

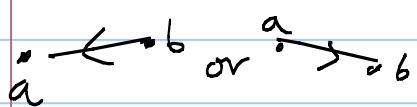
Primer: Structured Probability Distributions

Note Title

4/2/2006

Fundamentals of Bayesian Inference.

Two variables. joint probability $P(a,b)$
conditional prob $P(a|b)$ & $P(b|a)$
marginal probs $P(a)$ & $P(b)$



$$P(a,b) = P(a|b)P(b) = P(b|a)P(a)$$

Implies Bayes Rule

$$P(b|a) = \frac{P(a|b)P(b)}{P(a)}$$

Replace a by d (data)
 b by $h \in \mathcal{H}$ set of hypotheses.
Hypothesis Space.

$$P(h|d) = \frac{P(d|h)P(h)}{\sum_{h' \in \mathcal{H}} P(d|h')P(h')}$$

Compare Hypotheses.

E.G. Are data tosses from fair coin ($P_r(H) = 0.5$)
or a biased coin ($P_r(H) = 0.9$)

If $HH\dots H$ - probably think biased.
 $HHHTHTHTH$ - probably think fair.

(Page 2)

To address this problem formally
Let θ be the probability that the coin produces head.

Hypothesis

$h_0: \theta = 0.5$ unbiased

$h_1: \theta = 0.9$ biased.

Hypothesis Space $\mathcal{H} = \{h_0, h_1\}$

$$P(d|\theta) = \theta^{N_H} (1-\theta)^{N_T}$$

N_H no. of heads in data d

N_T no. of tails " " "

Posterior odds of the hypotheses

$$\frac{P(h_1|d)}{P(h_0|d)} = \frac{P(d|h_1) P(h_1)}{P(d|h_0) P(h_0)}$$

Gives 357:1 in favor of h_0 from HHH...H
165:1 in favor of h_1 from HHTHTHTHTT...

Combining Infinitely Many Hypotheses

Suppose $0 \leq \theta \leq 1$.

$$p(\theta|d) = P(d|\theta) p(\theta)$$

\leftarrow prior

$$p(d) = \int_0^1 \overbrace{P(d|\theta)}^{p(d)} p(\theta) d\theta.$$

(3)

Example: If $P(\theta) = 1$.

$$P(\theta|d) = \frac{\text{Beta}(N_H+1, N_T+1)}{N_H! N_T!} \theta^{N_H} (1-\theta)^{N_T}$$

$$\text{If } P(\theta) = \text{Beta}(V_H+1, V_T+1)$$

V_H & V_T
positive integer.

$$P(\theta|d) = \frac{\text{Beta}(N_H+V_H+1, N_T+V_T+1)}{(N_H+V_H)! (N_T+V_T)!} \theta^{N_H+V_H} (1-\theta)^{N_T+V_T}$$

Comparing hypotheses of different complexity.
~ Model Selection.

h_0 : hypothesis $\theta = 0.5$

h_1 : hypothesis θ drawn from uniform distribution over θ .

$$P(d|h_0) = (0.5)^{N_H+N_T}$$

$$P(d|h_1) = \int_0^1 P(d|\theta, h_1) P(\theta|h_1) d\theta$$

$$= \frac{N_H! N_T!}{(N_H+N_T+1)!}$$

$$P(\theta|h_1) = 1.$$

Can apply Bayes rule as before.:

Important: Occam's razor

complex hypotheses have more degrees of freedom & can be adapted to data. But integrating over θ prevents this

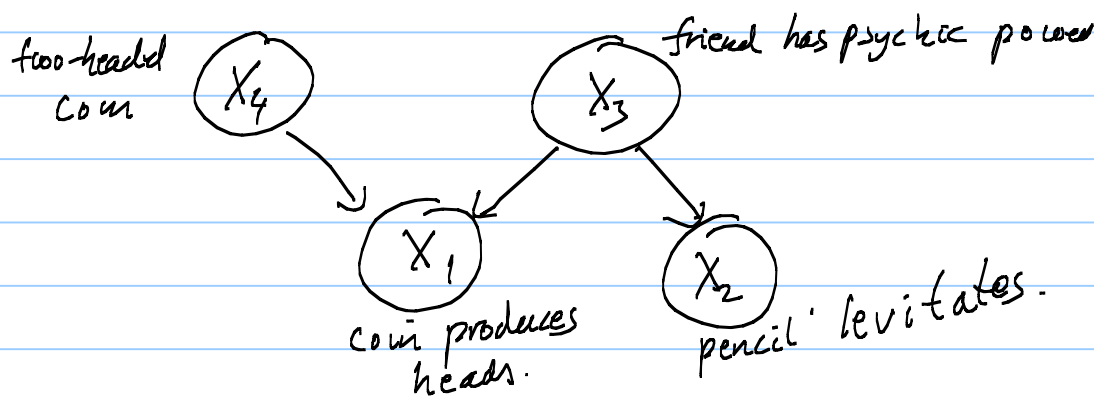
(4)

Representing Structural Probability Distributions

Probabilistic models can define the joint distribution for a set of random variables.

