# Decision Tree

Machine learning 2021
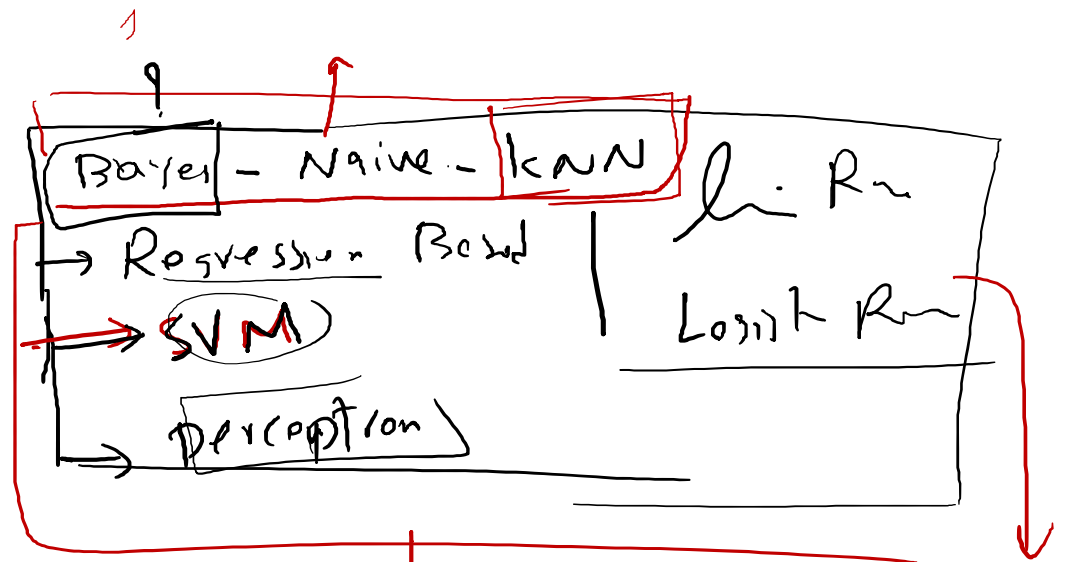
Mansoor Rezghi

# Ref:

Zaki

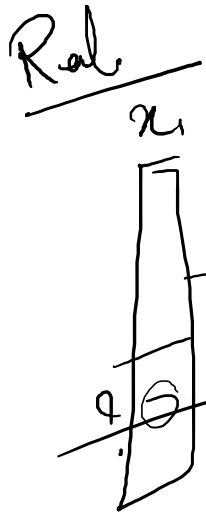Data mining Tan

Supervised
→ Regression
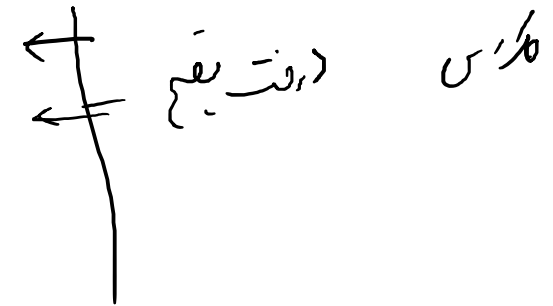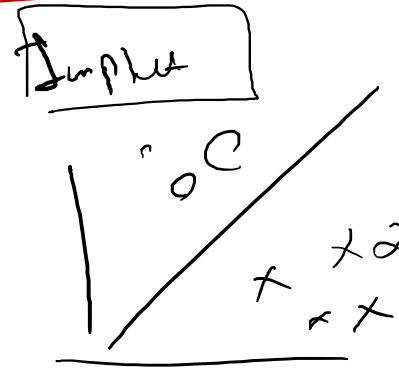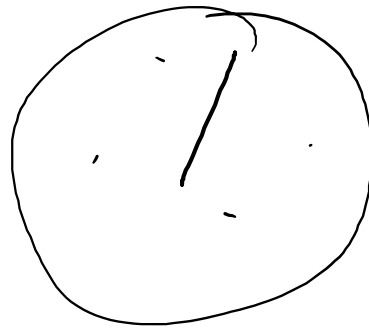→ Classification

Bayes — Naive — kNN

→ Regression Based

$SVM$

Perception

lin. Reg

Logist Reg

$$y = f(x)$$
$$w^T x + w_0$$

Explicite

Implicit

kNN → branch    lin

Real

$x_1$    $x_2$

$9.0$

Di

$x^2$
$x x$

$°C$

# Decision Tree classifiers

- A **root node** that has no incoming edges and zero or more outgoing edges.

- **Internal nodes**, each of which has exactly one incoming edge and two or more outg...

- **Leaf or term...** and no outgo...

Questions

Answer

Implicit discrimination

Body Temperature

Root node

Internal node

Warm

Cold

Gives Birth

Non-mammals

Yes

No

Mammals

Non-mammals

Leaf nodes

**Figure 4.4.** A decision tree for the mammal classification problem.

# Example of a Decision Tree

S  M  D

Depth.          Underfitting

$\leq 14$

No      Yes

| 4. | 3 |

| 5. | 2 |

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

categorical   categorical   continuous   class

Train

Training Data

*Splitting Attributes*

10

Refund

Yes        No

NO

MarSt

Single, Divorced        Married

TaxInc          NO

< 80K        > 80K

NO        YES

Model:  Decision Tree

# Another Example of Decision Tree



| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

There could be more than one tree that fits the same data!

