# CS60020: Foundations of Algorithm Design and Machine Learning

Sourangshu Bhattacharya

# NAÏVE BAYES

# Generative vs. Discriminative Classifiers

Discriminative classifiers (e.g. Logistic Regression)

- Assume some functional form for P(Y|X) or for the decision boundary
- Estimate parameters of P(Y|X) directly from training data

Generative classifiers (e.g. Naïve Bayes)

- Assume some functional form for P(X,Y) (or P(X|Y) and P(Y))
- Estimate parameters of P(X|Y), P(Y) directly from training data

arg max_Y P(Y|X) = arg max_Y P(X|Y) P(Y)

# A text classification task: Email spam filtering

```
From: ''' <takworlld@hotmail.com>
Subject: real estate is the only way... gem oalvgkay
Anyone can buy real estate with no money down
Stop paying rent TODAY !
There is no need to spend hundreds or even thousands for
similar courses
I am 22 years old and I have already purchased 6 properties
using the
methods outlined in this truly INCREDIBLE ebook.
Change your life NOW !
=================================================
Click Below to order:
http://www.wholesaledaily.com/sales/nmd.htm
=================================================
```

How would you write a program that would automatically detect and delete this type of message?

4

# Formal definition of TC: Training

Given:

- A document set X

    - Documents are represented typically in some type of high-dimensional space.

- A fixed set of classes C = $\{c_1, c_2, . . . , c_J\}$

    - The classes are human-defined for the needs of an application (e.g., relevant vs. nonrelevant).

- A training set D of labeled documents with each labeled document $<d, c> \in X \times C$

Using a learning method or learning algorithm, we then wish to

learn a classifier $\Upsilon$ that maps documents to classes:

$$\Upsilon : X \rightarrow C$$

# Formal definition of TC: Application/Testing

Given: a description $d \in X$ of a document Determine: $\Upsilon(d) \in C$, that is, the class that is most appropriate for $d$

# Examples of how search engines use classification

- Language identification (classes: English vs. French etc.)
- The automatic detection of spam pages (spam vs. nonspam)
- Topic-specific or *vertical* search – restrict search to a "vertical" like "related to health" (relevant to vertical vs. not)

# Derivation of Naive Bayes rule

We want to find the class that is most likely given the document:

$$c_{map} = \underset{c \in \mathbb{C}}{\arg\max}\ P(c|d)$$

Apply Bayes rule

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

$$c_{map} = \underset{c \in \mathbb{C}}{\arg\max}\ \frac{P(d|c)P(c)}{P(d)}$$

Drop denominator since P(d) is the same for all classes:

$$c_{map} = \underset{c \in \mathbb{C}}{\arg\max}\ P(d|c)P(c)$$

# Too many parameters / sparseness

$$c_{map} = \arg\max_{c \in \mathbb{C}} P(d|c)P(c)$$

$$= \arg\max_{c \in \mathbb{C}} P(\langle t_1, \ldots, t_k, \ldots, t_{n_d}\rangle | c)P(c)$$

▪There are too many parameters $P(\langle t_1, \ldots, t_k, \ldots, t_{n_d}\rangle | c)$ , one for each unique combination of a class and a sequence of words.

▪We would need a very, very large number of training examples to estimate that many parameters.

▪This is the problem of data sparseness.

9

# Naive Bayes conditional independence assumption

To reduce the number of parameters to a manageable size, we make the Naive Bayes conditional independence assumption:

$$P(d|c) = P(\langle t_1, \ldots, t_{n_d} \rangle | c) = \prod_{1 \leq k \leq n_d} P(X_k = t_k | c)$$

We assume that the probability of observing the conjunction of attributes is equal to the product of the individual probabilities $P(X_k = t_k | c)$.

# The Naive Bayes classifier

- The Naive Bayes classifier is a probabilistic classifier.

- We compute the probability of a document $d$ being in a class $c$ as follows:

$$P(c|d) \propto P(c) \prod_{1 \le k \le n_d} P(t_k|c)$$

- $n_d$ is the length of the document. (number of tokens)

- $P(t_k|c)$ is the conditional probability of term $t_k$ occurring in a document of class $c$

- $P(t_k|c)$ is a measure of how much evidence $t_k$ contributes that $c$ is the correct class.

- $P(c)$ is the prior probability of $c$.

- If a document's terms do not provide clear evidence for one class vs. another, we choose the c with highest $P(c)$.

