

Statistical Learning

1

Today

- Parameter Estimation:
 - Maximum Likelihood (ML)
 - Maximum A Posteriori (MAP)
 - Bayesian
 - Continuous case
- Learning Parameters for a Bayesian Network
- Naive Bayes
 - Maximum Likelihood estimates
 - Priors
- Learning Structure of Bayesian Networks

2

Coin Flip



$$P(H|C_1) = 0.1$$

$$P(H|C_2) = 0.5$$

$$P(H|C_3) = 0.9$$

Which coin will I use?

$$P(C_1) = 1/3$$

$$P(C_2) = 1/3$$

$$P(C_3) = 1/3$$

Prior: Probability of a hypothesis before we make any observations

3

Coin Flip



$$P(H|C_1) = 0.1$$

$$P(H|C_2) = 0.5$$

$$P(H|C_3) = 0.9$$

Which coin will I use?

$$P(C_1) = 1/3$$

$$P(C_2) = 1/3$$

$$P(C_3) = 1/3$$

Uniform Prior: All hypothesis are equally likely before we make any observations

4

Experiment I: Heads

Which coin did I use?

$$P(C_1|H) = ?$$

$$P(C_2|H) = ?$$

$$P(C_3|H) = ?$$

$$P(C_1|H) = \frac{P(H|C_1)P(C_1)}{P(H)} \quad P(H) = \sum_{i=1}^3 P(H|C_i)P(C_i)$$



$$P(H|C_1) = 0.1$$

$$P(C_1) = 1/3$$

$$P(H|C_2) = 0.5$$

$$P(C_2) = 1/3$$

$$P(H|C_3) = 0.9$$

$$P(C_3) = 1/3$$

5

Experiment I: Heads

Which coin did I use?

$$P(C_1|H) = 0.066$$

$$P(C_2|H) = 0.333$$

$$P(C_3|H) = 0.6$$

Posterior: Probability of a hypothesis given data



$$P(H|C_1) = 0.1$$

$$P(C_1) = 1/3$$

$$P(H|C_2) = 0.5$$

$$P(C_2) = 1/3$$

$$P(H|C_3) = 0.9$$

$$P(C_3) = 1/3$$

6

Experiment 2: Tails

Which coin did I use?

$$P(C_1|HT) = ? \quad P(C_2|HT) = ? \quad P(C_3|HT) = ?$$

$$P(C_1|HT) = \alpha P(HT|C_1)P(C_1) = \alpha P(H|C_1)P(T|C_1)P(C_1)$$



$$P(H|C_1) = 0.1 \\ P(C_1) = 1/3$$



$$P(H|C_2) = 0.5 \\ P(C_2) = 1/3$$



$$P(H|C_3) = 0.9 \\ P(C_3) = 1/3$$

7

Experiment 2: Tails

Which coin did I use?

$$P(C_1|HT) = 0.21 \quad P(C_2|HT) = 0.58 \quad P(C_3|HT) = 0.21$$

$$P(C_1|HT) = \alpha P(HT|C_1)P(C_1) = \alpha P(H|C_1)P(T|C_1)P(C_1)$$



$$P(H|C_1) = 0.1 \\ P(C_1) = 1/3$$



$$P(H|C_2) = 0.5 \\ P(C_2) = 1/3$$



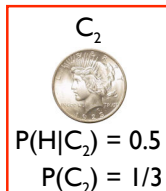
$$P(H|C_3) = 0.9 \\ P(C_3) = 1/3$$

8

Experiment 2: Tails

Which coin did I use?

$$P(C_1|HT) = 0.21 \quad P(C_2|HT) = 0.58 \quad P(C_3|HT) = 0.21$$



$$P(H|C_2) = 0.5 \\ P(C_2) = 1/3$$

9

Your Estimate?

What is the probability of heads after two experiments?

Most likely coin:



Best estimate for P(H)

$$P(H|C_2) = 0.5$$



$$P(H|C_2) = 0.5 \\ P(C_2) = 1/3$$

10

Your Estimate?

