

CS 2001 – Lecture 2

Bayesian belief networks

Milos Hauskrecht

milos@cs.pitt.edu

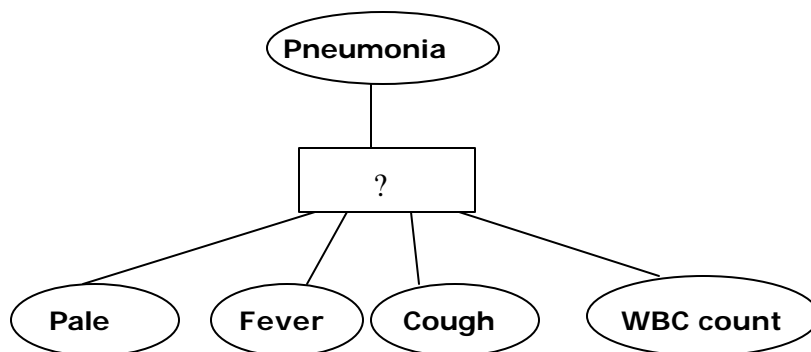
5329 Sennott Square

X4-8845

CS 2001 Bayesian belief networks

Modeling the uncertainty.

- How to describe, represent the relations in the presence of uncertainty?
- How to manipulate such knowledge to make inferences?



CS 2001 Bayesian belief networks

Probability theory

a well-defined coherent theory for representing uncertainty and for reasoning with it

Representation:

Proposition statements – assignment of values to random variables

$$Pneumonia = True \quad WBCcount = high$$

Probabilities over statements model the degree of belief in these statements

$$P(Pneumonia = True) = 0.001$$

$$P(WBCcount = high) = 0.005$$

$$P(Pneumonia = True, Fever = True) = 0.0009$$

$$P(Pneumonia = False, WBCcount = normal, Cough = False) = 0.97$$

CS 2001 Bayesian belief networks

Modeling uncertainty with probabilities

- **Full joint distribution:** joint distribution over all random variables defining the domain
 - it is sufficient to represent the complete domain and to do any type of probabilistic reasoning

Problems:

- **Space complexity.** To store full joint distribution requires to remember $O(d^n)$ numbers.
 - n – number of random variables, d – number of values
- **Inference complexity.** To compute some queries requires $O(d^n)$ steps.
- **Acquisition problem.** Who is going to define all of the probability entries?

CS 2001 Bayesian belief networks

Modeling uncertainty with probabilities

- **Knowledge based system era** (70s – early 80's)
 - Extensional non-probabilistic models
 - Probability techniques avoided because of space, time and acquisition bottlenecks in defining full joint distributions
 - Negative effect on the advancement of KB systems and AI in 80s in general
- Breakthrough (late 80s, beginning of 90s)
 - **Bayesian belief networks**
 - Give solutions to the space, acquisition bottlenecks
 - Significant improvements in the time cost of inferences

CS 2001 Bayesian belief networks

Bayesian belief networks (BBNs)

Bayesian belief networks.

- Represent the full joint distribution more compactly with smaller number of parameters.
- Take advantage of conditional and marginal independences among components in the distribution

