

## Variables

- ↳ Categorical → Nominal  
→ Ordinal
- ↳ Numerical → Discrete  
→ Continuous

### KNN Classifier Implementation:

It takes test data X as input and returns predictions for each data point in X. For each test data point x, the method computes the Euclidean distance between x and each data point in the training set X\_train.

Based on these distances, it identifies the indices of the k nearest neighbors. It sorts them according to their distances until the given k value (sorted\_indices[:self.k]). It then retrieves the corresponding labels from the training labels y\_train.

A majority vote is performed among these labels to determine the predicted label for the test data point x by counting the labels in loop.

This process is repeated for all test data points, and the predictions are collected into a list, which is then returned.

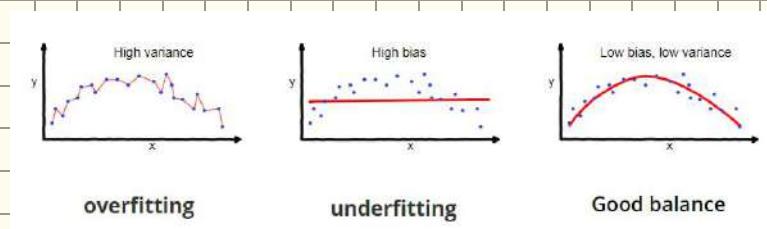
# Assessing and Comparing Classification Alg.

Input Encoding:

100 data  $\rightarrow$   $M_{100}$   $\rightarrow$  depends heavily on 100 data points

"  $\rightarrow$  SVM  $\rightarrow$  5 S.V. chosen from 100 d.p.

\* Changing data will change the model

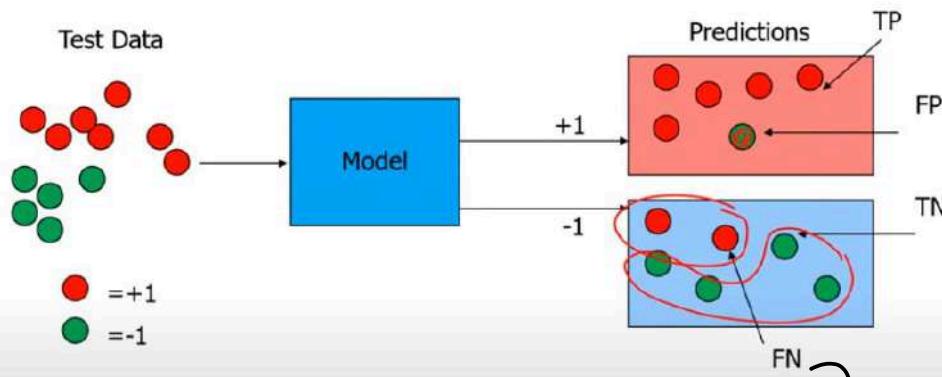


In 100 data you give someone 50 and you give other one 50 data. One of them result is higher because his data can be biased.

↓ Cross-validation

Cross-validation.

# Performance measurements -



$$\text{Accuracy} = \frac{TP + TN}{All}$$

↓  
FN sönemr kır əlavət.

$$\text{Precision} = \frac{TP}{TP + FP}$$

$\frac{\text{gercelən } +1}{\text{predicition } +1}$

+1  
Zəmənliklərinə kəm  
dəqiq?

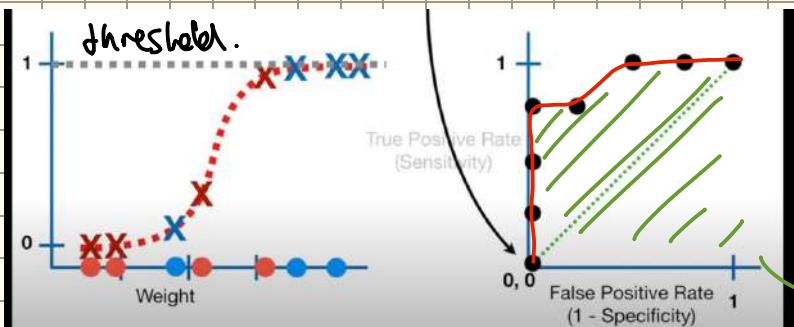
$$\text{Recall} = \frac{TP}{TP + FN}$$

$\frac{\text{gercelən } +1}{təm } +1$

} gercelən kəmərəkərək  
təməm dəqiq bildir?

F1-score

ROC-Curve



$$\text{Sensitivity} = \frac{TP}{TP + FN}$$

$$\text{Specificity} = \frac{TN}{FP + TN}$$

→ Summarizes all of the C.M.  
that each threshold produced.  
→ AUC

## Prediction vs actual performance

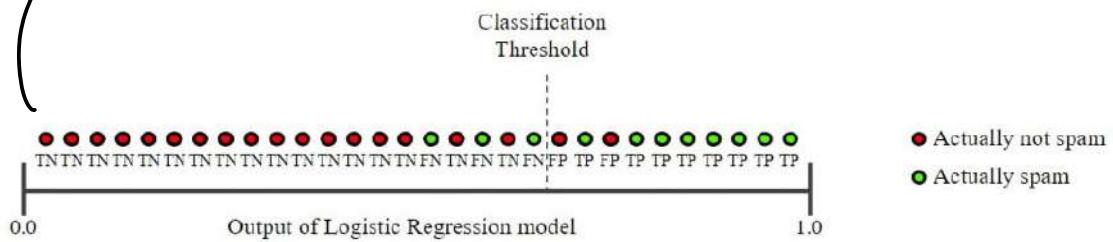


Figure 1. Classifying email messages as spam or not spam.

Let's calculate precision and recall based on the results shown in Figure 1:

True Positives (TP): 8	False Positives (FP): 2
False Negatives (FN): 3	True Negatives (TN): 17

Precision measures the percentage of **emails flagged as spam** that were correctly classified—that is, the percentage of dots to the right of the threshold line that are green in Figure 1:

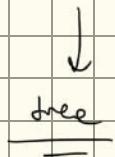
$$\text{Precision} = \frac{TP}{TP + FP} = \frac{8}{8 + 2} = 0.8$$

Recall measures the percentage of **actual spam emails** that were correctly classified—that is, the percentage of green dots that are to the right of the threshold line in Figure 1:

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{8}{8 + 3} = 0.73$$

# Decision Trees

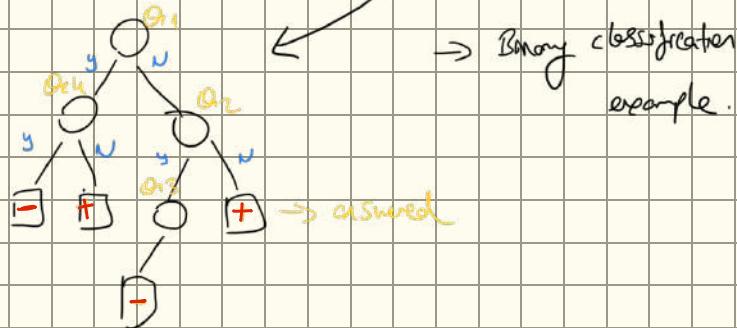
$x \rightarrow \text{Model} \rightarrow y$



cranny date

DT algorithm

" salary, savings  $\rightarrow$  credit worthiness  
not " "



→ Binary classification example.

- Can I have an algorithm to build a tree like this?

## Questions

$x \in \mathbb{R}^d$

$x_i > 10,000 \quad ?$  yes or no

↓  
salary

d-dimensional space

$x_i = \text{"March"}$

$x_i = \text{"March"} \parallel \text{"April"}$

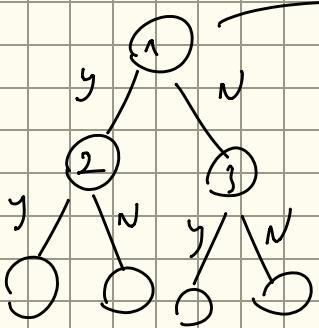
$x_i =$

:

:

- I can have infinite many questions.

- We have infinitely many possibilities.



Salary 1 M

Savings 10 M olsa, sadece root question iin brile

11 M tane sonum olacak.

(possibility)

$Q_1 \rightarrow 11M$

$Q_2 \rightarrow 11M$

$Q_3 \rightarrow 11M$

$11M^3$

impossible to solve.

We should find an optimal solution.

lets try with  
binary questions

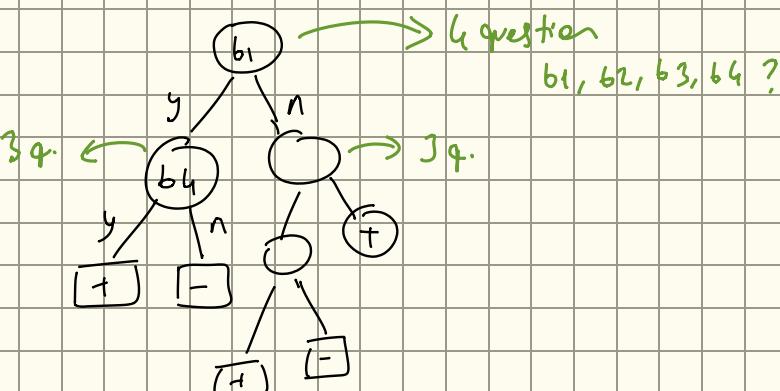
$$x_i = \{ b_1, b_2, b_3, b_4 \}$$

$\downarrow$        $\downarrow$        $\downarrow$        $\downarrow$   
 employed    children    server    citizen

try with categorical

$$\vec{x} = \{ c_1, c_2, c_3, c_4 \}$$

$\downarrow$        $\downarrow$   
 job degree ...



6 question  
 $b_1, b_2, b_3, b_4$  ?

is  $c_1$  teacher  
employer  
...  
proportional  
do date

binary / categorical  $\rightarrow$  decision  $\equiv$

integers  $\rightarrow$  decision  $\rightarrow <, =, > \dots$

real  $\rightarrow$  "  $\rightarrow <, > \dots$

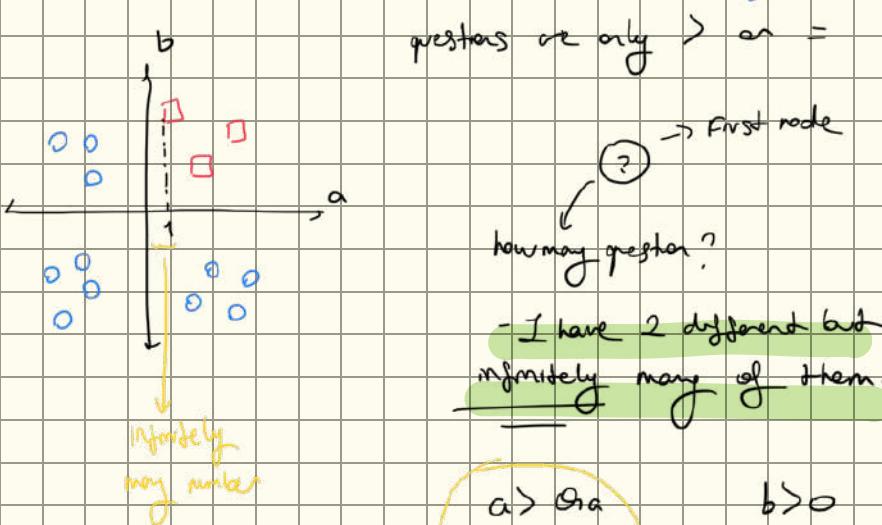
1-) What kind of questions? }  
 salary  $\rightarrow$  5000  
 savings  $\rightarrow$  10,000

2-) Limit those questions.

- For categorical only equal
- For numerical only greater.

Why limit?

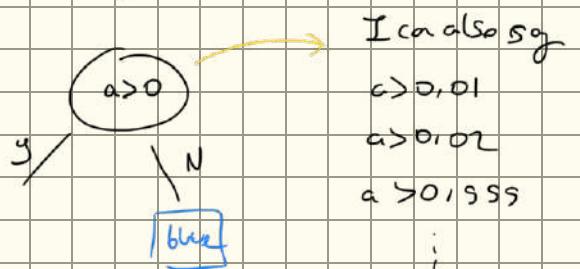
- Real world questions are NP-hard.
- (c) Helps to manage this computational complexity. Balance.



$$\begin{array}{ll} a > 0 \text{ or } a \\ b > 0 \text{ or } b \end{array}$$

$$\begin{array}{ll} b > 0 \\ b > 1 \\ b > -1 \\ \vdots \end{array}$$

lets say:

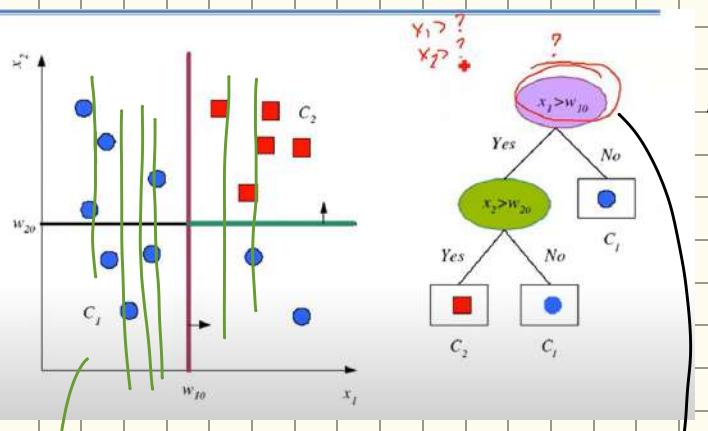


I need to look middle point

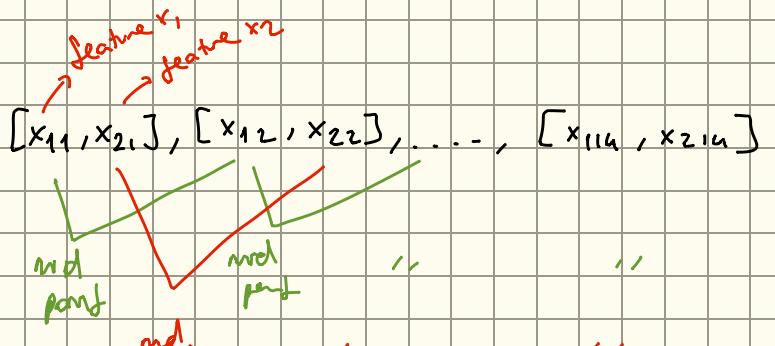
$\rightarrow$  Sort a's in ascending order.

Middle point of any pair is a candidate

for  $Q_C$ .



14 prediction  
red points



$$1^3 + 1^3 = 2^6 \text{ possibilities for } Q_{1,1} \cdot (2n-2)$$

$(n-1) \quad (n-1)$

Build trees  $\rightarrow$  "hierarchically"

Start from root, ask question

In a way that your job becomes easier in next one. Try to do this with fewer num. of questions.

Heuristic way, hope that previous choice was my best choice.

$$(2n-2) \cdot (2n-3) \Delta (2n-4)$$

$O(n^k)$

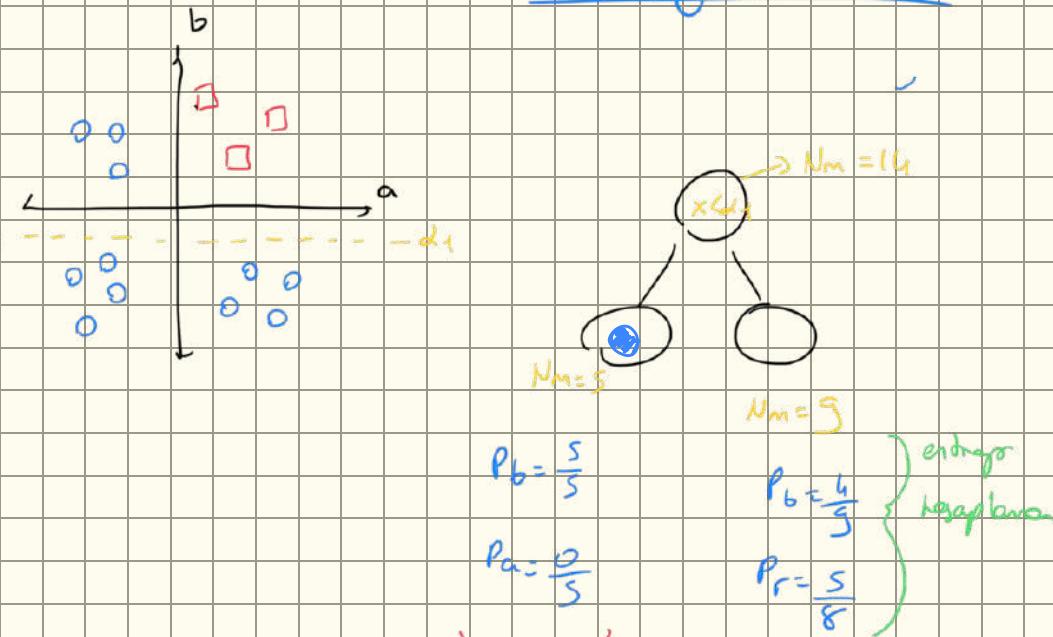
Alg  
 $= O(n^k)$

- Build all possible trees
- Assess performance of each
- Get the tree with best performance.

Choose the best split:

$$P(C_i | x_{im}) = p_{im}^i = \frac{N_m^i}{N_m} \rightarrow \begin{array}{l} \text{num of blue points} \\ \text{the num of points} \end{array}$$

Probability of blues.



This node is pure node

→ Measure of impurity is entropy

Entropy will be 0.

$$H_M = - \sum_{i=1}^k p_m^i \log_2 p_m^i$$

→ If you add information, entropy ↓

Best split?

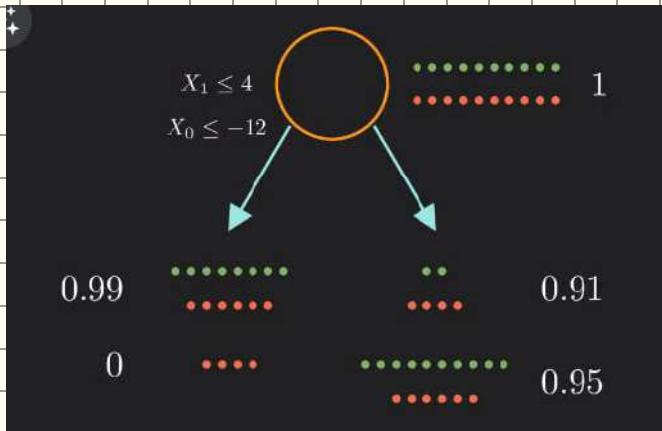
Maximize Information gain

$$\text{Entropy} = \sum -p_i \cdot \log(p_i)$$

→ possibility of  
ith class.

