# Lecture 2: Probability Distributions on Graphs

Graphs enable us to represent probability distributions compactly by exploiting the dependencies between variables.

- This is helpful for understanding the structure of the data described by the distribution — enables transferring distribution from one domain to another.
- It also makes it easier to perform inference and to do learning

Example:    $P(X_1, X_2, X_3, X_4)$    suppose $X_i \in \{0, 1\}$

this distribution has $2^4 - 1 = 15$ entries

A general distribution $P(X_1 .., X_N)$ has $2^N - 1$ entries, which are far too many to learn unless we have an enormous amount of data.

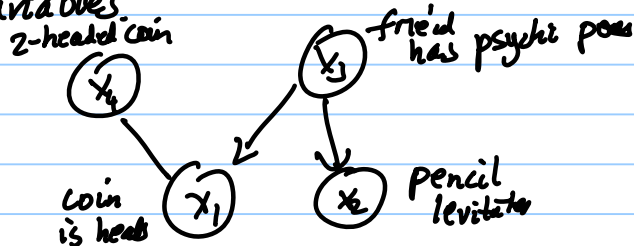But suppose we know which variables directly influence other variables.

Model the situation
friend claims psychic
powers which can predict
coin toss is head
→ but coin toss can
also be explained by 2-headed coin
→ an additional test — can friend
levitate pencil — if successful can
"explain away" the coin toss, and
justify the 2-headed coin explanation



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

~ specified by fewer
$(4 + 2 + 1 + 1 = 8 \text{ numbers})$

Knowing this structure gives:
- (i) knowledge about the problem — explaining away
- (ii) fewer numbers needed to describe distribution — less data needed to learn model
- (iii) faster inference.

(2)

<u>Inference</u> —

— normally takes $2^3 = 8$ operations.

to compute $P(X_1=1) = \sum_{X_2} \sum_{X_3} \sum_{X_4} P(X_1=1, X_2, X_3, X_4)$

exploiting graph structure.

$$= \sum_{X_2} \sum_{X_3} \sum_{X_4} P(X_1=1 \mid X_3, X_4) P(X_2 \mid X_3) P(X_3) P(X_4)$$
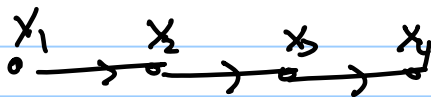
$$= \sum_{X_3} \sum_{X_4} P(X_1=1 \mid X_3, X_4) P(X_3) P(X_4) \underbrace{\sum_{X_2} P(X_2 \mid X_3)}_{1}$$

$$= \sum_{X_3} \sum_{X_4} P(X_1=1 \mid X_3, X_4) P(X_3) P(X_4)$$

only $2^2 = 4$ operations required

<u>General Result</u> if the graph structure has no closed loops, then we can use dynamic programming to compute any property of interest (e.g. $P(X_1)$) in polynomial time in no. of nodes and no. of states.

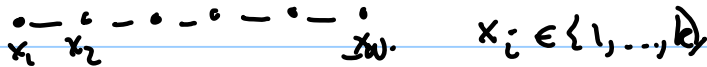<u>E.g.</u> $\quad X_1 \longrightarrow X_2 \longrightarrow X_3 \longrightarrow X_4$    rapid (polynomial) inference.

Intuition for dynamic programming — exploit the "linear structure" to break problem down into subcomponents.
E.G. to find the shortest path from Los Angeles to Boston which goes via Chicago, it is necessary only to find the shortest path from LA to Chicago and from Chicago to Boston <u>independently</u>.

More Technically — suppose we want to minimize
$$\varphi(X_1 \ldots X_N) = \varphi_{12}(X_1, X_2) + \varphi_{23}(X_2, X_3) + \ldots + \varphi_{N-1 N}(X_{N-1}, X_N)$$

$\underset{X_1}{\bullet} - \underset{X_2}{\bullet} - \bullet - \bullet - \bullet - \underset{X_N}{\bullet}$    $X_i \in \{1, \ldots, k\}$

Forward Pass of DP:

For each $X_2$, compute $h_2(X_2) = \min_{X_1} \varphi_{12}(X_1, X_2)$
        shortest cost to $X_2$

For each $X_3$, compute $h_3(X_3) = \min_{X_1 X_2} \{ \varphi_{12}(X_1, X_2) + \varphi_{23}(X_2, X_3) \}$

the clever bit $\longrightarrow$    $h_3(X_3) = \min_{X_2} \{ h_2(X_2) + \varphi_{23}(X_2, X_3) \}$

In general, $\quad h_m(X_m) = \min_{X_{m-1}} \{ h_{m-1}(X_{m-1}) + \varphi_{m-1 m}(X_{m-1}, X_m) \}$

So can compute $h_N(X_N)$ is poly time $\sim N k^2$ operators

(3)     Backward Pass of DP.

