

Summary

SUMMARY

(SKIPPED "COURSE INTRODUCTION" AND "INTRODUCTION TO DATA SCIENCE")

• DATA PREPROCESSING:

• DATA TYPES AND PROPERTIES:

• ATTRIBUTE TYPES:

- NOMINAL : es. ID numbers, ZIP codes → DISTINCTNESS $\rightsquigarrow \{/\}$
- ORDINAL : es. RANKINGS → DISTINCTNESS + ORDER $\rightsquigarrow \{/\}$
- INTERVAL : es. CALENDAR DATES → DISTINCTNESS, ORDER, AVERAGE $\rightsquigarrow \{/\}$
- RATIO : es. TEMPERATURE, LENGTH, WEIGHT → DISTINCTNESS, ORDER, ADDITION, MULTIPLICATION $\rightsquigarrow \{/\}$
- DISCRETE VS CONTINUOUS ATTRIBUTES

PROPERTIES :

- DISTINCTNESS $\rightsquigarrow \{/\}$
- DISTINCTNESS + ORDER $\rightsquigarrow \{/\}$
- DISTINCTNESS, ORDER, AVERAGE $\rightsquigarrow \{/\}$
- DISTINCTNESS, ORDER, ADDITION, MULTIPLICATION $\rightsquigarrow \{/\}$

• DATA SET TYPES:

• RECORDS :

STRUCTURED
UNSTRUCTURED
TEXT

- TABLES
- DOCUMENT DATA
- TRANSACTION DATA

• GRAPH:

\rightsquigarrow JSON

- VARIOUS WIDE GRAPHS
- MOLECULAR STRUCTURES

• ORDERED

- SPATIAL DATA
- MOLECULAR DATA
- SEQUENTIAL DATA

\rightsquigarrow TRANSACTIONS

• DATA QUALITY:

- NOISE : MODIFICATION OF ORIGINAL DATA \rightsquigarrow es. RECORDING OF A PERSON'S VOICE
- OUTLIERS : DATA OBJECT WITH CONSIDERABLY DIFFERENT CHARACTERISTICS IN DATA SET

↳ NOISE THAT INTERFERES WITH THE ANALYSIS

↳ OUT OF ANALYSIS \rightsquigarrow es. INTRUSION DETECTION

- MISSING VALUES:

- DUPLICATE DATA

• DATA PREPARATION:

• DATA PREPROCESSING:

- **AGGREGATION**: COMBINING 1 OR MORE ATTRIBUTES INTO A SINGLE ONE
→ DATA REDUCTION, CHANGE OF SCALE, LESS VARIABILITY

- **DATA REDUCTION**: REDUCED REPRESENTATION OF DATASET

→ SAMPLING, FEATURE SELECTION, DISCRETIZATION

• SAMPLING:

- SIMPLE RANDOM SAMPLING: $\forall \text{ item} \rightarrow \text{same } P$
- SAMPLING WITHOUT REPLACEMENT: $\forall \text{ chosen item} \rightarrow \text{removed from population}$
- SAMPLING WITH REPLACEMENT: $\forall \text{ item} \rightarrow \text{can be chosen more than 1 time}$
- STRATIFIED SAMPLING: SPLIT DATA IN PARTITIONS, THEN RANDOM SAMPLING

\hookrightarrow probability

- DIMENSIONALITY: $\uparrow \text{DIMENSIONALITY} \rightarrow \uparrow \text{SPARSITY}$
- DIMENSIONALITY REDUCTION:

- TECHNIQUES: PCA (PRINCIPAL COMPONENT ANALYSIS), SINGULAR VALUE DECOMPOSITION, OTHERS
 \hookrightarrow FIND A PROJECTION THAT UPGRADES LARGEST AMOUNT OF VARIATION IN DATA

• FEATURE SUBSET SELECTION

- FEATURE CREATION: CREATE NEW ATTRIBUTES THAT CAN CAPTURE MORE IMPORTANT INFORMATION
 \hookrightarrow e.g., FOURIER TRANSFORM: $t \mapsto f$

→ FEATURE EXTRACTION, MAPPING TO NEW SPACE, FEATURES CONSTRUCTION

• DISCRETIZATION:

- UNSUPERVISED DISCRETIZATION: FIND BREAKS IN DATA VALUES

- SUPERVISED DISCRETIZATION: USE CLASS LABELS TO FIND BREAKS

- BINARIZATION: MAPPING ATTRIBUTE \mapsto 1 OR MORE BINARY VARIABLES

• ATTRIBUTE TRANSFORMATION:

- NORMALIZATION: VALUES ARE SCALING IN ORDER TO FALL IN A RANGE:

$\textcircled{1}$ MIN-MAX NORMALIZATION: $V' = \frac{V - \min_A}{\max_A - \min_A} (\max_{N,A}^{\text{ANNO}} - \min_{N,A}) + \min_{N,A}$

$Z \sim N(0,1)$

$\textcircled{2}$ Z-SCORE NORMALIZATION: $V' = \frac{V - \mu_A}{\sigma_A} \rightarrow \text{MEAN VALUE}$

$\sigma_A \rightarrow \text{STAND. DEV.}$

• DECIMAL SCALING: $V' = \frac{1}{10^j} \cdot V \quad / \quad j: \text{smallest INT. such that } \max(|V'|) < 1$

• DATA PREPARATION FOR DOCUMENT DATA:

- DOCUMENT REPRESENTED AS A SET OF FEATURES: SET OF WORDS, CHARACTERS, TERM, SENTENCES
- DOCUMENT PROCESSING:

IT GENERATES A STRUCTURED DATA REPRESENTATION OF DOCUMENT DATA

\rightarrow 5 STEPS: DOCUMENT SPLITTING, TOKENISATION, CASE NORMALIZATION, STOPWORD REMOVAL, STEMMING

- DOCUMENT SPLITTING: DOCUMENT CAN BE SPLITTED IN PARAGRAPHS, SENTENCES, WORDS, SENTENCES, ETC.

- TOKENISATION: BREAKING TEXT IN SENTENCES (TOKENS)

- CASE NORMALIZATION: WORDS \mapsto UPPER/LOWER CASE CHAR

- STEMMING: REDUCE A WORD INTO ITS ROOT FORM

\rightarrow IDENTIFICATION OF PREFIXES, SUFFIXES, PLURALIZATION, ETC.

- STOPWORD ELIMINATION: ELIMINATION OF ARTICLES, PREPOSITIONS, ETC.

• WEIGHTED DOCUMENT REPRESENTATION:

- FEATURE VECTORS: DOCUMENT TEXT \mapsto FEATURES

- WEIGHTED SCHEMES:

- BINARY: A WORD $\rightarrow \exists / \nexists$

- FREQUENCY OCCURRENCE

- TERM FREQUENCY INVERSE DOCUMENT FREQUENCY: $t_f \cdot idf(t) = f(t, d) \log\left(\frac{m}{f(t, d)}\right)$

- DOCUMENT-TERM MATRIX X : $X_{ij} = f_{ij} \times g_j$

sd. N OF WORDS \forall DOCUMENT

	team	coach	y	pis	ball	snow	guitar	u	m	pink	purple	blue
Document 1	3	0	5	0	2	6	0	2	0	2	0	2
Document 2	0	7	0	2	1	0	0	3	0	0	0	0
Document 3	0	1	0	0	1	2	2	0	3	0	0	0



TERM t IN DOCUMENT d OF COLLECTION D

frequency

document

local weight

global weight

collection of documents

SIMILARITY AND DISSIMILARITY :

SIMILARITY :

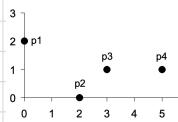
MEASURE OF HOW SIMILAR
2 DATA OBJECT ARE

DISSIMILARITY:

MEASURE OF HOW DIFFERENT
2 DATA OBJECT ARE

Attribute Type	Dissimilarity	Similarity
Nominal	$d = \begin{cases} 0 & \text{if } x = y \\ 1 & \text{if } x \neq y \end{cases}$	$s = \begin{cases} 1 & \text{if } x = y \\ 0 & \text{if } x \neq y \end{cases}$
Ordinal	$d = x - y /(n - 1)$ (values mapped to integers 0 to $n - 1$, where n is the number of values)	$s = 1 - d$
Interval or Ratio	$d = x - y $	$s = -d, s = \frac{1}{1+d}, s = e^{-d}, s = 1 - \frac{d - \min_d}{\max_d - \min_d}$

Ex.



point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

Distance Matrix

EUCLEDIAN DISTANCE :

$$d(x, y) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2}$$

$\rightarrow n$: N of dimensions, x_k, y_k are k^{th} ATTRIBUTES OF DATA OBJECTS X & Y

MINKOWSKI DISTANCE :

$\rightarrow r=2$

it's a generalization of EUCLIDEAN DISTANCE

$$d(x, y) = \left(\sum_{k=1}^n |x_k - y_k|^r \right)^{1/r}$$

$\rightarrow r$: parameter, n : N of dimensions x_k, y_k are k^{th} ATTRIBUTES

OF DATA OBJECTS X & Y

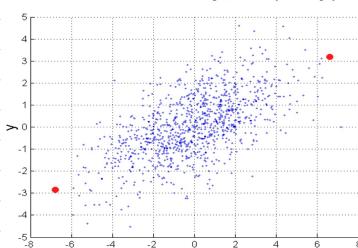
- $r = 1$: L₁ norm, CITY BLOCK DISTANCE \rightarrow MANHATTAN DISTANCE

- $r = 2$: L₂ norm, EUCLIDEAN DISTANCE

- $r = \infty$: L _{∞} norm \rightarrow MAXIMUM DISTANCE BETWEEN ANY COMPONENT OF THE VECTORS

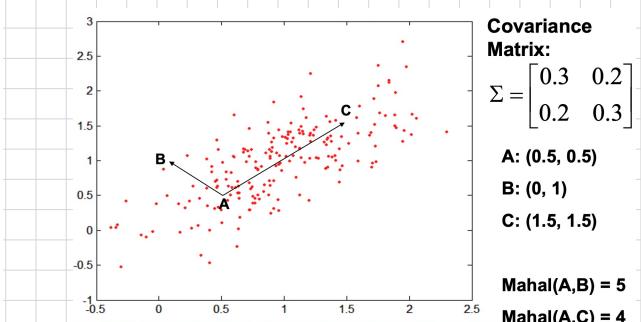
MAHALANOBIS DISTANCE.

$$\text{mahalanobis}(x, y) = (x - y)^T \Sigma^{-1} (x - y)$$



Σ is the covariance matrix

For red points, the Euclidean distance is 14.7, Mahalanobis distance is 6.



Covariance Matrix:

$$\Sigma = \begin{bmatrix} 0.3 & 0.2 \\ 0.2 & 0.3 \end{bmatrix}$$

A: (0.5, 0.5)

B: (0, 1)

C: (1.5, 1.5)

Mahal(A,B) = 5

Mahal(A,C) = 4

SIMILARITY BETWEEN BINARY VECTORS:

- CONSIDER BINARY VECTORS p, q

- COMPUTE SIMILARITY: $M_{00}, M_{01}, M_{10}, M_{11} / M_{x,y} = \begin{cases} n^{\circ} \text{ OF ATTRIBUTES} \\ \text{WHERE } p = X, q = Y \end{cases}$

- SIMPLE MATCHINR AND JACCARD COEFFICIENTS:

- SIMPLE MATCHINR: $SMC = \frac{n^{\circ} \text{ OF MATCHES}}{n^{\circ} \text{ OF ATTRIBUTES}} = \frac{M_{00} + M_{11}}{(M_{00} + M_{01} + M_{10} + M_{11})}$

- J = $\frac{n^{\circ} \text{ OF } 11 \text{ MATCHES}}{n^{\circ} \text{ OF NON-BOTH-ZERO ATTRIBUTES VALUE}} = \frac{M_{11}}{M_{01} + M_{10} + M_{11}} : SMC \quad | M_{00} = 0$

COSINE SIMILARITY:

- LET d_1, d_2 BE 2 DOCUMENT VECTORS:

- $\rightarrow \cos(d_1, d_2) = \frac{d_1 \cdot d_2}{\underbrace{\|d_1\| \cdot \|d_2\|}_{\sqrt{d_1^2}}$

- IN GENERAL:

$$\text{SIMILARITY}(x, y) = \frac{\sum_{k=1}^n w_k s_k(x, y)}{\sum_{k=1}^n w_k s_k}$$

DS.

$$d_1 = 3205000200$$

$$d_2 = 1000000102$$

$$d_1 \cdot d_2 = 3*1 + 2*0 + 0*0 + 5*0 + 0*0 + 0*0 + 2*1 + 0*0 + 0*2 = 5$$

$$\|d_1\| = (3^2 + 2^2 + 0^2 + 5^2 + 0^2 + 0^2 + 2^2 + 0^2 + 0^2)^{0.5} = (42)^{0.5} = 6.481$$

$$\|d_2\| = (1^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 1^2 + 0^2 + 2^2)^{0.5} = (6)^{0.5} = 2.245$$

$$\cos(d_1, d_2) = .3150$$

- $s_k(x, y)$: SIMILARITY FOR k^{th} ATTRIBUTES

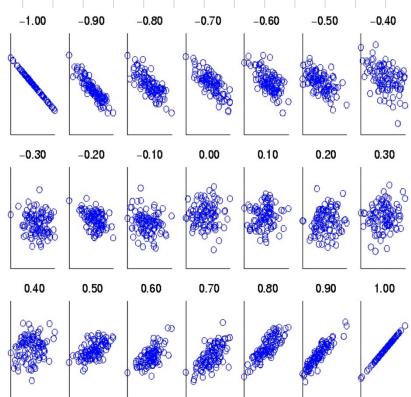
- $s_k = \begin{cases} 0, & \text{IF THE } k^{\text{TH}} \text{ IS ASYMETRICAL AND BOTH OBJECT HAVE VALUE } 0 \\ & \text{OR IF NONE OF THE OBJECTS HAS A MISSING VALUE FOR } k^{\text{TH}} \text{ ATTRIBUTE} \\ 1, & \text{otherwise} \end{cases}$

