**TON DUC THANG UNIVERSITY**
Faculty of Information Technology
Computer Science

# Naïve Bayes Classification

## LÊ ANH CƯỜNG

# Outline

- Classification problems
- What kind of machine learning type for NB classifier
- Text categorization with NB
- MultiNominal NB classification
- Gausian NB classification
- Benulli NB classification

# Classification Problem and Applications

| NAME | RANK | YEARS | TENURED |
|------|------|-------|---------|
| Mike | Assistant Prof | 3 | no |
| Mary | Assistant Prof | 7 | yes |
| Bill | Professor | 2 | yes |
| Jim | Associate Prof | 7 | yes |
| Dave | Assistant Prof | 6 | no |
| Anne | Associate Prof | 3 | no |

# Classification Problem and Applications

| age | income | student | credit_rating | buys_computer |
|---|---|---|---|---|
| <=30 | high | no | fair | no |
| <=30 | high | no | excellent | no |
| 31…40 | high | no | fair | yes |
| >40 | medium | no | fair | yes |
| >40 | low | yes | fair | yes |
| >40 | low | yes | excellent | no |
| 31…40 | low | yes | excellent | yes |
| <=30 | medium | no | fair | no |
| <=30 | low | yes | fair | yes |
| >40 | medium | yes | fair | yes |
| <=30 | medium | yes | excellent | yes |
| 31…40 | medium | no | excellent | yes |
| 31…40 | high | yes | fair | yes |
| >40 | medium | no | excellent | no |

# Classification Problem and Applications
# Is this spam?

**Subject:** **Important notice!**
**From:** Stanford University <newsforum@stanford.edu>
**Date:** October 28, 2011 12:34:16 PM PDT
**To:** undisclosed-recipients:;

Greats News!

You can now access the latest news by using the link below to login to Stanford University News Forum.

http://www.123contactform.com/contact-form-StanfordNew1-236335.html

Click on the above link to login for more information about this new exciting forum. You can also copy the above link to your browser bar and login for more information about the new services.

© Stanford University. All Rights Reserved.

# Positive or negative movie review?

- unbelievably disappointing
- Full of zany characters and richly applied satire, and some great plot twists
- this is the greatest screwball comedy ever filmed
- It was pathetic. The worst part about it was the boxing scenes.

# What is the subject of this article?

**MEDLINE Article**



?

**MeSH Subject Category Hierarchy**
- Antogonists and Inhibitors
- Blood Supply
- Chemistry
- Drug Therapy
- Embryology
- Epidemiology
- ...

# Outline

- Classification problems
- <mark>What kind of machine learning type for NB classifier</mark>
- Text categorization with NB
- MultiNominal NB classification
- Gausian NB classification
- Benulli NB classification

# Probabilistic models

- Probabilistic models see features and target variables as random variables.

- The process of modelling represents and **manipulates the level of uncertainty** with respect to these variables.

- There are two types of probabilistic models: **Predictive (or Discriminative) and Generative**.

- <mark>Predictive probability models</mark> use the idea of a **conditional probability** distribution $P(Y|X)$ from which $Y$ can be predicted from $X$.

- <mark>Generative models</mark> estimate the **joint distribution** $P(Y, X)$.

# Probabilistic models

- Once we know the joint distribution for the generative models, we can derive any conditional or marginal distribution involving the same variables.
- Thus, the generative model is capable of creating new data points and their labels, knowing the joint probability distribution.
- The joint distribution looks for a relationship between two variables. Once this relationship is inferred, it is possible to infer new data points.

# Naïve Bayes

- Naive bayes is a Generative model which is based on the joint probability, p( B, A), of the inputs B and the label A.
- It makes their predictions by using Bayes rules to calculate p(A | B), and then picking the most likely label A.
- Discriminative classifiers model the posterior p(A | B) directly, or learn a direct map from inputs B to the class labels .

