

Unsupervised Learning

Lecture 3: Naïve Bayes Classifier and Evaluation Scheme

Tho Quan
qttho@hcmut.edu.vn

Agenda

- Discriminative vs. Generative Classifiers
- Naïve Bayes Classifier
- Evaluation scheme: which method is better?
- NB for Text Classification and IR Metrics

Classification Methodologies

- There are three methodologies:

a) Model a classification rule directly

Examples: k-NN, linear classifier, SVM, neural nets, ..

b) Model the probability of class memberships given input data

Examples: logistic regression, probabilistic neural nets (softmax),...

c) Make a probabilistic model of data within each class

Examples: Naive Bayes....

Classification of Classifications

	Probabilistic	Non-Probabilistic
Discriminative	<ul style="list-style-type: none">• Logistic Regression• Probabilistic neural nets•	<ul style="list-style-type: none">• K-nn• Linear classifier• SVM• Neural networks•
Generative	<ul style="list-style-type: none">• Naïve Bayes•	N.A. (?)

Probability Basics

- Prior, conditional and joint probability for random variables
 - Prior probability: $P(x)$
 - Conditional probability: $P(x_1 | x_2), P(x_2 | x_1)$
 - Joint probability: $\mathbf{x} = (x_1, x_2), P(\mathbf{x}) = P(x_1, x_2)$
 - Relationship: $P(x_1, x_2) = P(x_2 | x_1)P(x_1) = P(x_1 | x_2)P(x_2)$
 - Independence:
$$P(x_2 | x_1) = P(x_2), P(x_1 | x_2) = P(x_1), P(x_1, x_2) = P(x_1)P(x_2)$$
- Bayesian Rule

$$P(c | \mathbf{x}) = \frac{P(\mathbf{x} | c)P(c)}{P(\mathbf{x})}$$

Discriminative

$$Posterior = \frac{Likelihood \times Prior}{Evidence}$$

Generative

Probabilistic Classification Principle

- **M**aximum **A** Posterior (**MAP**) classification rule
 - For an input \mathbf{x} , find the largest one from L probabilities output by a discriminative probabilistic classifier $P(c_1 | \mathbf{x}), \dots, P(c_L | \mathbf{x})$.
 - Assign \mathbf{x} to label c^* if $P(c^* | \mathbf{x})$ is the largest.
- Generative classification with the MAP rule
 - Apply Bayesian rule to convert them into posterior probabilities

$$P(c_i | \mathbf{x}) = \frac{P(\mathbf{x} | c_i)P(c_i)}{P(\mathbf{x})} \propto P(\mathbf{x} | c_i)P(c_i)$$

for $i = 1, 2, \dots, L$

Common factor for
all L probabilities

- Then apply the MAP rule to assign a label

Naïve Bayes

- Bayes classification

$$P(c | \mathbf{x}) \propto P(\mathbf{x} | c)P(c) = P(x_1, \dots, x_n | c)P(c) \text{ for } c = c_1, \dots, c_L.$$

Difficulty: learning the joint probability $P(x_1, \dots, x_n | c)$ is infeasible!

- Naïve Bayes classification

- Assume **all input features are class conditionally independent!**

$$\begin{aligned} P(x_1, x_2, \dots, x_n | c) &= \frac{P(x_1 | x_2, \dots, x_n, c)P(x_2, \dots, x_n | c)}{P(x_2, \dots, x_n | c)} \\ &= P(x_1 | c)P(x_2, \dots, x_n | c) \\ &= P(x_1 | c)P(x_2 | c) \cdots P(x_n | c) \end{aligned}$$

Applying the independence assumption

Summarization for Example

- Bayes' Theorem

$$P(H | X) = \frac{P(X | H)P(H)}{P(X)}$$

Training tuples

Class-labeled training tuples from the *AllElectronics* customer database.

<i>RID</i>	<i>age</i>	<i>income</i>	<i>student</i>	<i>credit_rating</i>	<i>Class: buys_computer</i>
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

Naïve Bayesian Classification

1. D is a training set.

$X = (\text{age} = \text{youth}, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit_rating} = \text{fair})$

2. m classes: C_i , $i=1..m$

$$P(C_i | X) = \frac{P(X | C_i)P(C_i)}{P(X)}$$

- $P(\text{buys_computer}=\text{yes} \mid \text{age}=\text{youth}, \text{income}=\text{medium}$
 $\text{student}=\text{yes}, \text{credit_rating}=\text{fair})$

→ We maximize $P(C_i | X)$.