Maximum Likelihood Estimate: The best hypothesis that fits observed data assuming uniform prior

Most likely coin:



Best estimate for P(H)

$$P(H|C_2) = 0.5$$



$$P(H|C_2) = 0.5 \\ P(C_2) = 1/3$$

11

Using Prior Knowledge

- Should we always use **Uniform Prior**?
- Background knowledge:
 - Heads => you go first in Abalone against TA
 - TAs are nice people
 - => TA is more likely to use a coin biased in your favor



$$P(H|C_1) = 0.1$$



$$P(H|C_2) = 0.5$$




$$P(H|C_3) = 0.9$$

12


Using Prior Knowledge

We can encode it in the **prior**:


$$P(C_1) = 0.05 \quad P(C_2) = 0.25 \quad P(C_3) = 0.70$$



$P(H|C_1) = 0.1$



$P(H|C_2) = 0.5$



$P(H|C_3) = 0.9$


13

Experiment 1: Heads


Which coin did I use?

$$P(C_1|H) = ? \quad P(C_2|H) = ? \quad P(C_3|H) = ?$$


$$P(C_1|H) = \alpha P(H|C_1)P(C_1)$$



$P(H|C_1) = 0.1$



$P(H|C_2) = 0.5$



$P(H|C_3) = 0.9$

$P(C_1) = 0.05$

$P(C_2) = 0.25$

$P(C_3) = 0.70$

14


Experiment 1: Heads

Which coin did I use?


$$P(C_1|H) = 0.006 \quad P(C_2|H) = 0.165 \quad P(C_3|H) = 0.829$$

ML posterior after Exp 1:


$$P(C_1|H) = 0.066 \quad P(C_2|H) = 0.333 \quad P(C_3|H) = 0.600$$



$P(H|C_1) = 0.1$



$P(H|C_2) = 0.5$



$P(H|C_3) = 0.9$

$P(C_1) = 0.05$

$P(C_2) = 0.25$

$P(C_3) = 0.70$


15

Experiment 2: Tails


Which coin did I use?

$$P(C_1|HT) = ? \quad P(C_2|HT) = ? \quad P(C_3|HT) = ?$$


$$P(C_1|HT) = \alpha P(HT|C_1)P(C_1) = \alpha P(H|C_1)P(T|C_1)P(C_1)$$



$P(H|C_1) = 0.1$



$P(H|C_2) = 0.5$



$P(H|C_3) = 0.9$

$P(C_1) = 0.05$

$P(C_2) = 0.25$

$P(C_3) = 0.70$


16

Experiment 2: Tails


Which coin did I use?

$$P(C_1|HT) = 0.035 \quad P(C_2|HT) = 0.481 \quad P(C_3|HT) = 0.485$$


$$P(C_1|HT) = \alpha P(HT|C_1)P(C_1) = \alpha P(H|C_1)P(T|C_1)P(C_1)$$



$P(H|C_1) = 0.1$



$P(H|C_2) = 0.5$



$P(H|C_3) = 0.9$

$P(C_1) = 0.05$

$P(C_2) = 0.25$


$P(C_3) = 0.70$

17

Experiment 2: Tails

Which coin did I use?

$$P(C_1|HT) = 0.035 \quad P(C_2|HT) = 0.481 \quad P(C_3|HT) = 0.485$$



$P(H|C_3) = 0.9$

$P(C_3) = 0.70$

18

Your Estimate?

What is the probability of heads after two experiments?

Most likely coin:



Best estimate for $P(H)$

$$P(H|C_3) = 0.9$$



$$P(H|C_3) = 0.9$$

$$P(C_3) = 0.70$$

19

Your Estimate?