- **A and B are independent**

$$P(A, B) = P(A)P(B)$$

- **A and B are conditionally independent given C**

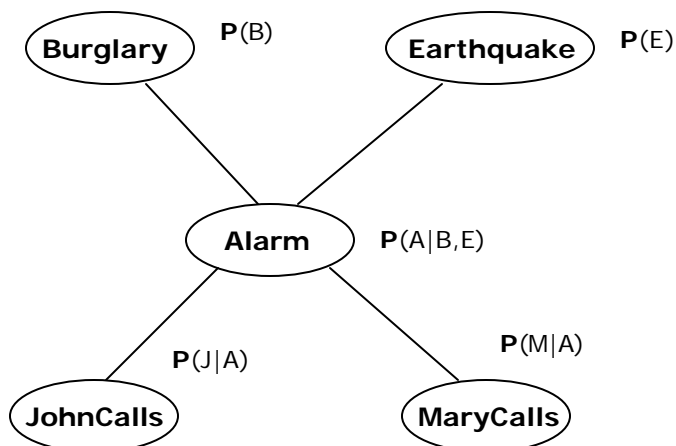
$$P(A, B | C) = P(A | C)P(B | C)$$

$$P(A | C, B) = P(A | C)$$

CS 2001 Bayesian belief networks

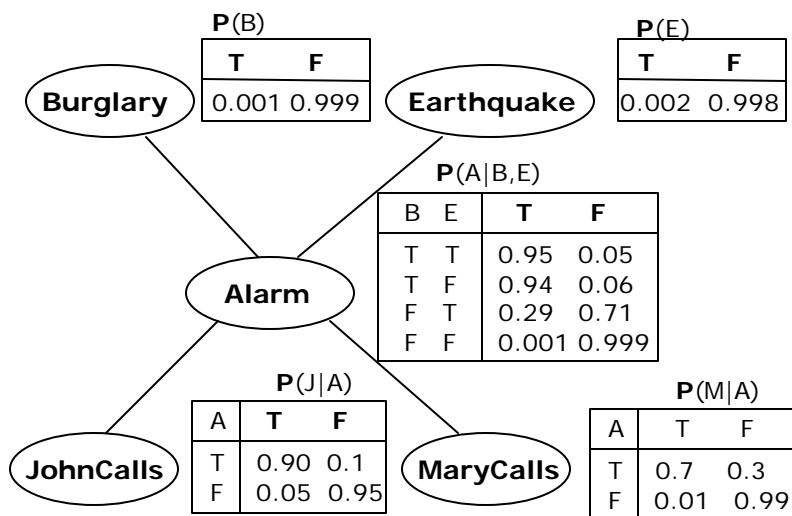
Bayesian belief network.

1. **Graph** represents marginal and conditional independences variables
2. **Parameters** defining local conditional distributions relating variables and their parents



CS 2001 Bayesian belief networks

Bayesian belief network.



CS 2001 Bayesian belief networks

Full joint distribution in BBNs

Full joint distribution is defined in terms of local conditional distributions (obtained via the chain rule):

$$\mathbf{P}(X_1, X_2, \dots, X_n) = \prod_{i=1, \dots, n} \mathbf{P}(X_i \mid pa(X_i))$$

Example:

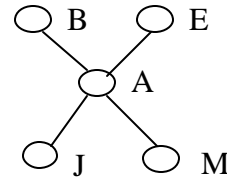
Assume the following assignment of values to random variables

$$B=T, E=T, A=T, J=T, M=F$$

Then its probability is:

$$P(B=T, E=T, A=T, J=T, M=F) =$$

$$P(B=T)P(E=T)P(A=T \mid B=T, E=T)P(J=T \mid A=T)P(M=F \mid A=T)$$



CS 2001 Bayesian belief networks

Parameter complexity problem

- In the BBN the full joint distribution is expressed as a product of conditionals (of smaller) complexity

$$\mathbf{P}(X_1, X_2, \dots, X_n) = \prod_{i=1, \dots, n} \mathbf{P}(X_i \mid pa(X_i))$$

Parameters:

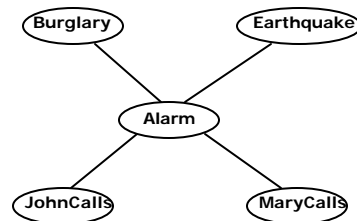
full joint: $2^5 = 32$

BBN: $2^3 + 2(2^2) + 2(2) = 20$

Parameters to be defined:

full joint: $2^5 - 1 = 31$

BBN: $2^2 + 2(2) + 2(1) = 10$



CS 2001 Bayesian belief networks

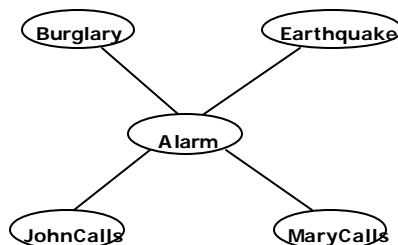
Inference in Bayesian networks

- BBN models compactly the full joint distribution by taking advantage of existing independences between variables
- Simplifies the acquisition of a probabilistic model
- But we are interested in solving various **inference tasks**:
 - **Diagnostic task. (from effect to cause)**
$$P(\text{Burglary} \mid \text{JohnCalls} = T)$$
 - **Prediction task. (from cause to effect)**
$$P(\text{JohnCalls} \mid \text{Burglary} = T)$$
 - **Other probabilistic queries** (queries on joint distributions).
$$P(\text{Alarm})$$
- **Question:** Can we take advantage of independences to construct special algorithms and speeding up the inference?