CORRELATION:

- $\text{Corr}(x, y) = \frac{G_{xy}}{G_x \cdot G_y}$

- $G_{xy} = \frac{1}{m-1} \sum_{k=1}^m (x_k - \bar{x})(y_k - \bar{y})$

- $G_x \text{ (ANALOG FOR } G_y) = \sqrt{\frac{1}{m-1} \sum_{k=1}^m (x_k - \bar{x})^2}$



ASSOCIATION RULES :

- GOAL : EXTRACTION OF FREQUENT CORRELATION OF PATTERNS FROM A TRANSACTIONAL DB

ASSOCIATION RULE MINING :

- COLLECTION OF TRANSACTION DATA

- ASSOCIATION RULE : $A, B \rightarrow C$
↳ NOT CAUSALITY

- IT IS AN EXPLORATORY TECHNIQUE THAT CAN BE APPLIED TO ANY DATA TYPE

DEFINITIONS :

- ITEMSET : SET INCLUDING 1 OR MORE ITEMS

ex. $\{ \text{BEER}, \text{DIAPERS} \}$

- K-ITEMSET : ITEMSET OF LENGTH K, K ITEMS CONTAINED

- SUPPORT COUNT (#) : ABSOLUTE # OF OCCURRENCE OF AN ITEMSET
ex. $\# \{ \text{BEER}, \text{DIAPERS} \} = 2$

- SUPPORT : RELATIVE # OF TRANSACTIONS THAT CONTAINS AN ITEM

ex. $\text{sup}(\{ \text{BEER}, \text{DIAPERS} \}) = \frac{2}{5}$

- FREQUENT ITEMSET : ITEM WHOSE $\text{sup}() > \text{minsup}$

- GIVEN THE ASSOCIATION RULE $A \rightarrow B$:

$$\textcircled{1} \text{ SUPPORT} = \frac{\#\{A, B\}}{|T|} \rightsquigarrow \text{CARTONALITY OR DB}$$

$$\textcircled{2} \text{ CONFIDENCE} = \frac{\text{sup}(\{A, B\})}{\text{sup}(A)} \rightsquigarrow \text{IT'S LIKE A CONDITIONAL PROBABILITY}$$

\rightarrow ASSOCIATION RULES MINING IS THE EXTRACTION OF RULES

Tickets at a supermarket counter

TID	Items
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diapers, Milk
4	Beer, Bread, Diapers, Milk
5	Coke, Diapers, Milk
...	...

Association rule

- diapers \Rightarrow beer
- 2% of transactions contains both items
- 30% of transactions containing diapers also contains beer

TID	Items
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diapers, Milk
4	Beer, Bread, Diapers, Milk
5	Coke, Diapers, Milk
...	...

ex.

From itemset {Milk, Diapers} the following rules may be derived

- Rule: Milk \Rightarrow Diapers
 - support
 - $\text{sup} = \#\{\text{Milk, Diapers}\}/\#\text{trans.} = 3/5 = 60\%$
 - confidence
 - $\text{conf} = \#\{\text{Milk, Diapers}\}/\#\{\text{Milk}\} = 3/4 = 75\%$
- Rule: Diapers \Rightarrow Milk
 - same support
 - $s = 60\%$
 - confidence
 - $\text{conf} = \#\{\text{Milk, Diapers}\}/\#\{\text{Diapers}\} = 3/3 = 100\%$

TID	Items
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diapers, Milk
4	Beer, Bread, Diapers, Milk
5	Coke, Diapers, Milk

$\left. \begin{array}{l} \text{SUPPORT} \geq \text{minsup} \\ \text{CONFIDENCE} \geq \text{minconf} \end{array} \right\}$

- BRUTE FORCE APPROACH:

- ENUMERATE ALL POSSIBLE PERMUTATION
- COMPUTE SUPPORT AND CONFIDENCE
- PRUNE BASED ON minsup & minconf

} → GIVEN d ITEMS → 2^d POSSIBLE ITEMSETS
 COMPUTATION VERY EXPENSIVE

② APRIORI ALGORITHM:

- APRIORI PRINCIPLE: "IF AN ITEMSET IS FREQUENT → ALL OF HIS SUBSETS MUST BE FREQUENT"

$$\rightarrow \text{IF } A \subseteq B \rightarrow \text{sup}(A) \geq \text{sup}(B)$$

→ IF ITEMSET S IS NOT FREQUENT → ALL SUBSETS GENERATED FROM S CAN BE DISCARDED

AGR '94 ALGORITHM:

PSEUDO CODE

- \forall ITERATION → EXTRACT ITEMSETS OF LENGTH K
- 2 STEPS:

1. CANDIDATE GENERATION: (1^{st} ITER; $K=1$)

1.1 GENERATE CANDIDATES OF LENGTH $K+1$ FROM FREQUENT ITEMSETS OF LENGTH K

1.2 PRUNE LENGTH $K+1$ ITEMSETS / CONTAIN AT MOST 1 K-ITEMSET THAT ARE NOT FREQUENT

2. FREQUENT ITEM GENERATION

2.1 SCAN DB TO GET SUPPORT FOR $K+1$ CANDIDATES

2.2 PRUNE CANDIDATES / $\text{sup}() \leq \text{minsup}$

so. SLIDES $21 \div 35$

FACTORS AFFECTING PERFORMANCE:

LOW MINIMUM SUPPORT THRESHOLD, DATASET DIMENSIONS, SIZE OF DB, AVG TRANSACTION WIDTH

IMPROVING APRIORI EFFICIENCY:

- 1995: HASH BASED ITEMSET COUNTING, TRANSACTION REJECTION
- 1996: PARTITIONING
- 1996: SAVING
- 1998: DYNAMIC ITEMSET COUNTING

C_k : Candidate itemset of size k
 L_k : frequent itemset of size k

```

 $L_f = \{\text{frequent items}\};$ 
for ( $k = 1$ ;  $L_k \neq \emptyset$ ;  $k++$ ) do
  begin
     $C_{k+1} = \{\text{candidates generated from } L_k\};$ 
    for each transaction  $t$  in database do
      increment the count of all candidates in  $C_{k+1}$ 
      that are contained in  $t$ 
     $L_{k+1} = \{\text{candidates in } C_{k+1} \text{ satisfying minsup}\}$ 
  end
return  $\cup_k L_k$ 

```

• FP - GROWTH Algorithm :

EXPLOIT A MAIN MEMORY COMPRESSED REPRESENTATION OF DB : THE FP-TREE

- FREQUENT PARTITION MIMING BY FP-Tree
 - ONLY 2 DB SCANS: COUNT ITYM SUPPORTS + BUILT FP-TREE
 - MAN 00 ALGORITHM:

- FP-TREE CONSTRUCTION: \sim SIMILAR TO HAMMING CODES TREE

1. COUNT ITEM SUPPORT AND PRUNE ITEMS WHERE SUP(.) <

1. COUNT ITEM SUPPORT AND PRUNE ITEMS WHERE $\text{sup}() \leq \text{minsup}$
 2. BUILD HEADER TABLE, SORTING ITEMS BY DECREASING $\text{sup}()$ ORDER ^{N²}
 3. CREATE FP-TREES:

HEADER :
 { B } -> 8
 { A } -> 7
 :
 :

\forall transaction t in DB :

- 3.7 ORDER TRANSACTION (IT IS IN DECREASING SUP() ORDER

- 3.2 INSERT E IN FP-TREE : USE COMMON PATH FOR COMMON PREFIX,

CREATE NEW BRANCH WITH PATH BECOMES DIFFERENT

• FP - Growth ALGORITHM :

1. SCAN MEADER TABLE FROM LOWEST SUPPORT ITEM

IN MEASURE TABLE

✓ ITEM IN HIGHER TABLE, EXTRACT FREQUENT ITEMSETS INCLUDING I AND ITEMS PRECEDING IT

- ## 2.1 BUILD CONDITIONAL PATTERN BASE FOR ITEM $i \rightarrow i - CPB$

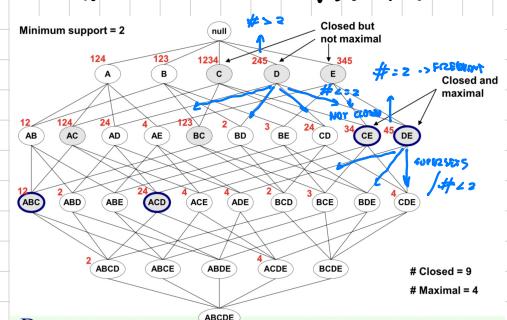
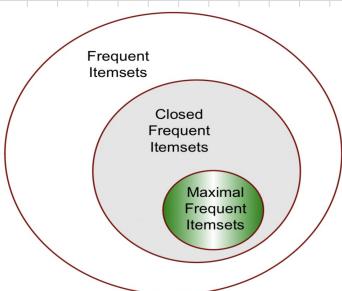
- SELECT PREFIX-PATHS OF ITEM 1 FROM FP-TREE

- ## 2.2 RECURSIVE INVOCATION OF FP-GROWTH ON i-CPB

20. SHOES $46 \div 53$ AND $56 \div 73$

④ MAXIMAL FREQUENT ITEMSET: AN ITEM IS FREQUENT MAXIMUS IF NONE OF ITS SUBSETS IS FREQUENT

• CLOSER ITEMSET : ITEMSET IS CLOSER IF NONE OF ITS IMMEDIATE SUPERSETS HAS THE



• EFFECT OF SUPPORT THRESHOLD:

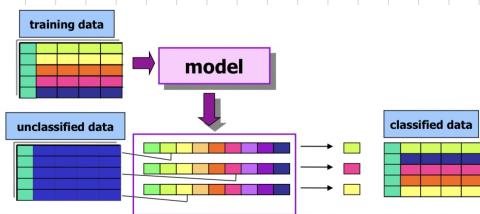
- TOO HIGH minsup: POSSIBLE LOSS FOR FAKE BUT INTERESTING ITEMSETS
- TOO LOW minsup: POSSIBLY TO BECOME COMPUTATIONALLY EXPENSIVE,
LARGE AMOUNT OF FREQUENT ITEMSETS
- CONFIDENCE MEASURE NOT ALWAYS RELIABLE. \rightarrow CORRELATION IS BETTER

$$\text{CORRELATION} = \frac{P(A, B)}{P(A) P(B)} = \frac{\text{CONF}(A \rightarrow B)}{\text{sup}(B)} \rightarrow \begin{cases} \text{CORR.} = 1 \rightarrow \text{STR. INDP.} \\ \text{CORR.} > 1 \rightarrow \text{POSITIVE CORR.} \\ \text{CORR.} < 1 \rightarrow \text{NEGATIVE CORR.} \end{cases}$$

- CONSIDERING ITEMSETS WEIGHT
- CONSIDERING HIERARCHIES \rightarrow AGGREGATE DATA: TAXONOMY
 \rightarrow GENERALIZING DATA CAN BE MORE USEFUL / MEANINGFUL SOMETIMES
 ↳ OR COMBINE DETAILED AND GENERALIZED DATA

CLASSIFICATION FUNDAMENTALS:

GOAL: PREDICTION OF A CLASS LABEL OF UNPREDICTED DATA
GIVEN TRAINING DATA



DEFINITIONS:

- TRAINING SET: COLLECTION OF LABELED DATA OBJECTS, USED TO LEARN THE CLASSIFICATION MODEL
- TEST SET: COLLECTION OF LABELED DATA OBJECTS TO VALIDATE THE CLASSIFICATION MODEL

CLASSIFICATION TECHNIQUES:

- Decision trees
- Classification rules
- Association rules
- Neural Networks
- Naïve Bayes and Bayesian Networks
- k-Nearest Neighbours (k-NN)
- Support Vector Machines (SVM)
- ...

