

UVA CS 6316

– Fall 2015 Graduate: Machine Learning

Lecture 14: Generative Bayes Classifier

Dr. Yanjun Qi

University of Virginia

Department of
Computer Science

10/21/15

1

Where are we ? →

Five major sections of this course

- Regression (supervised)
- Classification (supervised)
- Unsupervised models
- Learning theory
- Graphical models

10/21/15

2

Where are we ? →

Three major sections for classification

- We can divide the large variety of classification approaches into **roughly three major types**

1. Discriminative

- directly estimate a decision rule/boundary
- e.g., support vector machine, decision tree



2. Generative:

- build a generative statistical model
- e.g., **naïve bayes classifier**, Bayesian networks

3. Instance based classifiers

- Use observation directly (no models)
- e.g. K nearest neighbors

10/21/15

3

X_1	X_2	X_3	C

A Dataset for classification

$$f : X \rightarrow C$$

Output as Discrete
Class Label
 C_1, C_2, \dots, C_L

$$P(C | \mathbf{X})$$

- **Data/points/instances/examples/samples/records**: [rows]
- **Features/attributes/dimensions/independent variables/covariates/predictors/regressors**: [columns, except the last]
- **Target/outcome/response/label/dependent variable**: special column to be predicted [last column]

10/21/15

4

Bayes classifier

- Treat each attribute and class label as random variables.
- Given a sample \mathbf{x} with attributes (x_1, x_2, \dots, x_p) :
 - Goal is to predict class C .
 - Specifically, we want to find the value of C_i that maximizes $p(C_i | x_1, x_2, \dots, x_p)$.
- Generative Bayes classification

$$P(C | \mathbf{X}) \propto P(\mathbf{X} | C)P(C) = \underbrace{P(X_1, \dots, X_p | C)} P(C)$$

Difficulty: learning the joint probability $P(X_1, \dots, X_p | C)$

10/21/15

5

Probabilistic Models of text documents

$$\Pr(D | C = c) \quad ? \quad D = (w_1, w_2, \dots, w_k)$$



Two
Previous
models

$$\Pr(W_1 = \text{true}, W_2 = \text{false}, \dots, W_k = \text{true} | C = c)$$

Multivariate Bernoulli Distribution

$$\Pr(W_1 = n_1, W_2 = n_2, \dots, W_k = n_k | C = c)$$

Multinomial Distribution

10/21/15

6

Today : Generative Bayes Classifier

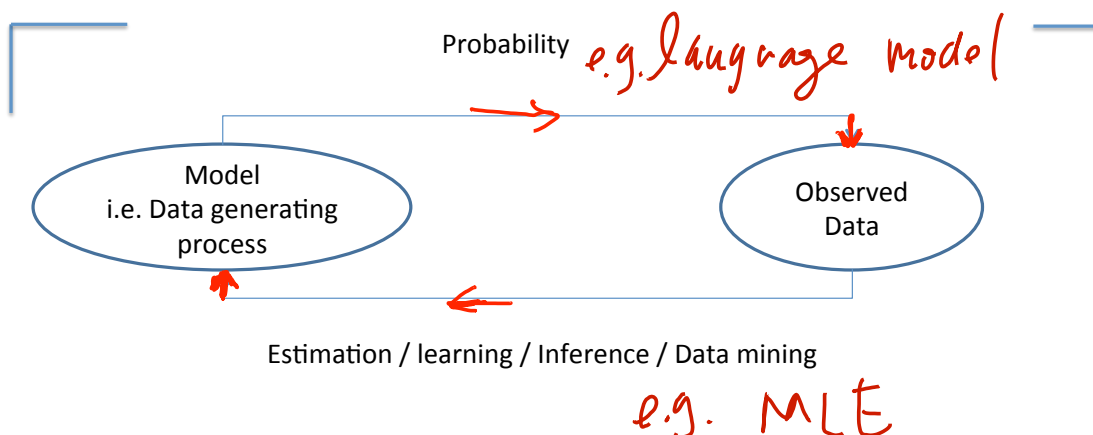
- ✓ Multinomial naïve Bayes classifier as **Conditional Stochastic Language Models**

- A unigram Language model approximates how a text document is produced.

$$\Pr(W_1 = n_1, \dots, W_k = n_k \mid C = c)$$

- ✓ Maximum Likelihood Estimation of parameters
- ✓ Gaussian Naïve Bayes Classifiers

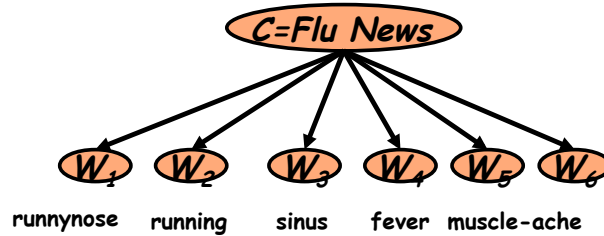
The Big Picture



But how to specify a model?

Build a *generative model* that approximates how data is produced.

Model 1: Multivariate Bernoulli



- **Conditional Independence Assumption:** Features (word presence) are *independent* of each other given the class variable:

this is naïve

$$\Pr(W_1 = \text{true}, W_2 = \text{false}, \dots, W_k = \text{true} \mid C = c)$$

$$= P(W_1 = \text{true} \mid C) \cdot P(W_2 = \text{false} \mid C) \cdot \dots \cdot P(W_k = \text{true} \mid C)$$

- Multivariate Bernoulli model is appropriate for **binary feature variables**

10/21/15

9

Review: Bernoulli Distribution e.g. Coin Flips

- You flip a coin
 - Head with probability p
 - Binary random variable
 - **Bernoulli trial** with success probability p

$$\Pr(W_i = \text{true} \mid C = c)$$

10/21/15

10

Model 2: Multinomial Naïve Bayes

- 'Bag of words' representation of text

word	frequency
grain(s)	3
oilseed(s)	2
total	3
wheat	1
maize	1
soybean	1
tonnes	1
...	...

$$\Pr(W_1 = n_1, \dots, W_k = n_k \mid C = c)$$

Can be represented as a multinomial distribution.