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

# Naïve Bayes

Probability of B occurring given evidence A has already occurred

Probability of A occurring

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Probability of A occurring given evidence B has already occurred

Probability of B occurring

# Text Classification: definition

- *Input*:
  - a document $d$
  - a fixed set of classes $C = \{c_1, c_2, ..., c_J\}$

- *Output*: a predicted class $c \in C$

# Classification Methods: Supervised Machine Learning

- *Input:*
  - a document *d*
  - a fixed set of classes $C = \{c_1, c_2, ..., c_J\}$
  - A training set of *m* hand-labeled documents $(d_1, c_1), ...., (d_m, c_m)$

- *Output:*
  - a learned classifier $\gamma : d \rightarrow c$

# Naïve Bayes Intuition

- Simple ("naïve") classification method based on Bayes rule

- Relies on very simple representation of document
  - Bag of words

# The bag of words representation

$$\gamma(\boxed{\begin{array}{l}\text{I love this movie! It's sweet,} \\ \text{but with satirical humor. The} \\ \text{dialogue is great and the} \\ \text{adventure scenes are fun...  It} \\ \text{manages to be whimsical and} \\ \text{romantic while laughing at the} \\ \text{conventions of the fairy tale} \\ \text{genre. I would recommend it to} \\ \text{just about anyone. I've seen} \\ \text{it several times, and I'm} \\ \text{always happy to see it again} \\ \text{whenever I have a friend who} \\ \text{hasn't seen it yet.}\end{array}}) = c$$

# Bayes' Rule Applied to Documents and Classes

- For a document *d* and a class *c*

$$P(c \mid d) = \frac{P(d \mid c)P(c)}{P(d)}$$

# Naïve Bayes Classifier (I)

$$c_{MAP} = \underset{c \in C}{\text{argmax}} \, P(c \mid d)$$

MAP is "maximum a posteriori" = most likely class

$$= \underset{c \in C}{\text{argmax}} \, \frac{P(d \mid c)P(c)}{P(d)}$$

Bayes Rule

$$= \underset{c \in C}{\text{argmax}} \, P(d \mid c)P(c)$$

Dropping the denominator

# Naïve Bayes Classifier (II)

$$c_{MAP} = \underset{c \in C}{\mathrm{argmax}}\, P(d \mid c)P(c)$$

$$= \underset{c \in C}{\mathrm{argmax}}\, P(x_1, x_2, \square\ , x_n \mid c)P(c)$$

Document d represented as features x1..xn

# Naïve Bayes Classifier (IV)

$$c_{MAP} = \underset{c \in C}{\text{argmax}}\, P(x_1, x_2, \square\ , x_n \mid c) P(c)$$

$O(|X|^n \bullet |C|)$ parameters

How often does this class occur?

Could only be estimated if a very, very large number of training examples was available.

We can just count the relative frequencies in a corpus

# Multinomial Naïve Bayes Independence Assumptions

$$P(x_1, x_2, \square, x_n \mid c)$$

- **Bag of Words assumption**: Assume position doesn't matter

- **Conditional Independence**: Assume the feature probabilities $P(x_i \mid c_j)$ are independent given the class $c$.

$$P(x_1, \square, x_n \mid c) = P(x_1 \mid c) \cdot P(x_2 \mid c) \cdot P(x_3 \mid c) \cdot \ldots \cdot P(x_n \mid c)$$

# Multinomial Naïve Bayes Classifier

$$c_{MAP} = \underset{c \in C}{\mathrm{argmax}}\, P(x_1, x_2, \ldots, x_n \mid c) P(c)$$

$$c_{NB} = \underset{c \in C}{\mathrm{argmax}}\, P(c_j) \widetilde{\bigcirc}_{x \in X} P(x \mid c)$$

# Applying Multinomial Naive Bayes Classifiers to Text Classification

positions ← all word positions in test document

$$c_{NB} = \underset{c_j \in C}{\text{argmax}} \, P(c_j) \prod_{i \in positions} P(x_i \mid c_j)$$

