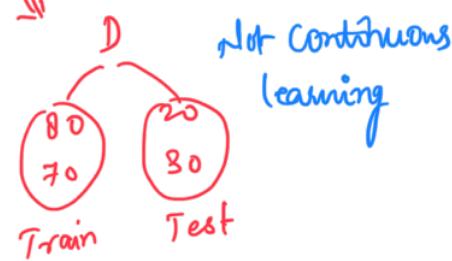
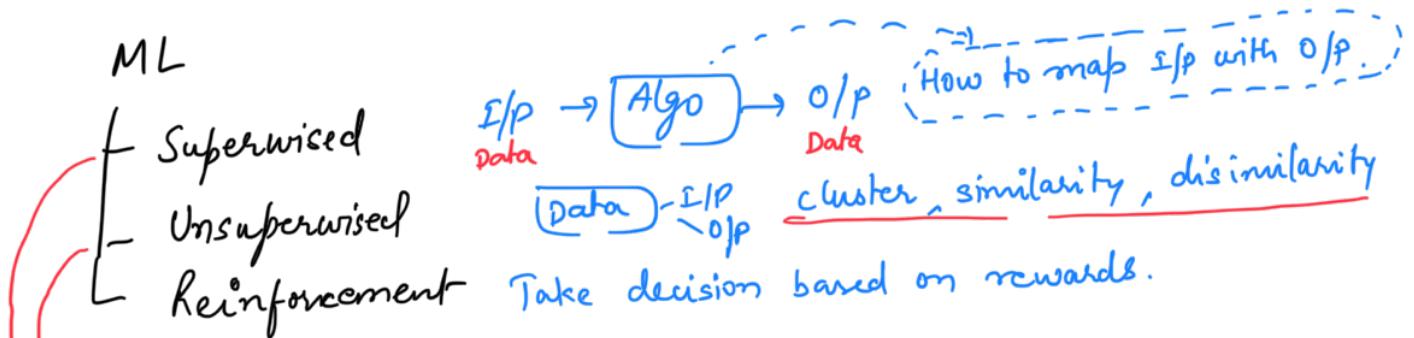


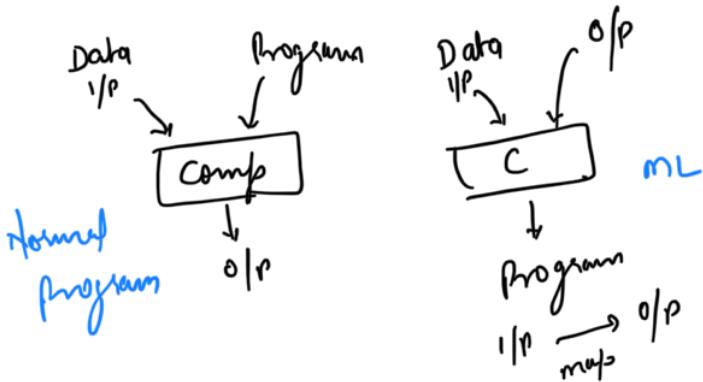
Machine learning - mid sem

Artificial Intelligence

- ✓ - Narrow Intelligence - specific task / weak AI
- ✗ - General " - Strong AI (Human Intelligence)
- ✗ - Super " - Hypothetical (No human intervention)



