

Naïve Bayes & Expectation Maximization

CSE 573

© Daniel S. Weld

1

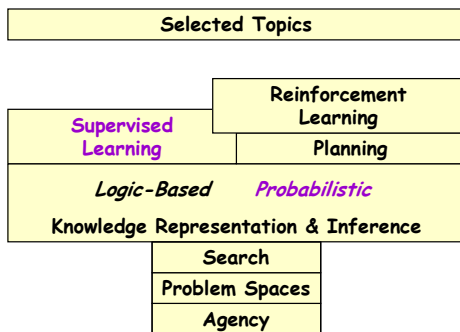
Logistics

- Team Meetings
- Midterm
 - Open book, notes
 - Studying
 - See AIMA exercises

© Daniel S. Weld

2

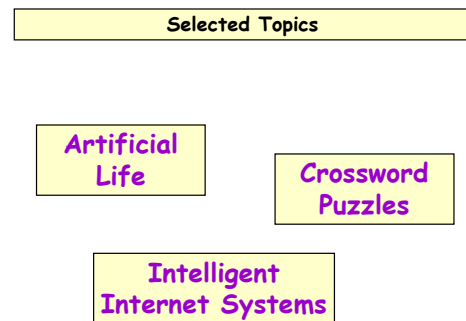
573 Schedule



© Daniel S. Weld

3

Coming Soon



© Daniel S. Weld

4

Topics

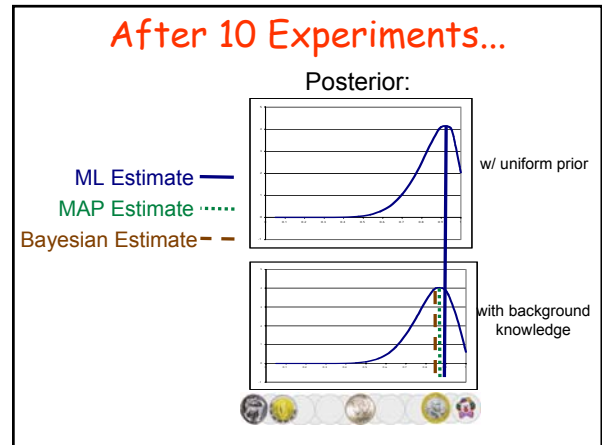
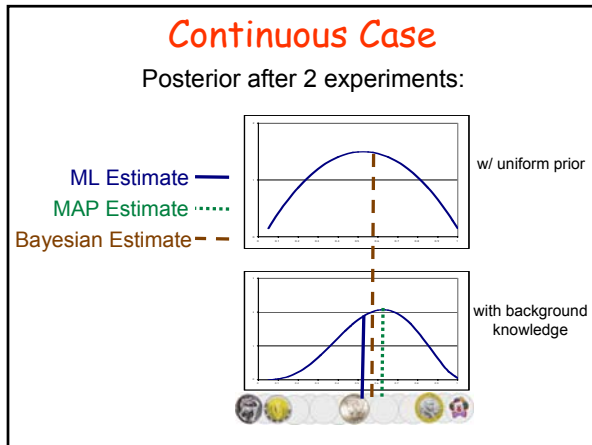
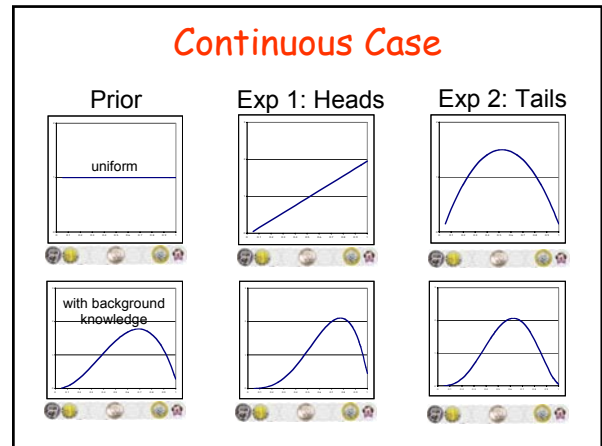
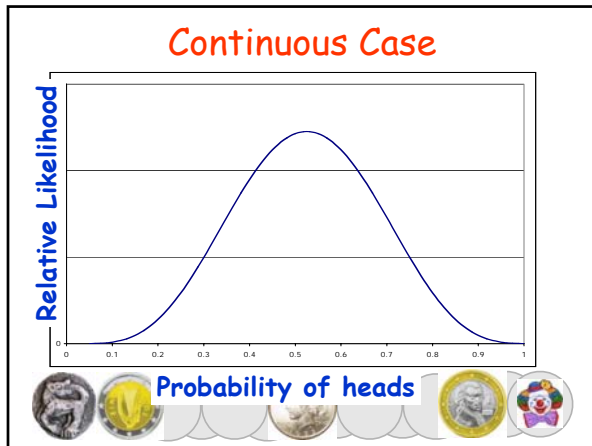
- Test & Mini Projects
- Review
- Naive Bayes
 - Maximum Likelihood Estimates
 - Working with Probabilities
- Expectation Maximization
- Challenge