$h_N(\hat{x}_N)$ is smallest cost of $P(x_1...x_N)$
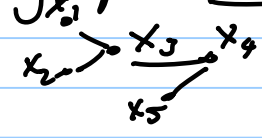
Solve   $\hat{x}_N = \arg\min h_N(x_N)$

$\hat{x}_{N-1} = \arg\min \{h_{N-1}(x_{N-1}) + \varphi_{N-1}(\hat{x}_N)\}$

recovers the states $\hat{x}_N, \hat{x}_{N-1}, ... \hat{x}_1$ which give the shortest cost.

Advantage : efficiency, $\sim k^2 N$ operators — instead of considering the total $k^N$ possible states of $P(x_1...x_N)$

Note ·(1) DP can be applied to any graph without closed loops → e.g extend to



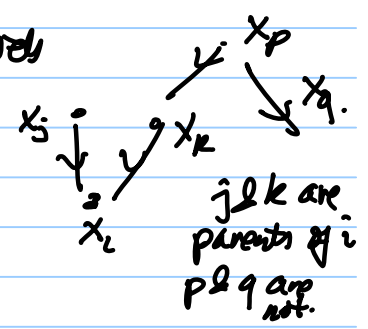(2) DP can be modified to compute other properties of interest.

lauritzen. (3) DP can be extended to graphs with
& closed loops to give the junction tree algorithm
spiegelhalter. — this includes a procedure for converting a graph to a tree. But for many graphs the tree is so large that computation on it is impractical.


Back to graphs


In general, Directed Graphs / Bayes Nets



$$P(x_1 ... x_N) = \prod_i P(x_i | Pa(x_i))$$

$Pa(x_i)$ are the parents of $x_i$, the nodes which have directed arcs directly into $x_i$

$j l k$ are parents of $i$
$p \& q$ are not.

DAG'S capture some of the causal structure of the variables in the problem.

This can include closed loops like this. →

Provided the arrows are consistent.



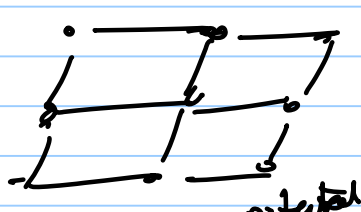not allowed

# (4) Undirected Graphs

$G = (V, E)$ — here the edges are not directed.

vertices edges.

$$P(\underline{x}) = \frac{1}{Z} \prod_{i,j \in E} \psi_{ij}(x_i, x_j) \prod_{i \in V} \psi_i(x_i)$$

$\chi$ normalization constant.

potentials. ← potentials

often write $\psi_{ij}(x_i, x_j) = e^{-\phi_{ij}(x_i, x_j)}$, etc.

$$P(\underline{x}) = \frac{1}{Z} e^{-\left( \sum_{i,j \in E} \phi_{ij}(x_i, x_j) + \sum_c \phi_i(x_i) \right)}$$

Undirected Graphs include Direct Graphs as a special case — just drop the arrows.

E.g. you can convert $A \rightarrow B \rightarrow C$ to $A - B - C$

to $A \rightarrow B \leftarrow C$.

For real causality — intervention by graph pruning can distinguish between them $A \rightarrow B \rightarrow C$ and $A \leftarrow B \leftarrow C$.

(1) $$P(x_A, x_B, x_C) = \frac{1 e^{-\left( \psi_A(x_A) + \psi_{AB}(x_A, x_B) + \psi_{BC}(x_B, x_C) \right)}}{Z}$$

(2) $P(x_C | x_B) P(x_B | x_A) P(x_A)$ directed, or $P(x_A | x_B) P(x_B | x_C) P(x_C)$

Can translate from (1) to (2) using dynamic programming
Translate from (1) to (2)   set $\psi_A(x_A) = -\log P(x_A)$
$\psi_{AB}(x_A, x_B) = -\log P(x_B | x_A)$, $\psi_{BC}(x_B, x_C) = -\log P(x_C | x_B)$

## Latent (Hidden) Variables

For both Directed and Undirected graphs some variables can be observed directly and so are 'observable', while the others are 'latent', 'hidden'.

Many, 'neural network' models can be expressed using hidden variables. Eg. Boltzmann Machine.

$x_i$ — observed
$y_j$ — hidden.

$$P(\underline{y}, \underline{x} | \underline{\omega}) = \frac{1}{Z} e^{-E[y, x | \omega]}$$

$$E[y, x, \omega] = \sum_{i,j} \omega_{ij}^o x_i y_j + \sum_{i,j} \omega_{ij}^h y_i y_j$$

(5.)    Graphical Models can also have variable
topology and no. of nodes.