Words = like colored balls, there are K possible type of them (i.e. from a dictionary of K words)

Document = contains N words, each word occurs n_i times (like a bag of N colored balls)

WHY is this naïve ???

multinomial coefficient, normally can leave out in practical calculations.

Why naïve ???

$$\frac{N!}{n_1! n_2! \dots n_k!} \theta_1^{n_1} \theta_2^{n_2} \dots \theta_k^{n_k}$$

10/21/15

Multinomial Naïve Bayes as \rightarrow a generative model that approximates how a text string is produced

- Stochastic Language Models:
 - Model probability of generating strings (each word in turn following the sequential ordering in the string) in the language (commonly all strings over dictionary Σ).
 - E.g., unigram model

$P(s \mid \text{model } C_i)$

Model C_1

0.2 the

0.1 a

0.01 boy

0.01 dog

0.03 said

0.02 likes

10/21/15

the	boy	likes	the	dog
0.2	0.01	0.02	0.2	0.01

Multiply all five terms

$$P(s \mid C_1) = 0.00000008$$

Adapt From Manning' textCat tutorial ¹²

Multinomial Naïve Bayes as Conditional Stochastic Language Models

- Model conditional *probability* of generating any string from two possible models

Model C1	Model C2	
0.2 the	0.2 the	
0.01 boy	0.0001 boy	
0.0001 said	0.03 said	
0.0001 likes	0.02 likes	
0.0001 black	0.1 black	
0.0005 dog	0.01 dog	
0.01 garden	0.0001 garden	

$P(C2) = P(C1)$

the	boy	likes	black	dog
0.2	0.01	0.0001	0.0001	0.0005
0.2	0.0001	0.02	0.1	0.01

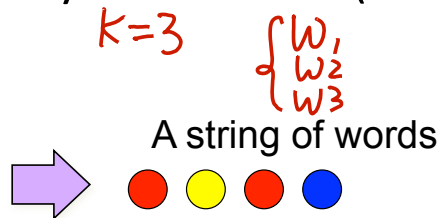
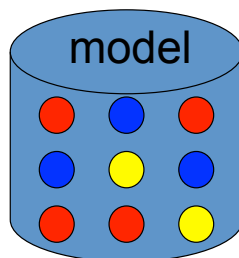
$\rightarrow P(d|C2)P(C2)$
 $P(s|C2)P(C2) > P(s|C1)P(C1)$
 $\rightarrow S$ is more likely to be from class C2

13

10/21/15

A Physical Metaphor

- Colored balls are randomly drawn from (with replacement)

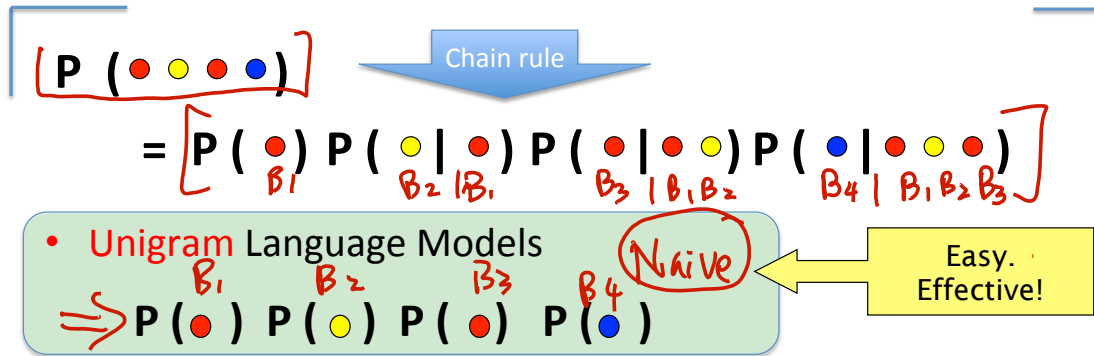


$$P(\text{red, yellow, red, blue}) = P(\text{red}) P(\text{yellow}) P(\text{red}) P(\text{blue})$$

10/21/15

14

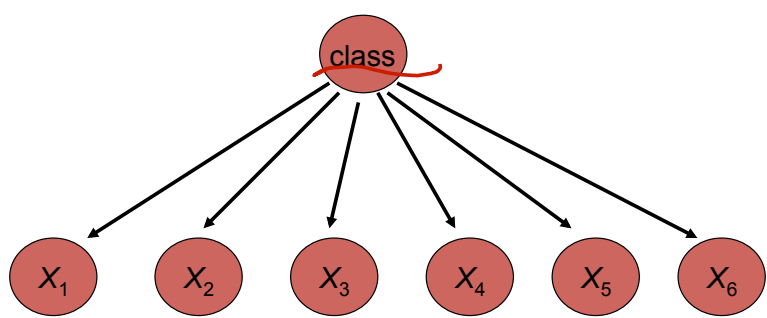
Unigram language model → More general: Generating language string from a probabilistic model



NAIVE: conditional independent on each position of the string

- Also could be **bigram** (or generally, *n*-gram) Language Models
 $P(b_1) P(b_2 | b_1) P(b_3 | b_2) P(b_4 | b_3) \dots P(b_j | b_{j-1})$

Multinomial Naïve Bayes = a class conditional unigram language model



- Think of X_i as the word on the i^{th} position in the document string
- Effectively, the probability of each class is done as a class-specific unigram language model

Using Multinomial Naive Bayes Classifiers to Classify Text: Basic method

- Attributes are text positions, values are words.