Friend claims psychic powers $\begin{cases} \text{test on coin tossing} \\ \text{test on pencil levitating} \end{cases}$

- X_1 represent truth of coin being flipped to give heads
- X_2 " pencil levitating.
- X_3 " friend having psychic powers
- X_4 " use of a "two-headed coin.



Represent the joint distribution:

$$P(X_1, X_2, X_3, X_4) = P(X_1 | X_3, X_4) P(X_2 | X_3) P(X_3) P(X_4)$$

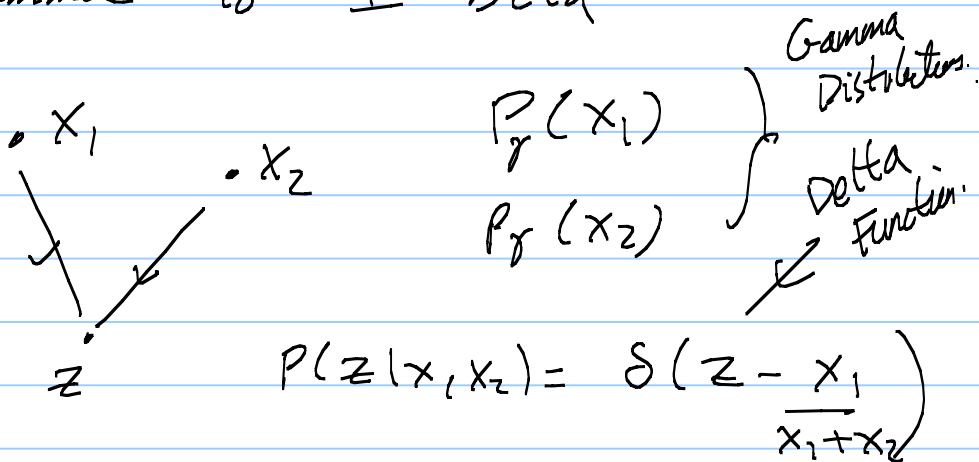
Markov Property: Graph represents the dependencies between variables. Direct & Indirection Relationship.

If you know X_3 , the knowing X_4 , won't give up on X_2 .
But, if you don't know X_3 - it will.

(Page 5)

Simple Example: Deterministic Special Case.

2 Gamma to 1 Beta



$$P(x_1, x_2, z) = P(z|x_1, x_2) P_\gamma(x_1) P_\gamma(x_2).$$

Marginal

$$P(z) = \int dx_1 \int dx_2 P(x_1, x_2, z)$$

β
Beta Distribution

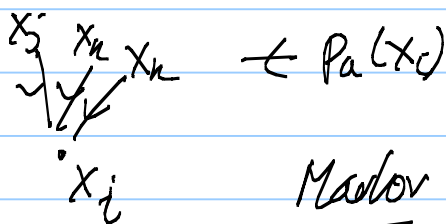
Relates to "special relationship"
between Beta & Gamma for sampling
(see lecture 2)

Page 6

Directed Bayes Network.

$$P(x_1, x_2, \dots, x_n) = \prod_i P(x_i | Pa(x_i))$$

where $Pa(x_i)$ are the parents of x_i



Martov Condition

conditioned on parents, each variable is independent of all other variables except its descendants

Factorizer \rightarrow allows us to use fewer numbers than directly specifying the full distribution $P(x_1, \dots, x_n) = 2^n - 1$

E.g. for psychic problem.

Need to specify 8 numbers, not $2^4 - 1 = 15$.

(Page 7)

Computations are also simplified by exploiting the structure.

$$\begin{aligned} \text{Eg. } P(X_1=1) &= \sum_{x_2} \sum_{x_3} \sum_{x_4} P(X_1=1, x_2, x_3, x_4) \\ &= \sum_{x_2} \sum_{x_3} \sum_{x_4} P(X_1=1 | x_3, x_4) P(x_2 | x_3) P(x_3) P(x_4) \\ &= \sum_{x_3} \sum_{x_4} P(X_1=1 | x_3, x_4) P(x_3) P(x_4). \end{aligned}$$

Sum over x_2 can be done automatically.