# Maximum a posteriori class

- Our goal in Naive Bayes classification is to find the "best" class.

- The best class is the most likely or maximum a posteriori (MAP) class $c_{map}$:

$$c_{\mathrm{map}} = \arg\max_{c \in \mathbb{C}} \hat{P}(c|d) = \arg\max_{c \in \mathbb{C}} \hat{P}(c) \prod_{1 \le k \le n_d} \hat{P}(t_k|c)$$

# Taking the log

- Multiplying lots of small probabilities can result in floating point underflow.

- Since log(*xy*) = log(*x*) + log(*y*), we can sum log probabilities instead of multiplying probabilities.

- Since log is a monotonic function, the class with the highest score does not change.

- So what we usually compute in practice is:

$$c_{\text{map}} = \arg\max_{c \in \mathbb{C}} \left[ \log \hat{P}(c) + \sum_{1 \leq k \leq n_d} \log \hat{P}(t_k|c) \right]$$

# Naive Bayes classifier

▪Classification rule:

$$c_{\mathrm{map}} = \arg\max_{c \in \mathbb{C}} \left[ \log \hat{P}(c) + \sum_{1 \leq k \leq n_d} \log \hat{P}(t_k|c) \right]$$

▪Simple interpretation:

  ▪Each conditional parameter log $\hat{P}(t_k|c)$ is a weight that indicates how good an indicator $t_k$ is for $c$.

  ▪The prior log $\hat{P}(c)$ is a weight that indicates the relative frequency of $c$.

  ▪The sum of log prior and term weights is then a measure of how much evidence there is for the document being in the class.

  ▪We select the class with the most evidence.

# Parameter estimation take 1: Maximum likelihood

- Estimate parameters $\hat{P}(c)$ and $\hat{P}(t_k|c)$ from train data: How?

- Prior:

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

- $N_c$ : number of docs in class $c$; $N$: total number of docs

- Conditional probabilities:

$$\hat{P}(t|c) = \frac{T_{ct}}{\sum_{t' \in V} T_{ct'}}$$

- $T_{ct}$ is the number of tokens of $t$ in training documents from class $c$ (includes multiple occurrences)

- We've made a Naive Bayes independence assumption here:

# The problem with maximum likelihood estimates: Zeros



$P(China|d) \propto P(China) \cdot P(\text{BEIJING}|China) \cdot P(\text{AND}|China)$
$\cdot P(\text{TAIPEI}|China) \cdot P(\text{JOIN}|China) \cdot P(\text{WTO}|China)$

- If WTO never occurs in class China in the train set:

$$\hat{P}(\text{WTO}|China) = \frac{T_{China,\text{WTO}}}{\sum_{t' \in V} T_{China,t'}} = \frac{0}{\sum_{t' \in V} T_{China,t'}} = 0$$

16

# The problem with maximum likelihood estimates: Zeros (cont)

▪If there were no occurrences of WTO in documents in class China, we'd get a zero estimate:

$$\hat{P}(\text{WTO}|China) = \frac{T_{China,\text{WTO}}}{\sum_{t' \in V} T_{China,t'}} = 0$$

▪→ We will get P(China|d) = 0 for any document that contains WTO!

▪Zero probabilities cannot be conditioned away.

# To avoid zeros: Add-one smoothing

▪Before:

$$\hat{P}(t|c) = \frac{T_{ct}}{\sum_{t' \in V} T_{ct'}}$$

▪Now: Add one to each count to avoid zeros:

▪B is the number of different words (in this case the size of the vocabulary: $|V| = B$)

$$\hat{P}(t|c) = \frac{T_{ct} + 1}{\sum_{t' \in V}(T_{ct'} + 1)} = \frac{T_{ct} + 1}{(\sum_{t' \in V} T_{ct'}) + B}$$

# To avoid zeros: Add-one smoothing

- Estimate parameters from the training corpus using add-one smoothing

- For a new document, for each class, compute sum of (i) log of prior and (ii) logs of conditional probabilities of the terms

- Assign the document to the class with the largest score

# Exercise

|  | docID | words in document | in $c = China$? |
|---|---|---|---|
| training set | 1 | Chinese Beijing Chinese | yes |
|  | 2 | Chinese Chinese Shanghai | yes |
|  | 3 | Chinese Macao | yes |
|  | 4 | Tokyo Japan Chinese | no |
| test set | 5 | Chinese Chinese Chinese Tokyo Japan | ? |

