

AS K INCREASES THE NUMERATOR
INCREASES OR DECREASES. AND
THIS ALSO HAPPENS FOR THE
DENOMINATOR.

(THEORETICALLY)
THE BEST IS THE ONE WITH
THE HIGHEST CH VALUE

S' IS THE IS STARTING ARRAY TO CALCULATE
THE CENTERS

CW IS THE ACTUAL CENTER,

A ARRAY TELLS YOU WHAT GROUP EACH
DATA GROUP WILL BE.

MACHINE LEARNING 9/16

- HOW TO DECIDE THE NUMBER OF K IN CLUSTERING

WITH CLUSTER VARIATION THERE IS A PROBLEM BECAUSE IT JUST KEEPS DECREASING.

BETWEEN CLUSTER VARIATION

$$B = \sum_{k=1}^K n_k \|\bar{x}_k - \bar{x}\|_2^2$$

n = CAN BE WEIGHTED TO MATCH THE # OF SAMPLES

PROBLEM BECAUSE AS K KEEPS INCREASING, SO DOES CLUSTER VARIATION

$w(k)$ - WITHIN CLUSTER VARIATION

$B(k)$
BETWEEN CLUSTER VARIATION

$$CH(k) = \frac{B(k)/(k-1)}{w(k)/(n-k)}$$

S O₁ O₂ O₃ ... O₁₀

| | | | | | |
|---|---|---|---|---|---|
| d | d | d | d | d | - |
| d | d | d | d | d | |

DOWNLOAD
IS A WRITE
OPERATION

OWNER CAN CHANGE THE PERMISSION OF A FILE. OWNER RIGHT TO CHANGE OF OWNER ONLY ROOT CAN CHANGE THE RIGHT OF OWNERS.

MAKE SURE CONCURRENCY COULD HAPPEN. RACE CONDITION : --- ETC. MANAGER CALLS INC-DE or DECL-ATR AND THOSE IMPLEMENT THE FUNCTIONS.

ACCESS CONTROL MECHANISM

ACCESS CONTROL LIST

CAPABILITIES

RIGHT BASED ACCESS CONTROL

IMPLEMENTING

- DEFINE SECURITY POLICY

- CREATE ACCESS CONTROL MATRIX

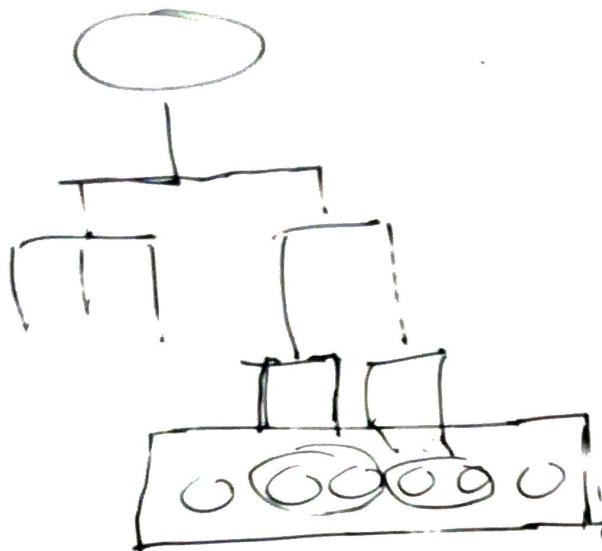
- IMPLEMENT AND PUT ACM IN STORAGE

- DEFINE PROCEDURE TO

ACCESS CONTROL LIST

EACH FILE HAS ITS OWN PERMISSIONS.

HIERARCHICAL CLUSTERING 9/6



MERGE THE CLOSEST TWO POINTS
TOGETHER THEN GROUPS = $m-1$

REPEAT IT AND YOU WILL HAVE
 $m-2$ GROUPS

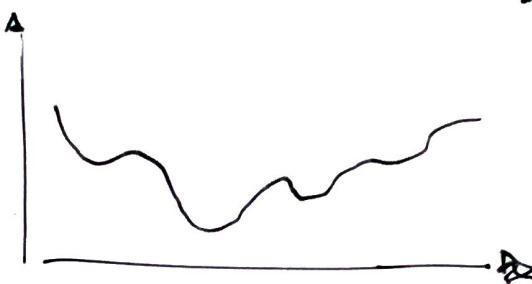
CONTINUING DOING THIS UNTIL
YOU HAVE ONE GROUP.

GIVE K AND WE WILL "CHUP"
THE RESULTS GIVING US
THE # OF GROUPS AND DATA
POINTS IN THEM.

CONVEX



HOW TO CHOOSE K
K MEANS



- FINDING THE LOCAL MINIMUM MIGHT BE RETURNED BY THE K MEANS ALGORITHM
- TO AVOID THIS RUN THE ALGORITHM MULTIPLE TIMES.

