

テキストのラベルづけ/レベルづけ  
のための特徴抽出の手法について  
(Feature Selection  
in Text Classification/Ranking)

奥村 学  
(東京工業大学)

# Background (1/2)

- Measuring text readability

(小学(1年-6年), 中学(1年-3年), 高校(1年-3年), 大学)

- Judging text level

(A1, A2, B1, B2, C1, C2 in CEFR)

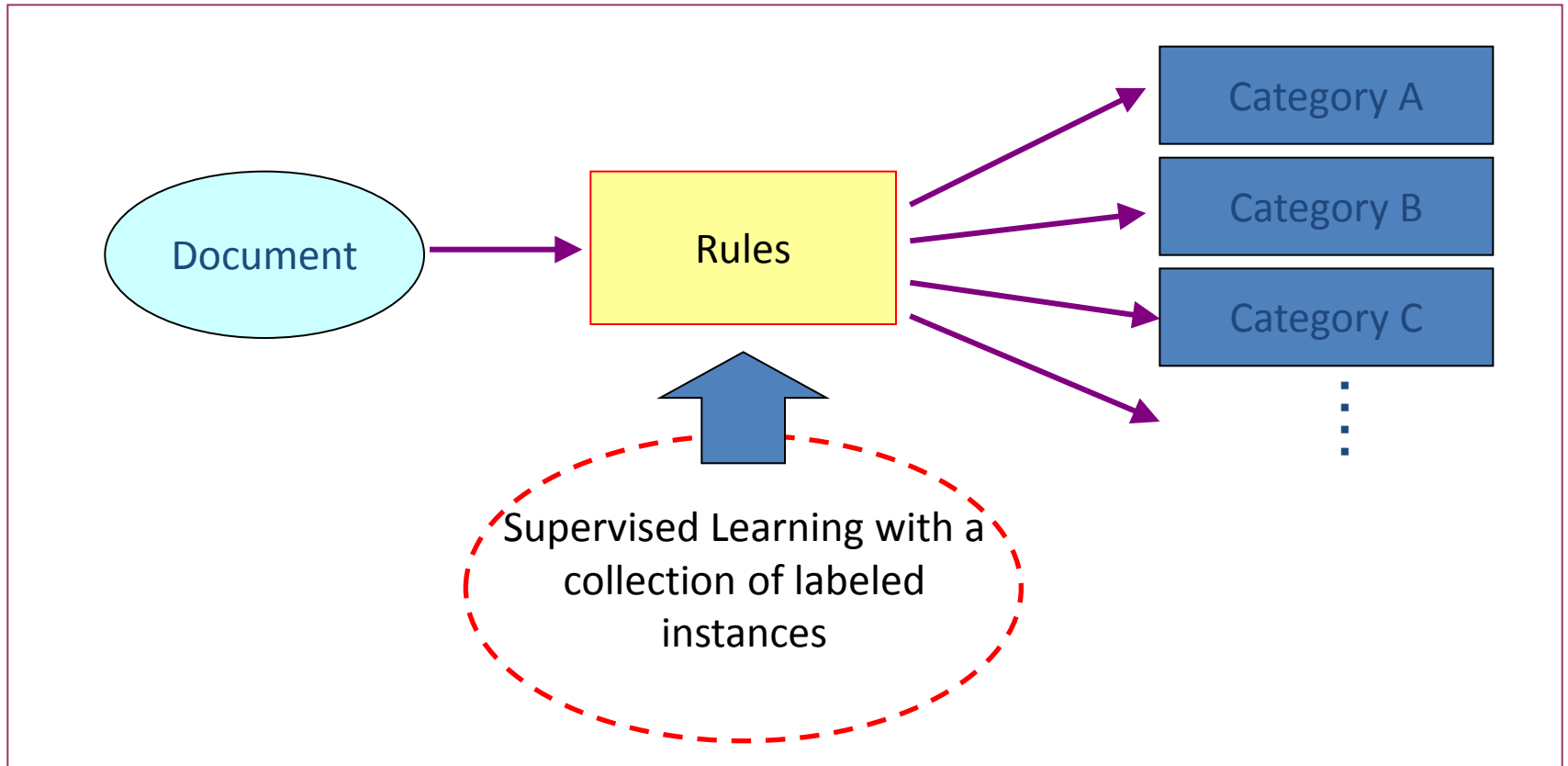
These tasks are considered as text classification/categorization or ranking in NLP (Natural Language Processing).