↑ Entropy ↑ uncertainty

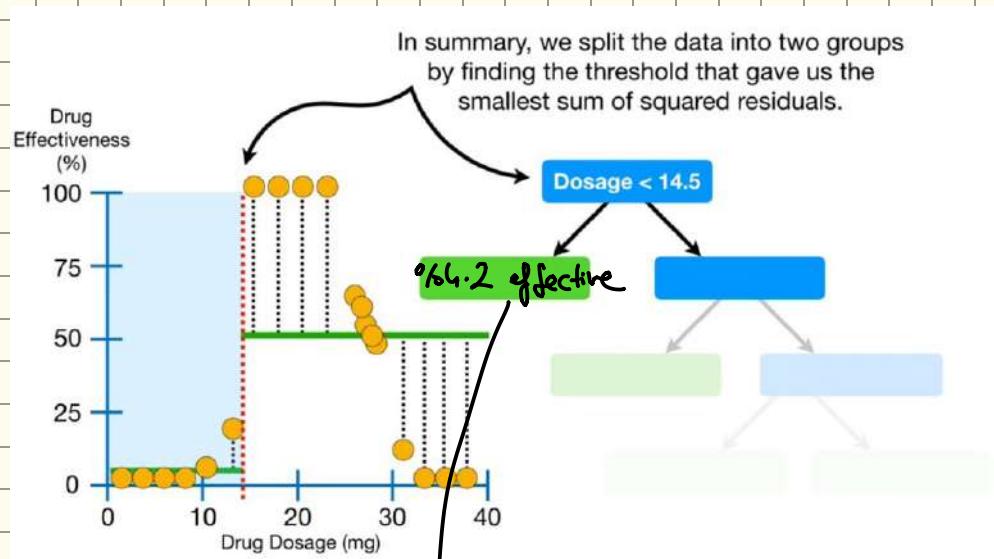
$$\text{Information Gain} = E(\text{parent}) - \sum_{\text{child}_i} \frac{E(\text{child}_i)}{\text{child num}} \downarrow \text{parent num}$$



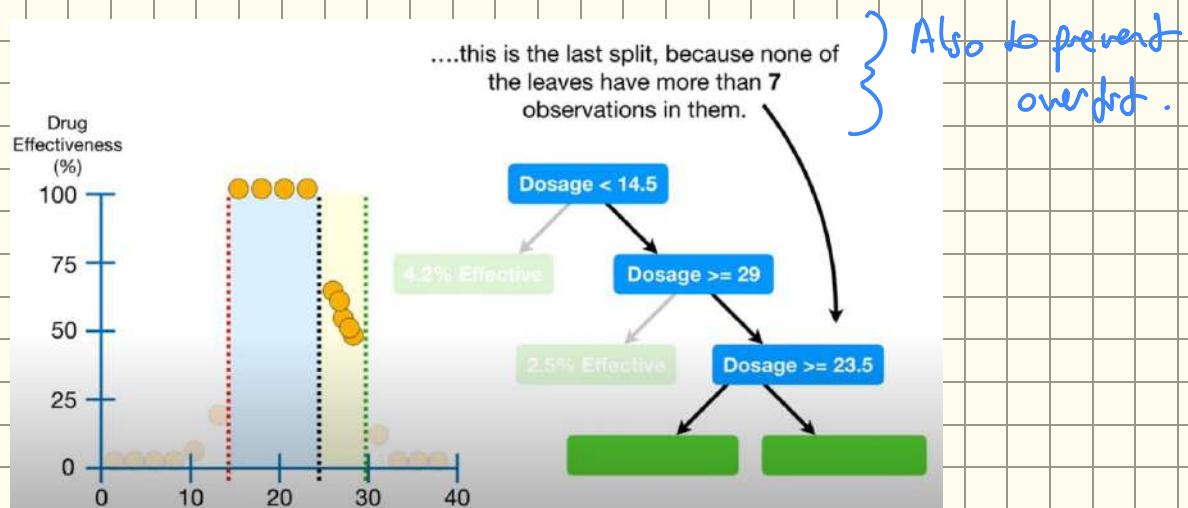
BIAS - VARIANCE TRADEOFF.

# Regression Trees.

Building trees:



We set `min_mim_split` as 20 to prevent overfit. So this node will be a leaf. (Left node's average is %64.2.)



→ Bir den fazla feature ile yapmak için her node için 16,5 gibi olen değer oluyor. (sum of squared residuals) - Min olan root obrak seçilir.

$\downarrow$   
min error.

"Better candidate for the root"

→ Overfitting entsteht im Pruning:

<u>Prepruning</u>	<u>Postpruning</u>	
↓		
- used before construction of tree	↳ used after construction	deeper vs bkt
- can be done using Hyperparameter tuning	↳ used when tree has very large depth	
(en: pruned in en: pruned parameters)	→ more accurate	stop test?
seiner (max_depth, min_sample_split) gibt		

→ faster

### Advantages of DT:

- Non parametric
- Classification and regression
- Classification fast once rules are developed.

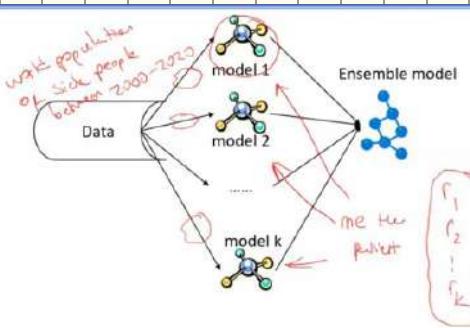
### Disadvantages of DT:

- Tend to overfit even pruning.
- Instability

# ENSEMBLE LEARNING

Combine multiple models together

- Committees – Bagging
  - Regression: Take an average of the predictions made by each model
  - Classification: Make classification by voting over a collection of classifiers
- Boosting – Adaboost
  - Train multiple models in sequence
- Decision trees
  - Different models are responsible for making predictions in different regions of input space



\* Majority vote.

\* Combine multiple models.

\* Combining diverse, independent opinions.

- How can I get different models?

→ different data + same models } Bagging (Bootstrapping)

→ same data + diff. models → not good idea } Boosting

→ change feature space for each model. (sampling features)

→ Randomness

↳ Random forest.

\* Choose random features and calculate the best split  
(gini)

If I have enough training data:

- How can I get different (independent) models?

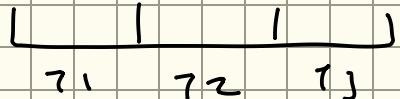
1-) Different datasets. (train data)

$$T_1 \rightarrow (\text{BML}) \rightarrow M_1$$

$$T_2 \rightarrow \dots \rightarrow M_2$$

$$T_3 \rightarrow \dots \rightarrow M_3$$

$$T_1 \cap T_2 = \emptyset \quad T_1 \cap T_3 = \emptyset \quad T_2 \cap T_3 = \emptyset$$



Performance of each models?

$$T_1 \cup T_2 \cup T_3 \rightarrow (\text{BML}) \rightarrow M$$

$P_1 > P_2 > P_3$  → because they have less data  
↓  
combine data.

$P_{1,2,3} > P_1 > P_2 > P_3$   
↓  
Majority voting

Trade-off

Majority vote gives better  $P_1$  % than  $P_2, P_3$

$P_2, P_3$  are also ok but  $P_1 > P_2, P_3$   
ok.

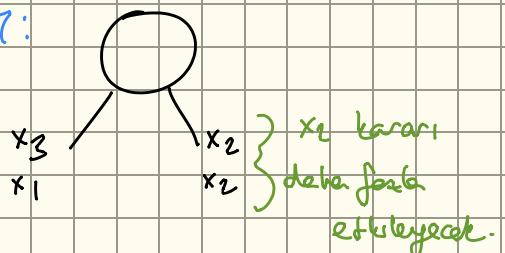
Leads to Subsampling

$x_1, x_2, x_3, x_4$   $M_1$

$x_3, x_1, x_2, x_2$   $M_2$

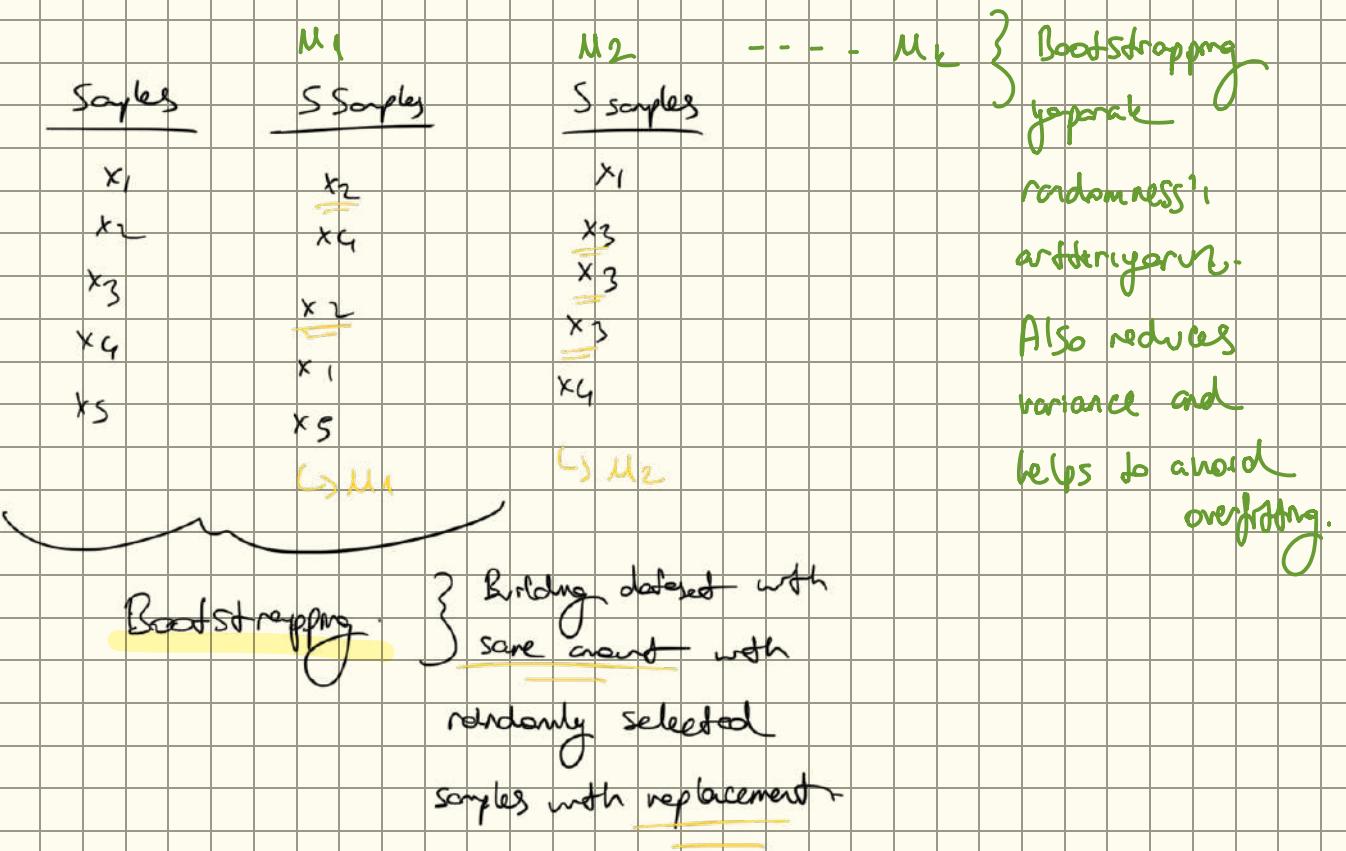
**SVM:** Data points are given  
soft margin belirleme mit iyi bir öreni  
olustur.

**DT:**



**1NN:** etkilevir.

**3NN:**  $x_2$ 'ye 2 tane hewarda mit data da  
test data  $x_2$ 'ye yeter gelirse etkilevir.



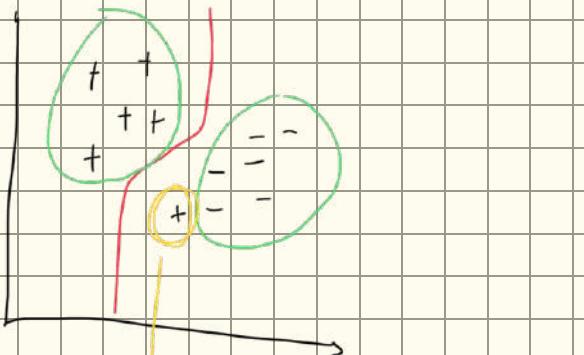
Bootstrapping  $\rightarrow$  Sampling (random) with replacement

## Bagging (Bootstrapping + Aggregation)

Regression: Take ave. of the predictions

Classification: Majority voting

## 2-) Use different models $\rightarrow$ Boosting



replicate this  
use for other bootstrapping.

It's gonna concentrate on his data  $\rightarrow$  boosting  
(problem areas)

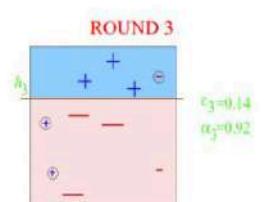
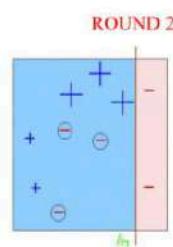
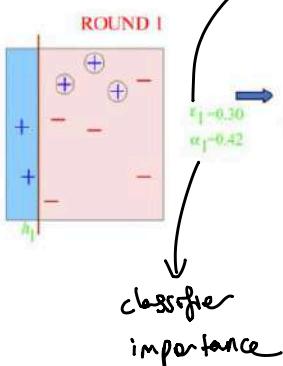
- Boost weak learners to strong learners.

Example: AdaBoost

### • AdaBoost

- Initially, set uniform weights on all the records
- At each round
  - Create a bootstrap sample based on the weights
  - Train a classifier on the sample and apply it on the original training set
  - Records that are wrongly classified will have their weights increased
  - Records that are classified correctly will have their weights decreased
  - If the error rate is higher than 50%, start over
- Final prediction is weighted average of all the classifiers with weight representing the training accuracy

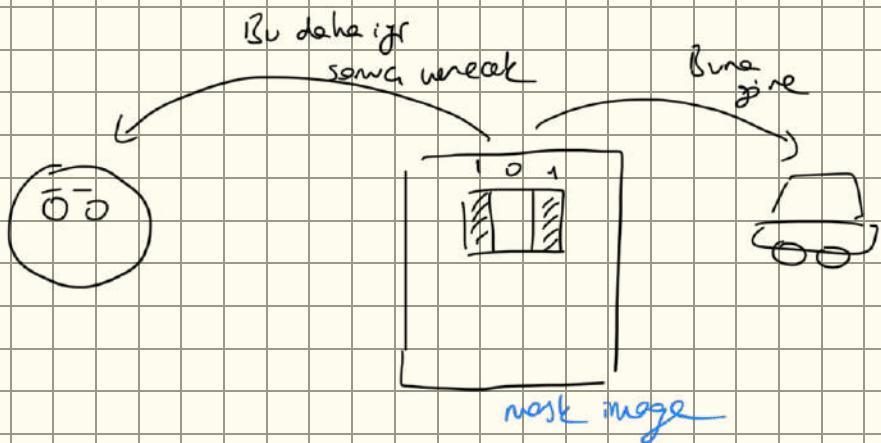
classifier error ( $\frac{3}{10}$ )



(training accuracy)

$$H_{\text{final}} = \text{sign} \left( 0.42 \cdot h_1 + 0.65 \cdot h_2 + 0.92 \cdot h_3 \right)$$

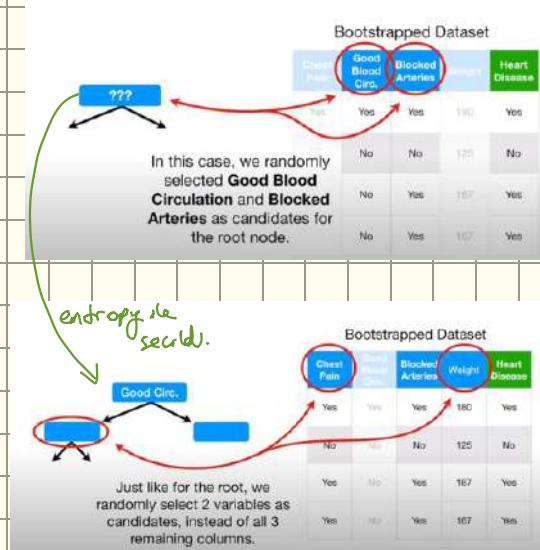
Boosting example : face detection.



# RANDOM FORESTS

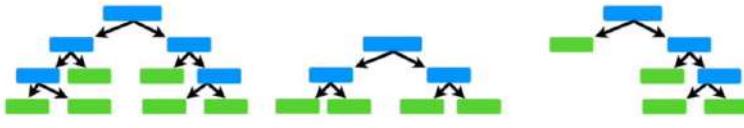
We built a tree:

- Bootstrapped dataset
- Only considering a random subset of features each step.

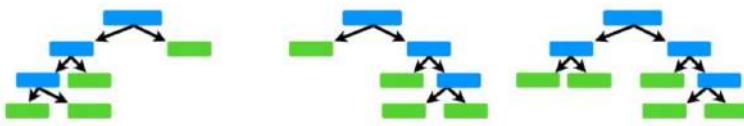


Repeat: Make a new bootstrapped dataset  
and build a tree with random features each step.

Using a bootstrapped sample and considering only a subset of the variables at each step results in a wide variety of trees.



The variety is what makes random forests more effective than individual decision trees.



→ then majority vote

- different dataset
- different models
- subset of features
- randomness

hyperparameters

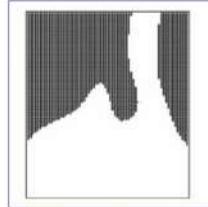
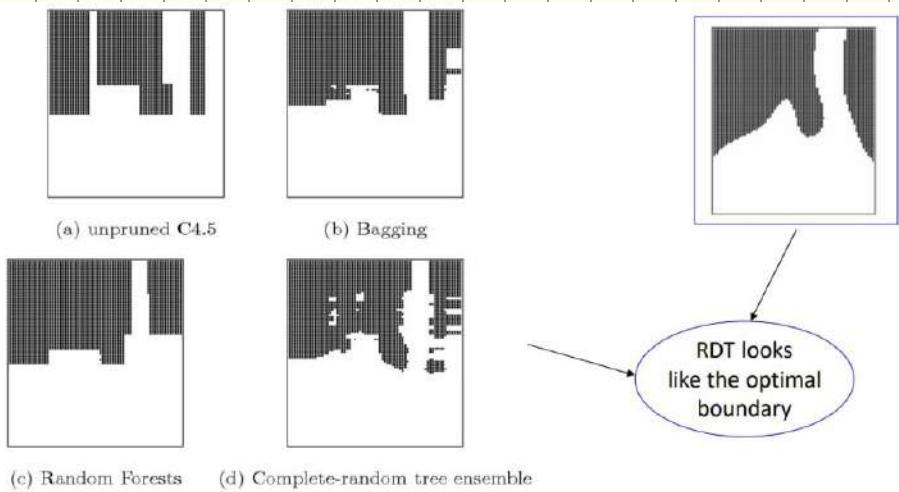
Algorithm

- Choose  $T$  — number of trees to grow
- Choose  $m < M$  ( $M$  is the number of total features) — number of features used to calculate the best split at each node (typically 20%)

- For each tree
  - Choose a training set by choosing  $N$  times ( $N$  is the number of training examples) with replacement from the training set
  - For each node, randomly choose  $m$  features and calculate the best split
  - Fully grown and not pruned → kind of want to overfit, achieve the independence.
  - Use majority voting among all the trees (or average for regression)

} Bootstrap

- Improve accuracy: Diversity, reduce variances
- Improve efficiency: Searching among subsets of features is much faster than searching among the complete set.