EVALUATION OF CLASSIFICATION TECHNIQUES:

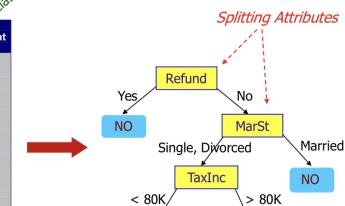
- Accuracy
 - quality of the prediction
- Interpretability
 - model interpretability
 - model compactness
- Incrementality
 - model update in presence of newly labelled record

- Efficiency
 - model building time
 - classification time
- Scalability
 - training set size
 - attribute number
- Robustness
 - noise, missing data

DECISION TREES:

- MANY ALGORITHMS: HUNT'S ALGORITHM, CART, ID3, C4.5, C5.0, SLIQ, SPRINT

Tid	Refund	Marital Status	Taxable Income	Cheat	class
1	Yes	Single	125K	No	
2	No	Married	100K	No	
3	No	Single	70K	No	
4	Yes	Married	120K	No	
5	No	Divorced	95K	Yes	
6	No	Married	60K	No	
7	Yes	Divorced	220K	No	
8	No	Single	85K	Yes	
9	No	Married	75K	No	
10	No	Single	90K	Yes	



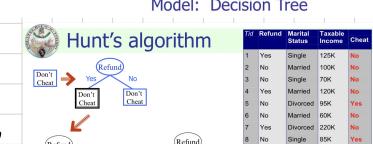
Training Data

HUNT'S ALGORITHM (IN BRIEF):

- GIVEN D_t : SET OF TRAINING RECORDS THAT REACH NODE t

STEPS:

- IF D_t CONTAINS RECORDS THAT BELONG TO MORE THAN 1 CLASS:
 - SELECT BEST ATTRIBUTE A ON WHICH TO SPLIT D_t AND CALL NODE t AS A
 - SPLIT D_t IN SMALLER SUBSETS AND RECURSIVELY APPLY THIS PROCEDURE: \forall subset
- IF D_t CONTAINS RECORDS THAT BELONG TO SAME CLASS $V_t \rightarrow t$ IS A LEAF NODE, LABELED AS THE MAJORITY CLASS V_t
- IF D_t IS AN EMPTY SET: t IS A LEAF NODE, LABELED AS THE MAJORITY CLASS V_d



Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

• DECISION TREE INDUCTION:

IT ADOPTS A GREEDY STRATEGY → BEST ATTRIBUTE IS SELECTED LOCALLY AT EACH STEP, IT'S NOT A GLOBAL

BINARY VS MULTIVALENT SPLIT OPTIMUM

↑

• ISSUES: STRUCTURES OF TEST CONDITIONS, BEST ATTRIBUTE SELECTION, STOPPING CONDITION

• SPLITTING ON COMMONS ATTRIBUTES :

• DISCRETIZATION: STATIC (ONCE AT BEGINNING) OR DYNAMIC (DURING TREES INDUCTION)

• BINARY DECISION: CONSIDER ALL POSSIBLE SPLITS AND FIND BEST CUT ($A < v$ or $A \geq v$)

• BEST ATTRIBUTE:

THE ONE WITH MOST HOMOGENOUS CLASS DISTRIBUTION, LESS NODE IMPURITY

• MEASURE OF NODE IMPURITY : GINI INDEX, ENTROPY, MISCLASSIFICATION ERROR

GINI INDEX:

$$GINI(t) = 1 - \sum_j [p(j|t)]^2 \quad / \quad p(j|t) \text{ IS THE RELATIVE FREQUENCY OF CLASS } j \text{ AT NODE } t$$

LEAF POUND

• $GINI(t) \in [0, 1 - \frac{1}{m_c}]$ → $\begin{cases} 0: \text{ALL RECORDS BELONG TO 1 CLASS} \\ 1 - \frac{1}{m_c}: \text{ALL RECORDS EQUALLY DISTRIBUTED} \end{cases}$ $\rightarrow m^o \text{ OF RECORDS}$

C1	0
C2	6

$$\cdot \begin{cases} P(C_1) = \frac{0}{6} = 0 \\ P(C_2) = \frac{6}{6} = 1 \end{cases} \rightarrow GINI = 1 - P(C_1)^2 - P(C_2)^2 = 0$$

C1	1
C2	5

$$\cdot \begin{cases} P(C_1) = \frac{1}{6} \\ P(C_2) = \frac{5}{6} \end{cases} \rightarrow GINI = 1 - \left(\frac{1}{6}\right)^2 - \left(\frac{5}{6}\right)^2 = 0.278$$

① SPLITTING BASED ON GINI:

GINI CAN BE USED AS A MEASURE FOR THE SPLITTING QUALITY

- WHEN A NODE P IS SPLITTED IN K PARTITIONS, QUALITY OF THE SPLIT:

$$GINI_{\text{split}} = \sum_{i=1}^K \frac{n_i}{n} GINI(i) \quad / \quad n_i : \text{n}^{\circ} \text{ OF RECORDS AT CHILD } i \\ n : \text{n}^{\circ} \text{ OF RECORDS AT NODE } p$$

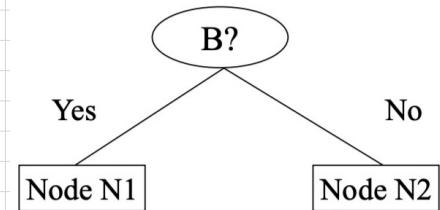
ex.

$$GINI(N_1) = 1 - \left(\frac{5}{7}\right)^2 - \left(\frac{2}{7}\right)^2 = 0.408$$

$$GINI(N_2) = 1 - \left(\frac{1}{5}\right)^2 - \left(\frac{4}{5}\right)^2 = 0.32$$

$$\rightarrow GINI(B \text{ split}) = \frac{7}{12} \cdot 0.408 + \frac{5}{12} \cdot 0.32 = 0.371$$

	Parent	
C1	6	
C2	6	
GINI = 0.500		
	N1	N2
C1	5	1
C2	2	4
GINI = ?		



- CATEGORICAL ATTRIBUTES:

COMPUTE MULTIWAY SPLIT

CarType		
	Family	Sports
C1	1	2
C2	4	1
GINI	0.393	

Multi-way split

Two-way split
(find best partition of values)

CarType		CarType	
(Sports, Luxury)	(Family)	(Sports, Family)	(Luxury)
C1	3	1	2
C2	2	4	1
GINI	0.400	0.419	

- CONTINUOUS ATTRIBUTES:

- BINARY DECISION ON 1 SPLITTING VALUE

- ALL POSSIBLE SPLITTING ON V: $A < V \wedge A > V$

Taxable Income													
→	60	70	75	85	87	92	90	95	97	100	120	125	220
→	55	65	72	80	87	92	90	95	97	100	122	172	230
Yes	0	3	0	3	0	3	1	2	2	1	3	0	3
No	0	7	1	6	2	5	3	4	3	4	3	5	2
GINI	0.420	0.400	0.375	0.343	0.417	0.400	0.300	0.343	0.375	0.400	0.420		

ENTROPY :

$$ENTROPY(t) = - \sum_j p(j|t) \log_2 p(j|t) \quad / \quad ENTROPY(t) \in [0; \log_2 n_c]$$

INFORMATION GAIN:

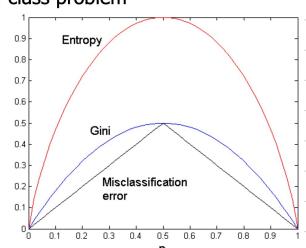
$$GAIN_{\text{split}} = ENTROPY(p) - \sum_{i=1}^k \frac{n_i}{n} ENTROPY(i)$$

\rightarrow Σ ; TENDS TO PREFER SPLITS WITH LARGE NUMBER OF SMALL PARTITIONS \rightarrow BIAS GAIN RATIO

GAIN RATIO:

$$GAIN RATIO_{\text{split}} = - \sum_{i=1}^k \frac{n_i}{n} \log_2 \frac{n_i}{n}$$

For a 2-class problem



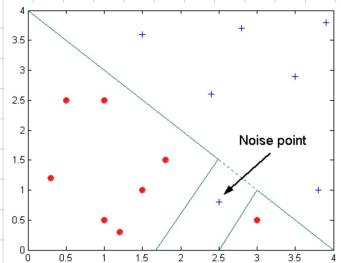
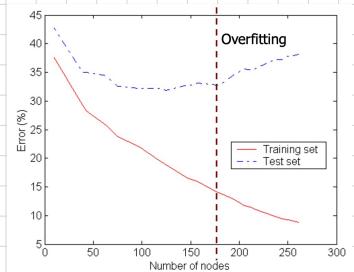
STOPPING CRITERIA:

- ALL RECORDS BELONG TO THE SAME CLASS
- ALL RECORDS HAVE SIMILAR ATTRIBUTE VALUES

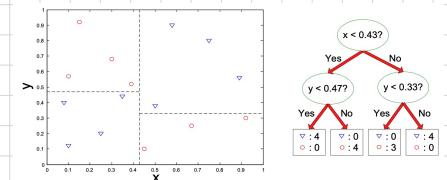
UNDERFITTING AND OVERFITTING:

- UNDERFITTING: WHEN MODEL IS TOO SIMPLE, BOTH TRAINING AND TEST RECORDS ARE WELT
- OVERFITTING: IT CAN BE CAUSED BY NOISE POINTS, DECISION BOUNDARY DISTORTIONS

→ OVERFITTING SOLUTION: PRE-PRUNING → EARLY STOPPING RULE
POST-PRUNING



- DATA FRAGMENTATION: NUMBER OF INSTANCES AT LEAF NODE COULD BE TOO SMALL TO MAKE ANY STATISTICAL SIGNIFICANT PRECISION
- MISSING ATTRIBUTE VALUES
- DECISION BOUNDARY: PARALLEL TO AXES, ∵ TEST CONDITION INVOLVES SINGLE ATTRIBUTE AT A TIME

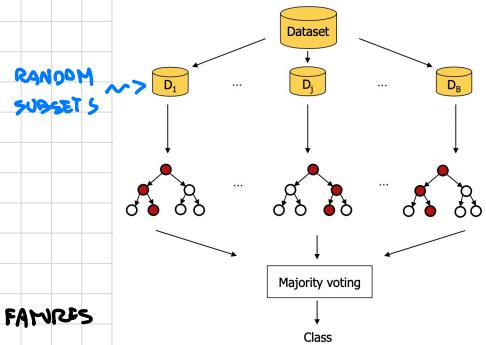


EVALUATION OF DECISION TREES:

- ACCURACY: COMPARES TO OTHER CLASSIFICATION TECHNIQUES FOR SIMPLE DATASETS
- INTERPRETABILITY: ONLY FOR SMALL TREES, SINCE PROGRAM ARE INTERPRETABLE
- INCREMENTALITY: NOT INCREMENTAL
- EFFICIENCY: FAST MODEL BUILDING, VERY FAST CLASSIFICATION
- SCALABILITY: SCALABLE BOTH IN TRAINING SET SIZE AND ATTRIBUTE NUMBER
- ROBUSTNESS: DIFFICULT HANDLING OF MISSING DATA

• RANDOM FOREST :

- MULTIPLE BASE MODELS ARE COMBINED
- > A NUMBER OF DECISION TREES BUILT AT SAME TIME
- CLASS ASSIGNMENTS BY MAJORITY VOTING
- \forall SUBSET \rightarrow TREE IS CREATED ON A RANDOM SET OF FEATURES



• BOOTSTRAP AGGREGATION :

GIVEN A TRAINING SET D OF m INSTANCES, IT SELECTS B TIMES A RANDOM SAMPLE (WITH REPLACEMENT) FROM D AND TRAIN TREES

- FOR $b: 1, \dots, B$: DATASET SUBSET D_b GENERATED, CLASSIFICATION TREE ON D_b
- FEATURE BAGGING :
 \forall CONDITION SPLIT, SELECT A RANDOM SUBSET OF THE FEATURES \rightarrow IF $P = n^0$ FEATURES

Typically \sqrt{P} ,

• N.B.: <https://data36.com/random-forest-in-python/>

- DATASET : IF N TREES \rightarrow N DATASET SUBSETS (WITH REPLACEMENT)
- \forall TREE : 1 SUBSET + ALL FEATURES CONSIDERED
↳ \forall SPLIT \rightarrow A RANDOM SUBSET OF FEATURE CHOOSEN

• EVALUATION :

- ↑ • ACCURACY : HIGHER THAN DECISION TREES
- INTERPRETABILITY : MODEL AND PREDICTION NOT INTERPRETABLE
- INCREMENTALITY : NOT INCREMENTAL
- EFFICIENCY : FAST MODEL BUILDING, VERY FAST CLASSIFICATION
- SCALABILITY : SCALABLE BOTH IN TRAINING SET SIZE AND ATTRIBUTE NUMBER
- ROBUSTNESS : ROBUST TO NOISE AND OUTLIERS

RULE-BASED CLASSIFICATION:

CATEGORIZE RECORDS BY USING: (CONDITION) → CLASS

Ex.

- (Blood Type=Warm) ∧ (Lay Eggs=Yes) → Birds
- (Taxable Income < 50K) ∧ (Refund=Yes) → Cheat=No

- RULE \vdash "COVERS" AN INSTANCE X : ATTRIBUTE OF INSTANCE SATISFY CONDITION OF THE RULE
- RULES CHARACTERISTICS.