CS 2001 Bayesian belief networks

Inference in Bayesian network

- **Bad news:**
 - Exact inference problem in BBNs is NP-hard (G. Cooper)
 - Approximate inference is NP-hard (Dagum, Luby)
- **But** very often we can achieve significant improvements
- Assume our Alarm network



- Assume we want to compute: $P(J = T)$

CS 2001 Bayesian belief networks

Inference in Bayesian networks

Computing: $P(J = T)$

Approach 1. Blind approach.

- Sum out all uninstantiated variables from the full joint,
- Express the joint distribution as a product of conditionals

$$\begin{aligned}
 P(J = T) &= \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(B = b, E = e, A = a, J = T, M = m) \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e)
 \end{aligned}$$

Computational cost:

Number of additions: **15**

Number of products: $16 * 4 = \mathbf{64}$

CS 2001 Bayesian belief networks

Inference in Bayesian networks

Approach 2. Interleave sums and products

- Combines sums and product in a smart way (multiplications by constants can be taken out of the sum)

$$\begin{aligned}
 P(J = T) &= \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \\
 &= \sum_{b \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(B = b) \left[\sum_{e \in T, F} P(A = a | B = b, E = e) P(E = e) \right] \\
 &= \sum_{a \in T, F} P(J = T | A = a) \left[\sum_{m \in T, F} P(M = m | A = a) \right] \left[\sum_{b \in T, F} P(B = b) \right] \left[\sum_{e \in T, F} P(A = a | B = b, E = e) P(E = e) \right]
 \end{aligned}$$

Computational cost:

Number of additions: $1 + 2 * (1) + 2 * (1 + 2 * (1)) = \mathbf{9}$

Number of products: $2 * (2 + 2 * (0) + 2 * (1 + 2 * (1))) = \mathbf{16}$

CS 2001 Bayesian belief networks

Inference in Bayesian networks

- The smart interleaving of sums and products can help us to speed up the computation of joint probability queries
- What if we want to compute: $P(B = T, J = T)$

$$P(B = T, J = T) = \sum_{a \in T, F} P(J = T | A = a) \left[\sum_{m \in T, F} P(M = m | A = a) \right] [P(B = T) \left[\sum_{e \in T, F} P(A = a | B = T, E = e) P(E = e) \right]]$$

Are there any similarities from the previous query ?

$$P(J = T) = \quad ?$$

Inference in Bayesian networks

- The smart interleaving of sums and products can help us to speed up the computation of joint probability queries
- What if we want to compute: $P(B = T, J = T)$

$$P(B = T, J = T) = \sum_{a \in T, F} P(J = T | A = a) \left[\sum_{m \in T, F} P(M = m | A = a) \right] [P(B = T) \left[\sum_{e \in T, F} P(A = a | B = T, E = e) P(E = e) \right]]$$

$$P(J = T) = \begin{matrix} \updownarrow & \updownarrow & \updownarrow & \updownarrow & \updownarrow \\ \sum_{a \in T, F} P(J = T | A = a) \left[\sum_{m \in T, F} P(M = m | A = a) \right] \left[\sum_{b \in T, F} P(B = b) \left[\sum_{e \in T, F} P(A = a | B = b, E = e) P(E = e) \right] \right] \end{matrix}$$

- A lot of shared computation
 - Smart caching of results can save the time for more queries

Inference in Bayesian networks

- The smart interleaving of sums and products can help us to speed up the computation of joint probability queries
- What if we want to compute: $P(B = T, J = T)$

$$\begin{aligned}
 P(B = T, J = T) &= \\
 &= \sum_{a \in T, F} P(J = T | A = a) \left[\sum_{m \in T, F} P(M = m | A = a) \right] \left[P(B = T) \left[\sum_{e \in T, F} P(A = a | B = T, E = e) P(E = e) \right] \right] \\
 \\
 P(J = T) &= \\
 &= \sum_{a \in T, F} P(J = T | A = a) \left[\sum_{m \in T, F} P(M = m | A = a) \right] \left[\sum_{b \in T, F} P(B = b) \left[\sum_{e \in T, F} P(A = a | B = b, E = e) P(E = e) \right] \right]
 \end{aligned}$$

- A lot of shared computation
 - Smart caching of results can save the time if more queries

