<u>Bag of words</u>: unordered collection of words {"John":1, "likes":7, "viagra":3}
<u>Term weightening</u>: tf-idf → How important a word is to a document. tf = # in document
$idf(t,D) = \log[|D|/|\{d \in D : t \in d\}|]$ → D = all docs ⇒ $tf\text{-}idf(t,d,D) = tf(t,d) \cdot idf(t,D)$ [Stop-Words]
<u>Similarity of vectors</u>: $sim(q_k, d_J) = \sum(t_{ik} \cdot t_{iJ})$. <u>Multinomial</u>: events are independent...
Doc is multinom outcome of words. $P(W_1=n_1, ..., W_k=n_k \mid N, \theta_1, ..., \theta_k) = \frac{N!}{n_1! \cdots n_k!} \cdot \theta_1^{n_1} \cdots \theta_k^{n_k}$
<u>Classification</u>: <u>KNN</u>: needs a distance metric. Determine K. → Majority vote...
between the K nearest. Possible weighting: $1/dist^2$. All computation deffered until classification...
A popular lable may dominate due to it's quantity → try to weight closest items to overcome.

<u>Association Rule Mining</u>: {Onions, potatos} ⇒ {burger}. $I = \{i_1, ..., i_n\}$ items (boolean)
$D = \{t_1, ..., t_m\}$ transactions. Rule: $X \Rightarrow Y$ where $X, Y \in I$ and $X \cap Y = \emptyset$. MinSup, MinConf.
<u>Support</u> $Supp(X)$ = Proportion of transactions containing itemset X. <u>Confidence</u> $Conf(X \Rightarrow Y) =$
$= Supp(X \cup Y)/Supp(X) = P(Y|X)$. <u>Lift</u> $= lift(X \Rightarrow Y) = Supp(X \cup Y)/(Supp(X) \cdot Supp(Y))$
Lift > 1 = מתאם חיובי, Lift < 1 = שלילי, Lift = 1 = חוסר תלות. <u>Apriori Principle</u>: if an itemset is frequent
then all of its subsets must also be frequent. {A,B,C,D} $Conf(ABC \rightarrow D) \geq Conf(AB \rightarrow CD) \geq ... (A \rightarrow \frac{BC}{D})$
$Lift(X \rightarrow Y) = P(Y|X)/P(X)$