# Decision Tree Classification Task

| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|-----|---------|---------|---------|-------|
| 1 | Yes | Large | 125K | No |
| 2 | No | Medium | 100K | No |
| 3 | No | Small | 70K | No |
| 4 | Yes | Medium | 120K | No |
| 5 | No | Large | 95K | Yes |
| 6 | No | Medium | 60K | No |
| 7 | Yes | Large | 220K | No |
| 8 | No | Small | 85K | Yes |
| 9 | No | Medium | 75K | No |
| 10 | No | Small | 90K | Yes |

Training Set

Tree Induction algorithm

Induction

Learn Model

Model

Decision Tree

| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|-----|---------|---------|---------|-------|
| 11 | No | Small | 55K | ? |
| 12 | Yes | Medium | 80K | ? |
| 13 | Yes | Large | 110K | ? |
| 14 | No | Small | 95K | ? |
| 15 | No | Large | 67K | ? |

Test Set

Apply Model

Deduction

# Apply Model to Test Data

Start from the root of tree.

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

# Apply Model to Test Data

Test Data

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

# Apply Model to Test Data

Test Data

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

# Apply Model to Test Data

Test Data

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

# Apply Model to Test Data

Test Data

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

# Apply Model to Test Data

**Test Data**

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ?  NO |



Assign Cheat to "No"

Refund
- Yes → NO
- No → MarSt
  - Single, Divorced → TaxInc
    - < 80K → NO
    - > 80K → YES
  - Married → NO

Imphte

# Decision Tree Classification Task

| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|-----|---------|---------|---------|-------|
| 1 | Yes | Large | 125K | No |
| 2 | No | Medium | 100K | No |
| 3 | No | Small | 70K | No |
| 4 | Yes | Medium | 120K | No |
| 5 | No | Large | 95K | Yes |
| 6 | No | Medium | 60K | No |
| 7 | Yes | Large | 220K | No |
| 8 | No | Small | 85K | Yes |
| 9 | No | Medium | 75K | No |
| 10 | No | Small | 90K | Yes |

Training Set

Tree Induction algorithm

Induction

Learn Model

Model

Decision Tree

Apply Model

Deduction

| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|-----|---------|---------|---------|-------|
| 11 | No | Small | 55K | ? |
| 12 | Yes | Medium | 80K | ? |
| 13 | Yes | Large | 110K | ? |
| 14 | No | Small | 95K | ? |
| 15 | No | Large | 67K | ? |

Test Set

# Decision Tree Induction

Many Algorithms:

Hunt's Algorithm (one of the earliest)
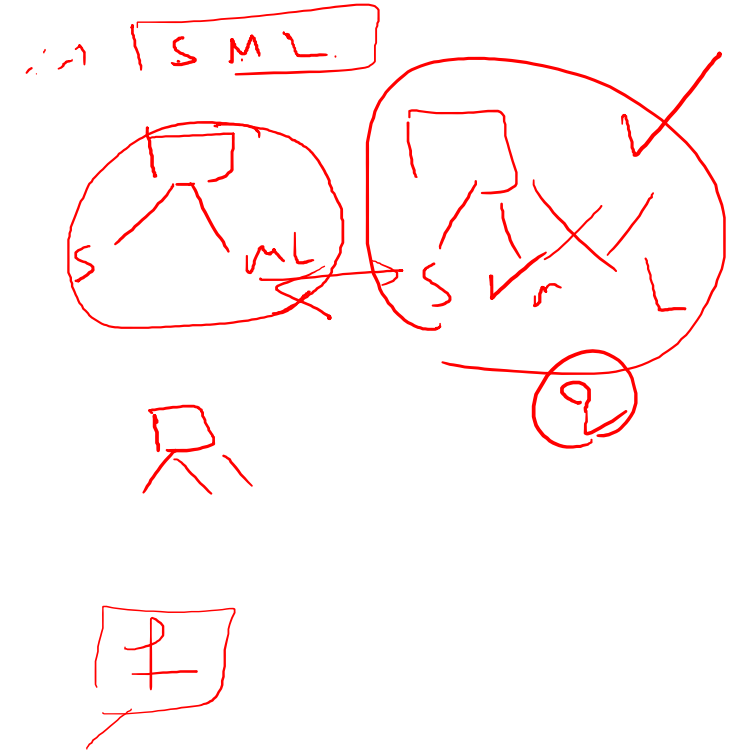
CART

ID3, C4.5

SLIQ, SPRINT

# Design Issues of Decision Tree Induction

## Design Issues of Decision Tree Induction

- Order of selecting the attributes (Evaluation )
- Designing the  Questions (Evaluation)