Inference in Bayesian networks

- When caching of results becomes handy?
- What if we want to compute a diagnostic query:

$$P(B = T | J = T) = \frac{P(B = T, J = T)}{P(J = T)}$$

- Exactly probabilities we have just compared !!
- There are other queries when caching and ordering of sums and products can be shared and saves computation

$$\mathbf{P}(B | J = T) = \frac{\mathbf{P}(B, J = T)}{P(J = T)} = \mathbf{aP}(B, J = T)$$

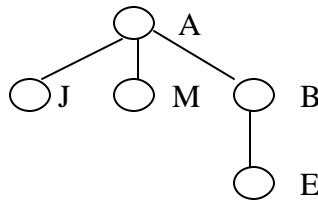
- General technique: Variable elimination

Inference in Bayesian networks

- General idea of variable elimination

$$\begin{aligned}
 P(\text{True}) &= 1 = \\
 &= \sum_{a \in T, F} \left[\underbrace{\sum_{j \in T, F} P(J=j | A=a)}_{f_J(a)} \right] \left[\underbrace{\sum_{m \in T, F} P(M=m | A=a)}_{f_M(a)} \right] \left[\underbrace{\sum_{b \in T, F} P(B=b) \left[\sum_{e \in T, F} P(A=a | B=b, E=e) P(E=e) \right]}_{f_E(a, b)} \right] \\
 &\hspace{15em} \underbrace{\hspace{15em}}_{f_B(a)}
 \end{aligned}$$

Variable order:



Results cached in the tree structure

CS 2001 Bayesian belief networks

Inference in Bayesian network

- **Exact inference algorithms:**
 - Symbolic inference (D'Ambrosio)
 - Recursive decomposition/variable elimination (Cooper, Dechter)
 - Message-passing algorithm (Pearl)
 - Clustering and joint-tree approach (Lauritzen, Spiegelhalter)
 - Arc reversal (Olmsted, Schachter)
- **Approximate inference algorithms:**
 - **Monte Carlo methods :**
 - Rejection sampling, Likelihood sampling
 - Variational methods

CS 2001 Bayesian belief networks

Learning Bayesian belief networks

- **Why learning?**
 - “subjective” estimates of conditional probability parameters by a human
 - need to adapt parameters in the light of observed data
 - large databases available
 - uncover important probabilistic dependencies from data and use them in inference tasks

CS 2001 Bayesian belief networks

Learning of BBN

Two learning tasks:

- Learning of the network structure
- Learning of parameters of conditional probabilities
- **Variables:**
 - **Observable** – values present in every data sample
 - **Hidden** – values are never in the sample
 - **Missing values** – values sometimes present, sometimes not
- **Here:**
 - Learning parameters of the fixed BBN structure
 - All variables are observable

CS 2001 Bayesian belief networks

Learning of BBNs.

Data: $D = \{D_1, D_2, \dots, D_n\}$
 $D_i = \mathbf{x}_i$ a vector of variable values

Random variables $\mathbf{X} = \{X_1, X_2, \dots, X_d\}$ with:

- **Continuous values**
- **Discrete values** ←

E.g. *blood pressure* with numerical values

or *chest pain* with discrete values

[no-pain, mild, moderate, strong]

Underlying true probability distribution:

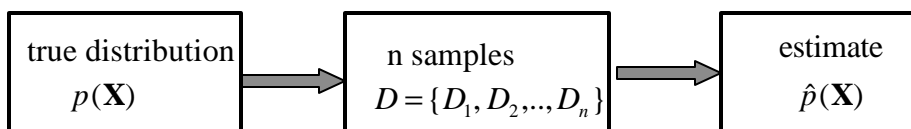
$$p(\mathbf{X})$$

CS 2001 Bayesian belief networks

Learning of BBNs.

Data: $D = \{D_1, D_2, \dots, D_n\}$
 $D_i = \mathbf{x}_i$ a vector of attribute values

Objective: try to estimate the underlying true probability distribution over variables \mathbf{X} , $p(\mathbf{X})$, using examples in D



Standard (iid) assumptions: Samples

- are independent of each other
- come from the same (identical) distribution (fixed $p(\mathbf{X})$)

CS 2001 Bayesian belief networks

Learning via parameter estimation

- For the fixed BBN we have:
 - **a model of the distribution** over variables in \mathbf{X}
with parameters \mathbf{T}
 - Parameters = parameters of local conditional probabilities

- **Objective:** find parameters \mathbf{T} that fit the data the best

$$\tilde{p}(\mathbf{X} | \mathbf{T}) \approx p(\mathbf{X})$$

- There are various criteria for defining the best set of parameters

CS 2001 Bayesian belief networks

Parameter estimation. Criteria.