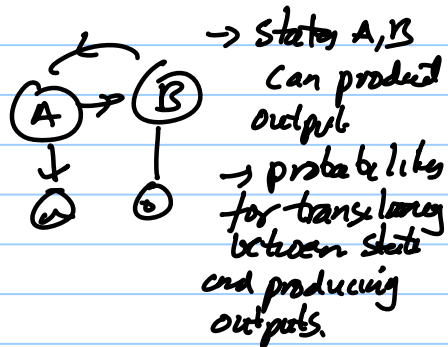E.g.    S C F G    Stochastic Context Free Grammar
Production Rules:         A -> (B, C)        A, B, C, .. Non-Trivial Nodes
                          A -> a             a, b, c terminal nodes

    assign probabilities to rules.

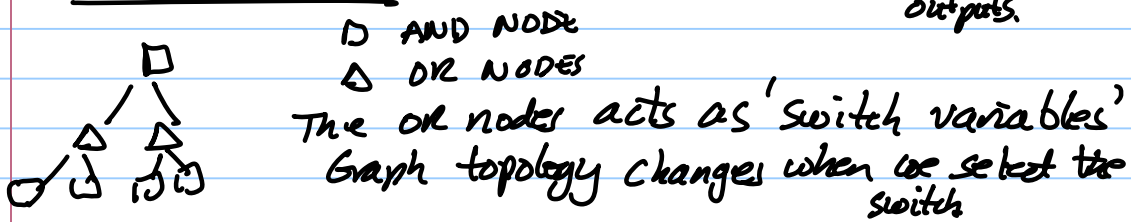root node      S          — Probability distribution over
            A.  B            the structure
          C   D   b          (see later in the course)
          |   |   b
          c   d

    Hidden Markov Models:                    -> States A, B
                                                can product
                           A -> B              output
                           |    |           -> probabilities
                          (a)  (b)             for transitioning
                                               between states
                                               and producing
                                               outputs.

    AND/OR Graphs :
                        □  AND NODE
                        △  OR NODES
                        The OR nodes acts as 'switch variables'
                        Graph topology changes when we select the
                                                            switch

For Brief Description of the
    models -> see Griffiths & Yuille handout
              -> or wait for descriptions later in the course

Inference Algorithms
    There are a range of inference algorithms that
we will describe in the next few lectures.
        Stochastic Sampling       )  these algorithms will
        Dynamic Programming       )  exploit the graph
        Steepest Descent.         )  structure
        Free Energy Methods.      )
        Graph Cuts.

## What do we want to compute?

Suppose we have $\quad P(X_1 \ldots X_N \mid data)$  ← data.

we may want to compute the marginal distribution

$$P(X_i \mid data) = \sum_{X_j \neq i} P(X_1 \ldots X_N \mid data).$$

or estimate $\quad \hat{x} = \arg\max_{x} P(\underline{x} \mid data)$

If there are hidden variables $\underline{y}$, we may want to compute

$$P(\underline{x} \mid data) = \sum_{\underline{y}} P(\underline{x}, \underline{y} \mid data)$$

or, compute $\hat{\underline{x}} = \text{ARGMAX } P(\underline{x} \mid data)$, knowing $P(\underline{x}, \underline{y} \mid data)$

→ requires the EM algorithm.

Also, we may want to estimate parameters of the model:

→ e.g. Boltzman Machine $P(\underline{x}, \underline{y} \mid \underline{\omega}) = e^{\sum_{ij} \omega_{ij}^0 x_i y_j + \sum_{ij} \omega_{ij}^0 y_i y_j}$

estimate: the $\underline{\omega}$ from a series of observations $\underline{x}^\mu$

maximize $\quad \prod_\mu P(\underline{x}^\mu \mid \omega) = \prod_\mu \sum_{\underline{y}} P(\underline{x}^\mu, \underline{y} \mid \omega)$
w.r.t. $\hat{\omega}$

sometimes called inference — but in the course we will call it learning.

easier if no hidden variables, e.g. estimate the probability that a coin yields "heads" from a set of sample coin tosses
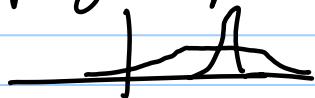
Note: most learning assumes that the form of the model is known — e.g. the graph structure. — it is usually much harder to learn the structure. But we will give some examples later in the course

Finally, model selection and Occam's factor.

Consider two models $\quad P(\underline{x}, \underline{h} \mid M_1) \quad P(\underline{x}, \tilde{\underline{h}} \mid M_2) \quad$ $\underline{h}, \tilde{\underline{h}}$ hidden variables

the prob of data $P(\underline{x}) = \sum_{\underline{h}} P(\underline{x}, \underline{h} \mid M_1)$ for $M_1$
$\qquad\qquad\qquad\qquad = \sum_{\underline{h}} P(\underline{x}, \underline{h} \mid M_2)$ for $M_2$.

This penalizes complex models and favours more precise models which fit the data.