# 18 SEP 2015     SEM. WEB MINING

## REAL VAL MEAS TO DIS. FEATURE

| REAL MEAS | FEAT |
|---|---|
| $m(x)$ | YES IF $m(x) > 50$ <br> NO OW |
|  | YES IF $m(x) < 20$ <br> NO OW |
| $m(x)$ | $m(x) \in [0, 10]$ <br> $m(x) \in [10, 50]$ <br> $\vdots$ |

# COSINE SIM.



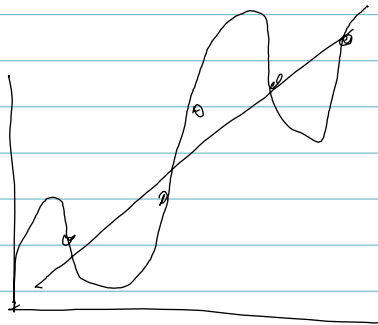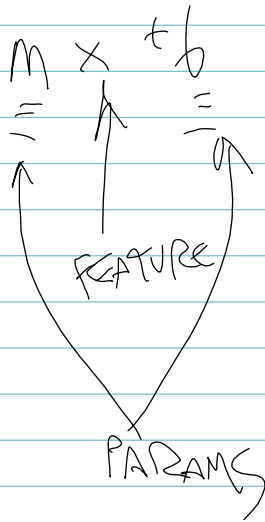$\cos\theta$ IS "SIMILARITY"
MEAS BETWN $S_a$ & $S_b$

---

UNSUPERVISED

DESCRIBING A DATASET

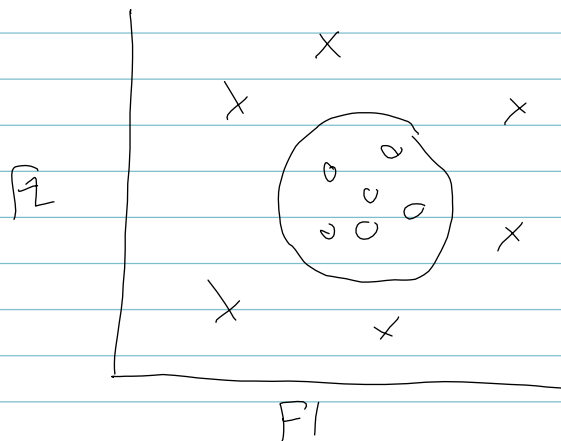SUPERVISED

TRAINING DATA — GIVES US AN
INTUITION
ON FUTURE DATA
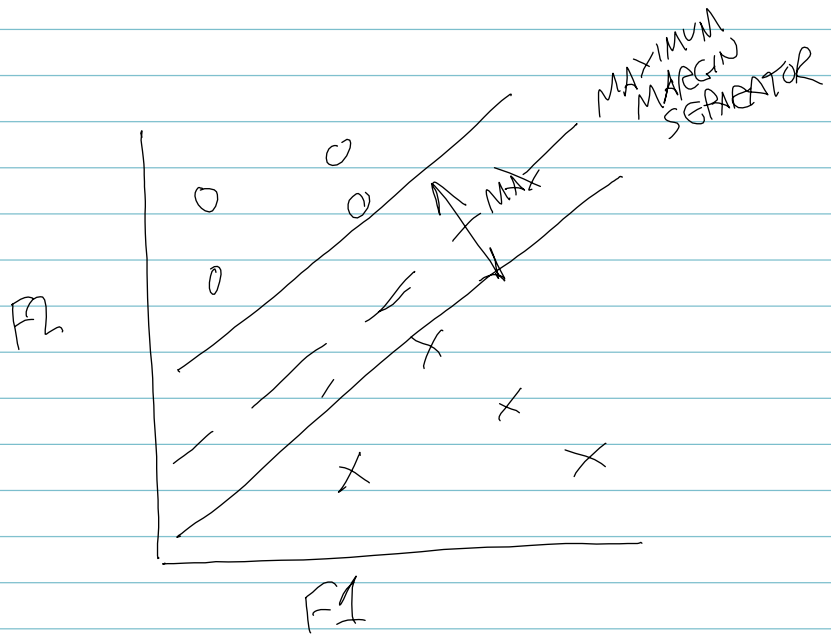
$$m \quad x \quad + b$$

$=$      $\uparrow$      $=$

FEATURE

PARAMS

## PARAMETRIC

NUMBER OF PARAMETERS
IS FIXED (NOT RELATED
TO x)

## NON-PARAMETRIC

NUMBER OF PARAMETERS
IS REL'D TO TNG DATA

# SVM

MAXIMUM MARGIN SEPARATOR
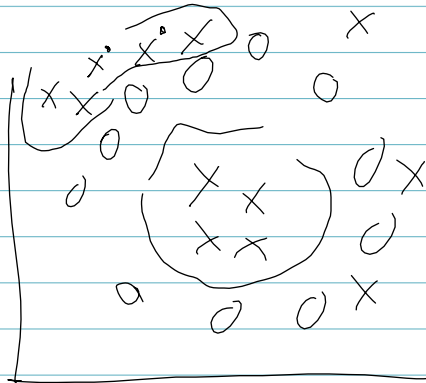
MAX

F2

F1

F2

F1

# DECISION STUMP

F

PICKE BASED
ON IG
FOR A
SAMPLE OF
THE TNG
DATA.

## DATA COMPLEXITY / IMBALANCE

$$P = \frac{\text{CORR CL. SAMPLES}}{\text{SAMPLE CLASS BY LEARNER}}$$

$$R = \frac{\text{CORR CL SAMPLES}}{\text{ALL SAMPLES OF CLASS}}$$

---

DATA SET $D$

LEAVE-ONE-OUT CROSS VAL.

$n$ SAMPLES IN $D$

EACH TEST IS DONE LEARNED

W. $D - \{x\}$

AND TESTED ON $X$

# 10-FOLD CROSS VAL

90% TNG
10% TEST

REPEAT 10x

REPORT AVG METRICS
(PREC, RECALL, AUC)
FOR EA. TRIAL

# VALIDATION EXPERIMENT

NEW DATA THAT
WASNT COLL W. INIT
DATASET