COMS4771, Columbia University

Machine Learning

4771

Instructors: Adrian Weller and Ilia Vovsha

Lecture 20

- •Undirected Graphs
- Undirected Separation
- •Inference: Marginals/Conditionals and MAP
- Moralization
- •Junction Trees
- •Triangulation

Undirected Graphs

- ~ conditional independence
- •Separation is *much easier* for undirected graphs
- •But, what are undirected graphs and why use them?
- Might be hard to call vars parent/child or cause/effect
- •Example: Image pixels
- •Each pixel is Bernoulli = $\{0,1\}$
- •Where 0=dark, 1=bright



Have probability over all pixels p(x₁₁,...,x_{1M},...,x_{M1},...,x_{MM})
We know Bright pixels tend to have Bright neighbors
Suggests adjacent pixels dependent, so connect with links
Obtain a graphical model that looks like a grid
But who is parent? No parents, just linked probabilities
Undirected models are called Markov Random Fields
Used in vision, physics (lattice, spin, or Ising models), etc.

Undirected Graphs

- •Undirected & directed not subsets,
- •Chain Graphs are a superset.
- Some distributions behave as undirected graphs, some as directed, some as both



Directed

Graphs

an undirected graph says that $p(x_1, ..., x_M)$ satisfies any statement $X_A \parallel X_B \mid X_C$ if no paths can go from X_A to X_B unless they go through X_C



Undirected

Graphs

Thus, undirected graphs obey the general Markov propertyRecall the simple Markov property

$$x_1 - x_2 - x_3 \qquad x_1 \parallel x_3 \mid x_2 \implies p\left(x_1 \mid x_2, x_3\right) = p\left(x_1 \mid x_2\right) \qquad 4$$

Hammersley Clifford Theorem

Theorem[HC]: any distribution that obeys the Markov property

$$p\left(x_{i} \mid X_{U \setminus i}\right) = p\left(x_{i} \mid X_{Ne(i)}\right) \quad \forall i \in U$$

can be written as a product of terms over each maximal clique $p(X_U) = p(x_1, ..., x_M) = \frac{1}{Z} \prod_{c \in C} \psi_c(X_c)$

Physics analogy ~Boltzmann distn: if all probs>0

$$\psi_{c}(X_{c}) = e^{-E_{c}(X_{c})}, p(X) = \frac{1}{Z}e^{-E(X)} \text{ where } E(X) = \sum_{c \in C} E_{c}(X_{c})$$

Cliques: subsets of variables that all connect to each other Maximal: cannot add any more variables and still be a clique Each c is a maximal clique of variables X_c in the graph C is the set of all maximal cliques





Undirected Graph Functions

•Probability for undirected factorizes as a product of mini non-negative Potential Functions over cliques in the graph

$$p(X) = p(x_1, \dots, x_M) = \frac{1}{Z} \prod_{c \in C} \psi_c(X_c)$$

- •Normalizing term $Z = \sum_{X} \prod_{c \in C} \psi_c(X_c)$ makes p(X) sum to 1 •Potentials ψ are non-negative un-normalized functions over cliques (subgroups of fully inter-connected variables)
- •Only maximal cliques since smaller ψ absorb into larger ψ

$$\psi(x_2, x_3)\psi(x_2) \to \psi'(x_2, x_3) = \begin{pmatrix} 1 & 2 \\ 5 & 0 \end{pmatrix}$$



$$p\left(X\right) = \frac{1}{Z}\psi\left(x_1, x_2\right)\psi\left(x_2, x_3\right)\psi\left(x_3, x_4, x_5\right)\psi\left(x_4, x_5, x_6\right)$$

Undirected Separation Examples



Logical Inference

Classic logic network: nodes are binary
Arrows represent AND, OR, XOR, NAND, NOR, NOT etc.
Inference: given observed binary variables, predict others



Problems: uncertainty, conflicts and inconsistency
Could get x₃=T and x₃=F following two different paths
We need a way to enforce consistency and combine conflicting statements via probabilities and Bayes rule!

Probabilistic Inference

Replace logic network with Bayesian network
Tables represent AND, OR, XOR, NAND, NOR, NOT etc.
Probabilistic Inference: given observed binary variables,



•Can also have soft versions of the functions

 $x_3 = f$ $x_1 = f 0$ $x_1 = t 1$

predict marginals over others

 $\begin{array}{c|c} x_1 = f & 1 & 0 \\ x_1 = t & 0 & 1 \\ x_3 = f & x_3 = t \end{array} \\ x_5 = f \\ x_5 = f \end{array}$

 $x_3 = f x_3 = t$

soft NOT

NOT

XOR

$$x_3 = f x_3 = t$$

 $x_1 = f .1.9$
 $x_1 = t .9.1$ 9

Probabilistic Inference

- •Two types of inference with a probability distribution:
 - $p\left(X\right) = p\left(x_{1}, \dots, x_{M}\right) \text{ with queries } X_{F} \subseteq X \text{ given evidence } X_{E} \subseteq X$

•Marginal Inference:

$$p\left(X_{F} \middle| X_{E}\right) = \frac{p\left(X_{F}, X_{E}\right)}{p\left(X_{E}\right)} = \frac{\sum_{X \setminus X_{F} \cup X_{E}} p\left(X\right)}{\sum_{X \setminus X_{E}} p\left(X\right)}$$

or... $p\left(x_i \middle| X_E\right) \forall x_i \in X_F$

•Maximum a posteriori (MAP) inference: $\arg \max_{X_F} p\left(X_F \middle| X_E\right)$

~ Energy minimization

...which is harder?



