

Data Mining 1/24 Lecture 2.

Data type: Categorical

Continuous.

Data Processing \Rightarrow Data Mining \Rightarrow Postprocessing

OSEMIN: Obtain. \rightarrow Scrub \rightarrow Explore \rightarrow Model \rightarrow Interpret

SQL

Clean

find max/min

mean

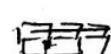
distribution



Machine/early

Handle missing values:

{ List wise deletion



{ Pair wise deletion



Imputation

Quality: Missing values.

Errors / inconsistent values.

Duplicate values

Noise / Outliers

typos

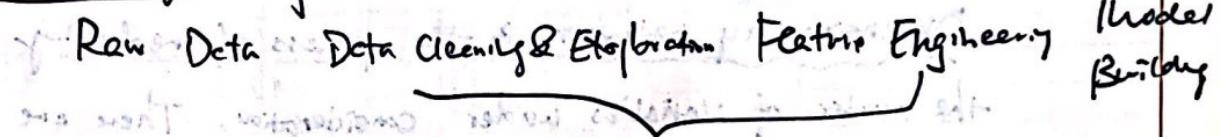
legitimate values but not in range

Exploratory Data Analysis (EDA)

Univariate Analysis : Categorical

Bivariate Analysis : Continuous Vs Continuous.

Feature Engineering



- Feature Transformation \Rightarrow (take log)

- Feature Creation

- Feature Scaling

(should be performed inside the cross-validation loop)

* Normalization : Scale all data to $[0-1]$

$$z = \frac{x - \min(x)}{\max(x) - \min(x)}$$

* Standardization : Rescale to have a mean of 0 and standard deviation of 1

$$z = \frac{x - \mu}{\sigma}$$

Noise/Outliers

Leave them as-is

Deletion

Transform \Rightarrow handle outliers

• Standardize or binning

Binning

Imputation

Separation : separate process

• Binning / Discretization : transform a cont. attribute into a categorical attribute (or to a binary, binarizing)

method examples : eg. k-means, Equal frequency.

Sampling : Reduce the dataset size

Aggregation : data by date \Rightarrow aggregate to by week/month

Dimensionality : The number of attributes in dataset

We can think of each row as a d-dimensional point

(e.g. show the 4th dimension by color)

Curse of Dimensionality

Dimensionality reduction is the process of reducing the number of variables under consideration. There are 2 methods:

feature selection by selecting a subset of attributes

feature extraction

Dom in knowledge

Missing Values ratio

Low variance filter

High correlation filter

Feature creation

Feature Selection Methods :

- Filter Methods : Score each feature based on statistical correlation with target value

- Wrapper Methods : Train a model using a subset of features. Based on the results, decide to add or remove features from our subset.

- Embedded Methods : Feature selection occurs as

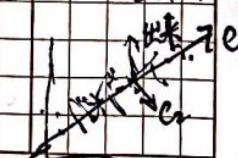
Feature Extraction / Dimensionality Reduction part of the data mining algorithm.

(PCA) Principle Component Analysis

(SVD) Singular Value Decomposition

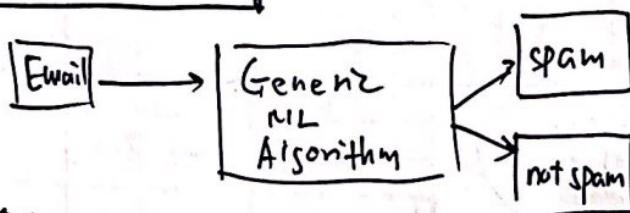
(LDA) Linear Discriminant Analysis

Reduction : Use techniques from Linear algebra to project the data from a high-dimensional space to a lower space in such a way that maximum variance of data is captured.



