

CS 3750 Machine Learning

Lecture 2

Advanced Machine Learning

Milos Hauskrecht

milos@cs.pitt.edu

5329 Sennott Square, x4-8845

<http://www.cs.pitt.edu/~milos/courses/cs3750/>

CS 3750 Advanced Machine Learning

Learning

Starts with data & prior knowledge

Typical steps in learning:

- Define a model space
- Define an objective criterion: criterion for measuring the goodness of a model (fit to data)
- Optimization: finding the best model

Alternative: optimization is replaced with the inference, e.g. Bayesian inference in the Bayesian learning

Evaluation/application:

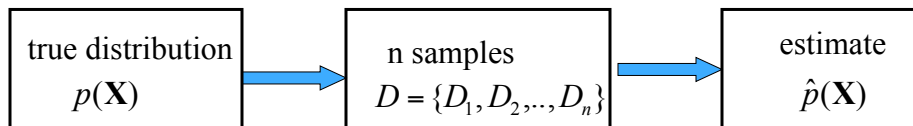
- Model learned from the training data
- generalization to the future (test) data

CS 3750 Advanced Machine Learning

Density estimation

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 3750 Advanced Machine Learning

Density estimation

Types of density estimation:

Parametric

- the distribution is modeled using a set of parameters Θ
 $p(\mathbf{X} | \Theta)$
- **Example:** mean and covariances of multivariate normal
- **Estimation:** find parameters $\hat{\Theta}$ that fit the data D the best

Non-parametric

- The model of the distribution utilizes all examples in D
- As if all examples were parameters of the distribution
- The density for a point \mathbf{x} is influenced by examples in its neighborhood

CS 3750 Advanced Machine Learning

Basic criteria

What is the best set of parameters?

- **Maximum likelihood (ML)**

$$\text{maximize } p(D | \Theta, \xi)$$

ξ - represents prior (background) knowledge

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

$$\text{maximize } p(\Theta | D, \xi)$$

Selects the mode of the posterior

$$p(\Theta | D, \xi) = \frac{p(D | \Theta, \xi) p(\Theta | \xi)}{p(D | \xi)}$$

CS 3750 Advanced Machine Learning

Example. Bernoulli distribution.

Outcomes: two possible values – 0 or 1 (head or tail)

Data: D a sequence of outcomes x_i with 0,1 values

Model: probability of an outcome 1 θ
probability of 0 $(1 - \theta)$

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

Objective:

We would like to estimate the probability of seeing 1:

$$\hat{\theta}$$

CS 3750 Advanced Machine Learning

Maximum likelihood (ML) estimate.

Likelihood of data:

$$P(D | \theta, \xi) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{(1-x_i)}$$

Maximum likelihood estimate

$$\theta_{ML} = \arg \max_{\theta} P(D | \theta, \xi)$$

Optimize log-likelihood

$$\begin{aligned} l(D, \theta) &= \log P(D | \theta, \xi) = \log \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{(1-x_i)} = \\ &= \sum_{i=1}^n x_i \log \theta + (1 - x_i) \log(1 - \theta) = \log \theta \underbrace{\sum_{i=1}^n x_i}_{N_1} + \log(1 - \theta) \underbrace{\sum_{i=1}^n (1 - x_i)}_{N_2} \end{aligned}$$

N_1 - number of 1s seen N_2 - number of 0s seen

CS 3750 Advanced Machine Learning

Maximum likelihood (ML) estimate.

Optimize log-likelihood

$$l(D, \theta) = N_1 \log \theta + N_2 \log(1 - \theta)$$

Set derivative to zero

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

Solving

$$\theta = \frac{N_1}{N_1 + N_2}$$

ML Solution: $\theta_{ML} = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2}$

CS 3750 Advanced Machine Learning

Maximum a posteriori estimate

Maximum a posteriori estimate

- Selects the mode of the posterior distribution

$$\theta_{MAP} = \arg \max_{\theta} p(\theta | D, \xi)$$

$$p(\theta | D, \xi) = \frac{P(D | \theta, \xi) p(\theta | \xi)}{P(D | \xi)} \quad (\text{via Bayes rule})$$

$P(D | \theta, \xi)$ - is the likelihood of data

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

$p(\theta | \xi)$ - is the prior probability on θ

How to choose the prior probability?