© Daniel S. Weld

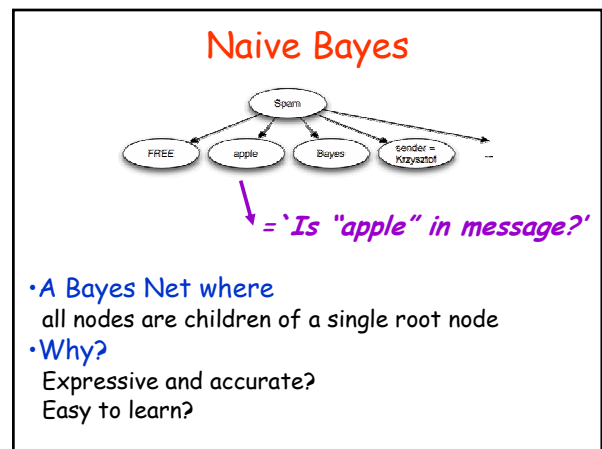
5

Estimation Models

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



- ### Topics
- Test & Mini Projects
 - Review
 - Naive Bayes
 - Maximum Likelihood Estimates
 - Working with Probabilities
 - Expectation Maximization
 - Challenge
- © Daniel S. Weld



Naive Bayes



- All nodes are children of a single root node
- Why? Expressive and accurate? **No** - why? Easy to learn?

Naive Bayes



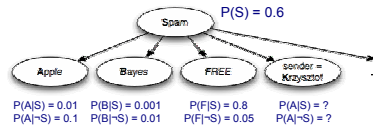
- All nodes are children of a single root node
- Why? Expressive and accurate? **No** Easy to learn? **Yes**

Naive Bayes



- All nodes are children of a single root node
- Why? Expressive and accurate? **No** Easy to learn? **Yes** Useful? **Sometimes**

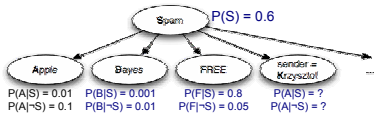
Inference In Naive Bayes



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

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

Inference In Naive Bayes



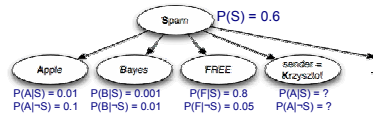
$$P(S|E) = \frac{P(E|S)P(S)}{P(E)}$$

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

Independence to the rescue!

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

Inference In Naive Bayes



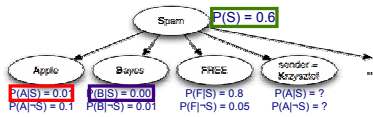
$$P(S|E) = \frac{P(A|S)P(\neg B|S)P(F|S)P(\neg K|S)P(\dots|S)P(S)}{P(E)}$$

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

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

But...