## How should the splitting procedure stop?

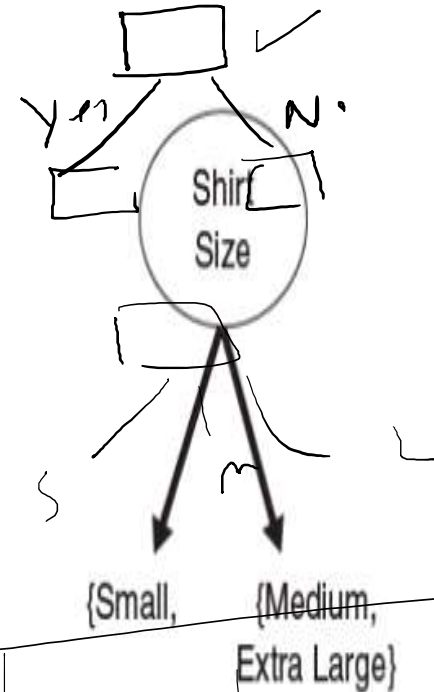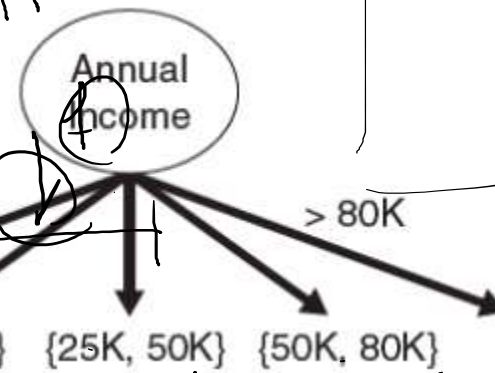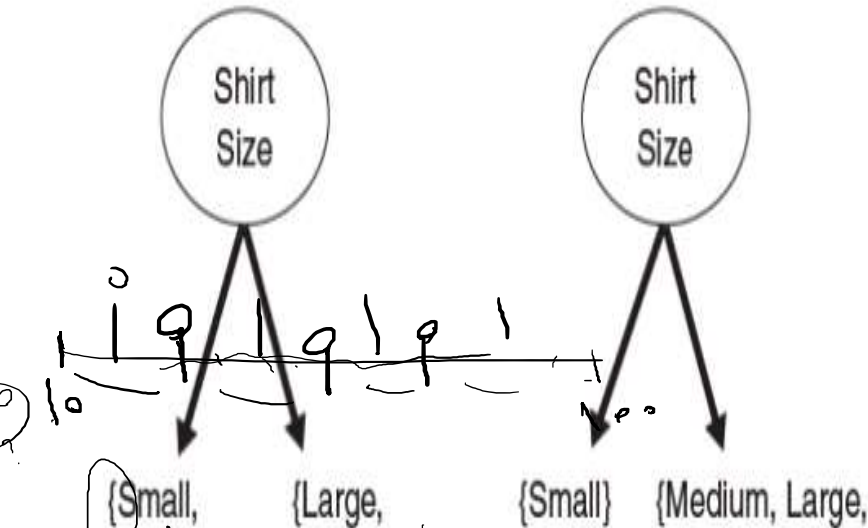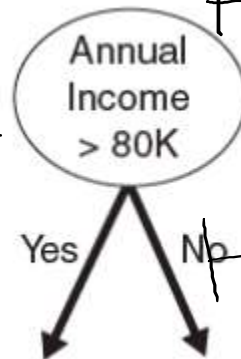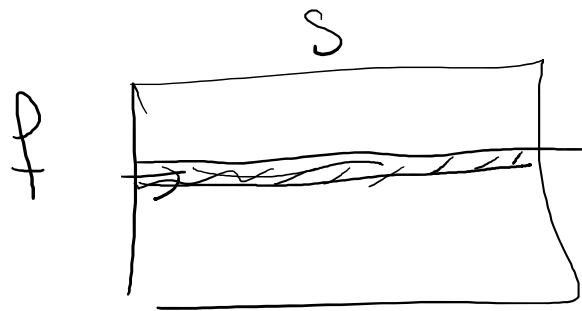- Under and over fitting!
- Evaluation is need

# Questions?

- Binary
- Nominal
- Ordinal
- Continuous

$\ne$ $\longrightarrow$ k $\ne$ $2^k$

k-1

k-2

k-3

2

③

C.⩽15 $\longrightarrow$ 2

○

①

V

100

$C_1$

$\vdots$

$C_k$

SM

$N_0$

$X_0$

V

V

pure $\longrightarrow$ Net $\longrightarrow$ Criten

S ... ML

V

V

Compare.

V

V

V

V

Node purity Evaluation

Tree Level Evaluation
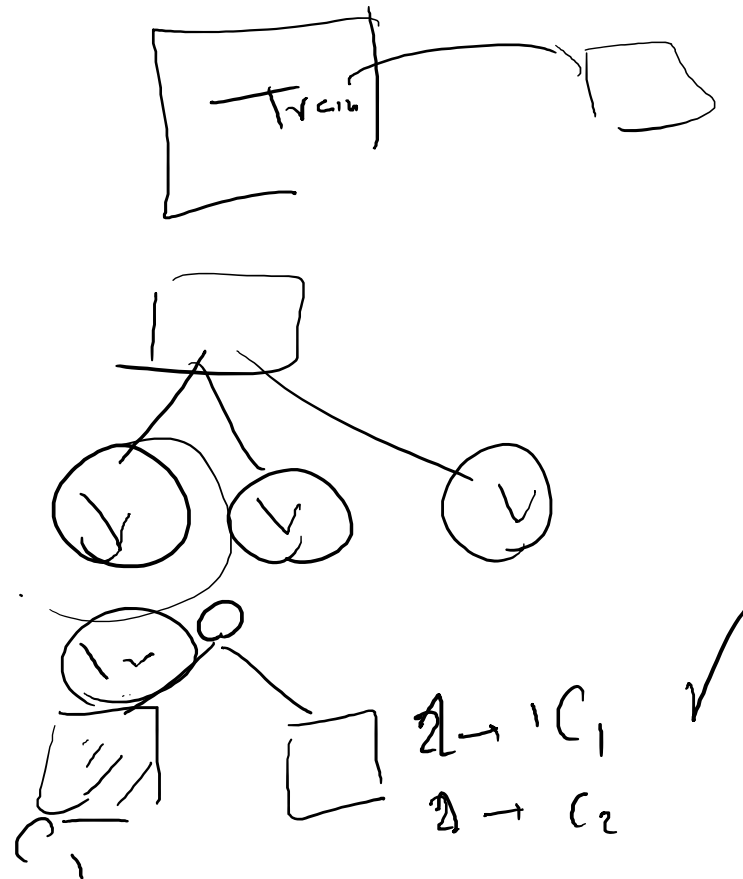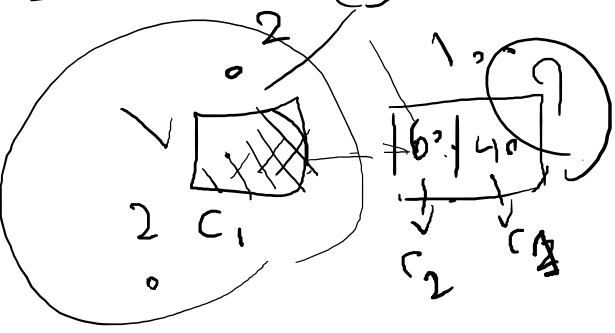
$t_1$