Date : 18-Aug-22

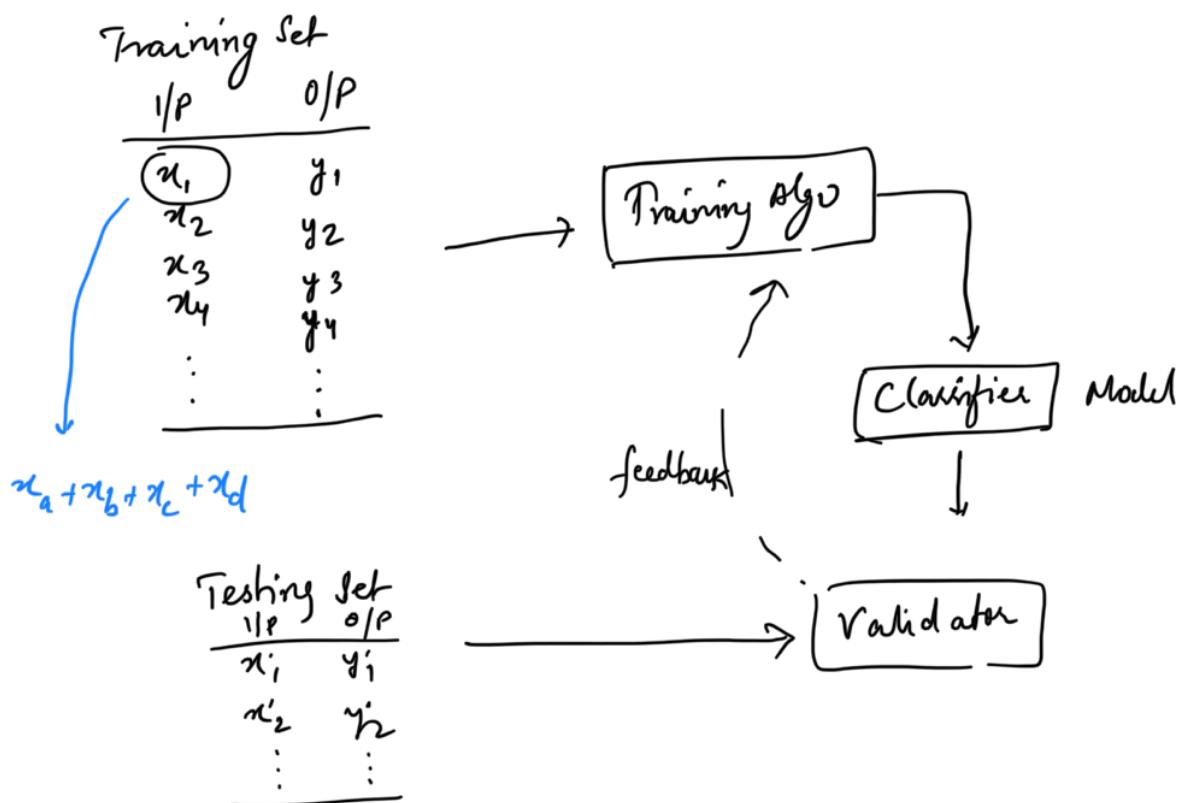


- ① Performance ↑
② Task (Experience) ↑
wrt.
Time
- Inductive learning
- ↳ ① Human
② Machine
③ Ant/Bee/Fungus

Supervised Learning

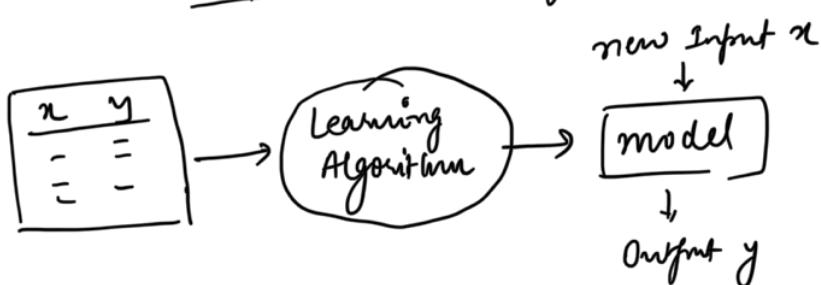
- Here question arises about: Complexity V/s Accuracy V/s Noise

- Assumptions we take in ML
(i.e. biases)
 - ↳ How much complexity we require to our model
 - ↳ which order we decide the complexities
 - i.e. n^0 to n^{10} so we can stop till
less complex More complex x^7 or x^8 .

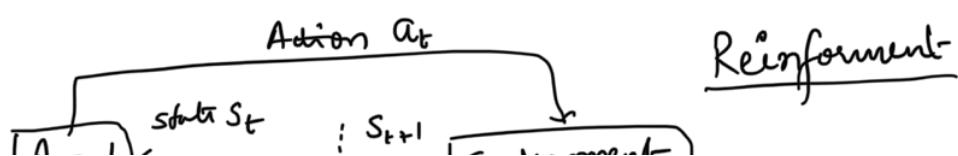
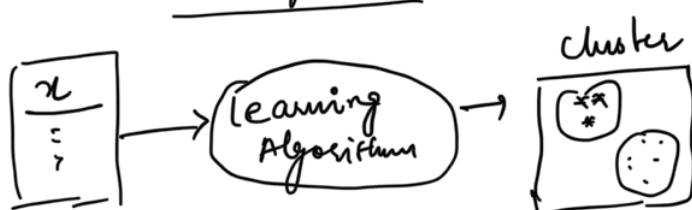


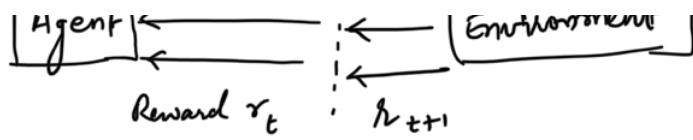
Date: 24-Aug-22

Supervised Learning



Unsupervised





S-L # Mathematical Explanation

Hypothesis space

- All possible solutions to particular problem.