- MUTUALLY EXCLUSIVE: 2 RULE CONDITIONS CAN'T BE TRUE AT THE SAME TIME

✓ RECORD → 1 RULE AT MOST EVERY RECORD IS COVERED BY AT MOST 1 RULE

- EXHAUSTIVE: RULES ACCOUNT \forall COMBINATIONS OF ATTRIBUTE VALUES

✓ RECORD → 1 rule at least \forall RECORDS → COVERED BY AT LEAST 1 RULE

- DECISION TREES CAN BE TRANSFORMED INTO RULES

- SOMETIMES RULES CAN BE SIMPLIFIED

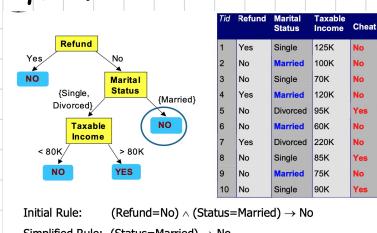
↳ RULES ARE NO LONGER MUTUALLY EXCLUSIVE

↳ SQL: ORDERED SET RULES

↳ RULES NO LONGER EXHAUSTIVE

↳ RECORDS MAY NOT TRIGGER ANY RULES

↳ SQL: USE A DEFAULT CLASS



BUILDING:

- DIRECT METHOD**: RULES FROM DATA

- INDIRECT METHOD**: RULES FROM OTHER CLASSIFICATION MODELS

EVALUATION:

- ACCURACY: HIGHER THAN DECISION TREES
- INTERPRETABILITY: MODEL AND PREDICTION INTERPRETABLE
- INCREMENTALITY: NOT INCREMENTAL
- EFFICIENCY: FAST MODEL BUILDING, VERY FAST CLASSIFICATION
- SCALABILITY: SCALABLE BOTH IN TRAINING SET SIZE AND ATTRIBUTE NUMBER
- ROBUSTNESS: ROBUST TO OUTLIERS

• ASSOCIATIVE CLASSIFIERS :

$$(\text{CONDITION}) \rightarrow \vee$$

• MODEL GENERATION :

- RULE SELECTING AND SORTING : BASED ON SUPPORT, CONFIDENCE AND CORRELATION THRESHOLDS

• RULE PRUNING : TRAINING SET COVERED BY SELECTING TOP MOST RULES ACCORDING TO PREVIOUS SORT

• EVALUATION :

- ↑ • ACCURACY : HIGHER THAN DECISION TREES AND RULE-BASED CLASSIFIERS
 ↗ CORRELATION IS CONSIDERED TOO
- ④ • INTERPRETABILITY : MODEL AND PREDICTION INTERPRETABLE
- INCREMENTALITY : NOT INCREMENTAL
- EFFICIENCY : RULE GENERATION MAY BE SLOW (e.g. SUPPORT VECTOR) , FAST CLASSIFICATION
- SCALABILITY : SCALABLE BOTH IN TRAINING SET SIZE BUT REQUIRES SCALABILITY IN ATTRIBUTE NUMBER
- ④ • ROBUSTNESS : ROBUST TO OUTLIERS, UNFFECTED BY MISSING DATA

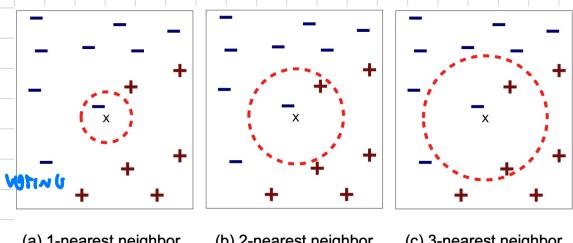
• K-NEAREST NEIGHBOR (KNN) :

USE TRAINING RECORDS TO PREDICT THE CLASS LABEL OF UNSEEN LABELS

• INSTANCE BASED CLASSIFIER ↗ NOTE LEARNER : REMEMBERS ALL TRAINING DATA, CLASSIFIES ONLY IF EXACT MATCH
NEAREST NEIGHBOR : USES K "CLOSEST" POINT

• CLASSIFICATION :

1. COMPUTE DISTANCE TO OTHER TRAINING RECORDS
2. IDENTIFY K NEAREST NEIGHBORS ↗
3. ASSIGN CLASS ACCORDING TO K NEAREST ↗ MAJORITY VOTING



• K=1 : VORONOI DIAGRAMS

• DISTANCE :

$$\text{ex. EUCLIDEAN: } d(p, q) = \sqrt{\sum_i (p_i - q_i)^2}$$

→ WEIGHT VOTE OR DISTANCE : $W = \frac{1}{d^2}$

- K
 - too small: sensitive to noise
 - too large: may include points from other classes
- SCALING ISSUE: ATTRIBUTES SHOULD BE NORMALIZED
 \hookrightarrow e.g. WEIGHT $\in [1.5 \text{ m} ; 2 \text{ m}]$, INCOME $\in [10 \cdot 10^3 ; 10 \cdot 10^6]$

EVALUATION:

- ACCURACY: COMPARABLE TO OTHER TECHNIQUES FOR SMALL DATASETS
- INTERPRETABILITY: MODEL AND PREDICTION NOT INTERPRETABLE
- INCREMENTALITY: INCREMENTAL, TRAINING SET MUST BE AVAILABLE
- EFFICIENCY: (ALMOST) NO MODEL BUILDING, SLOW FOR CLASSIFICATION (COMPLEX)
- SCALABILITY: WEAKLY SCALABLE IN TRAINING SET SIZE
- ROBUSTNESS: DEPENDS ON DISTANCE COMPUTATION

BAYESIAN CLASSIFICATION:

- BAYES THEOREM:

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

C: ANY CLASS LABEL
 X = $\langle x_1, \dots, x_n \rangle$ RECORD TO BE CLASSIFIED
 $\uparrow \Rightarrow = 1$

- BAYESIAN CLASSIFICATION:

1. \forall CLASS \rightarrow COMPUTE $P(C|X)$: p that X belongs to C
2. ASSIGN X TO CLASS WITH MAX $P(C|X)$

- suppose $\langle x_1, \dots, x_n \rangle$ STAT. INDIP.: $P(x_1, \dots, x_n | C) = P(x_1|C) \cdot \dots \cdot P(x_n|C)$

20.

Outlook	Temperature	Humidity	Windy	Class
sunny	hot	high	false	N
sunny	hot	high	true	N
overcast	hot	high	false	P
rain	mild	high	false	P
rain	cool	normal	false	P
rain	cool	normal	true	N
overcast	cool	normal	true	P
sunny	mild	high	false	N
sunny	cool	normal	false	P
rain	mild	normal	false	P
sunny	mild	normal	true	P
overcast	mild	high	true	P
overcast	hot	normal	false	P
rain	mild	high	true	N

outlook	
$P(\text{sunny} p) = 2/9$	$P(\text{sunny} n) = 3/5$
$P(\text{overcast} p) = 4/9$	$P(\text{overcast} n) = 0$
$P(\text{rain} p) = 3/9$	$P(\text{rain} n) = 2/5$
temperature	
$P(\text{hot} p) = 2/9$	$P(\text{hot} n) = 2/5$
$P(\text{mild} p) = 4/9$	$P(\text{mild} n) = 2/5$
$P(\text{cool} p) = 3/9$	$P(\text{cool} n) = 1/5$
humidity	
$P(\text{high} p) = 3/9$	$P(\text{high} n) = 4/5$
$P(\text{normal} p) = 6/9$	$P(\text{normal} n) = 2/5$
windy	
$P(\text{true} p) = 3/9$	$P(\text{true} n) = 3/5$
$P(\text{false} p) = 6/9$	$P(\text{false} n) = 2/5$

$$\begin{aligned} P(p) &= 9/14 \\ P(n) &= 5/14 \end{aligned}$$

$$P(C|\bar{x}) = P(\bar{x}|C)P(C)$$

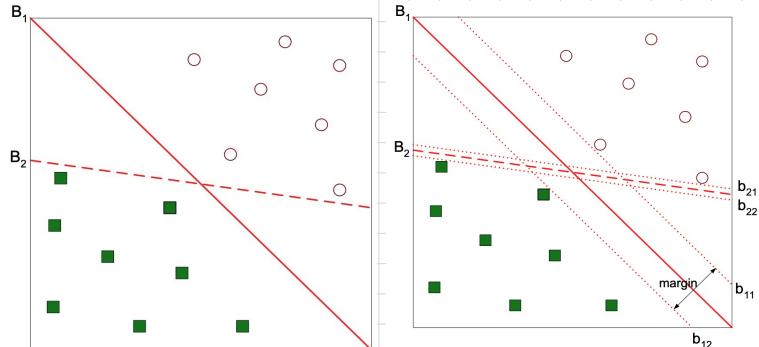
- Data to be labeled
 $X = \langle \text{rain, hot, high, false} \rangle$
- For class p
 $P(X|p) \cdot P(p) =$
 $= P(\text{rain}|p) \cdot P(\text{hot}|p) \cdot P(\text{high}|p) \cdot P(\text{false}|p) \cdot P(p) = 3/9 \cdot 2/9 \cdot 3/9 \cdot 6/9 \cdot 9/14 = 0.010582$
- For class n
 $P(X|n) \cdot P(n) =$
 $= P(\text{rain}|n) \cdot P(\text{hot}|n) \cdot P(\text{high}|n) \cdot P(\text{false}|n) \cdot P(n) = 2/5 \cdot 2/5 \cdot 4/5 \cdot 2/5 \cdot 5/14 = 0.018286$

EVALUATION :

- ACCURACY : SIMILAR / LOWER TO DECISION TREES
- INTERPRETABILITY : MODEL AND PREDICTION NOT INTERPRETABLE
- INCREMENTALITY : FULLY INCREMENTAL
- EFFICIENCY : FAST MODEL BUILDING, VERY FAST CLASSIFICATION
- SCALABILITY : SCALABLE IN TRAINING SET SIZE AND ATTRIBUTE NUMBER
- ROBUSTNESS : AFFECTED BY ATTRIBUTE CORRELATION

SUPPORT VECTOR MACHINES :

- FIND A LINEAR DECISION BOUNDARY THAT WILL SEPARATE THE DATA
- B_1 OR B_2 ?
 → THE ONE WHICH MAXIMIZES MARGIN
- BOUNDARY NOT LINEAR ? → TRANSFORM DATA INTO OTHERS DIMENSIONS



EVALUATION :

- ACCURACY : AMONG BEST PERFORMERS
- INTERPRETABILITY : MODEL AND PREDICTION NOT INTERPRETABLE
- INCREMENTALITY : NOT INCREMENTAL
- EFFICIENCY : MODEL BUILDING REQUIRES SIGNIFICANT PARAMETER TUNING, VERY FAST CLASSIFICATION
- SCALABILITY : MODERATE SCALABLE IN BOTH TRAINING SET SIZE AND ATTRIBUTE NUMBER
- ROBUSTNESS : ROBUST TO ATTRIBUTE CORRELATION

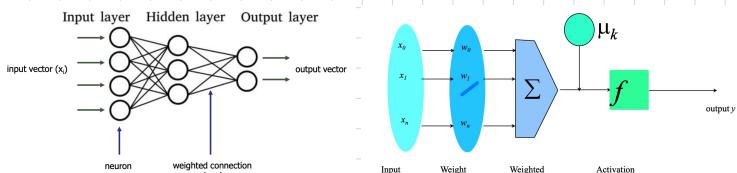
ARTIFICIAL NEURAL NETWORKS:

ACTIVATION FUNCTIONS:

SIMULATES ACTIVATION ON INPUT STIMULUS

NEURON OUTPUT IN FIXED DURATION,

REGULATE IF NEXT NEURON IS ACTIVATED OR NOT

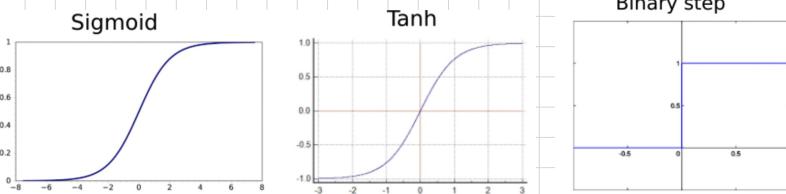


$$\text{softmax}(z_j) = \frac{e^{z_j}}{\sum_{i=0}^{N-1} e^{z_i}}$$

SOFTMAX

USED IN OUTPUT LAYER

GENERATE A FINAL PMF



BUILDING A FFNN: ~> FEED FORWARD NN

• ∀ NODE : $\{w_1, \dots, w_n\}$ WEIGHTS AND OFFSET VALUES

• ALGORITHM:

• INITIALLY: RANDOM VALUES w_i

• INSTANCES PROCESSING:

• ∀ NEURON: COMPUTE RESULT $\leftarrow \{w_1, \dots, w_n\}$, OFFSET, ACTIVATION FUNCTION

• FORWARD PROPAGATION UNTIL OUTPUT COMPUTED

• COMPARE OUTPUT WITH TRUE VALUES

• BACK PROPAGATION

\rightarrow IT STOPS $\leftarrow (ACCURACY > \epsilon_{THRESHOLD})$, ($\% \text{ ERROR ON PARTITION} > \epsilon$), MAX^{n^b} EPOCH

EVALUATION:

① ACCURACY: AMONG BEST PERFORMERS

• INTERPRETABILITY: MODEL AND PREDICTION NOT INTERPRETABLE (BLACK BOX)

• INCREMENTALITY: NOT INCREMENTAL

• EFFICIENCY: MODEL BUILDING REQUIRES SIGNIFICANT PARAMETER TUNING, VERY FAST CLASSIFICATION

• SCALABILITY: MODELS SCALABLE IN BOTH TRAINING SET SIZE AND ATTRIBUTES NUMBER

② ROBUSTNESS: ROBUST TO NOISE AND OUTLIERS, REQUIRE LARGE TRAINING SET

CNN: ~ convolutional NN

AUTOMATIC EXTRATION OF FEATURES

FROM IMAGES AND PERFORM CLASSIFICATION

CONVOLUTIONAL LAYER:

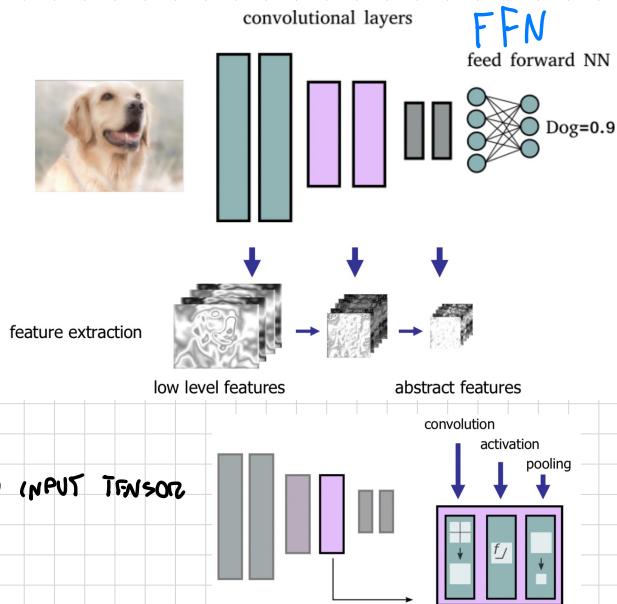
- CONVOLUTION: FEATURES EXTRACTION WITH SLIDING FILTERS
- ACTIVATION: APPLY ACTIVATION FUNCTIONS TO INPUT TENSOR
- POOLING: TENSOR DOWNSAMPLING