# Learning the Multinomial Naïve Bayes Model

- First attempt: maximum likelihood estimates
  - simply use the frequencies in the data

$$\hat{P}(c_j) = \frac{doccount(C = c_j)}{N_{doc}}$$

$$\hat{P}(w_i \mid c_j) = \frac{count(w_i, c_j)}{\sum_{w \in V} count(w, c_j)}$$

# Parameter estimation

$$\hat{P}(w_i \mid c_j) = \frac{count(w_i, c_j)}{\sum_{w \in V} count(w, c_j)}$$

fraction of times word $w_i$ appears among all words in documents of topic $c_j$

- Create mega-document for topic $j$ by concatenating all docs in this topic
  - Use frequency of $w$ in mega-document

# An Example

- I like this book: c1
- I do not like this movie: c2
- She loves this book: c1
- He do not like this book: c2
- I and her love the movie: c
- P(c1)= ?   P(c2)=?
- P(this|c1) = ?
- P(this|c2) = ?

$$\hat{P}(c_j) = \frac{doccount(C = c_j)}{N_{doc}}$$

$$\hat{P}(w_i \mid c_j) = \frac{count(w_i, c_j)}{\sum_{w \in V} count(w, c_j)}$$

# Problem with Maximum Likelihood

- What if we have seen no training documents with the word *fantastic* and classified in the topic **positive** (*thumbs-up*)?

$$\hat{P}(\text{"fantastic"} \mid \text{positive}) = \frac{count(\text{"fantastic", positive})}{\displaystyle\sum_{w \in V} count(w, \text{positive})} = 0$$

- Zero probabilities cannot be conditioned away, no matter the other evidence!

$$c_{MAP} = \text{argmax}_c \, \hat{P}(c) \prod_i \hat{P}(x_i \mid c)$$

# Laplace (add-1) smoothing for Naïve Bayes

$$\hat{P}(w_i \mid c) = \frac{count(w_i, c) + 1}{\sum_{w \in V} \left( count(w, c) + 1 \right)}$$

$$= \frac{count(w_i, c) + 1}{\left( \sum_{w \in V} count(w, c) \right) + |V|}$$

# Multinomial Naïve Bayes: Learning

- From training corpus, extract *Vocabulary*

- Calculate $P(c_j)$ terms
  - For each $c_j$ in $C$ do
    $$docs_j \leftarrow \text{all docs with class} = c_j$$
    $$P(c_j) \neg \frac{|docs_j|}{|\text{total \# documents}|}$$

- Calculate $P(w_k \mid c_j)$ terms
  - $Text_j \leftarrow$ single doc containing all $docs_j$
  - For each word $w_k$ in *Vocabulary*
    $$n_k \leftarrow \text{\# of occurrences of } w_k \text{ in } Text_j$$
    $$P(w_k \mid c_j) \neg \frac{n_k + a}{n + a \, |Vocabulary|}$$

# Laplace (add-1) smoothing: unknown words

Add one extra word to the vocabulary, the "unknown word" $w_u$

$$\hat{P}(w_u \mid c) = \frac{count(w_u, c) + 1}{\left(\sum_{w \in V} count(w, c)\right) + |V + 1|}$$

$$= \frac{1}{\left(\sum_{w \in V} count(w, c)\right) + |V + 1|}$$

$$\hat{P}(c) = \frac{N_c}{N}$$

$$\hat{P}(w \mid c) = \frac{count(w,c)+1}{count(c)+|V|}$$

| | Doc | Words | Class |
|---|---|---|---|
| Training | 1 | Chinese Beijing Chinese | c |
| | 2 | Chinese Chinese Shanghai | c |
| | 3 | Chinese Macao | c |
| | 4 | Tokyo Japan Chinese | j |
| Test | 5 | Chinese Chinese Chinese Tokyo Japan | ? |

