

CLASSIFICATION

- EVALUATION METRICS:

- Confusion matrix :

		predicted	
		Y	N
true	Y	TP	FN
	N	FP	TN

• accuracy = $\frac{\text{true overall}}{\text{overall}}$

• precision : $P = \frac{TP}{TP+FP}$

• recall : $R = \frac{TP}{TP+FN}$

• F1-measure : $= \frac{2RP}{R+P}$

• sensitivity = $\frac{TP}{TP+FN} = R$

• specificity = $\frac{TN}{TN+FP}$

= $\frac{\text{true positive}}{\text{classified positive}}$

= $\frac{\text{true positive}}{\text{positive}}$

↑ FP↓

↑ FN↓

↑ FP↓ FN↓

= true positive rate

= true negative rate

- PROBLEMS

- Bayes

- Decision tree

BAYES

$$\text{class} = \arg \max P(y|\vec{x}) \stackrel{\text{NB}}{=} \arg \max P(x_1|y) \cdots f(x_n|y) \cdot P(y)$$

Suppose the dataset :

	Outlook	Temp	Humidity	Windy	Play	target
	sunny/overcast/rainy	R	R	T/F	Y/N	

we create :

Outlook		Temp		Humidity		Windy		Play		
	Y	N	Y	N	Y	N	Y	N	Y	N
sunny	2+1	3+1	64.68	65.71	67.70	70.85	6+1	2+1	9	5
overcast	4+1	0+1	69.70	72.80	70.78	90.81	3+1	3+1		
rainy	3+1	2+1	72... ..	85... ..	80... ..	85... ..				
sunny	1/4	1/2	$\mu=73$ $\sigma=6.2$	$\mu=75$ $\sigma=7.9$	$\mu=79$ $\sigma=10.2$	$\mu=86$ $\sigma=9.2$	true false	7/11 4/11	3/7 4/7	$P(Y) = 9/14$ $P(N) = 5/14$
overcast	5/12	1/8								
rainy	1/3	3/8								



	outlook	Temp	Humidity	Windy	Play
	sunny	66	90	T	?

likelihood :

$$l(Y|\vec{x}) = P(\text{outlook}=\text{sunny}|Y) \cdot f_{\text{T}}(\text{temp}=66|Y) \cdot f_{\text{H}}(\text{humidity}=90|Y) \cdot P(\text{windy}=T|Y) \cdot P(Y)$$

$$l(N|\vec{x}) = P(\text{outlook}=\text{rainy}|N) \cdot f_{\text{T}}(66|N) \cdot f_{\text{H}}(90|N) \cdot P(T|N) \cdot P(N)$$



→ Normalization :

$$P(Y|\vec{x}) = \frac{l(Y|\vec{x})}{l(Y|\vec{x}) + l(N|\vec{x})} \quad ; \quad P(N|\vec{x}) = 1 - P(Y|\vec{x})$$

$$* f(\text{temp}=66|Y) = \frac{1}{\sqrt{2\pi(6.2)^2}} e^{-\frac{(66-73)^2}{2 \cdot (6.2)^2}}$$

$$\frac{1}{\sqrt{2\pi(10)^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

ASSOCIATION RULES

- Apriori
 - Eclat
 - FP-tree
 - Rule generation
 - Summarize itemsets
 - Thinning
 - Frequent sequences
- frequent itemset

$$\text{support} = \frac{\sigma(\{X, Y\})}{\# \text{transactions}} \sim \text{IP}(X, Y)$$

$$\text{confidence} = \frac{\sigma(\{X, Y\})}{\sigma(\{X\})} \sim \text{IP}(Y|X)$$

$$\sim \text{IP}(X \Rightarrow Y)$$

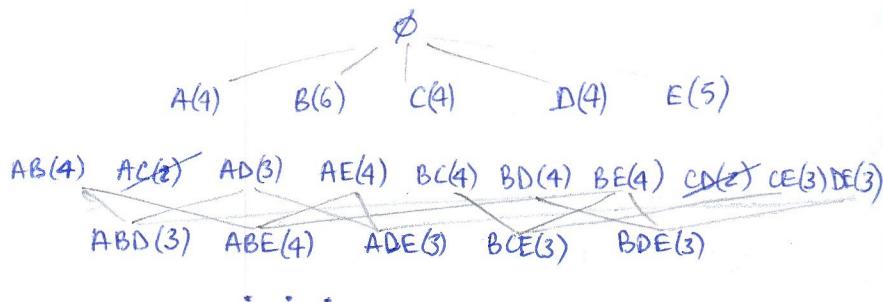
APRIORI

Threshold = 0.5

Transaction database

1	A	B	D	E
2		B	C	E
3	A	B	D	E
4	A	B	C	E
5	A	B	C	D
6		B	C	D

(threshold = 3/6)



ECLAT

Threshold = 0.5

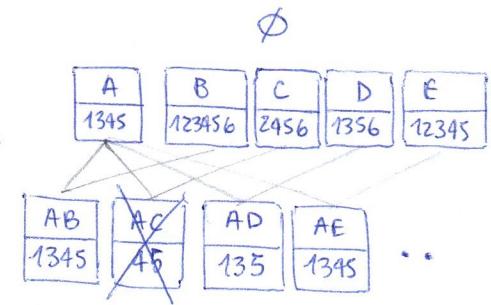
Transaction database

1	..
6	

(threshold = 3/6)