$= \text{argmax}_{c_j} P(C|X)$

$$c_{NB} = \text{argmax}_{c_j \in C} P(c_j) \prod_i P(x_i | c_j)$$

$$= \text{argmax}_{c_j \in C} P(c_j) P(x_1 = \text{"the"} | c_j) \cdots P(x_n = \text{"the"} | c_j)$$

the boy like the dog

?

- Still too many possibilities
 - Use same parameters for a word across positions
 - Result is bag of words model (over word tokens)

Multinomial Naïve Bayes: Classifying Step

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

Easy to implement, no need to construct bag-of-words vector explicitly !!!

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

Equal to, (leaving out of multinomial coefficient)

$$\Pr(W_1 = n_1, \dots, W_k = n_k | C = c_j)$$

the	boy	likes	black	dog
0.2	0.01	0.0001	0.0001	0.0005
0.2	0.0001	0.02	0.1	0.01

$P(s|C2) P(C2) > P(s|C1) P(C1)$

Unknown Words

- How to handle words in the test corpus that did not occur in the training data, i.e. *out of vocabulary* (OOV) words?
- Train a model that includes an explicit symbol for an unknown word (<UNK>).
 - Choose a vocabulary in advance and replace **other (i.e. not in vocabulary)** words in the training corpus with <UNK>.

10/21/15

19

Underflow Prevention: log space

- Multiplying lots of probabilities, which are between 0 and 1, can result in floating-point underflow.
- Since $\log(xy) = \log(x) + \log(y)$, it is better to perform all computations *by summing logs of probabilities rather than multiplying probabilities*.
- Class with highest final un-normalized log probability score is still the most probable.

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

- Note that model is now just max of sum of weights...

10/21/15

Adapt From Manning' textCat tutorial²⁰

Today : Generative Bayes Classifier

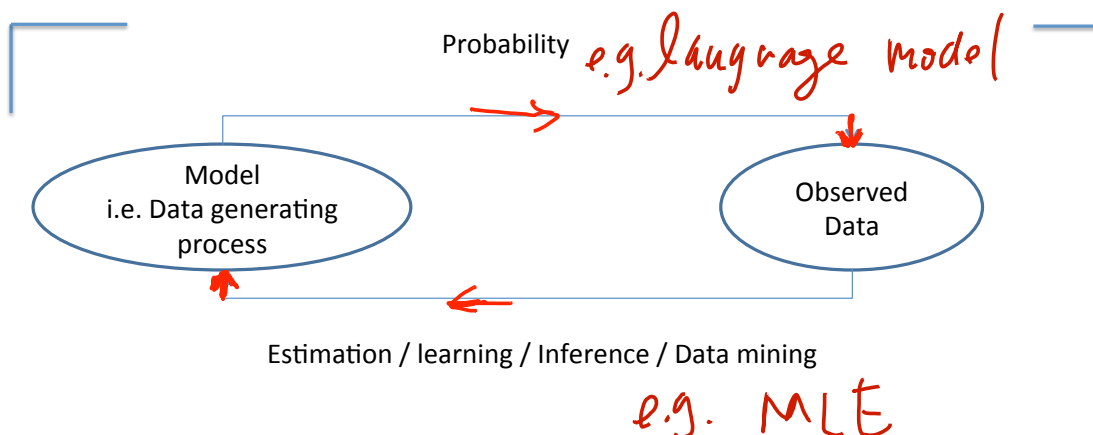
- ✓ Multinomial naïve Bayes classifier as **Conditional Stochastic Language Models**

- A unigram Language model approximates how a text document is produced.

$$\Pr(W_1 = n_1, \dots, W_k = n_k \mid C = c)$$

- ➔ ✓ Maximum Likelihood Estimation of parameters
- ✓ Gaussian Naïve Bayes Classifiers

The Big Picture



But how to specify a model?

Build a *generative model* that approximates how data is produced.

Generative Model & MLE

- Language model can be seen as a probabilistic automata for generating text strings

$$P(W_1 = n_1, \dots, W_k = n_k \mid N, \theta_1, \dots, \theta_k) \propto \theta_1^{n_1} \theta_2^{n_2} \dots \theta_k^{n_k}$$

- Relative frequency estimates can be proven to be *maximum likelihood estimates* (MLE) since they maximize the probability that the model M will generate the training corpus T .

$$\hat{\theta} = \operatorname{argmax}_{\theta} [P(\text{Train} \mid M(\theta))] \quad \text{likelihood}$$

MLE

Maximum Likelihood Estimation

A general Statement

Consider a sample set $T = (X_1, \dots, X_n)$ which is drawn from a probability distribution $P(X \mid \theta)$ where θ are parameters.

If the X s are independent with probability density function $P(X_i \mid \theta)$, the joint probability of the whole set is

$$P(X_1 \dots X_n \mid \theta) = \prod_{i=1}^n P(X_i \mid \theta)$$

this may be maximised with respect to θ to give the maximum likelihood estimates.

$$\hat{\theta} = \operatorname{argmax}_{\theta} [P(\text{Train} \mid M(\theta))] = \operatorname{argmax}_{\theta} P(X_1 \dots X_n \mid \theta)$$

The idea is to

- ✓ assume a particular **model with unknown parameters**, θ
- ✓ we can then define the probability of observing a given event conditional on a particular set of parameters. $P(X_i | \theta)$
- ✓ We have observed **a set of outcomes** in the real world. x_1, x_2, \dots, x_n
- ✓ It is then possible to choose a set of parameters which are most likely to have produced the observed results.

$$\hat{\theta} = \operatorname{argmax}_{\theta} P(X_1 \dots X_n | \theta) \quad \rightarrow \text{likelihood}$$

This is maximum likelihood. In most cases it is **both consistent and efficient**. It provides a standard to compare other estimation techniques.

$$\log(L(\theta)) = \sum_{i=1}^n \left[\log(P(X_i | \theta)) \right] \quad \rightarrow \text{log likelihood}$$

It is often convenient to work with the Log of the likelihood function.