Naïve Bayesian Classification

3. $P(C_i) = |C_{i,D}| / |D|$
- Or: $P(C_1) = P(C_2) = \dots = P(C_m)$

Naïve Bayesian Classification

4. Compute $P(X | C_i)$:

- Assume: Class conditional independence

$$P(X | C_i) = \prod_{k=1}^n P(x_k | C_i) = P(x_1 | C_i) * P(x_2 | C_i) * .. * P(x_n | C_i)$$

5. Class label of X is C_i

- $P(C_i | X) > P(C_j | X)$ với $1 \leq j \leq m, j \neq i$

Example

$\mathbf{X} = (\text{age} = \text{youth}, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit_rating} = \text{fair})$

$$C_1 = \{X' | X'.\text{buys_computer} = \text{yes}\}$$

$$C_2 = \{X'' | X''.\text{buys_computer} = \text{no}\}$$

$$P(\text{age} = \text{youth} \mid \text{buys_computer} = \text{yes}) = 2/9 = 0.222$$

$$P(\text{age} = \text{youth} \mid \text{buys_computer} = \text{no}) = 3/5 = 0.600$$

$$P(\text{income} = \text{medium} \mid \text{buys_computer} = \text{yes}) = 4/9 = 0.444$$

$$P(\text{income} = \text{medium} \mid \text{buys_computer} = \text{no}) = 2/5 = 0.400$$

$$P(\text{student} = \text{yes} \mid \text{buys_computer} = \text{yes}) = 6/9 = 0.667$$

$$P(\text{student} = \text{yes} \mid \text{buys_computer} = \text{no}) = 1/5 = 0.200$$

$$P(\text{credit_rating} = \text{fair} \mid \text{buys_computer} = \text{yes}) = 6/9 = 0.667$$

$$P(\text{credit_rating} = \text{fair} \mid \text{buys_computer} = \text{no}) = 2/5 = 0.400$$

$$\begin{aligned} P(\mathbf{X} \mid \text{buys_computer} = \text{yes}) &= P(\text{age} = \text{youth} \mid \text{buys_computer} = \text{yes}) \times \\ &\quad P(\text{income} = \text{medium} \mid \text{buys_computer} = \text{yes}) \times \\ &\quad P(\text{student} = \text{yes} \mid \text{buys_computer} = \text{yes}) \times \\ &\quad P(\text{credit_rating} = \text{fair} \mid \text{buys_computer} = \text{yes}) \\ &= 0.222 \times 0.444 \times 0.667 \times 0.667 = 0.044. \end{aligned}$$

$$P(\mathbf{X} \mid \text{buys_computer} = \text{no}) = 0.600 \times 0.400 \times 0.200 \times 0.400 = 0.019.$$

$$P(\mathbf{X} \mid \text{buys_computer} = \text{yes})P(\text{buys_computer} = \text{yes}) = 0.044 \times 0.643 = 0.028$$

$$P(\mathbf{X} \mid \text{buys_computer} = \text{no})P(\text{buys_computer} = \text{no}) = 0.019 \times 0.357 = 0.007$$

$$P(\text{buys_computer} = \text{yes}) = 9/14 = 0.643$$

$$P(\text{buys_computer} = \text{no}) = 5/14 = 0.357$$

$$\rightarrow \mathbf{X} \in C_1$$

Zero probability

- Zero probability elimination
 - 1,000 tuples: 0 tuples with *income = low*, 990 tuples with *income = medium*, and 10 tuples with *income = high*.
- Using the Laplacian correction
 - Assume that our training database, D, is **so large** that **adding one to each count** that we need would only make a negligible difference in the estimated probability value, yet would conveniently avoid the case of probability values of zero.
 - we pretend that we have 1 more tuple for each income-value pair.

Zero probability

Suppose a dataset with 1000 tuples, income=low (0), income= medium (990), and income = high (10)

Use **Laplacian correction** (or Laplacian estimator)

Adding 1 to each case

$$\text{Prob}(\text{income} = \text{low}) = 1/1003$$

$$\text{Prob}(\text{income} = \text{medium}) = 991/1003$$

$$\text{Prob}(\text{income} = \text{high}) = 11/1003$$

Model Evaluation and Selection

- Evaluation metrics: How can we measure accuracy? Other metrics to consider?
- Use **validation test set** of class-labeled tuples instead of training set when assessing accuracy
- Methods for estimating a classifier's accuracy:
 - Holdout method, random subsampling
 - Cross-validation

Classifier Evaluation Metrics: Confusion Matrix

Confusion Matrix:

Actual class\Predicted class	C_1	$\neg C_1$
C_1	True Positives (TP)	False Negatives (FN)
$\neg C_1$	False Positives (FP)	True Negatives (TN)

