

10-701  
**Machine Learning**

Naïve Bayes classifiers

# Types of classifiers

- We can divide the large variety of classification approaches into three major types
  1. Instance based classifiers
    - Use observation directly (no models)
    - e.g. K nearest neighbors
  2. Generative:
    - build a generative statistical model
    - e.g., Bayesian networks
  3. Discriminative
    - directly estimate a decision rule/boundary
    - e.g., decision tree

# Bayes decision rule

- If we know the conditional probability  $P(X | y)$  we can determine the appropriate class by using Bayes rule:

$$P(y = i | X) = \frac{P(X | y = i)P(y = i)}{P(X)} \stackrel{def}{=} q_i(X)$$

But how do we determine  $p(X|y)$ ?

# Computing $p(X|y)$

Recall...

$y$  – the class label

$X$  – input attributes  
(features)

- Consider a dataset with 16 attributes (lets assume they are all binary). How many parameters do we need to estimate to fully determine  $p(X|y)$ ?

age	employment	education	edun	marital	...	job	relation	race	gender	hour	country	wealth
39	State_gov	Bachelors	13	Never_mar	...	Adm_cleric	Not_in_fan	White	Male	40	United_Stat	poor
51	Self_emp	Bachelors	13	Married	...	Exec_man	Husband	White	Male	13	United_Stat	poor
39	Private	HS_grad	9	Divorced	...	Handlers_c	Not_in_fan	White	Male	40	United_Stat	poor
54	Private	11th	7	Married	...	Handlers_c	Husband	Black	Male	40	United_Stat	poor
28	Private	Bachelors	13	Married	...	Prof_speci	Wife	Black	Female	40	Cuba	poor
38	Private	Masters	14	Married	...	Exec_man	Wife	White	Female	40	United_Stat	poor
50	Private	9th	5	Married_sp	...	Other_serv	Not_in_fan	Black	Female	16	Jamaica	poor
52	Self_emp	HS_grad	9	Married	...	Exec_man	Husband	White	Male	45	United_Stat	rich
31	Private	Masters	14	Never_mar	...	Prof_speci	Not_in_fan	White	Female	50	United_Stat	rich
42	Private	Bachelors	13	Married	...	Exec_man	Husband	White	Male	40	United_Stat	rich
37	Private	Some_coll	10	Married	...	Exec_man	Husband	Black	Male	80	United_Stat	rich
30	State_gov	Bachelors	13	Married	...	Prof_speci	Husband	Asian	Male	40	India	rich
24	Private	Bachelors	13	Never_mar	...	Adm_cleric	Own_child	White	Female	30	United_Stat	poor
33	Private	Assoc_acc	12	Never_mar	...	Sales	Not_in_fan	Black	Male	50	United_Stat	poor
41	Private	Assoc_voc	11	Married	...	Craft_repair	Husband	Asian	Male	40	*MissingV	rich
34	Private	7th_8th	4	Married	...	Transport	Husband	Amer_Indic	Male	45	Mexico	poor
26	Self_emp	HS_grad	9	Never_mar	...	Farming_fi	Own_child	White	Male	35	United_Stat	poor
33	Private	HS_grad	9	Never_mar	...	Machine_c	Unmarried	White	Male	40	United_Stat	poor
38	Private	11th	7	Married	...	Sales	Husband	White	Male	50	United_Stat	poor
44	Self_emp	Masters	14	Divorced	...	Exec_man	Unmarried	White	Female	45	United_Stat	rich
41	Private	Doctorate	16	Married	...	Prof_speci	Husband	White	Male	60	United_Stat	rich

Learning the values for the full conditional probability table would require enormous amounts of data