TENSORS:

N-DIMENSIONAL VECTOR

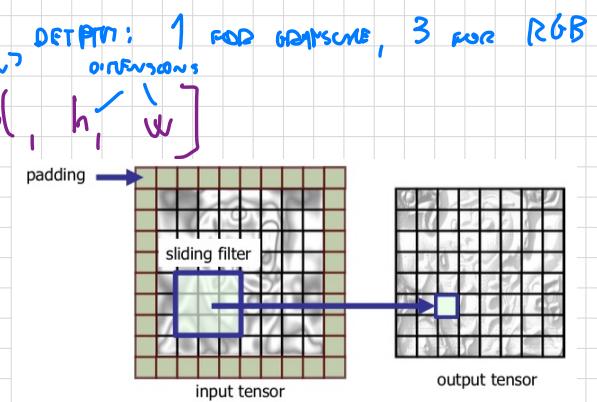
→ THEY FLOW THROUGH CNN LAYERS

so, IMAGES: RANK-3 TENSOR WITH



CONVOLUTION:

INPUT IMAGE/TENSOR $\xrightarrow{\quad}$ OUTPUT TENSOR WITH EXTRACTED FEATURES



• SLIDING FILTER CONTAIN THE TRAINABLE WEIGHTS OF NN

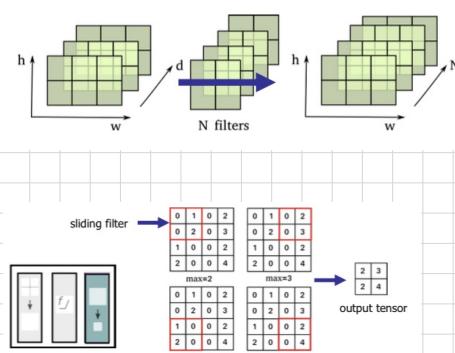
• EACH CONVOLUTIONAL LAYER → HUNDREDS OF FILTERS

$$\hookrightarrow [d, b, w] \mapsto [N, h, w]$$

POOLING:

• SLIDING FILTERS: REUSABLE TENSOR WITH SUMMARIZED STATISTICS OF NEARBY OUTPUTS

so, MAXPOOL : COMPUTE MAX AS STATISTIC (MOST COMMON)



- DURING TRAINING, EACH FILTER LEARN TO RECOGNIZE A PARTICULAR PATTERN

- SHALLOW LAYERS AND DEEPER LAYER

- SEMANTIC SEGMENTATION CNN'S :

- \forall PIXEL OF IMAGE \rightarrow CLASS
- ENCODER NETWORK + DECODER NETWORK

- RECURRENT NN :

SEQUENTIAL DATA PROCESSING $X(t)$

\rightarrow THEY KEEP A STATE AT EACH t

• INPUT . $x(t)$, PREVIOUS STATE $s(t-1) \rightarrow$ OUTPUT $v(t)$

• ERROR $v(t) - x(t)$ PROPAGATES \forall STEP

\hookrightarrow PROBLEM : VANISHING GRADIENT $\rightarrow \forall$ STEP \rightarrow LESS INPUT TO CHANGE NEXT WEIGHT

\hookrightarrow SOL : LSTM (long short term memory)

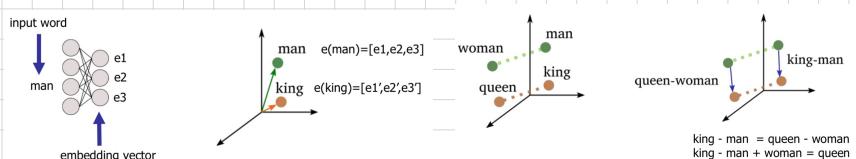
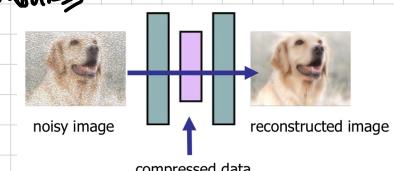
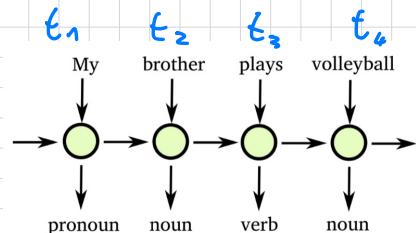
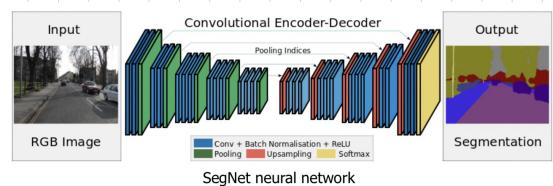
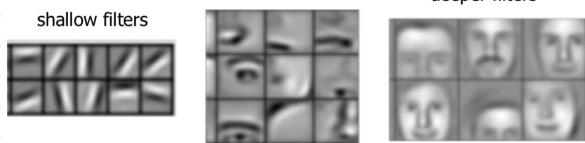
- AUTOENCODERS :

COMPRESS INPUT AND THEN RECONSTRUCT INPUT

- WORD EMBEDDINGS :

WORDS \mapsto VECTORS

\rightarrow MATH OPERATIONS CAN BE PERFORMED



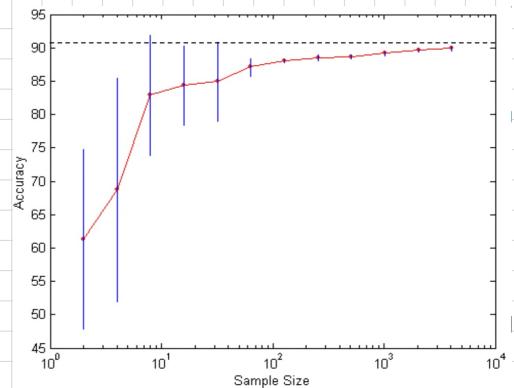
MODEL EVALUATION:

LEARNING CURVE:

$$\text{ACCURACY} = \frac{1}{N} (\text{SAMPLE SIZE})$$

IT NEEDS LARGE SAMPLES, OTHERWISE:

BIAS IN ESTIMATE, VARIANCE OF ESTIMATES



METHODS OF ESTIMATION: HOLDOUT, CROSS VALIDATION

HOLD OUT: 80% FOR TRAINING, 20% FOR TEST

CROSS VALIDATION: PARTITION DATA IN K DISTINCT SUBSETS

K-FOLD: \forall FOLD \rightarrow TRAIN ON K-1 PARTITIONS, TEST ON REMAINING ONE

MODEL VALIDATION:

TRAINING 60%, VALIDATION 20%, TEST 20%

METRICS FOR MODEL EVALUATION: \rightarrow es. FOR BINARY CLASSIFIER

$$\text{ACCURACY} = \frac{TP + TN}{TP + TN + FP + FN}$$



NOT ALWAYS RELIABLE:

$$\text{es. } \begin{cases} C: 0 \rightarrow \# : 990 \\ C: 1 \rightarrow \# : 100 \end{cases} \xrightarrow{\text{note: } 0 \mapsto C: 0} \text{Acc} = \frac{990}{1000} = 99\% (?)$$

		PREDICTED CLASS	
		Class=Yes	Class>No
ACTUAL CLASS	Class=Yes	a (TP)	b (FN)
	Class>No	c (FP)	d (TN)

\rightarrow INAPPROPRIATE FOR UNBALANCED LABEL DISTRIBUTION

$$\text{RECALL } (C: \text{Yes}) = \frac{TP}{TP + FN}$$



$$\text{PRECISION } (C: \text{Yes}) = \frac{TP}{TP + FP}$$



$$\text{F1-SCORE } (C: \text{Yes}) = 2 \cdot \frac{\text{PRECISION} \cdot \text{RECALL}}{\text{PRECISION} + \text{RECALL}} = 2 \cdot \frac{P \cdot R}{P + R}$$

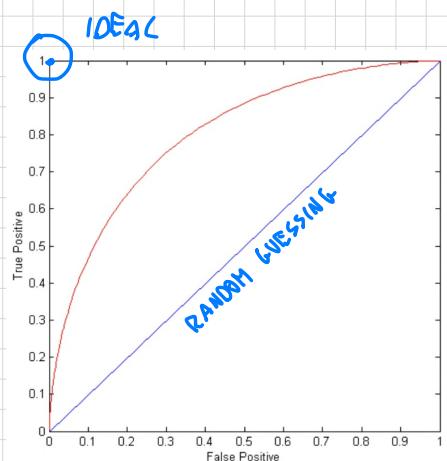
• ROC: \rightsquigarrow RECEIVER OPERATING CHARACTERISTICS

• GIVEN:

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}, \quad \text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

TP
True Positive Rate

FP
False Positive Rate



$\rightarrow \text{ROC} : \text{TPR} = f(\text{FPR})$

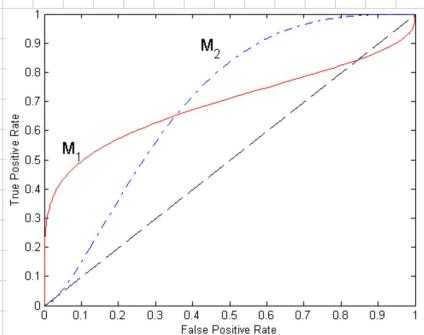
- $(0, 1)$: IDEAL
- $(0, 0)$: DECLARE EVERYTHING TO NEGATIVE CLASS
- $(1, 1)$: DECLARE EVERYTHING TO POSITIVE CLASS

\rightarrow DIAGONAL LINE : RANDOM GUESSING

\rightarrow BELOW DIAG. LINE: PREDICTION IS OPPOSITE OF THE TRUE CLASS

• BUILDING: USE CLASSIFIER THAT PRODUCES $P(+|A)$, \forall TEST INSTANCE
SORT INSTANCES DESC

APPLY THRESHOLD \forall value in $P(+|A)$



- M_1 : BETTER FOR small FPR
- M_2 : BETTER FOR large FPR

CLUSTERING FUNDAMENTALS:

PARTITIONING: DATA \rightarrow NON OVERLAPPING SET

HIERARCHICAL CLUSTERING: DATA \rightarrow NESTED CLUSTERS

CATEGORIES:

- WELL SEPARATED
- CENTER-BASED
- CONVEX VS CLUSTERS
- DENSITY BASED
- SHARED PROPERTY

K-MEANS CLUSTERING:

- K CLUSTERS, \forall CLUSTER \rightarrow CENTROID

S: CLUSTERS WITH DIFFERENT SHAPES, OBTINIES, NON HOMOGENEOUS SHAPES

ALGORITHM:

1. SELECT K

2. WHILE (CENTROIDS DON'T CHANGE):

 2.1 FORM K CLUSTER COMPUTING DISTANCE TO CLOSEST CENTROID

 2.2 RECOMPUTE CENTROID

CLUSTERS EVALUATION:

$$SSE = \sum_{i=1}^K \sum_{x_i \in C_i} \text{DIST}^2(m_i, x_i)$$

L: SUM OF
SQUARE ERRORS

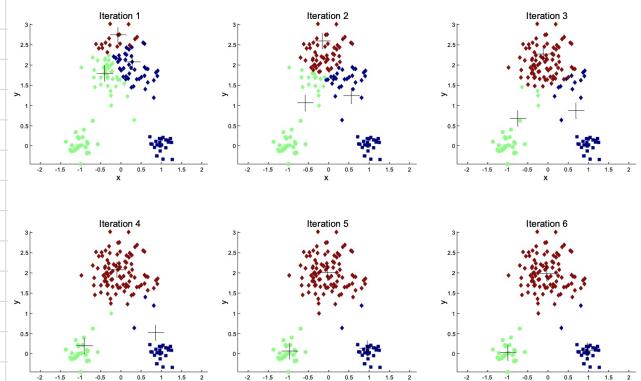
• PRE PROCESSING: NORMALIZE DATA, ELIMINATE OUTLIERS

• POST PROCESSING: DECIDE SMALL CLUSTERS, SPLIT CLUSTERS/HIGH SSE, MERGE CLUSTERS/LOW SSE

BISECTING K-MEANS:

K-MEANS, BUT 1 CLUSTER IN LIST IS BISECTED

THEN SELECT 2 CLUSTERS WITH LOWER SSE AND ADD TO LIST



```

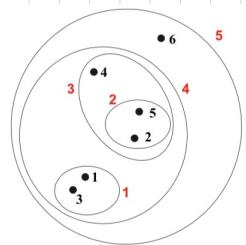
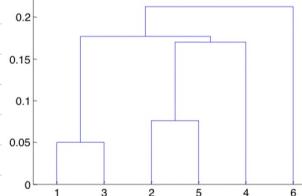
1: Initialize the list of clusters to contain the cluster containing all points.
2: repeat
3:   Select a cluster from the list of clusters
4:   for i = 1 to number_of_iterations do
5:     Bisect the selected cluster using basic K-means
6:   end for
7:   Add the two clusters from the bisection with the lowest SSE to the list of clusters.
8: until Until the list of clusters contains K clusters