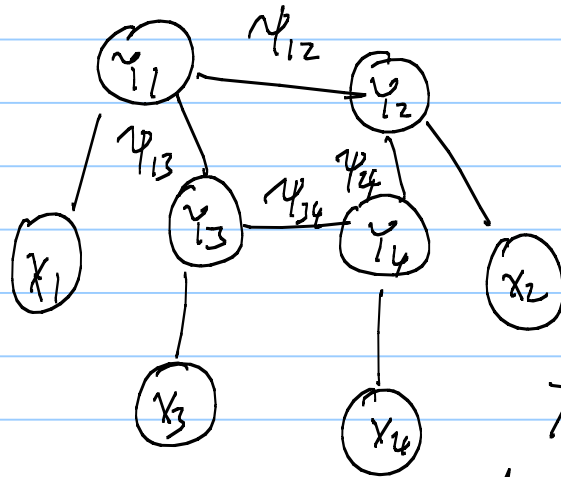
If the graph has no closed loops (tree) then dynamic programming (DP) can be used. (later in course).

Directed Graphs have many uses in Artificial Intelligence & Statistics Communities.

Page 8.

Undirected Graphical Models

Markov Random Fields (MRFs)



Undirected edges define neighbourhood structure on the graph.

These indicate the probability dependencies Markov Condition

Each fully connected neighbours is associated with a potential function

The distribution is the product of the potential functions:

$$P(\underline{x} | \underline{y}) = \prod_c P(x_i | y_i) P(\underline{y})$$

$$\text{with } P(\underline{y}) = \frac{1}{Z} \prod_{i,j \in E} \psi_{ij}(y_i, y_j) \prod_c \psi_c(\underline{x})$$

Note: Potentials not probabilities (unlike Directed Graphs).

Page 9

Example: Ising Spin Model.

$$x_1 \quad x_2 \quad x_3 \quad \dots \quad x_N \quad x_i \in \{+1, -1\}$$

$$E[x] = -J \sum_{i=1}^N x_i x_{i+1}$$

$$P[x] = \frac{1}{Z[J]} e^{J \sum_{i=1}^N x_i x_{i+1}} = \frac{1}{Z} e^{-E[x]}$$

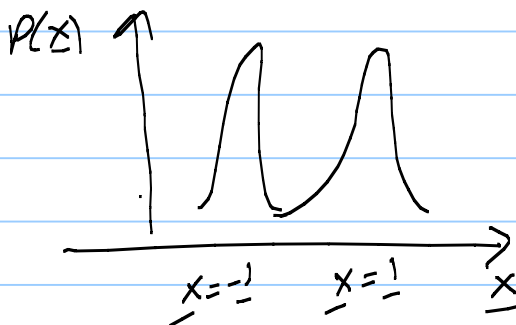
If $x_i = x_{i+1}$ $\begin{cases} (1)^2 = 1 \\ (-1)^2 = 1 \end{cases}$
Then $x_i x_{i+1} = 1$

Most probable states are lowest energy

$$x_1 = x_2 = \dots = x_N = 1 \quad \underline{x} = \underline{1}$$

$$\text{or } x_1 = x_2 = \dots = x_N = -1 \quad \underline{x} = \underline{-1}$$

The heights of the peaks depend on $J = K_T$
Large J means sharp peaks
Small J means soft peaks.



Properties. $\mu = \frac{1}{N} \sum_{i=1}^N x_i$ (net magnetization)

Physicists like to study the average magnetization $\int \mu(x) P[x]$

It has a phase transition at critical J_c .
Behaviour is quite different below J_c and above J_c

Page 10

Ising to Potts Model.

Extend to Potts Model.

$$x_i \in \{1, 2, \dots, M\}$$

$$E[x] = -J \sum_i \Phi(x_i, x_{i+1})$$

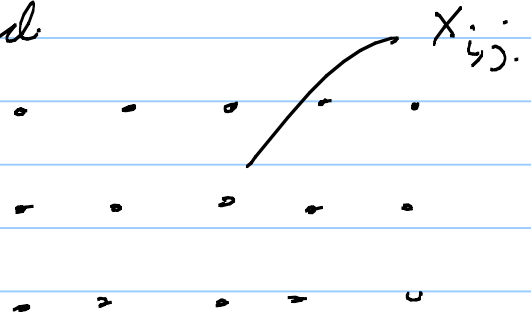
EG.

$$\Phi(x_i, x_{i+1}) = \delta_{x_i, x_{i+1}}$$

Kronecker delta.