PCA identifies the best hyperplane to project the data onto

Machine Learning



ML PK AI : ML is a subset of AI

Supervised learning : Labelled training data

Unsupervised learning : Unlabelled training data

Reinforcement learning :

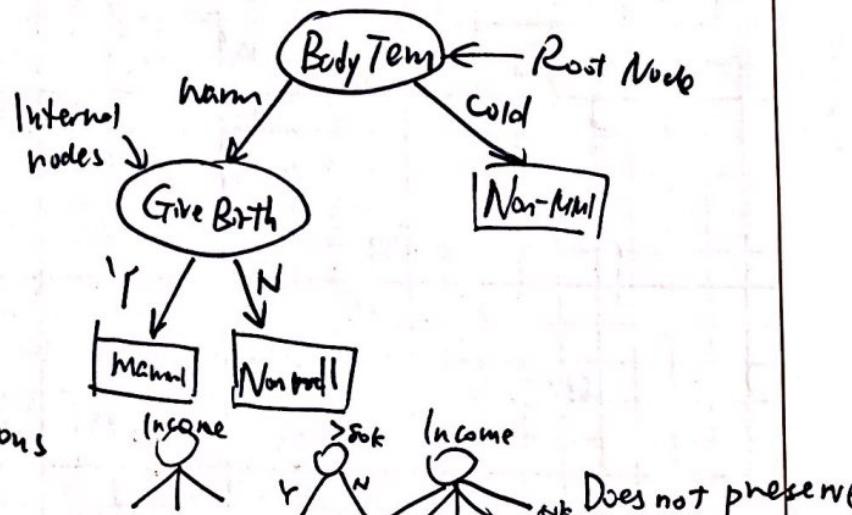
```

graph TD
    Env["Env"] -- "S, R" --> Model["Model"]
    Model -- "A" --> Game["Game: Chess - fail, try again"]
    Game -- "S, R" --> Env
    Env -- "S, R" --> FB["Reinforcement Feedback"]
    FB -- "3/5 are correct, try again" --> Model
  
```

Classification (Supervised learning)

Assigning:

Decision Trees



- Attribute Test Conditions

- Hunt's Algorithm

- All same class?