# Background (2/2)

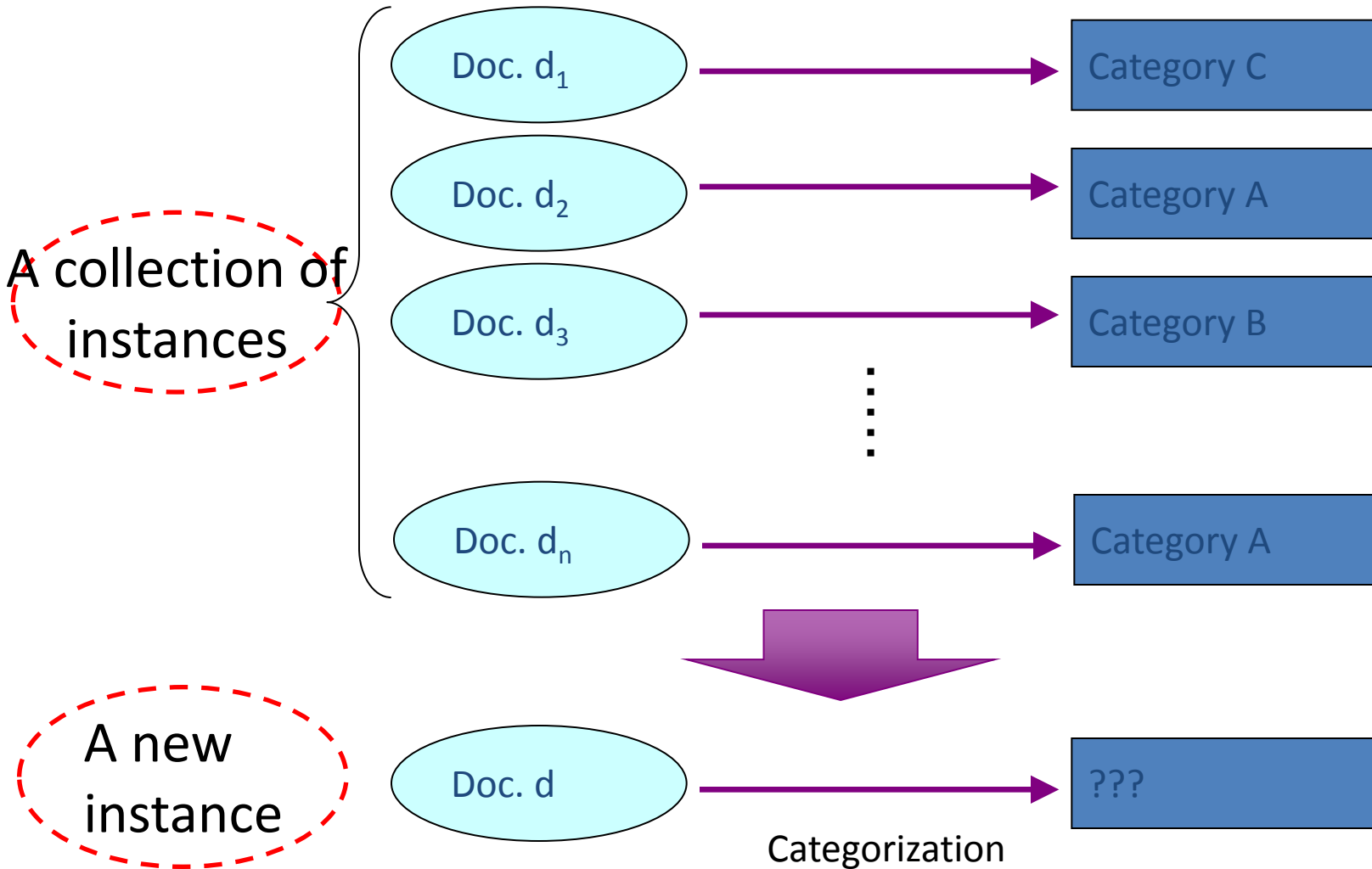
- What are good clues for measuring text readability/judging text level?

The task is considered as identification of effective features or feature selection in text classification or ranking.

# Text Classification (1/2)

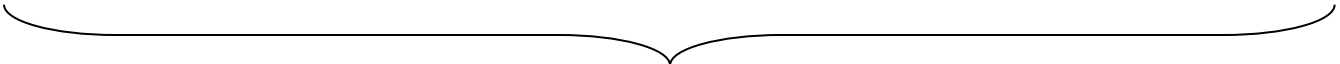


# Text Classification (2/2)



# Representation of documents

Terms	経営	情報	生産	パソコン	プリンタ	工学	品質	増加	減少	.....	楽し	述べ	する
$d_i =$	( 1 ,	0 ,	0 ,	0 ,	0 ,	1 ,	0 ,	2 ,	0 ,	..... ,	0 ,	1 ,	0 )



# Examples of Text Classification

- LABELS=BINARY
  - “spam” / “not spam”
- LABELS=TOPICS
  - “finance” / “sports” / “asia”
- LABELS=OPINION
  - “like” / “hate” / “neutral”
- LABELS=AUTHOR
  - “Shakespeare” / “Marlowe” / “Ben Jonson”
- ...

# Machine learning and our focus

- Like human learning from past experiences.
- A computer does not have “experiences”.
- A computer system learns from data, which represent some “past experiences” of an application domain.
- **Our focus:** learn a target function that can be used to predict the values of a discrete class attribute, e.g., approve or not-approved, and high-risk or low risk.
- The task is commonly called: **Supervised learning, classification.**



# The data and the goal

- **Data:** A set of data records (also called examples, instances or cases) described by
  - *k* attributes:  $A_1, A_2, \dots, A_k$ .
  - a class: Each example is labelled with a pre-defined class.
- **Goal:** To learn a **classification model** from the data that can be used to predict the classes of new (future, or test) cases/instances.

# An example: data (loan application)

Approved or not

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

# An example: the learning task

- Learn a classification model from the data
- Use the model to classify future loan applications into
  - Yes (approved) and
  - No (not approved)
- What is the class for following case/instance?