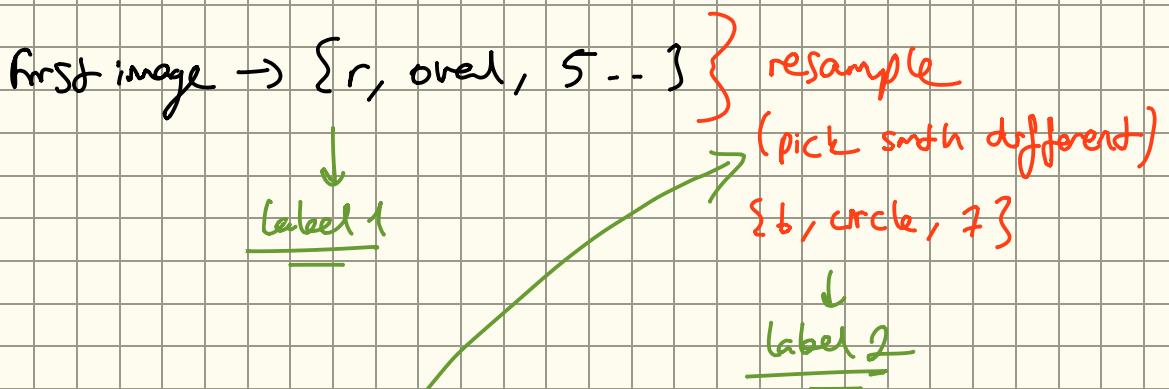
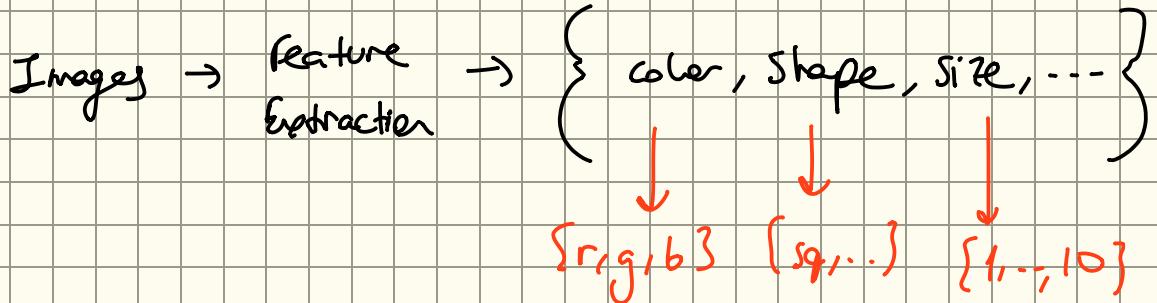
RDT looks  
like the optimal  
boundary

↓  
Decisions  
 $x_1 > 0$

↓  
you can make "multivariate"  
decisions!  
 $ax_1 + bx_2 + c > 0$

## RF with Clustering: (unsupervised learning)

apples  
pears



- RF to solve unsupervised learning problems, in particular, clustering problems and missing value imputation in the general sense
- RF generates a synthetic target variable in order to proceed with a regular run:
  - Give class label 1 to the original data
  - Create a copy of the data such that each variable is sampled independently from the values available in the original dataset
  - Give class label 2 to the copy of the data
  - RF on two class → dissimilarity measure for clustering...
  - Note that the second copy has marginal distributions identical to the first copy, whereas the possible dependency among predictors is completely destroyed
  - A necessary drawback is that the resulting dataset is twice as large as the original

## SUPERVISED ML

### Feature Reduction

Mapping high dimensional data to lower-dimensional space without losing essential information.

Some problem

100 features.

$$y = w^T x + b$$

100 parameters + 1

5 features.

$$y = w^T x + b$$

5 parameters + 1

simple

model (less noise)

I can have

- computational ↓

- more jobs to do

- less things to find

✓ optimization  
isn't ifr.

- simplicity

- faster

- noise can handle  
better.

- can improve accuracy  
(prevent overfit)

$f_1 \ f_2 \ - \ - \ - \ f_{100}$  sensors

↙

? is it needed

standard deviation =  $\sigma$  } it's same data  
or low

not useful  
because no  
variation

or  $F_2, F_3$  is highly  
correlated.

I don't need both of them.

→ Given a set of points, compute linear transformation

$$G \in \mathbb{R}^{d \times k}: x \in \mathbb{R}^d \rightarrow y = G^T x \in \mathbb{R}^k \quad (k \ll d)$$

↓  
originally were  
 $d$  dimensional  
space

↓  
reduced to  
 $k$  dimensional  
space.

## Why Reduce Dimensionality?

1. Reduces time complexity: Less computation
2. Reduces space complexity: Less parameters (compression)
3. Saves the cost of observing the feature
4. Simpler models are more robust on small datasets
5. More interpretable; simpler explanation
6. Noise removal
7. Data visualization (structure, groups, outliers, etc) if plotted in 2 or 3 dimensions

"Curse of dimensionality"

$x \xrightarrow{f} y \rightarrow \text{classification}$

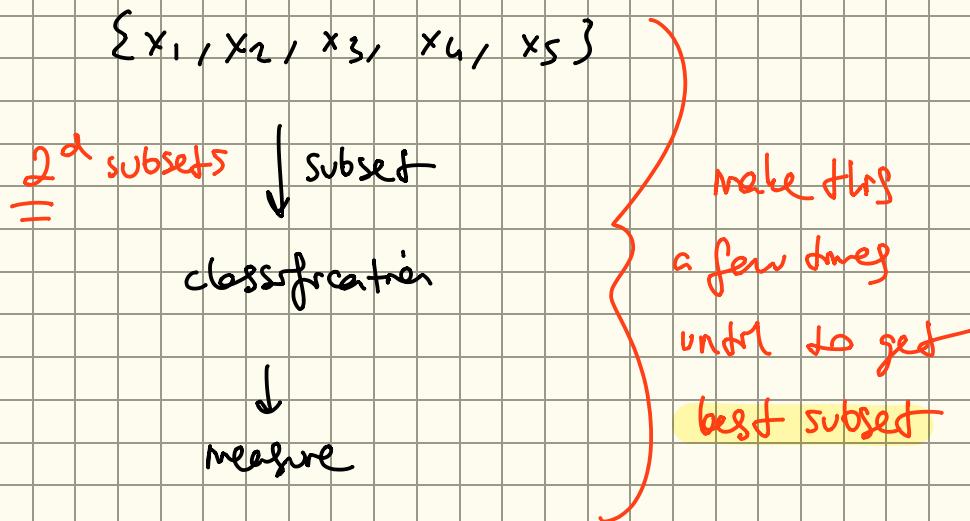


When your class. problem is in linear space but  
reduced one gives you non-linear space, it  
may not be a good idea to reduce it.

- **Feature selection:** Choosing  $k < d$  important features, ignoring the remaining  $d - k$ 
  - Subset selection algorithms
- **Feature extraction:** Project the original  $x_i, i = 1, \dots, d$  dimensions to new  $k < d$  dimensions,  $z_j, j = 1, \dots, k$ 
  - Principal components analysis (PCA), (unsupervised)
  - Linear discriminant analysis (LDA), (supervised)
  - Factor analysis (FA) ...

## Subset Selection:

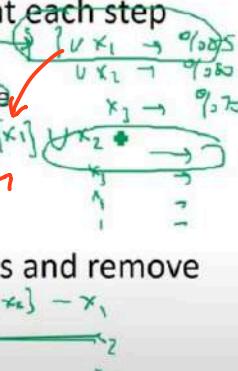
$$d=5$$



## Complexity:

$2^d$  work for a. Let's make with greedy approach!

- **Forward search:** Add the best feature at each step
  - Set of features  $F$  initially  $\emptyset$
  - At each iteration, find the best new feature  $j = \arg \min_i E(F \cup x_i)$
  - Add  $x_j$  to  $F$  if  $E(F \cup x_j) < E(F)$
  - Greedy  $O(d^2)$  algorithm → nested loop
- **Backward search:** Start with all features and remove one at a time, if possible  $\{x_1, x_2, x_3, x_4, x_5\} - x_1$
- **Floating search:** Add  $k$ , remove  $l$



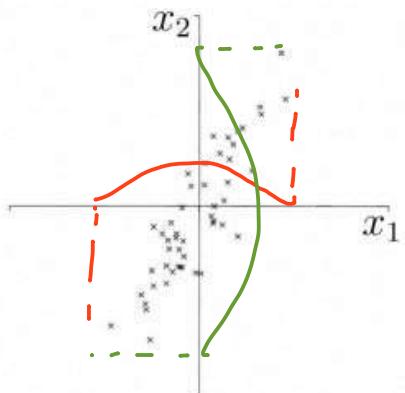
$x_1, x_2 \} x_3 \text{ ile bir boglarsa, greedy algo. buna } \cancel{x_1, x_2, x_3, x_4, x_5} \text{ feature selection cok dercih edilir. } \} \text{ gorenmeyeblir. There are non-homogeneous methods that can find those type of linear dependences } \Rightarrow$

# PCA (Principal Component Analysis):

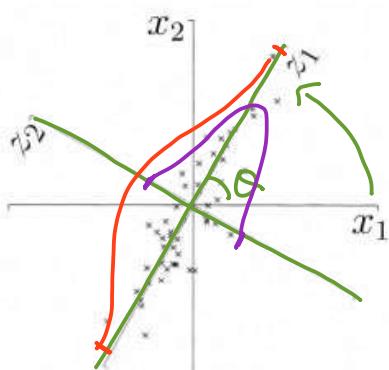


→ Reduce the dimensionality of a dataset by finding a new set of variables, smaller than the original.

↳ Still retains most of the information.



iki side genīe bir aclar kaptiyor. koyalice reduction yapanayiz.



new coordinate system

$$z_1 = \cos \theta \cdot x_1 - \sin \theta \cdot x_2 = w_1^T x$$

$$z_2 = \sin \theta \cdot x_1 + \cos \theta \cdot x_2 = w_2^T x$$

$$\delta(z_1) > \delta(x_1) \text{ or } \delta(x_2)$$

I have a new feature has larger variance than old one.

$$\delta(z_2) \ll \delta(x_1) \text{ or } \delta(x_2)$$

} variation is smaller.

$$\{x_1, x_2, \dots, x_n\} \in \mathbb{R}^d$$

$$z_1 = a_1^T \cdot x_i \quad (z_1 = w_1^T x)$$

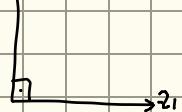
Linear Combination.

where  $\text{var}[z_1]$  is maximum.

To find next vector  $a_2$  maximizing  $\text{var}[z_2]$ , subject to  $\text{cov}[z_2, z_1] = 0$

} uncorrelated

$a_2$  also eigenvector of  $S$  whose eigenvalue  $\lambda = \lambda_2$  second largest.



In general;

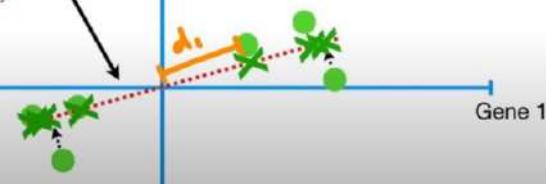
$$\text{var}[z_k] = a_k^T \cdot S \cdot a_k = \lambda_k$$

$a_1$  is an eigenvector of  $S$  corresponding to the largest eigenvalue  $\lambda = \lambda_1$ .

$$d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2 + d_6^2 = \text{sum of squared distances} = \text{SS}(distances)$$

Ultimately, we end up with this line. It has the largest SS(distances).

maximum var[z1]



To make PC1  
Mix 4 parts Gene 1 with 1 part Gene 2

When you do PCA with SVD, the recipe for PC1 is scaled so that this length = 1.

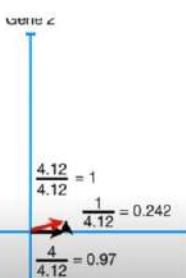


Terminology Alert!!!!  
Mathematicians call this cocktail recipe a "linear combination" of Genes 1 and 2.

The new values change our recipe...

To make PC1  
Mix 0.97 parts Gene 1 with 0.242 parts Gene 2

...but the ratio is the same: we still use 4 times as much Gene 1 as Gene 2.



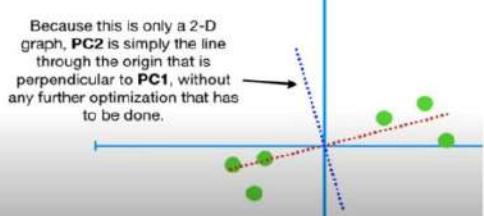
Gene 2

To make PC1  
Mix 0.97 parts Gene 1 with 0.242 parts Gene 2

...and the proportions of each gene are called "Loading Scores".

Gene 2

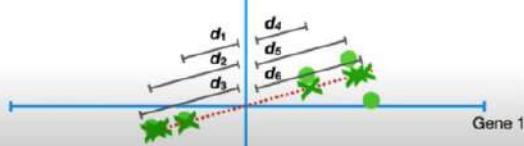
Terminology Alert!!! This 1 unit long vector, consisting of 0.97 parts Gene 1 and 0.242 parts Gene 2, is called the "Singular Vector" or the "Eigenvector" for PC1.



$$\frac{\text{SS}(distances for PC1)}{n-1} = \text{Eigenvalue for PC1}$$

(variation)

Also, while I'm at it, PCA calls the average of the SS(distances) for the best fit line the Eigenvalue for PC1...



PC2

0 0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

0 0

0

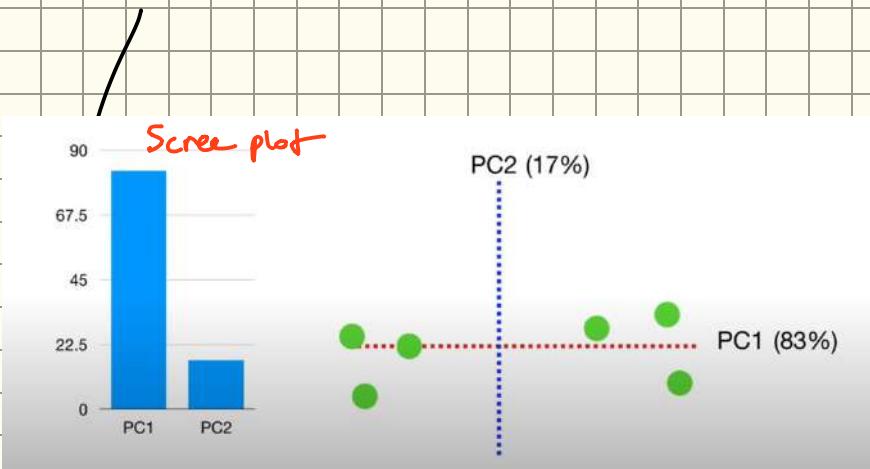
0 0

0

0 0

0

\* Eigenvalues are just measures of "variation"



Rotate  
Find variations

Covariance matrix

$$\text{var}[z_k] = \alpha_k^T S \alpha_k = \lambda_k$$

$$S = \frac{1}{n} \cdot X \cdot X^T \in \mathbb{R}^{d \times d}$$

$$X = [x_1, x_2, \dots, x_n] \in \mathbb{R}^{d \times n}$$

Obtain eigenvectors of  $S$  by computing the SVD of  $X$ :

MNIST digit recognition:  
 $10 \times 4k \rightarrow 10,000$   
 $30 \times 30 \rightarrow 900$  features

$$X = U \Sigma V^T$$

$$z_{k \times 1} = \begin{bmatrix} \alpha_1^T \\ \alpha_2^T \\ \vdots \\ \alpha_k^T \end{bmatrix} \in \mathbb{R}^{d \times 1} \quad (k \ll d)$$

All vectors

Assume  $\bar{x} = 0$   
Form the matrix:  $X = [x_1, x_2, \dots, x_n] \in \mathbb{R}^{d \times n}$

$$\text{then } S = \frac{1}{n} X X^T$$

Obtain eigenvectors of  $S$  by computing the SVD of  $X$ :

$$X = U \Sigma V^T$$

$$x_k' = \alpha_k \times d$$

orthonormal transformation

$$x_d' = T \times d$$

variations decreasing, when 0 ignore it.

How can I find k value?

When

$$\text{var}[z_{k+1}] = 0 \text{ stop.}$$

Ama her zero 0 olmazsa.

$$\vec{x} = [x_1, x_2, x_3]$$

$$x_3 = a x_1 + b x_2 \text{ oldugu durumda}$$

$\lambda_3 = 0$  olur, PCA turu hesaplayamaz.

Ama egere noise varsa?

$$x_3 = a x_1 + b x_2 + a \epsilon_1 + b \epsilon_2$$

Diger dogeler var.

PCA hesaplayamaz.

\* Step "elbow" of scree graph.

\* Proportion of Variance (PoV) explained

$$\frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_k + \dots + \lambda_d}$$

when  $\lambda_i$  are sorted in descending order

\* Typically, stop at PoV > 0.9

$$X \in \mathbb{R}^{d \times n} \rightarrow G^T \cdot X \in \mathbb{R}^{k \times n} \rightarrow \bar{X} = G(G^T \cdot X) \in \mathbb{R}^{d \times n}$$

↓  
Reconstruction

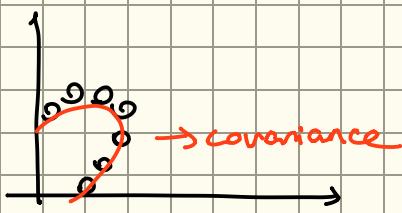
$$\vec{x} = [x_1, x_2, x_3]$$

$$x_3 = a x_1 + b x_2 \quad \text{odd gur durrada}$$

$$(x_1, x_2, x_3) \xrightarrow{\text{PCA}} (z_1, z_2)$$

Because there are linear relationship.

$$x_3 = x_1^2 + x_2 \rightarrow \text{not linear? } \quad \left. \begin{array}{l} \text{PCA} \\ \text{PCA} \end{array} \right\} (x_1, x_2, x_3) \rightarrow (x_1, x_1^2, x_2, x_2^2, x_3, x_3^2)$$



Now, I can reduce them.

$$\downarrow \text{PCA} \\ (x_1, x_1^2, x_2)$$

high dimension  
↓

apply non-linear trans. to high-dimensional space.