Yes → leaf

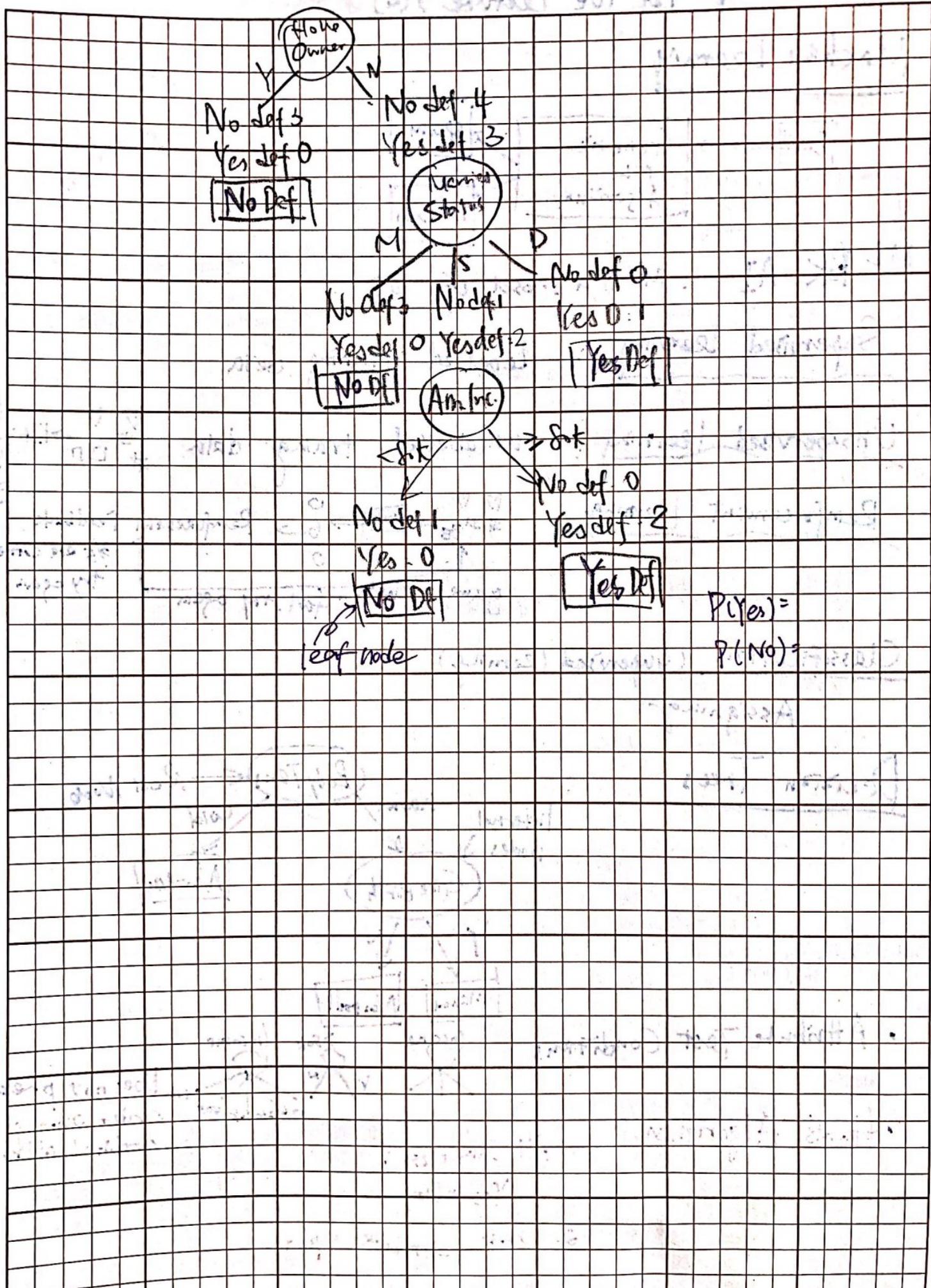
- Choose attribute split

No def: 3
Yes def: 0

No def

No def: 4
Yes def: 3

Married



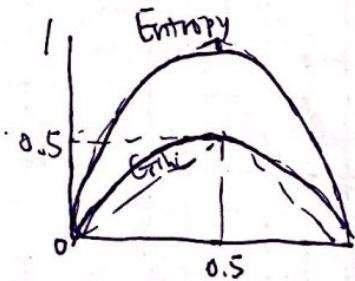
Greedy Algorithm: Which node is the best?

Selecting the best split.

Example: Impurity measures

- Entropy: $E(S) = \sum_{i=1}^C -p_i \log_2 p_i$

- Gini: $Gini = 1 - \sum_{i=1}^C (p_i)^2$



GAIN parent node \rightarrow child node the GAIN, the better split.

$$GAIN_{split} = Impurity(p) - \left(\sum_{i=1}^n \frac{n_i}{n} Impurity(i) \right)$$

Example impurity of parent $Gini_{parent} = 1 - \left(\frac{3}{10}\right)^2 - \left(\frac{7}{10}\right)^2 = 0.42$ [loan/class example Default]

* if split homeowner



Def No: 3 Def No: 4

Def Yes: 0

Def Yes: 3

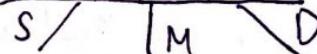
$$Gini = 1 - \left(\frac{0}{3}\right)^2 - \left(\frac{3}{3}\right)^2 = 0$$

$$Gini = 1 - \left(\frac{3}{7}\right)^2 - \left(\frac{4}{7}\right)^2 = 0.489$$

$$\sum Gini_{split} = \left(\frac{3}{10}\right) \cdot 0 + \left(\frac{7}{10}\right) \cdot 0.489 = 0.342$$

$$GAIN = 0.42 - 0.342 = 0.078$$

* if split on Marital status



Def No: 2 Def N: Y Def N: 1

Def Yes: 2 Def Y: 0 Def: 1

$$Gini = 0.5 : Gini = 0 : Gini = 0.5$$

$$\sum Impurity = \left(\frac{4}{10}\right) \cdot 0.5 + \left(\frac{4}{10}\right) \cdot 0 + \left(\frac{2}{10}\right) \cdot 0.5 = 0.3$$

$$GAIN = 0.42 - 0.3 = 0.12$$

(1) Standard Deviation Method

Split on Annual Income $\leq 55K$.