# Naïve Bayes Classifier

- Naïve Bayes classifiers assume that given the class label (Y) the attributes are **conditionally independent** of each other:

$$p(X | y) = \prod_j p_j(x^j | y)$$

Product of probability terms

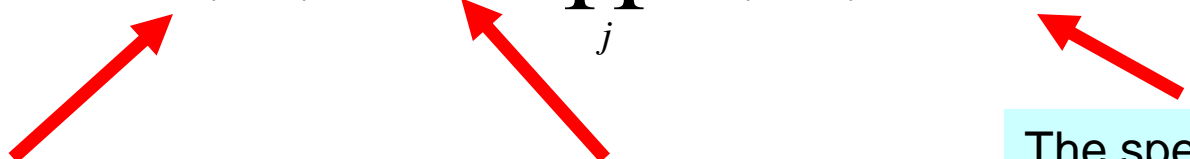
Specific model for attribute  $j$

- Using this idea the full classification rule becomes:

$$\begin{aligned}\hat{y} &= \arg \max_v p(y = v | X) \\ &= \arg \max_v \frac{p(X | y = v) p(y = v)}{p(X)} \\ &= \arg \max_v \prod_j p_j(x^j | y = v) p(y = v)\end{aligned}$$

$v$  are the classes we have

# Conditional likelihood: Full version

$$L(X_i | y_i = 1, \Theta) = \prod_j p(x_i^j | y_i = 1, \theta_1^j)$$


Vector of binary attributes for sample  $i$

The set of all parameters in the NB model

The specific parameters for attribute  $j$  in class 1

Note the following:

1. We assume conditional independence between attributes **given** the class label
2. We learn a **different** set of parameters for the two classes (class 1 and class 2).

# Learning parameters

$$L(X_i | y_i = 1, \Theta) = \prod_j p(x_i^j | y_i = 1, \theta_1^j)$$

- Let  $X_1 \dots X_{k_1}$  be the set of input samples with label 'y=1'
- Assume all attributes are **binary**
- To determine the MLE parameters for  $p(x^j = 1 | y = 1)$  we simply count how many times the j'th entry of those samples in class 1 is 0 (termed  $n_0$ ) and how many times its 1 ( $n_1$ ). Then we set:

$$p(x^j = 1 | y = 1) = \frac{n_1}{n_0 + n_1}$$

# Final classification

- Once we computed all parameters for attributes in both classes we can easily decide on the label of a **new** sample  $X$ .

$$\begin{aligned}\hat{y} &= \arg \max_v p(y = v | X) \\ &= \arg \max_v \frac{p(X | y = v)p(y = v)}{p(X)} \\ &= \arg \max_v \prod_j p_j(x^j | y = v)p(y = v)\end{aligned}$$

Perform this computation for both class 1 and class 2 and select the class that leads to a higher probability as your decision

Prior on the prevalence of samples from each class



# Example: Text classification

- What is the major topic of this article?

**boston.com**

HOME OBITUARIES SPORTS ENTERTAINMENT BUSINESS LIFESTYLE HEALTH TRAVEL CARS JOBS REAL ESTATE

NEWS [SHARE](#)

## The story behind Mitt Romney's loss in the presidential campaign to President Obama

By Michael Kranish  
Globe Staff  
DECEMBER 22, 2012 7:00 PM

...who believe that they are entitled to health care, to food, to housing, to you-name-it. That that's an entitlement. And the government should give it to them. And they will vote for this president no matter what.

**Mother Jones**  
VIDEO

A video from a May fund-raiser in Florida showed Romney characterizing nearly half of Americans as "victims" who want... Credit: Associated Press photo of Mother Jones video

# Example: Text classification

- Text classification is all around us

The screenshot shows a web browser window displaying the Google News homepage. The browser's address bar shows the URL <http://news.google.com/>. The page features the Google News logo and a search bar with the text "Search and browse 4,500 news sources".

On the left side, there is a navigation menu with categories: **Top Stories**, World, U.S., Business, Sci/Tech, Entertainment, Sports, Health, and Most Popular. Below these are links for "News Alerts", "Text Version", "Image Version", "RSS | Atom", "About Feeds", and "Mobile News".

