

# CS 188: Artificial Intelligence

## Spring 2011

Final Review  
5/2/2011

Pieter Abbeel – UC Berkeley

## Probabilistic Reasoning

---

- Probability
  - Random Variables
  - Joint and Marginal Distributions
  - Conditional Distribution
  - Inference by Enumeration
  - Product Rule, Chain Rule, Bayes' Rule
  - Independence
- Distributions over LARGE Numbers of Random Variables → Bayesian networks
  - Representation
  - Inference
    - Exact: Enumeration, Variable Elimination
    - Approximate: Sampling
  - Learning
    - Maximum likelihood parameter estimation
    - Laplace smoothing
    - Linear interpolation

2

# Probability recap

- Conditional probability  $P(x|y) = \frac{P(x,y)}{P(y)}$  ↗ ALWAYS
- Product rule  $P(x,y) = P(x|y)P(y)$  ↖
- Chain rule  $P(X_1, X_2, \dots, X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1, X_2) \dots$  ↖
- X, Y independent iff:  $\forall x, y : P(x, y) = P(x)P(y)$  ↖
- equivalently, iff:  $\forall x, y : P(x|y) = P(x)$  ↖
- equivalently, iff:  $\forall x, y : P(y|x) = P(y)$  ↖
- X and Y are conditionally independent given Z iff: ↖
- $\forall x, y, z : P(x, y|z) = P(x|z)P(y|z)$  ↖
- equivalently, iff:  $\forall x, y, z : P(x|y, z) = P(x|z)$  ↖
- equivalently, iff:  $\forall x, y, z : P(y|x, z) = P(y|z)$  ↖

3

# Inference by Enumeration

- P(sun)?

$$0.15 + 0.10 + 0.10 + 0.30$$

- P(sun | winter)?

$$\frac{P(\text{sun, winter})}{P(\text{winter})} = \frac{0.10 + 0.15}{0.10 + 0.20 + 0.15 + 0.20}$$

- P(sun | winter, hot)?

$$\frac{P(\text{sun, winter, hot})}{P(\text{winter, hot})}$$

S	T	W	P
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

4

# Chain Rule → Bayes net

- Chain rule: can **always** write **any** joint distribution as an incremental product of conditional distributions

$$P(x_1, x_2, x_3) = P(x_1)P(x_2|x_1)P(x_3|x_1, x_2)$$

$$P(x_1, x_2, \dots, x_n) = \prod_i P(x_i|x_1 \dots x_{i-1})$$

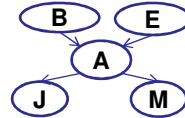
B, E  
A  
J, M

- Bayes nets: make conditional independence assumptions of the form:

$$P(x_i|x_1 \dots x_{i-1}) = P(x_i|\text{parents}(X_i))$$

giving us:

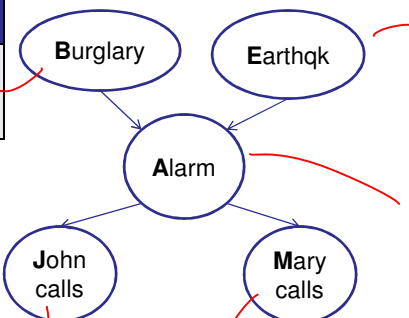
$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i|\text{parents}(X_i))$$



5

# Example: Alarm Network

B	P(B)
+b	0.001
-b	0.999



E	P(E)
+e	0.002
-e	0.998

A	J	P(J A)
+a	+j	0.9
+a	-j	0.1
-a	+j	0.05
-a	-j	0.95

A	M	P(M A)
+a	+m	0.7
+a	-m	0.3
-a	+m	0.01
-a	-m	0.99

B	E	A	P(A B,E)
+b	+e	+a	0.95
+b	+e	-a	0.05
+b	-e	+a	0.94
+b	-e	-a	0.06
-b	+e	+a	0.29
-b	+e	-a	0.71
-b	-e	+a	0.001
-b	-e	-a	0.999

## Size of a Bayes' Net for $P(X_1, X_2, \dots, X_n)$

- How big is a joint distribution over N Boolean variables?  
 $2^N$
- Size of representation if we use the chain rule  
 $2^N$
- How big is an N-node net if nodes have up to  $k$  parents?  
 $O(N * 2^{k+1})$
- Both give you the power to calculate
- BNs:
  - Huge space savings!
  - Easier to elicit local CPTs
  - Faster to answer queries

7

## Bayes Nets: Assumptions

- Assumptions we are required to make to define the Bayes net when given the graph:

$$P(x_i | x_1 \cdots x_{i-1}) = P(x_i | \text{parents}(X_i))$$

- Given a Bayes net graph additional conditional independences can be read off directly from the graph
- Question: Are two nodes *necessarily* independent given certain evidence?
- If no, can prove with a counter example
  - I.e., pick a set of CPT's, and show that the independence assumption is violated by the resulting distribution
- If yes, can prove with
  - Algebra (tedious)
  - D-separation (analyzes graph)

8



# D-Separation

?  $X \perp\!\!\!\perp Y \mid Z$

Question: Are X and Y conditionally independent given evidence vars {Z}?

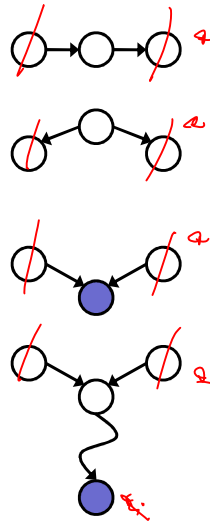
- Yes, if X and Y "separated" by Z
- Consider all (undirected) paths from X to Y
- No active paths = independence!