CS 3750 Advanced Machine Learning

Prior distribution

Choice of prior: Beta distribution

$$p(\theta | \xi) = \text{Beta}(\theta | \alpha_1, \alpha_2) = \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} \theta^{\alpha_1-1} (1 - \theta)^{\alpha_2-1}$$

Why?

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

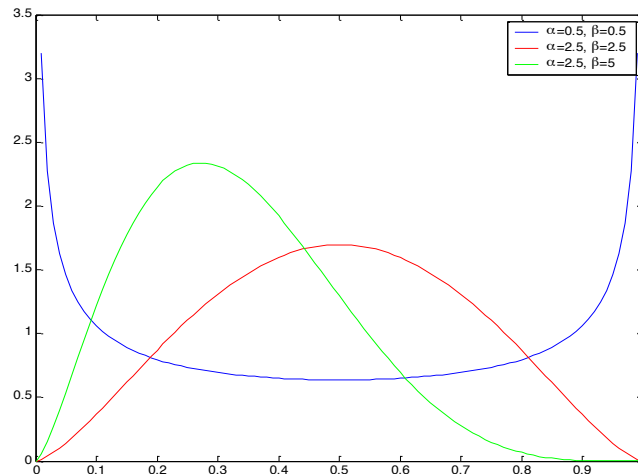
$$P(D | \theta, \xi) = \theta^{N_1} (1 - \theta)^{N_2}$$

$$p(\theta | D, \xi) = \frac{P(D | \theta, \xi) \text{Beta}(\theta | \alpha_1, \alpha_2)}{P(D | \xi)} = \text{Beta}(\theta | \alpha_1 + N_1, \alpha_2 + N_2)$$

MAP Solution:
$$\theta_{MAP} = \frac{\alpha_1 + N_1 - 1}{\alpha_1 + \alpha_2 + N_1 + N_2 - 2}$$

CS 3750 Advanced Machine Learning

Beta distribution



CS 3750 Advanced Machine Learning

Bayesian learning

- **Both ML or MAP pick one parameter value**
 - Is it always the best solution?
- **Full Bayesian approach**
 - Remedies the limitation of one choice
 - Keeps and uses a complete posterior distribution
- **How is it used? Assume we want: $P(\Delta | D, \xi)$**
 - Considers all parameter settings and averages the result
$$P(\Delta | D, \xi) = \int_{\theta} P(\Delta | \theta, \xi) p(\theta | D, \xi) d\theta$$
 - **Example:** predict the result of the next outcome
 - Choose outcome 1 if $P(x=1 | D, \xi)$ is higher

CS 3750 Advanced Machine Learning

Other distributions

The same ideas can be applied to other distributions

- Typically we choose distributions that behave well so that computations lead to a nice solutions

- **Exponential family of distributions**

Conjugate choices (sample – prior combinations) for some of the distributions from the exponential family:

- **Binomial – Beta**
- **Multinomial - Dirichlet**
- **Exponential – Gamma**
- **Poisson – Inverse Gamma**
- **Gaussian - Gaussian (mean) and Wishart (covariance)**

CS 2750 Machine Learning

Non-parametric density estimation

Parametric density estimation:

- A set of random variables $\mathbf{X} = \{X_1, X_2, \dots, X_d\}$
- **A model of the distribution** over variables in \mathbf{X}
with **parameters** $\Theta : \hat{p}(\mathbf{X} | \Theta)$
- **Data** $D = \{D_1, D_2, \dots, D_n\}$

Objective: find parameters Θ such that $p(\mathbf{X} | \Theta)$ models D

Parametric models are:

- restricted to specific forms, which may not always be suitable;
- **Nonparametric approaches:**
 - make few assumptions about the overall shape of the distribution being modelled.

CS 2750 Machine Learning

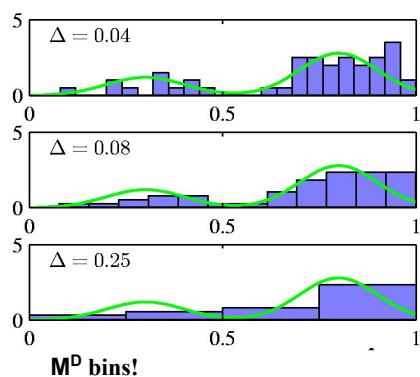
Nonparametric Methods

Histogram methods:

partition the data space into distinct bins with widths Δ_i and count the number of observations, n_i , in each bin.

$$p_i = \frac{n_i}{N\Delta_i}$$

- Often, the same width is used for all bins, $\Delta_i = \Delta$.
- Δ acts as a smoothing parameter.



