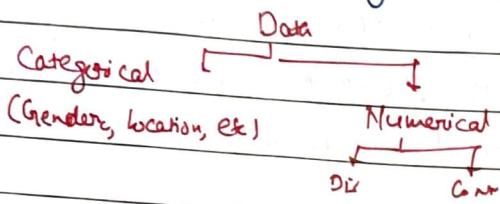


STATISTICAL METHODS FOR DECISIONMAKING.

Descriptive: - Concerned with data summarization, Graphs/Charts and tables.

Inferential: - Talk about a population parameter from a sample, point estimation and hypothesis testing.

Measure of central tendency:

1. AM:  $\bar{X} = \frac{\sum X}{n}$ , affected by extreme values, outliers

2. Median:  $(n+1)/2$ , great resistance to outliers (Can combine medians or mode of several groups)

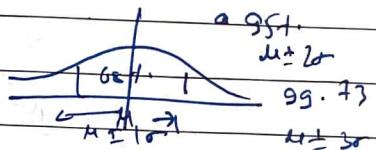
3. Mode: (Value which occurs most often).

Measure of dispersion: How large the spread is around central tendency.

1. Range:  $X_{\text{max}} - X_{\text{min}}$

2. IQR: - (less effected by outliers):  $Q_3 - Q_1$

3. SD: -  $\sqrt{\text{Var}}$



→ CV: (Coef of variation) :  $CV = \frac{SD}{\bar{X}} \times 100$

Five number summary: -  $X_{\text{smallest}}$ , First Quartile ( $Q_1$ ), Median ( $Q_2$ ),

$Q_3 \rightarrow X_{\text{largest}}$

Left skewed:

Median -  $X_{\text{smallest}}$

Symmetric

Median -  $X_{\text{smallest}}$

Right skewed:

Median -  $X_{\text{smallest}}$

$X_{\text{largest}} - \text{Median}$

$X_{\text{largest}} - \text{Median}$

$X_{\text{largest}} - \text{Median}$

$Q_1 - \text{Med} - Q_3$

Median -  $Q_1$

Median -  $Q_3$

$Q_3 - \text{Med}$

$Q_3 - \text{Med}$

$Q_3 - \text{Med}$

25.1. of date

25.1. of date

25.1. of date

$X_{\text{smallest}}$

$Q_1$

Med.

$Q_3$

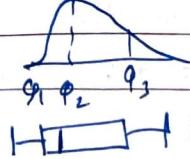
$X_{\text{largest}}$

Julian 19

1.5 times  
IQR



→ Left skewed



Right skewed

→ Covariance: Measures the strength of linear relationship between two numerical variables.

blw few num Var(X and Y).

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

(Doesn't measure the relative strength)

$$\rightarrow r = \text{Coef of correlation: } r = \frac{\text{Cov}(X, Y)}{S_x S_y} \quad -1 \leq r \leq 1$$

Gives relative strength of linear relationship between two numerical variables.

## PYTHON API'S

$$\rightarrow \text{Linear Reg: } \hat{y} = h_{\theta}(x) = \theta^T x$$

$\theta$  contains bias term and feature weights  $\theta_0, \theta_1, \dots, \theta_n$ .

$x$  = instance feature vector:  $x \in \mathbb{R}^n$  where  $x_0 = 1$  is equal to 1.

$$\rightarrow \text{Cost function: } \text{MSE}(\theta^T X, y) = \frac{1}{m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2$$

$$\hookrightarrow \text{Soln: } \hat{\theta} = (X^T X)^{-1} X^T y$$

$y \in (y^{(1)} \text{ to } y^{(n)})$ ,  $\hat{\theta}$  = value of  $\theta$  that min<sup>n</sup> cost function.

→ scikit-learn Linear Reg.

→ scipy.linalg.lstsq() [least sq].

↳ calculates:  $\hat{\theta} = X^T y$ ,  $X^T$  is pseudo inv. of  $X$

→ (The Moore-Penrose inv.).

→ np.linalg.pinv(X-b).dot(y) [directly calculate pseudoinv]  
pseudo inv is basically SVD calculation:  $U \Sigma V^T$

$X^T = V^T \Sigma^{-1} U^T$ ,  $\Sigma^{-1}$  = sets all other values

to zero than threshold and non-zero by their inv. and tr.