K MEANS IS SENSITIVE TO OUTLIERS.

K MEANS IS K-MEANS

K MEANS GETS NEW CENTER AS A POINT.

2) REPEAT UNTIL CONVERGENCE

{ FOR EVERY i , SET

$$c := \arg\min_{m_j} \|x^i - m_j\|^2$$

| x_i | $x_2 \dots x_n$ | c (cluster) |
|-------|-----------------|---------------|
| n | | |

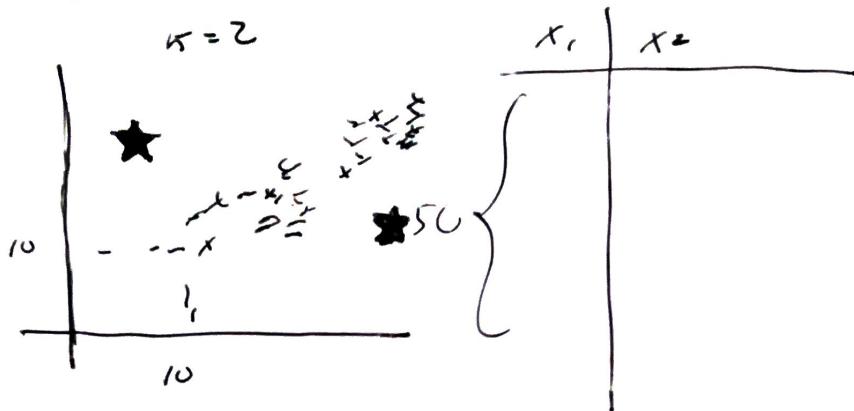
$$f = \arg\min (x - z)^2 \quad \text{WHAT VALUE OF } c \text{ WILL MAKE}$$

3) FOR EVERY j , SET

$$m_j := \frac{\sum_{i=1}^m 1 \cdot z \cdot c^{(i)} y_i x^{(i)}}{\sum_{j=1}^m 1 \cdot c^{(i)} y_i x^{(i)}}$$

MACHINE LEARNING

FOR THE K ALGORITHM YOU
MUST STATE THE NUMBER OF
CLUSTERS YOU WANT TO GROUP THE
DATA IN.



- EACH DATA POINT WILL CALCULATE THE DISTANCE TO EACH RANDOM POINT.
- THE CLOSEST RANDOM POINT WILL BELONG TO THAT GROUP.
- THE NEW CENTER WILL BE CALCULATED FOR EACH GROUP
- THEN THE PROCESS STARTS AGAIN AND CONTINUES UNTIL THE CENTERS ARE NO LONGER CHANGED.

K-MEANS:

INPUT $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

① INITIALIZE (CLUSTER CENTROIDS)

M_1, M_2, \dots, M_K [ERR]

| | (1,1) C1 | (1,2) C2 | GROUP |
|-------|-------------|-------------|-------|
| (1,1) | 0 | 1 | C1 |
| (2,1) | 1 | 0 | C2 |
| (1,3) | 3.61 | 2.81 | C2 |
| (5,4) | 5 | 4.24 | C2 |