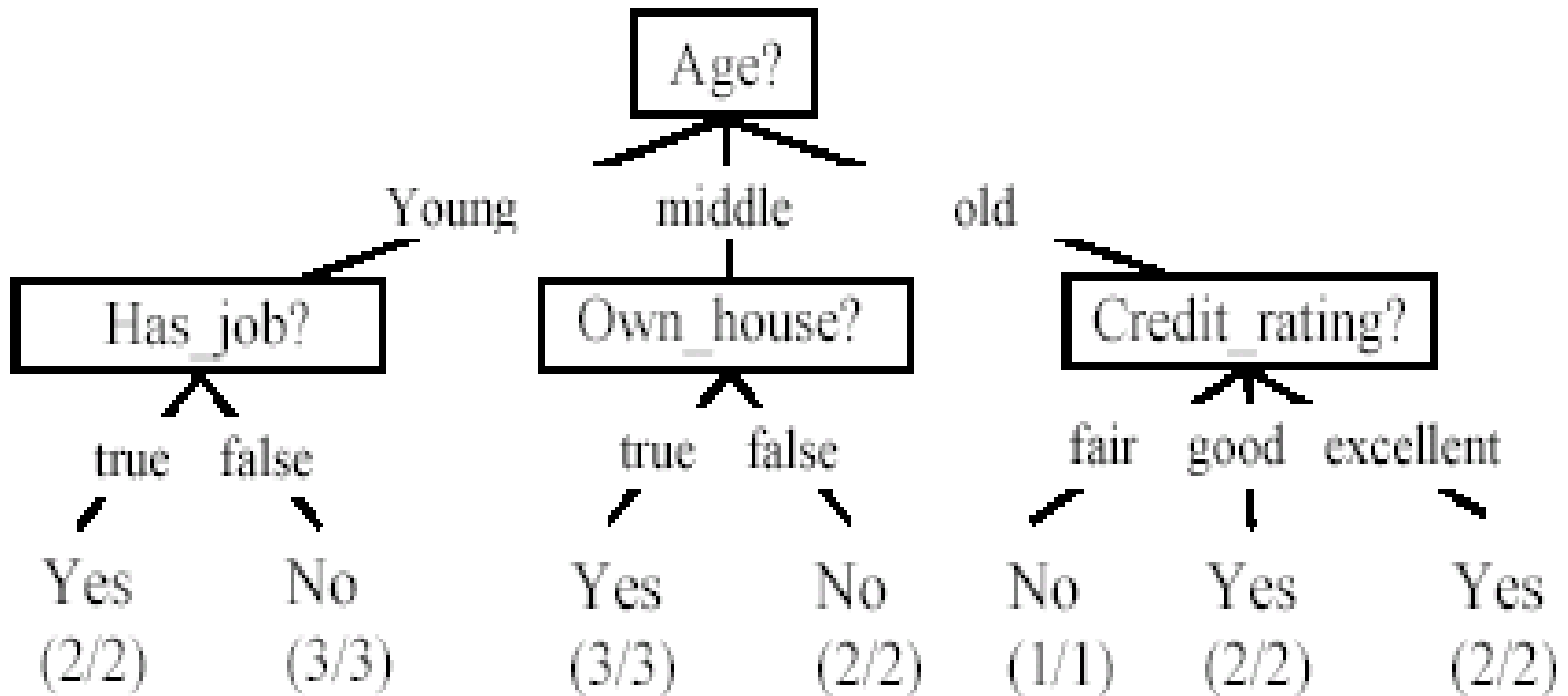
Age	Has_Job	Own_house	Credit-Rating	Class
young	false	false	good	?

# Decision Tree Learning

- Decision tree learning is one of the most widely used techniques for classification.
  - Its classification accuracy is competitive with other methods, and
  - it is very efficient.
- The classification model is a tree, called **decision tree**.
- **C4.5** by Ross Quinlan is perhaps the best known system. It can be downloaded from the Web.