CS 2750 Machine Learning

Nonparametric Methods

- Assume observations drawn from a density $p(x)$ and consider a small region R containing x such that

$$P = \int_R p(x) dx$$

- The probability that K out of N observations lie inside R is $\text{Bin}(K, N, P)$ and if N is large

$$K \approx NP$$

If the volume of R , V , is sufficiently small, $p(x)$ is approximately constant over R and

$$P \approx p(x)V$$

Thus

$$p(x) = \frac{P}{V}$$

$$p(x) = \frac{K}{NV}$$

CS 2750 Machine Learning

Nonparametric Methods: kernel methods

Kernel Density Estimation:

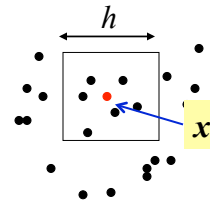
Fix \mathbf{V} , estimate \mathbf{K} from the data. Let \mathbf{R} be a hypercube centred on \mathbf{x} and define the kernel function (Parzen window)

$$k\left(\frac{x - x_n}{h}\right) = \begin{cases} 1 & |(x_i - x_{ni})| / h \leq 1/2 \\ 0 & \text{otherwise} \end{cases} \quad i = 1, \dots, D$$

- It follows that

- and hence
$$K = \sum_{n=1}^N k\left(\frac{x - x_n}{h}\right)$$

$$p(x) = \frac{1}{N} \sum_{n=1}^N \frac{1}{h^D} k\left(\frac{x - x_n}{h}\right)$$



CS 2750 Machine Learning

Nonparametric Methods: smooth kernels

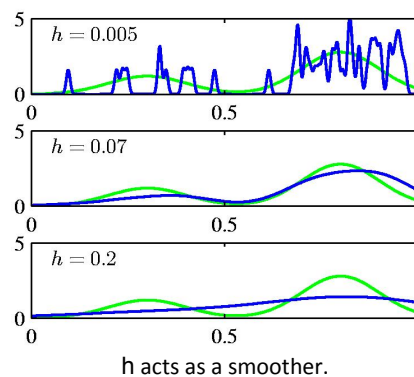
To avoid discontinuities in $p(x)$
because of sharp boundaries
use a **smooth kernel**, e.g. a
Gaussian

$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{(2\pi h^2)^{D/2}} \exp\left\{-\frac{\|\mathbf{x} - \mathbf{x}_n\|^2}{2h^2}\right\}$$

- Any kernel such that

$$\begin{aligned} k(u) &\geq 0, \\ \int k(u) du &= 1 \end{aligned}$$

- will work.



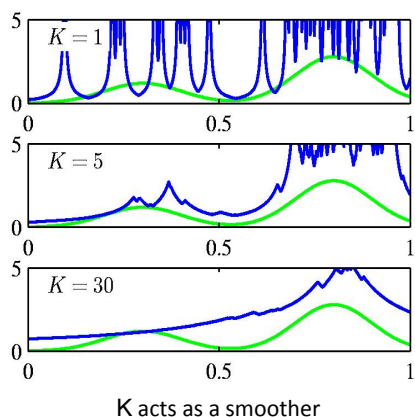
CS 2750 Machine Learning

Nonparametric Methods: kNN estimation

Nearest Neighbour Density Estimation:

fix K , estimate V from the data. Consider a hyper-sphere centred on \mathbf{x} and let it grow to a volume, V^* , that includes K of the given N data points. Then

$$p(\mathbf{x}) \simeq \frac{K}{NV^*}.$$



CS 2750 Machine Learning

Modeling complex multivariate distributions

How to model complex multivariate parametric distributions $\hat{p}(\mathbf{X})$ with large number of variables?

One solution:

- **Decompose the distribution. Reduce the number of parameters, using some form of independence.**

Two models:

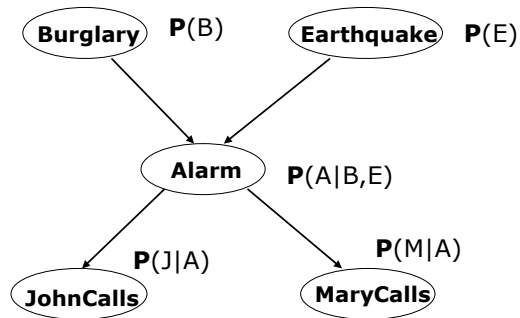
- **Bayesian belief networks (BBNs)**
- **Markov Random Fields (MRFs)**
- **Learning.** Relies on the decomposition.