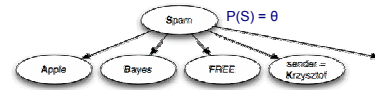
Inference In Naive Bayes



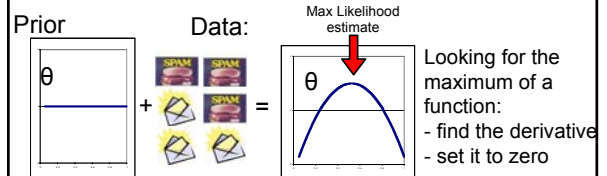
$$P(S|E) \propto P(A|S)P(B|S)P(F|S)P(K|S)P(\dots|S)P(S)$$

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

Parameter Estimation Revisited



Can we calculate Maximum Likelihood estimate of θ easily?



Topics

- Test & Mini Projects
- Review
- Naive Bayes
 - Maximum Likelihood Estimates
 - Working with Probabilities
 - Smoothing
 - Computational Details
 - Continuous Quantities
- Expectation Maximization
- Challenge

© Daniel S. Weld

21

Evidence is Easy?

$$P(X_i | S) = \frac{\# \text{spam}}{\# \text{spam} + \# \text{ham}}$$

• Or... Are their problems?

Smooth with a Prior

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

p = prior probability
 m = weight

Note that if $m = 10$, it means "I've seen 10 samples that make me believe $P(X_i | S) = p$ "

Hence, m is referred to as the **equivalent sample size**

Probabilities: Important Detail!

$$P(\text{spam} | X_1 \dots X_n) = \prod_i P(\text{spam} | X_i)$$

Any more potential problems here?

- We are multiplying lots of small numbers
 Danger of underflow!
 $0.5^{57} = 7 \text{ E } -18$
- Solution? Use logs and add!
 $p_1 * p_2 = e^{\log(p_1) + \log(p_2)}$
 Always keep in log form

P(S | X)

- Easy to compute from data if X discrete

Instance	X	Spam?
1	T	F
2	T	F
3	F	T
4	T	T
5	T	F

- $P(S | X) = \frac{1}{4}$ ignoring smoothing...

© Daniel S. Weld

25

P(S | X)

- What if X is real valued?

Instance	X	Spam?
1	-0.01 <T	False
2	-0.01- <T	False
3	-0.02- <T	False
4	-0.03- >T	True
5	-0.05- >T	True

- What now?

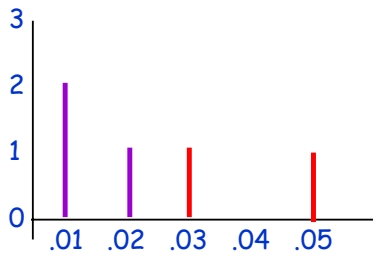
© Daniel S. Weld

26

Anything Else?

#	X	S?
1	0.01	F
2	0.01	F
3	0.02	F
4	0.03	T
5	0.05	T

$P(S|.04)?$

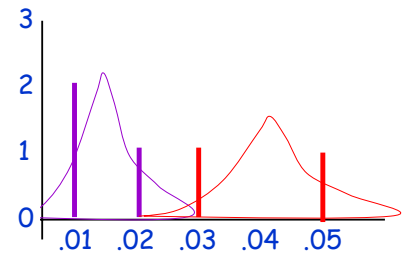


© Daniel S. Weld

27

Fit Gaussians

#	X	S?
1	0.01	F
2	0.01	F
3	0.02	F
4	0.03	T
5	0.05	T

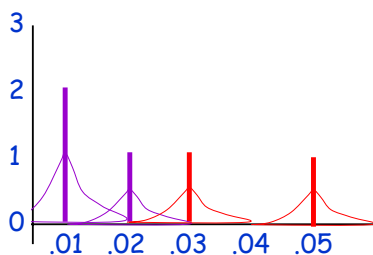


© Daniel S. Weld

28

Smooth with Gaussian then sum "Kernel Density Estimation"

#	X	S?
1	0.01	F
2	0.01	F
3	0.02	F
4	0.03	T
5	0.05	T



© Daniel S. Weld

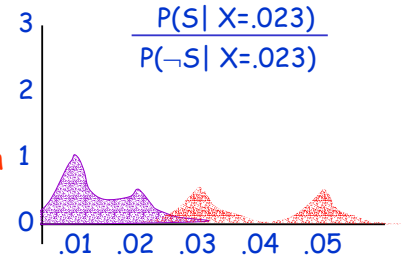
29

Spam?