# A decision tree from the loan data

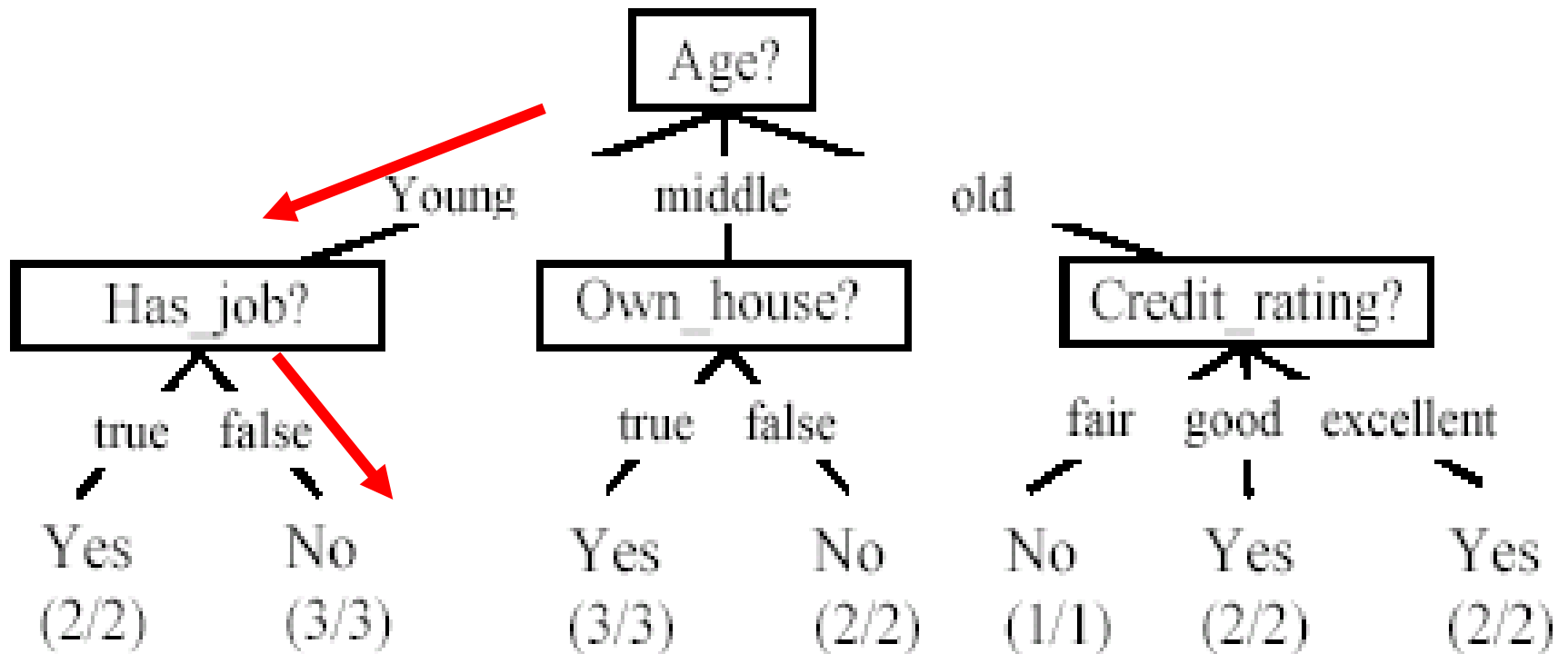
Decision nodes and leaf nodes (classes)



# Use the decision tree

Age	Has_Job	Own_house	Credit-Rating	Class
young	false	false	good	?

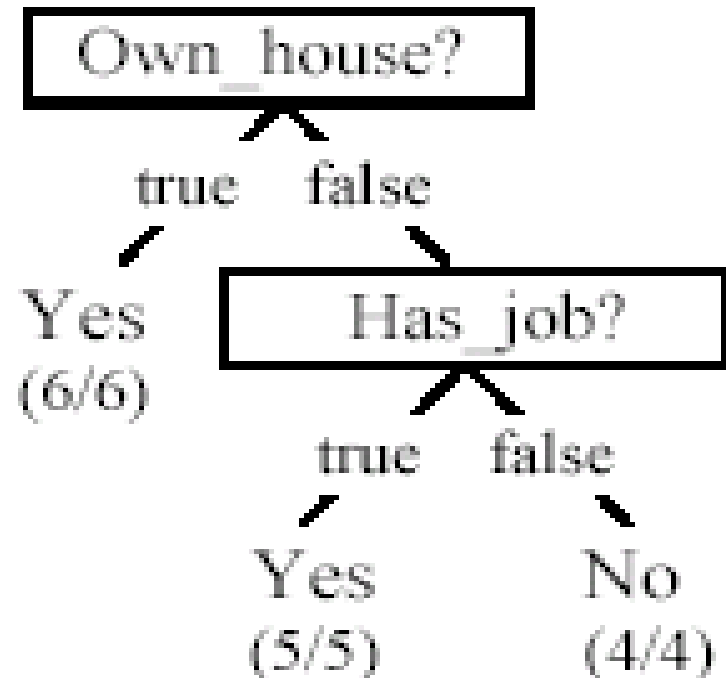
No



# From a decision tree to a set of rules

A decision tree can be converted to a set of rules

Each path from the root to a leaf is a rule.



Own\_house = true → Class = Yes [sup=6/15, conf=6/6]

Own\_house = false, Has\_job = true → Class = Yes [sup=5/15, conf=5/5]

Own\_house = false, Has\_job = false → Class = No [sup=4/15, conf=4/4]

# Algorithm for decision tree learning

- Basic algorithm (a greedy **divide-and-conquer** algorithm)
  - Assume attributes are categorical now (continuous attributes can be handled too)
  - Tree is constructed in a **top-down recursive manner**
  - At start, all the training examples are at the root
  - Examples are partitioned recursively based on selected attributes
  - Attributes are selected on the basis of an impurity function (e.g., **information gain**)
- Conditions for stopping partitioning
  - All examples for a given node belong to the same class
  - There are no remaining attributes for further partitioning – majority class is the leaf
  - There are no examples left



# Choose an attribute to partition data

- The *key* to building a decision tree - which attribute to choose in order to branch.
- The objective is to reduce impurity or uncertainty in data as much as possible.
  - A subset of data is *pure* if all instances belong to the same class.
- The *heuristic* in C4.5 is to choose the attribute with the maximum **Information Gain** or **Gain Ratio** based on information theory.