$$\begin{array}{ll}
 \text{Def } N: 0 & \text{Def } N: 7 \\
 \text{Def } Y: 0 & \text{Def } Y: 3 \\
 G_{\text{ini}} = 0 & G_{\text{ini}} = 1 - \left(\frac{3}{10}\right)^2 - \left(\frac{7}{10}\right)^2 = 0.42 \\
 \end{array}$$

$$\sum h_i = \frac{1}{10} \cdot 0 + \frac{7}{10} = 0.7$$

$$G_{\text{AN}} = 1 - \left(\frac{3}{7}\right)^2 - \left(\frac{4}{7}\right)^2 = 0.142$$

On Annual Income $\leq 65K$

$$\begin{array}{ll}
 \text{Def } N: 1 & \text{Def } N: 6 \\
 \text{Def } Y: 0 & \text{Def } Y: 3
 \end{array}$$

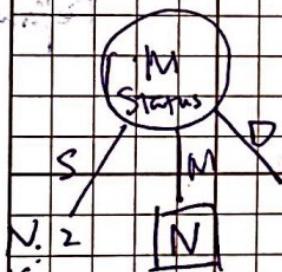
$$G_{\text{ini}} = 0 \quad G_{\text{ini}} = 0.445$$

$$\text{Impurity } \frac{h_i}{n} = \left(\frac{1}{10}\right) \cdot 0 + \left(\frac{9}{10}\right) \cdot 0.445 = 0.4$$

$$G_{\text{AN}} = 0.42 - 0.4 = 0.02$$

Example: - choose the least $\sum G_{\text{hi}}$. And split.

$$\begin{array}{c}
 Y: 3 | 0 \\
 N: 3 | 4 \\
 \hline
 0.34
 \end{array}$$



{ OMS

- (1) HO
- (2) Choose the highest
- (3) A Income

GAIN

90K

$GAIN = 0.445 - 0.4 = 0.045$

(A) 70K

Continuous Variables.

GAIN Ratio

GAIN RATIO = $\frac{GAIN_{\text{split}}}{SplitINFO}$

$$SplitINFO = - \sum_{i=1}^k \frac{n_i}{n} \log_2 \frac{n_i}{n}$$

Date 14th Jan 1/31 Thursday Lecture 4(2)

GAIN Ratio Continue

Example

$$SI(\text{marital status}) = -\frac{4}{10} \log_2 \left(\frac{4}{10}\right) - \frac{4}{10} \log_2 \left(\frac{4}{10}\right) - \left(\frac{2}{10}\right) \log_2 \left(\frac{2}{10}\right)$$

$$= 1.52$$

$$\Rightarrow GAINRATIO = \frac{0.12}{1.52} = 0.079 \quad \text{for Marital Status Attribute.}$$

$$SI(\text{Home owner}) = -\frac{3}{10} \log_2 \frac{3}{10} - \frac{7}{10} \log_2 \frac{7}{10} = 2.88$$

$$\Rightarrow GAINRATIO = \frac{0.078}{2.88} = 0.027 \quad \text{for Attribute Home Owner}$$

TIP: If all the split in the same data type (binary), choose GAIN
otherwise, choose GainRatio

Decision Boundaries: Use axis-parallel hyperplane to split data space into purer partitions.