→ $D(A,B) \leq D(A,C) + D(C,B)$. $Manhattan(a,b) = \sum(a_i - b_i)$
<u>Clustering</u>: <u>Distance Properties</u>: $D(A,B) = D(B,A)$; $D(A,A) = 0$; $D(A,B) = 0 \Rightarrow A = B$;
<u>Euclidean</u>: $d(X,Y) = \sqrt{\sum_{n=1}^{N}(x_n - y_n)^2}$; <u>Cosine Similarity</u>: $a \cdot b = |a| \cdot |b| \cdot \cos\theta$. Given $\vec{A}, \vec{B}$
Similarity $= \cos(\theta) = (A \cdot B)/(|A| \cdot |B|) = \sum_{1}^{n} A_i \cdot B_i / (\sqrt{\sum_{1}^{m} A_i^2} \cdot \sqrt{\sum_{1}^{n} B_i^2})$ range: $(-1, ..., 1)$.
<u>Jaccard Similarity</u>: $J(A,B) = |A \cap B| / |A \cup B|$. dist $= 1 - sim$... → A,B are groups.
In Dendograms: sim = height of lowest shared internal node. <u>Hierarchical Clustering</u>:
results in a dendogram. <u>K-means clustering</u>: aims to partition n observations
into K clusters. Given n observations $(X_1, ..., X_n)$ each x is d dimensional vector, partition to K sets
$S = \{S_1, ..., S_k\}$ minimize the "within cluster sum of squares" $\arg\min_S \sum_{j=1}^{S} \sum_{x_j \in S_i} |x_j - \mu_i|^2$. $\mu_i$ = mean of points in $S_i$
The algorithm alternates betweet two steps until convergence: ① Assignment step:
$S_i^t = \{x_p : |x_p - m_i^t|^2 \leq |x_p - m_J^t|^2 \; \forall \; 1 \leq J \leq K\}$ ② <u>Update step</u>: $m_i^{t+1} = \frac{1}{|S_i^t|} \cdot \sum_{x_J \in S_i^t} x_J$. Since the arithmetic
mean is a least-squares estimator, this also minimizes the → WCSS objective.
assumes clusters are of similar sizes. for different sizes the EM is better (genzation of Kmeans)
<u>K-medoids</u>: PAM ↔ Parartitioning Around Medoids: <u>Silhouette</u> - tool for determining K.
<u>Algorithm</u>: ① Initialize: randomly select K of the n data-points as the medoids
② Associate each data-point to the closest medoid. ③ for each medoid m:
   ③.1 For each non-medoid data point O ③.1.1 swap m and o and comput the
total cost of the configuration. ④ Select the configuration with the lowest cost.
⑤ Repeat steps 2-4 until there is no change in the medoid. ✱ cost = סכום המרחקים בין כל פריט למדואיד שלו
אבל <u>Silhouette (Clustering)</u>: How well each object lies within its cluster. $a(i)$ is average
distance of object i with all other objects within its own cluster. $b(i)$ is lowest
average distance of i to all other clusters. the cluster with distance $b(i)$ is the
neighboring cluster (best fit after current cluster) → $S(i) = [b(i) - a(i)]/\max\{a(i), b(i)\}$
$-1 \leq S(i) \leq 1$. 1 is great, (-1) is worst, 0 point is on the kiss-border between the neighboring clusters
<u>Agglomerative hierarchical Clustering</u>: in begining each element is a cluster of its own. ($x \in X, y \in Y$)
Single-linkage: $D(X,Y) = \min d(x,y)$. Complete linkage: $D(X,Y) = \max d(x,y)$. Average: $\frac{1}{|X| \cdot |Y|} \cdot \sum\sum d(x,y)$
$-P(yes)\log_2 P(yes) - P(no)\log_2 P(no)$ Y ה-מוצג פה, מוצאי, נכנסו האנטרופיה, זכור $-P(yes|r)\log_2 ...$ classify
<u>Decision Tree Learning</u>: Data comes in records of form $(X,Y) = (X_1, X_2, X_3, ..., X_k, Y)$
Work top-down by choosing a variable at each step that best splits the set of items. <u>Gini Impurity</u>:
is a measure of how often a randomly chosen element from the set would be incorrectly labled if it were
randomly labled according to the distribution of lables in the subset. suppose i takes on values {1,2,...,m} and
let $f_i$ be the fraction of items labled with value i in the set. $I_G(f) = \sum_{1}^{m} f_i \cdot (1-f_i) = \sum_{1}^{m}(f_i - f_i^2) =$
$= \sum_{1}^{m} f_i - \sum_{1}^{m} f_i^2 = 1 - \sum_{1}^{m} f_i^2$. <u>Information Gain</u>: $I_E(f) = -\sum_{1}^{m} f_i \cdot \log_2(f_i)$ כשהתוחלת ההפתעה מקסימלית
<u>Entropy</u>: E(S) is the measure of the amount of uncertainty in the set S. מחרת שווה/יותר מחומ
$E(S) = -P(yes|big)\log_2 P(yes|big) - P(no|big) \cdot \log_2 P(no|big)$ → big: אם הערך שלו גדול yes/no כמו
$IG(T, size) = E(T) - \frac{|S_{big}|}{|T|} \cdot E(S_{big}) - \frac{|S_{small}|}{|T|} \cdot E(S_{small})$ <u>Information Gain</u>: How much uncertainty in S was reduced
בוחרים את המאפיין after split אם יש יותר מאפיינים, בוחרים את המאפיין עם ה-IG הכי גבוה כי הוא מוריד הכי הרבה אי וודאות, כמובן
from כי את הגבוה ביותר של פיצול זה עדיף. אם IG=0 אז המאפיין לא עוזר בכלל. אחרת בוחרים אותו ומורידים
ומתחילים מהפיצול הבא. כמובן בוחרים כל פעם רק את המאפיין הכי טוב. <u>K-Medoids</u>