Biases

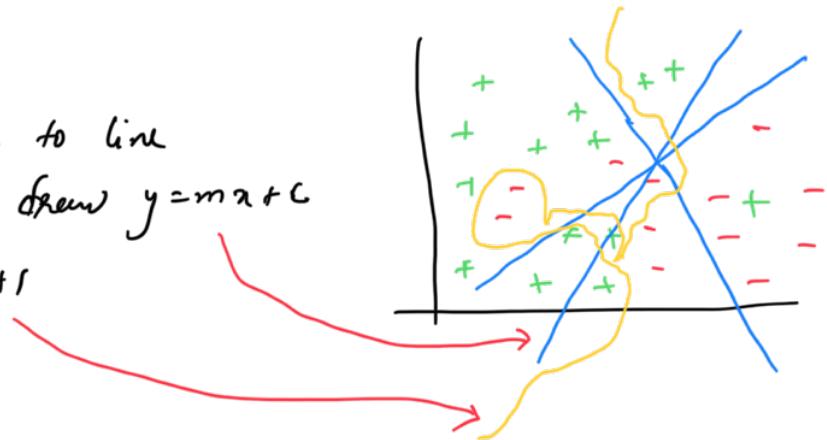
- To reduce the hypothesis space.

Next Class : Biases & Generalization

Date : 25-Aug-22

- Hypothesis h : fn. that approximates f .
- Hypo. space H : set of fn. we allow for approximating f .
- The set of hypothesis that can be produced, can be restricted further by specifying a language bias.
- Input : Training set $S \subseteq \mathcal{X}$
- Output : A hypothesis $h \in H$

Bias 1 : Restrict to line
i.e. we only draw $y = mx + c$
 $m = -2 + x^2 + x + 1$



Here, learning is Refining the Hypothesis Space.

Def: Bias : The tendency to prefer one hypothesis over another.

Ex: Occam's Razor

- A classical example of inductive bias.
- The simplest consistent hypothesis about the target function is actually the best.

Next Class : Problems in ML and Generalization

Date - 26-Aug-22

Important Issue in ML

..

Generalization
 components of gen" error < ^{Bias}
 variance : How much models estimated
 from different training sets differ
 from each other.

Underfitting and Overfitting

Underfitting < high bias and low variance
 High training error and high test error

Overfitting < low bias and high variance
 Low training error and high test error

NOTE : It is good if we are facing underfitting but overfitting
 is very bad always.

Experimental Evaluation of Learning Algorithms

→ Performance Evaluation:
 - Error : AE, MSE etc.
 - Accuracy
 - Precision / Recall

$$\text{Absolute Error} : \frac{1}{n} \sum |f(x_i) - \hat{y}_i| \quad \left. \right\} \text{Regression}$$

$$\text{MSE} : \frac{1}{n} \sum (f(x_i) - \hat{y}_i)^2 \quad \left. \right\} \text{Classification}$$

$$\frac{1}{n} \sum (f(x_i), y_i) \begin{cases} 0 \\ 1 \end{cases}$$

→ Sampling Methods :

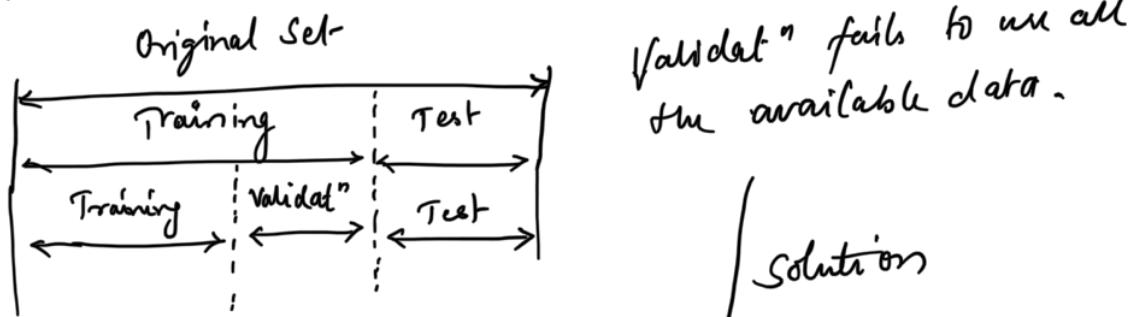
- Train / Test sets
- K-fold cross-validation (when dataset is very less)

Difficulties in evaluating hypotheses with limited data

① Bias in the estimate: The sample error is poor estimator
 of true error.