# Two possible roots, which is better?

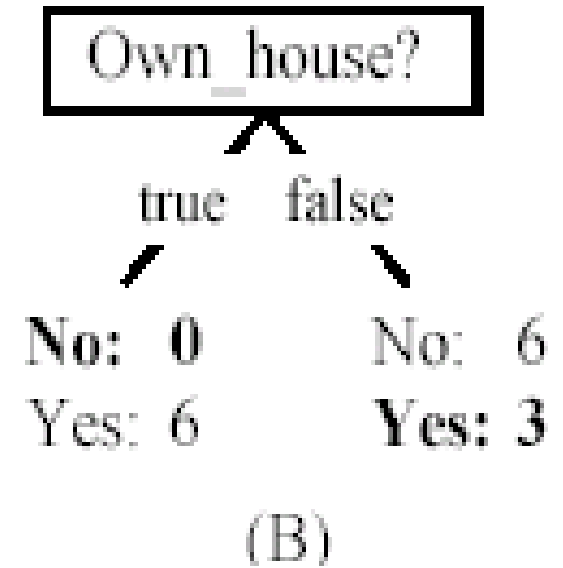
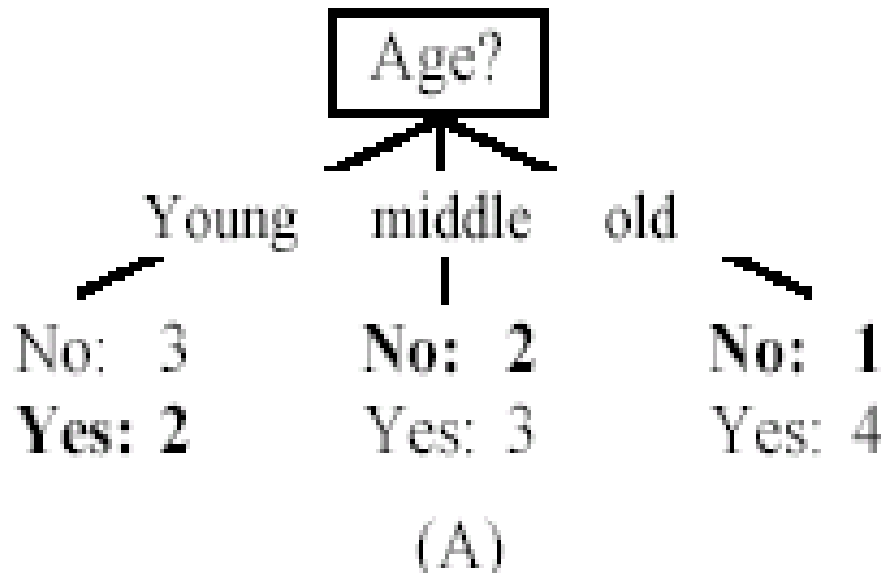


Fig. (B) seems to be better.

# Information theory: Entropy measure

- The entropy formula,

$$\text{entropy}(D) = -\sum_{j=1}^{|C|} \Pr(c_j) \log_2 \Pr(c_j)$$

$$\sum_{j=1}^{|C|} \Pr(c_j) = 1,$$

- $\Pr(c_j)$  is the probability of class  $c_j$  in data set  $D$
- We use entropy as a **measure of impurity or disorder** of data set  $D$ . (Or, a measure of information in a tree)

# Information gain

- Given a set of examples  $D$ , we first compute its entropy:

$$\text{entropy}(D) = - \sum_{j=1}^{|C|} \text{Pr}(c_j) \log_2 \text{Pr}(c_j)$$

- If we make attribute  $A_i$ , with  $v$  values, the root of the current tree, this will partition  $D$  into  $v$  subsets  $D_1, D_2, \dots, D_v$ . The expected entropy if  $A_i$  is used as the current root:

$$\text{entropy}_{A_i}(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times \text{entropy}(D_j)$$

# Information gain (cont ...)

- **Information gained** by selecting attribute  $A_i$  to branch or to partition the data is

$$\text{gain}(D, A_i) = \text{entropy}(D) - \text{entropy}_{A_i}(D)$$

- We choose the attribute with the highest gain to branch/split the current tree.

# An example

$$entropy(D) = -\frac{6}{15} \times \log_2 \frac{6}{15} - \frac{9}{15} \times \log_2 \frac{9}{15} = 0.971$$

$$\begin{aligned} entropy_{Own\_house}(D) &= -\frac{6}{15} \times entropy(D_1) - \frac{9}{15} \times entropy(D_2) \\ &= \frac{6}{15} \times 0 + \frac{9}{15} \times 0.918 \\ &= 0.551 \end{aligned}$$

$$\begin{aligned} entropy_{Age}(D) &= -\frac{5}{15} \times entropy(D_1) - \frac{5}{15} \times entropy(D_2) - \frac{5}{15} \times entropy(D_3) \\ &= \frac{5}{15} \times 0.971 + \frac{5}{15} \times 0.971 + \frac{5}{15} \times 0.722 \\ &= 0.888 \end{aligned}$$

Own\_house is the best choice for the root.