<u>Bayesian inference</u>: Bayes rule: $P(H|E) = \frac{P(E|H) \cdot P(H)}{P(E)}$ [likeliness / Evidence].
There are competing hypotheses from which one chooses the most probable.
<u>Naive Bayes Classifier</u>: assumes that the presence or absence of a particular
feature is unrelated to any other feature, given the class variable. The probability
model for a classifier is conditional: $P(C|F_1,...,F_n)$ ← Class; Feature. Posterior = $\frac{prior \cdot likelihood}{evidence}$
$P(C|F_1,...,F_n) = P(C) \cdot P(F_1,...,F_n|C) / P(F_1,...,F_n)$ ← Denom is Const.
$P(C|F_1,...,F_n) = \frac{1}{const} \cdot P(C) \cdot \prod_1^n P(F_i|C)$ ← Naïve assumption. All model parameters (Class
Priors and feature probability distributions) can be approximated with relative frequencies
from the training set. class prior → ① $\frac{1}{\#C}$ equally for all or ② $\#C_i/\#C$ ↔ Relative Frequency of $C_i$
For discrete features → Multinomial and Bernouli distributions are popular to assume.
Bernouli ↔ $X \sim Ber(p)$ ↔ $P(X=1)=P$. Multinomial ↔ $P(x_1,...,x_K; n, P_1,...P_n) = \frac{n!}{x_1!...x_K!} \cdot P_1^{x_1} \cdot ... \cdot P_K^{x_K}$.
For continuos features ↔ Gaussian dist. assumed. $P(X=v|C) = \frac{1}{\sqrt{2\pi\sigma_c^2}} \cdot exp(-[v-\mu_c]^2/[2\sigma_c^2])$ MAP
Sample-correction is needed when one of the fregs is 0 in training set. Maximum Aposteriori Decision Rule
<u>Classify</u> $(f_1,...,f_n) = argmax \; P(C=c) \cdot \prod P(F_i=f_i|C=c)$. $\mu = \frac{1}{N}\sum X_i$. $\sigma^2 = \frac{1}{N}\sum(X_i-\mu)^2$ $\frac{P(C)}{P(D)}$
Document Classification: $P(D|C) = \prod_i P(W_i|C)$. $P(D|C) = P(D \cap C)/P(C)$. $P(C|D) = P(D \cap C)/P(C) \Rightarrow P(C|D) = \frac{P(C)}{P(D)} \cdot P(D|C)$
Spam or ¬ Spam: $P(S|D) = \frac{P(S)}{P(D)} \cdot \prod P(W_i|S)$. Same for ¬S. $P(S|D)/P(¬S|D) = \frac{P(S)}{P(¬S)} \cdot \prod \frac{P(W_i|S)}{P(W_i|¬S)}$.
$ln[P(S|D)/P(¬S|D)] = \underrightarrow{\quad} ln\frac{P(S)}{P(¬S)} + \sum_i ln[P(W_i|S)/P(W_i|¬S)]$. Spam if $P(S|D) > P(¬S|D) \Rightarrow$
$\Rightarrow$ When $ln[P(S|D)/P(¬S|D)] > 0$. <u>LDA-Linear discriminant analysis</u>: find a linear
combination of features which characterizes or separates two or more classes of objects.
Fundemental assumption: the independent variables are normally distributed.
<u>For 2 classes</u>: Assumes $P(\vec{X}|y=0)$ and $P(\vec{X}|y=1)$ are both normally distributed
with mean and covariance $(\vec{\mu_0}, \Sigma_{y=0})$ and $(\vec{\mu_1}, \Sigma_{y=1})$. Under this assumption, the
Bayes optimal solution is to predict $y=1$ if the log of the likelihood ratios is below some
treshold $T$: $(\vec{X}-\vec{\mu_0})^T \Sigma_{y=0}^{-1}(\vec{X}-\vec{\mu_0}) + ln(\Sigma_{y=0}) - (\vec{X}-\vec{\mu_1})^T \Sigma_{y=1}^{-1}(\vec{X}-\vec{\mu_1}) - ln|\Sigma_{y=1}| < T$
Without further assumptions, the above is QDA. LDA assumes $\Sigma_{y=0} = \Sigma_{y=1} = \Sigma$ in this case
several terms cancel and the above becomes $\vec{W} \cdot \vec{X} > C$ for some threshold constant $C$, where
$\vec{W} \propto \Sigma^{-1}(\vec{\mu_1} - \vec{\mu_0}) \Rightarrow$ criterion of an input $\vec{X}$ being of class $y$ is purely a function of this
linear combination of the known observations. In geometrical terms we project a multidimention space point x