↓  
restricted set of trans. with good computational properties

↓  
"Kernel trick" ( $k_{\alpha} = \alpha \alpha$ )

PCA → No loss in lower eigenvalues

↓  
Pick a kernel function RBF (to high dimension)  
Does it now have lower eigenvalues?  
No

$$x_i = \{t, h, l, p, w\}$$

↓  
weekday, not gonna include it  
(categorical)

$$x'_i = \{t, h, l, p\}$$

(c) leave w aside and deal with real data.

$$[t, h, l, p, w] \rightarrow [f_1, f_2, w], y$$

push to  
me now  
*w*

$$\begin{bmatrix} f_1 \\ f_2 \end{bmatrix} = A \begin{bmatrix} t \\ h \\ l \\ p \end{bmatrix}$$

PCA does.

(c) I can't multiply with w  
Not include w in this process.

$$\text{New feature } i: \alpha_1 \cdot t + \alpha_2 \cdot h + \alpha_3 \cdot p + \alpha_4 \cdot l + \alpha_5 \cdot w$$

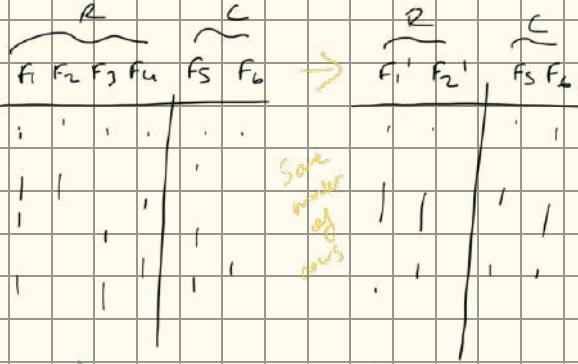
$$1, 2 + w \rightarrow \text{is this meaningful?}$$

Numerically it will give you smooth

But that doesn't mean anything  
Semantically



\* In PCA, we drop the categorical values, add them later.



2 features are enough.

g

SVM kNN DTV RFV

$$F_1 > 0, F_5 = C_1$$

2D is better

PCA says take this measurements

$$X^T \quad \text{1000x5 matrix} \quad X \cdot X^T \rightarrow \text{eigen values.}$$

$x_1 \dots x_5$

eigen vector

$e_1 \dots e_5$

PCA

$$X_i = \begin{bmatrix} e_1^T \\ e_2^T \\ e_3^T \\ e_4^T \\ e_5^T \end{bmatrix} \begin{bmatrix} t \\ h \\ l \\ p \\ w \end{bmatrix}$$

gentler  
say yes  
comparable

4 features  
(1000 data)  
 $x_5 = 0$   
we dropped it.

DT  $\rightarrow$   $k \approx (n-1)$  comparisons

iteration

$2^k$

number of iterations

for decision

1-) faster.

2-) When there is noise, PCA give clean data in next step.

3-) high dimensional overfit easily. prevent overfitting

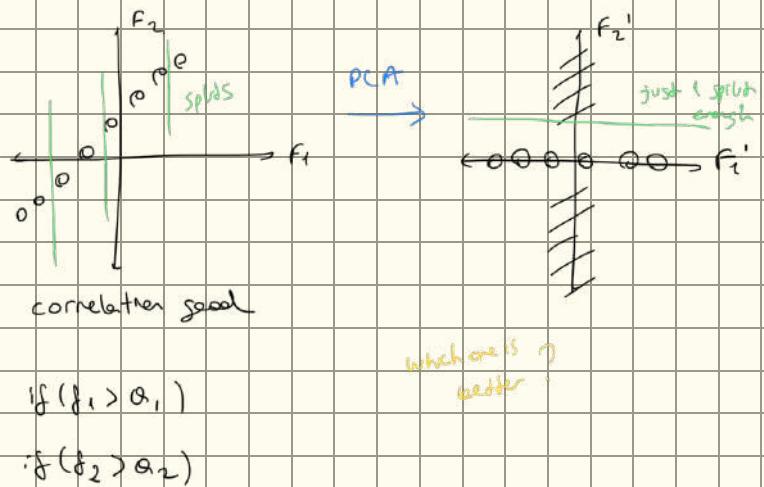
after PCA:

$$2^{(n-1)} + 2^k$$

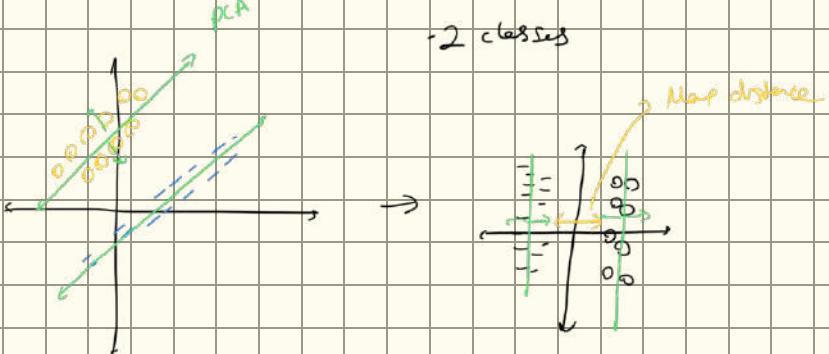
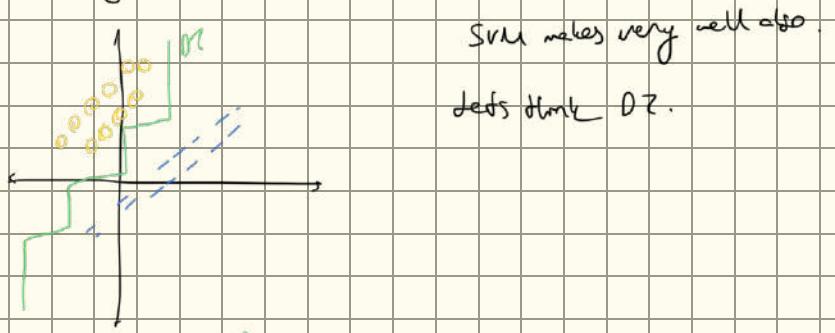
- computational advantage

- easier to model (6 features vs 5 features)

Can we do this with DT?

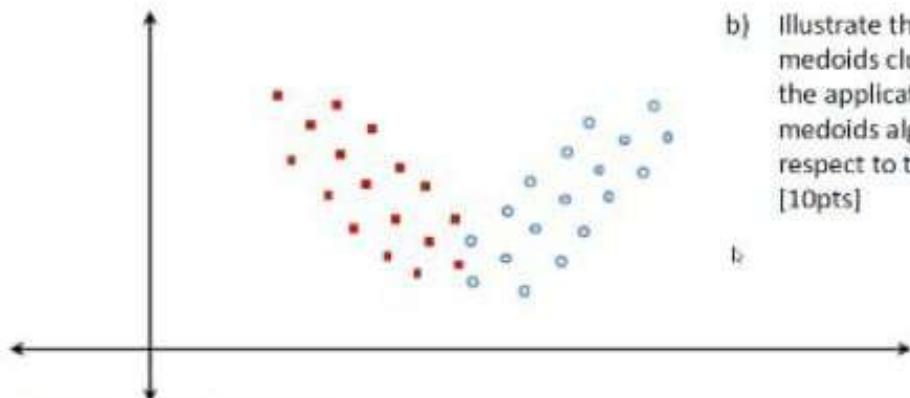


Binary classification problem:



If I rotate this, distance  
changes

- a) Given the following samples in a 2D Euclidean space, what does principal component (PCA) analysis reveal about the data? Show the principal components graphically. Explain your answer. [10pts]



- b) Illustrate the first two iterations of k-medoids clustering algorithm? Compare the application of k-means and k-medoids algorithms on this data with respect to their runtime complexities. [10pts]

CSF455 Final Exam June 13, 2023.

This is a handwritten question in this exam.

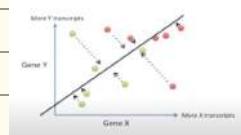
Find the  $a_1$  vector which represents  $\text{var}[z_1]$

# LDA (Linear Discriminant Analysis)

→ PCA → reduces dimensions with focuses on most variation ones.  
 - no class / labeled data (unsupervised)

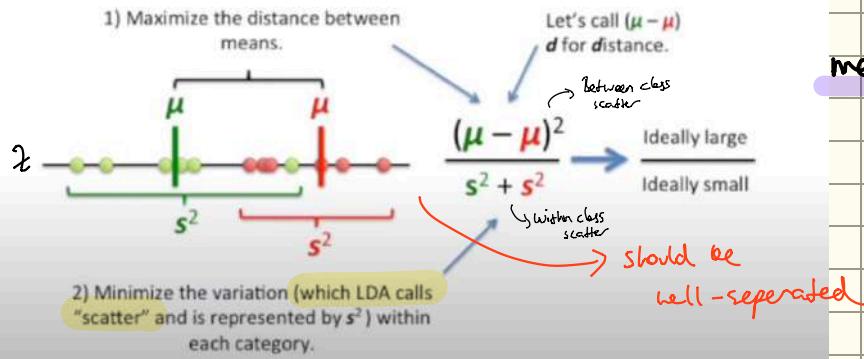
→ LDA → maximizing the separability between the known categories.  
 - labeled data (supervised)

LDA projects the data onto this new axis in a way to maximize the separation of 2 categories.



How LDA creates a new axis...

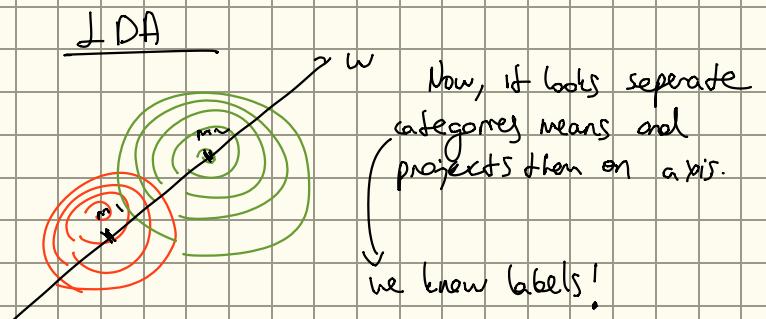
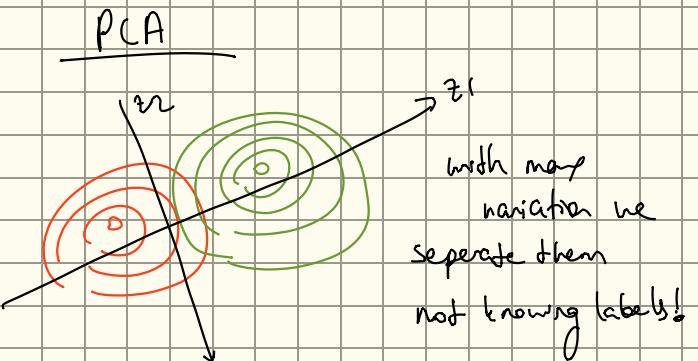
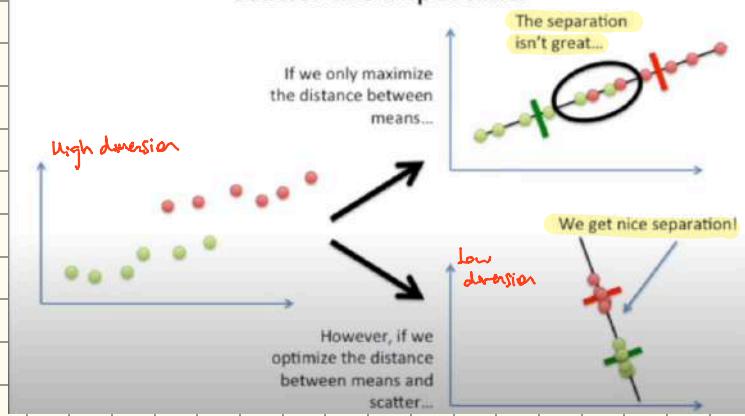
The new axis is created according to two criteria (considered simultaneously):



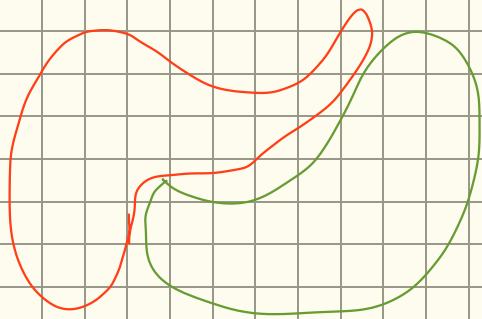
→ Maximize the distance between means, while minimizing the scatter.

$$z = (w_1, w_2) (x_1, x_2)$$

An example showing why both distance and scatter are important.



$\rightarrow$  LDA can't separate good in here.



More than 2 classes:

Fisher's linear discriminant: LDA with 2 classes.

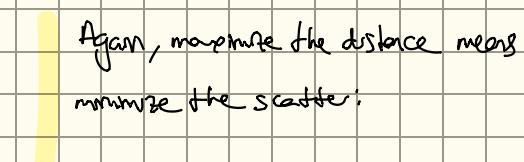
Find  $w$  that maps

$$J(w) = \frac{w^T \cdot S_{BS} \cdot w}{w^T \cdot S_w \cdot w} = \dots$$

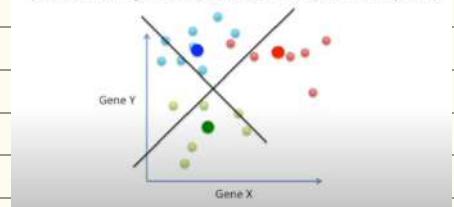
$\rightarrow$  max  
 $\hookrightarrow \min$

$$\boxed{w = c \cdot S_w^{-1} (m_1 - m_2)}$$

Again, maximize the distance means,  
minimize the scatter!



The second difference is LDA creates 2 axes to separate the data.  
This is because the 3 central points for each category define a plane.  
(Remember from high school: 2 points define a line, 3 points define a plane...)



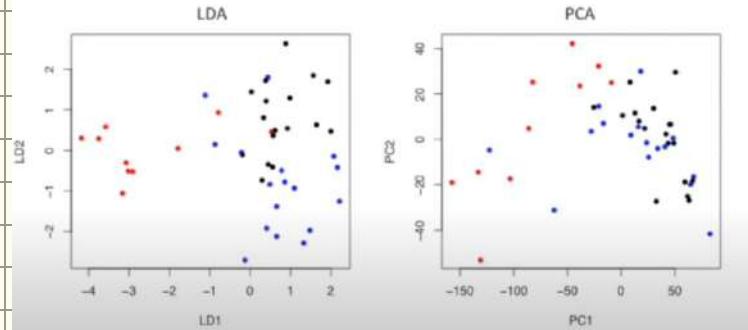
$\&$  Only applicable with linear problems.

$\rightarrow$  LDA doesn't have "soft margin" option like in SVM.

~~✗~~ PCA looks at the data with most variation

~~✗~~ LDA maximize the separation of known categories

Comparing LDA to PCA with 10,000 genes.



Eigenfaces: PCA application to specific case.

Calculation of eigenfaces:

- calculate average face  $\bar{v}$

Mean image:  $\bar{v}$

$$\bar{v} = \alpha_1 v_1 + \alpha_2 v_2 + \dots$$

Mean image  
Linear combination  
of images

$$v_i = \sum \alpha_i \phi_i$$

Some value.

recalculate  
images

→ If I have principal component of faces, then I can make your face with some parameters.

$\bar{v}$  is common to everything.

$[w_1, \dots, w_k]$  a signature of the faces.

$$\hat{x} = \bar{v} + w_1 u_1 + w_2 u_2 + \dots$$

right amount of each component.  
↓  
for face

$$= \bar{v} + \underbrace{w^T \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_k \end{bmatrix}}_{k < d \text{ feature extraction dimension}}$$

signature part  
other parts are common.

You can make a new face with this.

→ I'm just doing simple analysis of data.

$[1.2 \ -1.0 \ 1.5 \ 0.2] \rightarrow \text{this is my face}$

If I increase this a little

I can make new face similar to me.

PCA part

$x_1, \dots, x_N \rightarrow$  our data of faces.

$\downarrow$   
 $\Sigma \ \} \text{ covariance matrix}$

$u_1, \dots, u_k \} \text{ most } k \text{ eigen values relevant for us}$

We can represent new data as:

$$x_i \rightarrow \bar{v} + w_1 u_1 + w_2 u_2 + \dots$$

$\underbrace{[w_1, \dots, w_k]}_{100 \text{ dimension}}$

100 dimension (lower)

DT algo:

Partial tree

↳ soft

↳ take middle points

↳ choose best split by

Evaluate performance

Choose with best per-

Adaboost algo

Bootstrapped dataset

Train classifier on original dataset

Increase weight weak learners

Decrease " " good learners

Better or even better dev of weight-decay

Training accuracy is the model's error on training set.

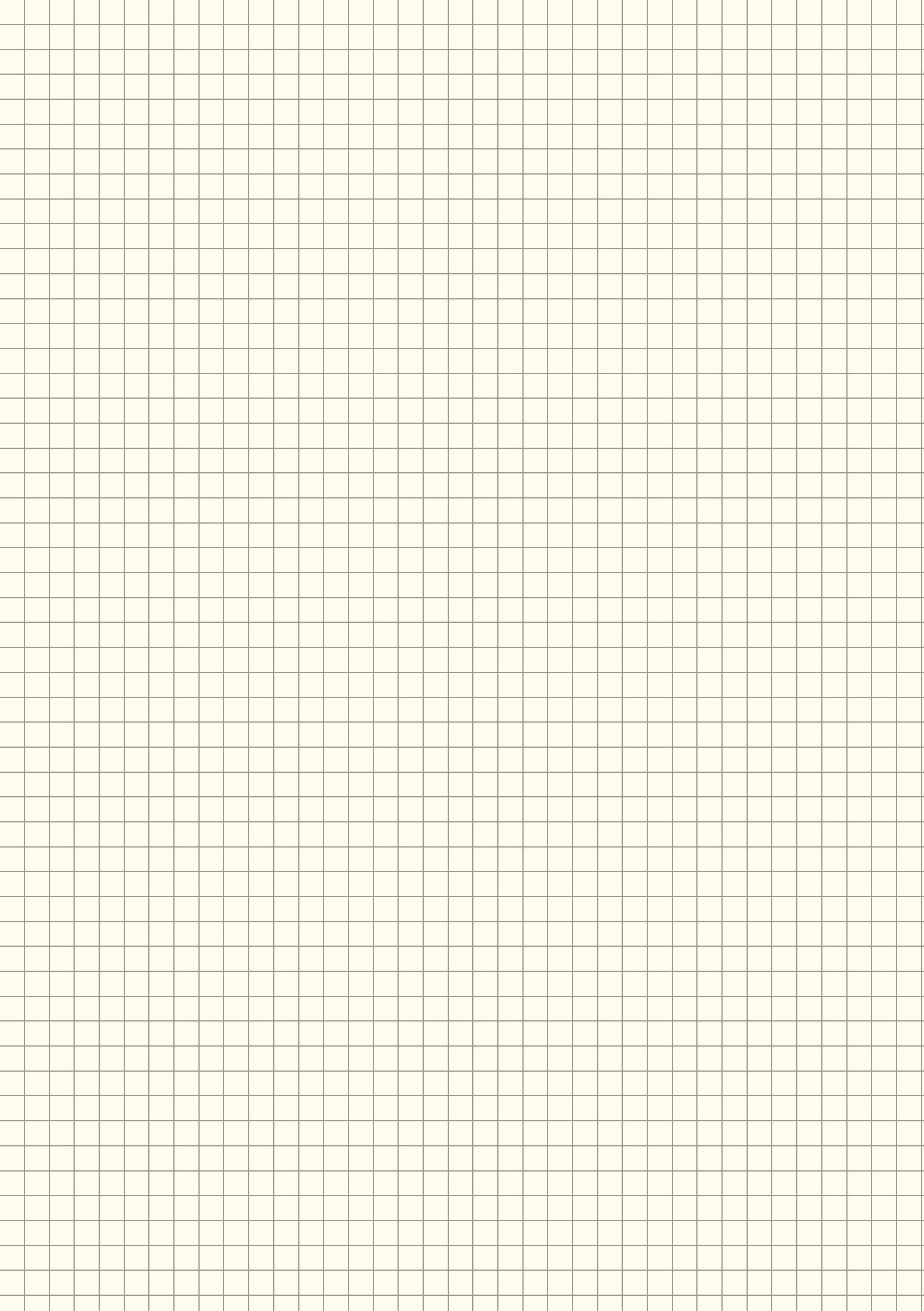
RF algo:

Partial trees!