Maximum A Posteriori (MAP) Estimate: The best hypothesis that fits observed data assuming a non-uniform prior

Most likely coin:



Best estimate for $P(H)$

$$P(H|C_3) = 0.9$$



$$P(H|C_3) = 0.9$$

$$P(C_3) = 0.70$$

20

Did We Do The Right Thing?

$$P(C_1|HT) = 0.035 \quad P(C_2|HT) = 0.481 \quad P(C_3|HT) = 0.485$$



C_1

$$P(H|C_1) = 0.1$$



C_2

$$P(H|C_2) = 0.5$$



C_3

$$P(H|C_3) = 0.9$$

21

Did We Do The Right Thing?

$$P(C_1|HT) = 0.035 \quad P(C_2|HT) = 0.481 \quad P(C_3|HT) = 0.485$$

C_2 and C_3 are almost equally likely



C_1

$$P(H|C_1) = 0.1$$



C_2

$$P(H|C_2) = 0.5$$



C_3

$$P(H|C_3) = 0.9$$

22

A Better Estimate

$$\text{Recall: } P(H) = \sum_{i=1}^3 P(H|C_i)P(C_i) = 0.680$$

$$P(C_1|HT) = 0.035 \quad P(C_2|HT) = 0.481 \quad P(C_3|HT) = 0.485$$



C_1

$$P(H|C_1) = 0.1$$



C_2

$$P(H|C_2) = 0.5$$



C_3

$$P(H|C_3) = 0.9$$

23

Bayesian Estimate

Bayesian Estimate: Minimizes prediction error, given data and (generally) assuming a non-uniform prior

$$P(H) = \sum_{i=1}^3 P(H|C_i)P(C_i) = 0.680$$

$$P(C_1|HT) = 0.035 \quad P(C_2|HT) = 0.481 \quad P(C_3|HT) = 0.485$$



C_1

$$P(H|C_1) = 0.1$$



C_2

$$P(H|C_2) = 0.5$$



C_3

$$P(H|C_3) = 0.9$$

24

Comparison

- **ML (Maximum Likelihood):**
 $P(H) = 0.5$
- **MAP (Maximum A Posteriori):**
 $P(H) = 0.9$
- **Bayesian:**
 $P(H) = 0.68$

25

Comparison

- **ML (Maximum Likelihood):**
 $P(H) = 0.5$
after 10 experiments (HTH⁸): $P(H) = 0.9$
- **MAP (Maximum A Posteriori):**
 $P(H) = 0.9$
after 10 experiments (HTH⁸): $P(H) = 0.9$
- **Bayesian:**
 $P(H) = 0.68$
after 10 experiments (HTH⁸): $P(H) = 0.9$

26

Comparison

- **ML (Maximum Likelihood):**
- **MAP (Maximum A Posteriori):**
- **Bayesian:**
 - Minimizes error => great when data is scarce
 - Potentially much harder to compute

27

Comparison

- **ML (Maximum Likelihood):**
- **MAP (Maximum A Posteriori):**
 - Still easy to compute
 - Incorporates prior knowledge
- **Bayesian:**
 - Minimizes error => great when data is scarce
 - Potentially much harder to compute

28

Comparison

- **ML (Maximum Likelihood):**
 - Easy to compute
- **MAP (Maximum A Posteriori):**
 - Still easy to compute
 - Incorporates prior knowledge
- **Bayesian:**
 - Minimizes error => great when data is scarce
 - Potentially much harder to compute

29

Summary For Now

- **Prior:**
- **Uniform Prior:**
- **Posterior:**
- **Likelihood:**

30

Summary For Now

- **Prior:** Probability of a hypothesis before we see any data
- **Uniform Prior:** A prior that makes all hypothesis equally likely
- **Posterior:** Probability of a hypothesis after we saw some data
- **Likelihood:** Probability of data given hypothesis

	Prior	Hypothesis	
Maximum Likelihood Estimate	Uniform	The most likely	Point
Maximum A Posteriori Estimate	Any	The most likely	Point
Bayesian Estimate	Any	Weighted combination	Average

31

Continuous Case

- In the previous example, we chose from a **discrete** set of three coins
- In general, we have to pick from a **continuous** distribution of biased coins