Can be two-dimensional

$$E[x] = -J \sum_{ij} \left(\delta_{x_{ij}, x_{i+1j}} + \delta_{x_{ij}, x_{i+1j+1}} \right)$$



Many applications of Potts models
- vision, encoding/decoding, etc.

Page 11

Hidden Markov Models (HMM)

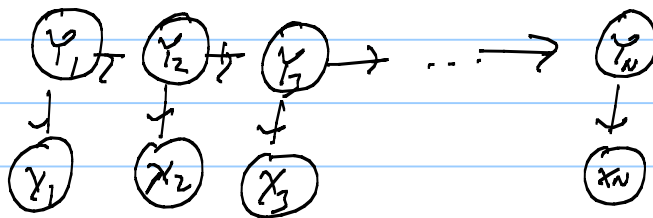
Used for speech and language processing

Sequence of T observations $\{x_t : t=1, \dots, T\}$
generated by hidden states $\{y_t : t=1, \dots, T\}$

Joint Distribution:

$$P(\{y_t\}, \{x_t\}, w) = P(w) P(y_1 | w) P(x_1 | y_1, w) \\ \prod_{t=2}^T P(y_t | y_{t-1}, w) P(x_t | y_t, w).$$

Word w .



DP Applying HMM's to recognize words requires algorithms to: (i) learn $P(x_t | y_t, w)$ & $P(y_t | y_{t-1}, w)$ for each w .

(ii) evaluate the probability:

$$P(\{x_t\}, w) = \sum_{\{y_t\}} P(\{y_t\}, \{x_t\}, w) \text{ for each word}$$

(iii) estimate $w^{\text{opt}} = \underset{w}{\text{ARG MAX}} \sum_{\{y_t\}} P(\{y_t\}, w | \{x_t\})$

Probabilistic Context-Free Grammars. (PCFG)

Define non-terminal nodes

$S, NP, VP, AT, NNS, VBD, PP, IN, DT, NO$

where S is a sentence.

VP is a verb phrase . . .

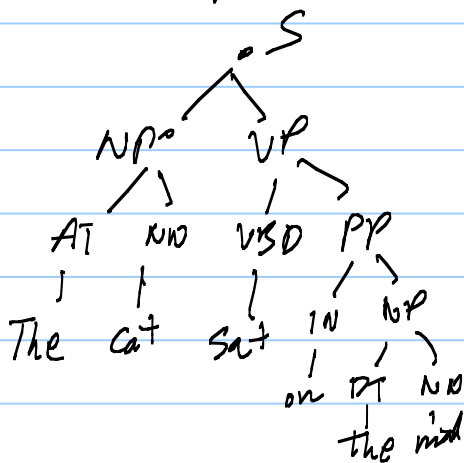
Terminal nodes are words from a dictionary

(eg. "the" "cat" "sat" "on" "the" "map").

Define production rules which are applied to non-terminal nodes to generate child nodes

eg. $S \rightarrow NP, VP$ or $NN \rightarrow \text{"cat"}$.

Define probability distributions for the production rules.



DP useful.

Generate a sentence by starting with the node S , and sampling the production rules.

Parse an input sentence by choosing the most probable parse tree.

Learn probabilities of rule

Page 13

Inference Algorithms

Want to infer values of latent (hidden variables) conditioned on data $P(y|x)$, calculate expectations, or sum out the hidden variables

$$P(\underline{x}|\underline{\theta}) = \sum_{\underline{y}} P(\underline{x}, \underline{y}|\underline{\theta})$$

Expectation Maximization (EM) Algorithms

Markov Chain Monte Carlo (MCMC)

Dynamic Programming (trees, DAGs)

This will be covered later in the course.

Converting Undirected Graphs to Directed.

If the graph structure is a tree (no closed loops) then it is possible to convert an undirected graph to a directed graph.

If graph structure has closed loops, then conversion is possible by augmenting variables. But may not be worthwhile

Decision Theory & Control Theory.

Bayes Decision Theory introduces loss function $L(h, d)$ for cost of making decision d when input is d and true hypothesis is h .

Select decision rule $d^*(\cdot)$ that minimizes risk, or expected loss,

$$R(d) = \sum_{h,d} L(h, d) P(h, d)$$

Basis of rational decision making

Loss function often set $L(h, d) = 1$, if $d \neq h$
 $L(h, d) = 0$, if $d = h$

Then best decision rule is

maximum a posteriori (MAP)

$$d^*(d) = \text{ARG MAX } P(h, d)$$

$$\text{If } L(h, d) = (h - d)^2$$

$$\text{then } d^*(d) = \sum_h h P(h, d)$$

Can extend Bayes risk to dynamical systems where decisions need to be made over time
 Gives. optimal control theory.