Vertical database

	A	B	C	D	E
1	1	2	1	1	1
2	2	4	3	2	
3	3	5	5	3	
4	3	6	6	4	
5	4				5
6	5				



FP-TREE

Threshold = 0.5

Transaction database

1	..
6	

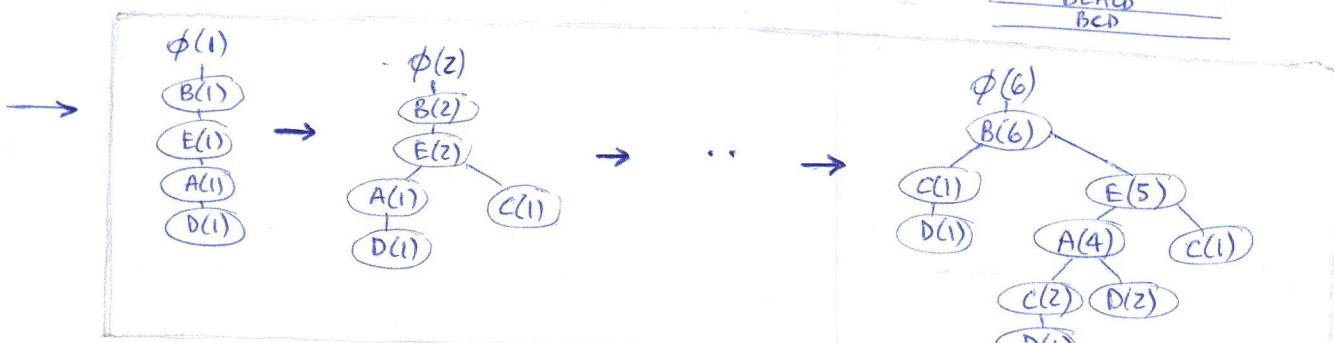
(threshold = 3/6)

Ordering:

B(6)
E(5)
A(4)
C(4)
D(4)

New transaction database:

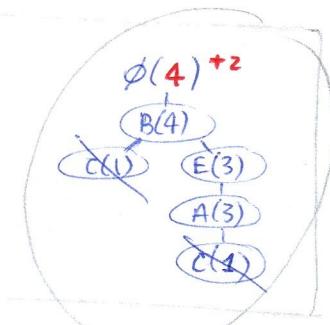
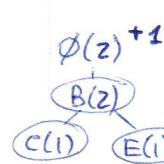
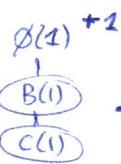
BEAD
BEC
BEAD
BEAC
BEACD
BCD



Conditional FP-tree for D:

Paths that end with D

BCD	count(D) = 1
BEACD	count(D) = 1
BEAD	count(D) = 2



Given D we have:

BEA, BE, EA, BE, B, E, A
(s(c) < minsup = 3)

Frequent itemsets

DBEA(3), DBE(3), DE(3), DBE(3), DA(3)

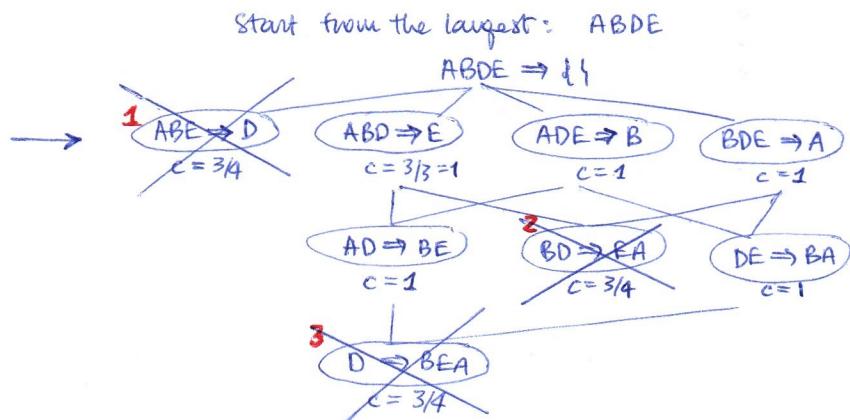
here c(2), still not enough

RULE GENERATION

min conf = 0.9

Report (frequent) itemsets

6	B
5	E, BE
4	A, C, D, AB, AE, BC, BD, ABE
3	AD, CE, DE, ABD, ADE, BCE, BDE, ABDE



We don't need to investigate:

1. ABE, AB, BE, AE, A, B, E
2. BD, B, D
3. D

Association rules:

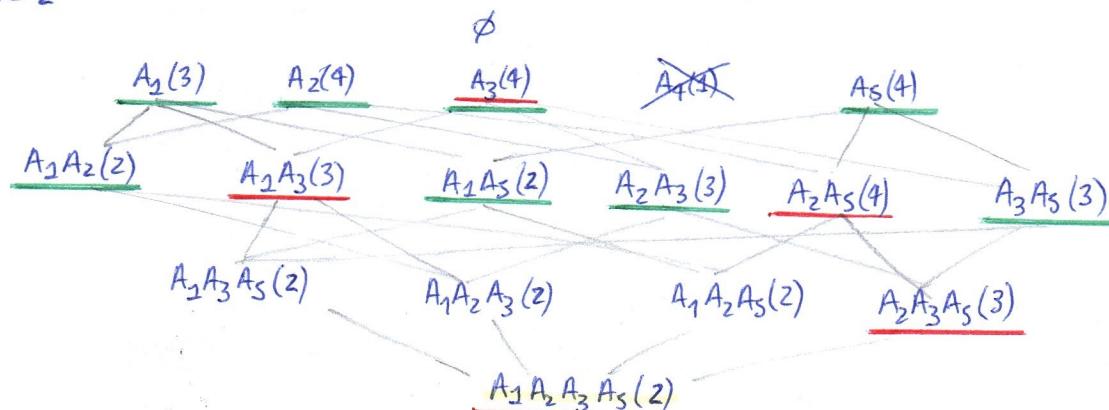
$$\begin{aligned} ABD &\Rightarrow E, \\ ADE &\Rightarrow B, \\ BDE &\Rightarrow A \end{aligned}$$