$t_2$

$$-\left[\left(\frac{1}{2}\log\frac{1}{2} + \frac{1}{2}\log\frac{1}{2}\right)\right]$$

$$= -\log\frac{1}{2} = \boxed{\log 2}$$

$$\boxed{\log K}$$

$$100 = N_t$$

$$N_t^1 = 100 \qquad P(1|t) = 1$$

$$N_t^2 = 0 \qquad P(2|t) = 0$$

$$-\sum P(i|t)\log P(i|t) = -(0\times\log 0 + 1\times\log 1)$$

$$= 0$$

Diagram labels:
$C_1$ $C_2$
$50$ $50$

$C_1$ $C_2$

$V$ $200$

$100 - 10$

$50$ $50$

Root

$N_1$ $N_2$ $N_3$

Internal

Leaf

$N$

$V_1$ $V_2$

$= I(v) - \sum I(v_i) \times \dfrac{N_i}{N}$

S

$f_1$
$f_2$

Evaluation of features. ✓

Evaluate Splits ✓

Evaluate Node ✓

$f_i$: Binary

Yes 0 1 No

Gini, Entropy & Error $e e_i$

$I(V_i)$

Categorical

$x_2$

$3$ | $c_1$ | $c_2$
$2$ | $r_2$ | $c_1$

$2$   $4$   $x_1$

if $x_1 < 4$

if $x_2 > 2$

if $x_1 > 2$

if $x_2 < 3$

$c_1$

## Measures for Selecting the Best Split

$o \leq \rightarrow \log k$

$\exists \ell \quad P(\ell|t) > P(i|t) \quad i \neq \ell$

$N_t = $ تعداد عناصر درنود t ام مستند

$t \ \square$

تعداد عناصر از نود t

$t$ ام که از کلاس i هستند $N_t^i$

$P(i|t) = \dfrac{N_t^i}{N_t}$

$-C_1 N_t^1$

$P(i|t) =$

$P(k|t) \quad C_{1k} \quad N_{tk}^k$

## Greedy approaches:

$$\text{Entropy}(t) = \boxed{-\sum_{i=0}^{c-1} p(i|t) \log_2 p(i|t),}$$

$$\text{Gini}(t) = 1 - \sum_{i=0}^{c-1} [p(i|t)]^2,$$

$$\text{Classification error}(t) = 1 - \max_i [p(i|t)],$$

Let $p(i|t)$ denote the fraction of records belonging to class $i$ at a given node $t$.

$N$

$N$

$\times \alpha$

$N_1 \quad N_2 \quad N_3$

Gain

$$\Delta = I(\text{parent}) - \sum_{j=1}^{k} \frac{N(v_j)}{N} I(v_j),$$

$$\text{Gain ratio} = \frac{\Delta_{\text{info}}}{\text{Split Info}}$$

$$\text{Split Info} = -\sum_{i=1}^{k} P(v_i) \log_2 P(v_i)$$

$P(v_i) = \dfrac{N_i}{N}$

| Node $N_1$ | Count |
|---|---|
| Class=0 | 0 |
| Class=1 | 6 |

$\text{Gini} = 1 - (0/6)^2 - (6/6)^2 = 0$
$\text{Entropy} = -(0/6)\log_2(0/6) - (6/6)\log_2(6/6) = 0$
$\text{Error} = 1 - \max[0/6, 6/6] = 0$

| Node $N_2$ | Count |
|---|---|
| Class=0 | 1 |
| Class=1 | 5 |

$\text{Gini} = 1 - (1/6)^2 - (5/6)^2 = 0.278$
$\text{Entropy} = -(1/6)\log_2(1/6) - (5/6)\log_2(5/6) = 0.650$
$\text{Error} = 1 - \max[1/6, 5/6] = 0.167$

| Node $N_3$ | Count |
|---|---|
| Class=0 | 3 |
| Class=1 | 3 |

$\text{Gini} = 1 - (3/6)^2 - (3/6)^2 = 0.5$
$\text{Entropy} = -(3/6)\log_2(3/6) - (3/6)\log_2(3/6) = 1$
$\text{Error} = 1 - \max[3/6, 3/6] = 0.5$