② Variance in the estimate: The smaller the test set the
 greater the expected variance.

Validation Set



K - Fold cross-validation

- 1) split the data into K equal subsets
- 2) perform K rounds of learning; on each round
 - $\frac{1}{K}$ of the data is held out as a test set and
 - the remaining examples are used as training data.
- 3) Compute the average test set score of the K rounds.

↓ So,

Trade-off

- In ML, there is always a trade-off b/w
 - complex hypotheses that fit the training data well
 - simpler hypothesis that may generalise better.
- As the amount of training data increases, the generalization error decreases.

Date : 27 - Aug - 22

Regression

Single variable LR

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

Population
y-intercept Population
slope Random
error

Variance

- Deterministic
- stochastic / probabilistic

Assumptions

- The data may not form a perfect line
- $E(\epsilon_i) = 0$ for $i = 1, 2, 3, \dots, n$. (Total error)
- $\sigma(\epsilon_i) = \sigma_\epsilon$
- The errors are independent
- The ϵ_i are normally distributed (with mean 0 and standard deviation σ_ϵ)

Multiply
varient LR

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

How do we "learn" parameters

- To find the values for the coeff. which minimize the objective function we take the partial derivatives of the objective function (sse) w.r.t the coefficients. set them to zero and solve

$$\beta_1 = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2} \quad \beta_0 = \frac{\sum y - \beta_1 \sum x}{n}$$

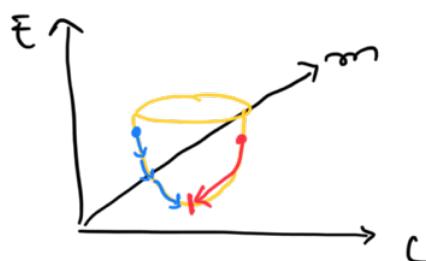
for $y = \beta_0 + \beta_1 x$

Minimize

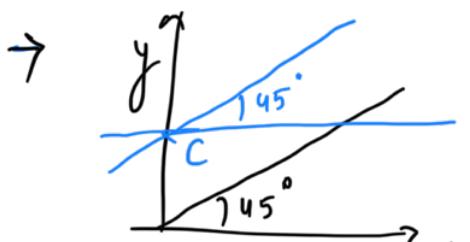
$$\sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)]^2$$

Gradient Descent Algorithm (LMS Algorithm)

least mean slope



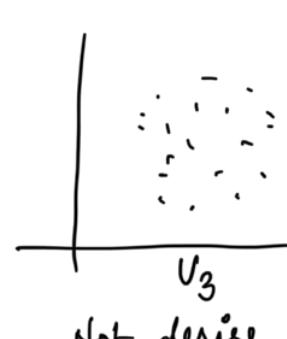
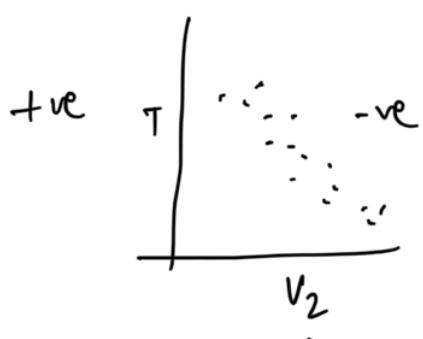
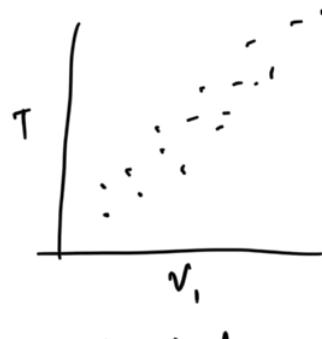
Data: 31-Aug



$$\begin{aligned} y &= x \\ &= 1 \cdot x \\ &= \tan 45^\circ \cdot x = mx \end{aligned}$$

$$y = mx + c$$

Coefficient of Variance (correlation)