Review: Binomial Distribution e.g. Coin Flips

- You flip n coins
 - How many heads would you expect
 - Head with probability p
 - Number of heads X : **discrete random variable**
- Following **Binomial distribution** with parameters n and p

Defining Likelihood

- Likelihood = $p(\text{data} \mid \text{parameter})$

→ e.g., for n independent tosses of coins, with **unknown** p

Observed data → x heads-up from n trials

10/21/15

Bernoulli
 $X_i \in \{0, 1\}$
 $p(X_1, X_2, \dots, X_n \mid p)$
 $= p^x (1-p)^{n-x}$

function of x_i

PMF: $f(x_i \mid p) = p^{x_i} (1-p)^{1-x_i}$

$$x = \sum_{i=1}^n x_i$$

LIKELIHOOD:

$$L(p) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} = p^x (1-p)^{n-x}$$

function of p

27

Deriving the Maximum Likelihood Estimate

maximize

$$L(p) = p^x (1-p)^{n-x}$$

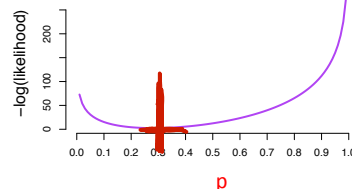
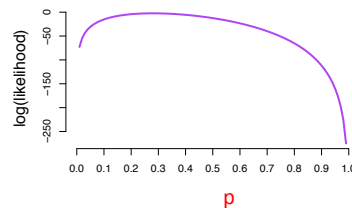
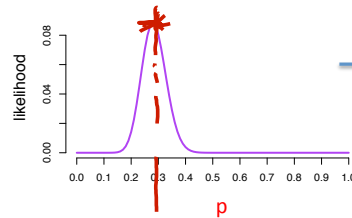
maximize

$$\log(L(p)) = \log[p^x (1-p)^{n-x}]$$

Minimize the **negative log-likelihood**

$$l(p) = -\log[p^x (1-p)^{n-x}]$$

10/21/15



28

Deriving the Maximum Likelihood Estimate

Minimize the negative log-likelihood

$$l(p) = -\log(L(p)) = -\log[p^x(1-p)^{n-x}]$$

$$l(p) = -\log(p^x) - \log((1-p)^{n-x})$$

$$l(p) = -x \log(p) - (n-x) \log(1-p)$$

Deriving the Maximum Likelihood Estimate

$$l(p) = -x \log(p) - (n-x) \log(1-p)$$

$$\frac{dl(p)}{dp} = -\frac{x}{p} - \frac{-(n-x)}{1-p} = 0$$

$$0 = -x + \hat{p}n$$

Minimize the negative log-likelihood

→ MLE parameter estimation

$$\hat{p} = \frac{x}{n}$$

i.e. Relative frequency of a binary event

$$0 = -\frac{x}{\hat{p}} + \frac{n-x}{1-\hat{p}}$$

$$0 = \frac{-x(1-\hat{p}) + \hat{p}(n-x)}{\hat{p}(1-\hat{p})}$$

$$0 = -x + \hat{p}x + \hat{p}n - \hat{p}x$$

Parameter estimation for homework3-Q1

- Multivariate Bernoulli model:

$$\hat{P}(X_w = \text{true} | c_j) = \begin{array}{l} \text{fraction of documents of topic } c_j \\ \text{in which word } w \text{ appears} \end{array}$$

- Multinomial model:

$$\hat{P}(X_i = w | c_j) = \begin{array}{l} \text{fraction of times in which} \\ \text{word } w \text{ appears} \\ \text{across all documents of topic } c_j \end{array}$$

- Can create a mega-document for topic j by concatenating all documents on this topic
- Use frequency of w in mega-document

10/21/15

Adapt From Manning' textCat tutorial ³¹

Deriving the Maximum Likelihood Estimate for multinomial distribution

LIKELIHOOD:

$$\arg \max_{\theta_1, \dots, \theta_k} P(d_1, \dots, d_T | \theta_1, \dots, \theta_k, C=C_j)$$

train T documents
w₁, ..., w_k

$$= \arg \max_{\theta_1, \dots, \theta_k} \prod_{t=1}^T P(d_t | \theta_1, \dots, \theta_k)$$

function of θ

$$= \arg \max_{\theta_1, \dots, \theta_k} \prod_{t=1}^T \frac{N_{d_t}!}{n_{1,d_t}! n_{2,d_t}! \dots n_{k,d_t}!} \theta_1^{n_{1,d_t}} \theta_2^{n_{2,d_t}} \dots \theta_k^{n_{k,d_t}}$$

s.t. $\sum_{i=1}^k \theta_i = 1$

$$= \arg \max_{\theta_1, \dots, \theta_k} \prod_{t=1}^T \theta_1^{n_{1,d_t}} \theta_2^{n_{2,d_t}} \dots \theta_k^{n_{k,d_t}}$$

data likelihood

10/21/15

32

Deriving the Maximum Likelihood Estimate for multinomial distribution

$\arg \max_{\theta_1, \dots, \theta_k} \log(L(\theta))$
Constrained optimization
 $s.t. \sum_{i=1}^k \theta_i = 1$

$= \arg \max_{\theta_1, \dots, \theta_k} \log\left(\prod_{t=1}^T \theta_1^{n_{1,d_t}} \theta_2^{n_{2,d_t}} \dots \theta_k^{n_{k,d_t}}\right)$

$= \arg \max_{\theta_1, \dots, \theta_k} \sum_{t=1, \dots, T} n_{1,d_t} \log(\theta_1) + \sum_{t=1, \dots, T} n_{2,d_t} \log(\theta_2) + \dots + \sum_{t=1, \dots, T} n_{k,d_t} \log(\theta_k)$