#	X	S?
1	0.01	F
2	0.01	F
3	0.02	F
4	0.03	T
5	0.05	T

$$\frac{P(S | X=.023)}{P(\neg S | X=.023)}$$

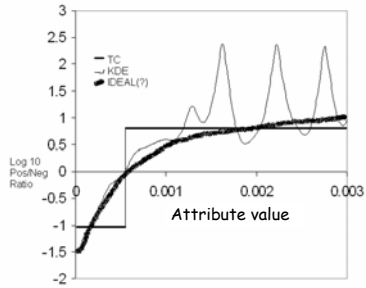
What's with
the shape?



© Daniel S. Weld

30

Analysis



© Daniel S. Weld

31

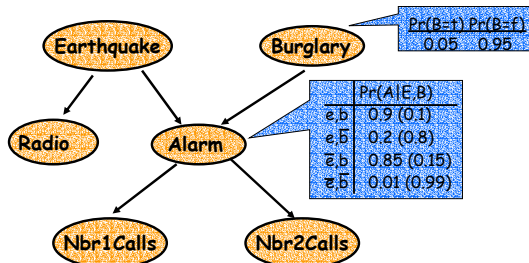
Topics

- Test & Mini Projects
- Review
- Naive Bayes
- Expectation Maximization
 - Review: Learning Bayesian Networks
 - Parameter Estimation
 - Structure Learning
 - Hidden Nodes
- Challenge

© Daniel S. Weld

32

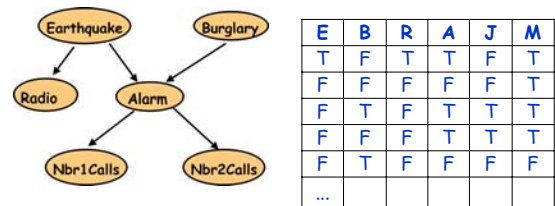
An Example Bayes Net



© Daniel S. Weld

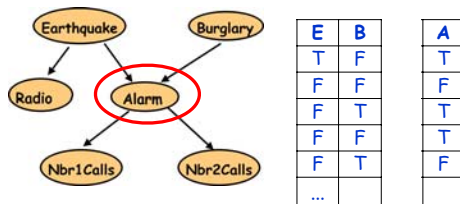
33

Parameter Estimation and Bayesian Networks



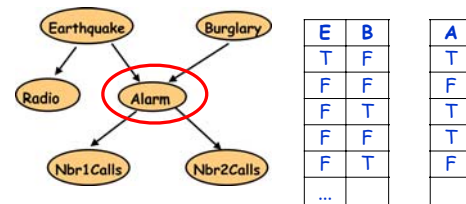
- We have:
- Bayes Net structure and observations
 - We need: Bayes Net parameters

Parameter Estimation and Bayesian Networks

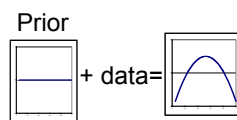


- $P(A|E,B) = ?$
 $P(A|E, \neg B) = ?$
 $P(A|\neg E,B) = ?$
 $P(A|\neg E, \neg B) = ?$

Parameter Estimation and Bayesian Networks



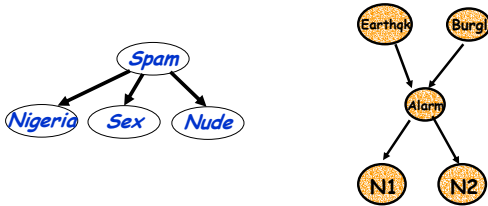
- $P(A|E,B) = ?$
 $P(A|E, \neg B) = ?$
 $P(A|\neg E,B) = ?$
 $P(A|\neg E, \neg B) = ?$



Now compute either MAP or Bayesian estimate

Recap

- Given a BN structure (with discrete or continuous variables), we can learn the parameters of the conditional prop tables.



© Daniel S. Weld

37

What if we *don't* know structure?