- Estimate parameters of Naive Bayes classifier

- Classify test document

# Example: Parameter estimates

Priors: $\hat{P}(c) = 3/4$ and $\hat{P}(\bar{c}) = 1/4$ Conditional probabilities:

$$\hat{P}(\text{CHINESE}|c) = (5+1)/(8+6) = 6/14 = 3/7$$
$$\hat{P}(\text{TOKYO}|c) = \hat{P}(\text{JAPAN}|c) = (0+1)/(8+6) = 1/14$$
$$\hat{P}(\text{CHINESE}|\bar{c}) = (1+1)/(3+6) = 2/9$$
$$\hat{P}(\text{TOKYO}|\bar{c}) = \hat{P}(\text{JAPAN}|\bar{c}) = (1+1)/(3+6) = 2/9$$

The denominators are (8 + 6) and (3 + 6) because the lengths of

$text_c$ and $text_{\bar{c}}$ are 8 and 3, respectively, and because the constant

$B$ is 6 as the vocabulary consists of six terms.

# Example: Classification

$$\hat{P}(c|d_5) \propto 3/4 \cdot (3/7)^3 \cdot 1/14 \cdot 1/14 \approx 0.0003$$

$$\hat{P}(\overline{c}|d_5) \propto 1/4 \cdot (2/9)^3 \cdot 2/9 \cdot 2/9 \approx 0.0001$$

Thus, the classifier assigns the test document to $c$ = *China*. The

reason for this classification decision is that the three occurrences

of the positive indicator $\mathrm{CHINESE}$ in $d_5$ outweigh the occurrences

of the two negative indicators $\mathrm{JAPAN}$ and $\mathrm{TOKYO}$.

# Class Conditional Probabilities

To compute, $P(x_k|C_i)$

- $A_k$ is categorical:

$$P(x_k|C_i) = \frac{\text{the number of tuples of class } C_i \text{ in D having the value } x_k \text{ for } A_k}{\text{the number of tuples of class } C_i \text{ in D.}}$$

- $A_k$ is continuous:

A continuous-valued attribute is typically assumed to have a Gaussian distribution with a mean $\mu$ and standard deviation $\sigma$

$$g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma^2}},$$

$$P(x_k|C_i) = g(x_k, \mu_{C_i}, \sigma_{C_i}).$$

# Generative model



$$P(c|d) \propto P(c) \prod_{1 \le k \le n_d} P(t_k|c)$$

- Generate a class with probability $P(c)$

- Generate each of the words (in their respective positions), conditional on the class, but independent of each other, with probability $P(t_k|c)$

- To classify docs, we "reengineer" this process and find the class that is most likely to have generated the doc.

# On naïve Bayesian classifier

- Advantages:
  - Easy to implement
  - Very efficient
  - Good results obtained in many applications
- Disadvantages
  - Assumption: class conditional independence, therefore loss of accuracy when the assumption is seriously violated (those highly correlated data sets)

# BAYESIAN LINEAR REGRESSION

# Maximum Likelihood and Least Squares

- Assume observations from a deterministic function with added Gaussian noise:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon \qquad \text{where} \qquad p(\epsilon|\beta) = \mathcal{N}(\epsilon|0, \beta^{-1})$$

- which is the same as saying,

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1}).$$

- Given observed inputs, $\mathbf{X} = \{\mathbf{x}_1, \ldots, \mathbf{x}_N\}$, and targets, $\mathbf{t} = [t_1, \ldots, t_N]^{\mathrm{T}}$, we obtain the likelihood function

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^{N} \mathcal{N}(t_n|\mathbf{w}^{\mathrm{T}}\boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}).$$

# Maximum Likelihood and Least Squares

- Taking the logarithm, we get

$$
\begin{aligned}
\ln p(\mathbf{t}|\mathbf{w}, \beta) &= \sum_{n=1}^{N} \ln \mathcal{N}(t_n|\mathbf{w}^{\mathrm{T}}\boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}) \\
&= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w})
\end{aligned}
$$

- where

$$
E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{t_n - \mathbf{w}^{\mathrm{T}}\boldsymbol{\phi}(\mathbf{x}_n)\}^2
$$

- is the sum-of-squares error.