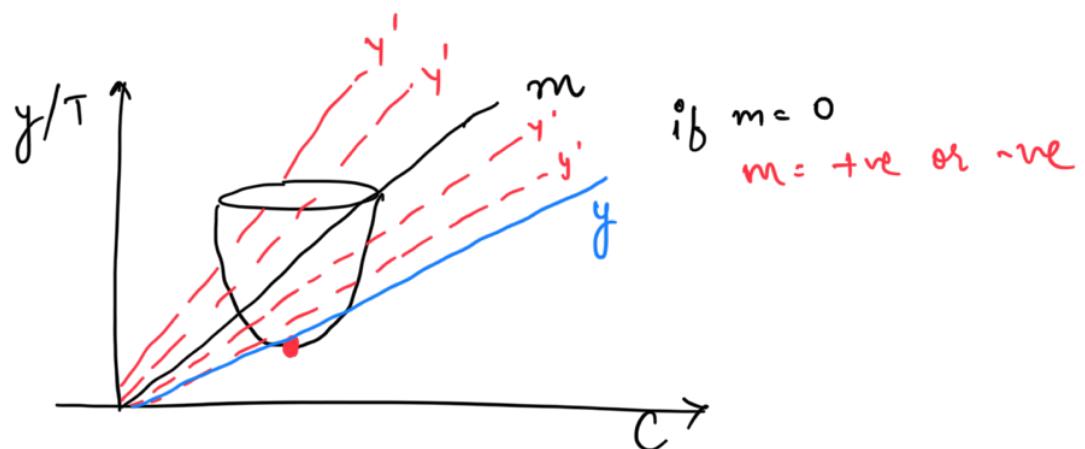


Zero

Not desire

Desired.

Desired



Coeff. of relation - Pearson's Coefficient $r(x, y)$

$$R = r_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}} \approx +1 \text{ or } -1$$

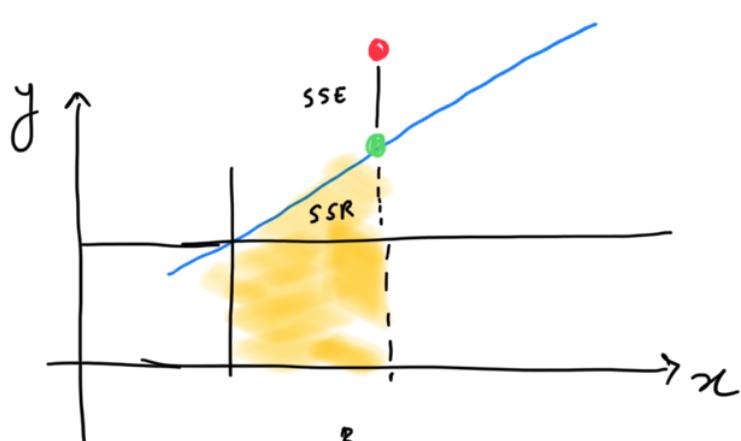
NOTE It is significant or we use only for finding coeff. of relation b/w dependant and independent variable.

Coeff. of Determinant-

NOTE In LR we try to minimize the stochastic variance.

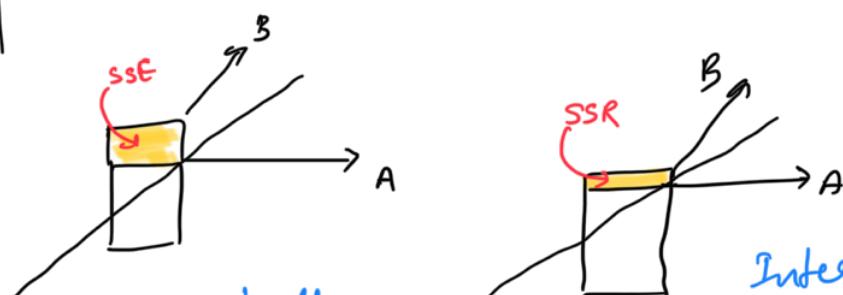
Variance

- 1) stochastic variance
- which we don't know
- 2) deterministic variance
- which we know



$$\left. \begin{array}{c} y \\ \hat{y} \\ \bar{y} \end{array} \right\} 6.1 \quad \left. \begin{array}{c} y \\ \hat{y} \\ \bar{y} \end{array} \right\} 5.9 \quad \left. \begin{array}{c} y \\ \hat{y} \\ \bar{y} \end{array} \right\} 5.5$$

Residual error (SSE)



we can handle this easily but not interested.

Interested to reduce this residual error.

Decision Tree

Issues

- what kind of data the D.T. will be generated.



- One proposal: prefer the smallest tree that is consistent with data (Bias)
i.e either less no. of leaves or low height.
- SPL: state space search

Choices

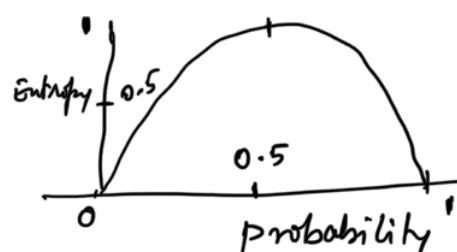
- When to stop
 - no more input features
 - All examples are classified the same
 - too few examples to make an informative split
- Which test to split on
 - split gives smallest error
 - with multi-valued features
 - split on all values or
 - split value into half