Data Fragmentation

Decision Tree Algorithm

Regression Trees

Practice Problem

$$\text{Gini (parent Impurity)} = 1 - \left(\frac{4}{7}\right)^2 - \left(\frac{3}{7}\right)^2 = 1 - \frac{16}{49} - \frac{9}{49} = 0.49$$

- ① if split by Pass all Assignments

Pass Y: 2 Pass N: 2

Pass Y: 2 Pass N: 2

$$\begin{aligned} GINI_1 &= 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 \\ &= 1 - \frac{1}{2} - \frac{1}{2} \\ &= \frac{1}{2} \end{aligned}$$

$$\begin{aligned} GINI_2 &= 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 \\ &= 1 - \frac{4}{9} - \frac{1}{9} \\ &= \frac{4}{9} \end{aligned}$$

$$\sum GINI(\text{split}) = \frac{4}{7} \cdot \frac{1}{2} + \frac{3}{7} \cdot \frac{4}{9} = 0.476$$

$$GAIN = 0.49 - 0.476 = 0.014$$

$$\begin{aligned} GAIN &= \frac{GAIN}{GI} = \frac{0.014}{-\frac{4}{7} \log_2 \frac{4}{7} - \frac{3}{7} \log_2 \frac{3}{7}} \\ &= 0.014 \end{aligned}$$

- ② if split by GPA < 3.1

Pass Y: 1 Pass Y: 3
Pass N: 2 Pass N: 1

$$\begin{aligned} GINI_1 &= 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 \\ &= 1 - \frac{1}{9} - \frac{4}{9} \\ &= \frac{4}{9} \end{aligned}$$

$$\begin{aligned} GINI_2 &= 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 \\ &= 1 - \frac{9}{16} - \frac{1}{16} = 1 - \frac{10}{16} = \frac{6}{16} \\ &= \frac{3}{8} \end{aligned}$$

$$\sum GINI(\text{split}) = \left(\frac{3}{7}\right) \cdot \frac{4}{9} + \left(\frac{4}{7}\right) \cdot \frac{3}{8} = 0.405$$

$$GAIN = 0.49 - 0.405 = 0.085$$

GPA ≤ 3

Y: 2 N: 5

Pass Y: 0 Pass Y: 4
Pass N: 2 Pass N: 1

$$GINI = 1 - \left(\frac{0}{7}\right)^2 - \left(\frac{2}{7}\right)^2 = 0$$

$$\begin{aligned} GINI &= 1 - \left(\frac{4}{5}\right)^2 - \left(\frac{1}{5}\right)^2 \\ &= 1 - \frac{16}{25} - \frac{1}{25} = \frac{8}{25} \end{aligned}$$

$$\sum GINI(\text{split}) = \left(\frac{2}{7}\right) \cdot 0 + \frac{5}{7} \cdot \frac{8}{25} = 0.229$$

$$GAIN = 0.49 - 0.229$$

$$\begin{aligned} &= 0.261 \quad GR = \frac{0.261}{\frac{2}{7} \log_2 \frac{2}{7} - \frac{5}{7} \log_2 \frac{5}{7}} \\ &= 0.302 \end{aligned}$$

③ If split on Language

Python 3	C++2	Java 2
Pass Y: 1	Pass Y: 2	Pass Y: 1
Pass N: 2	Pass N: 0	Pass N: 1
$G_{IN} = 1 - \left(\frac{1}{3}\right)^2 \cdot \left(\frac{1}{7}\right)^2$	$G_{IN} = 1 - \left(\frac{2}{7}\right)^2 \cdot \left(\frac{2}{7}\right)^2$	$G_{IN} = 1 - \left(\frac{1}{2}\right)^2 \cdot \left(\frac{1}{2}\right)^2$
$= 1 - \frac{1}{9} - \frac{4}{49}$	$= 0$	$= \frac{1}{4} - \frac{1}{4}$
$= \frac{4}{9}$		$= \frac{1}{2}$

$$\sum G_{IN}(\text{split}) = \left(\frac{3}{7}\right) \cdot \frac{4}{9} + \left(\frac{2}{7}\right) \cdot 0 + \left(\frac{2}{7}\right) \cdot \frac{1}{2} = 0.333$$

$$GAIN = 0.49 - 0.333 = 0.157$$

$$\begin{aligned} GAINRATIO &= \frac{GAIN}{SI(\text{lang})} = \frac{0.157}{\frac{3}{7} \log_2 \frac{3}{7} + \frac{2}{7} \log_2 \frac{2}{7} + \frac{2}{7} \log_2 \frac{2}{7}} \\ &= \frac{0.157}{1.557} = 0.1 \\ &= 0.1 \end{aligned}$$

Compare the 3 ways of GAINRATIO
choose Language to split on first.