# Bayesian Linear Regression (1)

- Define a conjugate prior over w

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0).$$

•Combining this with the likelihood function and using results for marginal and conditional Gaussian distributions, gives the posterior

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N)$$

- where

$$
\begin{aligned}
\mathbf{m}_N &= \mathbf{S}_N \left( \mathbf{S}_0^{-1}\mathbf{m}_0 + \beta\boldsymbol{\Phi}^{\mathrm{T}}\mathbf{t} \right) \\
\mathbf{S}_N^{-1} &= \mathbf{S}_0^{-1} + \beta\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi}.
\end{aligned}
$$

# Bayesian Linear Regression (2)

- A common choice for the prior is

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})$$

•for which

$$\mathbf{m}_N = \beta \mathbf{S}_N \boldsymbol{\Phi}^{\mathrm{T}} \mathbf{t}$$
$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \boldsymbol{\Phi}^{\mathrm{T}} \boldsymbol{\Phi}.$$

•Next we consider an example …

# Bayesian Linear Regression (3)

0 data points observed

# Bayesian Linear Regression (4)

1 data point observed



Likelihood

Posterior

Data Space

# Bayesian Linear Regression (5)

2 data points observed

# Bayesian Linear Regression (6)

20 data points observed



Likelihood      Posterior      Data Space

# Predictive Distribution (1)

- Predict t for new values of x by integrating over w:

$$
\begin{aligned}
p(t|\mathbf{t}, \alpha, \beta) &= \int p(t|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) \, \mathrm{d}\mathbf{w} \\
&= \mathcal{N}(t|\mathbf{m}_N^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}), \sigma_N^2(\mathbf{x}))
\end{aligned}
$$

- where

$$
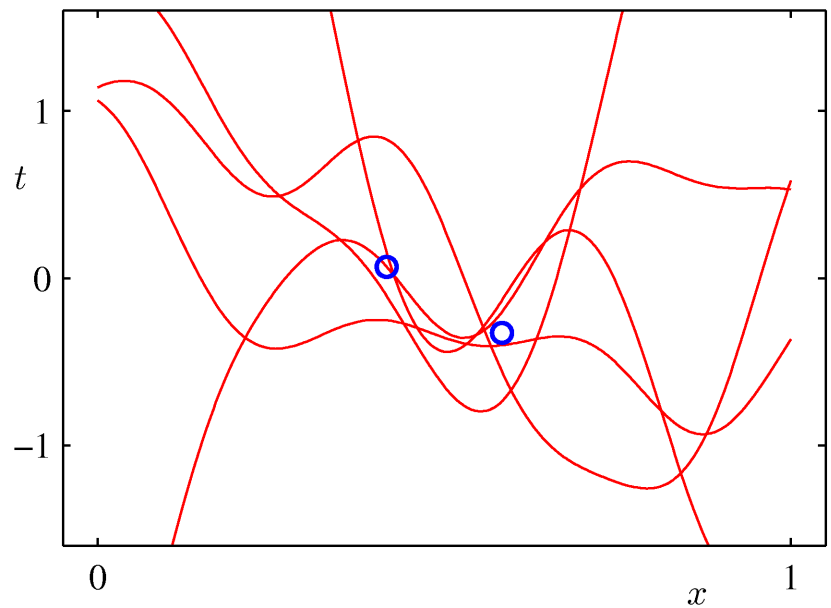\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \boldsymbol{\phi}(\mathbf{x})^{\mathrm{T}} \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x}).
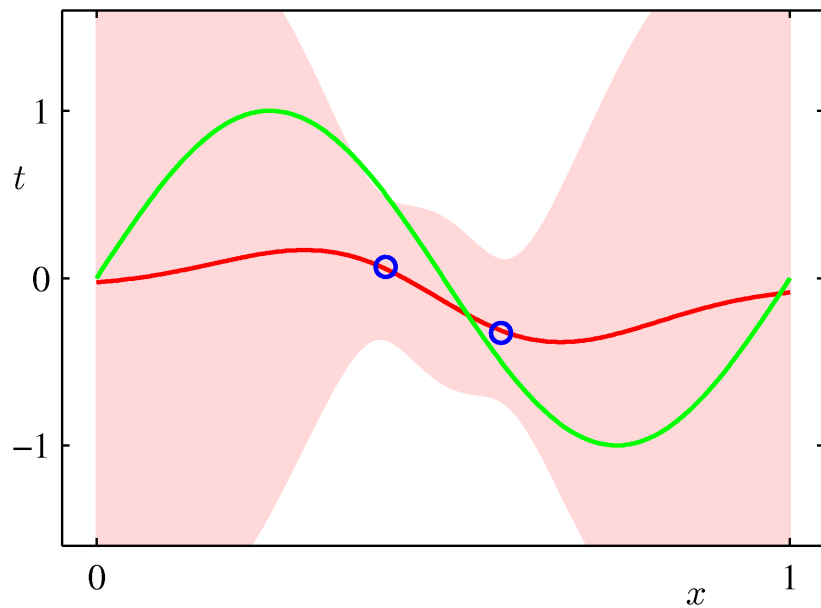$$