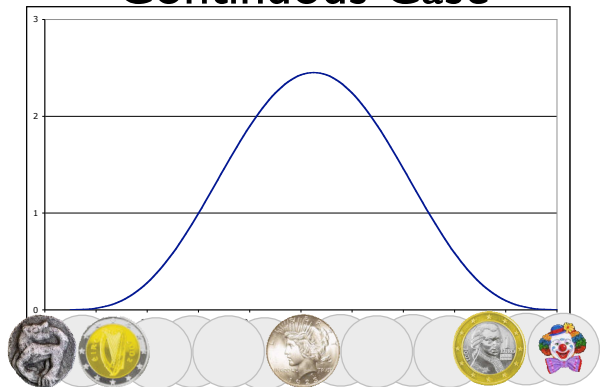
32

Continuous Case



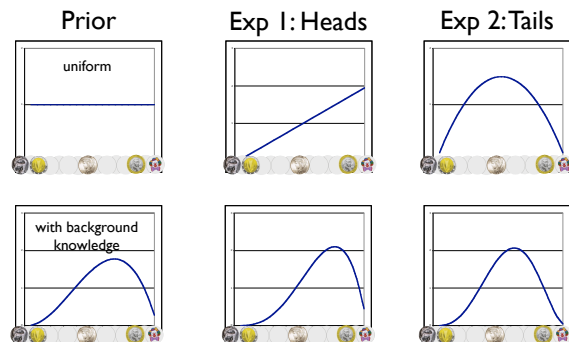
33

Continuous Case



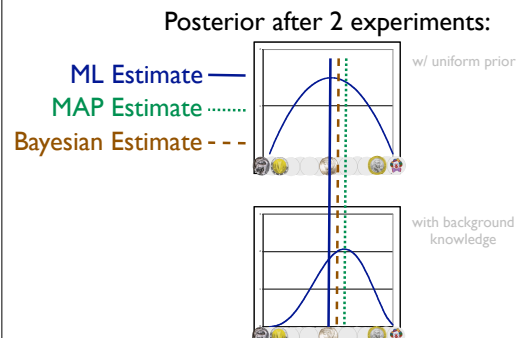
34

Continuous Case



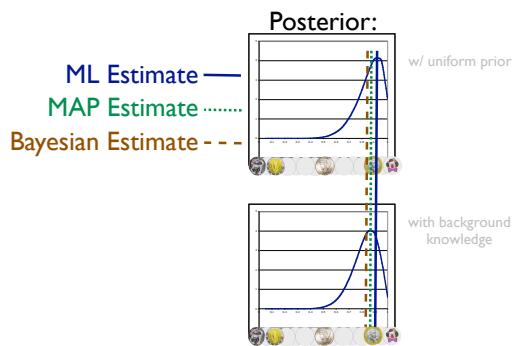
35

Continuous Case



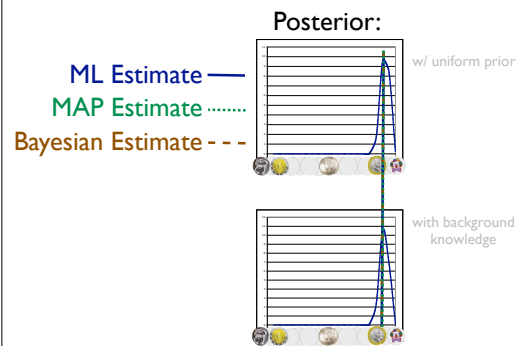
36

After 10 Experiments...



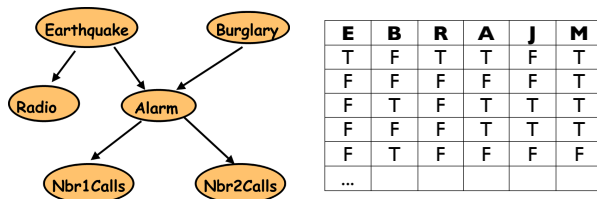
37

After 100 Experiments...