The main content area is titled "Top Stories" and includes a location filter set to "U.S." and a "Go" button. Three news items are listed:

- Israel Rejects Truce, Presses on With Gaza Strikes**  
Washington Post - 1 hour ago  
By Griff Witte and Sudarsan Raghavan JERUSALEM, Dec. 30 – Israel continued airstrikes against Gaza Strip targets for a fourth day on Tuesday, destroying civic and other buildings linked to the militant Hamas movement in a campaign Israeli leaders say ...  
Video: [Israel to fight Hamas to the end](#) RussiaToday  
[Egypt refuses full opening of Gaza crossing](#) The Associated Press  
[Times Online](#) - [Xinhua](#) - [AFP](#) - [BBC News](#)  
[all 22,637 news articles »](#)
- Pardo intended to kill others, police say**  
Los Angeles Times - 3 hours ago  
Residents gathered Monday in Covina to talk about the rampage in their community and, in the words of Police Chief Kim Raney, "begin the healing process."  
[Police: Santa shooter planned to kill divorce attorney, mother](#) CNN  
[Police: Santa gunman planned to kill more than 9](#) The Associated Press  
[San Francisco Chronicle](#) - [WQAD](#) - [Quad City Times](#) - [Houston Chronicle](#)  
[all 695 news articles »](#)
- No official word yet on weekend birth of Palin grandchild**  
The Miami Herald - 59 minutes ago  
Alaska Gov. Sarah Palin is a grandmother. The baby's name is Tripp, and he was born early Sunday morning, People magazine is reporting.  
Video: [Palin's Daughter Gives Birth to Son Named Tripp](#) AssociatedPress  
[Palin's teenaged daughter gives birth to son: Report](#) Indian Express  
[San Jose Mercury News](#) - [guardian.co.uk](#) - [The Week Magazine](#) - [BBC News](#)  
[all 869 news articles »](#)

At the bottom of the page, there are links for "Recommended for you" and "Local News".

# Feature transformation

- How do we encode the set of features (words) in the document?
  - What type of information do we wish to represent? What can we ignore?
  - Most common encoding: '**Bag of Words**'
  - Treat document as a collection of words and encode each document as a vector based on some dictionary
  - The vector can either be binary (present / absent information for each word) or discrete (number of appearances)
- 
- Google is a good example
  - Other applications include job search adds, spam filtering and many more.

# Feature transformation: Bag of Words

- In this example we will use a binary vector
- For document  $X_i$  we will use a vector of  $m^*$  indicator features  $\{\phi^j(X_i)\}$  for whether a word appears in the document
  - $\phi^j(X_i) = 1$ , if word  $j$  appears in document  $X_i$ ;
  - $\phi^j(X_i) = 0$  if it does not appear in the document
- $\Phi(X_i) = [\phi^1(X_i) \dots \phi^m(X_i)]^T$  is the resulting feature vector for the entire dictionary for document  $X_i$
- For notational simplicity we will replace each document  $X_i$  with a fixed length vector  $\Phi_i = [\phi^1 \dots \phi^m]^T$ , where  $\phi^j = \phi^j(X_i)$ .

\*The size of the vector for English is usually ~10000 words

# Example

Assume we would like to classify documents as election related or not.

Dictionary

- Washington
- Congress
- ...
- 54. Romney
- 55. Obama
- 56. Nader

$$\phi^{54} = \phi^{54}(X_i) = 1$$

$$\phi^{55} = \phi^{55}(X_i) = 1$$

$$\phi^{56} = \phi^{56}(X_i) = 0$$

**boston.com**

HOME OBITUARIES SPORTS ENTERTAINMENT BUSINESS LIFESTYLE HEALTH TRAVEL CARS JOBS REAL ESTATE

NEWS [SHARE](#)

## The story behind Mitt Romney's loss in the presidential campaign to President Obama

By Michael Kranish  
Globe Staff  
DECEMBER 22, 2012 7:00 PM

...who believe that they are entitled to health care, to food, to housing, to you-name-it. That that's an entitlement. And the government should give it to them. And they will vote for this president no matter what.