- **Maximum likelihood (ML)**

$$\text{maximize } P(D | \Theta, \mathbf{x})$$

\mathbf{x} - represents prior (background) knowledge
(structure of the BBN)

- **Maximum a posteriori probability (MAP)**

$$\text{maximize } P(\Theta | D, \mathbf{x})$$

$$P(\Theta | D, \mathbf{x}) = \frac{P(D | \Theta, \mathbf{x})P(\Theta | \mathbf{x})}{P(D | \mathbf{x})}$$

CS 2001 Bayesian belief networks

Parameter estimation. Example.

Coin example: we have a coin that can be biased

Outcomes: two possible values -- head or tail

Data: D a sequence of outcomes x_i such that

- **head** $x_i = 1$
- **tail** $x_i = 0$

Model: probability of a head q
 probability of a tail $(1-q)$

Objective:

We would like to estimate the probability of a head \hat{q}

Probability of an outcome x_i

$$P(x_i | q) = q^{x_i} (1-q)^{(1-x_i)} \quad \text{Bernoulli distribution}$$

Maximum likelihood (ML) estimate.

Likelihood of data:

$$P(D | q, \mathbf{x}) = \prod_{i=1}^n q^{x_i} (1-q)^{(1-x_i)}$$

Maximum likelihood estimate

$$q_{ML} = \arg \max_q P(D | q, \mathbf{x})$$

Optimize log-likelihood

$$l(D, q) = \log P(D | q, \mathbf{x}) = \log \prod_{i=1}^n q^{x_i} (1-q)^{(1-x_i)} =$$

$$\sum_{i=1}^n x_i \log q + (1-x_i) \log(1-q) = \log q \underbrace{\sum_{i=1}^n x_i}_{N_1} + \log(1-q) \underbrace{\sum_{i=1}^n (1-x_i)}_{N_2}$$

N_1 - number of heads seen N_2 - number of tails seen

Maximum likelihood (ML) estimate.

Optimize log-likelihood

$$l(D, \mathbf{q}) = N_1 \log \mathbf{q} + N_2 \log(1 - \mathbf{q})$$

Set derivative to zero

$$\frac{\partial l(D, \mathbf{q})}{\partial \mathbf{q}} = \frac{N_1}{\mathbf{q}} - \frac{N_2}{(1 - \mathbf{q})} = 0$$

Solving $\mathbf{q} = \frac{N_1}{N_1 + N_2}$

ML Solution: $\mathbf{q}_{ML} = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2}$

CS 2001 Bayesian belief networks

Maximum a posteriori estimate

Maximum a posteriori estimate

- Selects the mode of the posterior distribution

$$\mathbf{q}_{MAP} = \arg \max_{\mathbf{q}} p(\mathbf{q} | D, \mathbf{x})$$

$$p(\mathbf{q} | D, \mathbf{x}) = \frac{P(D | \mathbf{q}, \mathbf{x}) p(\mathbf{q} | \mathbf{x})}{P(D | \mathbf{x})} \quad (\text{via Bayes rule})$$

$P(D | \mathbf{q}, \mathbf{x})$ - is the likelihood of data

$$P(D | \mathbf{q}, \mathbf{x}) = \prod_{i=1}^n \mathbf{q}^{x_i} (1 - \mathbf{q})^{(1-x_i)} = \mathbf{q}^{N_1} (1 - \mathbf{q})^{N_2}$$

$p(\mathbf{q} | \mathbf{x})$ - is the prior probability on \mathbf{q}

How to choose the prior probability?

CS 2001 Bayesian belief networks

Prior distribution

Choice of prior: Beta distribution

$$p(\mathbf{q} | \mathbf{x}) = \text{Beta}(\mathbf{q} | \mathbf{a}_1, \mathbf{a}_2) = \frac{\Gamma(\mathbf{a}_1 + \mathbf{a}_2)}{\Gamma(\mathbf{a}_1)\Gamma(\mathbf{a}_2)} \mathbf{q}^{\mathbf{a}_1-1} (1-\mathbf{q})^{\mathbf{a}_2-1}$$

Why?