A path is active if each triple is active:

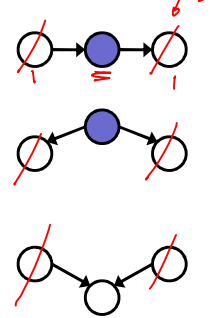
- Causal chain  $A \rightarrow B \rightarrow C$  where B is unobserved (either direction)
- Common cause  $A \leftarrow B \rightarrow C$  where B is unobserved
- Common effect (aka v-structure)  $A \rightarrow B \leftarrow C$  where B or one of its descendants is observed

All it takes to block a path is a single inactive segment

Active Triples



Inactive Triples



# D-Separation

Given query  $X_i \stackrel{?}{\perp\!\!\!\perp} X_j \mid \{X_{k_1}, \dots, X_{k_n}\}$

Shade all evidence nodes

For all (undirected!) paths between and

Check whether path is active

If active return  $X_i \not\perp\!\!\!\perp X_j \mid \{X_{k_1}, \dots, X_{k_n}\}$

(If reaching this point all paths have been checked and shown inactive)

Return  $X_i \perp\!\!\!\perp X_j \mid \{X_{k_1}, \dots, X_{k_n}\}$

## Example

$L \perp\!\!\!\perp T' | T$  ✓ Yes

$L \perp\!\!\!\perp B$  ✓ Yes

~~$L \perp\!\!\!\perp B | T$~~

~~$L \perp\!\!\!\perp B | T'$~~

$L \perp\!\!\!\perp B | T, R$  ✓ Yes

11

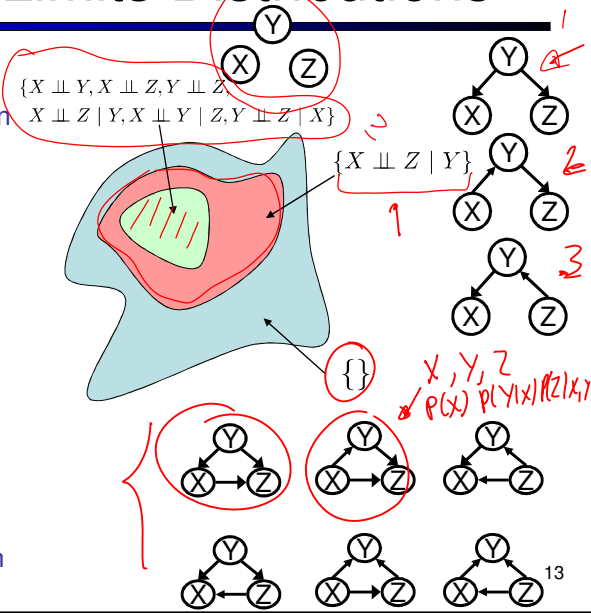
## All Conditional Independences

- Given a Bayes net structure, can run d-separation to build a complete list of conditional independences that are necessarily true of the form
 
$$X_i \perp\!\!\!\perp X_j | \{X_{k_1}, \dots, X_{k_n}\}$$
- This list determines the set of probability distributions that can be represented by Bayes' nets with this graph structure

12

# Topology Limits Distributions

- Given some graph topology  $G$ , only certain joint distributions can be encoded
- The graph structure guarantees certain (conditional) independences
- (There might be more independence)
- Adding arcs increases the set of distributions, but has several costs
- Full conditioning can encode any distribution



# Bayes Nets Status

## ✓ Representation

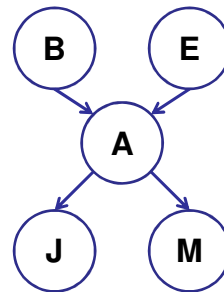
- Inference
- Learning Bayes Nets from Data

## Inference by Enumeration

---

- Given unlimited time, inference in BNs is easy
- Recipe:
  - State the marginal probabilities you need
  - Figure out ALL the atomic probabilities you need
  - Calculate and combine them
- Example:

$$P(+b | +j, +m) = \frac{P(+b, +j, +m)}{P(+j, +m)}$$



15

## Example: Enumeration

---

- In this simple method, we only need the BN to synthesize the joint entries

$$\begin{aligned}
 P(+b, +j, +m) = & \\
 & P(+b)P(+e)P(+a|+b, +e)P(+j|+a)P(+m|+a) + \\
 & P(+b)P(+e)P(-a|+b, +e)P(+j|-a)P(+m|-a) + \\
 & P(+b)P(-e)P(+a|+b, -e)P(+j|+a)P(+m|+a) + \\
 & P(+b)P(-e)P(-a|+b, -e)P(+j|-a)P(+m|-a)
 \end{aligned}$$

16



# Variable Elimination

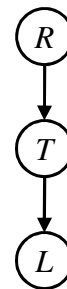
- Why is inference by enumeration so slow?
  - You join up the whole joint distribution before you sum out the hidden variables
  - You end up repeating a lot of work!
- Idea: interleave joining and marginalizing!
  - Called "Variable Elimination"
  - Still NP-hard, but usually much faster than inference by enumeration

17

# Variable Elimination Outline

- Track objects called **factors**
- Initial factors are local CPTs (one per node)

$P(R)$	$P(T R)$	$P(L T)$																												
<table border="1" style="display: inline-table;"><tr><td>+r</td><td>0.1</td></tr><tr><td>-r</td><td>0.9</td></tr></table>	+r	0.1	-r	0.9	<table border="1" style="display: inline-table;"><tr><td>+r</td><td>+t</td><td>0.8</td></tr><tr><td>+r</td><td>-t</td><td>0.2</td></tr><tr><td>-r</td><td>+t</td><td>0.1</td></tr><tr><td>-r</td><td>-t</td><td>0.9</td></tr></table>	+r	+t	0.8	+r	-t	0.2	-r	+t	0.1	-r	-t	0.9	<table border="1" style="display: inline-table;"><tr><td>+t</td><td>+l</td><td>0.3</td></tr><tr><td>+t</td><td>-l</td><td>0.7</td></tr><tr><td>-t</td><td>+l</td><td>0.1</td></tr><tr><td>-t</td><td>-l</td><td>0.9</td></tr></table>	+t	+l	0.3	+t	-l	0.7	-t	+l	0.1	-t	-l	0.9
+r	0.1																													
-r	0.9																													
+r	+t	0.8																												
+r	-t	0.2																												
-r	+t	0.1																												
-r	-t	0.9																												
+t	+l	0.3																												
+t	-l	0.7																												
-t	+l	0.1																												
-t	-l	0.9																												