38

Parameter Estimation and Bayesian Networks

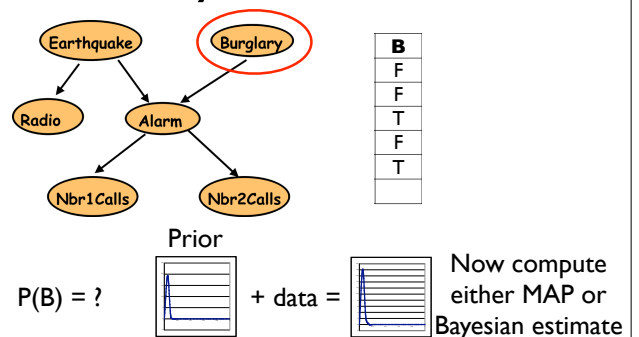


We have:

- Bayes Net **structure** and **observations**
- We need: Bayes Net **parameters**

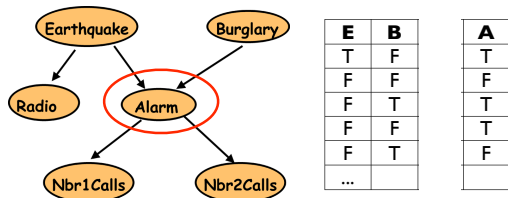
39

Parameter Estimation and Bayesian Networks



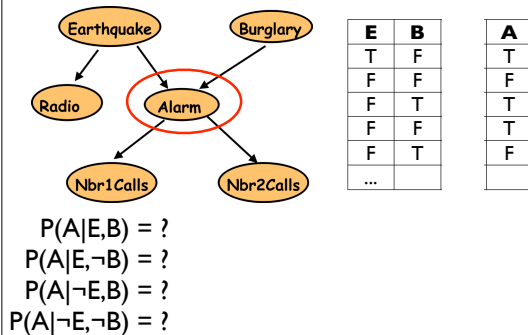
40

Parameter Estimation and Bayesian Networks



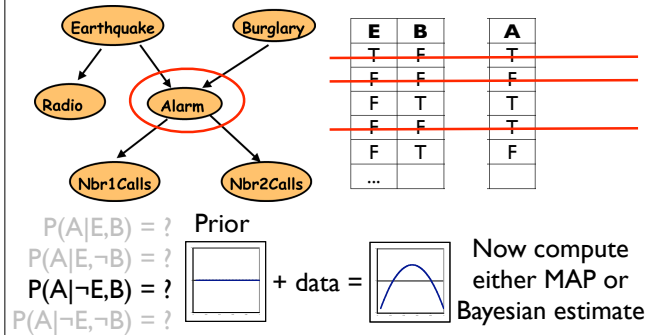
41

Parameter Estimation and Bayesian Networks



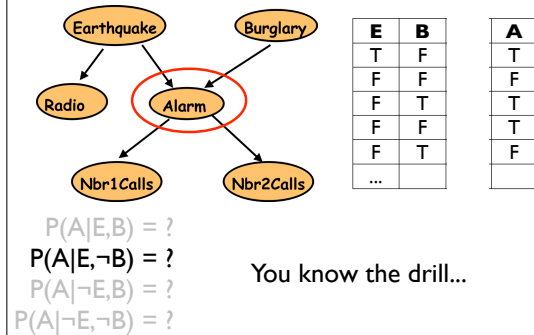
42

Parameter Estimation and Bayesian Networks



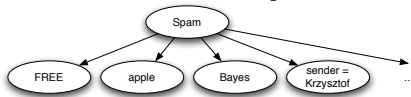
43

Parameter Estimation and Bayesian Networks



44

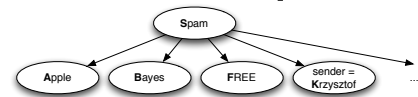
Naive Bayes



- A Bayes Net where all nodes are children of a single root node
- Why?
 - Expressive and accurate?
 - Easy to learn?