Learning The Structure of Bayesian Networks

- Search thru the space... of possible network structures! (for now, assume we observe all variables)
- For each structure, learn parameters
- Pick the one that fits observed data best
Caveat - won't we end up fully connected????
- When scoring, add a penalty \propto model complexity

problem!?!?

Learning The Structure of Bayesian Networks

- Search thru the space
- For each structure, learn parameters
- Pick the one that fits observed data best
- Problem?
 - Exponential number of networks!
 - And we need to learn parameters for each!
 - Exhaustive search out of the question!
- So what now?

Learning The Structure of Bayesian Networks

Local search!

Start with some network structure
Try to make a change
(add or delete or reverse edge)
See if the new network is any better

What should be the initial state?

Initial Network Structure?

- Uniform prior over random networks?
- Network which reflects expert knowledge?

Learning BN Structure

prior network+equivalent sample size

improved network(s)

data

X_1	X_2	X_3	
true	false	true	
false	false	true	
false	false	false	...
true	true	false	
⋮	⋮	⋮	⋮

© Daniel S. Weld 43

The Big Picture

- We described how to do MAP (and ML) learning of a Bayes net (including structure)
- How would Bayesian learning (of BNs) differ?
 - Find all possible networks
 - Calculate their posteriors
 - When doing inference, return weighed combination of predictions from all networks!

Hidden Variables

- But we can't observe the disease variable
- Can't we learn without it?

© Daniel S. Weld 45

We -could-

- But we'd get a fully-connected network

- With 708 parameters (vs. 78)
Much harder to learn!

© Daniel S. Weld 46

Chicken & Egg Problem

- If we knew that a training instance (patient) had the disease...
It would be easy to learn $P(\text{symptom} \mid \text{disease})$
But we can't observe disease, so we don't.
- If we knew params, e.g. $P(\text{symptom} \mid \text{disease})$ then it'd be easy to estimate if the patient had the disease.
But we don't know these parameters.

© Daniel S. Weld 47

Expectation Maximization (EM)

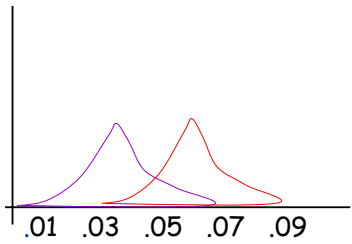
(high-level version)

- Pretend we **do** know the parameters
Initialize randomly
- **[E step]** Compute probability of instance having each possible value of the hidden variable
- **[M step]** Treating each instance as fractionally having **both** values compute the new parameter values
- Iterate until convergence!

© Daniel S. Weld 48

Simplest Version

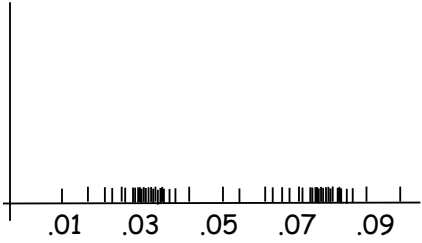
- Mixture of two distributions



- Know: form of distribution & variance, $\% = 5$
- Just need *mean* of each distribution

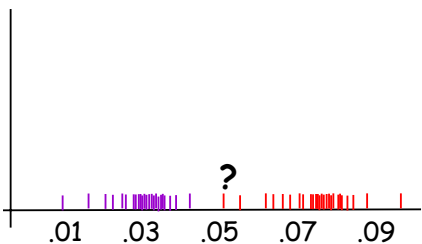
© Daniel S. Weld 49

Input Looks Like



© Daniel S. Weld 50

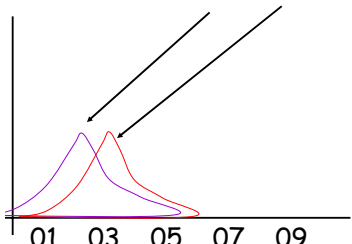
We Want to Predict



© Daniel S. Weld 51

Expectation Maximization (EM)