Principled Criterion

(1) Information gain

- Entropy = Impurity or randomness

Y/N	100% Y	$\rightarrow E = 0$
Case	100% N	$\rightarrow E = 0$
	50% Y & 50% N	$\rightarrow E = 1$



$\text{Gain}(S, A)$ = expected reduction in entropy due to

partitioning S on attribute A

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{\text{values}(r)} |S_r| / |S| \text{Entropy}(S_r)$$

(2.) Gini index

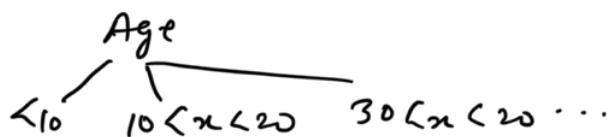
- measure of node impurity

(3.) Chi square

(4.) Reduction in Variance

Continuous Attribute :

- what could be the class / cluster size?



NOTE Problem in D.T. is that it works on greedy approach.

Random Forest -

- Builds multiple decision tree and merges them together.
- More accurate and stable prediction
- * Random decision forest correct for decision tree's habit of overfitting to their training set.
- Trained with the "bagging" method.

Overfitting

7 - Sep - 22

Reasons

- 1) Dataset set is small
- 2) Due to noise → remove the outlier
- 3) Data has high variance → cross validation
- 4) The model is too complex

Reason for Underfitting

- (4) The model is too simple
- (1) to (3) is same as above

NOTES on Overfitting

It happens when model is capture from limited setup.

K-NN

9-Sep-22

1 nearest neighbor

Training → Saving the instances
Testing / Predicting → instance x_j given
find (x_i, y_i) nearest to x_j
predicting → of y_i

KNN find $\{(x_1, y_1), (x_2, y_2) \dots (x_k, y_k)\}$
classification → $\{y_1, y_2, y_3 \dots y_k\} \leftarrow \begin{cases} \text{Mean} \\ \text{Median} \end{cases}$ y_j
Regression → $\left[\begin{array}{c} y_1 \\ y_2 \\ \vdots \\ y_k \end{array} \right] \leftarrow \begin{cases} \text{Mean} \\ \text{Median} \end{cases}$

Next :- locally weighted Average ...
class

Error

10-Sep-22

Feature Selection

16-Sep-22

- Univariate F-S. Methods
- Multivariate - "

Optimization in KNN

- ① Increase the accuracy
- ② Reduce the complexity

Feature Extraction

→ find a projection matrix w from N -dim to M -dimensional vectors that keeps

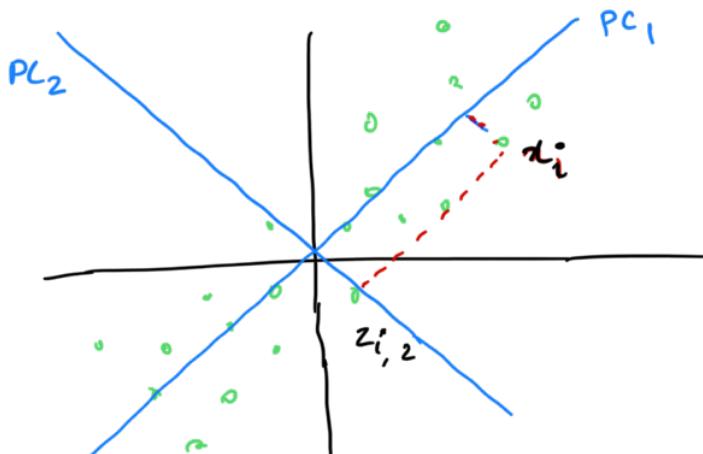
error low.

$$M < N$$

→ Expected:

- Uncorrelated, can't be reduced further $[\text{corr.}(x_1, x_2, \dots, x_n)]$
- Have large variance or otherwise bear no information $[\sigma^2(x_1, x_2, \dots, x_n)]$

Principle Components



if 2 variables (or principle component) are orthogonal, then there is no correlation b/w them.

Problem in PCA projection

- Data variation determines the projection
- class information will be missed?