SUMMARIZING ITEMSETS

- CLOSED ITEMSET X: all super-itemsets containing X have $<$ support
- MINIMAL GENERATOR X: all sub-sets of X have $>$ support
- MAXIMAL FREQUENT ITEMSET X: no frequent subset containing X



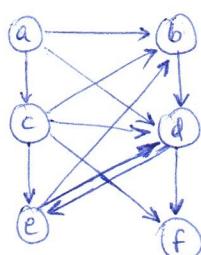
min sup = 2



Remember:

1st line: if frequent \Rightarrow —
 last line: if frequent \Rightarrow —

TRAWLING

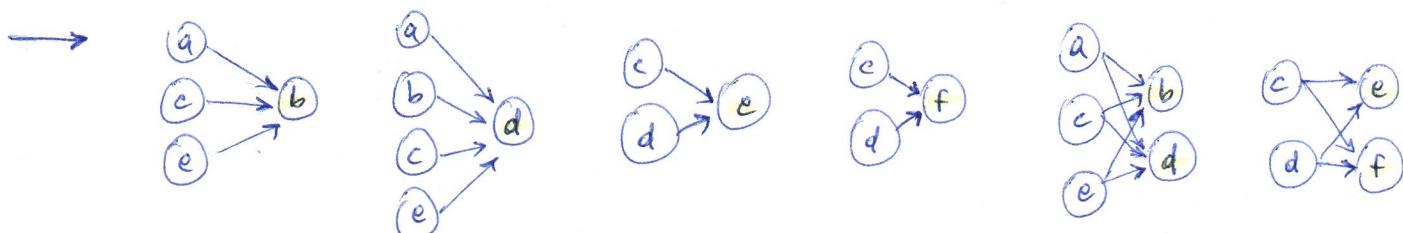


Itemsets:
 $a = \{b, c, d\}$
 $b = \{d\}$
 $c = \{b, d, e, f\}$
 $d = \{e, f\}$
 $e = \{b, d\}$
 $f = \{\}$

Frequent itemsets: min sup = 2

\emptyset
 $b(3)$ ~~x~~
 $d(4)$ ~~x~~
 $c(2)$
 $f(2)$
 $bd(3)$ ~~x~~
 $bc(2)$ ~~x~~
 $bd(2)$ ~~x~~
 $df(2)$ ~~x~~
 $ef(2)$
 \Rightarrow
 $b(3)$
 $d(4)$
 $e(2)$
 $f(2)$

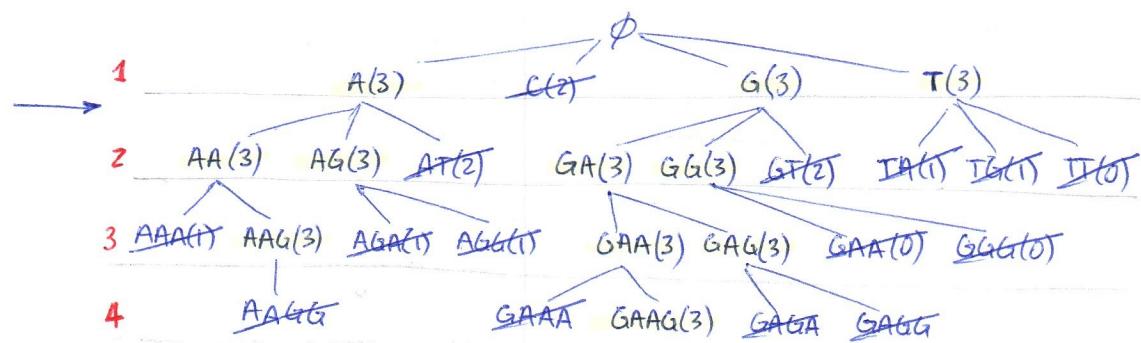
Communities:



FREQUENT SEQUENCES

minsup = 3

S1	CAGAAAGT
S2	TGACAG
S3	GAAGT



At level 1 dies C → we don't add it anywhere in 2, 3, 4
 At level 2 dies T → we don't add it anywhere in 2, 3, 4
 At level 4 dies A

CLUSTERING

- DISTANCES:

- Euclidean : $d(x, y) = \sqrt{\sum_i (x_i - y_i)^2}$
- Manhattan : $d(x, y) = \sum_i |x_i - y_i|$
- Jaccard : similarity : $J(x, y) = \frac{|x \cap y|}{|x \cup y|}$

- Hamming : distance : $d(x, y) = 1 - J(x, y)$
- d(x, y) = $\frac{\# \text{mismatch}}{\# \text{overall}}$

- Cosine : similarity = $\frac{\sum x_i y_i}{\sqrt{(\sum x_i^2)(\sum y_i^2)}}$
- distance : $d(x, y) = \arccos(\text{similarity})$

- Edit : $d(x, y) = |x| + |y| - 2|\text{LCS}|$ (longest common subsequence)