Beta distribution “fits” binomial sampling - **conjugate choices**

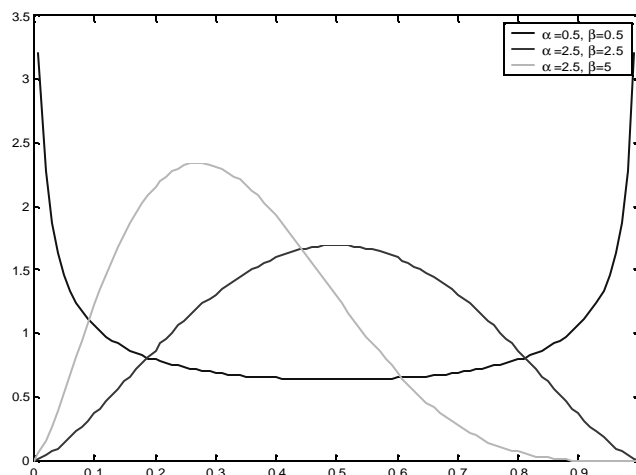
$$P(D | \mathbf{q}, \mathbf{x}) = \mathbf{q}^{N_1} (1-\mathbf{q})^{N_2}$$

$$p(\mathbf{q} | D, \mathbf{x}) = \frac{P(D | \mathbf{q}, \mathbf{x}) \text{Beta}(\mathbf{q} | \mathbf{a}_1, \mathbf{a}_2)}{P(D | \mathbf{x})} = \text{Beta}(\mathbf{q} | \mathbf{a}_1 + N_1, \mathbf{a}_2 + N_2)$$

MAP Solution:
$$\mathbf{q}_{MAP} = \frac{\mathbf{a}_1 + N_1 - 1}{\mathbf{a}_1 + \mathbf{a}_2 + N_1 + N_2 - 2}$$

CS 2001 Bayesian belief networks

Beta distribution



CS 2001 Bayesian belief networks

Bayesian learning

- **Both ML or MAP pick one parameter value**
 - Is it always the best solution?
 - **Assume:** there are two very different parameter settings that are close in terms of probability (ML or MAP). Using only one of them may introduce a strong bias, if we use them e.g. for predictions.
- **Bayesian approach**
 - Remedies the limitation of one choice
 - Considers all parameter settings and averages the result

$$P(\Delta | D, \mathbf{x}) = \int P(\Delta | \mathbf{q}, \mathbf{x}) p(\mathbf{q} | D, \mathbf{x}) d\mathbf{q}$$

- **Example:** predict the result of the next coin flip
 - Choose outcome 1 if $P(x=1 | D, \mathbf{x})$ is higher

CS 2001 Bayesian belief networks

Bayesian learning

- **Predictive probability of an outcome $x=1$ in the next trial**

$$P(x=1 | D, \mathbf{x})$$

$$\begin{aligned} P(x=1 | D, \mathbf{x}) &= \int_0^1 P(x=1 | \mathbf{q}, \mathbf{x}) p(\mathbf{q} | D, \mathbf{x}) d\mathbf{q} \\ &= \int_0^1 \mathbf{q} p(\mathbf{q} | D, \mathbf{x}) d\mathbf{q} = E(\mathbf{q}) \end{aligned}$$

- **Equivalent to the expected value of the parameter**
 - expectation is taken with regard to the posterior distribution

$$p(\mathbf{q} | D, \mathbf{x}) = \text{Beta}(\mathbf{q} | \mathbf{a}_1 + N_1, \mathbf{a}_2 + N_2)$$

CS 2001 Bayesian belief networks

Bayesian learning, expectation

The result is:

$$E(\mathbf{q}) = \int_0^1 \mathbf{q} \text{Beta}(\mathbf{q} | \mathbf{h}_1, \mathbf{h}_2) d\mathbf{q} = \frac{\mathbf{h}_1}{\mathbf{h}_1 + \mathbf{h}_2}$$

Bayesian estimate of the parameter

$$E(\mathbf{q}) = \frac{\mathbf{a}_1 + N_1}{\mathbf{a}_1 + N_1 + \mathbf{a}_2 + N_2}$$