- Any known values are selected
  - E.g. if we know  $L = +l$ , the initial factors are

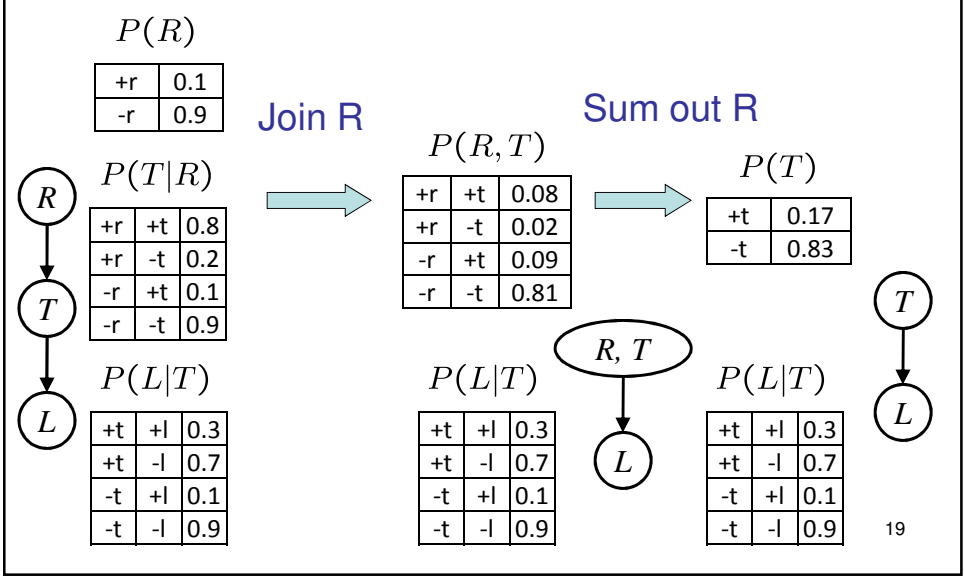
$P(R)$	$P(T R)$	$P(+l T)$																						
<table border="1" style="display: inline-table;"><tr><td>+r</td><td>0.1</td></tr><tr><td>-r</td><td>0.9</td></tr></table>	+r	0.1	-r	0.9	<table border="1" style="display: inline-table;"><tr><td>+r</td><td>+t</td><td>0.8</td></tr><tr><td>+r</td><td>-t</td><td>0.2</td></tr><tr><td>-r</td><td>+t</td><td>0.1</td></tr><tr><td>-r</td><td>-t</td><td>0.9</td></tr></table>	+r	+t	0.8	+r	-t	0.2	-r	+t	0.1	-r	-t	0.9	<table border="1" style="display: inline-table;"><tr><td>+t</td><td>+l</td><td>0.3</td></tr><tr><td>-t</td><td>+l</td><td>0.1</td></tr></table>	+t	+l	0.3	-t	+l	0.1
+r	0.1																							
-r	0.9																							
+r	+t	0.8																						
+r	-t	0.2																						
-r	+t	0.1																						
-r	-t	0.9																						
+t	+l	0.3																						
-t	+l	0.1																						

- VE: Alternately join factors and eliminate variables

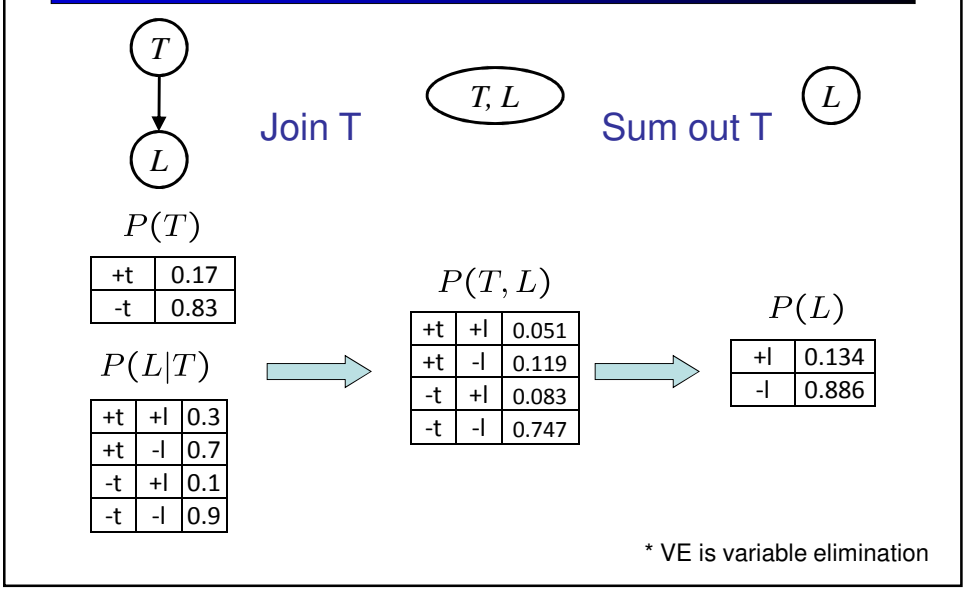
18

P(L)

# Variable Elimination Example



# Variable Elimination Example



## Example

$$P(B|j, m) \propto P(B, j, m)$$

$P(B)$	$P(E)$	$P(A B, E)$	$P(j A)$	$P(m A)$
--------	--------	-------------	----------	----------

Choose A

$$\begin{array}{l}
 P(A|B, E) \\
 P(j|A) \\
 P(m|A)
 \end{array}
 \xrightarrow{\times} P(j, m, A|B, E)
 \xrightarrow{\Sigma} P(j, m|B, E)$$