45

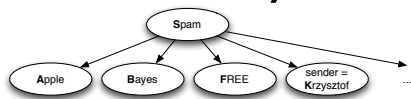
Naive Bayes



- A Bayes Net where all nodes are children of a single root node
- Why?
 - Expressive and accurate? **No** - why?
 - Easy to learn?

46

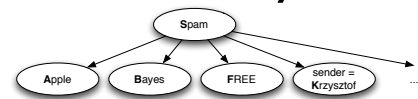
Naive Bayes



- A Bayes Net where all nodes are children of a single root node
- Why?
 - Expressive and accurate? **No**
 - Easy to learn? **Yes**

47

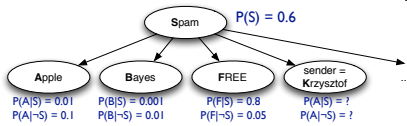
Naive Bayes



- A Bayes Net where all nodes are children of a single root node
- Why?
 - Expressive and accurate? **No**
 - Easy to learn? **Yes**
 - Useful? **Sometimes**

48

Inference In Naive Bayes

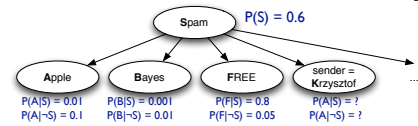


- Goal, given evidence (words in an email) decide if an email is spam

$$E = \{A, \neg B, F, \neg K, \dots\}$$

49

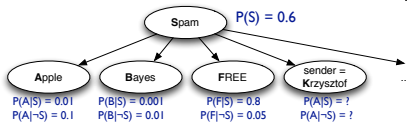
Inference In Naive Bayes



$$\begin{aligned} P(S|E) &= \frac{P(E|S)P(S)}{P(E)} \\ &= \frac{P(A, \neg B, F, \neg K, \dots | S)P(S)}{P(A, \neg B, F, \neg K, \dots)} \quad \text{Independence to the rescue!} \\ &= \frac{P(A|S)P(\neg B|S)P(F|S)P(\neg K|S)P(\dots|S)P(S)}{P(A)P(\neg B)P(F)P(\neg K)P(\dots)} \end{aligned}$$

50

Inference In Naive Bayes



$$P(S|E) = \frac{P(A|S)P(\neg B|S)P(F|S)P(\neg K|S)P(\dots|S)P(S)}{P(A)P(\neg B)P(F)P(\neg K)P(\dots)}$$

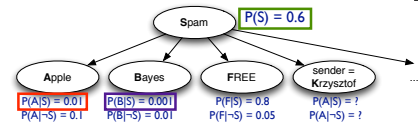
$$P(\neg S|E) = \frac{P(A|\neg S)P(\neg B|\neg S)P(F|\neg S)P(\neg K|\neg S)P(\dots|\neg S)P(\neg S)}{P(A)P(\neg B)P(F)P(\neg K)P(\dots)}$$

Spam if $P(S|E) > P(\neg S|E)$

But...

51

Inference In Naive Bayes

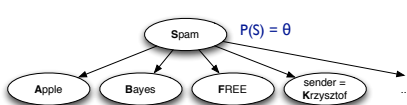


$$P(S|E) \propto P(A|S)P(\neg B|S)P(F|S)P(\neg K|S)P(\dots|S)P(S)$$

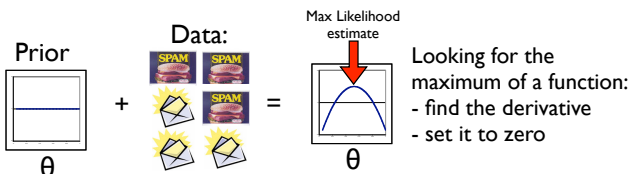
$$P(\neg S|E) \propto P(A|\neg S)P(\neg B|\neg S)P(F|\neg S)P(\neg K|\neg S)P(\dots|\neg S)P(\neg S)$$

52

Parameter Estimation Revisited

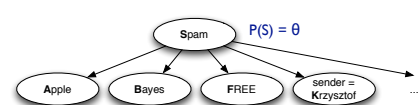


- Can we calculate Maximum Likelihood estimate of θ easily?