CS 3750 Advanced Machine Learning

Bayesian belief network.

1. Directed acyclic graph

- **Nodes** = random variables
- **Links** = direct (causal) dependencies between variables
 - Missing links encode independences

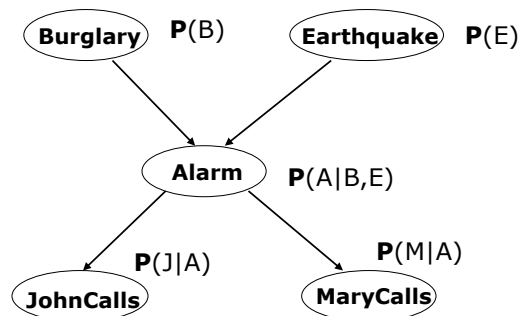


CS 3750 Advanced Machine Learning

Bayesian belief network.

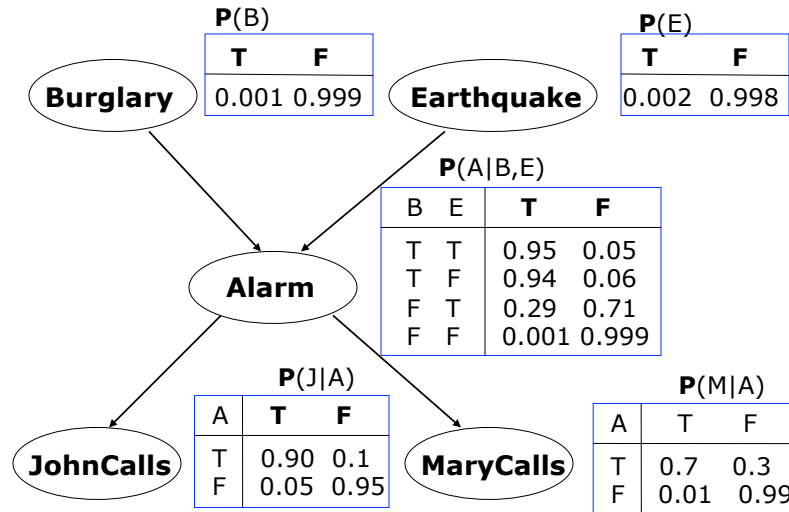
2. Local conditional distributions

- relate variables and their parents



CS 3750 Advanced Machine Learning

Bayesian belief network.



CS 3750 Advanced Machine Learning

Full joint distribution in BBNs

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

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

Example:

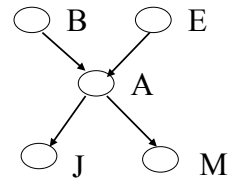
Assume the following assignment of values to random variables

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

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 3750 Advanced Machine Learning

Learning of BBN

Learning.

- **Learning of parameters of conditional probabilities**
- **Learning of the network structure**

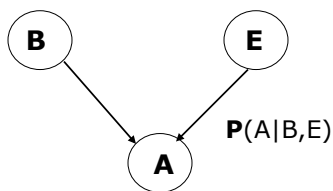
Variables:

- **Observable** – values present in every data sample
- **Hidden** – they values are never observed in data
- **Missing values** – values sometimes present, sometimes not

CS 2750 Machine Learning

Estimation of parameters of BBN

- **Idea:** decompose the estimation problem for the full joint over a large number of variables to a set of smaller estimation problems corresponding to local parent-variable conditionals.
- **Example:** Assume A,E,B are binary with *True*, *False* values



4 estimation problems

- $$\left\{ \begin{array}{l} P(A|B=T,E=T) \\ P(A|B=T,E=F) \\ P(A|B=F,E=T) \\ P(A|B=F,E=F) \end{array} \right.$$

- **Assumption that enables the decomposition:** parameters of conditional distributions are independent