$2 \times 2$

$P(B)$	$P(E)$	$P(j, m B, E)$
--------	--------	----------------

21

## Example

$P(B)$	$P(E)$	$P(j, m B, E)$
--------	--------	----------------

Choose E

$$\begin{array}{l}
 P(E) \\
 P(j, m|B, E)
 \end{array}
 \xrightarrow{\times} P(j, m, E|B)
 \xrightarrow{\Sigma} P(j, m|B)$$

$P(B)$	$P(j, m B)$
--------	-------------

Finish with B

$$\begin{array}{l}
 P(B) \\
 P(j, m|B)
 \end{array}
 \xrightarrow{\times} P(j, m, B)
 \xrightarrow{\text{Normalize}} P(B|j, m)$$

$P(B|j, m)$

22

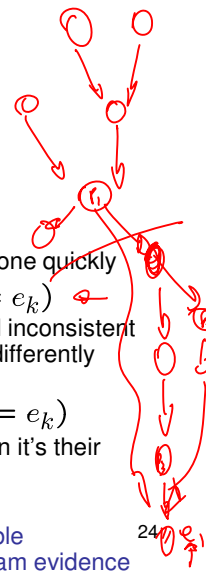
# General Variable Elimination

- Query:  $P(Q|E_1 = e_1, \dots, E_k = e_k)$
- Start with initial factors:
  - Local CPTs (but instantiated by evidence)
- While there are still hidden variables (not Q or evidence):
  - Pick a hidden variable H
  - Join all factors mentioning H
  - Eliminate (sum out) H
- Join all remaining factors and normalize

23

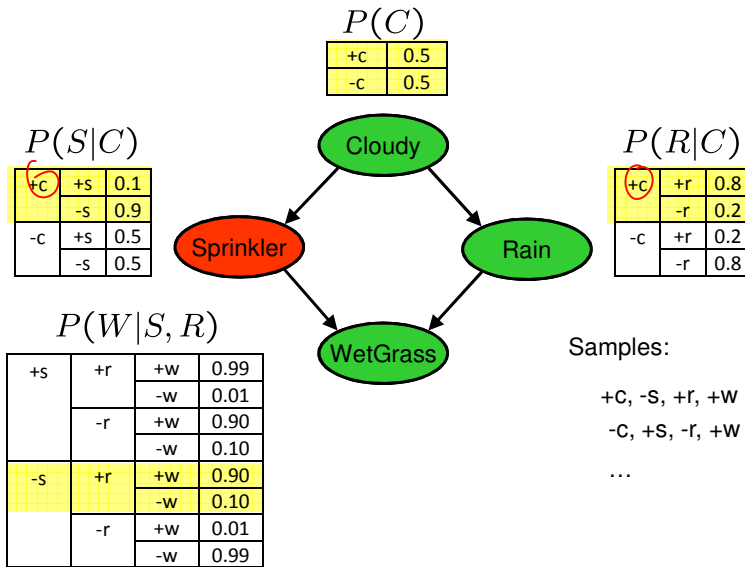
# Approximate Inference: Sampling

- Basic idea:
  - Draw  $N$  samples from a sampling distribution  $S$
  - Compute an approximate posterior probability
  - Show this converges to the true probability  $P$
- Why? Faster than computing the exact answer
- Prior sampling:
  - Sample ALL variables in topological order as this can be done quickly
- Rejection sampling for query  $P(Q|E_1 = e_1, \dots, E_k = e_k)$ 
  - = like prior sampling, but reject when a variable is sampled inconsistent with the query, in this case when a variable  $E_i$  is sampled differently from  $e_i$
- Likelihood weighting for query  $P(Q|E_1 = e_1, \dots, E_k = e_k)$ 
  - = like prior sampling but variables  $E_i$  are not sampled, when it's their turn, they get set to  $e_i$ , and the sample gets weighted by  $P(e_i | \text{value of parents}(e_i) \text{ in current sample})$
- Gibbs sampling: repeatedly samples each non-evidence variable conditioned on all other variables  $\rightarrow$  can incorporate downstream evidence



24

# Prior Sampling

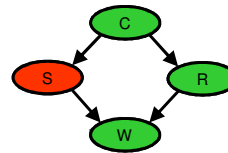


25

# Example

- We'll get a bunch of samples from the BN:

+C, -S, +r, +W  
 +C, +S, +r, +W  
 -C, +S, +r, -W  
 +C, -S, +r, +W  
 -C, -S, -r, +W



- If we want to know  $P(W)$

- We have counts  $\langle +w:4, -w:1 \rangle$
- Normalize to get  $P(W) = \langle +w:0.8, -w:0.2 \rangle$
- This will get closer to the true distribution with more samples
- Can estimate anything else, too
- What about  $P(C|+w)$ ?  $P(C|+r, +w)$ ?  $P(C|-r, -w)$ ?
- Fast: can use fewer samples if less time

$P(C|+w)$   
 $P(+w)$

26

# Likelihood Weighting

$P(C)$

+c	0.5
-c	0.5

$P(S|C)$

+c	+s	0.1
-c	-s	0.9
-c	+s	0.5
	-s	0.5

$P(R|C)$

+c	+r	0.8
-c	-r	0.2
-c	+r	0.2
	-r	0.8

$P(W|S, R)$

+s	+r	+w	0.99
	-r	-w	0.01
		+w	0.90
-s	+r	-w	0.10
		+w	0.90
	-r	+w	0.10
		-w	0.99

Samples:

+c, +s, +r, +w

...

$w = 1.0 \times 0.1 \times 0.99$

27

# Likelihood Weighting

- Sampling distribution if  $z$  sampled and  $e$  fixed evidence

$$S_{WS}(z, e) = \prod_{i=1}^l P(z_i | \text{Parents}(Z_i))$$

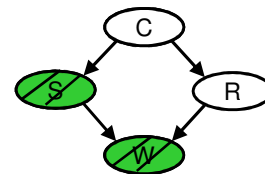
- Now, samples have weights

$$w(z, e) = \prod_{i=1}^m P(e_i | \text{Parents}(E_i))$$

- Together, weighted sampling distribution is consistent

$$S_{WS}(z, e) \cdot w(z, e) = \prod_{i=1}^l P(z_i | \text{Parents}(z_i)) \prod_{i=1}^m P(e_i | \text{Parents}(e_i))$$

$$= P(z, e)$$



28

# Gibbs Sampling

*+x, +y, +z, +w*  
*Resample X from  $P(X|+y, +z, +w) \rightarrow -x$*

- Idea: instead of sampling from scratch, create samples that are each like the last one.
- Procedure: resample one variable at a time, conditioned on all the rest, but keep evidence fixed.
- Properties: Now samples are not independent (in fact they're nearly identical), but sample averages are still consistent estimators!
- What's the point: both upstream and downstream variables condition on evidence.