- LINKAGE

- single (min) : $d(A, B) = \min d(x_i^A, x_j^B)$
- complete (max) : $d(A, B) = \max d(x_i^A, x_j^B)$
- mean/centroid
- group average

$$d(A, B) = \text{avg } d(x_i^A, x_j^B)$$

- PROBLEMS

- Hierarchical clustering
- DBSCAN
- HDBSCAN
- silhouette coefficient & k-means

HIERARCHICAL CLUSTERING

single linkage

Distance matrix : D =

1	/			
2	1			
3	6	5	/	
4	10	8	4	/
5	9	8	5	3
	1	2	3	4

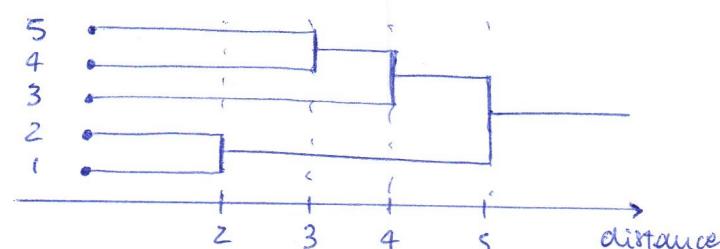
→ Merge {1}, {2}, d = 2
 $d(\{1, 2\}, \{3\}) = \min(d(1, 3), d(2, 3)) = \min(6, 5) = 5$
 $d(\{1, 2\}, \{4\}) = 9$
 $d(\{1, 2\}, \{5\}) = 8$

$D = \begin{matrix} & \{1, 2\} & - & & \\ \{1, 2\} & 3 & 5 & - & \\ & 4 & 3 & 4 & - \\ & 5 & 3 & 3 & - \\ & 1, 2 & 3 & 4 & 5 \end{matrix}$

Merge {4}, {5}, d = 3

...

→ dendrogram :

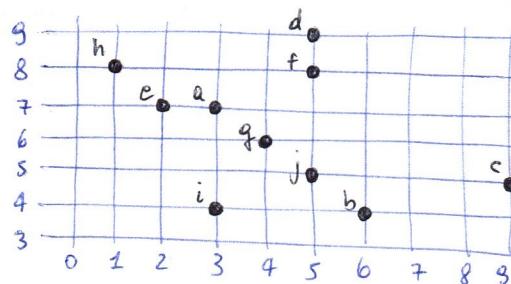
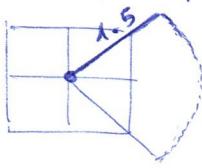


DBSCAN

- CORE POINT x : $|N_\varepsilon(x)| \geq \text{minpts}$
- BORDER POINT x : $|N_\varepsilon(x)| < \text{minpts}$, $x \in N_\varepsilon(\text{core point})$
- NOISE POINT x : $|N_\varepsilon(x)| < \text{minpts}$, $x \notin N_\varepsilon(\text{core point})$

! The point itself is part of the count

$$\varepsilon = 1.5, \text{minpts} = 3$$

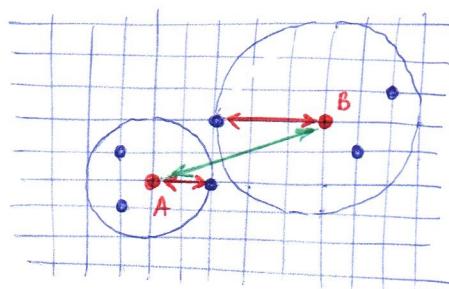


	# $N_\varepsilon(\cdot)$	
a	3	core
b	2	border
c	1	noise
d	2	noise
e	3	core
f	2	noise
g	3	core
h	2	border
i	1	noise
j	3	core

HDBSCAN

- CORE DISTANCE OF x : max distance of x to its k -nearest neig.
- MUTUAL REACHABILITY DISTANCE : $d_{\text{reach}}(a, b) = \max \{ \text{core}_k(a), \text{core}_k(b), d(a, b) \}$

$$k=3$$



$$\begin{aligned} \text{core}_3(A) &= 2 \\ \text{core}_3(B) &= 3 \\ d(A, B) &= \sqrt{29} = 5.38 \\ d_{\text{reach}}(A, B) &= \sqrt{29} = 5.38 \end{aligned}$$

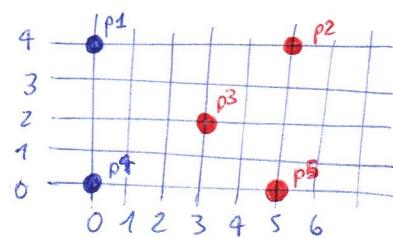
SILHOUETTE COEFFICIENT

- For every point x_i :
- a_i = average distance with points inside x_i 's cluster
 - b_i = average distance with points of the closest cluster
- $$s_i = \frac{b - a}{\max\{a, b\}}$$

Then average on i .