Constrained optimization MLE estimator

How optimize ? See Handout - EXTRA

$$\theta_i = \frac{\sum_{t=1, \dots, T} n_{i,d_t}}{\sum_{t=1, \dots, T} n_{1,d_t} + \sum_{t=1, \dots, T} n_{2,d_t} + \dots + \sum_{t=1, \dots, T} n_{k,d_t}} = \frac{\sum_{t=1, \dots, T} n_{i,d_t}}{\sum_{t=1, \dots, T} N_{d_t}}$$

→ i.e. We can create a mega-document by concatenating all documents d_1 to d_T
 → Use relative frequency of w in mega-document

for c_j

Naïve Bayes: Learning Algorithm for parameter estimation with MLE

- From training corpus, extract *Vocabulary*
- Calculate required $P(c_j)$ and $P(w_k | c_j)$ terms
 - For each c_j in C do
 - $docs_j \leftarrow$ subset of documents for which the target class is c_j

$$P(c_j) \leftarrow \frac{|docs_j|}{|\text{total \# documents}|}$$

- $Text_j \leftarrow$ is length n and is a single document containing all $docs_j$
- for each word w_k in *Vocabulary*
 - $n_k \leftarrow$ number of occurrences of w_k in $Text_j$; n is length of $Text_j$

$$P(w_k | c_j) \leftarrow \frac{n_k + \alpha}{n + \alpha |Vocabulary|} \quad \text{e.g., } \alpha = 1$$

Relative frequency of word w_k appears across all documents of class c_j

Naive Bayes: Time Complexity

- **Training Time:** $O(|D|L_d + |C||V|)$

where L_d is the average length of a document in D .

- Assumes V and all D_i , n_i , and n_{ij} pre-computed in $O(|D|L_d)$ time during one pass through all of the data.
- Generally just $O(|D|L_d)$ since usually $|C||V| < |D|L_d$

- **Test Time:** $O(|C| L_t)$

where L_t is the average length of a test document.

- **Very efficient overall**, linearly proportional to the time needed to just read in all the data.
- Plus, **robust** in practice

$|D|$ num. doc
 $|V|$ dict size
 $|C|$ class size

Naive Bayes is Not So Naive

- **Naïve Bayes: First and Second place in KDD-CUP 97 competition, among 16 (then) state of the art algorithms**

Goal: Financial services industry direct mail response prediction model: Predict if the recipient of mail will actually respond to the advertisement – 750,000 records.

- **Robust to Irrelevant Features**
Irrelevant Features cancel each other without affecting results
Instead Decision Trees can **heavily** suffer from this.
- **Very good in domains with many equally important features**
Decision Trees suffer from *fragmentation* in such cases – especially if little data
- **A good dependable baseline for text classification (but not the best)!**
- **Optimal if the Independence Assumptions hold:** If assumed independence is correct, then it is the Bayes Optimal Classifier for problem
- **Very Fast:** Learning with one pass of counting over the data; testing linear in the number of attributes, and document collection size
- **Low Storage requirements**

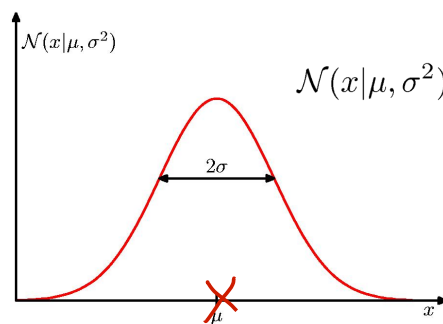
Today : Generative vs. Discriminative

- ✓ Multinomial naïve Bayes classifier as Stochastic Language Models
 - ✓ a unigram Language model approximates how a text document is produced.
- ✓ Maximum Likelihood Estimation of parameters
- ➔ ✓ Gaussian Naïve Bayes Classifiers
 - Gaussian distribution
 - Gaussian NBC
 - LDA, QDA

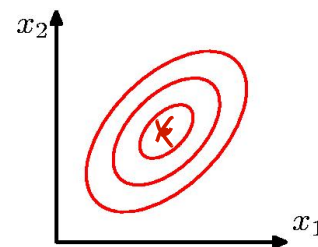
10/21/15

37

The Gaussian Distribution



$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\}$$



$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}$$

Mean
Covariance Matrix

10/21/15

Courtesy: <http://research.microsoft.com/~cmbishop/PRML/index.htm>

38

Multivariate Gaussian Distribution

- A multivariate Gaussian model: $\mathbf{x} \sim N(\boldsymbol{\mu}, \Sigma)$ where

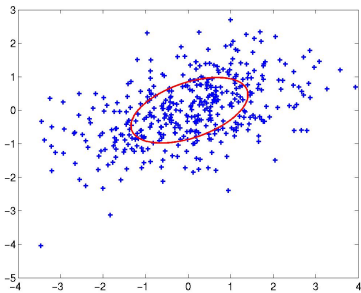
Here $\boldsymbol{\mu}$ is the mean vector and Σ is the covariance matrix,
if $p=2$

$$\boldsymbol{\mu} = \{\mu_1, \mu_2\}$$

$$\Sigma = \begin{bmatrix} \text{var}(x_1) & \text{cov}(x_1, x_2) \\ \text{cov}(x_1, x_2) & \text{var}(x_2) \end{bmatrix}$$

- The covariance matrix captures linear dependencies among the variables

MLE Estimation for Multivariate Gaussian



- We can fit statistical models by maximizing the probability / likelihood of generating the observed samples:

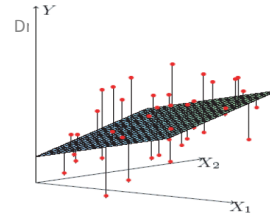
$$L(x_1, \dots, x_n | \theta) = p(x_1 | \theta) \dots p(x_n | \theta)$$

(the samples are assumed to be IID)

- In the Gaussian case, we simply set the mean and the variance to the sample mean and the sample variance:

$$\bar{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad \bar{\boldsymbol{\sigma}}^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\boldsymbol{\mu}})^2$$

Probabilistic Interpretation of Linear Regression



- Let us assume that the target variable and the inputs are related by the equation:

$$y_i = \theta^T \mathbf{x}_i + \varepsilon_i$$

where ε is an error term of unmodeled effects or random noise

- Now assume that ε follows a Gaussian $N(0, \sigma)$, then we have:

$$p(y_i | x_i; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - \theta^T \mathbf{x}_i)^2}{2\sigma^2}\right)$$

- By IID assumption:

$$L(\theta) = \prod_{i=1}^n p(y_i | x_i; \theta) = \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n \exp\left(-\frac{\sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2}{2\sigma^2}\right)$$

Other prob. models.

10/21/15

41

Dr. Yanjun Qi / UVA CS 6316 / f15

Probabilistic Interpretation of Linear Regression (cont.)

- Hence the log-likelihood is:

$$l(\theta) = n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \frac{1}{2} \sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2$$

- Do you recognize the last term?

Yes it is:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \theta - y_i)^2$$

- Thus under independence assumption, residual square error (**RRS**) is equivalent to **MLE** of ϑ !

10/21/15

42

Today : Generative vs. Discriminative

- ✓ Multinomial naïve Bayes classifier as Stochastic Language Models
 - ✓ a unigram Language model approximates how a text document is produced.
- ✓ Maximum Likelihood Estimation of parameters
- ✓ Gaussian Naïve Bayes Classifiers
 - Gaussian distribution
 - Gaussian NBC
 - LDA, QDA

10/21/15

43

Gaussian Naïve Bayes Classifier

$$\operatorname{argmax}_C P(C | X) = \operatorname{argmax}_C P(X, C) = \operatorname{argmax}_C P(X | C)P(C)$$

Naïve
Bayes
Classifier

$$\begin{aligned} P(X | C) &= P(X_1, X_2, \dots, X_p | C) \\ &= P(X_1 | X_2, \dots, X_p, C)P(X_2, \dots, X_p | C) \\ &= P(X_1 | C)P(X_2, \dots, X_p | C) \\ &= P(X_1 | C)P(X_2 | C) \cdots P(X_p | C) \end{aligned}$$

$$\hat{P}(X_j | C = c_i) = \frac{1}{\sqrt{2\pi}\sigma_{ji}} \exp\left(-\frac{(X_j - \mu_{ji})^2}{2\sigma_{ji}^2}\right)$$

μ_{ji} : mean (average) of attribute values X_j of examples for which $C = c_i$
 σ_{ji} : standard deviation of attribute values X_j of examples for which $C = c_i$

10/21/15

Gaussian Naïve Bayes Classifier

- Continuous-valued Input Attributes

- Conditional probability modeled with the normal distribution

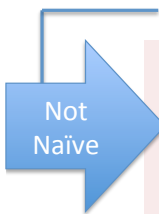
$$\hat{P}(X_j | C = c_i) = \frac{1}{\sqrt{2\pi}\sigma_{ji}} \exp\left(-\frac{(X_j - \mu_{ji})^2}{2\sigma_{ji}^2}\right)$$

μ_{ji} : mean (average) of attribute values X_j of examples for which $C = c_i$

σ_{ji} : standard deviation of attribute values X_j of examples for which $C = c_i$

- **Learning Phase:** for $\mathbf{X} = (X_1, \dots, X_p)$, $C = c_1, \dots, c_L$
Output: $p \times L$ normal distributions and $P(C = c_i)$ $i = 1, \dots, L$
- **Test Phase:** for $\mathbf{X}' = (X'_1, \dots, X'_p)$
 - Calculate conditional probabilities with all the normal distributions
 - Apply the MAP rule to make a decision

Naïve Gaussian means ?



$$P(X_1, X_2, \dots, X_p | C) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}$$

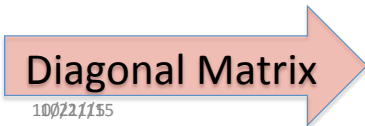


P-dim random vector

$$P(X_1, X_2, \dots, X_p | C = c_j) = P(X_1 | C)P(X_2 | C) \cdots P(X_p | C)$$

$$= \prod_i \frac{1}{\sqrt{2\pi}\sigma_{ji}} \exp\left(-\frac{(X_j - \mu_{ji})^2}{2\sigma_{ji}^2}\right)$$