Actual class\Predicted class	buy_computer = yes	buy_computer = no	Total
buy_computer = yes	6954	46	7000
buy_computer = no	412	2588	3000
Total	7366	2634	10000

- Given m classes, an entry, $\mathbf{CM}_{i,j}$ in a **confusion matrix** indicates # of tuples in class i that were labeled by the classifier as class j
- May have extra rows/columns to provide totals

Classifier Evaluation Metrics: Accuracy, Error Rate, Sensitivity and Specificity

A\P	C	¬C	
C	TP	FN	P
¬C	FP	TN	N
	P'	N'	All

- **Classifier Accuracy**, or recognition rate: percentage of test set tuples that are correctly classified

$$\text{Accuracy} = (TP + TN) / \text{All}$$

- **Error rate**: $1 - \text{accuracy}$, or
 $\text{Error rate} = (FP + FN) / \text{All}$

- **Class Imbalance Problem:**

- One class may be *rare*, e.g. fraud, or HIV-positive
- Significant *majority of the negative class* and minority of the positive class
- **Sensitivity**: True Positive recognition rate
 - **Sensitivity** = TP / P
- **Specificity**: True Negative recognition rate
 - **Specificity** = TN / N

Classifier Evaluation Metrics: Precision and Recall, and F-measures

- **Precision:** exactness – what % of tuples that the classifier labeled as positive are actually positive

$$precision = \frac{TP}{TP + FP}$$

- **Recall:** completeness – what % of positive tuples did the classifier label as positive?

$$recall = \frac{TP}{TP + FN}$$

- Perfect score is 1.0
- Inverse relationship between precision & recall
- **F measure (F_1 or F-score):** harmonic mean of precision and recall,

$$F = \frac{2 \times precision \times recall}{precision + recall}$$

- **F_β :** weighted measure of precision and recall
 - assigns β times as much weight to recall as to precision

$$F_\beta = \frac{(1 + \beta^2) \times precision \times recall}{\beta^2 \times precision + recall}$$

Classifier Evaluation Metrics: Example

Actual Class\Predicted class	cancer = yes	cancer = no	Total	Recognition(%)
cancer = yes	90	210	300	30.00 (<i>sensitivity</i>)
cancer = no	140	9560	9700	98.56 (<i>specificity</i>)
Total	230	9770	10000	96.40 (<i>accuracy</i>)