Bootstrapped data

Random feature selection

Many弱者



# NEURAL NETWORKS

- The curved called "Activation functions".
  - Another "Hidden Layer"
  - Parameters we multiply } weights  
" " " add } biases.
  - Neural networks starts with identical Activation functions.  
But using weights and biases, it flips and stretches the Act. func. into new shapes.  
which are then added together to get a squiggle.
- sigmoid  
ReLU  
softmax

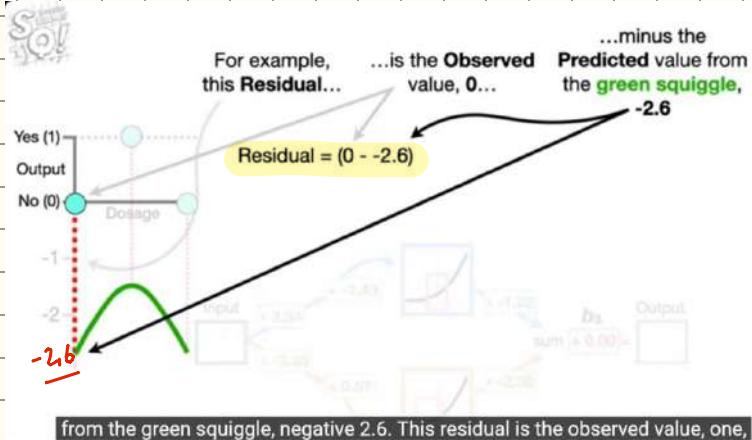
**Back propagation:** Predict weights and biases.

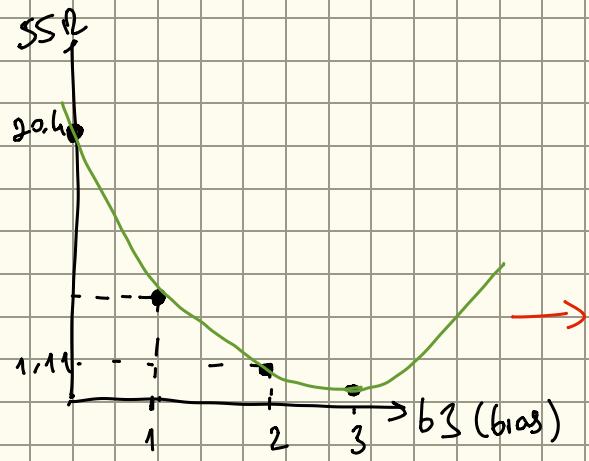
- To quantify how good is the squiggle } Sum of the Squared Residuals.

$$\text{Residual} = \text{Observed} - \text{Predicted}$$

So, when  $b_3 = 0$ , the SSR = 20.4...

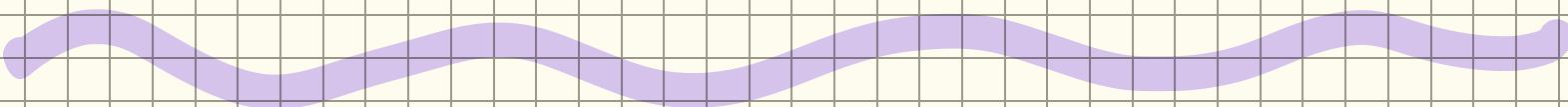
$$\begin{aligned} \text{SSR} &= (0 - -2.6)^2 \\ &\quad + (1 - -1.61)^2 \\ &\quad + (0 - -2.61)^2 = 20.4 \end{aligned}$$





Take the smallest and  
find it with gradient descent.

$$SSR = \sum (\text{observed}_i - \text{predicted}_i)^2$$



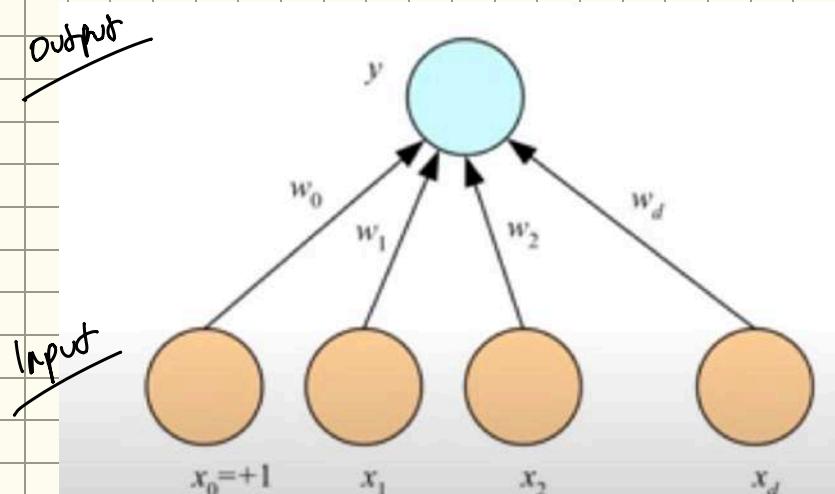
Large numbers of neurons:  $10^{10}$

Large connectivity:  $10^5$  as add.  $(\frac{1}{3}) + (\frac{1}{3}) + (\frac{1}{3})$

Parallel processing

Distributed computation / memory: They have to remember

Perceptron!



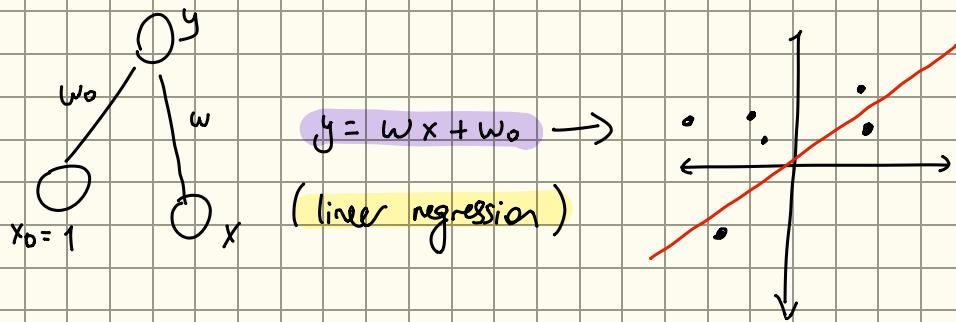
$$y = \mathbf{w}^T \mathbf{x} + b \quad (\text{lets assume this})$$

$$= \mathbf{w}^T \mathbf{x}$$

$$\mathbf{w} = [w_0, w_1, \dots, w_d]^T$$

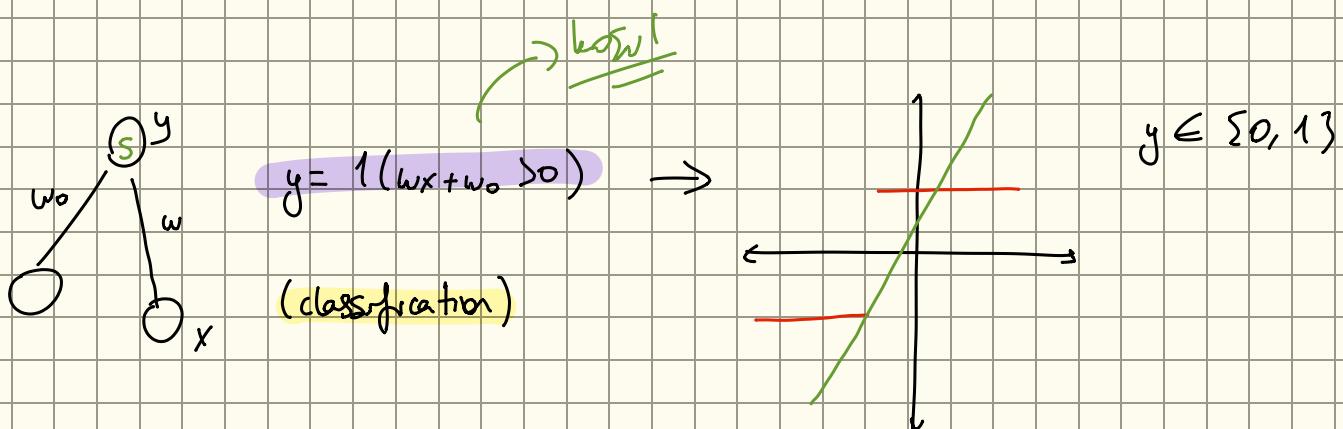
$$\mathbf{x} = [1, x_1, \dots, x_d]^T$$

Can I use it for regression?  
" " classification?

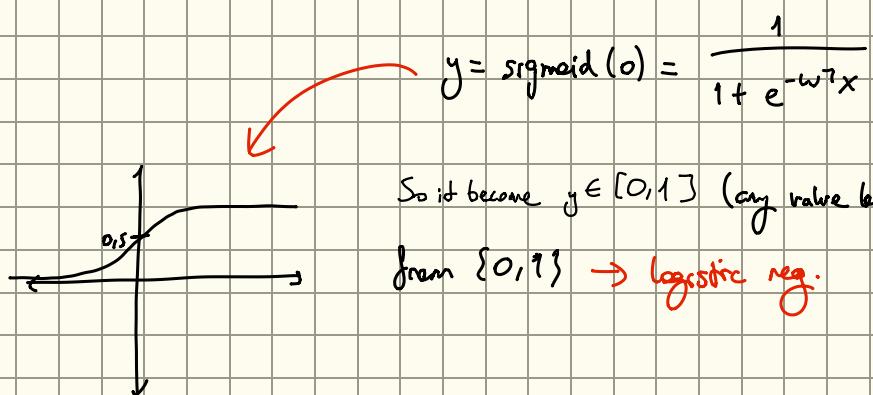


$$(x_1, x_2, x_3) \rightarrow w_1 x_1 + w_2 x_2 + w_3 x_3 + w_0$$

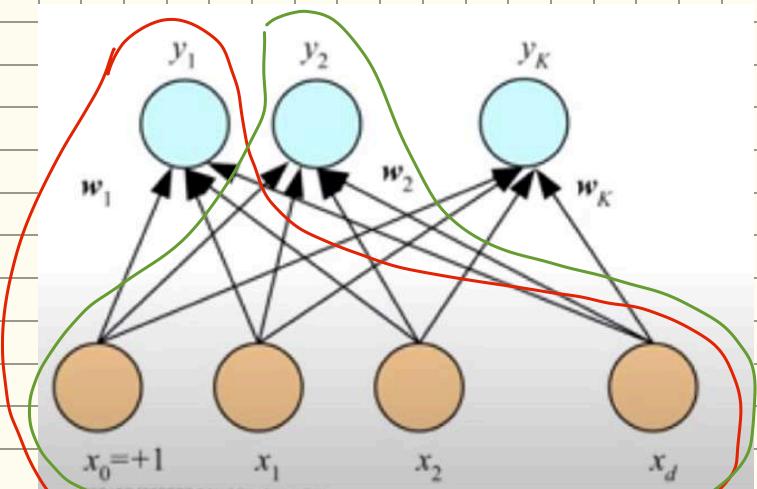
So, we can do linear regression ✓



↓ Can be continuity problem  
Use sigmoid



## Regression with K outputs:



$$\begin{aligned} y_1 &= w_1^T x + w_{10} \\ y_2 &= w_2^T x + w_{20} \\ &\vdots \\ y_k &= w_k^T x + w_{k0} \end{aligned} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} y = w x$$

$$\vec{y}_{(k)} = w_{k \times d} \cdot x^{(d)}$$

## Classification with K outputs:

More than 2 classes  $\rightarrow$  generate a categorical value

$\hookrightarrow$  Map these categorical values to a binary vector

Image  $\rightarrow$  Classifier  $\rightarrow \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$   
(Digit problem)

0 ? 8      } Same composition  
0 ? 1      } Not ordinal,  
1 ? 8      } Categorical numbers.

bcs they are just labels  
in this problem.

One hot encoding:

$$y' = [ \quad | \quad | \quad | \quad | \quad | \quad | \quad ]$$

binary  $\{0, 1\}$

So,

$$\text{"0"} = \{1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0\}$$

$$\text{"8"} = \{0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0\}$$

:

If you apply this:

$$1 \quad 0 \quad \dots \quad 0 \rightarrow \text{digit 1}$$

$$0 \quad 1 \quad \dots \quad 0 \rightarrow \text{digit 2}$$

:

You can do the same thing  
with regression:

$$o_i = w_i^T x$$

$$y_i = \frac{e^{o_i}}{\sum_k e^{o_k}}$$

If one of  $y_1, y_2, \dots, y_k$  are much

larger than others  $e^{o_i}$  will take them

↓

$$y = \frac{e^{o_i}}{e^{o_1} + e^{o_2} + e^{o_3} + \dots}$$

A way to use sigmoid function in

one-hot encoding.

↓

?? Softmax function (as opposed to sigmoid.)

For regression:  $y = w^T x$

"classification":  $y = \sigma(w^T x)$

## Training :

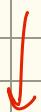
- Online learning (instances seen one by one) vs Batch learning (whole sample)
  - ↳ red iron vs red apple

$$\arg \min_w \sum_{i=1}^n |y_i - w^T x_i|^2 \rightarrow \text{Batch}$$

↓ simpler.

$$\arg \min_w |y_i - w^T x_i|^2 \rightarrow \text{online}$$

↳ Current model explains the data perfectly.



Calculate the difference  
btw. apple and iron.

Next time you see apple or iron, try to minimize the difference?

## Regression

$$E(w|x^t, r^t) = \frac{1}{2} (r^t - y^t)^2 = \frac{1}{2} (r^t - (w^T x))^2$$

$$\Delta w_{ij}^t = M (r_i^t - y_i^t) \cdot x_j^t \quad \text{per observation}$$

↳ hyperparameter

Update = Learning factor \* (Desired output - Actual output) \* Input.

↳ learning factor      ↳ error

Gradient Descent Update.



How big our update should be to correct that mistake.

- If you overreact : (making learning factor too large) not desirable

↳ prevent you from learning.

$$\Delta w_{ij}^t = M(r_i^t - y_i^t) \cdot x_j^t$$

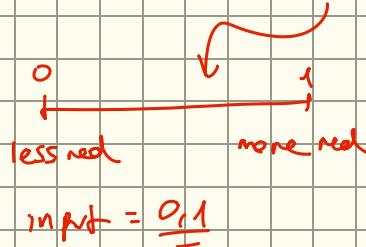
→ Delta to your weight is proportional to your mistake.

(error you make with your current  $w$ )

→ If mistake is 0, you learnt it.

→ If .. .. big, you should update your mistake.

↳ If error (mistake) is big but input is small, you should make small update.



Because it probably doesn't effect that much.

↳ If error and input is big, I should make big update.

### Algorithm

Step 0: Start with a random  $w$ . (small)

Step 1: Take a sample  $(x_i, y_i)$  error random

$$-y_i = w^T x$$

$$-\text{loss} = \text{error} = r_i - y_i$$

Update  $w$  with  $w + \mu(r_i - y_i)x_i$

Step 2: Repeat Step 1 as long as necessary.

If converges, you can stop ( $w$  doesn't change anymore)

randomness

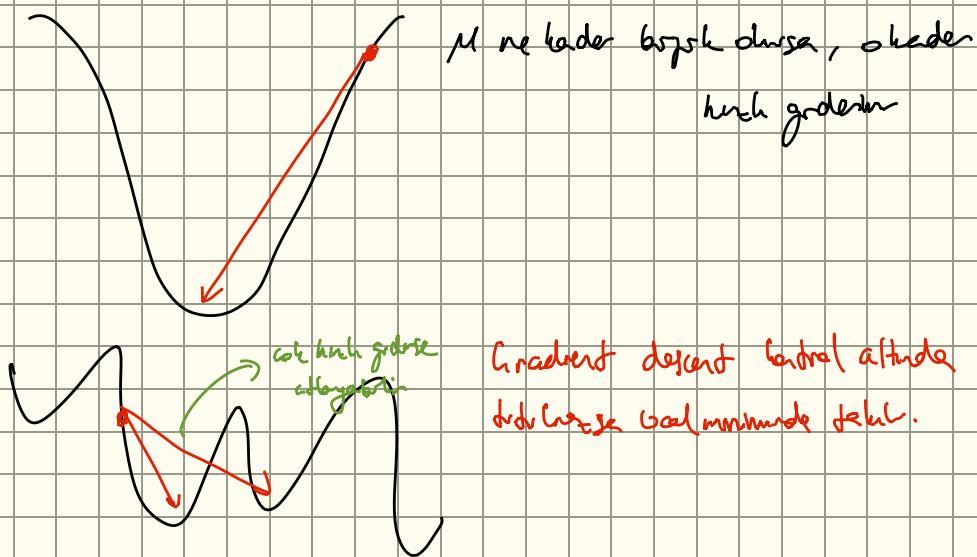
Stochastic

Gradient-Descent

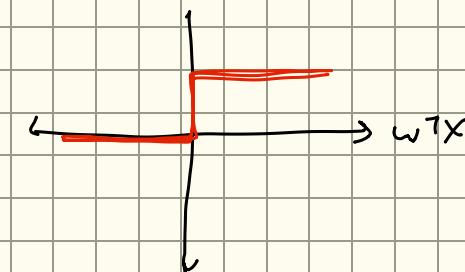
Alg.

or whenever you get tired

↳ Regression is easy bcs we are just trying to minimize that value.  
(gradient-descent minimization)



## Classification

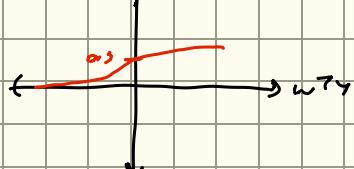


→ We try to make sure that it is mapped to a logical value (with sigmoid)

↳ Logistic regression

$$\begin{cases} 0 & \text{if } (w^T x) < 0 \\ 1 & \text{otherwise} \end{cases}$$

Then we applied some approximation to make continuity.