```

HIERARCHICAL CLUSTERING:

SET OF NESTED CLUSTERS ORGANIZED AS
HIERARCHICAL TREE

DENDROGRAM



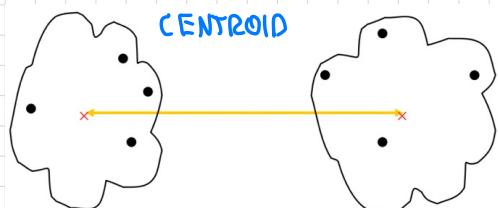
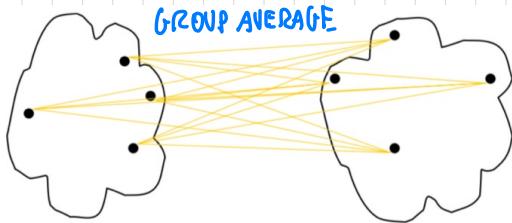
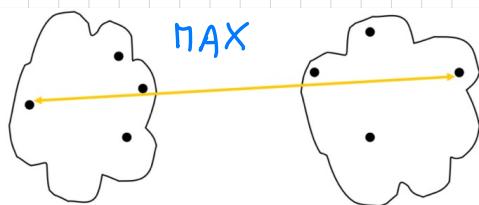
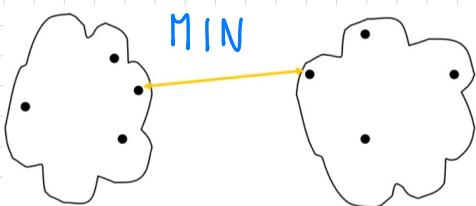
V: DO NOT HAVE TO ASSUME A N° OF CLUSTER \rightarrow JUST NEED TO "WT" AT A LEVEL

- TYPES:
 - AGGLOMERATIVE: START AS ALL-INCLUSIVE CLUSTER \rightarrow JOIN 2 CLOSEST CLUSTERS UNTIL $1/K$ REMAIN
 - DIVISIVE: START AS ALL-INCLUSIVE CLUSTER \rightarrow SPLIT CLUSTER UNTIL $1/K$ REMAIN

\rightarrow USES OF A SIMILARITY OR DISTANCE MATRIX

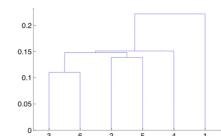
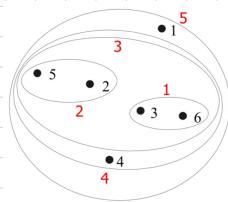
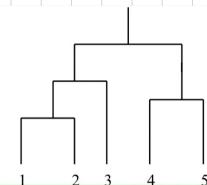
- ALGORITHM (most popular):

1. COMPUTE PROXIMITY MATRIX
 2. LET EACH DATA POINT BE A CLUSTER
 3. while (until 1 cluster remain):
 - 3.1 MERGE THE 2 CLOSEST CLUSTERS
 - 3.2 UPDATE PROXIMITY MATRIX \rightarrow DIFFERENCE IN DEFINING CLUSTERS DISTANCE
- HOW TO DEFINE CLUSTERS SIMILARITY?



MIN :

	I1	I2	I3	I4	I5
I1	1.00	0.90	0.10	0.65	0.20
I2	0.90	1.00	0.70	0.60	0.50
I3	0.10	0.70	1.00	0.40	0.30
I4	0.65	0.60	0.40	1.00	0.80
I5	0.20	0.50	0.30	0.80	1.00



V: NON-EUCLIDEAN

SMOKE HANDLING

S: SENSITIVE TO NOISE

AND OUTLIERS

Original Points

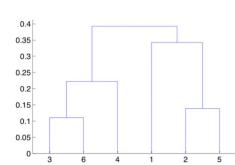
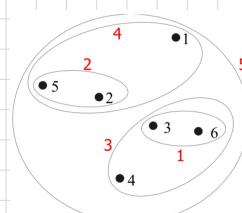
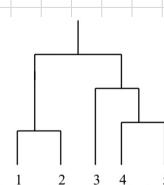
Two Clusters

Original Points

Two Clusters

MAX :

	I1	I2	I3	I4	I5
I1	1.00	0.90	0.10	0.65	0.20
I2	0.90	1.00	0.70	0.60	0.50
I3	0.10	0.70	1.00	0.40	0.30
I4	0.65	0.60	0.40	1.00	0.80
I5	0.20	0.50	0.30	0.80	1.00



V: LESS SUSCEPTIBLE TO NOISE AND OUTLIERS

S: TENDS TO BREAK LARGE CLUSTERS

BIASED TOWARDS SUBULAR CLUSTERS

Original Points

Two Clusters

Original Points

Two Clusters

GROUP AVERAGE:

CLUSTER

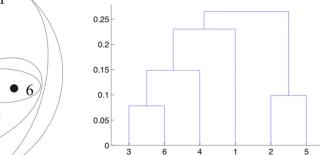
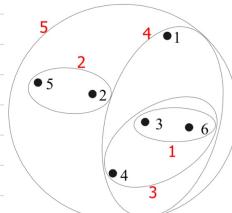
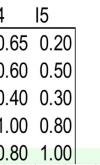
$$\text{proximity } (C_i, C_j) = \frac{\sum_{P \in C_i} \sum_{P' \in C_j} \text{proximity}(P, P')}{|C_i| \cdot |C_j|}$$

↳ CLUSTER orientation

$P \in C_i$

r

	I1	I2	I3	I4	I5
I1	1.00	0.90	0.10	0.65	0.20
I2	0.90	1.00	0.70	0.60	0.50
I3	0.10	0.70	1.00	0.40	0.30
I4	0.65	0.60	0.40	1.00	0.80
I5	0.20	0.50	0.30	0.80	1.00



V: LESS SUSCEPTIBLE TO NOISE AND OUTLIERS

S: BIASED TOWARDS SUBULAR CLUSTERS

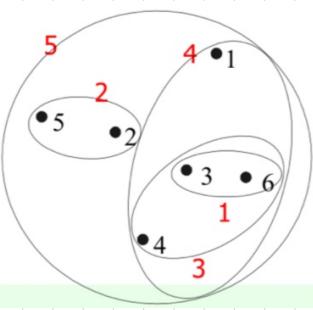
WARD'S METHOD:

SIMILARITY ↗ INCREASE OF MSE WHEN 2 CLUSTERS ARE MERGED

→ HIERARCHICAL ANALOGUE OF K-MEANS

V: LESS SUSCEPTIBLE TO NOISE AND OUTLIERS

S: BIASED TOWARDS SUBULAR CLUSTERS



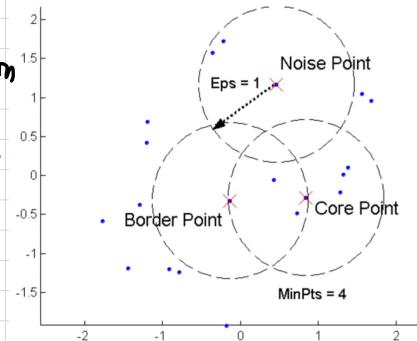
• COMPUTATION:

- $O(n^2)$ IN SPACE: USE OF A MATRIX $n \times n$
- $O(n^2)$ IN TIME: N STEPS / 1 STEP $\rightarrow N^2$ MATRIX TO BE UPDATED
 \rightarrow CAN BE OPTIMIZED TO $O(N^2 \log(N))$

• DBSCAN:

IT IS A DENSITY BASED ALGORITHM / DENSITY: n^o POINTS
 WITHIN A SPECIFIED RADIUS $\approx \text{EPS}$
 C_{EPS} : CIRCLE OF RADUS EPS

\rightarrow THRESHOLD FOR MINIMUM POINTS

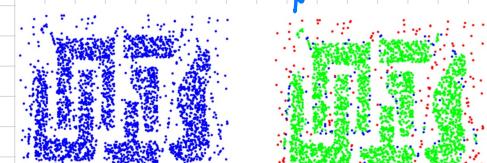


- CORE POINT $\Leftrightarrow \# \text{POINTS} > \text{MINPTS}$ IN C_{EPS}
- BORDER POINT $\Leftrightarrow (\# \text{POINTS} \in C_{\text{EPS}} < \text{MIN PTS})$
 AND $(\in \text{NEIGHBORHOOD OF A CORE POINT})$
- NOISE POINT $\Leftrightarrow (\text{! CORE POINT})$ AND (! BORDER POINT)

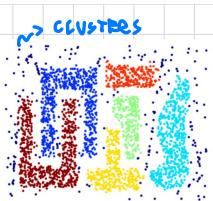
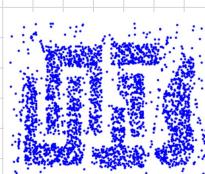
• ALGORITHM:

```

current_cluster_label ← 1
for all core points do
    if the core point has no cluster label then
        current_cluster_label ← current_cluster_label + 1
        Label the current core point with cluster label current_cluster_label
    end if
    for all points in the Eps-neighborhood, except  $i^{th}$  the point itself do
        if the point does not have a cluster label then
            Label the point with cluster label current_cluster_label
        end if
    end for
end for
    
```



CORE, BORDER, NOISE POINTS

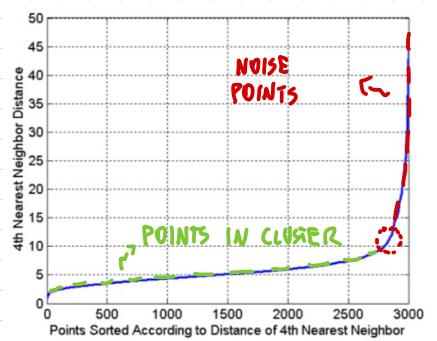


V: NOISE RESISTANT

CAN HANDLE DIFFERENT CLUSTER SHAPES

• DETERMINING EPS AND MIN PTS:

- \forall POINTS IN CLUSTER $\rightarrow k^{th}$ NEAREST NEIGHBOR \approx SAME DISTANCE
 \rightarrow NOISE POINTS HAVE k^{th} NEAREST NEIGHBOR AT DISTANCE $\gg 0$



CLUSTER VACUITY:

- DETERMINE CLUSTERING TENDENCY OF A SET OF DATA $\rightarrow \exists$ NON RANDOM DATA?
- COMPARISON OF RESULT WITH WELL KNOWN RESULTS
- EVALUATING CLUSTER ANALYSIS RESULTS
- COMPARING RESULTS OF 2 DIFFERENT SET OF CLUSTERS ANALYSIS
- DETERMINING "CORRECT" # OF CLUSTERS
- MEASURE OF CLUSTER VACUITY:

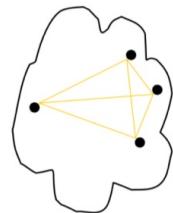
- EXTERNAL INDEX: MEASURE WHICH CLUSTER LABELS MATCH EXTERNAL SUPPLIED LABELS
 \rightarrow ex. ENTROPY, PURITY
- INTERNAL INDEX: MEASURE OF A CLUSTER HOMOGENEITY /S/ EXTERNAL INFO
- RELATIVE INDEX: COMPARE 2 DIFFERENT CLUSTERS /CLUSTERING/
 \rightarrow ex. SSE, ENTROPY
- INTERNAL MEASURE:

CLUSTER COMBINATION:

MEASURE OF HOW CLOSELY RELATED ARE OBJECT IN CLUSTER

$$WSS = \sum_i \sum_{x \in C_i} (x - m_i)^2$$

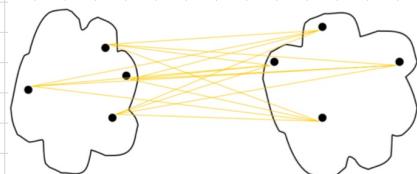
L>CENTROID



CLUSTER SEPARATION:

MEASURE OF HOW WELL DISTANCED CLUSTERS ARE

$$BSS = \sum_i |C_i| (m - m_i)^2$$



SILHOUETTE: HOW WELL OBJECT LIES IN ITS CLUSTER, (-) COMBINATION, SEPARATION

$\in [-1, 1]$

\rightarrow COMPUTABLE FOR INDIVIDUAL POINTS, INDIVIDUAL C_i , ALL CLUSTERS (C)

$$s(i) = \frac{b(i) - a(i)}{\max\{b(i), a(i)\}}$$

a(i): AVG dissimilarity of i with all objects in C of point i
b(i): min { AVG dissimilarity of i to any point not in C of i, but in other C }

- $\cdot \text{AVG}(s_i)$ ON ALL p IN CLUSTER: HOW WELL TIGHT DATA ARE IN C
- $\cdot \text{AVG}(s_i)$ ON ALL DATASET: CORRECTNESS ON THIS CLUSTERING

EXTERNAL MEASURES:

ENTROPY AND PURITY

entropy For each cluster, the class distribution of the data is calculated first, i.e., for cluster j we compute p_{ij} , the 'probability' that a member of cluster j belongs to class i as follows: $p_{ij} = m_{ij}/m_j$, where m_j is the number of values in cluster j and m_{ij} is the number of values of class i in cluster j . Then using this class distribution, the entropy of each cluster j is calculated using the standard formula $e_j = \sum_{i=1}^L p_{ij} \log_2 p_{ij}$, where L is the number of classes. The total entropy for a set of clusters is calculated as the sum of the entropies of each cluster weighted by the size of each cluster, i.e., $e = \sum_{j=1}^K \frac{m_j}{m} e_j$, where m_j is the size of cluster j , K is the number of clusters, and m is the total number of data points.

purity Using the terminology derived for entropy, the purity of cluster j , is given by $\text{purity}_j = \max p_{ij}$ and the overall purity of a clustering by $\text{purity} = \sum_{j=1}^K \frac{m_j}{m} \text{purity}_j$.