onto vector $\vec{W}$ (thus we only consider its direction). In other words, the observation belongs to
$y$ if coresponding $\vec{X}$ is located on a certain side of a hyperplane perpendicular to $\vec{W}$. the
location of the plane is defined by the threshold $c$. <u>Fisher's linear discriminant</u>: Suppose 2
Classes of observations have means $\vec{\mu}_{y=0}, \vec{\mu}_{y=1}$ and covariances $\Sigma_{y=0}, \Sigma_{y=1}$ then the linear
combination of features $\vec{W} \cdot \vec{X}$ will have means $\vec{W} \cdot \vec{\mu}_{y=i}$ and variances $\vec{W}^T \Sigma_{y=i} \vec{W}$ for $i = 0,1$
Fisher defines the separation as the ratio of the variance between the classes to the
variance within the classes: $S = \sigma_{Between}^2 / \sigma_{within}^2 = [\vec{W} \cdot \vec{\mu}_{y=1} - \vec{W} \cdot \mu_{y=0}]^2 / [\vec{W}^T \Sigma_{y=1} \vec{W} + \vec{W}^T \Sigma_{y=0} \vec{W}]$
$= [\vec{W}(\vec{\mu}_{y=1} - \vec{\mu}_{y=0})]^2 / [\vec{W}^T(\Sigma_{y=0} + \Sigma_{y=1})\vec{W}]$ optimal when $\vec{W} \propto (\Sigma_{y=0} + \Sigma_{y=1})^{-1}(\vec{\mu}_{y=1} - \vec{\mu}_{y=0})$
$\vec{W}$ is <u>normal</u> to the discriminant hyperplane. If projections of points from both classes have
approx same distr, a good choice of $C$ will be a hyperplane between the two means:
$C = \vec{W} \cdot \frac{1}{2}(\vec{\mu}_{y=0} + \vec{\mu}_{y=1}) = \frac{1}{2}\vec{\mu}_{y=1}^t \Sigma^{-1}\vec{\mu}_{y=1} - \frac{1}{2}\vec{\mu}_{y=0}^t \Sigma^{-1}\vec{\mu}_{y=0}$.
<u>Logistic Regression</u>: measures the relationship between a categorical dependant variable and
one or more independant variables (usualy continuos), by using probability score as the predicted
values of the dependent variable. Logistic regression is used to predict the <u>odds</u> of result success
based on the values of the independant variables (predictors). The odds are defined as the probabil
that a particular outcome is success devided by the probability that it's failure. <u>Logistic Function</u>:
$F(t) = e^t/(e^t+1) = 1/(1+e^{-t}) \in [0...1]$ where $t$ is a linear function of an explanatory variable $X$.
It can be written as: $\pi(x) = exp(\beta_0 + \beta_1 x)/[exp(\beta_0 + \beta_1 x) + 1]$ ⇐ Probability for success.
The inverse of the logisic function: $g(x) = ln[\pi(x)/(1-\pi(x))] = \beta_0 + \beta_1 X$
נמצאים ברווח של שוליים... ואפשר להכניס נקודות "מחיר" אל תוך הרווח. $C$ קובע כמה. Soft Margin SVM נתון: סיווג בינארי של $m$ דוגמאות.
$\vec{x_i}$ ו$y_i$ את הוקטור מסווג $K$ בוחרים נניח. אחרי שבחרנו $K$ נגדיר ווקטור... שמכיל את המשקל
$\vec{x_i}$ ורוצים שאותו ווקטור $d$ ובכל נקודה בסיווג הגיאומטריה. ה... את ... ומגדירים את ...
$P_e = \int P(X|yes) \cdot P(yes) dx + \int P(X|no) \cdot P(no) dx$ בעיית אפשר של הסתברות שגיאה: Bayes סיווג
Ryes <u>ANN-Artificial Neural Networks</u>: <u>perceptrons</u>: output $\{-1,1\}$
$O(X_1,...,X_n) = 1$ if $W_0 + W_1 X_1 + W_2 X_2 + ... + W_n X_n > 0$ and $(-1)$ otherwise. $(-W_0 \cdot 1)$ is the threshhold.
We imagine $X_0 = 1 \Rightarrow O(X_1...X_n) = 1$ if $\sum W_i X_i > 0 \Rightarrow O(\vec{X}) = sgn(\vec{W} \cdot \vec{X})$. Perceptron training rule:
$W_i = W_i + \Delta W_i$ where $\Delta W_i = \eta(t-0)X_i$ $t$=target; $0$=output; $\eta$=learning rate-must be small! (step)