NEW CENTER IS AVERAGE OF EACH FEATURE DIVIDED BY # OF POINTS

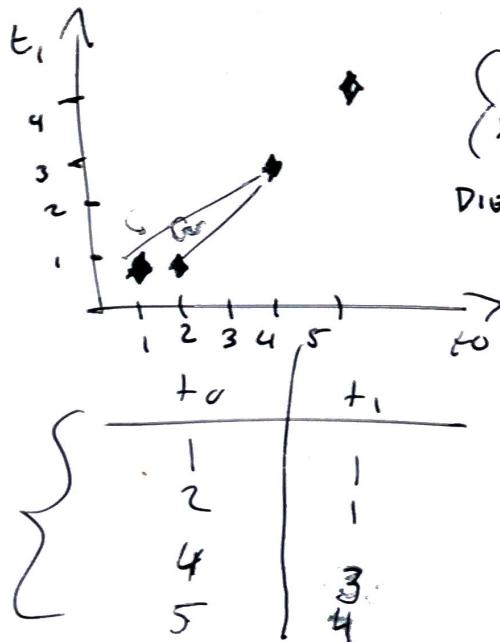
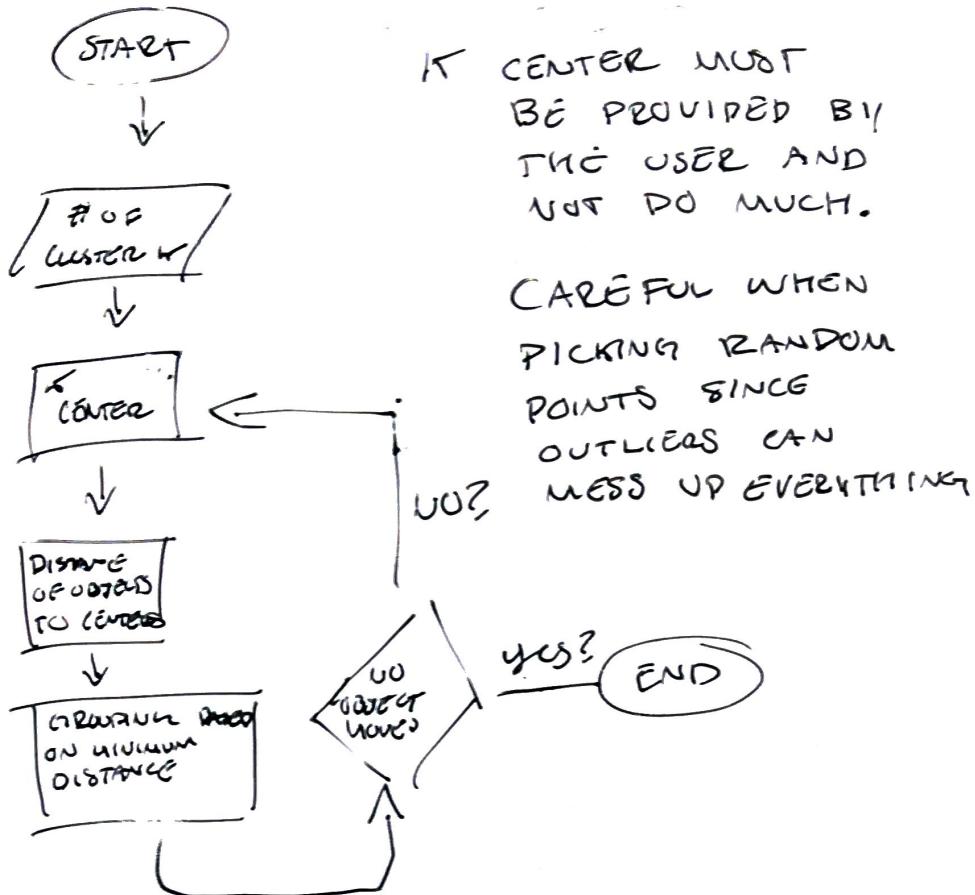
$$C_2 = \left(\frac{2+4+5}{3}, \frac{1+3+6}{3} \right)$$

$$C_1 = \left(\frac{1}{3}, \frac{1}{3} \right) = (1, 1)$$

THE PROCESS REPEATS WITH THE NEW CENTERS.

UNTIL THE CENTERS DO NOT CHANGE.

THIS METHOD MAY NOT WORK WITH DATA THAT HAS MULTIPLE DIMENSIONS.



$$k = 2$$

$$(x_1, y_1)$$

$$(x_2, y_2)$$

$$Dist = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

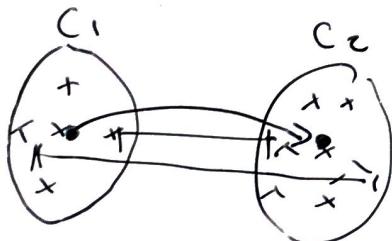
APPLICATIONS

- PATTERN RECOGNITION

WHAT IS GOOD CLUSTERING?

- HIGH QUALITY CLUSTERING
 - HIGH INTER CLASS SIMILARITY
 - LOW BETWEEN CLASS SIMILARITY
- QUALITY IS A RESULT ON SIMILARITY MEASURE USED AND ITS IMPLEMENTATION.
- ALSO ABILITY TO DISCOVER SOME OR ALL PATTERNS.

CALCULATE THE CENTER OF THE CLUSTERING.



DISTANCE BETWEEN CLUSTERS

- CENTER OF BOTH
- THE TWO CLOSEST POINTS
- THE TWO FURTHEST POINTS

PARTITIONING METHOD

K-MEANS ALGORITHM.

K - # OF CLUSTERS OR GROUPS

MACHINE LEARNING 914

TRAINING DATA

| YEAR | MILEAGE | ... |
|------|---------|-----|
| 1998 | 200,000 | |

$n = 108$

SUPERVISED

- PROVIDE THE ANSWER

UNSUPERVISED

- NO ANSWER IS PROVIDED

m - # OF TRAINING SAMPLES

n - DIMENSION OF DATA SAMPLES. (NO ANSWER)

X - TRAINING DATA SET

~~$X^{(i)}$~~ - TRAINING SET FOR A SAMPLE

Y - LABEL / ANNOTATION

$X \rightarrow Y$ RELATIONSHIP
(SUPERVISED)

X_j - j^{th} FEATURE

WHAT IS CLUSTERING?