$O(kp + kp)$

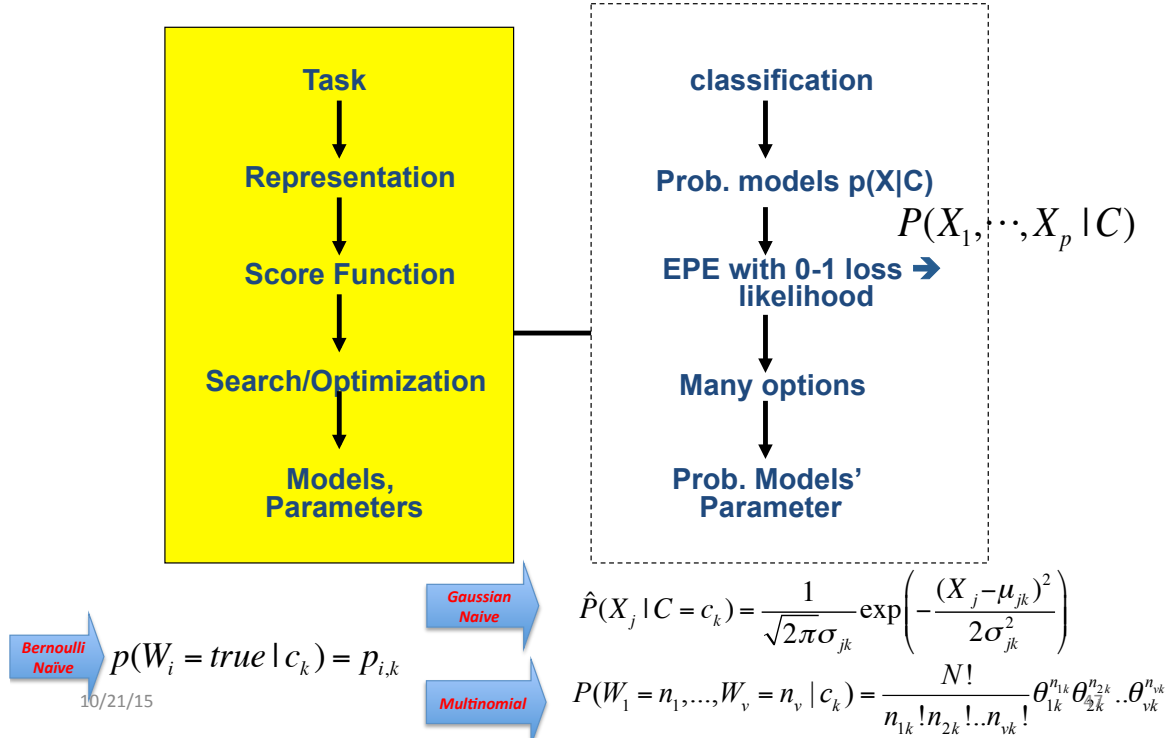


$$\boldsymbol{\Sigma}_{c_k} = \boldsymbol{\Lambda}_{c_k}$$

Each class' covariance matrix is diagonal

$$\operatorname{argmax}_k P(C = k | X) = \operatorname{argmax}_k P(X, C) = \operatorname{argmax}_k P(X|C)P(C)$$

Generative Bayes Classifier



Today : Generative vs. Discriminative

- ✓ Multinomial naïve Bayes classifier as Stochastic Language Models
 - ✓ a unigram Language model approximates how a text document is produced.
 - ✓ Maximum Likelihood Estimation of parameters
- ➔ ✓ Gaussian Naïve Bayes Classifiers
 - Gaussian distribution
 - Gaussian NBC
 - Not-naïve Gaussian BC \rightarrow LDA, QDA

(1) covariance matrix are the same across classes

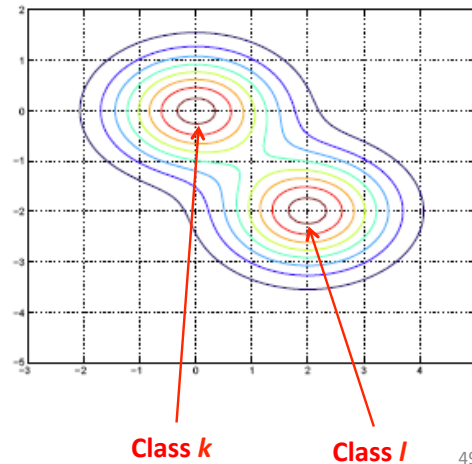
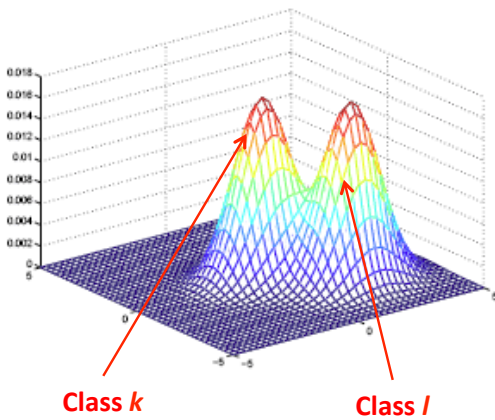
→ LDA (Linear Discriminant Analysis)

$$O(p^2 + kp)$$

Linear Discriminant Analysis : $\Sigma_k = \Sigma, \forall k$

Each class' covariance matrix is the same

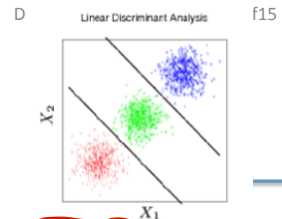
The Gaussian Distribution are shifted versions of each other



10/21/15

49

Optimal Classification



$$\operatorname{argmax}_k P(C = k | X) = \operatorname{argmax}_k P(X, C) = \operatorname{argmax}_k P(X | C) P(C)$$

$$= \operatorname{argmax}_k \left[-\log\left(\frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}}\right) - \frac{1}{2}(x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log(\pi_k) \right]$$

$$= \operatorname{argmax}_k \left[-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log(\pi_k) \right]$$

- Note

Linear Discriminant Function for LDA

$$-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1} (x - \mu_k) = \left[x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k - \frac{1}{2} x^T \Sigma^{-1} x \right]$$

10/21/15

50

$$\operatorname{argmax}_k P(C_k | X) = \operatorname{argmax}_k P(X, C_k) = \operatorname{argmax}_k P(X | C_k) P(C_k)$$

$$= \operatorname{argmax}_k \log \{ \underbrace{P(X | C_k)} \underbrace{P(C_k)} \}$$

$$= \operatorname{argmax}_k \log P(X | C_k) + \log P(C_k) \Rightarrow \pi_k$$

Decision Boundary points,

$$\begin{aligned} \log \frac{P(C_k | X)}{P(C_l | X)} &= 0 = \log \frac{P(X | C_k)}{P(X | C_l)} + \log \frac{\pi_k}{\pi_l} \\ &= \log P(X | C_k) - \log P(X | C_l) + \log \frac{\pi_k}{\pi_l} \end{aligned}$$