distance: Euclidean

	x	y	label
p1	0	4	0
p2	5	4	1
p3	3	2	1
p4	0	0	0
p5	5	0	1



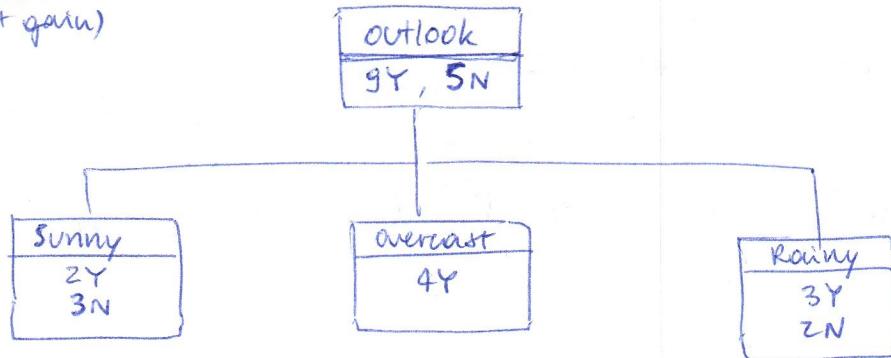
$$\begin{aligned} \rightarrow p1) \quad a &= d(p1, p4) = 4 \\ p2) \quad a &= \frac{1}{2} (d(p2, p3) + d(p2, p5)) = 2 + \sqrt{2} = 3.4 \\ p3) \quad a &= \frac{1}{2} (d(p3, p2) + d(p3, p5)) = 2\sqrt{2} = 2.83 \\ p4) \quad a &= a(p1) \\ p5) \quad a &= a(p2) \end{aligned} \quad \begin{aligned} b &= \frac{1}{3} (d(p1, p2) + d(p1, p3) + d(p1, p5)) = 5 \\ b &= \frac{1}{2} (d(p2, p1) + d(p2, p4)) = 5.7 \\ b &= \frac{1}{2} (d(p3, p4) + d(p3, p1)) = 3.6 \\ b &= b(p1) \\ b &= b(p2) \end{aligned}$$

	a	b	s_i
p1	4	5	0.2
p2	3.4	5.7	0.4
p3	2.83	3.6	0.21
p4	4	5	0.2
p5	3.4	5.7	0.4

$$s_{\text{TOT}} = \frac{\sum_i s_i}{n} = 0.282$$

DECISION TREE

(gain \rightarrow we take the attribute with highest gain)



1. Entropy before splitting

$$\text{info}([9, 5]) = \text{entropy}([9/14, 5/14]) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.94$$

2. Informations for each value:

$$\begin{aligned} \text{sunny} &: \text{info}([2, 3]) = \text{entropy}([2/5, 3/5]) = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right) = 0.971 \\ \text{overcast} &: \text{info}([4, 0]) = \text{entropy}([1, 0]) = 0 \\ \text{rainy} &: \text{info}([3, 2]) = \text{info}([2, 3]) = 0.971 \end{aligned}$$

3. Entropy after the splitting

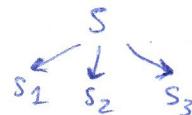
$$\begin{aligned} \text{info(split)} &= \frac{5}{14} \text{info}_{\text{sunny}} + \frac{4}{14} \text{info}_{\text{overcast}} + \frac{5}{14} \text{info}_{\text{rainy}} \\ &= \frac{5}{14} \cdot 0.971 + \frac{4}{14} \cdot 0 + \frac{5}{14} \cdot 0.971 = 0.693 \end{aligned}$$

4. Gain

$$\text{gain}(\text{outlook}) = \text{entropy}_{\text{before}} - \text{entropy}_{\text{after}} = 0.94 - 0.693 = 0.247$$

Intrinsic information of a split:

$$\text{IntrinsicInfo}(S, \text{split}) = - \sum_i \frac{|S_i|}{|S|} \log_2\left(\frac{|S_i|}{|S|}\right)$$



Information gain ratio:

$$\text{GainRatio}(S, \text{split}) = \frac{\text{gain}(S)}{\text{IntrinsicInfo}(S, \text{split})}$$

$$gini(D) = 1 - \sum_{j=1}^c p_j^2$$

$c = \# \text{ classes in } D$

$p_j = \% \text{ of elements of class } j$

$$gini(\text{split}) = \sum \frac{|D_i|}{|D|} gini(D)$$

$$\Delta gini(\text{split}) = gini(D) - gini(\text{split})$$

We look for the lower gini:

Split: {sunny}, {overcast, rainy}

$$1. gini(\text{outlook}) = 1 - (9/14)^2 - (5/14)^2 = 0.46$$

$$2. gini(\text{sunny}) = 1 - (2/5)^2 - (3/5)^2 = 0.48$$

$$gini(\text{overcast}) = 1 - (4/9)^2 - (1/9)^2 = 0.364$$

$$3. \underline{gini(\text{split})} = \frac{5}{14} 0.48 + \frac{9}{14} 0.364 = 0.394$$

$$4. \underline{\Delta gini(\text{split})} = gini(\text{outlook}) - gini(\text{split}) \\ = 0.46 - 0.394 = 0.066$$

DECISION TREE - NUMERICAL

1. sort the values (with class labels)
2. check feasible cut points
3. choose the one with best information gain

those that have a change in the output

Ex.

1	1	2	3	3	3	3	4	4	4	4	4	4
N	Y	N	Y	N	N	Y	N	Y	Y	Y	Y	N

(1) (2)

(3)

6Y
6N

1. Entropy before :

$$\text{info}(D) = \text{info}([6, 6]) = \text{entropy}([Y_2, Y_2]) = 1$$

2. Informations of the values

$$(1) \quad \text{info}(\leq 1) = \text{info}([1, 1]) = \text{entropy}([Y_2, Y_2]) = 1$$

$$\text{info}(> 1) = \text{info}([5, 5]) = \text{entropy}([Y_2, Y_2]) = 1$$

$$(2) \quad \text{info}(\leq 2) = \text{info}([1, 2]) = \text{entropy}([Y_3, 2/3]) = -\frac{1}{3} \log_2(Y_3) - \frac{2}{3} \log_2(2/3) = 0.918$$

$$\text{info}(> 2) = \text{info}([5, 4]) = \text{entropy}([5/8, 4/8]) = 0.981$$

$$(3) \quad \text{info}(\leq 3) = \text{info}([3, 4]) = \text{entropy}([3/7, 4/7]) = 0.985$$

$$\text{info}(> 3) = \text{info}([3, 2]) = \text{entropy}([3/5, 2/5]) = 0.971$$

3. Entropy after

$$\text{info}((1)) = \frac{1}{6} \cdot 1 + \frac{1}{6} \cdot 1 = 1$$

$$\text{info}((2)) = \frac{1}{4} \cdot 0.918 + \frac{3}{4} \cdot 0.981 = 0.973$$

$$\text{info}((3)) = \frac{7}{12} \cdot 0.985 + \frac{5}{12} \cdot 0.971 = 0.979$$