**Algorithm 4.1** A skeleton decision tree induction algorithm.

TreeGrowth (E, F)
 1: **if** stopping_cond(E,F) = *true* **then**
 2:    *leaf* = createNode().
 3:    *leaf.label* = Classify(E).
 4:    return *leaf*.
 5: **else**
 6:    *root* = createNode().
 7:    *root.test_cond* = find_best_split(E, F).
 8:    let $V = \{v|v$ is a possible outcome of *root.test_cond* $\}$.
 9:    **for each** $v \in V$ **do**
10:       $E_v = \{e \mid root.test\_cond(e) = v$ and $e \in E\}$.
11:       *child* = TreeGrowth($E_v$, F).
12:       add *child* as descendent of *root* and label the edge (*root* → *child*) as *v*.
13:    **end for**
14: **end if**
15: return *root*.

# Characteristics of Decision Tree Induction

1-

Decision tree induction is a nonparametric approach :

Finding an optimal decision tree is an NP-complete problem.

Since most decision tree algorithms employ a top-down, recursive partitioning approach, the number of records becomes smaller as we traverse down the tree. At the leaf nodes, the number of records may be too small to make a statistically significant decision about the class representation of the nodes. This is known as the **data fragmentation** problem. One possible solution is to disallow further splitting when the number of records falls below a certain threshold.

Tree replication problem.

# Linear or nonlinear



An **oblique decision tree** can be used to overcome this limitation because it allows test conditions that involve more than one attribute. The data set given in Figure 4.21 can be easily represented by an oblique decision tree containing a single node with test condition

$$x + y < 1.$$

Although such techniques are more expressive and can produce more compact trees, finding the optimal test condition for a given node can be computationally expensive.

Solution

# Overfitting

Training set

$x_1$

$x_2$

Training Error
Test Error

Error Rate

Number of Nodes

Depth

$f_1$
$f_2$
$f_2$

overfit

- Presence of noise
- Lake of representative samples

Decisio Tree

Train

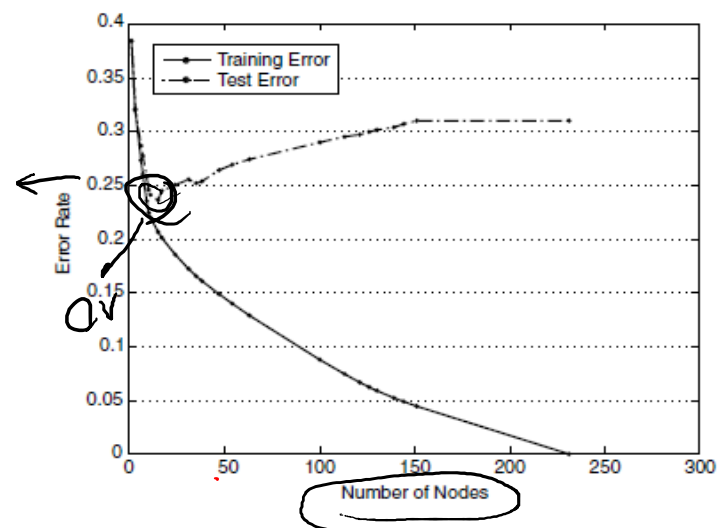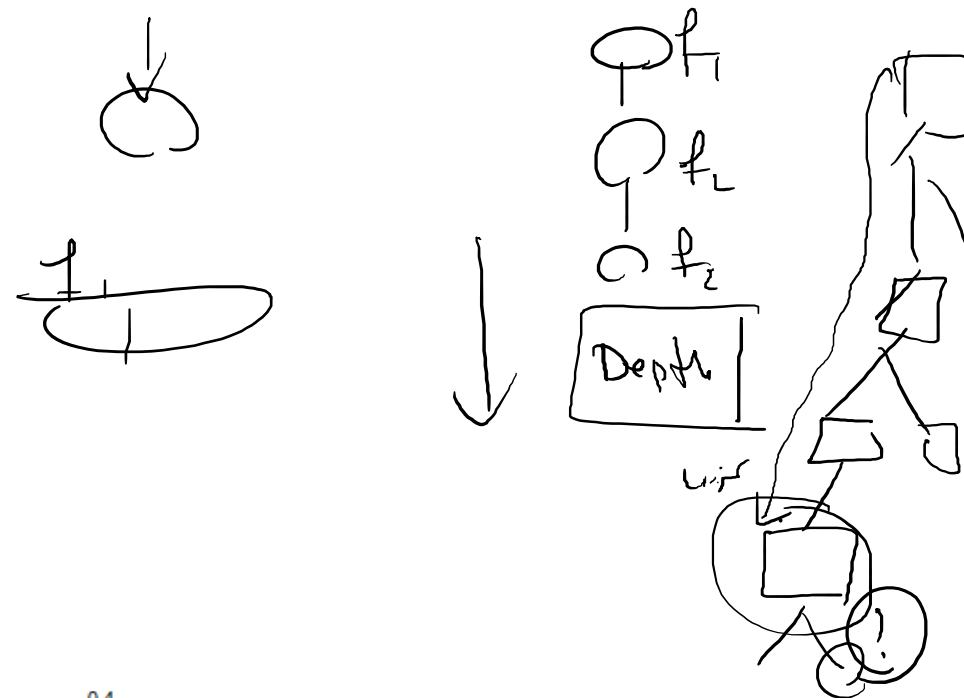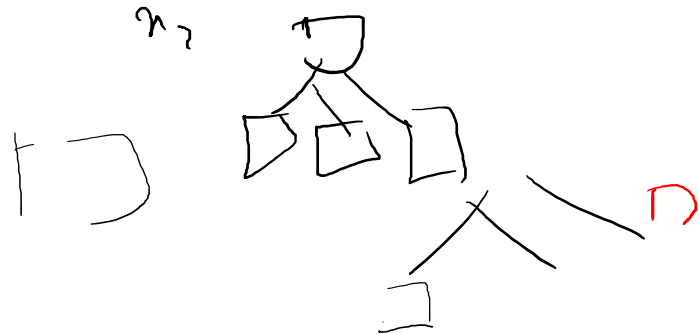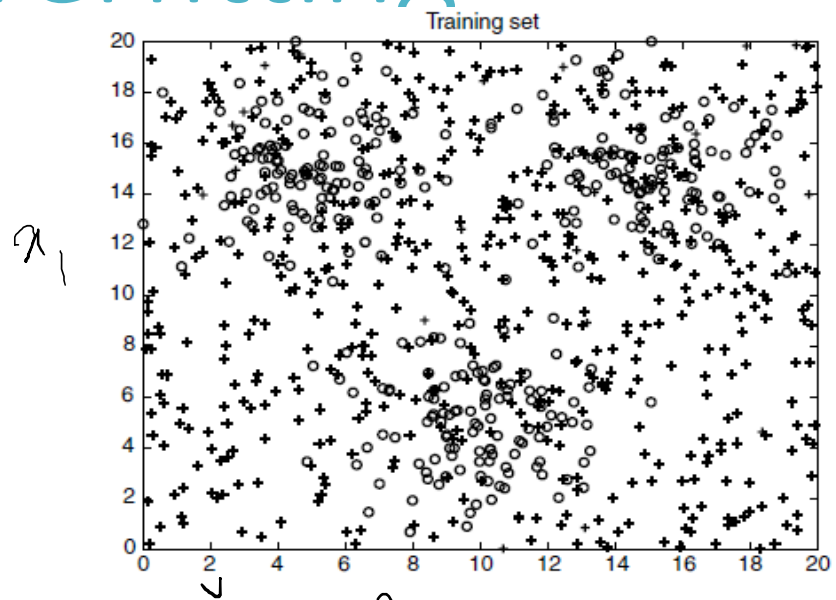| Name | Body Temperature | Gives Birth | Four-legged | Hibernates | Class Label |
|---|---|---|---|---|---|
| porcupine | warm-blooded | yes | yes | yes | yes |
| cat | warm-blooded | yes | yes | no | yes |
| bat | warm-blooded | yes | no | yes | no* |
| whale | warm-blooded | yes | no | no | no* |
| salamander | cold-blooded | no | yes | yes | no |
| komodo dragon | cold-blooded | no | yes | no | no |
| python | cold-blooded | no | no | yes | no |
| salmon | cold-blooded | no | no | no | no |
| eagle | warm-blooded | no | no | no | no |
| guppy | cold-blooded | yes | no | no | no |

