction Multi Class

OvA → One vs All

OvO → One vs One

Cross validation →

train - validate - test split

K Fold cross validation

Stratified K fold cross validation

Leave one out cross validation

68 - 95 - 99 - 7

|||||

Imbalanced Dataset →

Underamplify

oversampling

K Means, near Miss,
Deletion

SMOTE, Adasyn,
Duplication

~~Multicollinearity, Dummy variable, Bias-Variance, Linear Transform, Discriminability, K-M, DBSCAN, correlation~~

SMOTE - Synthetic Minority Over-sampling 2002

"SMOTE works by selecting examples close in feature space (KNN kinda), drawing line b/w them & drawing new sample along that line"

ADASYN - Adaptive Synthetic Sampling

↳ generate more samples where the density of minority class is less

↳ ignores majority class hence may have issues at border

Borderline SMOTE - Border of classes amplified

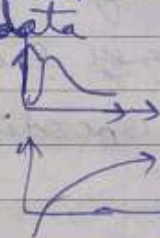
~~Linear Transform~~
Lower Transform

log transform, Box-Cox, yeo-johnson

to make data gaussian

why?

log(x)



$$\frac{x^2 - 1}{\lambda}$$

done in algorithm data gaussian error gaussian, makes data smoother hence decision boundaries easier to make

Correlations

Pearson, Spearman, Point-Biserial, Phi coefficient

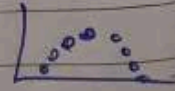
Pearson $\rightarrow \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$

Linear Correlations

Sensitive to outliers

Correlation does not imply causation

Dependency \rightarrow Pearson correlation



Spearman correlation \rightarrow captures monotonicity well
 ↓
 doesn't work for (discrete variables)

$$r = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

 ↓
 d_i = difference in ranks of variables

Multicollinearity

↓
 reduces explainability
 VIF (Variable Inflation Factor)

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

$$1 - R^2$$

↓
 fix by dropping correlated columns

$$VIF \in [1, \infty)$$

1 \rightarrow not correlated

1-5 \rightarrow moderately correlated

VIF > 5 \rightarrow highly

$R^2 \rightarrow$ how much variance of y explained by x

Dummy Variable trap
 Happens in OLS

GLM (Generalized Linear Model)

Linear Regression $\rightarrow y_i \sim N(\mu_i, \epsilon)$; $\mu_i = b_0 + b_1 x_i$

Poisson Regression $\rightarrow y_i \sim \text{Poisson}(\lambda_i)$; $\ln \lambda_i = b_0 + b_1 x_i$

↓
 for discrete y
 count data

Distribution link fn
 for

Logistic Regression $\rightarrow y_i \sim \text{Bern}(p_i)$

$$p_i = \frac{1}{1 + e^{-(b_0 + b_1 x_i)}}$$

$$\text{Poisson} \rightarrow P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

$$\ln\left(\frac{p_i}{1-p_i}\right) = b_0 + b_1 x_i$$

Bernoulli

$$C_x p^x (1-p)^{1-x}$$

$$C_x p^x (1-p)^{1-x}$$

p for $x=1$, $1-p$ for $x=0$ x = outcome

n = no. of predictors
 $\logit(p_i)$
 \rightarrow logistic link fn

Linear Regression \rightarrow Independent observations, Normality of errors, Linearity, Homoscedasticity, Non-multicollinearity

Logistic Regression loss fn = log loss, - its cross entropy loss

$$J(\theta) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))$$

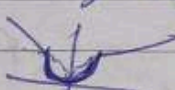
$h_{\theta}(x)$ = predicted value = Sigmoid / Logistic fn

$$\text{Ridge Regression } \rightarrow \lambda \sum |a_i|^2 \quad P(Y=1|x, \theta)$$

Lasso Regression (L1) $\rightarrow \lambda |a_i|$ \rightarrow can eliminate variables

Elastic Net Regression $\rightarrow \lambda_1 |a_i| + \lambda_2 |a_i|^2$

OLS Regression \rightarrow Ordinary Least Squares

Huber Loss \rightarrow  \rightarrow Huber regression anti outliers

Multinomial Logistic Regression \rightarrow Softmax fn

$$P(\text{Class } i | X) = \frac{e^{w_i \cdot x + b_i}}{\sum_j e^{w_j \cdot x + b_j}}$$

Naive Bayes \rightarrow Bernoulli NB (predictors are boolean)

Gaussian NB (predictors are continuous, gaussian)

Multinomial NB (predictors discrete, eg)

Requires independent predictors!