**Mother Jones**  
VIDEO

A video from a May fund-raiser in Florida showed Romney characterizing nearly half of Americans as "victims" who want associated Press photo of Mother Jones Video Credit: Associated Press photo of Mother Jones Video

# Example: cont.

We would like to classify documents as election related or not.

- Given a collection of documents with their labels (usually termed ‘training data’) we learn the parameters for our model.
- For example, if we see the word ‘Obama’ in  $n_1$  out of the  $n$  documents labeled as ‘election’ we set  $p(\text{‘obama’} | \text{‘election’}) = n_1/n$
- Similarly we compute the priors ( $p(\text{‘election’})$ ) based on the proportion of the documents from both classes.

**boston.com**

HOME OBITUARIES SPORTS ENTERTAINMENT BUSINESS LIFESTYLE HEALTH TRAVEL CARS JOBS REAL ESTATE

NEWS SHARE

## The story behind Mitt Romney's loss in the presidential campaign to President Obama

By Michael Kranish  
Globe Staff  
DECEMBER 22, 2012 7:00 PM

...who believe that they are entitled to health care, to food, to housing, to you-name-it. That that's an entitlement. And the government should give it to them. And they will vote for this president no matter what.

**Mother Jones**  
VIDEO

A video from a May fund-raiser in Florida showed Romney characterizing nearly half of Americans as "victims" who want...  
Associated Press photo of Mother Jones video Credit: Associated Press photo of Mother Jones video

# Example: Classifying Election (E) or Sports (S)

Assume we learned the following model

$$P(\phi^{\text{romney}} = 1 | E) = 0.8, \quad P(\phi^{\text{romney}} = 1 | S) = 0.1 \quad P(S) = 0.5$$

$$P(\phi^{\text{obama}} = 1 | E) = 0.9, \quad P(\phi^{\text{obama}} = 1 | S) = 0.05 \quad P(E) = 0.5$$

$$P(\phi^{\text{clinton}} = 1 | E) = 0.9, \quad P(\phi^{\text{clinton}} = 1 | S) = 0.05$$

$$P(\phi^{\text{football}} = 1 | E) = 0.1, \quad P(\phi^{\text{football}} = 1 | S) = 0.7$$

For a specific document we have the following feature vector

$$\phi^{\text{romney}} = 1 \quad \phi^{\text{obama}} = 1 \quad \phi^{\text{clinton}} = 1 \quad \phi^{\text{football}} = 0$$

$$P(y = E | 1, 1, 1, 0) \propto 0.8 * 0.9 * 0.9 * 0.9 * 0.5 = 0.5832$$

$$P(y = S | 1, 1, 1, 0) \propto 0.1 * 0.05 * 0.05 * 0.3 * 0.5 = 0.000075$$

So the document is classified as 'Election'

# Naïve Bayes classifiers for continuous values

- So far we assumed a binomial or discrete distribution for the data given the model ( $p(X_i|y)$ )
- However, in many cases the data contains continuous features:
  - Height, weight
  - Levels of genes in cells
  - Brain activity
- For these types of data we often use a Gaussian model
- In this model we assume that the observed input vector  $X$  is generated from the following distribution

$$X \sim N(\mu, \Sigma)$$



# Gaussian Bayes Classifier Assumption

- The  $i$ 'th record in the database is created using the following algorithm
  1. Generate the output (the “class”) by drawing  $y_i \sim \text{Multinomial}(p_1, p_2, \dots, p_{N_y})$
  2. Generate the inputs from a Gaussian PDF that depends on the value of  $y_i$ :

$$\mathbf{x}_i \sim N(\mathbf{m}_i, \mathbf{S}_i).$$

# Gaussian Bayes Classification

- To determine the class when using the Gaussian assumption we need to compute  $p(X|y)$ :

$$P(y = v | X) = \frac{p(X | y = v)P(y = v)}{p(X)}$$

$$P(X | y) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left[ -\frac{1}{2} (X - \mu)^T \Sigma^{-1} (X - \mu) \right]$$