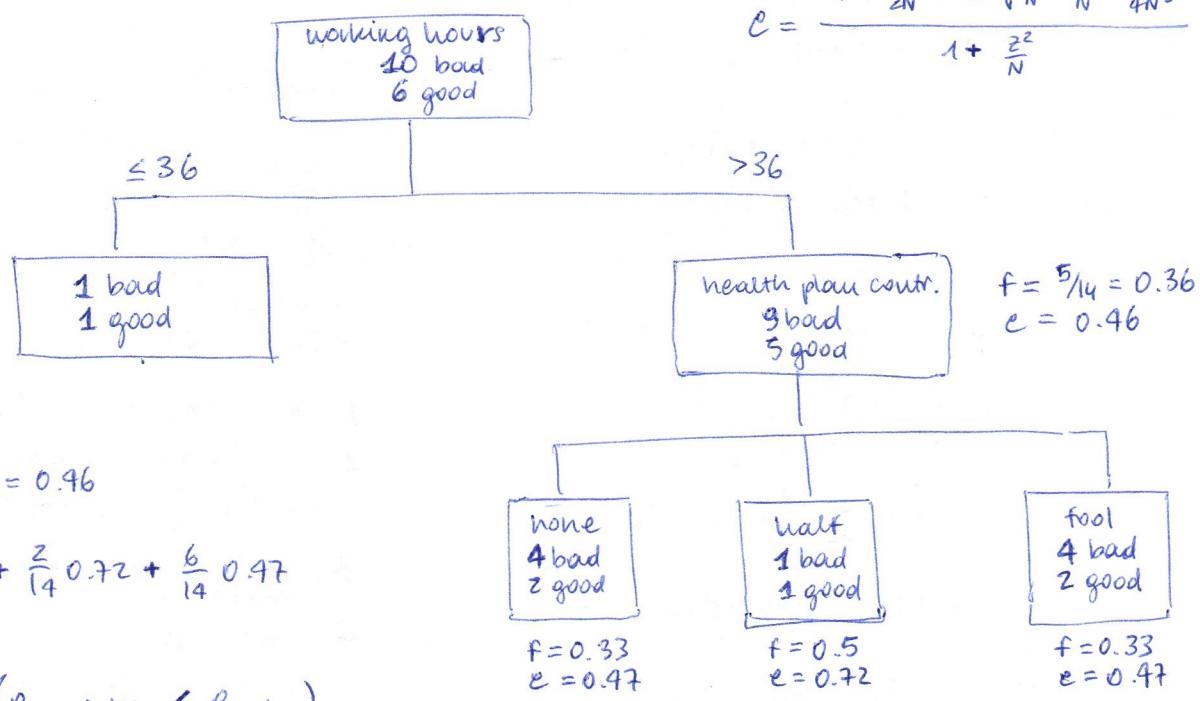
4. Gain

$$\text{gain}((1)) = 1 - 1 = 0$$

$$\text{gain}((2)) = 1 - 0.973 = 0.027 \quad \leftarrow$$

$$\text{gain}((3)) = 1 - 0.979 = 0.021$$

PRUNING



ADABOOST (ENSEMBLE)

Ex.

Dataset :

x	0	1	2	3	4	5	6	7	8	9
y	+1	+1	+1	-1	-1	-1	-1	-1	+1	+1

Initial weights \rightarrow uniform :

w	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
x	0	1	2	3	4	5	6	7	8	9
y	+1	+1	+1	-1	-1	-1	-1	-1	+1	+1

$$\varepsilon_i = \sum_j w_j \frac{1}{2} \delta_{y_j \neq h_i(x_j)}$$

$$\alpha_i = \frac{1}{2} \ln \left(\frac{1 - \varepsilon_i}{\varepsilon_i} \right)$$

$$H(x) = \text{sign} \left(\sum_i \alpha_i h_i(x) \right)$$

We bootstrap the data with the assigned weights.

We generate h_1 that returns +1 if $x \leq 2$

w	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
x	0	1	2	3	4	5	6	7	8	9
y	+1	+1	+1	-1	-1	-1	-1	-1	+1	+1
y-pred	+1	+1	+1	-1	-1	-1	-1	-1	-1	-1
									x	x

$$\Rightarrow \varepsilon_1 = \frac{2}{10}, \quad \alpha_1 = 0.69$$

The new weights are : $w_2 = w_1 e^{-\alpha_1}$ if correct, $w_2 = w_1 e^{\alpha_1}$ if wrong

We bootstrap and we generate h_2 that returns +1 if $x \geq 8$

w	0.0625	0.0625	0.25	0.25
x	0	1	2	3	4	5	6	7	8
y	+1	+1	+1	-1	-1	-1	-1	-1	+1
y-pred	-1	-1	-1	-1	-1	-1	-1	-1	+1
	x	x	x						

$$\Rightarrow \varepsilon_2 = 0.0625 + 0.0625 + 0.0625 = 0.1875$$

$$\alpha_2 = 0.7332$$

The new weights are : ..