SOLUTION

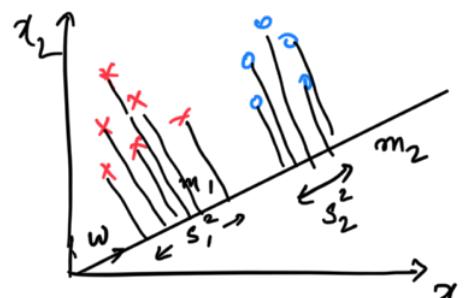


LDA

- Maximize between-class distance.
- Minimize within-class distance

Fisher linear discriminant-

$$j(w) = \frac{(m_1 - m_2)^2}{S_1^2 + S_2^2}$$



Imp

Recommendation System

17-Sept-22

Types
1) - Popularity RS 20%

2) - Content-based RS 30%

3) - Collaborative RS 50% User based RS

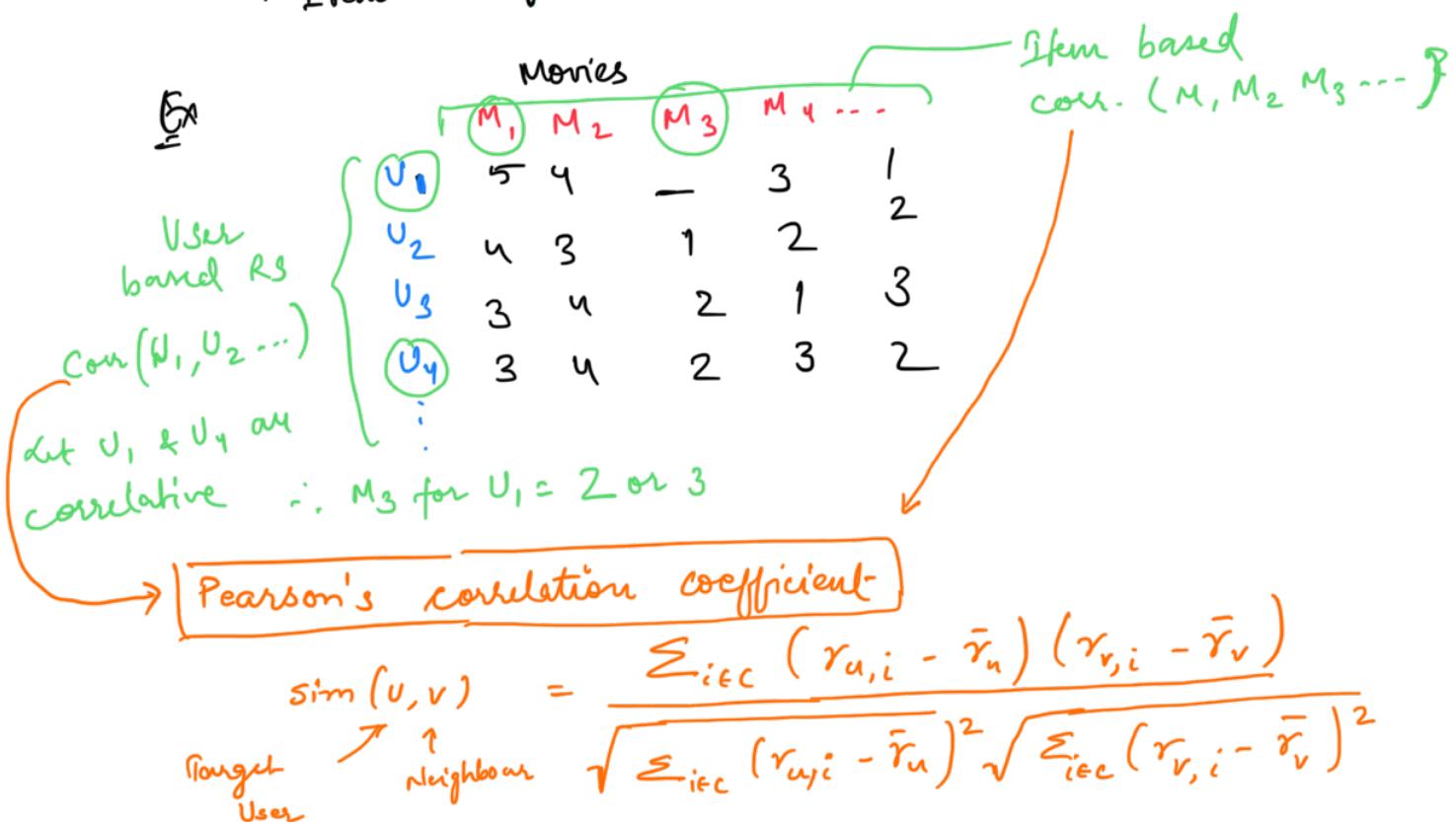
4) - Hybrid (100% used)

} All are based on KNN

3) Collaborative filtering for Rating Prediction

- User based :- Nearest similar user (Neibour)

- Item based :- Nearest neighbour's item (similar)



- Method/Process name : Neighborhood formation phase

Now Recommendation Phase

- Use the following formula to compute the rating prediction of item i for target user u

$$P(u, i) = \bar{r}_u + \frac{\sum_{v \in V} \text{sim}(u, v) \times (r_{v,i} - \bar{r}_v)}{\sum_{v \in V} |\text{sim}(u, v)|}$$

Average \nearrow

where V is the set of K similar users, $r_{v,i}$ is the rating of user v given to item i ,

Issue

- lack of scalability.