Pseudocount / Laplace Smoothing

$$P(E_i | A) \propto \underbrace{P(A | E_i)}_{\text{probability of observing such input given output } E_i} \underbrace{P(E_i)}_{\text{prior probability}} = \frac{P(A | E_i) P(E_i)}{P(A)}$$

probability of observing such input given output E_i

$$P(y | x_1, x_2, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i | y)$$

$$\text{Precision} = \frac{TP}{TP + FP} \quad \begin{matrix} \text{P} \\ \text{Recall} / \\ \text{Sensitivity} / \\ \text{True + ve rate} \end{matrix} \Rightarrow \frac{TP}{TP + FN}$$

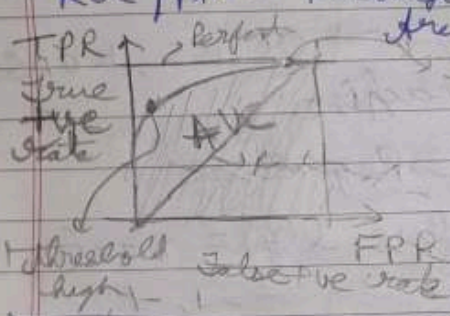
$$F1 \text{ score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

Threshold ↑
Precision ↑
Recall ↓

$$\text{Specificity} \quad \begin{matrix} \text{TN} \\ \text{True - ve rate} \end{matrix} = \frac{TN}{TN + FP}$$

$$\text{False - ve Rate} = \frac{FP}{FP + TN} = \text{Just opposite of Recall}$$

ROC / AUC Receiver Operating Characteristics



Area under curve

Hughes phenomenon

→ even for fixed no of data pts

→ curve of dimensionality

No of dimensions

Clustering → dist →

K-Means

Hierarchical Clustering → Dendrograms

Elbow method

Silhouette score = $\frac{b-a}{\max(a, b)}$

$b = \text{avg inter cluster dist}$

$a = \text{avg intra-cluster dist}$

0 → insignificant
1 → perfect

K-mode clustering → categorical data

K-medoids → helps against outliers

K-prototype clustering → (num + cat) data

K-means ++ → Better initializations

Density Based Clustering

DBSCAN → Density Based Spatial Clustering

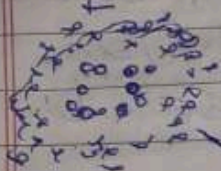
Application with Noise

~~DBSCAN~~
 Hierarchical clustering, SVC, SVM, Hyperparameters & Curbs of DBSCAN

DBSCAN → 2 parameters

→ epsilon → min dist to consider 2 pts as neighbors

minpoints → no. of points needed within epsilon dist



→ no need to specify no. of clusters

Disadv → some data pts may remain unclustered

parameter sensitive
 Cannot handle variable densities

Other density clustering algs → OPTICS, BOPSCAN

SVM → Support Vector machine

SVC → classifier SVR → Regressor

SVM loss fn → Hinge loss $L(y) = \max(0, 1 - ty)$
 $t = 1, -1$

Soft margin SVM

$$\min_w \frac{1}{2} \sum_{i=1}^n w_i^2 + C \sum_{i=1}^n \max(0, 1 - t_i y_i)$$

C ↓ → regularization ↑

Hard margin SVM → no regularization

SVM → large margin classifier

Support vectors → data pts closest to decision boundary
 dist b/w hyperplane & support vectors → margin

$$SVM_{eqn} \rightarrow F(x) = \text{sign}(w^T x + b)$$

Support vector regression (SVR) → minimizes pts b/w data boundary

kernels → Gaussian, polynomial, RBF (Radial Basis fn)

