

# Machine Learning

- Covered topic
- Identification trees
  - Arranging tests
  - Measures of disorder
  - Numeric tests
  - Rules
  - Transformations
- Some test are expensive
- Note: a good 2D tree has lower depth and simple structure
- Disorder Set
- positive total negative

$$D(\text{set}) = -\frac{P}{T} \log_2 \frac{P}{T} - \frac{N}{T} \log_2 \frac{N}{T}$$

If we have a unbiased data set:

$$\frac{P}{T} = \frac{1}{2} : D(\text{set}) = -(\frac{1}{2} \log_2 \frac{1}{2}) \times 2 \quad \text{e.g. } \begin{matrix} \square \\ \square \end{matrix}$$

If we have a fully total data:

$$\frac{P}{T} = 1 : D(\text{set}) = -1 \log_2 1 - 0 \log_2 0 = 0, \text{ e.g. } \begin{matrix} \square \\ \square \end{matrix}$$

Disorder test:

$$D(\text{test}) = \sum \text{Dis}(\text{set}) \times \frac{\text{set size}}{\text{Total}}$$

Branches

— O'Cam's razor: the simplest explanation is often the best explanation

$D(\text{set})$  or  $H(X) = -[\sum_{i=1}^N P(i) \ln(P(i))]$

Entropy in information theory

Note: P 越大, 熵 越小  
即事件经常发生, 因此对应的信息量越少

1.0 0.5 信总 = 2 4 bit

1111, N 种可能 = 2

$\Rightarrow$  信总 =  $\log_2(N)$

$\Rightarrow$   $= \log_2(2)$

## Decision Tree Example

Data: Vampire Data in Romanian

Vampire	Shadow	Garlic	Complexion	Accent
No	?	Y	Pale	None
No	Yes	Y	Ruddy	None
Yes	?	N	Ruddy	None
Yes	N	N	Average	Heavy
Yes	?	N	Average	Odd
No	Y	N	Pale	Heavy
No	?	N	Average	Heavy
No	?	Y	Ruddy	Odd

Goal: Use some test result of characteristic to determine whether it's an Vampire or not.

Step1: Compute the disorder of each feature test. (Actual, try to use your intuition at first).

Shadow

Disorder Score

Yes/No

Yes: 3/8, No: 1/8

Disorder Score:  $\frac{3}{8} \log_2 \frac{3}{8} + \frac{1}{8} \log_2 \frac{1}{8} = \frac{1}{2} = 0.5$

$D(\text{test}) = \frac{3}{8} \log_2 \frac{3}{8} + \frac{1}{8} \log_2 \frac{1}{8} = 0.5$

Complexion

Pale/Ruddy

Pale: 1/4, Ruddy: 3/4

Disorder Score:  $\frac{1}{4} \log_2 \frac{1}{4} + \frac{3}{4} \log_2 \frac{3}{4} = 0.915$

$D(\text{test}) = \frac{1}{4} \log_2 \frac{1}{4} + \frac{3}{4} \log_2 \frac{3}{4} = 0.915$

## Step2: Eliminating the root feature and build the next level for decision tree:

Garlic

Yes/No

Yes: 3/8, No: 1/8

Disorder Score:  $\frac{3}{8} \log_2 \frac{3}{8} + \frac{1}{8} \log_2 \frac{1}{8} = 0.5$

$D(\text{test}) = 0.5$

Step3: Construct the Complete disorder-minimizing Identification/Decision tree (Note: In practice, the tree can be deeper than 2 levels, so you might need to take more step to reach here).

Shadow

Yes/No

Yes: 3/8, No: 1/8

Disorder Score:  $\frac{3}{8} \log_2 \frac{3}{8} + \frac{1}{8} \log_2 \frac{1}{8} = 0.5$

$D(\text{test}) = 0.5$

Complexion

Pale/Ruddy

Pale: 1/4, Ruddy: 3/4

Disorder Score:  $\frac{1}{4} \log_2 \frac{1}{4} + \frac{3}{4} \log_2 \frac{3}{4} = 0.915$

$D(\text{test}) = 0.915$

Boundary threshold for KNN

Note: After you have a Decision tree, you can use it to develop some rules/policies with if-else clauses.

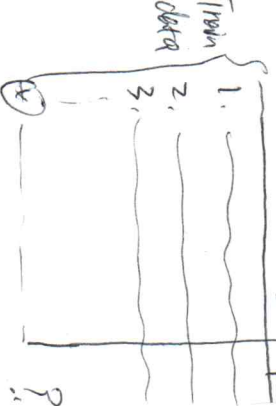
# Machine learning

## Supervised Machine learning

Features  $X_1, X_2, X_3$

Classification

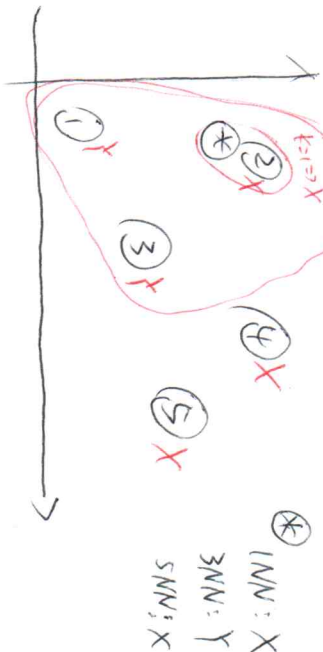
- Boolean: Y/N
- Discrete: e.g. color, party, ...
- Continuous: Height, weight