**Priors:**

$P(c)= \dfrac{3}{4}$

$P(j)= \dfrac{1}{4}$

**Conditional Probabilities:**

P(Chinese|$c$) =   (5+1) / (8+6) = 6/14 = 3/7

P(Tokyo|$c$)   =   (0+1) / (8+6) = 1/14

P(Japan|$c$)    =   (0+1) / (8+6) = 1/14

P(Chinese|$j$) =   (1+1) / (3+6) = 2/9

P(Tokyo|$j$)    =   (1+1) / (3+6) = 2/9

P(Japan|$j$)     =   (1+1) / (3+6) = 2/9

**Choosing a class:**

P(c|d5)  $\propto$  $3/4 * (3/7)^3 * 1/14 * 1/14$

$\approx 0.0003$

P(j|d5)  $\propto$  $1/4 * (2/9)^3 * 2/9 * 2/9$

$\approx 0.0001$

# Evaluation: Classic Reuters-21578 Data Set

- Most (over)used data set, 21,578 docs (each 90 types, 200 toknens)

- 9603 training, 3299 test articles (ModApte/Lewis split)

- 118 categories
  - An article can be in more than one category
  - Learn 118 binary category distinctions

- Average document (with at least one category) has 1.24 classes

- Only about 10 out of 118 categories are large

Common categories
(#train, #test)

- Earn (2877, 1087)
- Acquisitions (1650, 179)
- Money-fx (538, 179)
- Grain (433, 149)
- Crude (389, 189)

- Trade (369,119)
- Interest (347, 131)
- Ship (197, 89)
- Wheat (212, 71)
- Corn (182, 56)

# Reuters Text Categorization data set (**Reuters-21578**) document

<REUTERS TOPICS="YES" LEWISSPLIT="TRAIN" CGISPLIT="TRAINING-SET" OLDID="12981" NEWID="798">

<DATE> 2-MAR-1987 16:51:43.42</DATE>

<TOPICS><D>livestock</D><D>hog</D></TOPICS>

<TITLE>AMERICAN PORK CONGRESS KICKS OFF TOMORROW</TITLE>

<DATELINE>    CHICAGO, March 2 - </DATELINE><BODY>The American Pork Congress kicks off tomorrow, March 3, in Indianapolis with 160 of the nations pork producers from 44 member states determining industry positions on a number of issues, according to the National Pork Producers Council, NPPC.

    Delegates to the three day Congress will be considering 26 resolutions concerning various issues, including the future direction of farm policy and the tax law as it applies to the agriculture sector. The delegates will also debate whether to endorse concepts of a national PRV (pseudorabies virus) control and eradication program, the NPPC said.

    A large trade show, in conjunction with the congress, will feature the latest in technology in all areas of the industry, the NPPC added. Reuter

&#3;</BODY></TEXT></REUTERS>

# Micro- vs. Macro-Averaging

- If we have more than one class, how do we combine multiple performance measures into one quantity?

- **Macroaveraging**: Compute performance for each class, then average.

- **Microaveraging**: Collect decisions for all classes, compute contingency table, evaluate.

# Micro- vs. Macro-Averaging: Example

| Class 1 | | |
|---|---|---|
| | Truth: yes | Truth: no |
| Classifier: yes | 10 | 10 |
| Classifier: no | 10 | 970 |

| Class 2 | | |
|---|---|---|
| | Truth: yes | Truth: no |
| Classifier: yes | 90 | 10 |
| Classifier: no | 10 | 890 |

| Micro Ave. Table | | |
|---|---|---|
| | Truth: yes | Truth: no |
| Classifier: yes | 100 | 20 |
| Classifier: no | 20 | 1860 |

- Macroaveraged precision: $(0.5 + 0.9)/2 = 0.7$
- Microaveraged precision: $100/120 = .83$
- Microaveraged score is dominated by score on common classes