- SIMILAR DATA GROUPED
TOGETHER

X_j^i - CONSTANTLY USED
TOGETHER
TO PINPOINT
A CERTAIN PIECE OF
DATA

FIND THE BEST FIT LINE TO FIND EQUATION AND USE THAT TO PREDICT SELLING PRICE.

"REGRESSION PROBLEM" USED WHEN PREDICTED VALUE IS CONTINUOUS
CAN IT USE WITH DISCRETE VALUES

CLASSIFICATION PROBLEM

CANCER - PROVIDE DATA WITH RIGHT ANSWER AND THEN USE THAT TO PREDICT WHAT KIND OF CANCER.

+ BENIGN - MALIGNANT
DISCRETE VALUES

UNSUPERVISED LEARNING

- DATA WITH NO LABEL BUT THE ALGORITHM WILL CLUSTER TOGETHER SO IT MAKES SENSE.
- NO ANSWER IS GIVEN

ROLE OF STATISTICS: INFERENCE FROM A SAMPLE

ROLE OF COMPUTER SCIENCE: EFFICIENT ALGORITHMS TO MAKE INFERENCES.

SOME TYPES

LEARNING ASSOCIATIONS: RELATIONSHIP IN THE DATA

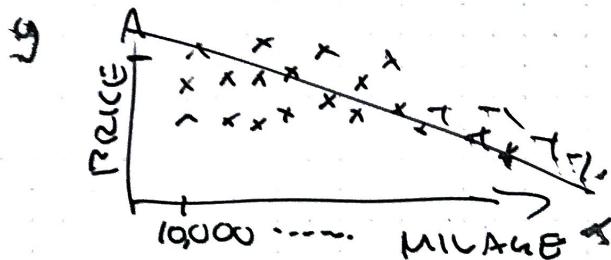
SUPERVISED: GIVE EXAMPLES WITH LABELS

UNSUPERVISED: GIVE DATA AND THE MACHINE FINDS PATTERNS.

REINFORCEMENT LEARNING: LEARNS A POLICY THAT REWARDS FOR BEING CORRECT.

SUPERVISED

- SELL CAR: COLLECT SELLING DATA FROM PAST (RIGHT ANSWER MUST BE PROVIDED)



CS 4347: MACHINE LEARNING

MACHINE LEARNING WHEN TO USE IT?

CAN USE TO PREDICT STOCKS TAKING INTO ACCOUNT PAST STATES.

WHY LEARN?

- PROGRAM TO OPTIMIZE A PERFORMANCE CRITERION USING DATA OR PAST EXPERIENCE
- WE NEED TO LEARN TO CALCULATE PAYROLL
- LEARNING USED
 - HUMAN EXPERTISE DOES NOT EXIST
 - HUMANS ARE UNABLE TO EXPLAIN THEIR EXPERTISE
 - SOLUTION CHANGES IN TIME
 - SOLUTION NEEDS TO ADAPT TO PARTICULAR CASES. (FACE RECOGNITION)

WHAT WE TALK ABOUT WHEN SAYING LEARNING

- DATA IS CHEAP + ABUNDANT
 - KNOWLEDGE IS EXPENSIVE + SCARCE
 - LEARNING GENERAL MODELS FROM A DATA OF PARTICULAR EXAMPLES.
 - BUILD MODEL THAT IS GOOD AND USEFUL APPROXIMATION TO THE DATA.
- DATA MINING - APPLICATION OF MACHINE LEARNING

MACHINE LEARNING AUG, 28

DEEP LEARNING - IAN GOODFELLOW
PATTERN RECOGNITION AND MACHINE LEARNING

MACHINE LEARNING - ARTHUR SAMUEL (1959)

FIELD OF STUDY THAT GIVES COMPUTERS THE ABILITY TO LEARN WITHOUT BEING EXPLICITLY PROGRAMMED.

TOM MITCHELL (1998) WELL-POSED LEARNING PROBLEM: A COMPUTER PROGRAM IS SAID TO LEARN FROM EXPERIENCE E WITH RESPECT TO SOME TASK T AND SOME PERFORMANCE MEASURE P , IF IT IS PERFORMANCE ON T AS MEASURED BY P , IMPROVE WITH EXPERIENCE E .

SUPERVISED LEARNING

- MOST POPULAR
- LOTS OF DATA WITH THE ANSWER.

UNSUPERVISED LEARNING

- DATA BUT NO RIGHT ANSWER IS PROVIDED