$W \leftarrow \vec{W} + \Delta W$ . where $\Delta W = -\eta \nabla E(\vec{W})$. $\Delta W_i = -\eta \cdot \frac{\partial E}{\partial W_i}$ . $\boxed{\frac{\partial E}{\partial W_i} = \sum_j (t_d - O_d) \cdot (-x_{id})}$ $_{d \in D}$ (where $\sigma(y) = 1/(1 + exp(-y))$]
When we need to ensure a differentiable threshold unit : $O = \sigma(\vec{W} \cdot \vec{X})$ .
<u>Back-Propagation</u> : Multilayer-Network. We redefine E to sum the errors over all of the network output units : $E(\vec{W}) \equiv \frac{1}{2} \sum_{d \in D} \sum_{k \in output} (t_{kd} - O_{kd})^2$ . input/weight from $i$ to $J$ is denoted as $x_{Ji}, W_{Ji}$ . ⊜ for each $\langle \vec{x}, \vec{t} \rangle$ in training examples Do : ① input the instance $\vec{X}$ to the network and Compute the output $O_u$ for every unit $u$ in the network. ② Propagate the errors back through the network by : For each network output unit $K$, calculate its error term $\delta_K$
$\delta_K \leftarrow O_K (1 - O_K)(t_K - O_K)$ ③ For each hidden unit $h$, calculate its error term $\delta_h$
$\delta_h \leftarrow O_h (1 - O_h) \cdot \sum_{k \in outputs} W_{kh} \delta_K$ ④ Update each network weight $W_{Ji} \leftarrow W_{Ji} + \Delta W_{Ji}$ where
$\Delta W_{Ji} = \eta \delta_J x_{Ji}$ . מנטרפ ים <u>The EM Algorithm</u> : Used when only a subset of the relevant instance features are observable. Example : finding the means $\mu_1, \mu_2$ of $K = 2$ Gaussians where $\sigma_i^2$ are known... The task is to output hypothesis $h = \langle \mu_1 ... \mu_K \rangle$ We want $h$ that maximizes $P(D|h)$. We can think of the full description of each instance as $\langle x_i, z_{i1}, z_{i2} \rangle$ where $x_i$ is the observed value and $z_j$ are indicators of wher $x_i$ came from (which gaussian) $z_j$ are hidden variables. EM algorithm searches for a maximum likelihood hypothesis by repeatedly re-estimating the expected values of the hidden variables $z_{iJ}$ given the current hypothesis $\langle \mu_1, ..., \mu_K \rangle$ then recalculating the ML-hypothesis using these expected values for the hidden variables. First set $h = \langle \mu_1, \mu_2 \rangle$ arbitrarly. Then, iterativly re-estimate $h$ by repeating 2-steps until convergence. Step 1 : Calculate expected value $E[z_{iJ}]$ for each hidden variable $z_{iJ}$ assuming $h = \langle \mu_1, \mu_2 \rangle$. Step 2 : Calculate a new ML-hypo $h' = \langle \mu_1', \mu_2' \rangle$ assuming value taken by each hidden variable $z_{iJ}$ is $E[z_{iJ}]$ calculated in step 1, then replace $h$ with $h'$ and iterate. In our example, $E[z_{iJ}]$ is just the probability that instance $x_i$ was generated by the $J$th Gaussian : $E[z_{iJ}] = \frac{P(X = x_i | \mu = \mu_J)}{\sum_{n=1}^{2} P(X = x_i | \mu = \mu_n)}$
$= exp(-\frac{1}{2\sigma^2} \cdot (x_i - \mu_J)^2) / \sum_{n=1}^{2} exp(-\frac{1}{2\sigma^2}(x_i - \mu_n)^2)$ .
Step 2 : $\mu_J \leftarrow \frac{1}{m} \sum_i^m E[z_{iJ}] \cdot x_i$ . Converges to a local maximum likelihood hypothesis for $\mu_1, \mu_2$ in general case : $X = \{x_1, ..., x_m\}$ observed, $Z = \{z_1, ..., z_m\}$ unobserved, $Y = X \cup Z$. $h$ current $h'$ revised. EM searches for ML-hypo $h'$ by seeking $h'$ that maximizes $E[\ln P(Y|h')]$. We define a function $Q(h'|h) = E[\ln P(Y|h')|h, X]$. Step 1 : Estimation (E) step : Calculate $Q(h'|h)$ using the current hypo $h$ and observed data $X$ to estimat prob-dist over $Y$. $Q(h'|h) \leftarrow E[\ln P(Y|h')|h, X]$ Step 2 : Maximization (M) step : replace $h$ with $h'$ that maximizes this $Q$ function $h \leftarrow \underset{h'}{argmax} Q(h'|h)$