CS 2750 Machine Learning

Estimates of parameters of BBN

- Two assumptions that permit the decomposition:
 - **Sample independence**

$$P(D | \Theta, \xi) = \prod_{u=1}^N P(D_u | \Theta, \xi)$$

- **Parameter independence**

$$p(\Theta | D, \xi) = \prod_{i=1}^n \prod_{j=1}^{q_i} p(\theta_{ij} | D, \xi)$$

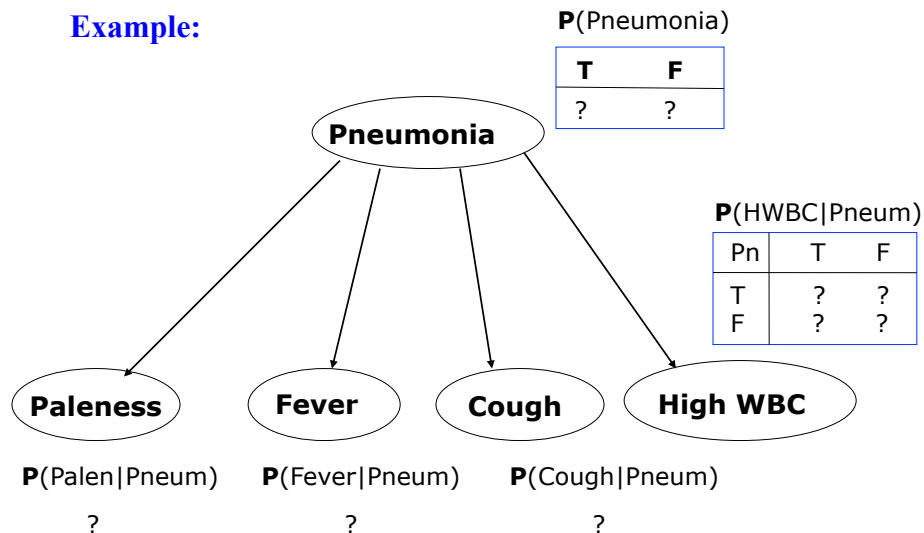
of nodes
 # of parents values

Parameters of **each conditional** (one for every assignment of values to parent variables) can be learned independently

CS 2750 Machine Learning

Learning of BBN parameters. Example.

Example:



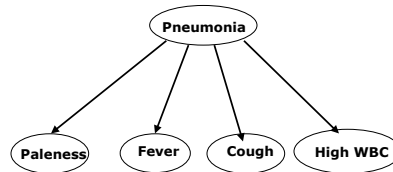
CS 2750 Machine Learning

Learning of BBN parameters. Example.

Data D (different patient cases):

Pal Fev Cou HWB Pneu

T	T	T	T	F
T	F	F	F	F
F	F	T	T	T
F	F	T	F	T
F	T	T	T	T
T	F	T	F	F
F	F	F	F	F
T	T	F	F	F
T	T	T	T	T
F	T	F	T	T
T	F	F	T	F
F	T	F	F	F



CS 2750 Machine Learning

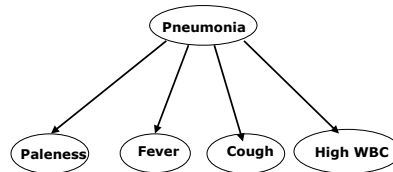
Learning of BBN parameters. Example.

Learn: $P(\text{Fever} \mid \text{Pneumonia} = T)$

Step 1: Select data points with Pneumonia=T

Pal Fev Cou HWB Pneu

T	T	T	T	F
T	F	F	F	F
F	F	T	T	T
F	F	T	F	T
F	T	T	T	T
T	F	T	F	F
F	F	F	F	F
T	T	F	F	F
T	T	T	T	T
F	T	F	T	T
T	F	F	T	F
F	T	F	F	F



CS 2750 Machine Learning

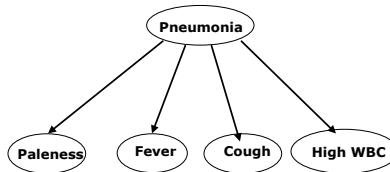
Learning of BBN parameters. Example.

Learn: $P(\text{Fever} \mid \text{Pneumonia} = T)$

Step 1: Ignore the rest

Pal Fev Cou HWB Pneu

F	F	T	T	T
F	F	T	F	T
F	T	T	T	T
T	T	T	T	T
F	T	F	T	T



CS 2750 Machine Learning

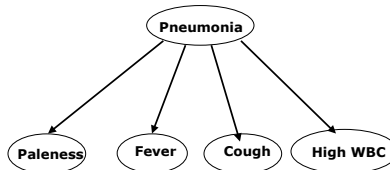
Learning of BBN parameters. Example.