53

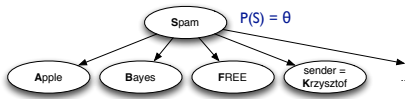
Parameter Estimation Revisited



- What function are we maximizing?
 $P(\text{data}|\text{hypothesis})$

54

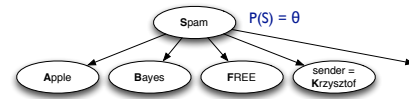
Parameter Estimation Revisited



- What function are we maximizing?
 $P(\text{data}|\text{hypothesis})$
- hypothesis = h_θ (one for each value of θ)

55

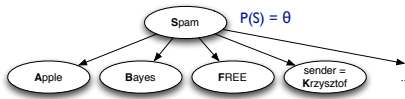
Parameter Estimation Revisited



- What function are we maximizing?
 $P(\text{data}|\text{hypothesis})$
- hypothesis = h_θ (one for each value of θ)
- $P(\text{data}|h_\theta) = P(\text{Apple}|h_\theta)P(\text{Bayes}|h_\theta)P(\text{FREE}|h_\theta)P(\text{sender = Krzysztof}|h_\theta)$

56

Parameter Estimation Revisited

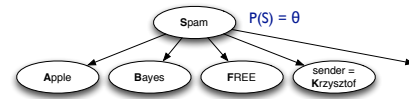


- What function are we maximizing?
 $P(\text{data}|\text{hypothesis})$
- hypothesis = h_θ (one for each value of θ)
- $P(\text{data}|h_\theta) = P(\text{Apple}|h_\theta)P(\text{Bayes}|h_\theta)P(\text{FREE}|h_\theta)P(\text{sender = Krzysztof}|h_\theta)$

$$= \theta (1-\theta) (1-\theta) \theta$$

57

Parameter Estimation Revisited



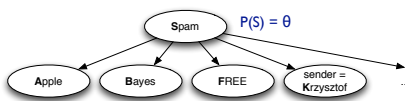
- What function are we maximizing?
 $P(\text{data}|\text{hypothesis})$
- hypothesis = h_θ (one for each value of θ)
- $P(\text{data}|h_\theta) = P(\text{Apple}|h_\theta)P(\text{Bayes}|h_\theta)P(\text{FREE}|h_\theta)P(\text{sender = Krzysztof}|h_\theta)$

$$= \theta (1-\theta) (1-\theta) \theta$$

$$= \theta^{\# \text{Apple}} (1-\theta)^{\# \text{Bayes}}$$

58

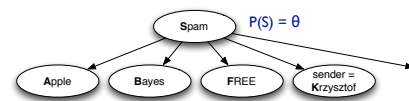
Parameter Estimation Revisited



- To find θ that maximizes $\theta^{\# \text{Apple}} (1-\theta)^{\# \text{Bayes}}$ we take a derivative of the function and set it to 0. And we get:

59

Parameter Estimation Revisited



- To find θ that maximizes $\theta^{\# \text{Apple}} (1-\theta)^{\# \text{Bayes}}$ we take a derivative of the function and set it to 0. And we get:
- $P(S) = \theta = \frac{\# \text{Apple}}{\# \text{Apple} + \# \text{Bayes}}$
- You knew it already, right?

60

Problems With Small Samples

- What happens if in your training data apples are not mentioned in any spam message?
- $P(A|S) = 0$
- Why is it bad?

$$P(S|E) \propto \mathbf{0} \quad P(\neg B|S)P(F|S)P(\neg K|S)P(\dots|S)P(S) = \mathbf{0}$$

61

Smoothing

- **Smoothing** is used when samples are small
- **Add-one smoothing** is the simplest smoothing method: just add 1 to every count!