*Resample Y from  $P(Y|-x, +z, +w) \rightarrow +y$*

*Resample Z from  $P(Z|-x, +y, +w) \rightarrow -z$*

*Resample X from  $P(X|+y, -z, +w) \rightarrow +x$*

*Handwritten notes and diagrams include:*  
 $P(-x, +y, +z, +w)$   
 $\frac{P(-x, +y, +z, +w)}{P(-x, +y, +w)}$   
 $\frac{\sum P(-x|+y)}{\sum P(z|-x, +y)}$   
 $\frac{P(+w|z)}$   
 $P(-x)P(+y)$   
 $P(z|-x, +y)$   
 $P(+w|z)$

29

# Markov Models

- A Markov model is a chain-structured BN
  - Each node is identically distributed (stationarity)
  - Value of X at a given time is called the **state**
  - As a BN:

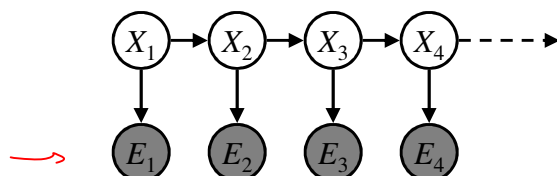
$P(X_1)$        $P(X|X_{-1})$

- The chain is just a (growing) BN
  - We can always use generic BN reasoning on it if we truncate the chain at a fixed length
- Stationary distributions
  - For most chains, the distribution we end up in is independent of the initial distribution
  - Called the **stationary distribution** of the chain
- Example applications: Web link analysis (Page Rank) and Gibbs Sampling

*Handwritten notes include:*  
 $P_{\infty}$   
 $P_{\infty}(X_{\infty}) = \sum_x P(X_{\infty}|x) P(x)$

## Hidden Markov Models

- Underlying Markov chain over states  $S$
- You observe outputs (effects) at each time step



- Speech recognition HMMs:
  - $X_i$ : specific positions in specific words;  $E_i$ : acoustic signals
- Machine translation HMMs:
  - $X_i$ : translation options;  $E_i$ : Observations are words
- Robot tracking:
  - $X_i$ : positions on a map;  $E_i$ : range readings

## Online Belief Updates

- Every time step, we start with current  $P(X | \text{evidence})$
- We update for time:

$$P(x_t | e_{1:t-1}) = \sum_{x_{t-1}} P(x_{t-1} | e_{1:t-1}) \cdot P(x_t | x_{t-1})$$

The diagram shows two white circles,  $X_1$  and  $X_2$ , with a solid arrow pointing from  $X_1$  to  $X_2$ .

- We update for evidence:

$$P(x_t | e_{1:t}) \propto_X P(x_t | e_{1:t-1}) \cdot P(e_t | x_t)$$

The diagram shows a white circle  $X_2$  with a solid arrow pointing down to a grey circle  $E_2$ .

- The forward algorithm does both at once (and doesn't normalize)

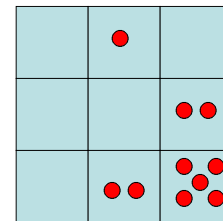
*↳ = variable elimination in  $X_1, \dots, X_n$*



# Particle Filtering

- = likelihood weighting + resampling at each time slice
- Why: sometimes  $|X|$  is too big to use exact inference
- Particle is just new name for sample

0.0	0.1	0.0
0.0	0.0	0.2
0.0	0.2	0.5

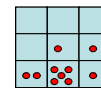
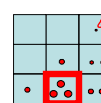
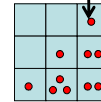
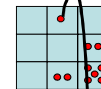


# Particle Filtering

- Elapse time:
  - Each particle is moved by sampling its next position from the transition model

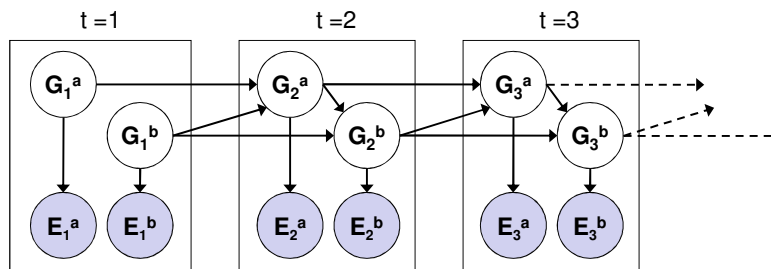
$$x' = \text{sample}(P(X'|x))$$

- Observe:
  - We don't sample the observation, we fix it and downweight our samples based on the evidence
  - This is like likelihood weighting, so we multiply with  $P(e_t|x')$
- Resample:
  - Rather than tracking weighted samples, we resample
  - N times, we choose from our weighted sample distribution



## Dynamic Bayes Nets (DBNs)

- We want to track multiple variables over time, using multiple sources of evidence
- Idea: Repeat a fixed Bayes net structure at each time
- Variables from time  $t$  can condition on those from  $t-1$



- Discrete valued dynamic Bayes nets are also HMMs

## Bayes Nets Status

- ✓ Representation
- ✓ Inference
- Learning Bayes Nets from Data

X	P
r	$\theta_r$
g	$\theta_g$
b	$\theta_b$

## Parameter Estimation

X	P
r	$\theta$
g	$1-\theta$

- Estimating distribution of random variables like X or X | Y
- **Empirically:** use training data
  - For each outcome x, look at the *empirical rate* of that value:  $P_{\theta}(x_i)$

$$P_{ML}(x) = \frac{\text{count}(x)}{\text{total samples}}$$

r	g	g
---	---	---

$P_{ML}(r) = 1/3$

- This is the estimate that maximizes the *likelihood of the data*

$$L(x, \theta) = \prod_i P_{\theta}(x_i) \quad \theta \cdot (1-\theta) \cdot (1-\theta)$$