COMS4771, Columbia University

MAP Example: Image De-Noising



Original Image

Noisy Image

MAP Example: Image De-Noising



Observe $\mathbf{y} = \{y_i\}$ Original(unobserved) $\mathbf{x} = \{x_i\}$ $\mathbf{x}_i, y_i \in \{-1, +1\}$ $E(\mathbf{x}, \mathbf{y}) = h \sum_{i} x_i - \beta \sum_{\{i,j\}} x_i x_j$ Signs? $-\eta \sum_{i} x_i y_i$ $p(\mathbf{x}, \mathbf{y}) = \frac{1}{Z} \exp\{-E(\mathbf{x}, \mathbf{y})\}$

COMS4771, Columbia University

MAP Example: Image De-Noising





Restored Image (ICM)

Early approximate MAP algorithm

Restored Image (Graph cuts)

More recent, exact MAP algorithm

Denoising Results



Pairwise strengths increasing

Initially we'll focus on one approach to marginal inference, Then later show how we can use the same technique for MAP...

Traditional Marginal Inference

- •Marginal inference problem: given graph and probability function $p(X) = p(x_1, ..., x_M)$ for any subsets of variables find $p(X_F | X_E) = \frac{p(X_F, X_E)}{p(X_E)}$
- So, we basically compute both marginals and divide
 But finding marginals can take exponential work!
 A problem for both directed & undirected graphs:

$$p\left(x_{j}, x_{k}\right) = \sum_{x_{1}} \sum_{x_{2}} \cdots \sum_{x_{M}} \prod_{i=1}^{M} p\left(x_{i} \mid \pi_{i}\right)$$
$$p\left(x_{j}, x_{k}\right) = \sum_{x_{1}} \sum_{x_{2}} \cdots \sum_{x_{2}} \frac{1}{\sum_{x_{M}} \sum_{x_{M}} \sum_{x_{L}} \sum_{c \in C} \psi_{c}\left(X_{c}\right)}$$
Sum over all vars other than $\mathbf{x}_{j}, \mathbf{x}_{k}$

- •Graphs gave efficient storage, learning, Bayes Ball...
- •Graphs can also be used to perform efficient inference!
- •Junction Tree Algorithm: method to efficiently find marginals

Traditional Marginal Inference

- •Example: brute force inference on a directed graph...
- •Given a directed graph structure & *filled-in* CPTs
- •We would like to efficiently compute arbitrary marginals
- •Or we would like to compute arbitrary conditionals $p(X) = p(x_{1})p(x_{2} | x_{1})p(x_{3} | x_{1})p(x_{4} | x_{2})p(x_{5} | x_{3})p(x_{6} | x_{2}, x_{5})$ $p(x_{1}, x_{3}) = p(x_{1})p(x_{3} | x_{1})$ $p(x_{1}, x_{6}) = \sum_{\substack{x_{2}, x_{3}, x_{4}, x_{5} \\ x_{2}, x_{3}, x_{4}, x_{5}}} p(x_{1})p(x_{2} | x_{1})p(x_{3} | x_{1})p(x_{4} | x_{2})p(x_{5} | x_{3})p(x_{6} | x_{2}, x_{5})$ $p(x_{1} | x_{6}) = \frac{\sum_{\substack{x_{2}, x_{3}, x_{4}, x_{5} \\ \sum_{x_{2}, x_{3}, x_{4}, x_{5}}} p(x)}{\sum_{x_{1}, x_{2}, x_{3}, x_{4}, x_{5}}} p(x)$

•For example, we may have some evidence, i.e. x_6 =TRUE

$$p(x_1 \mid x_6 = TRUE) = \frac{\sum_{x_2, x_3, x_4, x_5} p(x_{U\setminus 6}, x_6 = TRUE)}{\sum_{x_1, x_2, x_3, x_4, x_5} p(x_{U\setminus 6}, x_6 = TRUE)}$$

•This is tedious & does not exploit the graph's efficiency

Efficient Marginals & Inference

Another idea is to use some efficient graph algorithmTry sending messages (small tables) around the graph



•Hopefully these somehow settle down and equal marginals $\hat{p}(x_1, x_6) = \sum_{x_2, x_3, x_4, x_5} p(X)$

•AND marginals are self-consistent •Note: can't just return conditionals since they can be inconsistent •Junction Tree Algorithm must find consistent marginals $\sum_{x_1} \hat{p}(x_1, x_6) = \sum_{x_2} \hat{p}(x_2, x_6)$ $\sum_{x_1} \hat{p}(x_6 | x_1) \neq \sum_{x_2} \hat{p}(x_2 | x_6)$