Predictive probability of event $x=1$

$$E(\mathbf{q}) = P(x=1 | \mathbf{q}, \mathbf{x}) = \frac{\mathbf{a}_1 + N_1}{\mathbf{a}_1 + N_1 + \mathbf{a}_2 + N_2}$$

CS 2001 Bayesian belief networks

Example. Multi-way coin toss

Multi-way coin toss (roll of dice)

- **Data:** a set of N outcomes (multi-set)

N_i - a number of times an outcome i has been seen

Model parameters: $\mathbf{?} = (\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_k)$ s.t. $\sum_{i=1}^k \mathbf{q}_i = 1$
 \mathbf{q}_i - probability of an outcome i

Probability of data (likelihood)

$$P(N_1, N_2, \dots, N_k | \mathbf{?}, \mathbf{x}) = \frac{N!}{N_1! N_2! \dots N_k!} \mathbf{q}_1^{N_1} \mathbf{q}_2^{N_2} \dots \mathbf{q}_k^{N_k} \quad \text{Multinomial distribution}$$

ML estimate:

$$\mathbf{q}_{i,ML} = \frac{N_i}{N}$$

CS 2001 Bayesian belief networks

MAP estimate

Choice of prior: Dirichlet distribution

$$Dir(? | \mathbf{a}_1, \dots, \mathbf{a}_k) = \frac{\Gamma(\sum_{i=1}^k \mathbf{a}_i)}{\prod_{i=1}^k \Gamma(\mathbf{a}_i)} \mathbf{q}_1^{\mathbf{a}_1-1} \mathbf{q}_2^{\mathbf{a}_2-1} \dots \mathbf{q}_k^{\mathbf{a}_k-1}$$

Dirichlet is the conjugate choice for multinomial

$$P(D | ?, \mathbf{x}) = P(N_1, N_2, \dots, N_k | ?, \mathbf{x}) = \frac{N!}{N_1! N_2! \dots N_k!} \mathbf{q}_1^{N_1} \mathbf{q}_2^{N_2} \dots \mathbf{q}_k^{N_k}$$

Posterior distribution

$$p(? | D, \mathbf{x}) = \frac{P(D | ?, \mathbf{x}) Dir(? | \mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k)}{P(D | \mathbf{x})} = Dir(? | \mathbf{a}_1 + N_1, \dots, \mathbf{a}_k + N_k)$$

MAP estimate:
$$\mathbf{q}_{i,MAP} = \frac{\mathbf{a}_i + N_i - 1}{\sum_{i=1, \dots, k} (\mathbf{a}_i + N_i) - k}$$

CS 2001 Bayesian belief networks

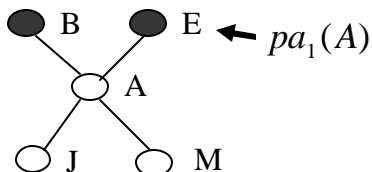
Learning of parameters of BBNs

- We need to estimate parameters of conditional distributions

$$\mathbf{P}(X_i | pa(X_i))$$

- **Idea:**

- Fix an assignment of the values to parent variables of X_i



- Pull out of the dataset D only data that are in agreement with this assignment -- $pa_1(A)$
- Learn the parameters of the conditional $\mathbf{P}(A | pa_1(A))$
- If discrete-valued variable A this is like a multiway coin toss

CS 2001 Bayesian belief networks

Estimates of parameters of BBN

- **Two assumptions to make this possible:**
 - **Sample independence**

$$P(D | \mathbf{T}, \mathbf{x}) = \prod_{u=1}^n p(D_u | \mathbf{T}, \mathbf{x})$$

- **Parameter independence**

$$P(\mathbf{T} | D, \mathbf{x}) = \prod_{i=1}^d \prod_{j=1}^{q_i} p(q_{ij} | D, \mathbf{x})$$

Parameters of each node-parents conditional can be optimized independently

CS 2001 Bayesian belief networks

ML Course

CS2750 Machine Learning, Spring 2003

Instructor: **Milos Hauskrecht**

web page: <http://www.cs.pitt.edu/~milos/courses/cs2750/>

- Covers modern machine learning techniques, including learning of BBNs, their structures and parameters in different settings, as well as, many other learning frameworks, such as neural networks, support vector machines etc.

CS 2001 Bayesian belief networks