$y^t = \text{sigmoid}(w^T x) \rightarrow$  Single sigmoid output.

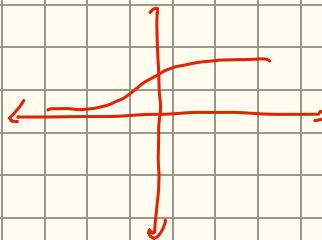
$$E'(w | x^t, r^t) = -r^t \cdot \log y^t - (1-r^t) \cdot \log(1-y^t)$$

loss func. I want to minimize

$$\Delta w_j^t = M(r^t - y^t)x_j^t$$

↳ same formula with regression

logistic selector: 0,1 we 0,2 argandeler hiz, 0,2 we 0,3 'den data yarang.



Bu yuzden sadece  $r-y$  bolunak yeter

$(1-r, 1-y)$ 'de batıylarız.

↓  
cross entropy.

↳ 2 softmax outputs

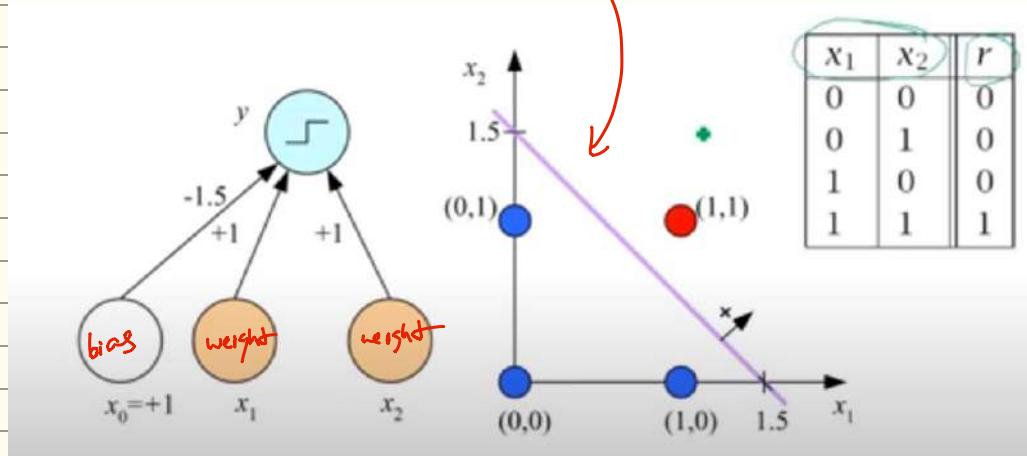
$$y^t = -\sum r_i^t \log y_i^t$$

$$\Delta w_{ij}^t = M(r_i^t - y_i^t)x_j^t$$

## Learning Boolean AND

→ Let's try Basic linear classifications

$$x_1 + x_2 - 1.5 = 0 \quad \left. \begin{array}{l} \text{line} \\ \text{y} = \text{sign}(w^T x) \end{array} \right\}$$

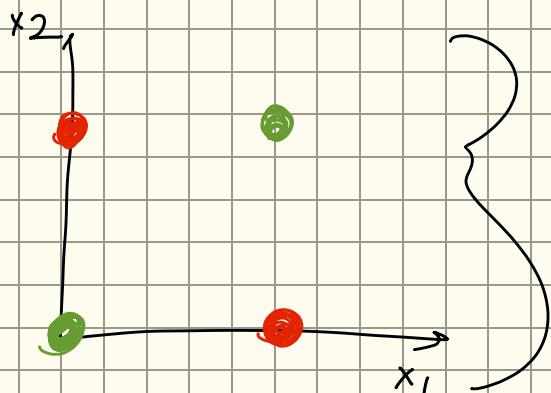


- 2 input neural network

XOR

$$w_1 x_1 + w_2 x_2 + w_0 \geq 0$$

How to draw this line? Intuitively



Positive negative  
sign. teraz direkt  
berlin cübeln.

Not a linear problem.

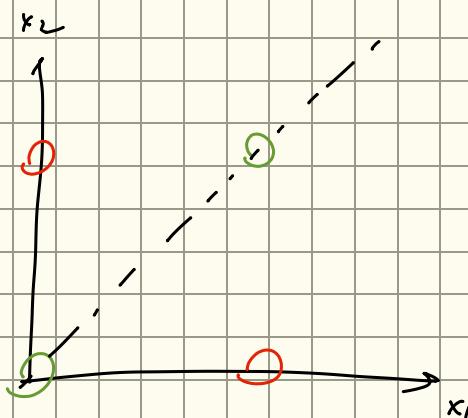
} A perceptron can't learn this  
with any linear classifier

1. ✓

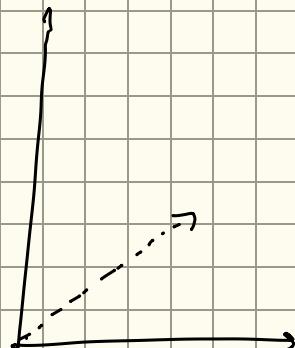
Upgrade to  
higher dimension

2. Non-linear method

1-)



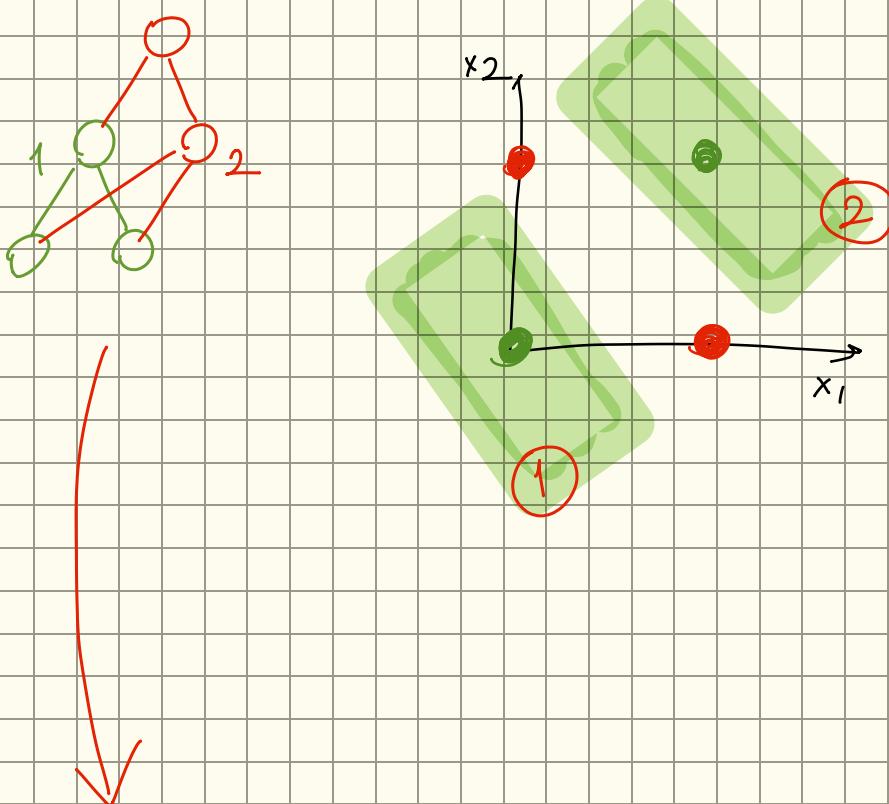
2D space



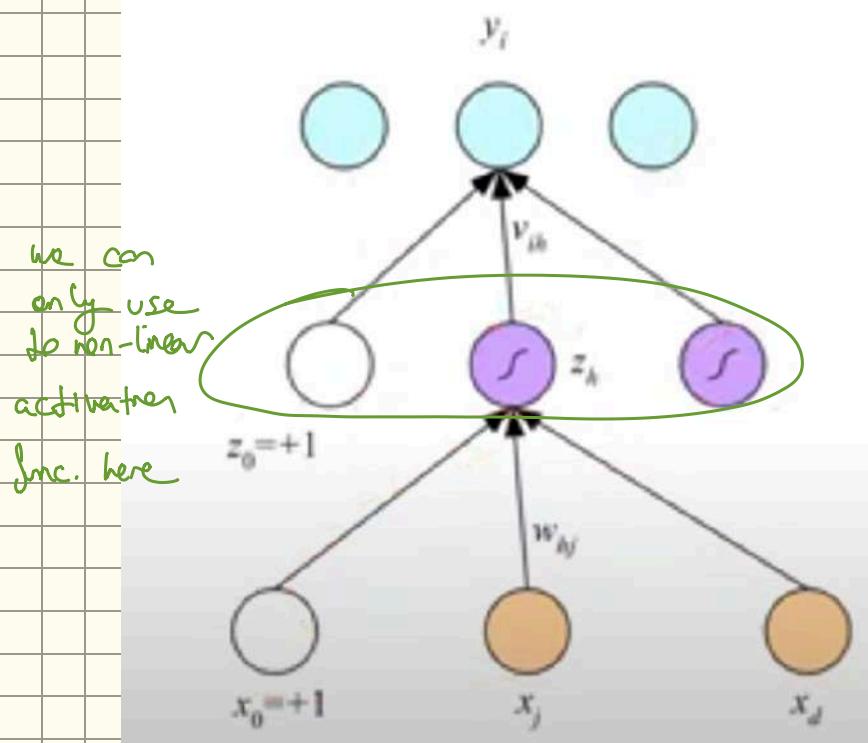
$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \Phi(x_1, x_2)$$

Yuksek boyutlu dairenin  
etrafında.

2-) Perceptron itself is a single linear classifier. So, let's combine 2 linear classifiers. (perceptron)

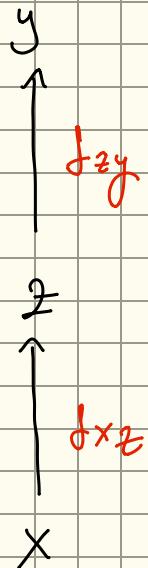


# Multi-layer Perceptron

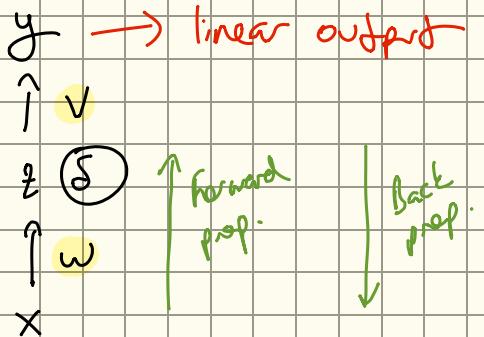


Multiple classifier combined

to get a full classifier  
and you are doing this at  
the same time. If it's  
converges, you're done.



How can I minimize the loss in this?



$\rightarrow$  Take  $x$  and calculate  $z$  and  $y$  with  
your current w, v values.  
(Initialized random)

Activation functions: It maps the resulting values in between 0 to 1 or -1 etc.

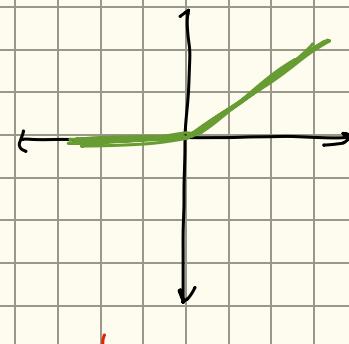
- Sigmoid : Differentiable

$$\frac{1}{1 + e^{-x}}$$

- Tanh : Similar to sigmoid

$$- \text{ReLU } (\omega^T x) = \begin{cases} \text{if } \omega^T x < 0, & 0 \\ \text{else} & , \omega^T x \end{cases}$$

↓  
linear



very simple non-linearity

Sketches for Gradient - Descent:

Initialize all  $v_{ih}$  and  $w_{hj}$  to  $\text{rand}(-0.01, 0.01)$

Repeat

For all  $(\mathbf{x}^t, r^t) \in \mathcal{X}$  in random order

For  $h = 1, \dots, H$

$z_h \leftarrow \text{sigmoid}(\mathbf{w}_h^T \mathbf{x}^t)$

For  $i = 1, \dots, K$

$y_i = \mathbf{v}_i^T z$

} Forward prop.  
Now I have  $y, z$ .

For  $i = 1, \dots, K$

$\Delta \mathbf{v}_i = \eta(r_i^t - y_i^t) z$

For  $h = 1, \dots, H$

$\Delta \mathbf{w}_h = \eta(\sum_i (r_i^t - y_i^t) v_{ih}) z_h (1 - z_h) \mathbf{x}^t$

For  $i = 1, \dots, K$

$\mathbf{v}_i \leftarrow \mathbf{v}_i + \Delta \mathbf{v}_i$  // add v to your existing ones

For  $h = 1, \dots, H$

$\mathbf{w}_h \leftarrow \mathbf{w}_h + \Delta \mathbf{w}_h$  // add w //

Back prop.

Until convergence

First randomness!: Choose small numbers bcs big numbers means you know smth.

Second .. : No batch learning. So, we need to decide how to feed the data, and we feed random.

Forward prop + Back prop = Iteration

Iteration

Use all data in exactly 1 iteration = Epoch

→ Repeat epochs until convergence.



loss function for input  $(x_i, r_i)$

→ The loss function will be different for input  $(x_j, r_j)$



We hope that our local minimums are closer.



But if my loss function is changing so much depending inputs, I have problem.

( $\hookrightarrow$ ) Maybe I can use batch learning → But it takes much time

( $\hookrightarrow$ ) Alternatively, mini batch }  $n$  data → take randomly  $k$  of it



non stable loss functions

$k=n$  } epoch = 1 iteration

$$n=100$$

Nasıl iteration yapmak istiyorum?

$$k=5$$

$$\# \text{epoch} = \frac{100}{5}$$

$k=1$  ise n tane dolanıştır

1 epoch tane belli olur.

$$\# \text{iterations} = \frac{100}{5} \times 100$$

$k$  seciwi böyle close data GPU'ya gider.

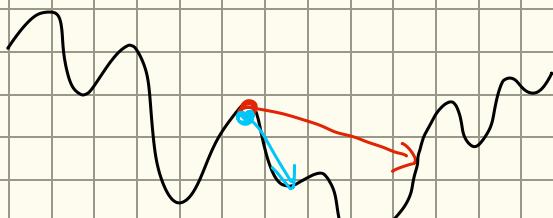
$M$  is hyperparameter.

Larger  $M$

- Faster convergence
- less stable
- Better to be avoid local min.

Smaller  $M$

- Slower "
- Stable
- not good with local min.

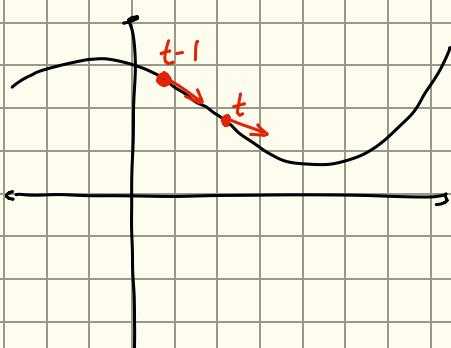


Strategies to set  $M$ :

1-) Start with big steps ( $M \uparrow$ ) and slow it down as you go

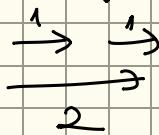


2-) Momentum:



- If both of them going same way, go faster
- " " opposite way, go slower.

Gradient yourself backward momentum ile keyeberiniz.



### 3-) Adaptive learning rate:

$$\Delta \mu = \begin{cases} +\alpha & \text{if } E^{t+1} < E^t \\ -bM & \text{otherwise} \end{cases} \quad // \text{If previous step has more error than current Step, go faster}$$

ilk basa etrafda dolayim.

$\gamma$ : Mercurt oldugum yeme ve gradient tarihe giderken  
guncelliyor

Cok hizli giderken hizlaniyor

Cek yemek -> yemeklar



Convergence iyi olmamak icin



M ne kadar boyuk olursa, o kadar  
hizli giderim

Gradient descent hatali altinda  
dolayla local minimumde durul.

#### 4-) ADAM

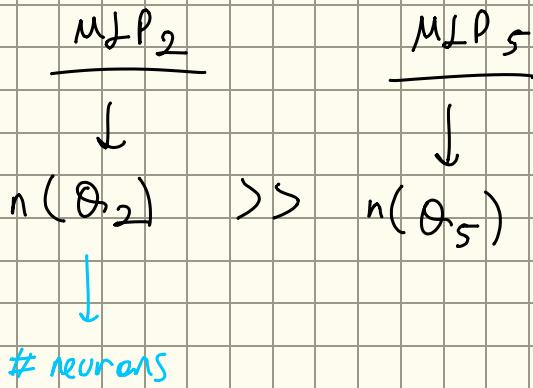
:

↳ They all depending on your loss function and your model shape

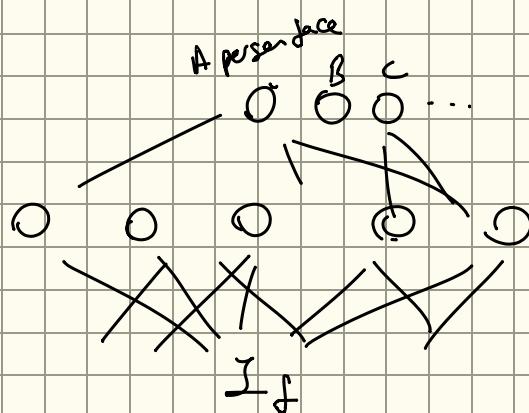
Multiple Hidden Layers :

Theoretically, you can solve any problem with 2 hidden layer MLP  
as opposed to K layer MLP.

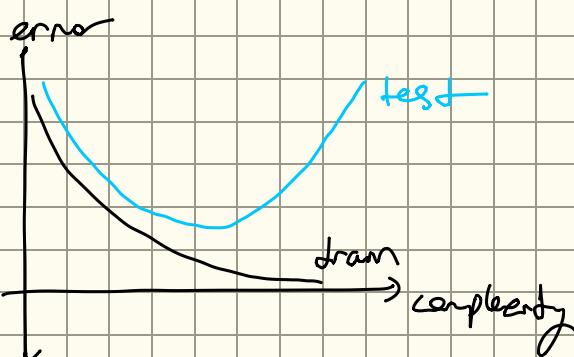
↓  
You would need less  
# of parameters.



face detection problem }



- I can make this single hidden layer to memorize the faces.



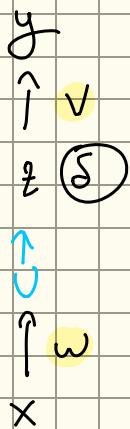
→ Test data can give me high variance  
in this case.

→ It is not practical.

\* So, you can use simpler networks but it will require more complex networks.

\* Using multiple layers may lead to simpler networks.

↳ But # of multiplication will increase in back prop.



→ If you multiply number  $> 1$ , results gonna increase

↳  $2 \times 2 \times 2 = 1$  very big number

↳ Exploding gradient } If you add very big number

$0, 1 \times m \dots$

You cannot learn

$\rightarrow \text{--} \text{--} \text{--} < 1, \text{--} \text{--} \text{--}$  decrease

$$0,2 \times 0,2 \times 0,2 = \downarrow \text{closer to } 0.$$

( $\hookrightarrow$  Diminishing gradient) When weights become 0, you cannot update the weight



You cannot learn.