\*  $(X^T X)$  has an issue it may or may not work due to invertibility such that  $m < n$  but pseudoinv will always work.

SVD ( $O(n^2)$ ):

$$\rightarrow \text{DMS E}(\theta) = \frac{1}{m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2$$

$$\nabla_{\theta} \text{DMS E}(\theta) = \frac{2}{m} X^T (X\theta - y)$$

→ Adding More feat to underfitting model will not help  
use diff model.

In case of overfitting add more data

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

Model Complexity ↑ → overfitting

↓ → underfitting

→ Ridge ( $L^2$ ): -  $\alpha \sum_{i=1}^n \theta_i^2 = \rightarrow$  add to cost function only.

ds' hyp: - ,  $\alpha = 0$ , it's linear,  $\propto \sim \omega$ , all are  $\omega$  weights  $\rightarrow 0$ . flat line going through mean.

$$J(\theta) = \text{MSE}(\theta) + \frac{\lambda}{2} \sum_{i=1}^n \theta_i^2$$

→ I'd the more the value:

$$\|V\|_k = (|V_1|^k + |V_2|^k + |V_3|^k + \dots + |V_n|^k)^{\frac{1}{k}}$$

more  $k \uparrow$  neglects small values and put emphasis on outliers so RMSE is sensitive to MAE.

- Standards scalar should be used before  $L_2$  or Ridge Norm.

$$\hat{\theta} = (X^T X + \lambda I)^{-1} X^T y$$

(A is  $(n+1) \times (n+1)$  identity Mat.)  
↳ uses Andre-Louis Cholesky method.

→ Lasso ( $L^1$ ): -  $J(\theta) = \text{MSE} + \alpha \sum_{i=1}^n |\theta_i|$

[It tends to eliminate weight of last imp feature make it 0. i.e. all high degree poly are zero, - Automatically perform feature selection and outlier sparse.

→ Elastic Net (Middle ground): -  $\gamma = \text{mix ratio}$ .

$$J(\theta) = \text{MSE} + r \sum_{i=1}^n |\theta_i| + (1-r) \frac{1}{2} \sum_{i=1}^n \theta_i^2$$

↳ preferred over Lasso.

→ Logistic or Logit Reg: Binary classifier.  $\text{logit } (\beta) = \log \frac{p}{1-p}$

$$\hat{p} = h_\theta(x) = \sigma(X^T \theta) \text{ sigmoid. } \frac{1}{1+e^{-x^T \theta}}$$

$$c(\theta) = \begin{cases} -\log(p) = c_f & y=1 \\ -\log(1-p) = c_o & y=0 \end{cases}$$

$$\rightarrow J(\theta) = \frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(\hat{p}^{(i)}) + (1-y^{(i)}) \log(1-\hat{p}^{(i)})]$$

↳ convex function.

$$\rightarrow \text{Softmax: } S_k(x) = \frac{e^{x^T \theta^k}}{\sum_j e^{x^T \theta^j}}$$

$$P_k = \sigma(S_k(x))_k = \frac{e^{S_k(x)}}{\sum_{j=1}^k e^{S_j(x)}} \quad k = \text{no. of class}$$

$S(x)$  = Vector containing the scores.

$$y^* = \arg \max_k \sigma(S(x))_k = \arg \max_k S_k(x) = \arg \max_k (\theta^{(k)})^T x$$

$$\rightarrow \text{Cross entropy: } -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(\hat{P}_k^{(i)})$$

$y_k^{(i)}$  → target prob.  $i^{\text{th}}$  in  $k^{\text{th}}$  class.

$$J_0 = \frac{1}{m} \sum_{i=1}^m (\hat{P}_k^{(i)} - y_k^{(i)}) x^{(i)}$$

→ SVM: - Sensitive to feature scale.

→ Hard Margin: - Sensitive to outliers and works with only linearly separable data. how value of  $C \rightarrow$  slack in SVM.  
(is Regularisation factor).