GMM → Gaussian Mixture Model

Soft clustering, Hard clustering

gaussian
 →

used to
 - into
 to an
 covered
 for

GMM \rightarrow EM (Expectation Maximization algorithm)

$$p(x) = \sum_{i=1}^K \phi_i \cdot \mathcal{N}(x | \mu_i, \sigma_i^2)$$

E step \rightarrow

M step \rightarrow

mean \rightarrow mean of data pts

weight \rightarrow avg of probs

variance \rightarrow weighted variance probability; i^{th} data pt in i^{th} gaussian

weight of each data pt

K-means loss fun \rightarrow SSE sum sq errors

GMM \rightarrow maximises log-likelihood

$$\log P(x|D) = \sum_{i=1}^n \log \left(\sum_{j=1}^K \pi_j \mathcal{N}(x_i | \mu_j, \Sigma_j) \right)$$

PCA (Principal Component Analysis) \rightarrow eigenvalues, eigenvectors, score plot, SVD

Dimensionality redn \rightarrow PCA, LDA, t-SNE (Singular Value Decomposition)

LDA \rightarrow Linear Discriminant Analysis

Fisher LDA \rightarrow assumes gaussian, linearly separable

$$J(w) = \frac{\| \mu_1 - \mu_2 \|^2}{S_1^2 + S_2^2}$$

PCA done using SVD $\rightarrow X = U S V^T$

t-SNE \rightarrow stochastic Neighbor Embeddings

SNE \rightarrow 2002

tSNE \rightarrow 2008

plexity \rightarrow effective no of neighbours for which we calculate similarity

$$p(j|i) = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k=1}^n (1 + \|y_i - y_k\|^2)^{-1}}$$

low dim similarity

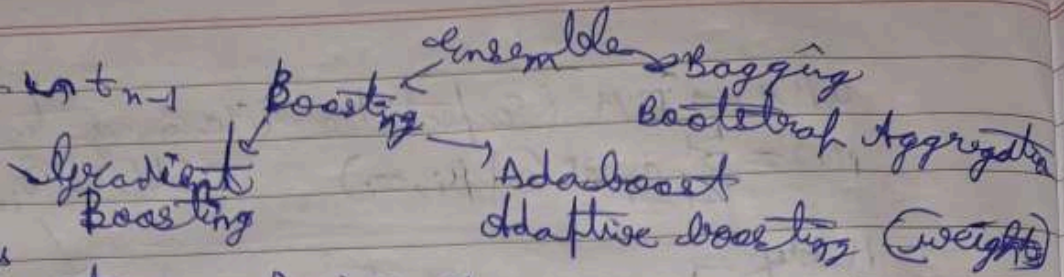
$$C_2 \sum_i KL(P_i || Q_i) = \sum_i P_i \log \frac{P_i}{Q_i}$$

KL Divergence

optimize this

Non-deterministic algorithm Preserves neighbours / local relations

$$T \rightarrow \frac{\bar{x} - \mu}{s/\sqrt{n}} \sim t_{n-1}$$



Classifiers
↓
Discriminative
Generative

Decision Tree
CART (Classification & Regression Trees)

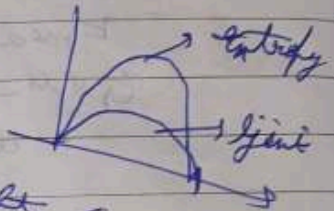
Entropy = $-\sum_{i=1}^n p_i \log_2(p_i)$

Information gain → categorical
Gini index → continuous

Cost fn $J(K, t_K) = \frac{m_{\text{left}}}{m}$

single feature threshold

$$1 - \sum_{i=1}^n p_i^2$$



$G_{\text{left}} + \frac{m_{\text{right}}}{m}$

MSE for continuous

impurity

Boosting
Majority voting

Bias Variance Tradeoff