RAND INDEX :

• MAIN IDEA: 2 OBJECT IN SAME CLUSTER \Rightarrow SAME CLASS

$$\text{RAND INDEX} = \frac{f_{00} + f_{11}}{f_{00} + f_{01} + f_{10} + f_{11}}$$

PAIRS OF 2 POINTS: $(P_1, P_2) \rightarrow$ CONSIDER THESE COMBINATIONS

- f_{00} : # PAIRS / CLASS IS \neq , $C_i \neq C_j$
- f_{01} : # PAIRS / CLASS IS \neq , $C_i \neq C_j$
- f_{10} : # PAIRS / CLASS IS \neq , $C_i \neq C_j$
- f_{11} : # PAIRS / CLASS IS $=$, $C_i = C_j$

Table 5.9. K-means Clustering Results for LA Document Data Set

Cluster	Entertainment	Financial	Foreign	Metro	National	Sports	Entropy	Purity
1	3	5	40	506	96	27	1.2270	0.7474
2	4	7	280	29	39	2	1.1472	0.7756
3	1	1	1	7	4	671	0.1813	0.9796
4	10	162	3	119	73	2	1.7487	0.4390
5	331	22	5	70	13	23	1.3976	0.7134
6	5	358	12	212	48	13	1.5523	0.5525
Total	354	555	341	943	273	738	1.1450	0.7203

• REGRESSION ANALYSIS :

DEF.: REGRESSION . PREDICTION OF A VALUE $\in \mathbb{R}$ ↳ PREDICTION OF A FINITE CLASS

- REGRESSION FINDS A MODEL THAT ALLOWS PREDICTING THE TARGET VARIABLE OF NEW OBJECT USING $Y = f(X_1, \dots, X_m)$

• SIMPLE LINEAR REGRESSION :

$$Y = \beta_0 + \beta_1 \cdot X$$

X : INDEPENDENT VARIABLE

Y : ESTIMATED / PREDICTED VARIABLE

β_1 : ESTIMATION OF REGRESSION SLOPE $(Y = mX + q)$

β_0 : REGRESSION INTERCEPT $(Y = mX + q)$
↳ ESTIM. Y when $X = 0$

- HOW TO FIND β_0 AND β_1 ?

→ ONE METHOD IS MINIMIZING THE RESIDUAL SUM OF SQUARES (RSS) :

- GIVEN Y : ACTUAL VALUES, \hat{Y} : PREDICTED VALUES :

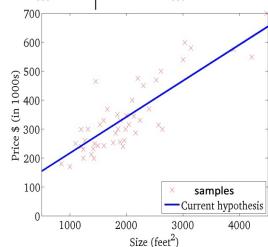
$$\min(RSS) = \min \left(\sum_i (Y_i - \hat{Y}_i)^2 \right) = \min \left[\sum_i \left(Y_i - (\beta_0 + \beta_1 \cdot X_i) \right)^2 \right]$$

- ESTIMATION OF β_0 AND β_1 :

$$\begin{aligned} \cdot Y = \beta_0 + \beta_1 \cdot X &\longrightarrow \begin{cases} \beta_1 = \frac{\sum_i (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_i (X_i - \bar{X})^2} \\ \beta_0 = \bar{Y} - \beta_1 \cdot \bar{X} \end{cases} & \begin{aligned} \bar{Y} &= \frac{1}{n} \sum_i Y_i \\ \bar{X} &= \frac{1}{n} \sum_i X_i \end{aligned} \end{aligned}$$

∴ ESTIM. HOME \$ = f(SIZE IN feet²)

Size in feet ²	Price (\$) in 1000's
2104	460
1416	232
1534	315
852	178
...	...



· X : SIZE IN feet², Y : PRICE IN \$ IN 1000'S

$$\rightarrow Y = \beta_0 + \beta_1 \cdot X$$

· β_0 MEANING: PRICE WHEN $X = 0$

→ β_0 IS THE PARTITION OF HOUSE PRICE

WHICH IS INDEPENDENT FROM SQUARE FEET

MULTIPLE LINEAR REGRESSION:

$$Y = f(\vec{X}) = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \dots + \beta_n \cdot X_n + \epsilon$$

\nearrow n-dimensional plane equation

RESIDUAL ϵ : PORTION OF Y WHICH IS NOT EXPLAINED BY THE MODEL

β_i : AVG EFFECT ON Y OF A UNIT INCREASE IN X_i , HOLDING FIXED ALL OTHER PREDICTORS

- CORRELATION AMONG X_i CAN CAUSE PROBLEMS
 - ↳ COLLINEARITY SHOULD BE AVOIDED AMONG X_i 'S
- IF $\# X_i > 0 \rightarrow$ POSSIBILITY OF PROVIDING REDUNDANT INFORMATION

→ FEATURES X_i SELECTION AND REMOVAL (BASED ON CORRELATION)
 → FEATURES SELECTION BASED ON CORRELATION TEST

POLYNOMIAL REGRESSION:

IT IS USED WHEN THE RELATION $Y = f(\vec{X})$ IS NON-LINEAR:

$$Y = \beta_0 + \beta_1 \vec{x} + \beta_2 \vec{x}^2 + \dots + \beta_k \vec{x}^k + \epsilon$$

- IT CONSISTS OF:
 - COMPUTING NEW FEATURES AS A POWER FUNCTION OF INPUT FEATURES
 - APPLY LINEAR REGRESSION ON THESE NEW FEATURES

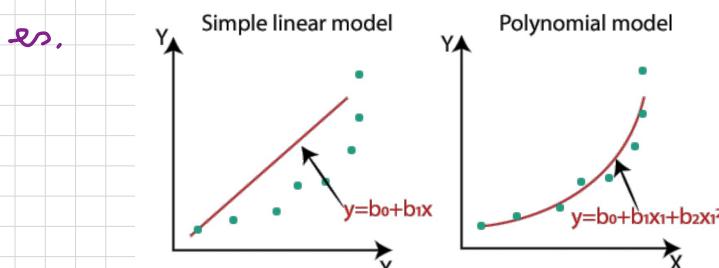
Ex. $n=2 \rightarrow \vec{x} = (x_1, x_2)$:

$$Y = \beta_0 + \beta_1 \vec{x} + \beta_2 \vec{x}^2 + \epsilon \longrightarrow Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{11} x_1^2 + \beta_{22} x_2^2 + \beta_{12} x_1 x_2 + \epsilon$$

- THE MODEL CAN BE WRITTEN AS $Y = X \cdot \beta + \epsilon$

$$X = [x_1, \dots, x_n]^T$$

$$\beta = [\beta_0, \dots, \beta_k]^T$$



• CONSIDERATIONS :

- KEEP THE ORDER OF THE POLYNOMIAL AS LOW AS POSSIBLE
usuallly $\rightarrow \text{MAX ORDER} = 2$
- HIGHER ORDER CAN LEAD TO OVERFITTING
- TRY MODEL FITTING FOR INCREASING ORDER UNTIL NONE ORDERS IS NON SIGNIFICANT

V : POLYNOMIAL FIT WIDE RANGE OF CURVATURE

USUALLY THEY ARE THE BEST APPROXIMATION

S : TOO SENSITIVE TO OUTLIERS

HIGH DEGREE \rightarrow POSSIBLY DATA OVERFITTING

\hookrightarrow AVOID DATA OVERFITTING : • USE A LARGE TRAINING DATA

• USE LOWER MODEL COMPLEXITY

• USE REGULARIZATION TECHNIQUES

REGULARIZATION TECHNIQUES

THEY PERFORM VARIABLE SELECTION AND REGULARIZATION TO ENHANCE ACCURACY AND INTERPRETABILITY

→ REDUCTION IN MODEL COMPLEXITY AND OVERFITTING WITH :

(# VARIABLES > # OBSERVATIONS TO BE DESCRIBED) OR (LEARNED MODEL IS POORLY GENERALIZABLE)

→ THEY PERFORM CHANGES ON β_j

• NORMAL LINEAR REGRESSION : $RSS = \sum_i (y_i - \hat{y}_i)^2 = \sum_i (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2$

• RIDGE REGRESSION : $RSS + \lambda \sum_{j=1}^p \beta_j^2$
IT ADDS L2 AS PENALTY : $\sum_j \beta_j^2$

→ IT'S LIKE RSS, UNDERR THE CONDITION

■ Ridge tends to lower uniformly all the coefficients

- Coefficients already close to 0 have little effect on the sum of squares (if $x \approx 0$, $x^2 < x$)



$\sum_{j=1}^p \beta_j^2 \leq C$, FOR $C > 0$

FACT ABSOLUTE SHRINKAGE AND SELECTION OPERATOR

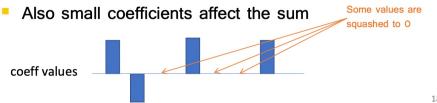
• LASSO REGRESSION : $RSS + \lambda \sum_{j=1}^p |\beta_j|$

IT ADDS L1 AS PENALTY : $\sum_j |\beta_j|$

→ IT APPLIES TRUE CONSTRAINT $\sum_{j=1}^p |\beta_j| \leq C$, FOR $C > 0$

■ Lasso tends to assign the value 0 to some coefficients (feature selection)

- Also small coefficients affect the sum



/ λ : AMOUNT OF SHRINKAGE

• RIDGE vs LASSO :

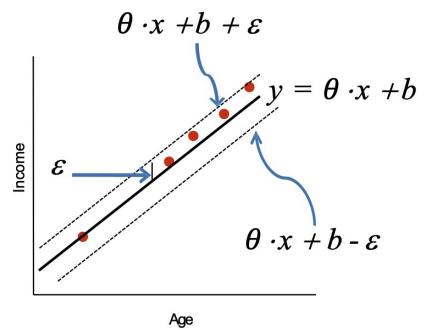
• LASSO SIMPLIFIES A SIMPLER MODEL : LEADS $\beta_j = 0$ COEFFICIENT

. { RIDGE : SHRINK SIZE OF SOME β_j
{ LASSO : SETS SOME $\beta_j := 0$

SUPPORT VECTOR REGRESSION:

DEFINE HOW MUCH ERROR IS ACCEPTABLE IN MODEL

→ IT FINDS AN APPROPRIATE HYPERPLANE TO FIT DATA



TRAINING PROBLEM IS A CONVEX OPTIMIZATION PROBLEM:

$$\min \left\{ \frac{1}{2} \|\theta\|^2 \right\}, \text{ CONSTRAINTS: } \begin{cases} y^i - \theta \cdot x^i - b \leq \varepsilon \\ \theta \cdot x^i + b - y^i \leq \varepsilon \end{cases}$$

SOFT MARGIN:

SOMETIMES, GIVEN AN ε

→ PROBLEM NOT FEASIBLE

→ REFORMULATE PROBLEM CONSIDERING ERRORS

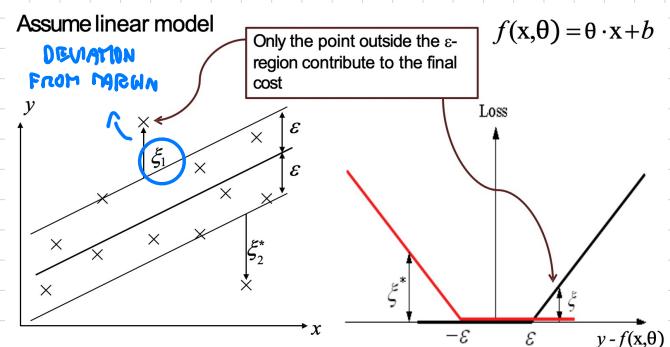
/ DON'T SATISFY ε

→ OPTIMIZATION PROBLEM:

$$\min \left\{ \frac{1}{2} \|\theta\|^2 \right\} + C \sum_{i=1}^m (\xi_i + \xi_i^*)$$

$$\text{CONSTRAINTS: } \begin{cases} y^i - \theta \cdot x^i - b \leq \varepsilon + \xi_i \\ \theta \cdot x^i + b - y^i \leq \varepsilon + \xi_i^*, \quad \xi_i, \xi_i^* \geq 0 \end{cases}$$

→ WE MINIMIZE DEVIATION ξ_i FROM THE MARGIN



NON-LINEAR CASE:

USE A MAPPING FOR A HIGHER ORDER DIMENSIONAL SPACE:

KERNEL TRANSFORMATION (POLYNOMIAL, GAUSSIAN RADIAL, ...)

