

DMLL, 22 Jan 2019

Frequent itemsets with multiple min support

Item  $i \rightarrow \text{min-sup}(i)$

Set  $\{i_1, i_2, \dots, i_k\}$  - min support is

$\min_{j \in \{1, 2, \dots, k\}} (\text{min-sup}(i_j))$

Destroys a priori property wrt smallest min-sup

$\{a, b, c\} \xrightarrow{\text{set has min-sup 0.01}}$

$0.01 \ 0.1 \ 0.2 \xrightarrow{\{b, c\} \text{ has min-sup 0.1}}$

Adopt level-by-level algorithm

$C_1 \rightarrow F_1 \rightarrow C_2 \rightarrow F_2 \rightarrow \dots \rightarrow F_{K-1} \rightarrow C_K \rightarrow F_K$

$$F_i = \{i_j \mid \text{sup}(i_j) \geq \text{min-sup}(i_j)\}$$



$$C_2 \neq F_1 \times F_1 \rightarrow \{(i_j, i_k) \mid \text{sup}(i_j) \geq \text{min-sup}(i_j), \text{sup}(i_k) \geq \text{min-sup}(i_k)\}$$

$C_2$  should have

$$\underline{(i_j, i_k)}, \text{sup}(i_j) \geq \text{min-sup}(i_j)$$

$$\text{ass. } \text{sup}(i_j) \leq \text{sup}(i_k) \quad \text{sup}(i_k) \geq \text{min-sup}(i_j)$$

## Convenhn

Enumerate any  $X \subseteq I$  in ascending order of min-sup.

$$\{i_1, i_2, \dots, i_k\}$$

$$\text{min-sup}(i_1) \leq \text{min-sup}(i_2) \leq \dots \leq \text{min-sup}(i_k)$$

Correct  $C_2$

$$C_2 = \{(i_j, i_k) \mid \text{sup}(i_j), \text{sup}(i_k) \geq \text{min-sup}(i_j)\}$$

Count & get  $F_2$  as usual

$F_{k-1} \rightarrow C_k$ 

Merge

 $\{i_1, i_2, \dots, i_{k-2}, i_{k-1}\}$  $\in F_{k-1}$  $\{i_1, i_2, \dots, i_{k-2}, i'_{k-1}\}$  $\{i_1, i_2, \dots, i_{k-2}, i_{k-1}, i'_{k-1}\} \in C_k$ 

↓ Prune using a-prion

 $\{i_2, \dots, i_{k-2}, i_{k-1}, i'_{k-1}\}$  Check all  $k-1$  subsets

has a higher mm-sup

in  $F_{k-1}$  [uniform case]

- not consider this  $k-1$  subset when pruning

## Generating rules

$$X = \{i_1, i_2, \dots, i_k\} \in F_k$$

$$X \setminus \alpha \rightarrow \alpha$$

$$\frac{\text{sup}(X)}{\text{sup}(X \setminus \alpha)} \geq \text{min-conf}$$

↳ What if  $\alpha = \{i_1\}$ ?

Modify level-by-level to keep track of this

If  $\{i_1, i_2, \dots, i_k\} \in C_k$

Count occurrences of  $\{i_1, i_2, \dots, i_k\}$

and

$\{i_2, \dots, i_k\}$

need not qualify for  $F_{k-1}$ ,

but available for rule gen.

Other variations - Sequential rules

Buy product today  
accessory later

## Special Case

Last attribute is a category — topic  
of doc

“Class Association Rules”

$$X \rightarrow \{c\}$$

---

General problem of supervised learning

Build a model to predict a category  
given some attributes

# Bank loans

Outcome

ID	Age	Has_job	Own_house	Credit_rating	Class
1	young	false	false	fair	No
2	young	false	false	good	No
3	young	true	false	good	Yes
4	young	true	true	fair	Yes
5	young	false	false	fair	No
6	middle	false	false	fair	No
7	middle	false	false	good	No
8	middle	true	true	good	Yes
9	middle	false	true	excellent	Yes
10	middle	false	true	excellent	Yes
11	old	false	true	excellent	Yes
12	old	false	true	good	Yes
13	old	true	false	good	Yes
14	old	true	false	excellent	Yes
15	old	false	false	fair	No

Training data - historical data with  
class labels

↓  
Model

"labelled data"

↑  
Supervised learning

Attributes of  
unknown  
item

→ Predict  
class

Implicit assumption

Training data is  
statistically representative of future data

How to evaluate the model?