Learn: $P(\text{Fever} \mid \text{Pneumonia} = T)$

Step 2: Select values of the random variable defining the distribution of Fever

Pal Fev Cou HWB Pneu

F	F	T	T	T
F	F	T	F	T
F	T	T	T	T
T	T	T	T	T
F	T	F	T	T



CS 2750 Machine Learning

Learning of BBN parameters. Example.

Learn: $P(\text{Fever} \mid \text{Pneumonia} = T)$

Step 2: Ignore the rest

Fev

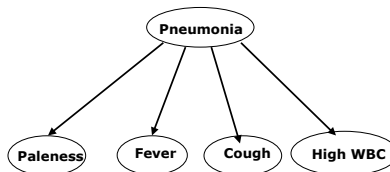
F

F

T

T

T



CS 2750 Machine Learning

Learning of BBN parameters. Example.

Learn: $P(\text{Fever} \mid \text{Pneumonia} = T)$

Step 3a: Learning the ML estimate

Fev

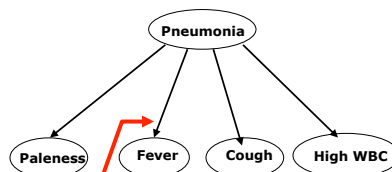
F

F

T

T

T



$P(\text{Fever} \mid \text{Pneumonia} = T)$

T	F
0.6	0.4

CS 2750 Machine Learning

Learning of BBN parameters. Bayesian learning.

Learn: $P(\text{Fever} \mid \text{Pneumonia} = T)$

Step 3b: Learning the Bayesian estimate

Assume the prior

$$\theta_{\text{Fever} \mid \text{Pneumonia} = T} \sim \text{Beta}(3, 4)$$

Fev

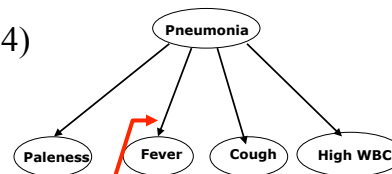
F

F

T

T

T



Posterior:

$$\theta_{\text{Fever} \mid \text{Pneumonia} = T} \sim \text{Beta}(6, 6)$$

CS 2750 Machine Learning

Hidden variables

Modeling assumption:

Variables $\mathbf{X} = \{X_1, X_2, \dots, X_n\}$

- Additional variables are hidden – never observed in data

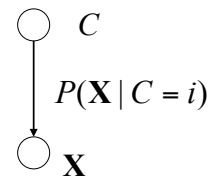
Why to add hidden variables?

- **More flexibility in describing the distribution** $P(\mathbf{X})$
- **Smaller parameterization of** $P(\mathbf{X})$
 - New independences can be introduced via hidden variables

Example:

- Latent variable models
 - hidden classes (categories)

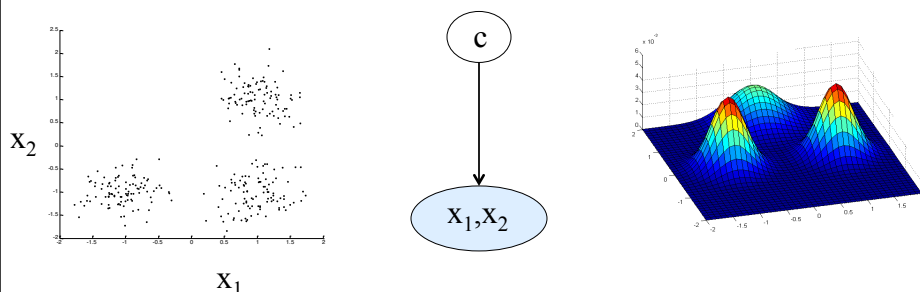
Hidden class variable



CS 2750 Machine Learning

Latent variable models

- We can have a model with hidden variables
- Hidden variables may help us to induce the decomposition of a complex distribution



CS 3750 Advanced Machine Learning

Learning with hidden variables and missing values

Goal: Find the set of parameters $\hat{\Theta}$

Estimation criteria:

- **ML** $\max_{\Theta} p(D | \Theta, \xi)$
- **Bayesian** $p(\Theta | D, \xi)$

Optimization methods for ML: gradient-ascent, conjugate gradient, Newton-Rhapson, etc.