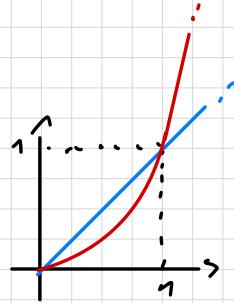
$$x_i \mapsto \phi(x_i)$$

$$\min \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*)$$

$$\text{s.t. } y^i - \theta \cdot \phi(x^i) - b \leq \varepsilon + \xi_i; \\ \theta \cdot \phi(x^i) + b - y^i \leq \varepsilon + \xi_i^* \\ \xi_i, \xi_i^* \geq 0, i = 1, \dots, m$$

EVALUATING REGRESSION:

- MEAN ABSOLUTE ERROR: $MAE = \frac{1}{n} \sum_i |Y_i - \hat{Y}_i|$



- MEAN SQUARED ERROR: $MSE = \frac{1}{n} \sum_i (Y_i - \hat{Y}_i)^2$

\hookrightarrow PENALIZES LESS ERRORS CLOSER TO 0

$\rightarrow MAE, MSE > 0$, $\downarrow MAE, MSE \rightarrow$ BETTER MODEL

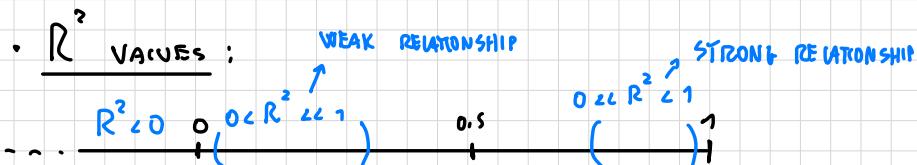
- RESIDUAL STANDARD ERROR: $RSE = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2}$ / n : # SAMPLES

R-SQUARED:

IT MEASURES THE GOODNESS OF FIT OF A MODEL

$$R^2 = 1 - \frac{RSS}{TSS} = 1 - \left[\frac{MSE}{S^2} \right] \quad / TSS = S^2 = \sum_i (Y_i - \bar{Y}_i)^2$$

• FVU: RATIO BETWEEN UNEXPLAINED VARIANCE AND TOTAL VARIANCE



• $R^2 = 1$: PERFECT LINEAR RELATIONSHIP BETWEEN X AND Y

• $R^2 = 0$: NO LINEAR RELATIONSHIP BETWEEN X AND Y

• $R^2 < 0$: MODEL IS PREDICTING WORSE THAN $\hat{Y} = \bar{Y}$

R^2 ADJUSTED:

IF MULTIPLE LINEAR REGRESSION $\rightarrow \uparrow \# X_i \rightarrow \downarrow R^2$ NEVER REMAIN CONSTANT
 n^2 OR INCREASES ONLY

\rightarrow USING THE ADJUSTED \bar{R}^2 : $\uparrow \bar{R}^2 \Leftrightarrow$ NEW X_i ADDED IMPROVE THE MODEL

- $\bar{R}^2 = 1 - \frac{n-1}{n-p-1} \left(1 - R^2 \right)$ / n : # SAMPLES
 p : # X_i

TIME SERIES ANALYSIS :

TIME SERIES : SEQUENTIAL SET OF DATA POINTS MEASURED OVER SUCCESSIVE TIMES $\rightarrow S = f(t)$

COMPONENTS :

TREND COMPONENT :

LONG-TERM MOVEMENT, INCREASING / DECREASING OVER t

UPWARD OR DOWNWARD TENDENCY

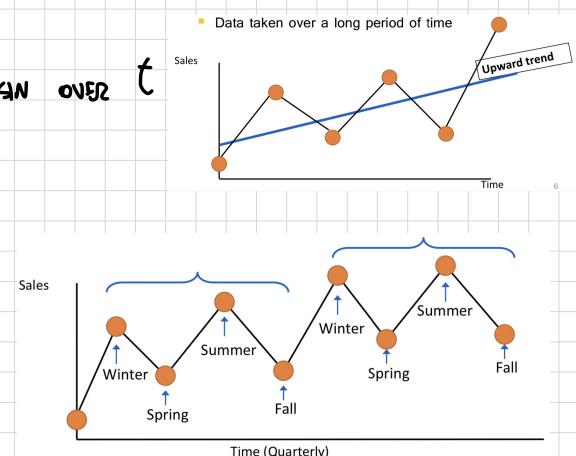
CAN BE LINEAR OR NON-LINEAR



SEASONAL COMPONENT :

REGULAR PERIODIC FLUCTUATIONS

REGULAR WAVE-LIKE PATTERN

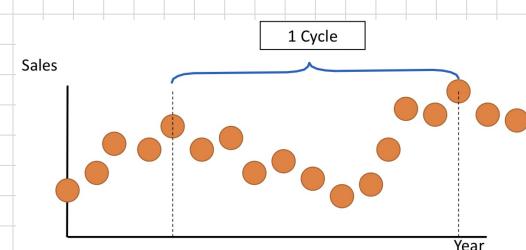


CYCICAL COMPONENT :

REPEATING CYCLES OR RHYTHMS

REGULARLY OCCUR BUT MAY VARY IN LENGTH

MEASURED PEAK TO PEAK



IRREGULAR / RANDOM COMPONENTS :

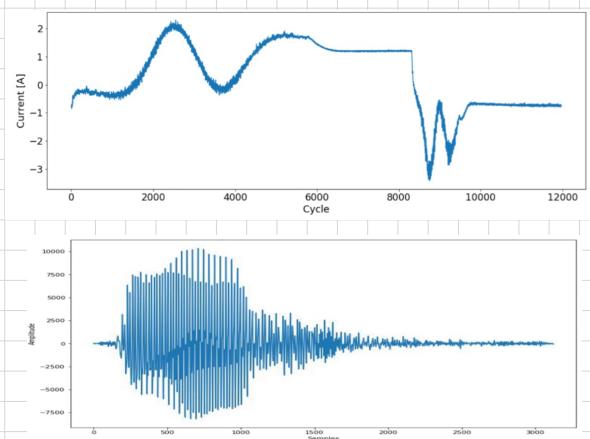
CAUSED BY UNPREDICTABLE INFLUENCES

\rightarrow THIS VARIANCE REPRESENTS NOISE IN TIME SERIES

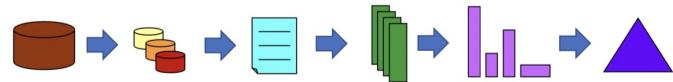
DISCRETE TIME SERIES



CONTINUOUS TIME SERIES



- KDD . KNOWLEDGE DISCOVERY IN DB



- PRE-PROCESSING**
- TIME SERIES ALIGNMENT TECHNIQUES

RIGHT RE REQUIRED

$$S_1, \dots, S_m = f(t)$$



TRANSFORMATION

- FEATURE ENGINEERING , FEATURE EMBEDDING

- ANALYTICS TASKS :

CLASSIFICATION AND FORECASTING

- ARTIFICIAL NEURAL NETWORKS CAN BE USED

- CNN :

- VGGISH IS A CNN TO EXTRACT FEATURES FROM AUDIO SIGNALS
- Log Mel Spectrogram Audio \mapsto Audio Tag

- RNN :

CONNECTION BETWEEN NODES FORM A DIRECTED GRAPH ALONG A TEMPORAL SEQUENCE

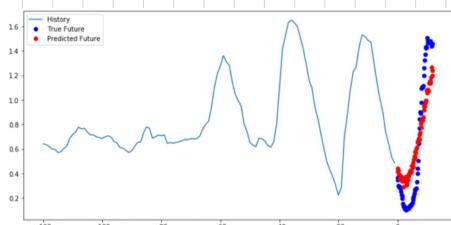
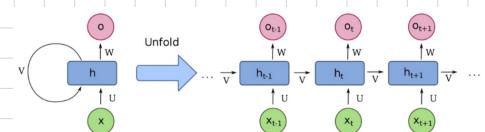
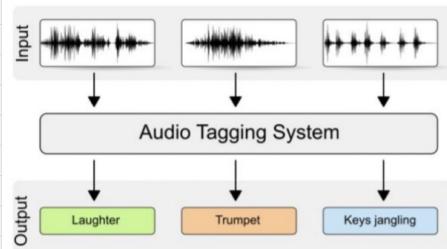
\rightarrow RNN CAN USE THEIR INTERNAL STATE MEMORY TO PROCESS SEQUENCES OF INPUT

CONVOLUTIONAL



CNN : CLASSIFICATION TASK

RNN : FORECASTING TASK
 \rightsquigarrow RECURRENT



\checkmark : POWERFUL ALGORITHMS TO TRAIN ACCURATE MODELS

ABLE TO DEAL WITH DIFFERENT COMPLEX ANALYTICS TASKS

S: HUGE AMOUNT OF DATA REQUIRED

TRAINING IS AN HEAVY TASK COMPUTATIONALLY , IN TIME AND MEDIUM
FEATURE LEARNING STEP HIDDEN IN THE NETWORK

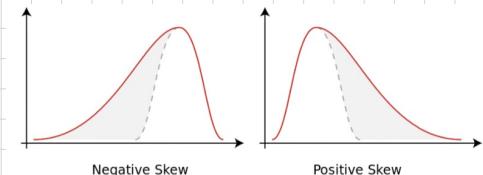
FEATURE ENGINEERING :

FEATURE COMPUTATION OVER A TIME SERIES ;

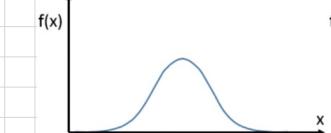
- BASIC STATISTICS (MIN, MAX, etc.), INDICES (KURTOSIS, SKEWNESS), SUMMARIZATION (PERCENTILE, CDF + PERCENTILE, DEVIATE + CDF + PERCENTILE)
- , LINEAR REGRESSION
- SKEWNESS :

MEASURE OF ASYMMETRY OF PROB. DISTRIBUTION OF A REAL VALUED RV ABOUT ITS MEAN

$$\gamma_1 = E \left[\left(\frac{X - \mu}{\sigma} \right)^3 \right] = \frac{\mu_3}{\sigma^3} = \frac{k_3}{k_2^{3/2}}$$



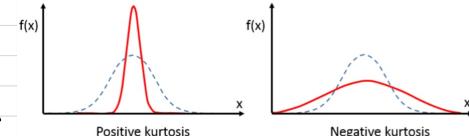
k_t : t-th cumulant



- KURTOSIS :

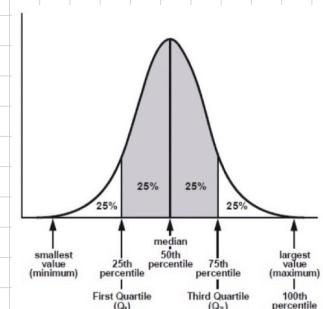
MEASURE OF "TAILEDNESS" OF PROB. DISTRIBUTION OF A REAL VALUED RV

$$KURT[X] = E \left[\left(\frac{X - \mu}{\sigma} \right)^4 \right] = \frac{\mu_4}{\sigma^4}$$



- SUMMARIZATION :

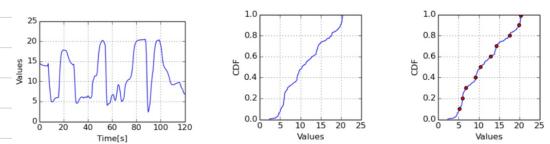
• PERCENTILE: VALUE BELOW WHICH A GIVEN % OF OBSERVATIONS IN A GROUP OF OBSERVATIONS FALL
→ OUTLIERS CAN BE REMOVED, REMOVING LAST PERCENTILES OF THE DISTRIBUTION



• PERCENTILES ARE FEATURES DESCRIBING THE TIME SERIES

• CDF: $F_X(x) = P(X \leq x)$

$$\rightarrow P(a < X < b) = F_X(b) - F_X(a)$$



- LINEAR REGRESSION:



- FEATURE ENGINEERING CAN ALSO BE COMPUTED IN LOCAL PARTS OF ENTIRE TIME SERIES, i.e. TIME WINDOW

TIME WINDOW:

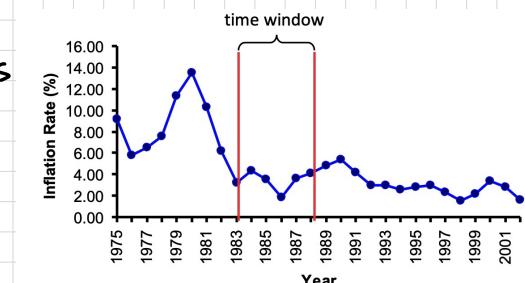
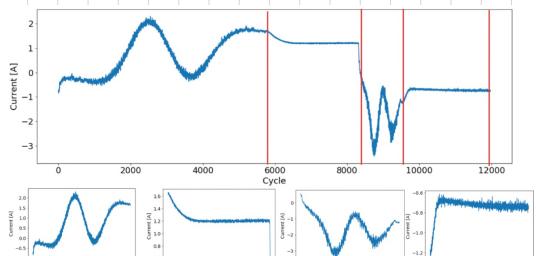
TIME SERIES PORTION ON WHICH COMPUTE FEATURES

- WINDOW LENGTH: SIZE IN t UNITS
- DOMAIN - DRIVEN / DATA - DRIVEN

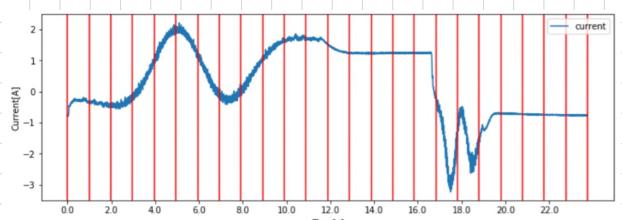
- WINDOW SHIFT: POSITION IN RESPECT TO CONSECUTIVE WINDOWS

- WINDOWS CAN BE NON-OVERLAPPED / OVERLAPPED

DOMAIN - DRIVEN



DATA - DRIVEN



- TIME SERIES TREND CAN BE CAPTURED THROUGH THE FEATURES EXTRACTED FROM EACH SUB-CYCLE OF EACH TIME SERIES

- TO DEAL WITH LONG HORIZON PREDICTION

→ AVERAGING OVER A TIME WINDOW

- IN SOME CASES PURE FEATURES CAN BE CORRELATED

→ REVERSE CORRELATED FEATURES

- TIME SERIES IN FREQUENCY DOMAIN;