<u>שיטת EM</u> : Unsupervised. c שיטה לא-מונחית, מניחה שהנתונים נמצאו מ א'ב גאוסיאני כאשרה הפרמטרים שלו אינם כזהוי ותמצאה ה'כ בגבלאנט
$N(X|\mu_K, \Sigma_K) = \frac{1}{(2\pi)^{d/2} \cdot |\Sigma_K|^{1/2}} \cdot exp(-\frac{1}{2}(X - \mu_K)^T \cdot \Sigma_K^{-1}(X - \mu_K))$ הפרמטרים של הגאוסיאן. נרצה $\gamma_{ik}$ כי לא נ'כ בשיבורי שאנקפקד כל פעם 1 את התצפית $x_i$ לקוחה מהרכיב ה-K הסתברות רכיבה באופן שאנקפקד כל $\gamma_{ik}$ נתון בגבהוב $\gamma_{ik} \in [0, 1], \sum_K \gamma_{ik} = 1$
$P(x_i) = \sum_K P(x_i, r_{ik}) = \sum_K P(r_{ik}) \cdot P(x_i|r_{ik}) = \sum_K \pi_K N(x_i|\mu_K, \Sigma_K)$. $\gamma(r_{ik}) = P(r_{ik} = 1|x_i) = $ :נחשב
$= \frac{P(r_{ik}) \cdot P(x_i|r_{ik})}{\sum_{J=1}^{K} P(r_{iJ}) \cdot P(x_i|r_{iJ})} = \frac{\pi_K N(x_i|\mu_K, \Sigma_K)}{\sum_{J=1}^{K} \pi_J N(x_i|\mu_J, \Sigma_J)}$ $\}\rightarrow$ $\gamma(r_{ik})$ נקראת ה-Posterior הוא כי $\gamma(r_{ik})$ Responsabl-שלו את מידה כמה הרכיב ה-K לקיחת להתהסתבירי את התצפית ה-$i$ .
$\{\pi_K, \mu_K, \Sigma_{ik}\}_{K=1}^{K}$ הם הפרמטרים את נרצה למצוא באופן ש'גדיל את ה Dataset log-likelihood : $\ln P(X|\pi, \mu, \Sigma) = \sum_{i=1}^{N} \ln \{\sum_K \pi_K N(x_i|\mu_J\} $ $K = \underset{K}{argmax} \gamma(r_{iJ})$ : שמקיים K-ה Cluster-ל שייך $x_i$ נשבץ את התצפית
כללי האלגוריתם ה-EM ⓛ מאתחלים את הפרמטרים הגאוסיאנים. ⓶ שלב (E-step) מחשב לכל אחד מהם כי הסכום הוא. ⓷ $\gamma(r_{ik}) = \frac{\pi_K N(x_i|\mu_K, \Sigma_K)}{\sum_{J=1}^{K} \pi_J N(x_i|\mu_J, \Sigma_J)}$ (M-step) מחשב ערכי הפרמטרים החדשים כאופן האומדנים המקסימלי
$\mu_K^{new} = \frac{1}{N_K} \sum_{i=1}^{N} \gamma(r_{ik}) x_i$ , $\pi_K^{new} = \frac{N_K}{N}$ $N_K = \sum_{i=1}^{N} \gamma(r_{ik})$ : התפלגות
⓸ $\Sigma_K^{new} = \frac{1}{N_K} \sum_{i=1}^{N} \gamma(r_{ik})(x_i - \mu_K^{new})(x_i - \mu_K^{new})^T$ עוצרים כאשר האלגוריתם מתכנס חוזר לשלב ה-$\delta$-2.