Once again, we need lots of data to compute the values of the mean  $\mu$  and the covariance matrix  $\Sigma$

# Gaussian Bayes Classification

- Here we can also use the Naïve Bayes assumption: Attributes are independent given the class label
- In the Gaussian model this means that the covariance matrix becomes a **diagonal matrix** with zeros everywhere except for the diagonal
- Thus, we only need to learn the values for the variance term for each attribute:  $x^j \sim N(\mu^j, \sigma^j)$

$$P(X | y = v) = \prod_j \frac{1}{(2\pi)^{1/2} \sigma_v^j} \exp \left[ -\frac{(\mathbf{x}_j - \mu_v^j)^2}{2(\sigma_v^j)^2} \right]$$

Separate means and variance for each class

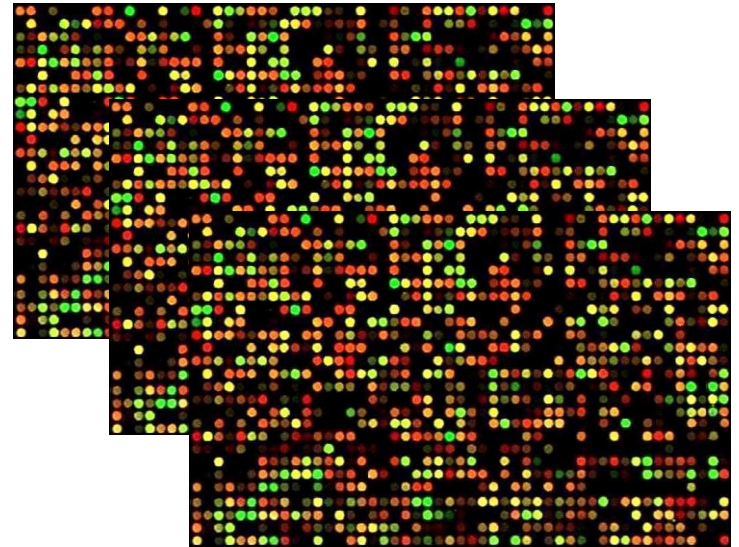
# MLE for Gaussian Naïve Bayes Classifier

- For each class we need to estimate one global value (prior) and two values for each feature (mean and variance)
- The prior is computed in the same way we did before (counting) which is the MLE estimate For each feature
- Let the numbers of input samples in class 1 be  $k_1$ . The MLE for mean and variance is computed by setting:

$$\mu_1^j = \frac{1}{k_1} \sum_{X_i | s.t. y_i = 1} x_i^j \quad \sigma_1^{j^2} = \frac{1}{k_1} \sum_{X_i | s.t. y_i = 1} (x_i^j - \mu_1^j)^2$$

# Example: Classifying gene expression data

- Measures the levels (up or down) of genes in our cells
- Differs between healthy and sick people and between different disease types
- Given measurement of patients with two different types of cancer we would like to generate a classifier to distinguish between them



# Classifying cancer types

- We select a subset of the genes (more in our 'feature selection' class later in the course).
- We compute the mean and variance for each of the genes in each of the classes
- Compute the class priors based on the input samples