Junction Tree Algorithm

An algorithm that achieves fast inference, by passing messages on undirected graphs.
We first convert a directed graph to an undirected one



- •Then apply the efficient Junction Tree Algorithm:
 - 1) Moralization
 - 2) Introduce Evidence
 - 3) Triangulate
 - 4) Construct Junction Tree
 - 5) Propagate Probabilities (Junction Tree Algorithm)

Moralization

 \mathcal{X}

specific

•Converts directed graph into undirected graph •By moralization, marrying the parents:

Why? 1) Connect nodes that have common children

2) Drop the arrow heads to get undirected

 $p\left(x_{1}\right)p\left(x_{2} \mid x_{1}\right)p\left(x_{3} \mid x_{1}\right)p\left(x_{3} \mid x_{1}\right)p\left(x_{4} \mid x_{2}\right)p\left(x_{5} \mid x_{3}\right)p\left(x_{6} \mid x_{2}, x_{5}\right)$

 $\rightarrow \frac{1}{Z} \psi(x_1, x_2) \psi(x_1, x_3) \psi(x_2, x_4) \psi(x_3, x_5) \psi(x_2, x_5, x_6)$



most

general

 x_6

•Note: moralization resolves *coupling* due to marginalizing moral graph is more general (loses some independencies) most

 x_3



Introducing Evidence

 \mathcal{X}_{\neg}

- •Given moral graph, note what is observed $X_E \to \overline{X}_E$ $p\left(X_F \mid X_E = \overline{X}_E\right) \equiv p\left(X_F \mid \overline{X}_E\right)$
- •If we know this is *always* observed at $X_E \rightarrow \overline{X}_E$, simplify... •Reduce the probability function since those X_E fixed •Only keep probability function over remaining nodes X_F •Only get marginals and conditionals with subsets of X_F

$$\begin{array}{ccc} & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & \\$$

Replace potential functions with slices

0.3	0.13
0.12	0.1

$$p\left(X_{F} \mid \bar{X}_{E}\right) \propto \frac{1}{Z} \psi\left(x_{1}, x_{2}, x_{3} = \bar{x}_{3}, x_{4} = \bar{x}_{4}\right) \psi\left(x_{4} = \bar{x}_{4}, x_{5}\right) \psi\left(x_{4} = \bar{x}_{4}, x_{6}\right) \psi\left(x_{4} = \bar{x}_{4}, x_{7}\right) \\ \propto \frac{1}{Z} \tilde{\psi}\left(x_{1}, x_{2}\right) \tilde{\psi}\left(x_{5}\right) \tilde{\psi}\left(x_{6}\right) \tilde{\psi}\left(x_{7}\right)$$

But, need to recompute different normalization Z... 21

Introducing Evidence

•Recall undirected separation, observing X_E separates others



•But, need to recompute new normalization ...

•Just avoid Z & normalize at the end when we are querying individual marginals and conditionals as subsets of X_F

$$\tilde{p}\left(x_{2}\right) = \frac{\sum_{x_{1},x_{5},x_{6},x_{7}}\tilde{\psi}\left(x_{1},x_{2}\right)\tilde{\psi}\left(x_{5}\right)\tilde{\psi}\left(x_{6}\right)\tilde{\psi}\left(x_{7}\right)}{\sum_{x_{2}}\sum_{x_{1},x_{5},x_{6},x_{7}}\tilde{\psi}\left(x_{1},x_{2}\right)\tilde{\psi}\left(x_{5}\right)\tilde{\psi}\left(x_{6}\right)\tilde{\psi}\left(x_{7}\right)}$$
22

23

Junction Trees

 Given moral graph want to build Junction Tree: each node is a clique (ψ) of variables in moral graph edges connect cliques of the potential functions unique path between nodes & root node (tree) between connected clique nodes, have separators (φ) separator nodes contain intersection of variables



 $p(X) = \frac{1}{Z}\psi(A, B, D)\psi(B, C, D)\psi(C, D, E)$

Triangulation

• Problem: imagine the following undirected graph



- •To ensure Junction Tree is a tree (no cycles) / R.I.P. before forming it must first Triangulate moral graph before finding the cliques...
- •Triangulating gives more general graph (like moralization)
- •Adds links to get rid of cycles or loops
- Triangulation: Connect nodes in moral graph until no chordless cycle of 4 or more nodes remains in the graph 24

Triangulation

•Triangulation: Connect nodes in moral graph such that no cycle of 4 or more nodes remains in graph



So, *add edges*, but many possible choices...
HINT: Try to keep largest clique size small

(makes junction tree algorithm more efficient)

- •Sub-optimal triangulations of moral graph are Polynomial
- •Triangulation that minimizes largest clique size is NP-hard
- •But, OK to use a suboptimal triangulation (slower JTA...)

Triangulation

•Triangulation: Connect nodes in moral graph such that no cycle of 4 or more nodes remains in graph



So, add edges, but many possible choices...
HINT: Try to keep largest clique size small (makes junction tree algorithm more efficient)
Sub-optimal triangulations of moral graph are Polynomial
Triangulation that minimizes largest clique size is NP-hard
But, OK to use a suboptimal triangulation (slower JTA...)