Test

| Name | Body Temperature | Gives Birth | Four-legged | Hibernates | Class Label |
|---|---|---|---|---|---|
| human | warm-blooded | yes | no | no | yes |

Train

| Name | Body Temperature | Gives Birth | Four-legged | Hibernates | Class Label |
|------|------------------|-------------|-------------|------------|-------------|
| salamander | cold-blooded | no | yes | yes | no |
| guppy | cold-blooded | yes | no | no | no |
| eagle | warm-blooded | no | no | no | no |
| poorwill | warm-blooded | no | no | yes | no |
| platypus | warm-blooded | no | yes | yes | yes |

# Overfitting Due to Preser

### Train

| Name | Body Temperature | Gives Birth | Four-legged | Hibernates | Class Label |
|---|---|---|---|---|---|
| porcupine | warm-blooded | yes | yes | yes | yes |
| cat | warm-blooded | yes | yes | no | yes |
| bat | warm-blooded | yes | no | yes | no* |
| whale | warm-blooded | yes | no | no | no* |
| salamander | cold-blooded | no | yes | yes | no |
| komodo dragon | cold-blooded | no | yes | no | no |
| python | cold-blooded | no | no | yes | no |
| salmon | cold-blooded | no | no | no | no |
| eagle | warm-blooded | no | no | no | no |
| guppy | cold-blooded | yes | no | no | no |