- **Laplace smoothing**
  - Pretend saw every outcome k extra times  $P_{LAP,k}(x) = \frac{c(x) + k}{N + k|X|}$
  - Smooth each condition independently:  $P_{LAP,k}(x|y) = \frac{c(x, y) + k}{c(y) + k|X|}$

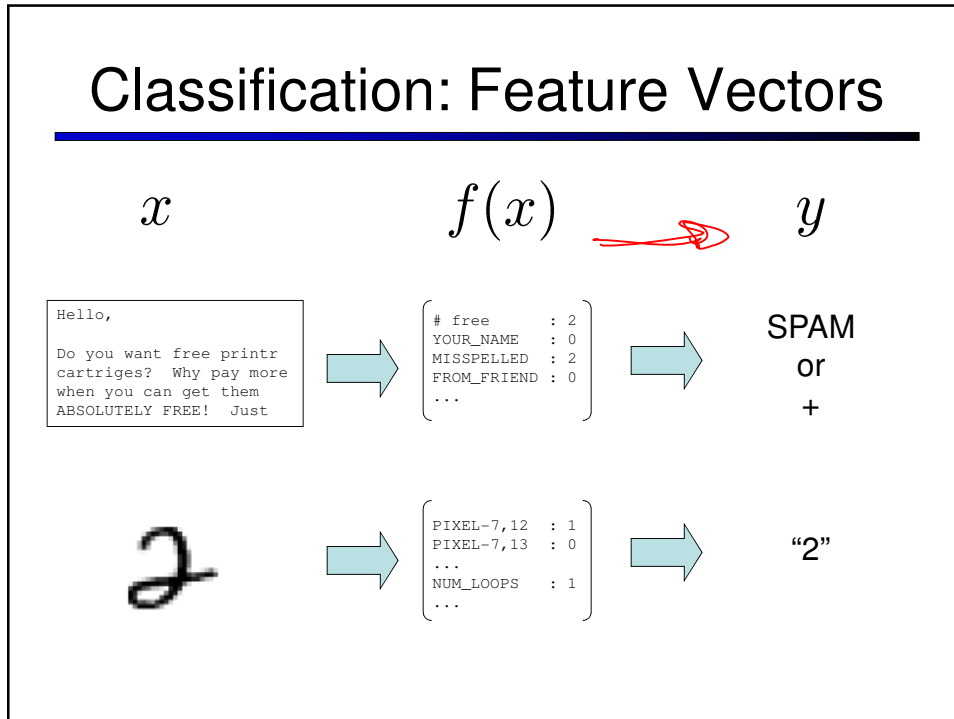
## Bayes Nets Status

---

- ✓ Representation
- ✓ Inference
- ✓ Learning Bayes Nets from Data

38

# Classification: Feature Vectors



# Classification overview

- **Naïve Bayes:**
  - Builds a model training data
  - Gives prediction probabilities
  - Strong assumptions about feature independence
  - One pass through data (counting)
- **Perceptron:**
  - Makes less assumptions about data
  - Mistake-driven learning
  - Multiple passes through data (prediction)
  - Often more accurate
- **MIRA:**
  - Like perceptron, but adaptive scaling of size of update
- **SVM:**
  - Properties similar to perceptron
  - Convex optimization formulation
- **Nearest-Neighbor:**
  - Non-parametric: more expressive with more training data
- **Kernels**
  - Efficient way to make linear learning architectures into nonlinear ones

# Bayes Nets for Classification

- One method of classification:
  - Use a probabilistic model!
  - Features are observed random variables  $F_i$
  - $Y$  is the query variable
  - Use probabilistic inference to compute most likely  $Y$

$$y = \operatorname{argmax}_y P(y|f_1 \dots f_n)$$

- You already know how to do this inference

# General Naïve Bayes

- A general *naive Bayes* model:

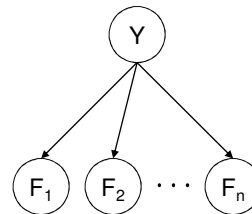
$|Y| \times |F|^n$   
parameters

$$P(Y, F_1 \dots F_n) =$$

$$P(Y) \prod_i P(F_i|Y)$$

$|Y|$  parameters

$n \times |F| \times |Y|$   
parameters



- We only specify how each feature depends on the class
- Total number of parameters is *linear* in  $n$
- Our running example: digits

# Bag-of-Words Naïve Bayes

- Generative model

$$P(Y, W_1 \dots W_n) = P(Y) \prod_i P(W_i | Y)$$

Word at position  $i$ , not  $i^{\text{th}}$  word in the dictionary!

- Bag-of-words

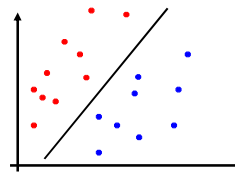
- Each position is identically distributed
- All positions share the same conditional probs  $P(W|C)$
- When learning the parameters, data is shared over all positions in the document (rather than separately learning a distribution for each position in the document)

- Our running example: spam vs. ham

# Linear Classifier

- Binary linear classifier:

$$y = \begin{cases} +1 & \text{if } w \cdot f(x) \geq 0 \\ -1 & \text{if } w \cdot f(x) < 0 \end{cases}$$



- Multiclass linear classifier:

- A weight vector for each class:

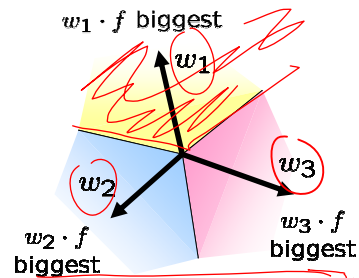
$$w_y$$

- Score (activation) of a class  $y$ :

$$w_y \cdot f(x)$$

- Prediction highest score wins

$$y = \arg \max_y w_y \cdot f(x)$$



Binary = multiclass where the negative class has weight zero

## Perceptron = algorithm to learn weights $w$

- Start with zero weights
- Pick up training instances one by one
- Classify with current weights

$$\begin{aligned} \rightarrow y &= \arg \max_y w_y \cdot f(x) \\ &= \arg \max_y \sum_i w_{y,i} \cdot f_i(x) \end{aligned}$$

- If correct, no change!
- If wrong: lower score of wrong answer, raise score of right answer

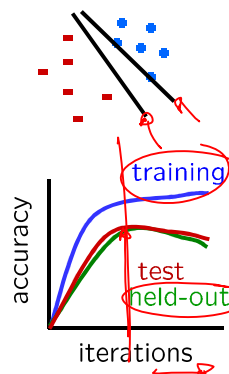
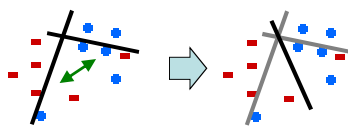
$$w_y = w_y - f(x)$$

$$w_{y^*} = w_{y^*} + f(x)$$

$$\begin{aligned} (w_y - f(x)) \cdot f(x) &\leq w_y \cdot f(x) \\ (w_{y^*} + f(x)) \cdot f(x) &\geq w_{y^*} \cdot f(x) \end{aligned}$$

## Problems with the Perceptron

- Noise: if the data isn't separable, weights might thrash
  - Averaging weight vectors over time can help (averaged perceptron)
- Mediocre generalization: finds a "barely" separating solution
- Overtraining: test / held-out accuracy usually rises, then falls
  - Overtraining is a kind of overfitting



## Fixing the Perceptron: MIRA

- Update size that fixes the current mistake and also minimizes the change to  $w$

Guessed  $y$  instead of  $y^*$  on example  $x$  with features  $f(x)$

- Update  $w$  by solving:

$$\min_{w, \tau} \frac{1}{2} \sum_y \|w_y - w'_y\|^2$$

$$\rightarrow w_{y^*} \cdot f(x) \geq w_y \cdot f(x) + 1$$

$$\begin{cases} w_y = w'_y - \tau f(x) \\ w_{y^*} = w'_y + \tau f(x) \end{cases}$$

$$0 \leq \tau \leq C$$

$$\|w\|^2 = w \cdot w$$

$$\|w\| = \sqrt{w \cdot w}$$

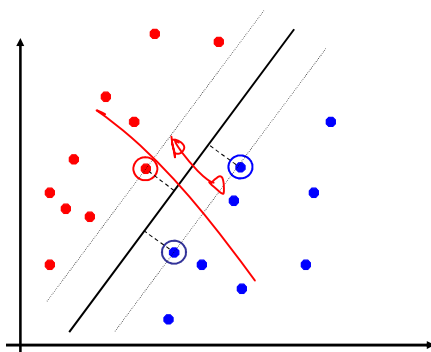
SVM

$$\min_w \frac{1}{2} \|w\|^2$$

$$\forall i, y \quad w_{y^*} \cdot f(x_i) \geq w_y \cdot f(x_i) + 1$$

## Support Vector Machines

- Maximizing the margin:** good according to intuition, theory, practice
- Support vector machines (SVMs) find the separator with max margin
- Basically, SVMs are MIRA where you optimize over all examples at once



SVM

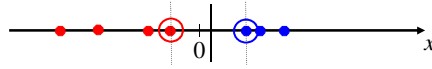
$$\min_w \frac{1}{2} \|w\|^2$$

$$\forall i, y \quad w_{y^*} \cdot f(x_i) \geq w_y \cdot f(x_i) + 1$$

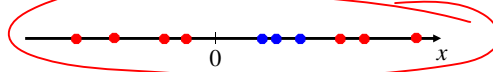


# Non-Linear Separators

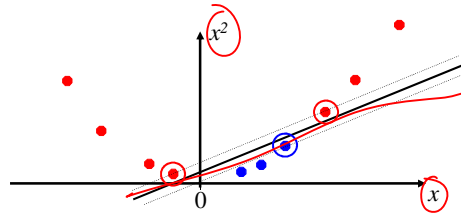
- Data that is linearly separable (with some noise) works out great:



- But what are we going to do if the dataset is just too hard?



- How about... mapping data to a higher-dimensional space:

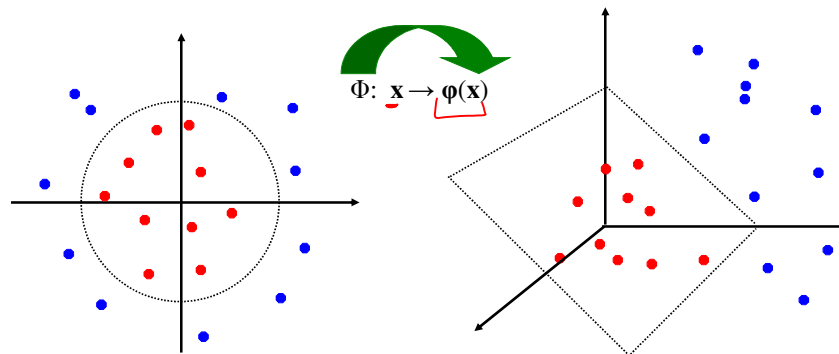


49

This and next few slides adapted from Ray Mooney, UT

# Non-Linear Separators

- General idea: the original feature space can always be mapped to some higher-dimensional feature space where the training set is separable:



50

## Some Kernels

- Kernels **implicitly** map original vectors to higher dimensional spaces, take the dot product there, and hand the result back