ID	Age	Has_Job	Own_House	Credit_Rating	Class
1	young	false	false	fair	No
2	young	false	false	excellent	No
3	young	true	false	good	Yes
4	young	true	true	good	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

Age	Yes	No	entropy(D <sub>i</sub> )
young	2	3	0.971
middle	3	2	0.971
old	4	1	0.722

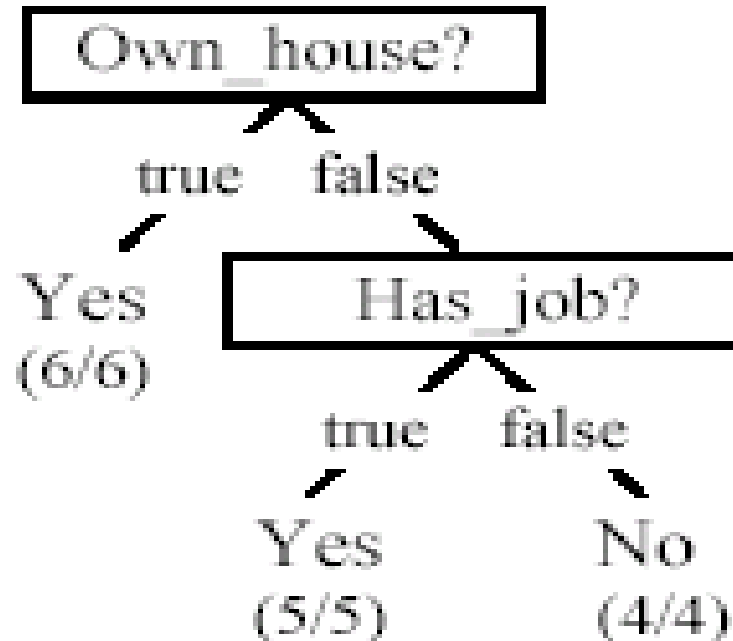
$$gain(D, Age) = 0.971 - 0.888 = 0.083$$

$$gain(D, Own\_house) = 0.971 - 0.551 = 0.420$$

$$gain(D, Has\_Job) = 0.971 - 0.647 = 0.324$$

$$gain(D, Credit\_Rating) = 0.971 - 0.608 = 0.363$$

# We build the final tree



# Bayesian classification

- **Probabilistic view:** Supervised learning can naturally be studied from a probabilistic point of view.
- Let  $A_1$  through  $A_k$  be attributes with discrete values. The class is  $C$ .
- Given a test example  $d$  with observed attribute values  $a_1$  through  $a_k$ .
- Classification is basically to compute the following posteriori probability. The prediction is the class  $c_j$  such that

$$\Pr(C = c_j \mid A_1 = a_1, \dots, A_{|A|} = a_{|A|})$$

is maximal



# Apply Bayes' Rule

$$\begin{aligned} & \Pr(C = c_j | A_1 = a_1, \dots, A_{|A|} = a_{|A|}) \\ = & \frac{\Pr(A_1 = a_1, \dots, A_{|A|} = a_{|A|} | C = c_j) \Pr(C = c_j)}{\Pr(A_1 = a_1, \dots, A_{|A|} = a_{|A|})} \\ = & \frac{\Pr(A_1 = a_1, \dots, A_{|A|} = a_{|A|} | C = c_j) \Pr(C = c_j)}{\sum_{r=1}^{|C|} \Pr(A_1 = a_1, \dots, A_{|A|} = a_{|A|} | C = c_r) \Pr(C = c_r)} \end{aligned}$$

$\Pr(C=c_j)$  is the class *prior* probability: easy to estimate from the training data.

# Computing probabilities

- The denominator  $P(A_1=a_1, \dots, A_k=a_k)$  is irrelevant for decision making since it is the same for every class.
- We only need  $P(A_1=a_1, \dots, A_k=a_k \mid C=c_i)$ , which can be written as
$$\Pr(A_1=a_1 \mid A_2=a_2, \dots, A_k=a_k, C=c_j) * \Pr(A_2=a_2, \dots, A_k=a_k \mid C=c_j)$$
- Recursively, the second factor above can be written in the same way, and so on.
- Now an assumption is needed.

# Conditional independence assumption

- All attributes are conditionally independent given the class  $C = c_j$ .
- Formally, we assume,

$$\Pr(A_1=a_1 \mid A_2=a_2, \dots, A_{|A|}=a_{|A|}, C=c_j) = \Pr(A_1=a_1 \mid C=c_j)$$

and so on for  $A_2$  through  $A_{|A|}$ . I.e.,

$$\Pr(A_1 = a_1, \dots, A_{|A|} = a_{|A|} \mid C = c_j) = \prod_{i=1}^{|A|} \Pr(A_i = a_i \mid C = c_j)$$

# Final naïve Bayesian classifier

$$\begin{aligned} & \Pr(C = c_j \mid A_1 = a_1, \dots, A_{|A|} = a_{|A|}) \\ &= \frac{\Pr(C = c_j) \prod_{i=1}^{|A|} \Pr(A_i = a_i \mid C = c_j)}{\sum_{r=1}^{|C|} \Pr(C = c_r) \prod_{i=1}^{|A|} \Pr(A_i = a_i \mid C = c_r)} \end{aligned}$$

- We are done!
- How do we estimate  $P(A_i = a_i \mid C = c_j)$ ? Easy!.