Enhancement-

- Item based CF

Bayesian Learning

21- Sep - 22

Probability in Learning

- Probability
 - Bayesian Estimation
 - Bayesian Probability

Bayes Theorem

- $$\text{Bayes Rule: } p(h|D) = \frac{p(D|h) p(h)}{p(D)}$$

$$P(D) = \text{prior probability of hypothesis, } h$$

$$P(h) = \text{" " " of training data, } D$$

$$P(D) = \text{Prbs. of } h \text{ given } D \text{ (Posterior density)}$$

$$P(D|h) = \text{Prbs. of } D \text{ given } h \text{ (likelihood of } D \text{ given } h)$$

Gramph

Dialerma

1

M/F ?

→ Given Out of 100 women at the movie

50 - long hair
50 - short hair
96 - short hair
89 - long hair

Prior knowledge

$$\text{and } 98 \text{ men } " " = \frac{98}{2} - S.H$$

To find : $P(\text{man} | \text{long hair}) = ?$

Concepts

Conditional

Marginal

Joint

Probability

$$\begin{aligned} \text{Suppose } p(\text{man with LH}) &= p(LH) \approx p(\text{man} | LH) \\ p(LH \text{ and man}) &= p(m) \approx p(\text{man} | LH) \\ \therefore p(m \& LH) &= p(LH \& m) \end{aligned}$$

1

$$p(m|w) = \dots$$

Maximum A Posteriori (MAP) hypothesis

$$p(h|D) = \frac{p(D|h) p(h)}{p(D)}$$

The goal of Bayesian learning: the most probable hypothesis given the training data (max. A posteriori hypothesis)

$$h_{MAP} = \arg \max_{h \in H} P(h|D)$$

$$= \arg \max_{h \in H} \frac{p(D|h) p(h)}{p(D)}$$

$$= \arg \max_{h \in H} p(D|h) p(h)$$

#1 Compute ML hypo (ML - Maximum likelihood)

$$h_{ML} = \arg \min_{h \in H} \sum_{i=1}^m (d_i - h(x_i))^2$$

Bayes Optimal

22-Sep-22

NOTE finding $p(y|x)$ is hard, while it is easy to find $p(x|y)$.

$$\hat{P} \quad p(y|x)$$

$$p(1|x=(0,2))$$

$$p(y=0) = \frac{6}{10}$$

$$p(y=1) = \frac{4}{10}$$

$$p(x=(0,2)|y=1) = \frac{1}{4}$$

x_1	x_2	y
0	0	0
0	1	1
1	2	1
0	0	1
2	2	0
1	1	0
0	2	1
2	0	0

$$P(x=(0,2) | y=0) = 0 \quad \begin{matrix} 2 & 1 & 0 \\ 1 & 0 & 0 \end{matrix}$$

$$P(x=(0,2) | y=0) * P(y=0) = 0 * \frac{6}{10} = 0$$

$$P(x=(0,2) | y=1) * P(y=1) = \frac{1}{4} * \frac{4}{10} = \frac{1}{10}$$

Estimated value of y is 1 given $x=(0,2)$

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

likelihood prior
 ~ Evidence

Problems

- Prob. of particular $x=\{\cdot\}$ set in very large dataset - is very complex.

Solution → Naive Bayes

- Assume all x are independent.

$$\therefore P(x=(0,2) | y=1) = P(x_1=0 | y=1) * P(x_2=2 | y=1)$$

$$= \frac{3}{4} * \frac{2}{4} = \frac{3}{8}$$

$$\& P(x=(0,2) | y=0) = P(x_1=0 | y=0) * P(x_2=2 | y=0)$$

$$= \frac{1}{6} * \frac{1}{6} = \frac{1}{36}$$

$$\Rightarrow y = 2 - 2 \rightarrow 1 \perp \perp \frac{6}{36}$$

- To $\bar{q} \bar{u} / \bar{b} b^{10}$

Don't forget to multiply the priors probability. i.e $P(y=1)$ and $P(y=0)$.

- - -

Problem in Naive Bayes

Smoothing is the solution to problem where y^{new} is zero.

* Text recognition is the practical example of Naive Bayes classifier or regression.

Note Even if the variable are not independent Naive Bayes classifier works well.