Train error: 0



30%

(a) Model M1

### Test

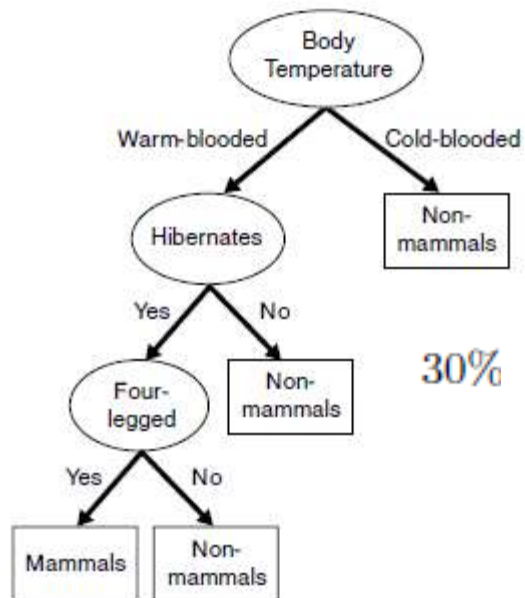| Name | Body Temperature | Gives Birth | Four-legged | Hibernates | Class Label |
|---|---|---|---|---|---|
| human | warm-blooded | yes | no | no | yes |
| pigeon | warm-blooded | no | no | no | no |
| elephant | warm-blooded | yes | yes | no | yes |
| leopard shark | cold-blooded | yes | no | no | no |
| turtle | cold-blooded | no | yes | no | no |
| penguin | cold-blooded | no | no | no | no |
| eel | cold-blooded | no | no | no | no |
| dolphin | warm-blooded | yes | no | no | yes |
| spiny anteater | warm-blooded | no | yes | yes | yes |
| gila monster | cold-blooded | no | yes | yes | no |

Train error: 20%



error rate (10%)

# small number of training records

| Name | Body Temperature | Gives Birth | Four-legged | Hibernates | Class Label |
|------|------------------|-------------|-------------|------------|-------------|
| salamander | cold-blooded | no | yes | yes | no |
| guppy | cold-blooded | yes | no | no | no |
| eagle | warm-blooded | no | no | no | no |
| poorwill | warm-blooded | no | no | yes | no |
| platypus | warm-blooded | no | yes | yes | yes |



30%

lin Regression

logisten

SVM | Bayes | Naive Bayes | 1cNN , Perceptron ( _Varial_ )  Variants  → Decision Trees

$P(x|C_i)P(C_i)$

→ Evaluation

→ combine :   Ensemble method

$x_1$  $x_c$

$\beta$

$z_{ki}$

→ 0

( — )

Root N

$N_1$  $N_2$

→ Interkr

Leaf

S

$f_1$ p

$p_2$

V

$f_1$ vs $f_2$

Feature Evaluation

Split Evaluation

Gin
Entropy
Err rate.

Node Evaluation

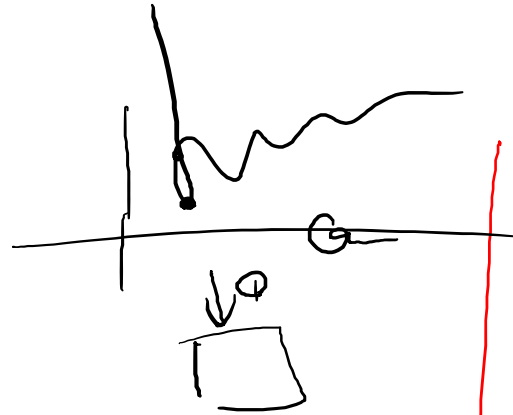$$Gain = f(v) - \sum \frac{N_i}{N} (v_i)$$

Split

# Decision Trees

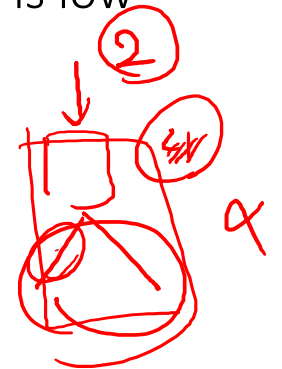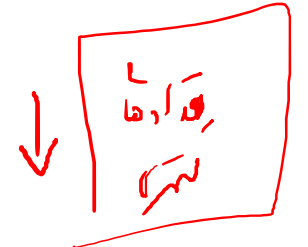**Prepruning :Early Stopping Rule**

a more restrictive stopping condition
stop expanding a leaf node when the observed gain in impurity measure is low

**Post-pruning**

decision tree is initially grown to its maximum size
tree-pruning step
replacing a subtree with a new leaf node

# Evaluating the Performance of a Classifier

accuracy or error rate computed from the test set can used to compare different classifiers

class labels of test records must be known

**Holdout Method**

1. labeled examples partitioned into two disjoint sets: training and the test sets
2. classification model is then induced from the training set
3. its performance is evaluated on the test set

✓ smaller training set size, larger variance of the model
✓ training set is too large, then the estimated accuracy computed from the
smaller test set is less reliable