The final model is :

$$H(x) = \text{sign} (0.69 h_1(x) + 0.73 h_2(x))$$

TEXT MINING

$$f(q, d) = q_1 d_1 + \dots + q_n d_n = q_1 (w_1 d_1) + \dots + q_n (w_n d_n)$$

$q = "news about presidential campaign"$

$d_1 = \dots news \dots$

$d_2 = \dots news \dots campaign \dots$

$d_3 = \dots news \dots presidential \dots campaign \dots$

$d_4 = \dots news \dots presidential \dots campaign \dots presidential \dots$

$d_5 = \dots news \dots campaign \dots campaign \dots campaign \dots campaign$

$$f(q, d_1) = 2, 2$$

$$f(q, d_2) = 3, 3$$

$$f(q, d_3) = 3, 3$$

$$f(q, d_4) = 3, 4$$

$$f(q, d_5) = 2, 5$$

count of
matchings

Count of
matchings
+ Frequency

To penalize words that frequently occur in many documents:

$$\text{IDF}(w) = \log_2 \left(\frac{M}{k} \right)$$

$M = \# \text{ documents in the corpus}$

$k = \# \text{ documents that contain word } w$

$c(w_i, d) = \text{frequency of word } w_i \text{ in } d$

$$d_i = c(w_i, d) \cdot \text{IDF}(w_i)$$

$$q_i = c(w_i, q) \cdot \text{IDF}(w_i)$$

$$f(q, d) = \sum_{i=1}^N d_i q_i = \sum_{w \in \{q\} \cap d} c(w, d) \cdot c(w, q) (\text{IDF}(w))^2$$

We are given 40 documents:

- in 20 there is the word "set"
- in 10 there is the word "computing"
- in 5 there is the word "mining"

Which keyword is the most important according to IDF?

$$\text{IDF}("set") = \log_2 \left(\frac{40}{20} \right) = 1$$

$$\text{IDF}("computing") = \log_2 \left(\frac{40}{10} \right) = 2$$

$$\text{IDF}("mining") = \log_2 \left(\frac{40}{5} \right) = 3$$

most important since it appears in fewer documents

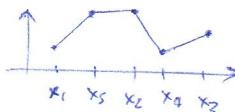
FEATURES SELECTION

- embedded approaches: feature selection is part of the algorithm (e.g. Lasso)
- filter approaches: features are selected \Downarrow specific problem
- wrapper approaches: apply data mining algorithms as black box to find the best subset of features

FILTER APPROACHES:

- reduced variance feature selection:
eliminate the variables that have a variance $<$ threshold
(low variance \rightarrow similar feature values \rightarrow not informative (?)
(however a low-variance feature may be important)
- univariate feature selection:
select the best features based on a statistical test
(stat. test: how much does a variable explain the target (\sim supervised method))
- PCA (normalize!):
features that best capture the variance of the data
(remember: we use only variables, not the target)
- recursive features elimination:
we start with all and eliminate one-at-the-time the ones that bring less information about the target

WRAPPER APPROACHES: (wrapped around target)

- random forests:
random forest has a function that, based on the forest, returns a score for each feature/variable (importance)
 - hill climbing:
 - start from a random set of features
 - provide an evaluation of the performance (algorithm-dependent)
 - generate a perturbation in the solution
 - evaluate the performance of the perturbed sol.
 - take the best
 - genetic algorithms:
as hill climbing but with "genetic variation"
 - permutation importance
(how important is a single feature)
 - train the model on a dataset
 - shuffle the values in a single column
 - apply the model to both datasets
 - Feature importance = loss in performance in shuffled data
 - partial dependence plot
(how the target depends on average on one feature)
 - variable $X = \{x_1, \dots, x_n\}$
 - train a model
 - for each $x_i \in X$:
 - substitute x_i in all the columns
 - predict the output of the model for each sample
 - average the outputs
 - plot:
- 
- } how much a feature is important for the model, not the target

DISCRETIZATION? continuous \rightarrow nominal

- supervised: use informations about the target to generate intervals that preserve relevant informations (e.g. decision trees)
- unsupervised: do not \uparrow
 - equal width intervals
 - equal frequency intervals
 - k-means (unsupervised clustering)

PCA? Standardize (StandardScaler().fit_transform(X)) → Standardize almost always (PCA, clustering, regression, ...)

Exploration?

- barplot for nominal
- heatmap for correlation: darker → + correlated (~1) (.heatmap())
- clustermap : - annot=True → clusters of columns (attributes)
- \emptyset → clusters of both attributes and values (hierarchical)
- distplot: histogram (if kde=False: no kernel density estimation) / map()
- scatterplots (.pairplot())
- scatterplots + histograms (.jointplot())
- boxplot (if jitter=True → points) / violinplot
- missing values (drop/mean/"None")
- skewness
- one-hot-encoding categorical / ordinal instead of categorical
- Feature creation from date stamps
- data.describe() → mean, std, min, 25%, 50%, 75%, max

Clusters?

- inconsistency = $\frac{h - avg}{std}$ compares the current high h of the merge with the average height of the previous merges
- acceleration: we can search the highest acceleration of merge distance growth (in terms of distance)

Internal measures

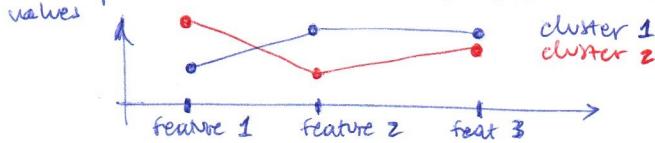
→ WSS/BSS

/ External measures

→ already existing labels

Represent clusters?