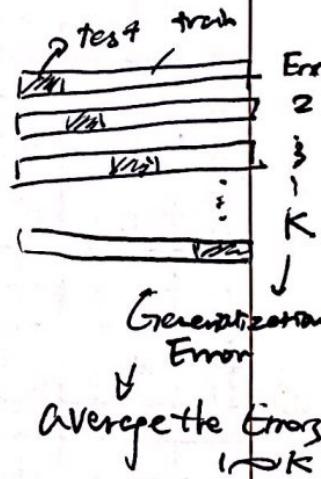
Data Mining 2/5 The Lecture 5 (2)

Types of Error

- Training Error in Training Phase
the percent of misclassification errors on the training set.
 - Testing Error/Generalization Error

Partitioning Data to predict Error on the test data.

K-Fold Cross-Validation : split data into K folds



Overfitting

An over-fitted model : too specific to the training data and will not generalize well to new data.

Pre-pruning: stop grow the tree before it is fully grown.

e.g. χ^2 (Chi-square) pruning example.

$$\chi^2 = \sum_{\substack{\text{all classes } i \\ \text{children } j}} \frac{(O_{i,j} - E_{i,j})^2}{E_{i,j}}$$

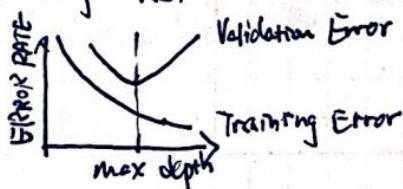
hull : you split on random data. correctly.

Post-pruning Trim the fully grown tree from the bottom up.

-Reduced Error Pruning (REP)

Model Selection

" hyperparameter tuning.



Re select another ~~completely~~ relevant ~~data~~ with.

Data Mining 3/12 Tuesday Lecture

Nearest Neighbor Classifiers → Eager Learner
→ Lazy Learner

K-nearest Neighbors: use euclidian distance, $d = \sqrt{\sum (p_i - q_i)^2}$ Scaling
for nominal attributes: $\begin{cases} \text{distance} = 0 & \text{if they same} \\ \text{distance} = 1 & \text{if different} \end{cases}$

Voting Schemes

Weight ~~vectors~~ factor: $w = \frac{1}{d^2}$

Decision Boundaries

$k=1$, overfitting

$k=20$

\vdots underfitting
 \Downarrow

choose the right k , using cross-validation

Instance Reduction Algorithms

In Class Activity -

$$d_{12} = \sqrt{1 + (7-7)^2 + (7-4)^2} = \sqrt{1+9} = \sqrt{10}$$

$$d_{13} = \sqrt{0 + (7-5)^2 + (7-4)^2} = \sqrt{16+9+1} = \boxed{5} = 5$$

$$d_{14} = \sqrt{1 + (7-1)^2 + (7-4)^2} = \sqrt{1+36+9+1} = \boxed{8} = \sqrt{46}$$

$$d_{23} = \sqrt{1 + 4^2 + 0^2} = \sqrt{1+16+0} = \sqrt{17}$$

$$d_{24} = \sqrt{1 + 6^2 + 0^2} = \sqrt{36} = 6$$

$$d_{34} = \sqrt{1 + 2^2 + 0^2} = \sqrt{5}$$

Data Mining 2/14 Thursday Lecture 8

Bayesian Classifiers

Naive Bayes

: Bayes Theorem
and a "naive" assumption

$$P(A, B) = P(A) \times P(B|A) \quad \text{dependent}$$

$$P(A, B) = P(A) \times P(B) \quad \text{Independent.}$$

Bayes Theorem

$$P(B|A) = \frac{P(A|B) \cdot P(B)}{P(A)} \quad \begin{matrix} \text{conditional} \\ \text{class probability} \end{matrix}$$