By defn, we don't know correct answers for unseen data

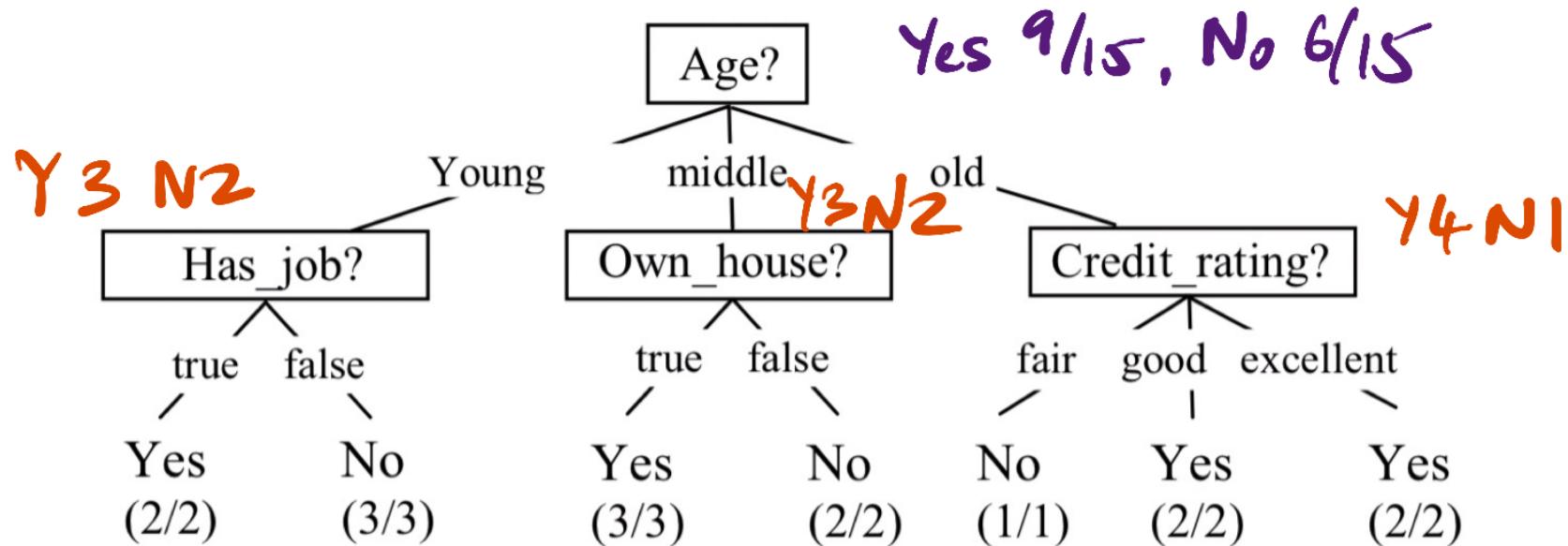
What does it mean "to predict"?

Data is not always "complete"