# Classify a test instance

- If we only need a decision on the most probable class for the test instance, we only need the numerator as its denominator is the same for every class.
- Thus, given a test example, we compute the following to decide the most probable class for the test instance

$$c = \arg \max_{c_j} \Pr(c_j) \prod_{i=1}^{|A|} \Pr(A_i = a_i \mid C = c_j)$$

# An example

Compute all probabilities required for classification

A	B	C
m	b	t
m	s	t
g	q	t
h	s	t
g	q	t
g	q	f
g	s	f
h	b	f
h	q	f
m	b	f

$$\Pr(C = t) = 1/2,$$

$$\Pr(C = f) = 1/2$$

$$\Pr(A = m \mid C = t) = 2/5$$

$$\Pr(A = g \mid C = t) = 2/5$$

$$\Pr(A = h \mid C = t) = 1/5$$

$$\Pr(A = m \mid C = f) = 1/5$$

$$\Pr(A = g \mid C = f) = 2/5$$

$$\Pr(A = h \mid C = f) = 2/5$$

$$\Pr(B = b \mid C = t) = 1/5$$

$$\Pr(B = s \mid C = t) = 2/5$$

$$\Pr(B = q \mid C = t) = 2/5$$

$$\Pr(B = b \mid C = f) = 2/5$$

$$\Pr(B = s \mid C = f) = 1/5$$

$$\Pr(B = q \mid C = f) = 2/5$$

Now we have a test example:

$$A = m \quad B = q \quad C = ?$$

# An Example (cont ...)

- For  $C = t$ , we have

$$\Pr(C = t) \prod_{j=1}^2 \Pr(A_j = a_j | C = t) = \frac{1}{2} \times \frac{2}{5} \times \frac{2}{5} = \frac{2}{25}$$

- For class  $C = f$ , we have

$$\Pr(C = f) \prod_{j=1}^2 \Pr(A_j = a_j | C = f) = \frac{1}{2} \times \frac{1}{5} \times \frac{2}{5} = \frac{1}{25}$$

- $C = t$  is more probable.  $t$  is the final class.

# Support Vector Machines

- Support vector machines were invented by V. Vapnik and his co-workers in 1970s in Russia and became known to the West in 1992.
- SVMs are **linear classifiers** that find a hyperplane to separate **two class** of data, positive and negative.
- **Kernel functions** are used for nonlinear separation.
- SVM not only has a rigorous theoretical foundation, but also performs classification more accurately than most other methods in applications, especially for high dimensional data.
- It is perhaps the best classifier for text classification.



# So far ...

- Text classification based on decision tree learning and Bayesian classification
- Effective features are easily identified in decision tree learning and Bayesian classification.

# Feature Selection: Why?

- Text collections have a large number of features
  - 10,000 – 1,000,000 unique words ... and more
- May make using a particular classifier feasible
  - Some classifiers can't deal with 1,000,000 features
- Reduces training time
  - Training time for some methods is quadratic or worse in the number of features
- Makes runtime models smaller and faster
- Can improve generalization (performance)
  - Eliminates noise features
  - Avoids overfitting

# Feature Selection: Frequency

- The simplest feature selection method:
  - Just use the commonest terms
  - No particular foundation
  - But it make sense why this works
    - They're the words that can be well-estimated and are most often available as evidence
  - In practice, this is often 90% as good as better methods

# Ranking Problem

- Consider a problem of ranking essays
  - Three ranking categories: good, ok, bad
  - Given a input document, predict its ranking category
- How should we formulate this problem?
- A simple multiple class solution
  - Each ranking category is an independent class
- But, there is something missing here ...
  
- We miss the ordinal relationship between classes !

# “Learning to rank” (1/2)

- Classification probably isn't the right way to think about approaching ranking:
  - Classification problems: Map to a unordered set of classes
  - Regression problems: Map to a real value
  - Ordinal regression problems: Map to an ordered set of classes
    - A fairly obscure sub-branch of statistics, but what we want here
- This formulation gives extra power:
  - Relations between levels are modeled

# “Learning to rank” (2/2)

- Assume a number of categories  $\mathbf{C}$  exist
  - These are totally ordered:  $c_1 < c_2 < \dots < c_j$
  - This is the ordinal regression setup
- Assume training data is available consisting of documents represented as feature vectors  $\psi_i$  and ranking  $c_i$
- We could do **point-wise learning**, where we try to map items of a certain rank to a subinterval (e.g, Crammer et al. 2002 PRank)
- But most work does **pair-wise learning**, where the input is a pair of results, and the class is the ordering relationship between them

# Pairwise learning: The Ranking SVM

[Herbrich et al. 1999, 2000; Joachims et al. 2002]

- Aim is to classify instance pairs as correctly ranked or incorrectly ranked
  - This turns an ordinal regression problem back into a binary classification problem