Problem: No or very small advantage from the structure of the corresponding belief network when unobserved variable values

Expectation-maximization (EM) method

- An alternative optimization method
- Suitable when there are missing or hidden values
- **Takes advantage of the structure of the belief network**

CS 2750 Machine Learning

General EM

The key idea of a method:

Compute the parameter estimates iteratively by performing the following two steps:

Two steps of the EM:

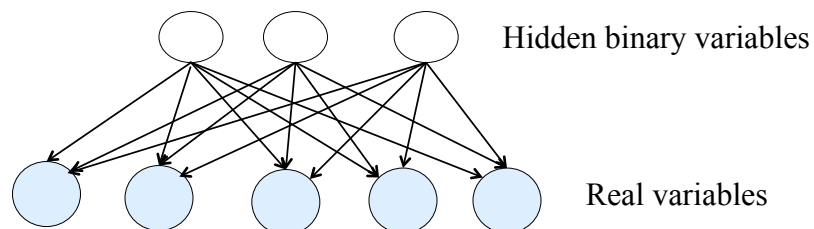
1. **Expectation step.** Complete all hidden and missing variables with expectations for the current set of parameters Θ'
2. **Maximization step.** Compute the new estimates of Θ for the completed data

Stop when no improvement possible

CS 2750 Machine Learning

Latent variable models

- More general latent variable models
- Various relations in between hidden and observable variables
- **Example:** Continuous vector quantizer (CVQ) model



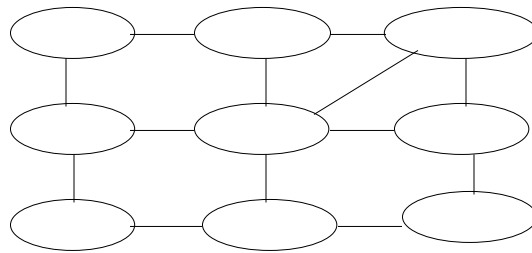
- **Possible uses:**
- A probabilistic model
- A low dimensional representation of observable data

CS 3750 Advanced Machine Learning

Markov Random Fields (MRFs)

Undirected graph

- **Nodes** = random variables
- **Links** = direct relations between variables
- BBNs used to model **asymmetric** dependencies (most often causal),
- MRFs model **symmetric** dependencies (bidirectional effects) such as spatial dependences

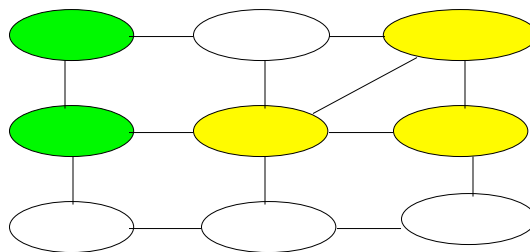


CS 3750 Advanced Machine Learning

Markov Random Fields (MRFs)

A probability distribution is defined in terms of potential functions defined over cliques of the graph

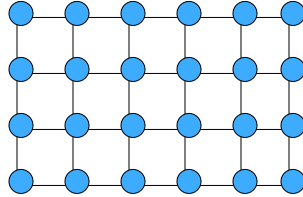
$$P(X_1, X_2, \dots, X_n) = \frac{1}{Z} \prod_{C_i \in \text{cliques}(G)} \Psi(C_i)$$



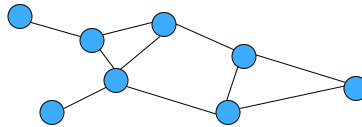
CS 3750 Advanced Machine Learning

Markov random fields

- regular lattice (Ising model)



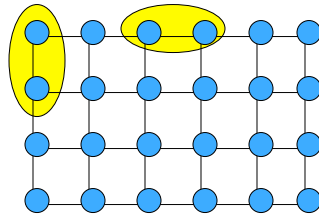
- Arbitrary graph



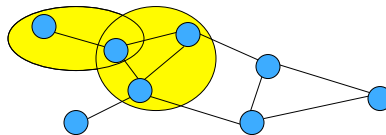
CS 3750 Advanced Machine Learning

Markov random fields

- regular lattice (Ising model)



- Arbitrary graph



CS 3750 Advanced Machine Learning