\* In both cases, these are disadvantages of adding more layers.

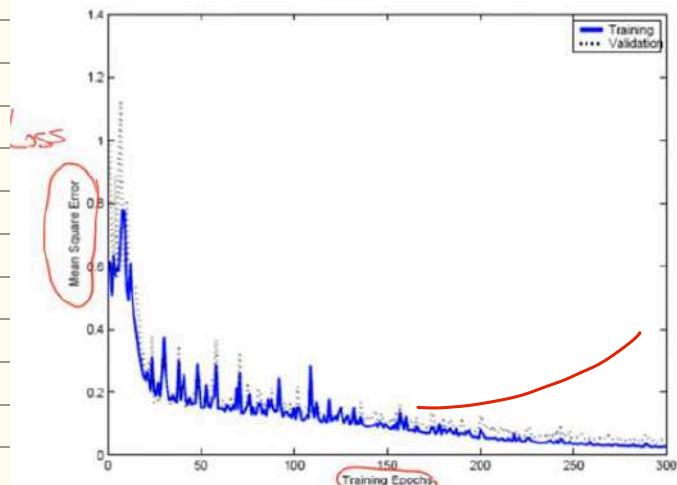
$\downarrow$  solution

Deep learning = Artificial neural networks with multiple layers.

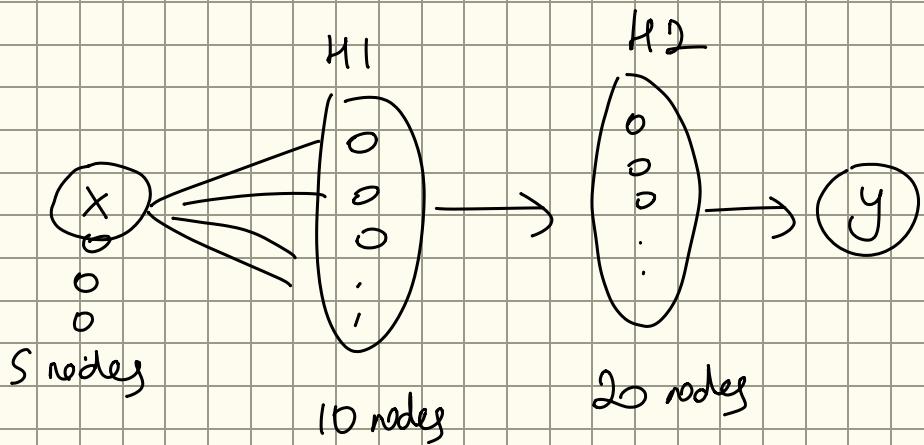
( $\hookrightarrow$ ) Solutions to the problem of diminishing gradient, regularization...

AI  $\rightarrow$  ML  $\rightarrow$  NN  $\rightarrow$  DL  
SS  
AI

\* The best regularization is adding more data.



- Validate like cross-validation } overfitting



$5 \times 10$

connections

$10 \times 20$

/

Solution:



- Structured MLP:

Not fully but repeated connections       $10^9$  neurons  
 $10^5$  connections

- Weight Sharing:

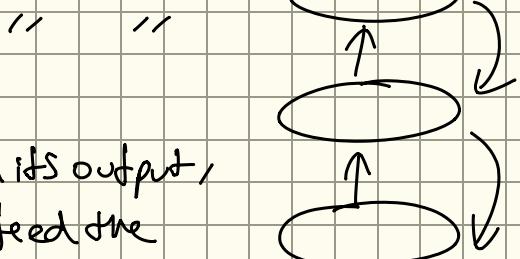
.

.

- Feed Forward Network = MLP

Back prop.

Updated weights.



find its output,  
 feed the  
 next layer.

# CLUSTERING

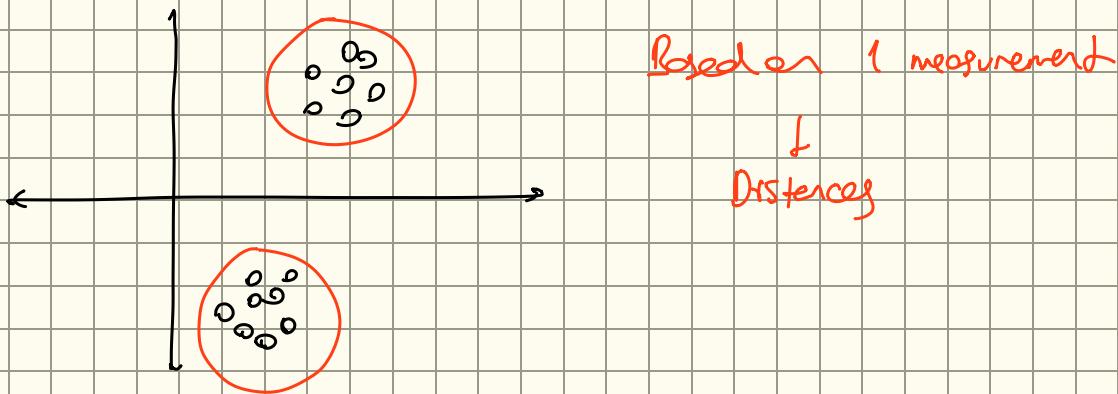
$$x_i \in \mathbb{R}^d \quad i = 1, \dots, N$$

↙

# classes

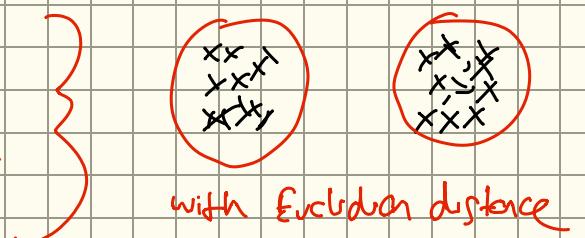
- I don't know how many classes. (fruit types)

Cluster analysis : finding similarities in data



Good clustering: High quality clusters

- High intra-class similarity → in cluster
- Low inter-class ↵ → between =



## Quality Measures:

Data matching

Dissimilarity measures.

$$d(i,j) = d(j,i)$$

Nominal Attributes: Categorical without ordering (red, blue, ...)

Simple Matching:

$$d(i,j) = \frac{p-m}{p} \rightarrow \# \text{ of matches}$$

$\hookrightarrow \# \text{ total variables}$

Binary Attributes:

Contingency Table for Binary Attributes

		Object <i>j</i>		sum
		1	0	
Object <i>i</i>	1	<i>q</i>	<i>r</i>	<i>q+r</i>
	0	<i>s</i>	<i>t</i>	<i>s+t</i>
sum		<i>q+s</i>	<i>r+t</i>	<i>p</i>

7 hrs will be  
big number generally.

→ for symmetric binary attr. each state is equally valuable.

$$d(i,j) = \frac{r+s}{q+r+s+t}$$

Symmetric binary dissimilarity

→ Assymetric //, positive and negative result of disease.  
(not equally important)

$$d(i,j) = \frac{r+s}{q+r+s}$$

} 1-1 matching is more important than 0-0.

**Example 2.18** Dissimilarity between binary attributes. Suppose that a patient record table (Table 2.4) contains the attributes *name*, *gender*, *fever*, *cough*, *test-1*, *test-2*, *test-3*, and *test-4*, where *name* is an object identifier, *gender* is a symmetric attribute, and the remaining attributes are asymmetric binary.

For asymmetric attribute values, let the values *Y* (*yes*) and *P* (*positive*) be set to 1, and the value *N* (*no* or *negative*) be set to 0. Suppose that the distance between objects

**Table 2.4** Relational Table Where Patients Are Described by Binary Attributes

<i>name</i>	<i>gender</i>	<i>fever</i>	<i>cough</i>	<i>test-1</i>	<i>test-2</i>	<i>test-3</i>	<i>test-4</i>
Jack	M	Y	N	P	N	N	N
Jim	M	Y	Y	N	N	N	N
Mary	F	Y	N	P	N	P	N
:	:	:	:	:	:	:	:

Sym.

asym.

$$d(\text{Jack}, \text{Mary}) = d_S(J, M) + d_A(J, M)$$



We can make normalization bcs

$d_S$  have 1 entry,  $d_A$  have 6.

## Chapter 2 Getting to Know Your Data

(patients) is computed based only on the asymmetric attributes. According to Eq. (2.14), the distance between each pair of the three patients—Jack, Mary, and Jim—is

$$d(\text{Jack}, \text{Jim}) = \frac{1+1}{1+1+1} = 0.67,$$

$$d(\text{Jack}, \text{Mary}) = \frac{0+1}{2+0+1} = 0.33,$$

$$d(\text{Jim}, \text{Mary}) = \frac{1+2}{1+1+2} = 0.75.$$

dissimilarity  
with  
binary

These measurements suggest that Jim and Mary are unlikely to have a similar disease because they have the highest dissimilarity value among the three pairs. Of the three patients, Jack and Mary are the most likely to have a similar disease. ■

Combining binary and real number?

$$\Delta w + \Delta c$$



$$[0, 1] \quad \{0, 1\}$$

Interval features: age [0, 20], height [20, 150] etc.

Let assume there are 2 features in set.

$[-50, 20]$  } They are not comparable.

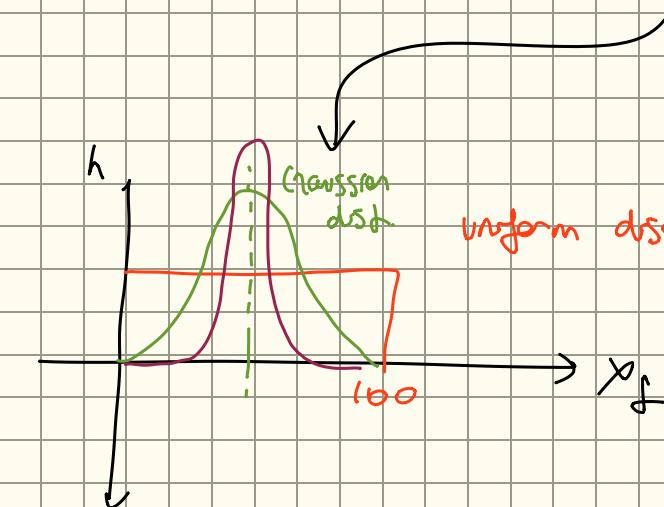
$[700, 1200]$  We should standardize data

### (1) z-score

Why not this?

$$z = \frac{x - m}{s} \quad \text{or} \quad \frac{x - m}{\text{dist}(\text{max} - \text{min})}$$

(mean absolute deviation)

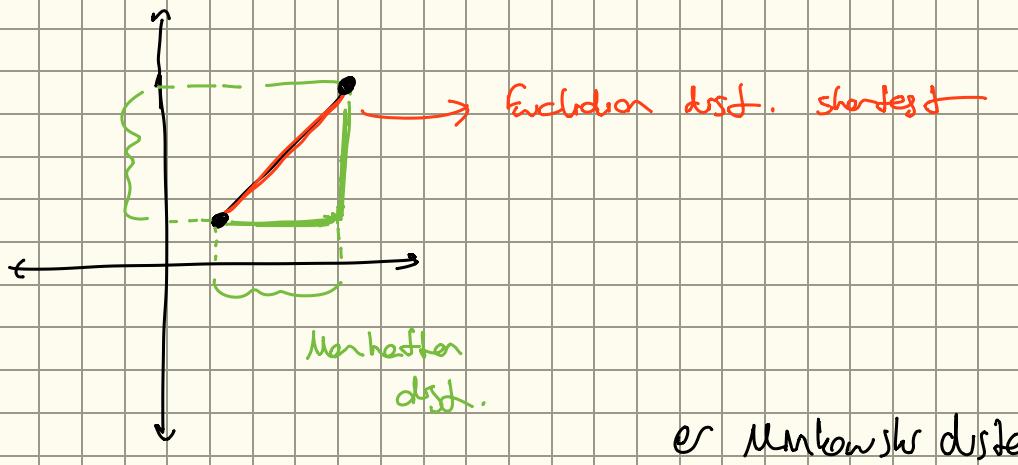


How can I normalize it between -1, 1?

$$\frac{x_f - 50}{100} \rightarrow \text{mean}$$
$$100 \rightarrow \text{distribution}$$

Gaussian from 100 normalizable elements. Center exist  
discontinuity.

By z-score  $\frac{x - m}{s}$  will make  $\frac{x - m}{\text{dist}(\text{max} - \text{min})}$  i.e. same  
data type our.



What happens on binary features? (Not intervals)

$T \rightarrow$  is more "important" than  $F$ :

$\begin{matrix} T \\ F \end{matrix}$  is more important to contribute similarity.

} asimetrik

Build contingency matrix

Dissimilarity for symmetric:

" " asymmetric:

Jaccard distance: important is  $T-T$

$$sim_j(i, j) = \frac{a}{a+b+c}$$

Asymmetric binary similarity

Decide symmetric or asymmetric?

↳ look at the data

→ Meget "cough" bør danne andre data fåle obbetur.

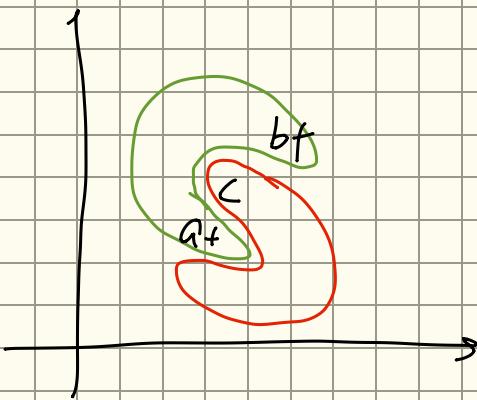
→ "gender" her varian symmetric altså bør → meget engerning aabenbundet

Categorical features? (Nominal variables)

Simple Matching:

$$d(i,j) = \frac{p-m}{p} \rightarrow \# \text{ of matches}$$

$\hookrightarrow \# \text{ total variables}$



$$d(a,b) < d(a,c)$$

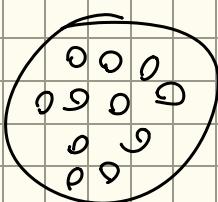
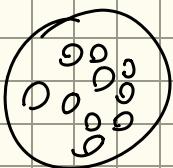
$$< d(b,c)$$

$$d_e(a,b) \gg d(a,c)$$

$$d(b,c)$$

No way to use  
these metrics.

Calculating distances between Clusters



Single Link: Shortest dist

(Complexity:  $n^2$ )

Complete Link: Longest dist,  $n^2$

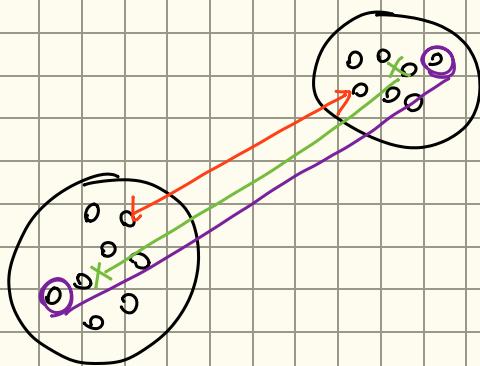
Average Link: Avg dist.  $n^2$

Centroid:  $\frac{1}{n} \sum_{i=1}^n x_i$

Medoid: An actual point  $\underline{\underline{n^2}}$

one closer, centrally located object.

How can we measure the quality of within cluster?



1-) average distance to center } works only  
when  $x_i \in \mathbb{R}^d$

gender, age, place, salary

2-)  $\frac{1}{n \times n} \sum_{i,j} d(x_i, x_j)$  } dissimilating

3-) maximum distance } K-means

Evaluate clusters } Goodness within + Goodness across

↓  
depends on shape of clusters  
Distance / dissim. metrics

Supervised:

$$\text{avg map } \sum_{i=1}^N g(y_i, f(x_i))$$

(↳ This goodness is difference b/w. what you generate and the actual.)

Unsupervised:  $\rightarrow$  # of clusters

$$\arg \min_{k, c_k \in F} \sum_{k=1}^K \sum_{i=1}^N g(c_k(x_i))$$

some mapping func.

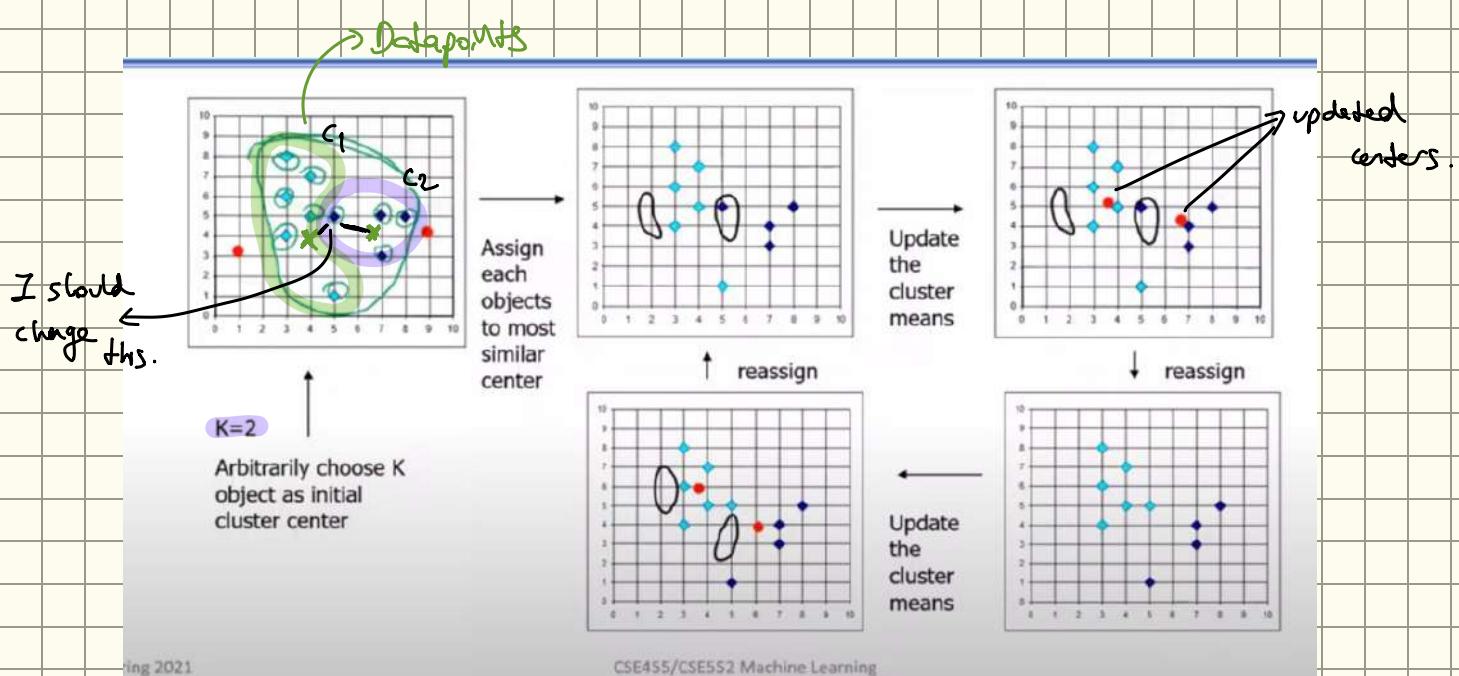
(↳ This goodness is actually calculating.)