# Evaluating the Performance of a Classifier
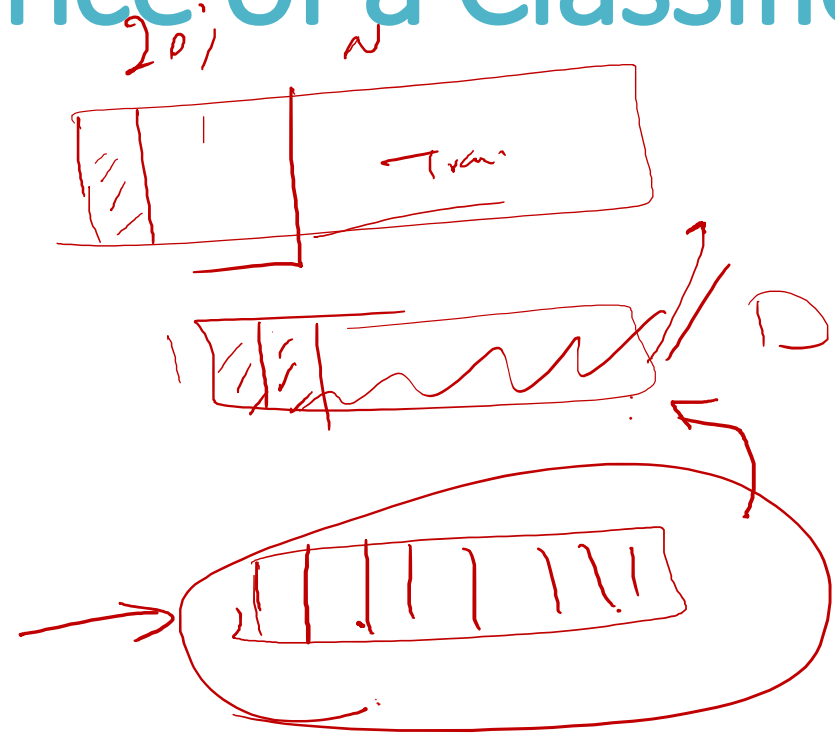
**Random Subsampling**

Repeated holdout

**Bootstrap**

Sampling with replacement

**Cross-Validation**

each record is used the same number of times for training and exactly once
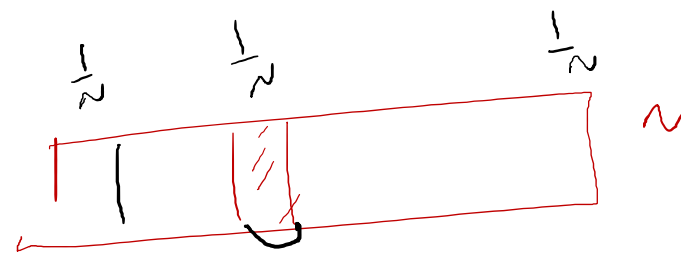for testing
K-fold Cross-Validation

# sampling

$\rightarrow$

Cross-Validation

Bootstrap

$$1 - (1 - 1/N)^N. \implies 1 - e^{-1} = 0.632.$$

Leave one out

$$\frac{1}{N} \qquad \frac{1}{N} \qquad \frac{1}{N}$$

$N$

$$1 - \left(1 - \frac{1}{N}\right)^N$$

$$\left(1 - \frac{1}{N}\right)\left(1 - \frac{1}{N}\right) \cdots \left(1 - \frac{1}{N}\right) =$$

$50 \qquad 50$

$$\frac{1}{2} = \int$$

$\frac{2a}{|?| \; |?|}$

# Metrics for class imbalance problem

# Imbalance

✓ Data sets with imbalanced class distributions
✓ in credit card fraud detection, fraudulent transactions are outnumbered by legitimate transactions
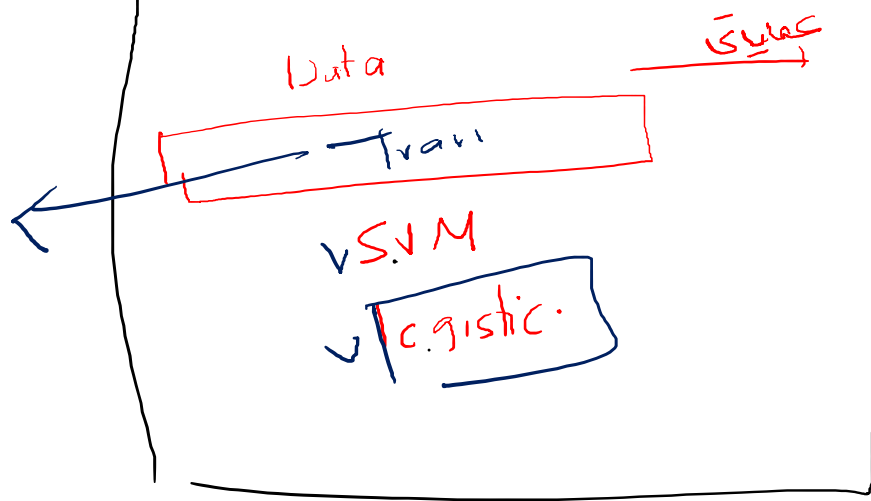✓ accuracy measure, used extensively for classifiers, may not be well suited for evaluating models derived from imbalanced data sets

example : 1% of the credit card transactions fraudulent,
a model that predicts every transaction as legitimate
accuracy 99%
it fails to detect any of the fraudulent activities.
binary classification, the rare class is often denoted as the positive class against negative class

|        |   | Predicted Class |              |
|--------|---|-----------------|--------------|
|        |   | +               | −            |
| Actual | + | $f_{++}$ (TP)   | $f_{+-}$ (FN) |
| Class  | − | $f_{-+}$ (FP)   | $f_{--}$ (TN) |

confusion matrix

# Imbalance

Precision : fraction of records that actually turns out to be positive in the group the classifier has declared as a positive class

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

Recall measures the fraction of positive examples correctly predicted by the classifier

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

maximizes both precision and recall

# Imbalance

Precision and recall can be summarized into another metric known as the F1 measure

$$F_1 = \frac{2}{\frac{1}{r} + \frac{1}{p}}.$$

tends to be closer to the smaller of the two numbers
a high value of *F1*-measure ensures that both precision and recall are reasonably high

$$\text{Weighted accuracy} = \frac{w_1 TP + w_4 TN}{w_1 TP + w_2 FP + w_3 FN + w_4 TN}.$$