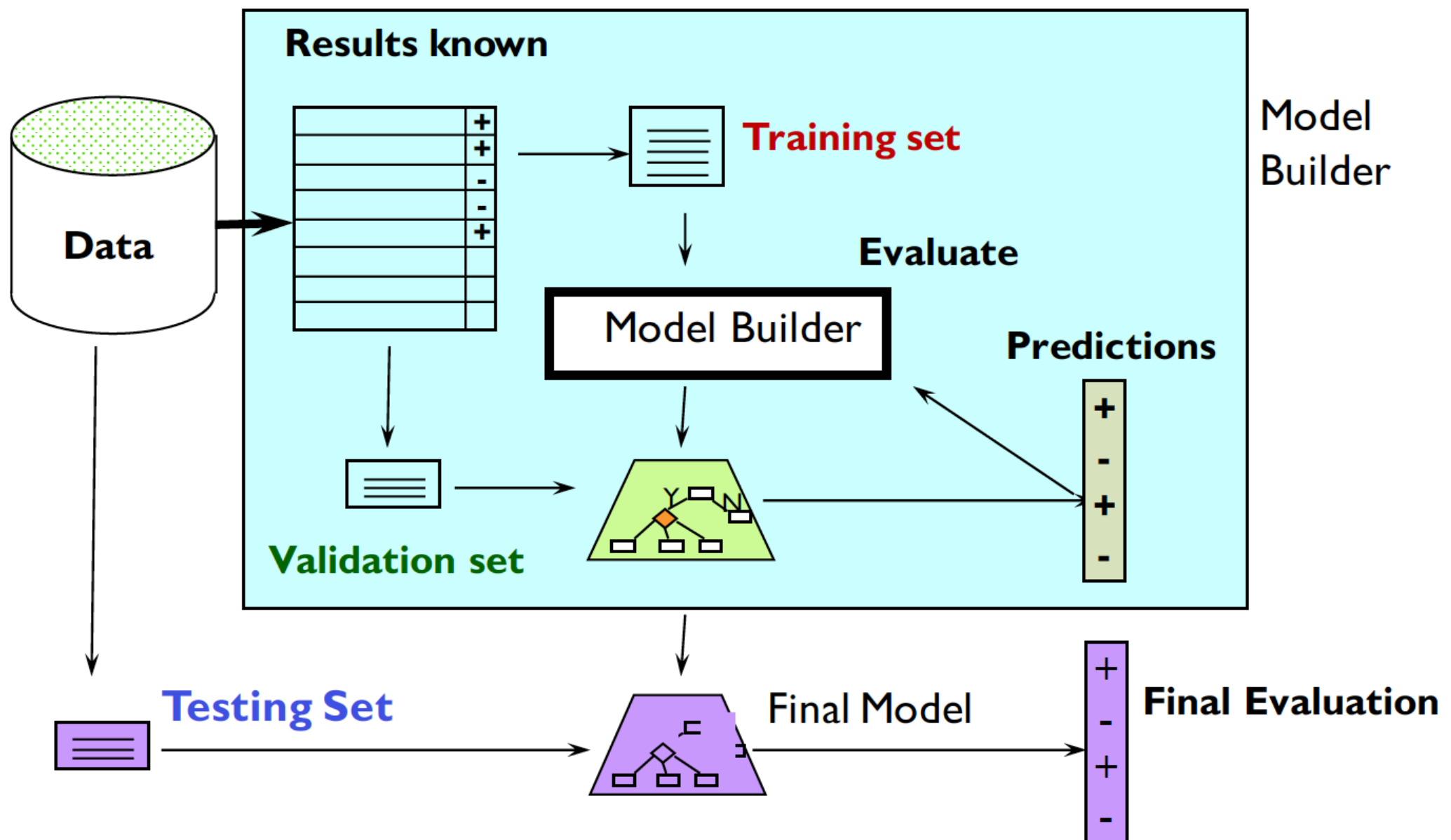
- $Precision = 90/230 = 39.13\%$

$$Recall = 90/300 = 30.00\%$$

Evaluating Classifier Accuracy: Holdout & Cross-Validation Methods

- **Holdout method**
 - Given data is randomly partitioned into two independent sets
 - Training set (e.g., 2/3) for model construction
 - Test set (e.g., 1/3) for accuracy estimation
 - Random sampling: a variation of holdout
 - Repeat holdout k times, accuracy = avg. of the accuracies obtained
- **Cross-validation** (k -fold, where $k = 10$ is most popular)
 - Randomly partition the data into k *mutually exclusive* subsets, each approximately equal size
 - At i -th iteration, use D_i as test set and others as training set
 - Leave-one-out: k folds where $k = \#$ of tuples, for small sized data
 - *Stratified cross-validation*: folds are stratified so that class dist. in each fold is approx. the same as that in the initial data

Train, Validation, Test



Uses of NB classification

- Text Classification
- Spam Filtering
- Hybrid Recommender System
 - Recommender Systems apply machine learning and data mining techniques for filtering unseen information and can predict whether a user would like a given resource
- Online Application
 - Simple Emotion Modeling

Examples of Text Classification

- CLASSES=BINARY
 - “spam” / “not spam”
- CLASSES =TOPICS
 - “finance” / “sports” / “politics”
- CLASSES =OPINION
 - “like” / “hate” / “neutral”
- CLASSES =TOPICS
 - “AI” / “Theory” / “Graphics”
- CLASSES =AUTHOR
 - “Shakespeare” / “Marlowe” / “Ben Jonson”

Naive Bayes for Text Categorization

- Attributes are text positions, values are words.

$$\begin{aligned}c_{NB} &= \operatorname{argmax}_{c_j \in C} P(c_j) \prod_i P(x_i | c_j) \\ &= \operatorname{argmax}_{c_j \in C} P(c_j) P(x_1 = \text{"our"} | c_j) \cdots P(x_n = \text{"text"} | c_j)\end{aligned}$$

- Still too many possibilities
- Assume that classification is *independent* of the positions of the words
 - Use same parameters for each position
 - Result is bag of words model (over tokens not types)

Naïve Bayes: Learning

- From training corpus, extract *Vocabulary*
- Calculate required $P(c_j)$ and $P(x_k \mid c_j)$ terms
 - For each c_j in C do
 - $docs_j \leftarrow$ subset of documents for which the target class is c_j
 - $$P(c_j) \leftarrow \frac{|docs_j|}{|\text{total \# documents}|}$$
 - $Text_j \leftarrow$ single document containing all $docs_j$
 - for each word x_k in *Vocabulary*
 - $n_k \leftarrow$ number of occurrences of x_k in $Text_j$

$$P(x_k \mid c_j) \leftarrow \frac{n_k + \alpha}{n + \alpha |Vocabulary|}$$

Naïve Bayes: Classifying

- positions \leftarrow all word positions in current document which contain tokens found in *Vocabulary*
- Return c_{NB} , where

$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_{i \in \text{positions}} P(x_i | c_j)$$