- We want a ranking function  $f$  such that

$$c_i > c_k \text{ iff } f(\psi_i) > f(\psi_k)$$

- ... or at least one that tries to do this with minimal error

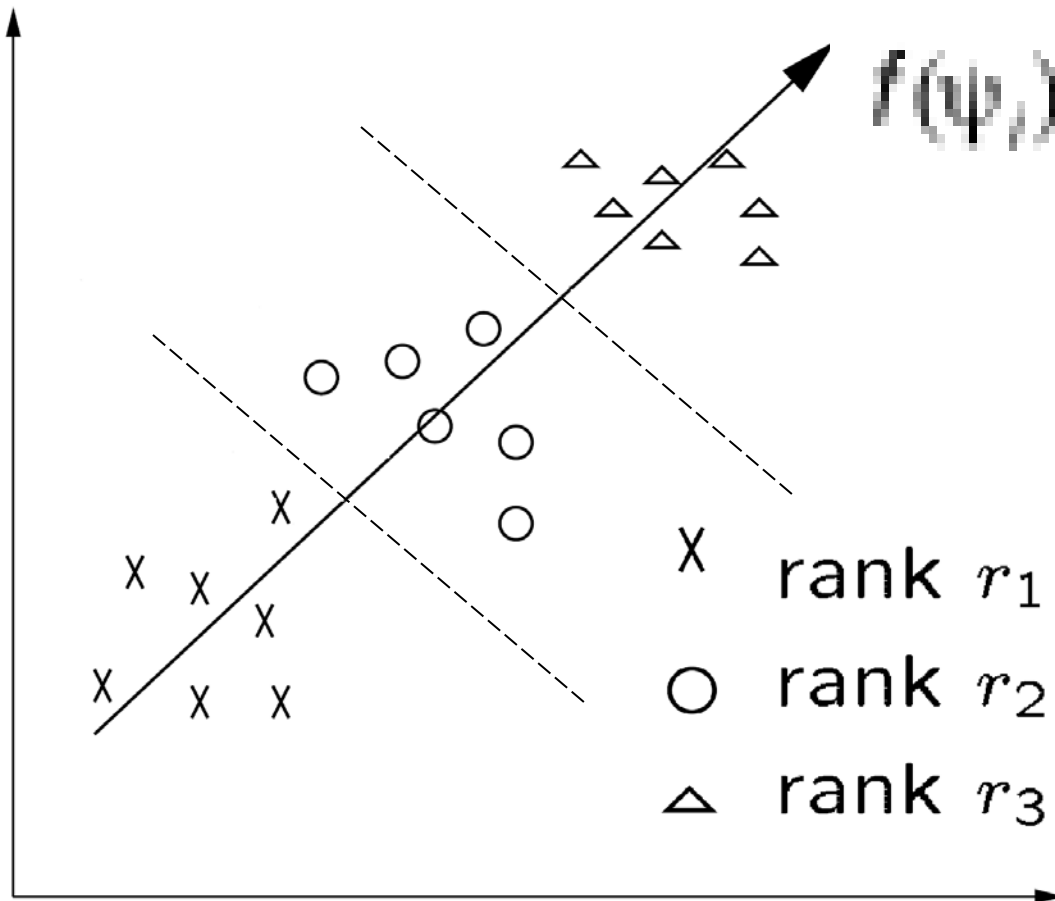
- Suppose that  $f$  is a linear function

$$f(\psi_i) = \mathbf{w} \bullet \psi_i$$

# The Ranking SVM

[Herbrich et al. 1999, 2000; Joachims et al. 2002]

- Ranking Model:  $f(\psi_i)$





# Other machine learning methods for learning to rank

- Of course!
- I've only presented the use of SVMs for machine learned ranking, but other machine learning methods have also been used successfully
  - Boosting: RankBoost
  - Ordinal Regression loglinear models
  - Neural Nets: RankNet
  - (Gradient-boosted) Decision Trees

# Summary

- Text classification based on decision tree learning and Bayesian classification
- Effective features are easily identified in decision tree learning and Bayesian classification.
- Feature selection based on mutual information and Chi-square
- Methods for learning to rank



# Combining Cues via Naive Bayes (1/2)

## Authorship ID: Who Wrote a Student's Term Paper?

Word in Text	Frequency as Student A	Frequency as Student B
optimally	97	1
certainly	84	3
typically	46	4
perspicuous	26	0
actually	13	4
whilst	6	0
the	241	229
awesome	0	63
totally	0	40
wonderful	0	26
incredibly	0	13

these stats  
come from term  
papers of *known*  
authorship

(i.e., supervised  
training)

$$\frac{P(\text{optimally}|\text{Student A})}{P(\text{optimally}|\text{Student B})} = \frac{97}{1}$$

$$\frac{P(\text{the}|\text{Student A})}{P(\text{the}|\text{Student B})} = \frac{1.1}{1}$$

# Combining Cues via Naive Bayes (2/2)

$$\frac{P(\textit{optimally}|\textit{StudentA})}{P(\textit{optimally}|\textit{StudentB})} = \frac{97}{1} \qquad \frac{P(\textit{the}|\textit{StudentA})}{P(\textit{the}|\textit{StudentB})} = \frac{1.1}{1}$$

$$\frac{P(\textit{awesome}|\textit{StudentA})}{P(\textit{awesome}|\textit{StudentB})} = \frac{0}{63}$$

$$\frac{P(\textit{StudentA})}{P(\textit{StudentB})} \times \frac{P(w_1 | \textit{StudentA})}{P(w_1 | \textit{StudentB})} \times \frac{P(w_2 | \textit{StudentA})}{P(w_2 | \textit{StudentB})} \times \dots$$

“Naive Bayes” model for classifying text

Would this kind of sentence be more typical of a student A paper or a student B paper?