Attributes not captured

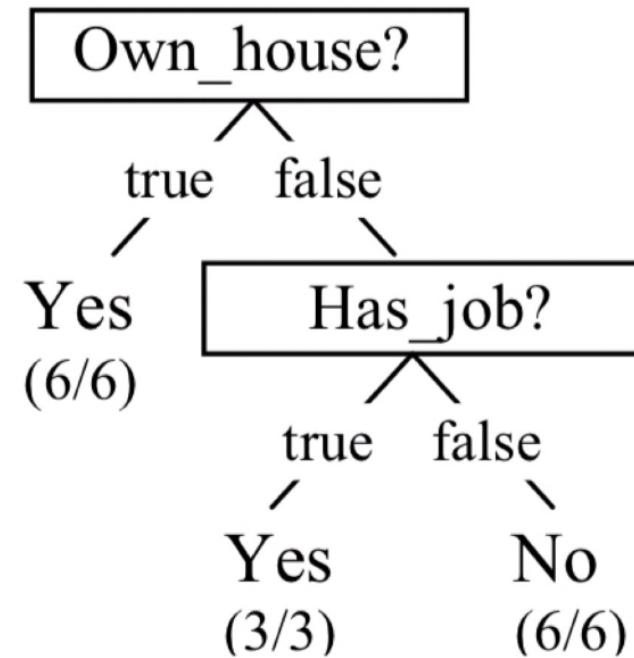
# Decision Tree

ID	Age	Has_job	Own_house	Credit_rating	Class
1	young	false	false	fair	No
2	young	false	false	good	No
3	young	true	false	good	Yes
4	young	true	true	fair	Yes
5	young	false	false	fair	No
6	middle	false	false	fair	No
7	middle	false	false	good	No
8	middle	true	true	good	Yes
9	middle	false	true	excellent	Yes
10	middle	false	true	excellent	Yes
11	old	false	true	excellent	Yes
12	old	false	true	good	Yes
13	old	true	false	good	Yes
14	old	true	false	excellent	Yes
15	old	false	false	fair	No



## 1. How to decide what question to ask next?

Answer is not unique.

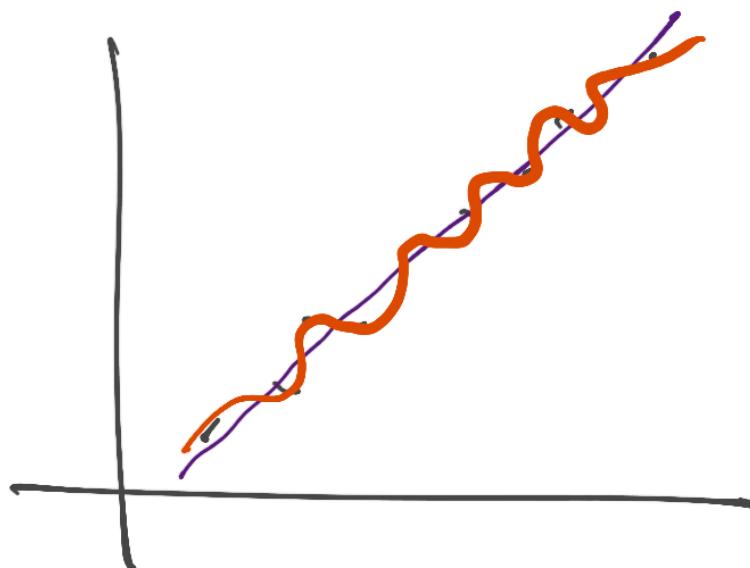


“Smaller” tree

2. Which tree are we looking for?

Small trees are better

- └ Simpler explanation
- └ Shorter trees generalize the training data better



By any defn of "smallest", finding the smallest tree is computationally hard

NP-complete

"Greedy" strategy - heuristic to get as good a tree as we can

Goal  $\rightarrow$  Narrow down candidates to all Yes or No  $\xrightarrow{\text{Assume}}$  binary classification