# Partitioning Algorithms

Kmeans : Start with random center  
Update center in each iteration

$$l \times (n \times 2 + n) \rightarrow O(n \times l \times k)$$

clusters  
↓  
data      ↓  
          num of iterations



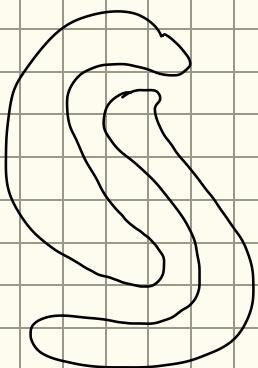
→ Greedy algo.

→ find means, label it with its closer mean.

→ Algo starts with random 2 points as centers.



fire with kmeans, **ignores shape**



### Wealthiness

- Applicable only when mean is defined
- Determining k
- Sensitive to outliers
- Only with numeric variables (can calculate distance)

↳ one hot encoding

**k-medoids**: Instead of taking the mean, we can use medoids,

↓  
most centrally located object

kmeans → calculate mean  $O(n)$ , fast on certain problems

kmedoids → " medoid  $O(n^3)$  ?

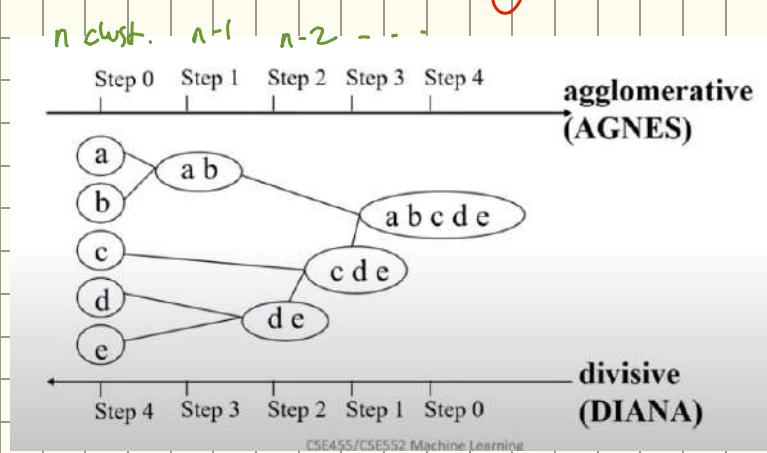
|| apl

**PAM**:  $O(k(n-k)^2)$

aims to minimise the sum of dissimilarities  
btw points and their medoids.

**CLARA**:  $O(kS^2 + k(n-k))$  large datasets.

# Hierarchical Clustering!



- In the beginning, Each data point 1 cluster.

- Step 0  $\rightarrow n$

- Step 1  $\rightarrow n^2$  combinations

$\otimes$  Doesn't require any number  $k$ , but needs termination conditions.

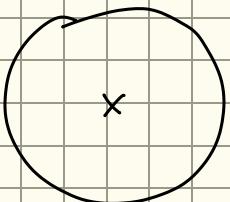
- Don't scale well.  $\rightarrow O(n^2)$  at best

- When you make decision, you can't undo (like decision trees)

$\otimes$  Combine hierarchical with distance-based clustering } - BIRCH

} tree like search

- CHAMELEON



outlier

- K-means sensitive to outliers
- ↓
- k-medoid  $\checkmark \rightarrow O(n^3)$  to find medoid.

Density - Based Methods:  $O(n^2)$

### DBSCAN

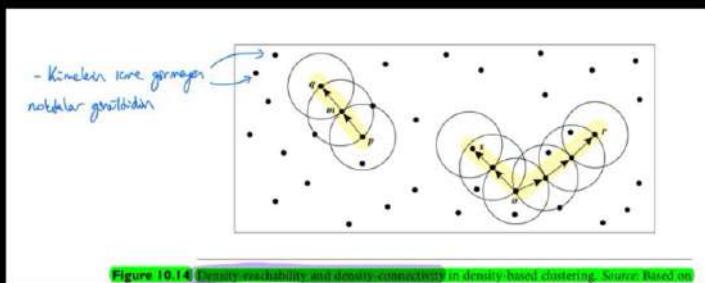


Figure 10.14 Density-reachability and density-connectivity in density-based clustering. Source: Based on

### Parameters:

Eps: Map radius of neighborhood

MinPts: Min num of points in an Eps - neighborhood of that point

(According to this I can say this border is dense enough.)

Algorithm :

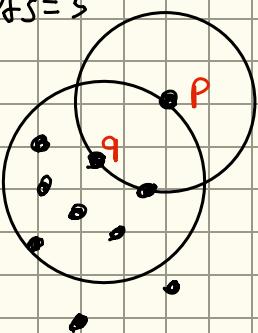
- 1-) Select point  $P \rightarrow$  visited
- 2-) density - reachable
- 3-) If  $p$  is core point , a cluster is formed
- 4-) If not , no points density reachable,  
find new point
- 5-) Continue until visit all points

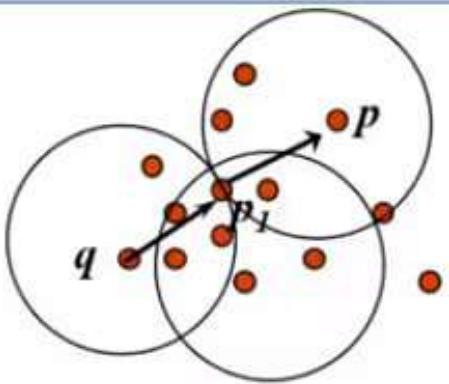
$$\text{Eps} = 1 \text{ cm}$$

$$\text{MinPts} = 5$$

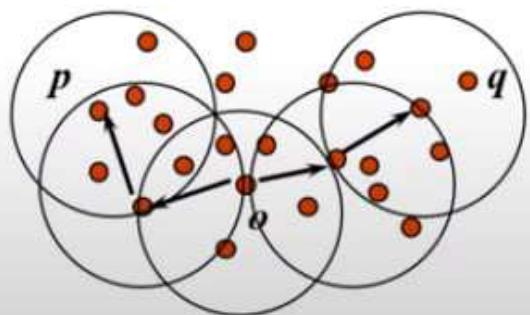
Density reachable if  $p$  belongs to  $N_{\text{Eps}}(q)$

then, Core point if  $|N_{\text{Eps}}(q)| \geq \text{MinPts}$





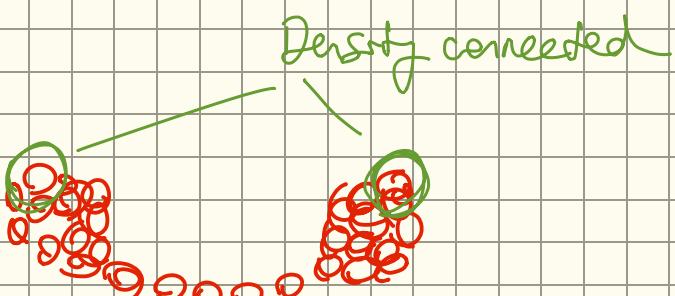
$q - p_1 - p \}$  density reachable  
chain of points



Density connected

Both  $p$  and  $q$  are density reachable  
from point  $o$ .

So  $p$  and  $q$  are density connected.



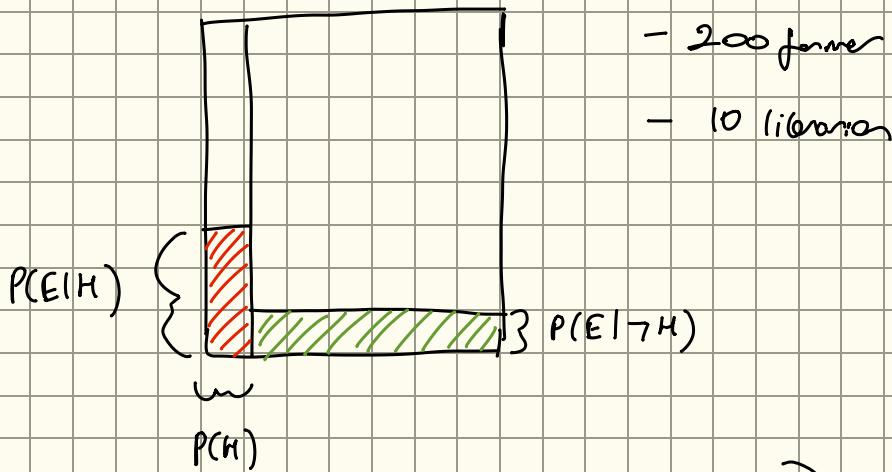
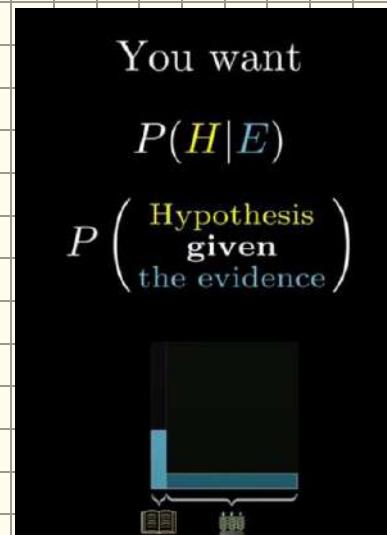
# BAYES

- Verenegenim herarı koşulu doğrudanın zemini, ökülmess en kontroversel oluyor.

$$P(\text{ciftler} \mid \text{süren}) = \frac{20}{24}$$

$$P(\text{kıraklıclar} \mid \text{süren}) = \frac{4}{24}$$

$$P(H|E) = \frac{P(H) \cdot P(E|H)}{P(E)}$$



Prior

$$P(H) = \frac{\text{librarians}}{\text{All people}} = \frac{1}{21}$$

$$P(E|H) = \frac{4}{10}$$

} The prob. of Steven in given librarians

$$P(E|\neg H) = \frac{20}{200}$$

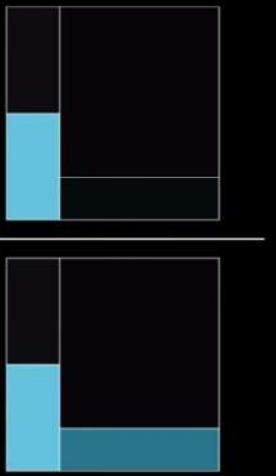
} The prob. of Steven in given farmers.

$$P(H|E) = \frac{\frac{4}{10}}{\frac{4}{10} + \frac{20}{200}}$$

The prob. of librarians in given Steven.

# Bayes' theorem

$$P(H|E) = \frac{P(H)P(E|H)}{P(E)}$$



Binary Classification

$$\begin{aligned} P(+|x) &+ \\ P(-|x) &= 1 \end{aligned}$$

log max  $P(C|x)$

Problem :

$T$	$P$	$H$	Conj / Not C.
high			
low			
medium			

$P(C|T, P(H))$  ?

Now, let's just look at  $T$ .

(to bayes formula inc.)

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

$$P(C|T = \text{high}) = \frac{P(C) \cdot P(T_h|C)}{P(T = \text{high})}$$

$$= \frac{\frac{4}{5} \cdot \frac{2}{4}}{\frac{3}{5}}$$

$T$	$C$
h	+
h	+
h	-
l	+
l	+

$$P(C|T = \text{med}) = \frac{\frac{4}{5} \cdot 0}{?}$$

$$P(C|T = \text{low}) = \frac{\frac{4}{5} \cdot \frac{2}{4}}{\frac{2}{5}}$$

$x_1 = \text{t hrs about leg circle}$

$$Dg = 100 \circ$$

$$\rho_g = 110 \circ$$

$x_2 = \text{" " " single 1"}$

:

:

:

:

$x_3 = \text{" " " 2 single 1"}$

:

$x_4 : \text{" " " 2 " " "}$

for all O's what's  
the prob. of the see this  
pattern?

$$\rightarrow \frac{1}{10}$$

$$\uparrow P(d_0 | x_1 x_2 x_3 x_4) = \frac{P(x_1 x_2 x_3 x_4 | d_0) P(d_0)}{P(x_1 x_2 x_3 x_4)}$$

$$P(h_1) \quad P(h_2) \rightarrow P(h_1 | \dots) \uparrow$$

bigg      or

Naive Bayes

$$P(x_1=+, x_2=+ | \text{cancer})$$

$$P(x_1=-, x_2=+ | \text{cancer})$$

;

$$\text{So, } P(x_1, x_2, x_3 | v)$$

=

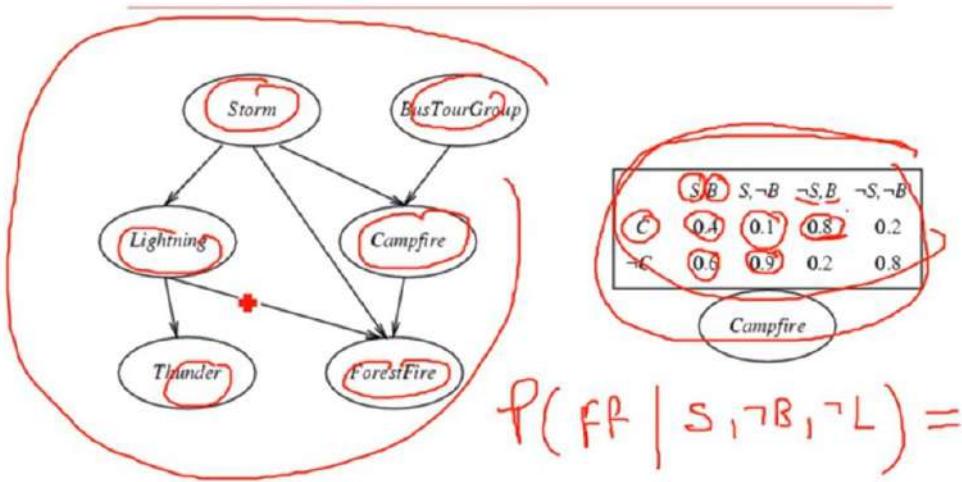
$$P(x_1 | v) \cdot P(x_2 | v) \cdot P(x_3 | v)$$

conditionally independence.

,  $\rightarrow$  complete

$$P(h | x_1 \dots x_d) = \arg \max_{h \in H} P(h) \cdot \prod_{i=1}^d P(x_i | h)$$

# Bayesian Belief Network ; No loops



CSE455/CSE552 Machine Learning

In continuous data prob  $\rightarrow$  1.) Discretize it

↳ if age make groups

↳ if salary .. ..

$$P(x_1 | ?_{+/-})$$



$2 \times 100$  num of prob.

$\begin{matrix} + \\ - \end{matrix}$  hypothesis       $\begin{matrix} + \\ - \end{matrix}$  group num  
 $\begin{matrix} + \\ - \end{matrix}$  for salary

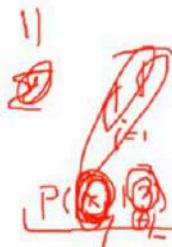
discretize  $x_1 = [0, 100000]$   
 $\downarrow = \{0, 1, \dots, 100\}$

Insurance mt:

$$\left\{ \begin{array}{l} x_1 = \text{Balance} \\ x_2 = \text{Savers} \\ x_3 = \text{Age} \\ x_4 = \text{Occupation} \end{array} \right.$$

no explicit linear  
variables.

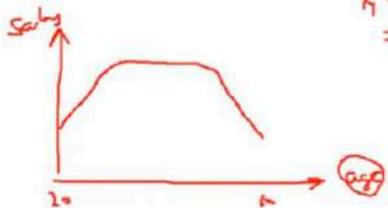
$$P(?) \mid x_1 = 5k, x_2 = 10k, x_3 = 20-25, x_4 = s$$



Naive Bayes  
(non-conditioned)

$$P(x_1 | ?) P(x_2 | ?) P(x_3 | ?) P(x_4 | ?) P(?)$$

# Thanks for listening!



$$x_1 = [0, 100, 200] \\ = \{0, 1, 2, \dots, 100\}$$

$$2 \times 100 \quad 2 \times 100 \quad 2 \times 20 \quad 2 \times 10$$

$$(100) \times 100$$

If conditioned, they will grow very quickly.

## Final Aitomsgen

Q2)

Dissim

Asym

$$\begin{array}{r} 3 \\ - \\ 3 \end{array}$$

Sym

$$\begin{array}{r} 1 \\ - \\ 3 \end{array}$$

Paccord

$$0$$

O(h) < No, Merged, Middle >

$$P(\text{No} | \langle \text{No}, m-, m \rangle)$$



$$P(\text{No}) \cdot P(\langle \dots \rangle | \text{No})$$

$$P(\dots)$$

$$\frac{7}{10} \cdot$$

$$P(\text{Yes} | \langle \dots \rangle)$$

$$\frac{P(\text{Yes}) \cdot P(\langle \dots \rangle | \text{Yes})}{P(\langle \dots \rangle)}$$

$$\frac{3}{10} \cdot$$

$$\frac{1}{3} \cdot \frac{1}{3} \cdot 0$$

$$\frac{4}{7} \cdot \frac{4}{7} \cdot \frac{3}{2}$$

$$0$$

b) Gaussian distribution

Low | Middle: 2 | High: 3

~~1-~~) Calculate mean and variance for Yes and No of income.

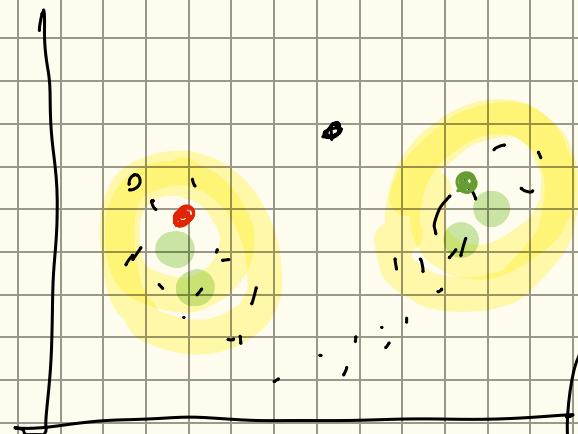
2-) Calculate  $P(\text{Middle} = 2 | \text{Yes})$  with prob. density function.

$$P(\text{Middle} = 2 | \text{Yes})$$

3-)  $P(\text{No} | \langle \text{No}, \text{Med.} \rangle) =$  Dög-sykerligini, mane yine  
forlich yoruduk.

$$\frac{P(\text{No}) \cdot P(\langle \dots \rangle | \text{No})}{P(\langle \dots \rangle)}$$

Q3) b)

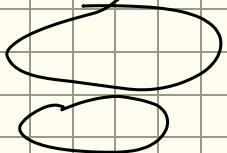


1-) Decide k and choose random  
2 medord.

2-) Find clusters

3-) Look data points, if quality better  
swap medord

k means

Basiswieder her ren ,  
 esig le skeletohs .