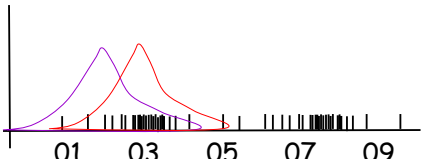
- Pretend we *do* know the parameters
Initialize randomly: set $\theta_1 = ?$; $\theta_2 = ?$



© Daniel S. Weld 52

Expectation Maximization (EM)

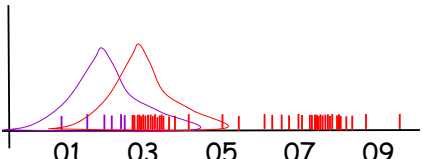
- Pretend we *do* know the parameters
Initialize randomly
- [E step]** Compute probability of instance having each possible value of the hidden variable



© Daniel S. Weld 53

Expectation Maximization (EM)

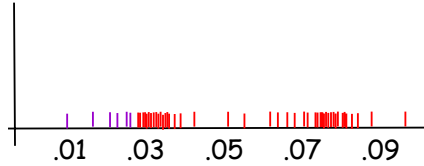
- Pretend we *do* know the parameters
Initialize randomly
- [E step]** Compute probability of instance having each possible value of the hidden variable



© Daniel S. Weld 54

Expectation Maximization (EM)

- Pretend we **do** know the parameters
Initialize randomly
- **[E step]** Compute probability of instance having each possible value of the hidden variable
- **[M step]** Treating each instance as fractionally having **both** values compute the new parameter values

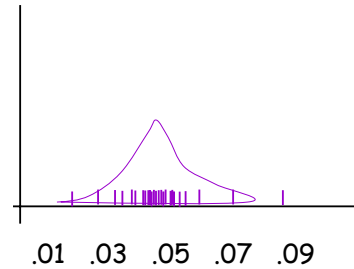


© Daniel S. Weld

55

ML Mean of Single Gaussian

$$U_{ml} = \operatorname{argmin}_u \sum_i (x_i - u)^2$$

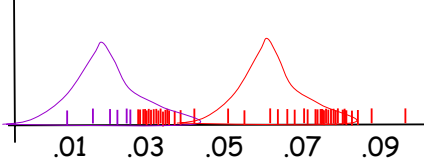


© Daniel S. Weld

56

Expectation Maximization (EM)

- Pretend we **do** know the parameters
Initialize randomly
- **[E step]** Compute probability of instance having each possible value of the hidden variable
- **[M step]** Treating each instance as fractionally having **both** values compute the new parameter values



© Daniel S. Weld

57

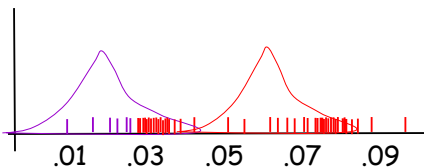
Iterate

© Daniel S. Weld

58

Expectation Maximization (EM)

- **[E step]** Compute probability of instance having each possible value of the hidden variable

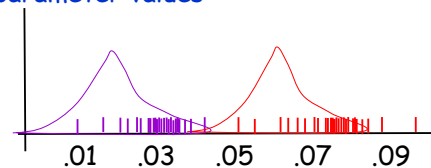


© Daniel S. Weld

59

Expectation Maximization (EM)

- **[E step]** Compute probability of instance having each possible value of the hidden variable
- **[M step]** Treating each instance as fractionally having **both** values compute the new parameter values

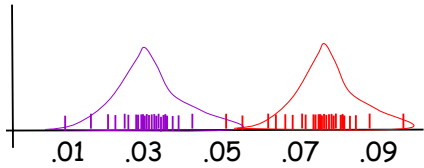


© Daniel S. Weld

60

Expectation Maximization (EM)

- [E step] Compute probability of instance having each possible value of the hidden variable
- [M step] Treating each instance as fractionally having *both* values compute the new parameter values



© Daniel S. Weld

61

Until Convergence

- Problems
 - Need to assume form of distribution
 - Local Maxima
- But
 - It really works in practice!
 - Can easily extend to multiple variables
 - E.g. Mean & Variance
 - Or much more complex models...

© Daniel S. Weld

62

Crossword Puzzles

© Daniel S. Weld

63