62

Priors!

- Recall that $P(S) = \frac{\# \text{spam}}{\# \text{spam} + \# \text{ham}}$

63

Priors!

- Recall that $P(S) = \frac{\# \text{spam}}{\# \text{spam} + \# \text{ham}}$
- If we have a slight hunch that $P(S) \approx p$

$$P(S) = \frac{\# \text{spam} + p}{\# \text{spam} + \# \text{ham} + 1}$$

64

Priors!

- Recall that $P(S) = \frac{\# \text{spam}}{\# \text{spam} + \# \text{ham}}$
- If we have a slight hunch that $P(S) \approx p$

$$P(S) = \frac{\# \text{spam} + p}{\# \text{spam} + \# \text{ham} + 1}$$

- If we have a **big** hunch that $P(S) \approx p$

$$\frac{\# \text{spam} + mp}{\# \text{spam} + \# \text{ham} + m}$$

where m can be any number > 0

65

Priors!

$$P(S) = \frac{\# \text{spam} + mp}{\# \text{spam} + \# \text{ham} + m}$$

- Note that if $m = 10$ in the above, it is like saying "I have seen 10 samples that make me believe that $P(S) = p$ "
- Hence, m is referred to as the **equivalent sample size**

66

Priors!

$$P(S) = \frac{\# \text{spam} + mp}{\# \text{spam} + \# \text{not spam} + m}$$

- Where should p come from?
- No prior knowledge $\Rightarrow p=0.5$
- If you build a personalized spam filter, you can use $p = P(S)$ from some body else's filter!

67

Inference in Naive Bayes Revisited

- Recall that

$$P(S|E) \propto P(A|S)P(\neg B|S)P(F|S)P(\neg K|S)P(\dots|S)P(S)$$

Is there any potential for trouble here?

68

Inference in Naive Bayes Revisited

- Recall that

$$P(S|E) \propto P(A|S)P(\neg B|S)P(F|S)P(\neg K|S)P(\dots|S)P(S)$$

- We are multiplying lots of small numbers together \Rightarrow danger of underflow!
- Solution? Use logs!

69

Inference in Naive Bayes Revisited

$$\log(P(S|E)) \propto \log(P(A|S)P(\neg B|S)P(F|S)P(\neg K|S)P(\dots|S)P(S))$$

$$\propto \log(P(A|S)) + \log(P(\neg B|S)) + \log(P(F|S)) + \log(P(\neg K|S)) + \log(P(\dots|S)) + \log(P(S))$$

- Now we add “regular” numbers -- little danger of over- or underflow errors!

70

Learning The Structure of Bayesian Networks

- General idea: look at all possible network structures and pick one that fits observed data best
- Impossibly slow: exponential number of networks, and for each we have to learn parameters, too!
- What do we do if searching the space exhaustively is too expensive?

71

Learning The Structure of Bayesian Networks

- Local search!
 - Start with some network structure
 - Try to make a change (add, delete or reverse node)
 - See if the new network is any better

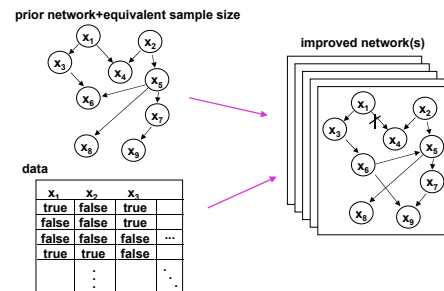
72

Learning The Structure of Bayesian Networks

- What network structure should we start with?
- Random with uniform prior?
- Networks that reflects our (or experts') knowledge of the field?

73

Learning The Structure of Bayesian Networks



74

Learning The Structure of Bayesian Networks

- We have just described how to get an ML or MAP estimate of the structure of a Bayes Net
- What would the Bayes estimate look like?
 - Find all possible networks
 - Calculate their posteriors
 - When doing inference: result weighed combination of all networks!

75