Go from "impure" mix to "pure" set

Purity as a measure to make a locally  
optimal choice — which question to ask

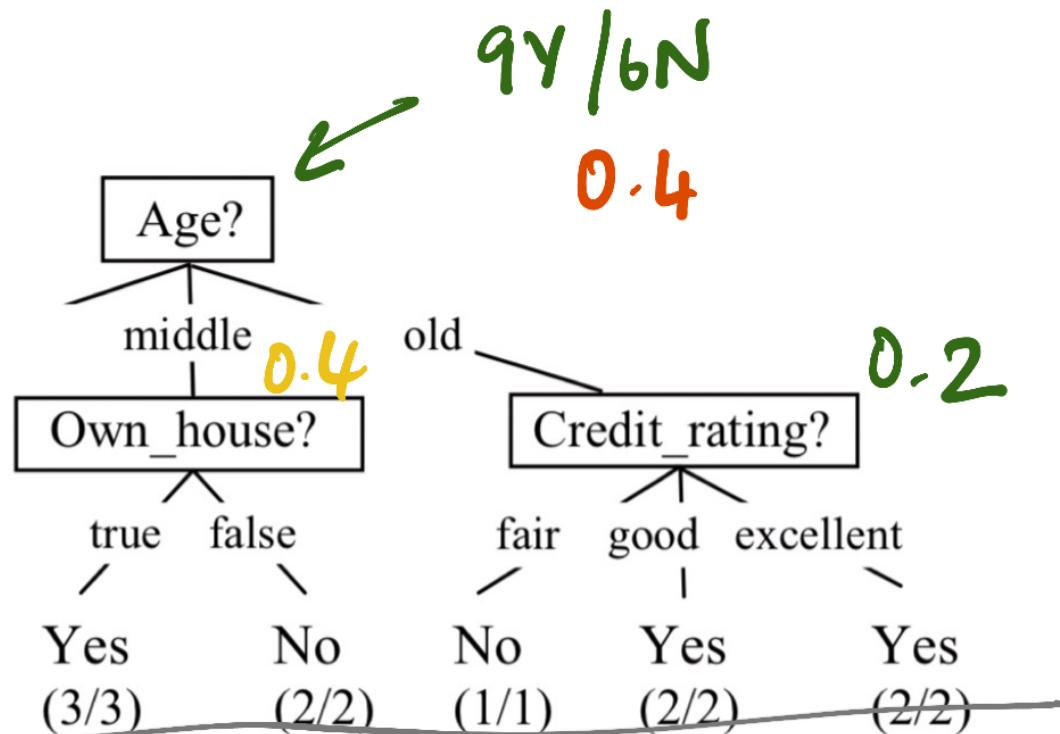
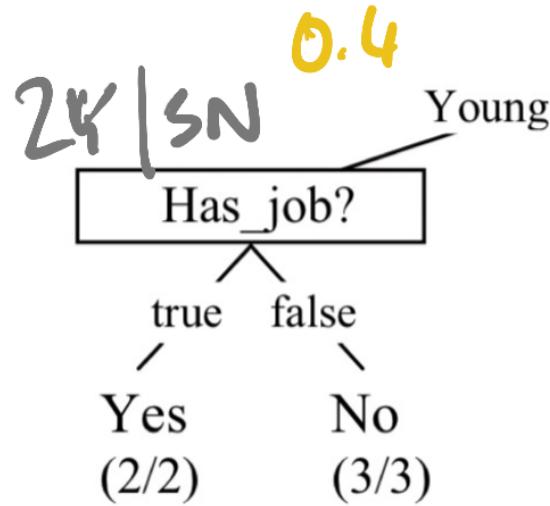
Suppose we want  $\delta$  questions to ask

End up  
with



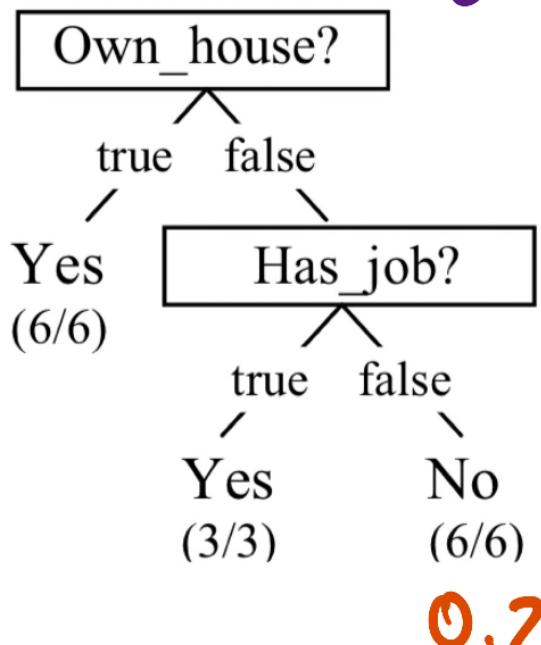
Predict  
majority  
value = Y

→ Error rate is 37%



Weighted avg

$$\begin{aligned}
 &= 5 \cdot (0.4) \\
 &+ 5(0.2) \\
 &+ 5(0.2) \\
 &= \frac{0.33}{15}
 \end{aligned}$$



$$\begin{aligned}
 &6 \cdot 0 + 1 \cdot 0.33 \\
 &= \frac{0.2}{15}
 \end{aligned}$$

## Algorithm

Among the attributes still to be explored,  
choose the one that improves purity  
by largest amount

## Stopping criterion

- Impurity = 0
- Run out of questions — use majority value as prediction

A "better" notion of impurity

Based on empirical observations