→ Also instead of using Linear SVC, SGD can be used with loss = 'hinge',  $\alpha = 1/(m+c)$  applies to Stochastic GD, to train SVM classifier. (used for huge datasets)

$$\text{Non-linear SVM: } \text{Gauss: } e^{-\gamma \|x - x_i\|^2} \quad (\gamma \text{ and } c \text{ hyp})$$

Complexity:  $\theta = \text{tol insensitivity}$  only SVC supports Kernel trick.

|            | Scaling                                    | Kernel |
|------------|--|--------|
| Linear SVC | ○ (m+n)                                    | ✗      |
| SGD Class. | ○ (m+n)                                    | ✗      |
| SVC        | ○( $m^2 \times n$ ) to ○( $m^3 \times n$ ) | ✗      |

Regression:  $t$  is more insensitive to data

$$\rightarrow \hat{y} = \begin{cases} 0 & \text{if } S^T X + b < 0 \\ 1 & \text{if } S^T X + b \geq 0 \end{cases}$$

$$\min_{w, b} \frac{1}{2} w^T w \approx \frac{1}{2} \|w\|^2$$

$$\text{Subject: } t^{(i)} (w^T x^{(i)} + b) \geq 1 \quad \text{for } i = 1, 2, \dots, n$$

## Softmargin

$$\min_{w, b, \xi} \frac{1}{2} w^T w + C \sum_{i=1}^n \xi^{(i)}$$

$$\text{subject } t^{(i)} (w^T x_i + b) \geq 1 - \xi^{(i)}, \xi^{(i)} \geq 0, i=1 \dots n$$

Both are convex quadratic optimization (Rg Richard Breiman  
→ formal normal Boyd)

dual for SVC.

→ Kernel: → in ML means a function capable of computing the dot product  $\phi(a)^T \phi(b)$  based only on original vector  $a$  and  $b$ , without having to compute the transformation  $\phi$ .

$$\rightarrow \text{Linear: } K(a, b) = a^T b$$

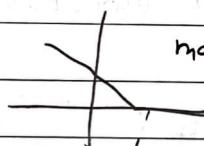
$$\text{Poly: } K(a, b) = (\gamma a^T b + r)^d$$

$$\text{Gaussian RBF: } K(a, b) = \exp(-\gamma \|a - b\|^2)$$

$$\text{Sigmoid: } K(a, b) = \tanh(\gamma a^T b + r)$$

→ online SVM: Online means learning incrementally, typically as new instances arrive.

$$\max(0, 1-t) = \text{large loss.}$$



$$J(w, b) = \frac{1}{2} w^T w + C \sum_{i=1}^n \max(0, 1-t^{(i)}) (w^T x^{(i)} + b)$$

### \* Decision Trees: (both clas. and reg. task)

Adv: - Less data prep, Almost no scaling or centring.

Gini calculates impurity, pure  $\Rightarrow$  gini = 0

$$G_C = 1 - \sum_{k=1}^K p_k k^2$$

→ CART is used by Scikit. they have always 2 child unlike older ones ID3 which have more than 2.

↳ Decision tree is called white box Model

because it's easy to interpret unlike Neural and Random forest.

$$\text{Depth} = \log_2 [n]^2$$

## CART (Classification and Regression tree)

CLASSMATE

Date \_\_\_\_\_

Page \_\_\_\_\_

nug instances Gini or Measure of impurity

$$\rightarrow J(k, t_k) = \frac{m_{\text{left}}}{m} G_{\text{left}} + \frac{m_{\text{right}}}{m} G_{\text{right}}$$

→ Greedy algo. ( $O(\exp(m))$  time.)

→  $O(\log_2 m)$  nodes. Also gini is faster to compute than Entropy.

→ Non-parametric model.  $\because$  no ~~no~~ parameters is not clear prior to training.

→ max depth is hyperparam. And used L2 regularization too.

→ Cart for Reg.  $J(k, t_k) = \frac{m_{\text{left}} \text{MSE}_{\text{left}}}{m} + \frac{m_{\text{right}} \text{MSE}_{\text{right}}}{m}$

$$\text{where } \text{MSE} = \sum_{i \in \text{node}} (y_{\text{node}}^{(i)} - \bar{y}^{(i)})^2$$

$$\bar{y}^{(i)} = \frac{1}{m_{\text{node}}} \sum_{i \in \text{node}} y^{(i)}$$

- Cons:
- They love to be orthogonally decision boundary, makes sensitive to training set rotation.
  - Small variation will cause changes.