Classification

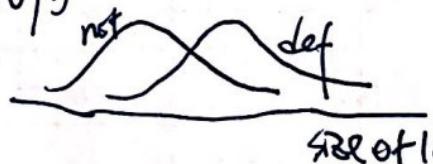
$$P(C|x) = \frac{P(x|C) P(C)}{P(x)} \quad \begin{matrix} \leftarrow \text{Prior} \\ \downarrow \text{posterior} \end{matrix}$$

Example:

$$\begin{aligned} P(\text{not def} | \{9, 5000, \text{own}\}) &\propto P(T \leq 10 | \text{not}) \times P(S \leq 10000 | \text{not}) \\ &\quad \times P(\text{own} | \text{not}) \times P(\text{NOT DEF}) \\ &\propto \frac{4}{6} \times \frac{3}{6} \times \frac{4}{6} \times \frac{6}{10} \\ &= 0.1334 \end{aligned}$$

$$\begin{aligned} P(\text{DEF} | \{9, 5000, \text{own}\}) &\propto P(T \leq 10 | \text{DEF}) \times P(S \leq 10000 | \text{DEF}) \\ &\quad \times P(\text{own} | \text{DEF}) \times P(\text{DEF}) \\ &\propto \frac{3}{4} \times \frac{4}{4} \times \frac{1}{4} \times \frac{4}{10} \\ &= 0.075 \end{aligned}$$

Continuous Attributes : Draw the pdf



Laplace Smoothing

if? zero oft an attribute value we've never seen

Add 1 to the numerator and denominator of our class-conditional Pwb.

$$P(\text{Green} | \text{Apple}) = \frac{1}{4} \left(\frac{1}{4} + \frac{0}{3} \right)$$

$$P(\text{Green} | \text{Banana}) = \frac{0}{1+0} \xrightarrow{\text{Laplace Smoothing}} \frac{1}{2} \quad \begin{matrix} \text{to avoid zeroing the} \\ \text{attribute out.} \end{matrix}$$

Decision Boundaries / Probability Distributions.

Characteristics: Naive Bayes Attributes (independent).

Text Classification

Bag of words model:

D₁: Each state has its own laws.

D₂: Every country has its own culture

Class

D₁

D₂

0 1 1 0 0

- 0 1 -

Example:

$$P(C|X) = \frac{P(X|C) \cdot P(C)}{P(X)}$$

$$P(\text{Spam} | R, H) = P(R|\text{SPAM}) \times P(H|\text{SPAM}) \times P(\text{SPAM})$$

$$\begin{aligned} &= \frac{400}{500} \times \frac{10}{500} \times \frac{500}{1500} \\ &= 0.0053 \end{aligned}$$

$$P(\text{Ham} | R, H) = P(R|\text{Ham}) \times P(H|\text{Ham}) \times P(\text{Ham})$$

$$\begin{aligned} &= \frac{25}{1000} \times \frac{50}{1000} \times \frac{1000}{1500} \\ &= 0.00083 \end{aligned}$$

No class on Thursday

Ensemble Methods: Improve classification accuracy by aggregating the predictions of multiple classifiers.

◆ Requirements: base classifiers better than 50% accurate & All base classifiers be independent. Diverse

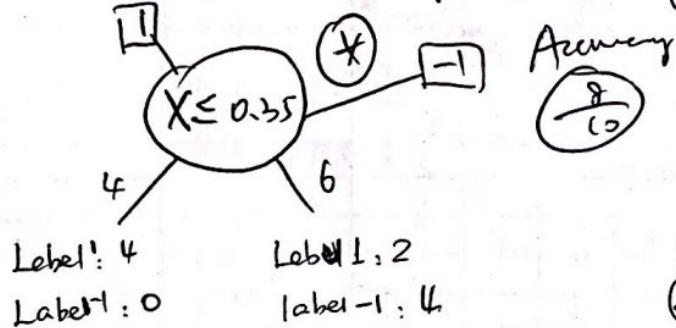
◆ Types of EM: Homogeneous → ensemble classifier
Heterogeneous → voting classifier