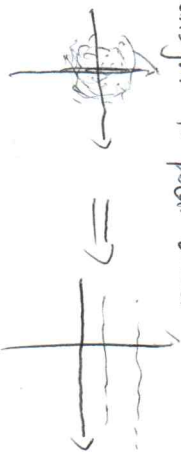
## Supervised ML objective:

Given training data with known features ( $X$ ) and known classification ( $Y$ ), and produce some rule for classifying new points with known features and unknown classification.

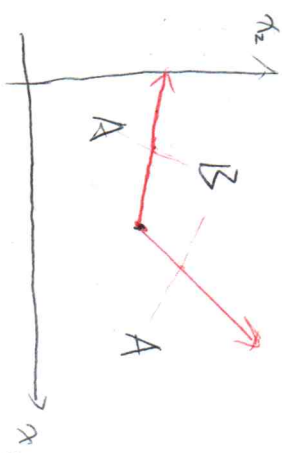
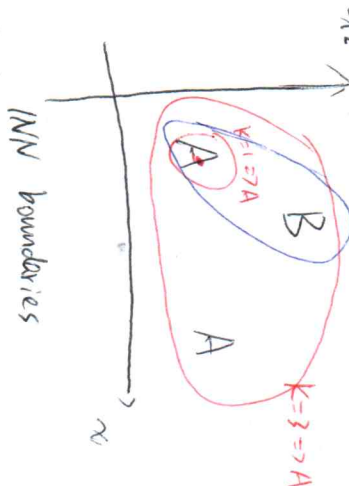
You can also working backward to figure out the training data:



5) Transform to polar coord



K-Nearest Neighbor  $K=2 \Rightarrow$  tie



Other thing we have control over for

K-NN is Distance Metric:

1) Euclidean distance:  $D(\vec{u}, \vec{v}) = \sqrt{\sum (u_i - v_i)^2}$

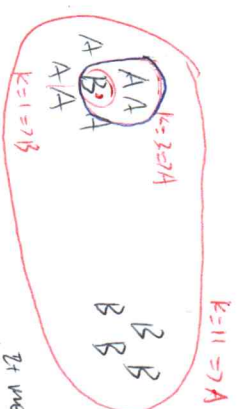
(Note: this is the default method, which simply find the straight line distance between two point, aka geometric distance)

2) Manhattan distance:  $D(\vec{u}, \vec{v}) = \sum |u_i - v_i|$

(Note: aka taxicab distance, used in case like classifying people based on what grocery store one they lived closest to, for advertising purpose, so people can't go across the block but have to go along the streets.)

3) Hamming distance:  $D(\vec{u}, \vec{v}) = \sum \text{diff?}(u_i, v_i)$

(Note: It's useful for comparing symbolic things. It's just looking at whether the vector components are same or not. e.g. like comparing same children's toys:  $\vec{u} = (\Delta, \text{red}, \text{soft}) \Rightarrow 0+1+0=1$   
 $\vec{v} = (\Delta, \text{blue}, \text{soft})$



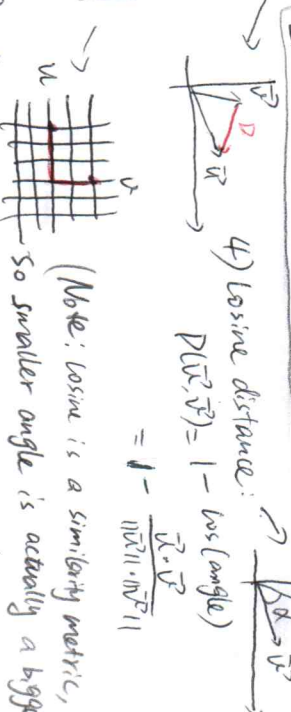
K=1 Overfitting

It means you overly rely on your data, so you're assuming even small fluctuations in the data are significant.

How to find a good value of K? much, so you're not relying enough on the data. As a result you're ignoring the difference in the data should be

Idea: set aside part of training data, then training using the remaining data, then we'll test whether the data we set aside gets classified correctly. e.g. leave-one out CV, K-CV.

4) Cosine distance:  $D(\vec{u}, \vec{v}) = 1 - \cos(\text{angle})$   
 $= 1 - \frac{\vec{u} \cdot \vec{v}}{\|\vec{u}\| \cdot \|\vec{v}\|}$



(Note: cosine is a similarity metric, so smaller angle is actually a bigger number, so we use the negation, so that's effective larger angle means large dist. It's useful if you care about the ratio between the features than actual quantities. E.g. Recipe.