**Class 1  
(ALL)**

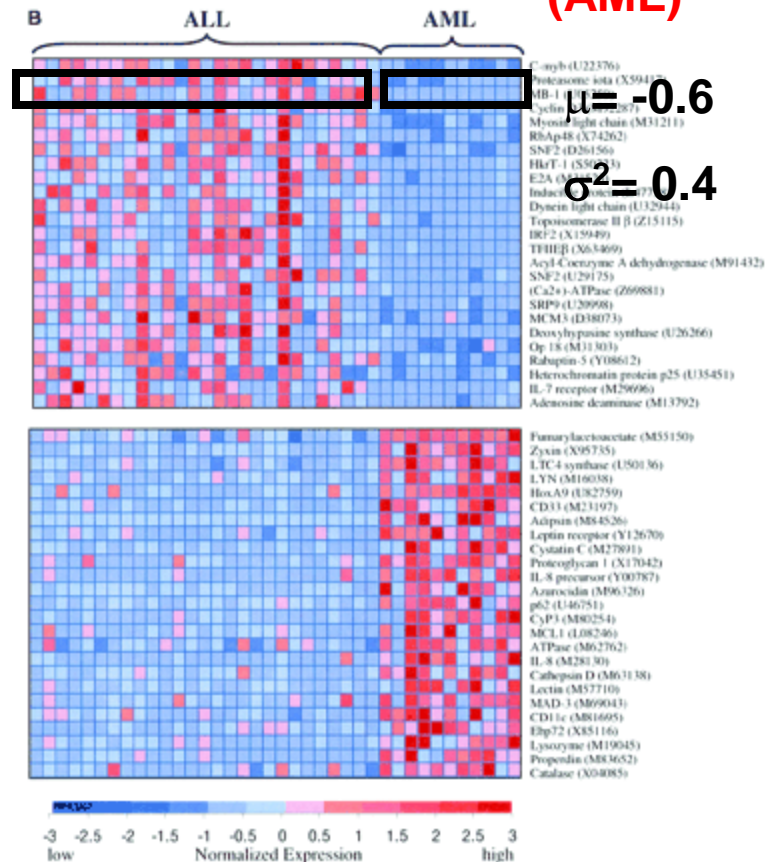
$$\mu = 1.8$$

$$\sigma^2 = 1.1$$

**Class 2  
(AML)**

$$\mu = -0.6$$

$$\sigma^2 = 0.4$$



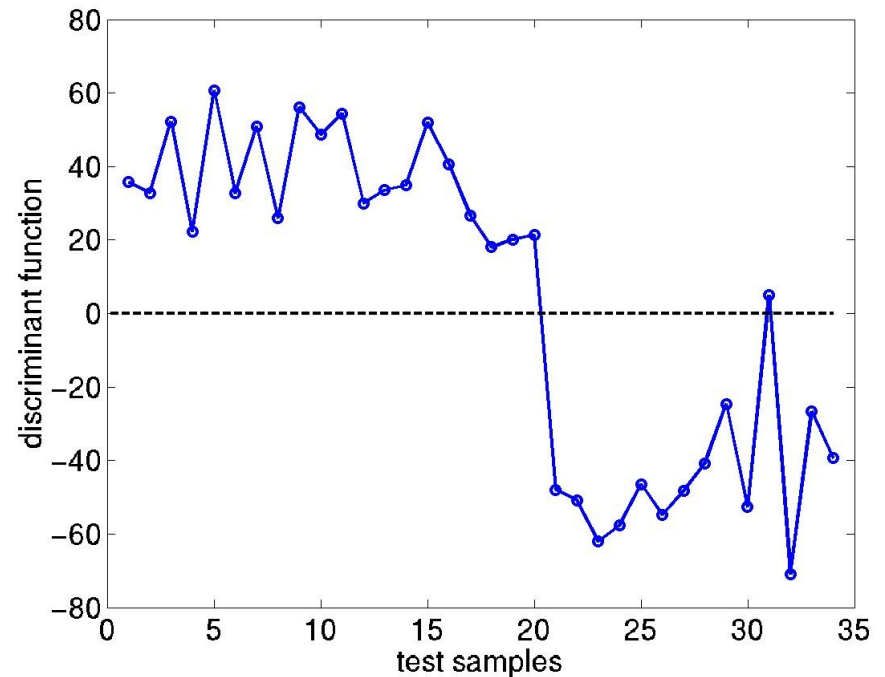
# Classification accuracy

- The figure shows the value of the discriminate function

$$f(x) = \log \frac{p(y = 1 | X)}{p(y = 0 | X)}$$

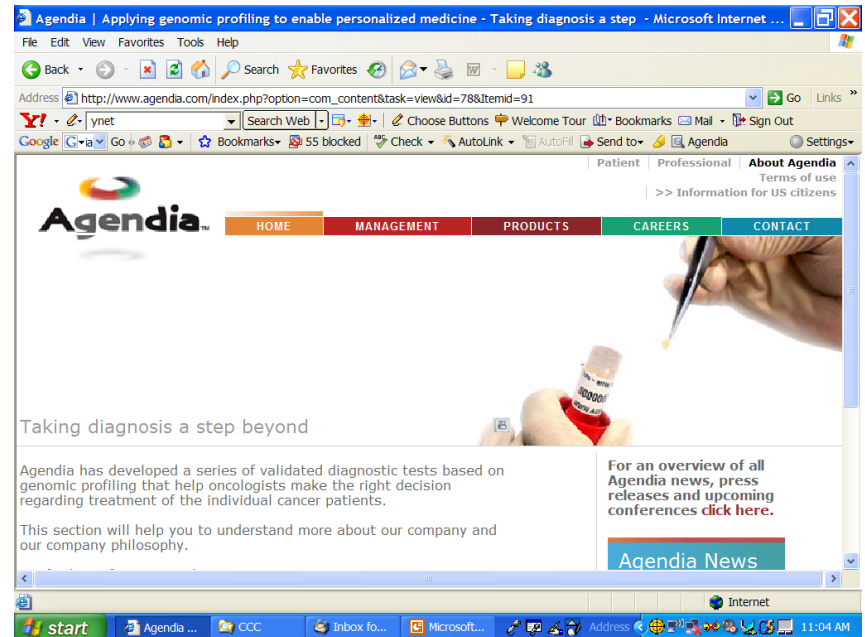
across the test examples

- The only test error is also the decision with the lowest confidence



# FDA Approves Gene-Based Breast Cancer Test\*

“ MammaPrint is a DNA microarray-based test that measures the activity of 70 genes... The test measures each of these genes in a sample of a woman's breast-cancer tumor and then uses a specific formula to determine whether the patient is deemed low risk or high risk for the spread of the cancer to another site.”



\*Washington Post, 2/06/2007



# Possible problems with Naïve Bayes classifiers: Assumptions

- In most cases, the assumption of conditional independence given the class label is violated
  - much more likely to find the word 'Barack' if we saw the word 'Obama' regardless of the class
- This is, unfortunately, a major shortcoming which makes these classifiers inferior in many real world applications (though not always)
- There are models that can improve upon this assumption without using the full conditional model (one such model are Bayesian networks which we will discuss later in this class).

# Possible problems with Naïve Bayes classifiers: Parameter estimation

- Even though we need far less data than the full Bayes model, there may be cases when the data we have is not enough
- For example, what is  $p(S=1, N=1 | E=2)$ ?
- This can get worst. Assume we have 20 variables, almost all pointing in the direction of the same class except for one for which we have no record for this class.
- Solutions?

Summer?	Num > 20	Evaluation
1	1	3
1	0	3
0	1	2
0	1	1
0	0	3
1	1	1

# Decision trees and Naïve Bayes

- What are the relationships between the assumptions the two classifiers make?
- How does this affect their ability to model different input datasets?
  - Number of feature?
  - Number of samples?
- How does this affect the way they handle the different features?

# Important points

- Problems with estimating full joints
- Advantages of Naïve Bayes assumptions
- Applications to discrete and continuous cases
- Problems with Naïve Bayes classifiers