Overfitting & max-depth, independent of scaling.

## \* Ensemble Learning and Random Forests:-

→ Due to law of large numbers, work best on methods that are independent of each other

→ Voting Classifier (V).

Bootstrap agg.

Bagging:- Sample is performed with replacement, this is called bagging. And Pasting when sampling is performed without replacement. [Pg 193 diagram].

- Allows training instances to be sampled several times for same predictor.
- Predictors can be trained in parallel, and prediction can be made in parallel.

Random  
Patches  
Random  
Subspace

→ Bagging ends up slightly higher bias than boosting, but small var. & overall bagging results in better model.

However ultimately cross val should be used to test if bagging or boosting is efficient.

→ Out-of-bag eval. Due to bootstrapping, with resp. 63.1% are sampled for each predictor  $(1 - e^{-1}) \approx 63.2\%$

Remaining 37% are not sampled called out-of-bag (OOB) instances.

Random patches → 1 max. feat., bootstrap features are two hyperparameters.

Subspaces → Good for high-D pb. Sampling both training ins. and feat. is called Random patches.

Random forest: Ensemble of D.T. generally trained via bagging method. Rand Forest Cfgs - key parameter is DecisionTree Cfg.

→ Instead of searching for the best feat. when splitting a node it ~~performs~~ searches best features among a random subset of feats. → Results in greater div → Greater bias → Low Var. → Overall better model.

→ Extra tree: - Feature + Feature Imp: - using feature imp. var.

→ RF are very handy to get quick understanding of what feat. actually matter, good for feat. selection.

Boosting: Ensemble method that combine several weak learners into a strong learner. → Train predictor sequentially.

→ Ada Boost: [Pg 200]. Drawback: - Not parallelized.

Does not scale as well as Bagging or boosting.

$$T_j = \sum_{i=1}^m w^{(i)} y_j^{(i)} \neq g^{(i)}$$

Hyp. ( $\eta$ )

$$\begin{aligned} y(x) &= \text{argmax}_{j=1}^n a_j \\ \hat{y}_j(x) &= K \end{aligned}$$

Final n

$$\text{Predictor weight: } \alpha = \eta \log \frac{1 - r_j}{r_j}$$

Weight update

$$w^{(i)} \leftarrow \begin{cases} w^{(i)} & \text{if } \hat{y}_j^{(i)} = g^{(i)} \\ w^{(i)} \exp(\alpha_j) & \text{if } \hat{y}^{(i)} \neq y^{(i)} \end{cases}$$

Split has SAME.

- queries then fly & learning rate classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

(iii) Gradient Boosting: Works by sequentially adding predictors to an ensemble; each one correcting its pred.

Uses residual error made by pre. unlike ada boost. tweaks instance weight of predecessors.

→ import XGBoost (optimised implementation of Gradient boost)

Stacking: [Page 206].

DIM REDUCTION:- (i) Projection (ii) Manifold learning,

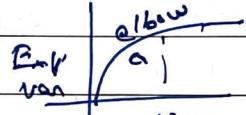
(iii) PCA → AIM FINDING MAX VAR

→ works on CVD! -  $U \in V^T$  contains unit vect that defines all PC.

$$U = \begin{bmatrix} 1 & 1 & 1 \\ \vdots & c_2 & c_n \\ 1 & 1 & 1 \end{bmatrix}$$

(Data should be centered around origin)

→  $X_{\text{proj}} = XW_d$ , to find value of each var ratio.  
explained var ratio.



→ use randomized PCA for faster one dim

$O(m+d^2) + O(d^3)$  vs rather than  $O(m \times n^2) + O(n^3)$

→ Incremental PCA to solve instead of whole dataset

Split into mini batch and feed 1 PCA along one on one-mini batch.

→ Kernel PCA:-

LLC Locally linear embedding → Non-linear dim. Red.

→ Manifold technique. → don't rely on projection

→ Note Not good for large dataset due to  $m^2$  complexity