Snake plots of standardized values



K-MEANS ASSUMPTIONS

- problems with:
- wrong k
 - different variance
 - unevenly sized clusters
 - different densities

Connectivity constraints

constraints that we can put on hierarchical agg. clustering.
for example we can say:

"a point can be merged only with points in the 10 nearest neig."
This prevents to merge clusters that are too far away.

SILHOUETTE

measures cohesion and separation, it's based on the difference between the average distance to points in the closest cluster and to points in the same cluster

$$s_i = \frac{b_i - a_i}{\max\{b_i, a_i\}}$$

$b_i = \text{avg } d(x_i, \text{points of the closest cluster})$
 $a_i = \text{avg } d(x_i, \text{points in the same cluster})$

$$s_i \sim \begin{cases} +1 & x_i \text{ close to its cluster points more than the closer cluster's points} \\ 0 & x_i \text{ close to boundary} \\ -1 & x_i \text{ closest to the closer cluster's points than its own} \end{cases}$$

ENSEMBLES

- generate models difficult to analyze
- provide ways to score the variables used by the models (.feature-importances-)

TIME SERIES

- trend : long term increase/decrease
- seasonal pattern: if \exists seasonal factors
- cyclic: rises/falls that are not of fixed frequency
- time series: additive/multiplicative
- $\cdot \text{head}()$, $\cdot \text{describe}()$, $\cdot \text{describe}(\text{include} = [\text{np.objects}])$, $\cdot \text{info}()$, $\cdot \text{plot}()$
- sample at a given frequency (! N)
- filling strategies:
 - pad/fill : use last valid obs
 - backfill/bfill : use next valid obs
 - fill-value : we put the value
- shift:
 - .shift : shifting ahead of 1 value in the index
 - .shift(n) : shift ahead of n (if $n < 0 \Rightarrow$ shift back)
 - asfreq('3D').shift(n) : shift ahead of n but take the frequency of 3 days
- percentage change: $\frac{x_t - x_{t-1}}{x_{t-1}}$ (pct_change())

- correlation between pairs of series (CF)
measures how much two series vary together

$$r_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sigma_x \sigma_y}$$

- high positive $\rightarrow (\downarrow\downarrow), (\uparrow\uparrow)$
high negative $\rightarrow (\downarrow\uparrow), (\uparrow\downarrow)$
low \rightarrow weak association

autocorrelation: (.auto corr()) (ACF)

$$r_k = \text{corr}(x_t, x_{t+k})$$

: how a series is correlated with a lagged copy of itself
(any substantial non-zero autocorrelation implies that series can be forecast from past)
 \rightarrow we can plot it with correlograms

- random noise: $E = 0, \sigma^2 = 1$, zero autocorrelation
- random walk:
- with drift: $P_{t+1} = p_t + \varepsilon_t$
- $P_{t+1} = \mu + p_t + \varepsilon_t$

} unpredictable

$$P_{t+1} = \alpha + \beta P_t + \varepsilon_t$$

$H_0: \beta = 1 \rightarrow$ random walk
 $H_1: \beta < 1 \rightarrow$ not

$$P_{t+1} - P_t = \alpha + \beta P_t + \varepsilon_t$$

$H_0: \beta = 0 \rightarrow$ r.w. (Dickey-Fuller)
 $H_1: \beta < 0 \rightarrow$ not

- stationary: statistical properties are constant over time

(zero trend, const variance, const autocorrelation)

non-stationary? \rightarrow diff (first difference) (r.w. non-stat, its diff. stationary)

statistical test: augmented Dickey-Fuller : $H_0: \text{non-stat.}$

- ARMA(p,q) : $R_t = \mu + \phi R_{t-1} + \varepsilon_t + \theta \varepsilon_{t-1}$ ($p=q=1$)

today_return = mean + ϕ yesterday_return + today_noise + θ yesterday_noise

- ARMAX(1,1): $R_t = \mu + \phi R_{t-1} + \varepsilon_t + \theta \varepsilon_{t-1} + \gamma d_{t-1}$

external variable
to bring extra (11 series) info
 \leftarrow non-stat

- ARIMA(1,1,1): $\Delta R_t = \alpha \Delta R_{t-1} + \varepsilon_t + \beta \varepsilon_{t-1}$

$$\Delta R_t = R_t - R_{t-1}$$

- Diagnostic : analysis of residuals (real values - predicted) \sim white noise?

- plot
- histogram
- qqplot gaussian
- correlogram

- Model selection :
 - partial autocorrelation function : evaluates the benefit of adding lags
 - Bayesian information criterion :
- $$BIC = T \log\left(\frac{SSE}{T}\right) + (k+2) \log(T)$$
- $T = \# \text{ obs for the estimate}$
 $k = \# \text{ model parameters}$
- auto.arima(data)
 - BOX-JENKINS
 1. identification : stationary? \rightarrow transform , AR/MA ..?
 2. estimation : data \rightarrow estimate parameters
 3. diagnostic
 4. decision : good model?
 - Evaluation metrics :
 - single error: $e_t = y_t - \hat{y}_t$
 - Mean Absolute Error: $MAE = \frac{1}{N} \sum |e_t|$
 - Root Mean Square Error: $RMSE = \sqrt{\frac{1}{N} \sum e_t^2}$
 - Mean Absolute percentage error: $MAPE = \frac{1}{N} \sum |p_t| , p_t = 100 \frac{|e_t|}{y_t}$