<u>Decision Trees</u> : $N$-Number of observations, $K$-Number of classes, $N_t$-Number of observations at node $t$, $N_t^K$-Number of observations from class $K$ at node $t$, $\hat{P}(K|t)$-proportion of observations from class $K$ at node $t$. $\hat{P}(K, t) = N_t^K / N_t$. $Y(t)$-Class assigned to the terminal node $t$. $Entropy(t) = -\sum_{K=1}^{K} \hat{P}(K|t) \log_2 (\hat{P}(K|t))$. $Gini\ Index(t) = 1 - \sum_{K=1}^{K} \hat{P}^2(K|t)$
$Misclassification\ error(t) = 1 - \underset{L=1...K}{max} \{\hat{P}(K|t)\}$, $IG(S, A) = Entropy(S) - \sum_{v \in Val(A)} \frac{|S_v|}{|S|} Entropy(S_v)$.

שיטת היה'ך/שונות : מרכבות Overfitting שונות-כ נמוך, הטיה-אם גבוהה, שונות-גבוה הטיה-אם אותה
אווריאנס=כמה המודל ישתנה אם נאמץ אותו ... הטיה=עד כמה 'סטה תחזית המודל מהנתונים... הב'ך: כאשר
על פרוס של המשוחזן. (במקרה של קבוצת הסבורים)

$P(Y=1|X=x) = h(W^Tx) = \frac{1}{1+e^{-W^Tx}}$  (W^Tx)  :לפי P(Y|X) אינו ישירות מתאים ? לוגיס רגרסיה מסווג

$P(Y=1|X)/P(Y=0|X) \geq 1 \Rightarrow \log[P(Y=1|x)/(1-P(Y=1|x))]$ כלוגר כם $P(Y=1|X=x) \geq P(Y=0...$ אם X במחלקה 1

We determine $\hat{w}$ with MLE: $L(w) = \log \prod_{i=1}^n P(Y=1|X_i, w)^{y_i} \cdot P(Y=0|X_i, w)^{1-y_i}$  שערך־ עולה כמונוט

$g(x) = 1-h$ $P(Y=1|X)=h(W^Tx)$ $\tau \geq$  אם כמה סיווגים שם כפונקציה של $X_{n+1} \leftarrow X_n - \frac{f(x)}{f'(x)}$

Specificity & Specifity יש סטטיסטי, סיבוכי, אבל ? מתסמ" . מסתדר יותר אחרי הודעות יש 2 קטגורים

Specificity= True negative/(True neg + false pos), Sensitivity = True positive/(true positive + false negative) for the

K-medoids ← ... סוגי מגבירים, הטרוכינא על טוריא, המבל סבה ... תגויות

MBA ← ... ויחס הקשר "חבר" מיעש אישי

KNN ← כלומר יבחן ... נתוני סימון ואת בסיור אחסון ... כי מיסוה כי

Gaussian Mixture ← ... ודומה ... רק תקון

Soft limit Cluster ← ... המשמת ... היא המסמרות

Hierarch-Clustering ← ...

אם: התמלה ← ... הטנרונניא ... האטרונרים של ... הטנרונריא הטכנה

$E(yes)+E(no)$. Y ... של Attr ...

size=big קיום ... כמה תוכמונות Y=yes ... size=big

2 ... של Y ... גינורי מגריב $\hat{y}$ ...

Bayes Classifier ← ... $P(1|X) > \frac{2+4}{2+4+5}$

Naive Bayes ← $\prod_i$ ... $\beta \cdot P(1|X) > \alpha \cdot P(2|X)$

$\prod_i f(x_1,x_2)$ : ...

ID3 ← ... יש ה $IG$ ...

ML/MAP ← ML ... P(Data|Model) ... MAP ... $P(\text{model params}|\text{data})$

רגרסיה לוגיסטית ← ... $\tau$ ...

Specificity= $\frac{d}{b+d}$ , Sensitivity= $\frac{a}{a+c}$   1-Specificity :

| Pred | True1 | True0 |
|------|-------|-------|
| Pred1 | a | b |
| Pred0 | c | d |

Sensitivity : Y ... ROC ...

SVM ← ...

PCA ← ...

LDA ← ... $L(Y|X)$ ... $L(Y|X)$ מקסימיא $L(Y|X)$

LDA ← מהליכות: ... נצמות  MLE ... $\log$ או $\log$

ניוטון : ... ראשית נצ' ...

① ... $X_{n+1} = X_n - f(x_n)/f'(x_n)$ ... $Y = 0-$ ... פונק:

$x_1$ →