# Predictive Distribution (2)

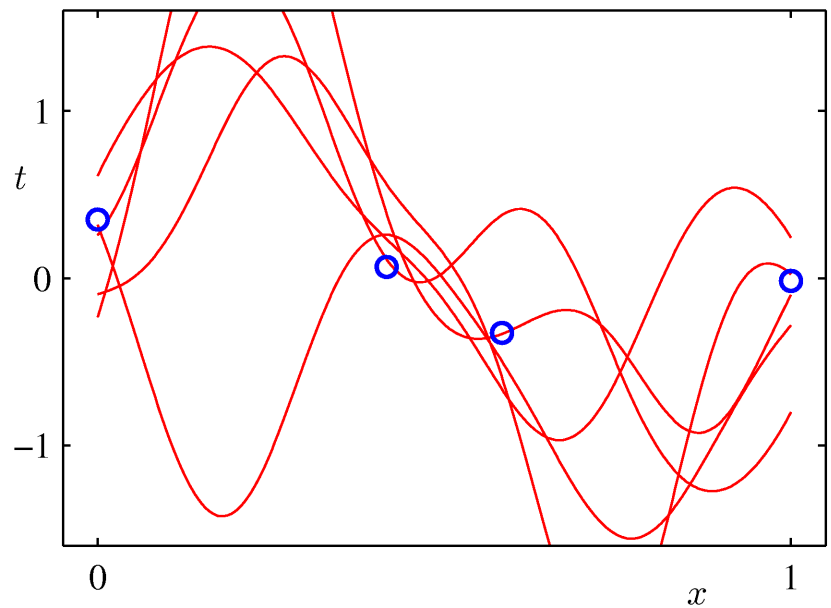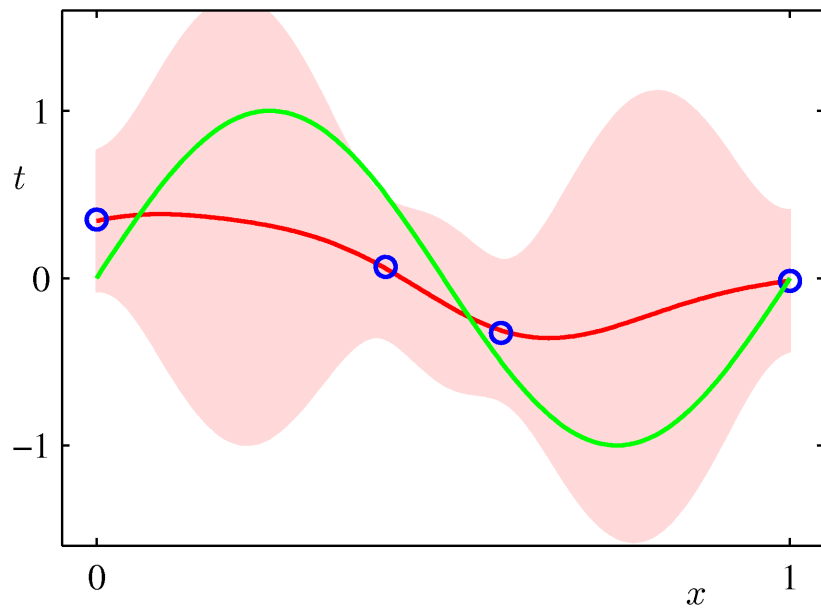- Example: Sinusoidal data, 9 Gaussian basis functions, 1 data point

# Predictive Distribution (3)

- Example: Sinusoidal data, 9 Gaussian basis functions, 2 data points

# Predictive Distribution (4)

- Example: Sinusoidal data, 9 Gaussian basis functions, 4 data points

# Predictive Distribution (5)

- Example: Sinusoidal data, 9 Gaussian basis functions, 25 data points