- Linear kernel:  $K(x, x') = \underline{x \cdot x'} = \sum_i x_i x'_i$

$$\phi(x) = x$$

- Quadratic kernel:  $K(x, x') = \underline{(x \cdot x' + 1)^2} = \underline{\phi(x) \cdot \phi(x')}$

$$= \sum x_i x_j x'_i x'_j + 2 \sum x_i x'_i + 1$$

For  $x \in \mathbb{R}^3$ :

$$\rightarrow \phi(x) = [x_1 x_1 \quad x_1 x_2 \quad x_1 x_3 \quad x_2 x_1 \quad x_2 x_2 \quad x_2 x_3 \quad x_3 x_1 \quad x_3 x_2 \quad x_3 x_3 \quad \sqrt{2} x_1 \quad \sqrt{2} x_2 \quad \sqrt{2} x_3 \quad 1]$$

## Some Kernels (2)

- Polynomial kernel:  $K(x, x') = (x \cdot x' + 1)^d$

$$= \phi(x) \cdot \phi(x')$$

For  $x \in \mathbb{R}^3$ :

$$\phi(x) = [x_1^d \quad x_2^d \quad x_3^d \quad \sqrt{d} x_1^{d-1} x_2 \quad \sqrt{d} x_1^{d-1} x_3 \quad \dots \quad \sqrt{d} x_1 \quad \sqrt{d} x_2 \quad \sqrt{d} x_3 \quad 1]$$

For  $x \in \mathbb{R}^n$  the  $d$ -order polynomial kernel's implicit feature space is  $\binom{n+d}{d}$  dimensional.

By contrast, computing the kernel directly only requires  $O(n)$  time.

# Why and When Kernels?

- Can't you just add these features on your own (e.g. add all pairs of features instead of using the quadratic kernel)?
  - Yes, in principle, just compute them
  - No need to modify any algorithms
  - But, number of features can get large (or infinite)
- Kernels let us compute with these features implicitly
  - Example: implicit dot product in polynomial, Gaussian and string kernel takes much less space and time per dot product
- When can we use kernels?
  - When our learning algorithm can be reformulated in terms of only inner products between feature vectors
    - Examples: perceptron, support vector machine

# K-nearest neighbors

- 1-NN: copy the label of the most similar data point
- K-NN: let the k nearest neighbors vote (have to devise a weighting scheme)



2 Examples



10 Examples



100 Examples



10000 Examples



Truth

- Parametric models:
  - Fixed set of parameters
  - More data means better settings
- Non-parametric models:
  - Complexity of the classifier increases with data
  - Better in the limit, often worse in the non-limit
- (K)NN is non-parametric

# Basic Similarity

- Many similarities based on **feature dot products**:

$$\text{sim}(x, x') = f(x) \cdot f(x') = \sum_i f_i(x) f_i(x')$$

*K(x, x')*

- If features are just the pixels:

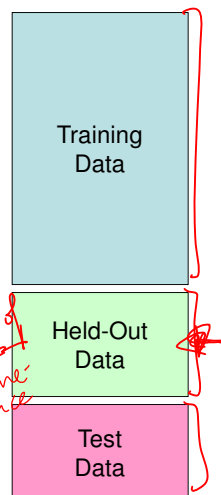
$$\text{sim}(x, x') = x \cdot x' = \sum_i x_i x'_i$$

- Note: not all similarities are of this form

55

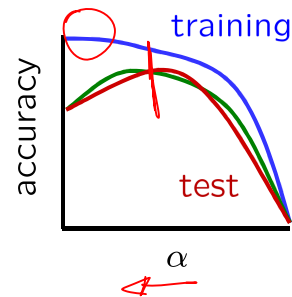
# Important Concepts

- Data: labeled instances, e.g. emails marked spam/ham
  - Training set
  - Held out set
  - Test set
- Features: attribute-value pairs which characterize each  $x$
- Experimentation cycle
  - Learn parameters (e.g. model probabilities) on training set
  - (Tune hyperparameters on held-out set)
  - Compute accuracy of test set
  - Very important: never “peek” at the test set!
- Evaluation
  - Accuracy: fraction of instances predicted correctly
- Overfitting and generalization
  - Want a classifier which does well on *test* data
  - Overfitting: fitting the training data very closely, but not generalizing well
  - We’ll investigate overfitting and generalization formally in a few lectures



# Tuning on Held-Out Data

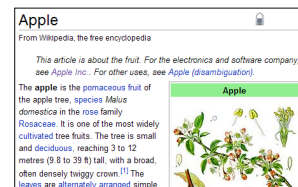
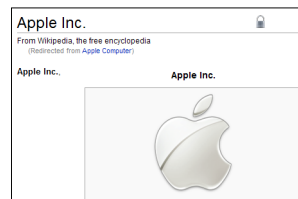
- Now we've got two kinds of unknowns
  - Parameters: the probabilities  $P(Y|X)$ ,  $P(Y)$
  - Hyperparameters:
    - Amount of smoothing to do:  $k$ ,  $\alpha$  (naïve Bayes)
    - Number of passes over training data (perceptron)
- Where to learn?
  - Learn parameters from training data
  - Must tune hyperparameters on different data
  - For each value of the hyperparameters, train and test on the held-out data
  - Choose the best value and do a final test on the test data



# Extension: Web Search

- Information retrieval:
  - Given information needs, produce information
  - Includes, e.g. web search, question answering, and classic IR
- Web search: not exactly classification, but rather ranking


$x = \text{"Apple Computers"}$



# Feature-Based Ranking


$x = \text{"Apple Computers"}$

$f(x, \text{Apple}) = [0.3 \ 5 \ 0 \ 0 \ \dots]$



PR #1 #2 #3 #4  
 ↓ ↓ ↓ ↓

$f(x, \text{Apple Inc.}) = [0.8 \ 4 \ 2 \ 1 \ \dots]$



# Perceptron for Ranking

- Inputs  $x$
- Candidates  $y$
- Many feature vectors:  $f(x, y)$
- One weight vector:  $w$

- Prediction:

$$y = \arg \max_y \underline{w} \cdot \underline{f(x, y)}$$

- Update (if wrong):

$$w = \underline{w} + \underline{f(x, y^*)} - \underline{f(x, y)}$$