Define **Linear Discriminant Function**

$$\delta_k(x) = -\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log \pi_k$$

$\delta_k(x) = \delta_g(x)$

→ The **Decision Boundary Between class k and l**, $\{x : \delta_k(x) = \delta_l(x)\}$, is linear

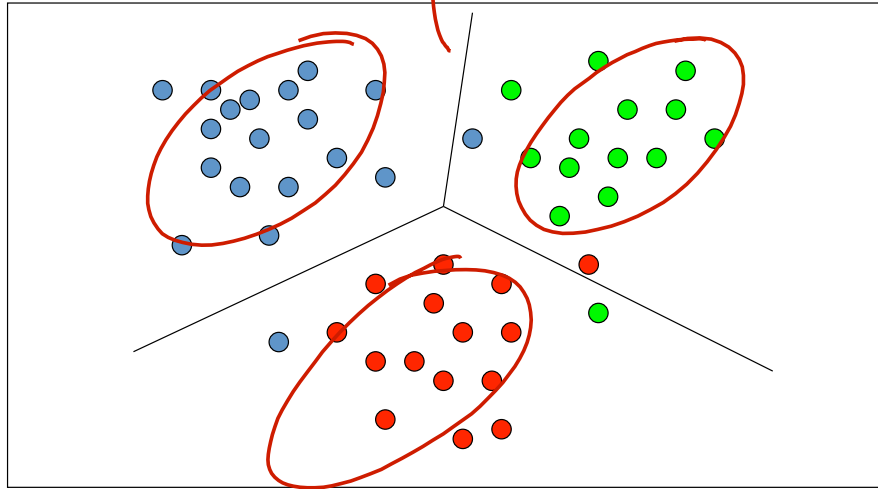
$$\log \frac{P(C_k | X)}{P(C_l | X)} = \log \frac{P(X | C_k)}{P(X | C_l)} + \log \frac{P(C_k)}{P(C_l)}$$

$$= \log \frac{\pi_k}{\pi_l} - \frac{1}{2} (\mu_k + \mu_l)^T \Sigma^{-1} (\mu_k - \mu_l) + x^T \Sigma^{-1} (\mu_k - \mu_l) \quad (4.9)$$

Equals to zero

Boundary points X : when $P(c_k|X) = P(c_l|X)$, the left linear equation $= 0$, a linear line / plane

Visualization (three classes)



10/21/15

53

(2) If covariance matrix are not same

e.g. → **QDA (Quadratic Discriminant Analysis)**

$O(KP^2 + KP)$

- ▶ Estimate the covariance matrix Σ_k separately for each class k , $k = 1, 2, \dots, K$.

- ▶ Quadratic discriminant function:

$$\delta_k(x) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log \pi_k .$$

- ▶ Classification rule:

$$\hat{G}(x) = \arg \max_k \delta_k(x) .$$

$$\log \frac{P(q|x)}{P(r|x)} = 0$$

- ▶ Decision boundaries are quadratic equations in x .
- ▶ QDA fits the data better than LDA, but has more parameters to estimate.

10/21/15

54

LDA on Expanded Basis

- ▶ Expand input space to include X_1X_2 , X_1^2 , and X_2^2 .
- ▶ Input is five dimensional: $X = (X_1, X_2, X_1X_2, X_1^2, X_2^2)$.

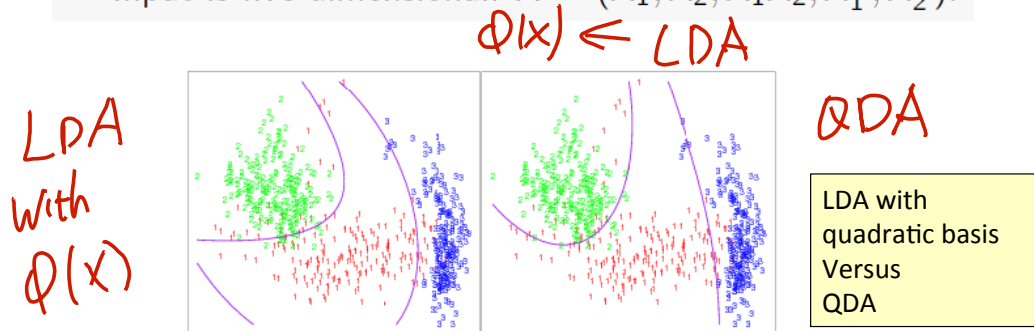


Figure 4.6: Two methods for fitting quadratic boundaries. The left plot shows the quadratic decision boundaries for the data in Figure 4.1 (obtained using LDA in the five-dimensional space $x_1, x_2, x_{12}, x_1^2, x_2^2$). The right plot shows the quadratic decision boundaries found by QDA. The differences are small, as is usually the case.

10/21/15

55

(3) Regularized Discriminant Analysis

- ▶ A compromise between LDA and QDA.
- ▶ Shrink the separate covariances of QDA toward a common covariance as in LDA.
- ▶ Regularized covariance matrices:

$$\hat{\Sigma}_k(\alpha) = \alpha \hat{\Sigma}_k + (1 - \alpha) \hat{\Sigma}.$$

- ▶ The quadratic discriminant function $\delta_k(x)$ is defined using the shrunken covariance matrices $\hat{\Sigma}_k(\alpha)$.
- ▶ The parameter α controls the complexity of the model.

10/21/15

56

References

- Prof. Andrew Moore's review tutorial
- Prof. Ke Chen NB slides
- Prof. Carlos Guestrin recitation slides
- Prof. Raymond J. Mooney and Jimmy Lin's slides about language model