◆ Bias and Variance

◆ Methods for constructing an Ensemble Classifier

- By manipulating
 - ① training set (bagging & boosting)
 - ② input features (random features)
 - ③ class labels (multi-class partitioning)
 - ④ learning algorithm

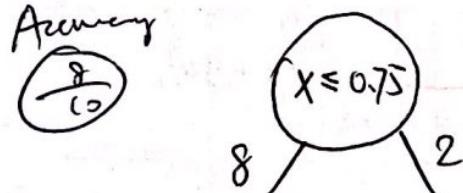
① Bagging: Bootstrap aggregating



$$GINI(Split) = 0 \cdot \frac{4}{10} + 0 \cdot \frac{4}{9} \cdot \frac{6}{10}$$

$$= \frac{24}{90} = \frac{12}{45} = \frac{4}{15}$$

$$= 0.2667$$



$$\begin{aligned} & \text{Left child: } 1 : 4, -1 : 4 \\ & \text{Right child: } 1 : 2, -1 : 0 \\ & GINI = 1 - \left(\frac{4}{8}\right)^2 - \left(\frac{4}{8}\right)^2 \\ & = 1 - \frac{1}{4} - \frac{1}{4} \\ & = \frac{1}{2} \end{aligned}$$

$$GINI(Split) = \frac{4}{10} \cdot \frac{1}{2} + \left(\frac{2}{10}\right) \cdot 0 = 0.4$$

② Random Forests : An extensive extension of bagged decision trees
Randomly select a subset of input features.

① → Boosting : Using weights when doing bootstrap bagging

updated weight according to previous correction

Class Imbalance Problem: Accuracy doesn't work well.

Confusion Matrix (Binary Classification)

		Predicted Class	
		+	-
Actual Class	+	142	8
	-	1	989

Precision and Recall Trade-off

Confusion

F-measure

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

ROC Curves (Receiver Operating Characteristic)

ROC

ACT

60	50
2	17

True Positives = 60

False Positives = 2

True Negatives = 50

False Negatives = 17

TPR = $\frac{60}{60+17} = 0.79$

TNR = $\frac{50}{50+2} = 0.96$

FPR = $\frac{2}{2+50} = 0.04$

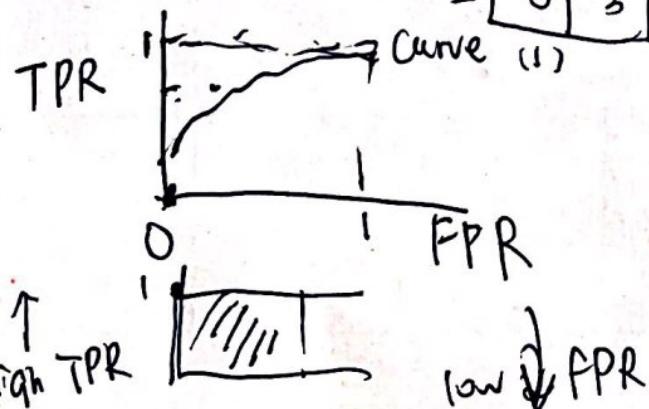
Q: Example

		+	-
+	+	2	1
	-	0	3

$$\begin{aligned} \text{TPR} &= \frac{2}{2+1} = \frac{2}{3} \\ \text{FPR} &= \frac{1}{1+3} = \frac{1}{4} \end{aligned}$$

		+	-
+	+	2	1
	-	1	2

Curve (2)



Mitigating Class Imbalances

- Sampling based approaches:
 - * Undersampling
 - * oversampling

→ Sampling Algorithms - SMOTE: Synthetic Minority Over-Sampling Technique