# Development Test Sets and Cross-validation

| Training set | Development Test Set | Test Set |
|---|---|---|

- Metric: P/R/F1 or Accuracy

- Unseen test set
  - avoid overfitting ('tuning to the test set')
  - more conservative estimate of performance
- Cross-validation over multiple splits
  - Handle sampling errors from different datasets
  - Pool results over each split
  - Compute pooled dev set performance

| Training Set | Dev Test | |
|---|---|---|

| Training Set | | Dev Test |
|---|---|---|

| Dev Test | | Training Set |
|---|---|---|

| Test Set |
|---|

# Code

```
1  from sklearn.feature_extraction.text import CountVectorizer
2  #Feature extraction, vectorize text features
3  vec=CountVectorizer()
4  X_train=vec.fit_transform(X_train)
5  X_test=vec.transform(X_test)
```

```
1  from sklearn.naive_bayes import MultinomialNB
2  mnb=MultinomialNB()
3  mnb.fit(X_train,y_train)
4  y_predict=mnb.predict(X_test)
```
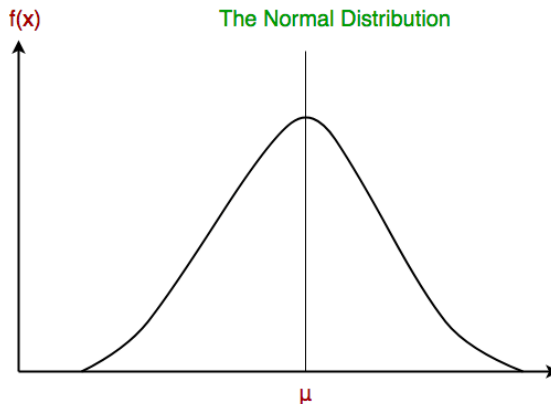
```
1  from sklearn.metrics import classification_report
2  print('The accuracy of Navie Bayes Classifier is',mnb.score(X_test,y_test))
3  print(classification_report(y_test,y_predict,target_names=news.target_names))
```

https://www.programmersought.com/article/10835149414/

# Gaussian Naive Bayes

- In Gaussian Naive Bayes, continuous values associated with each feature are assumed to be distributed according to a Gaussian distribution.
- A Gaussian distribution is also called Normal distribution. When plotted, it gives a bell shaped curve which is symmetric about the mean of the feature values as shown below:

# Gaussian Naive Bayes

f(x)

**The Normal Distribution**

μ

The likelihood of the features is assumed to be Gaussian, hence, conditional probability is given by:

$$P(x_i|y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} exp\left(-\frac{(x_i-\mu_y)^2}{2\sigma_y^2}\right)$$

```python
# load the iris dataset
from sklearn.datasets import load_iris
iris = load_iris()

# store the feature matrix (X) and response vector (y)
X = iris.data
y = iris.target

# splitting X and y into training and testing sets
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test

# training the model on training set
from sklearn.naive_bayes import GaussianNB
gnb = GaussianNB()
gnb.fit(X_train, y_train)

# making predictions on the testing set
y_pred = gnb.predict(X_test)

# comparing actual response values (y_test) with predicted res
from sklearn import metrics
print("Gaussian Naive Bayes model accuracy(in %):", metrics.ac
```

# Bernoulli Naive Bayes

•**Bernoulli Naive Bayes**: In the multivariate Bernoulli event model, features are independent booleans (binary variables) describing inputs. Like the multinomial model, this model is popular for document classification tasks, where binary term occurrence(i.e. a word occurs in a document or not) features are used rather than term frequencies(i.e. frequency of a word in the document).

# Code

```python
from sklearn.naive_bayes import BernoulliNB
BNB = BernoulliNB()
BNB.fit(X_train, Y_train)
accuracy_score_bnb = metrics.accuracy_score(BNB.predict(X_test),Y_test)
print('BNB accuracy = ' + str('{:4.2f}'.format(accuracy_score_bnb*100))+'%')
```

```
BNB accuracy = 60.61%
```

https://towardsdatascience.com/sentiment-analysis-introduction-to-naive-bayes-algorithm-96831d77ac91