nhen Sr troductior Definitions pplications raphical lodels raining

CSCE 496/896 Lecture 11: Structured Prediction and Probabilistic **Graphical Models**

Stephen Scott and Vinod Variyam

(Adapted from Sebastian Nowozin and Christoph H. Lampert)

sscott@cse.unl.edu

Introduction Nebraska Out with the old ...

We've long known how to answer the question: Does this picture contain a cat?



E.g., convolutional layers feeding connected layers feeding softmax

2/74

Definitions

Graphical Models

raining



Nebraska Lincoln	Introduction and in with the new.	Nebraska	Outline
CSCE 496/896 Lecture 11: Structured Prediction and Probabilistic Graphical Models Stephen Scott and Vinod Vanyam htroduction Definitions Applications Graphical Models Iraining	What we want to know now is: Where are the cats?	CSCE 496/896 Lecture 11: Structured Prediction and Probabilistic Graphical Models Stephen Scott and Vinod Variyam Introduction Definitions Applications Graphical Models Training	 Definitions Applications Graphical modeling of probability distributions Training models Inference
3/74	・ロ・・(型・・ミ・・ミ・ 多くの)	4/74	<ロ> <日> <日> <日> <日> <日> <日> <日> <日> <日> <日



on	CSCE 496/896 Lecture 11: Structured Prediction and Probabilistic Graphical Models
ssion,	Stephen Scot and Vinod Variyam
ti a .a	Introduction
iction	Definitions
	Applications

Definitions Nebraska Structured Outputs (2)

Can think of structured data as consisting of parts, where each part contains information, as well as how they fit together

- Text: Word sequence matters
- Hypertext: Links between documents matter
- Chemical structures: Relative positions of molecules matter
- Images: Relative positions of pixels matter



Nebraska Lincoln	Applications Image Processing (2)
CSCE 496/896 Lecture 11: Structured Probabilistic Graphical Models Stephen Scott and Vinodes Stephen Scott and Vinodes Perintions Applications Graphical Models Training	Pose estimation: $f: \{mages\} \rightarrow \{K \text{ positions & angles}\}$ $a_{i} \rightarrow a_{i} \rightarrow \{K \text{ positions & angles}\}$
8/74	- ロ > (型 > (注) (注) (三) の(の



Nebraska Lincoln	Applications Image Processing (4)	
CSCE 496/896 Lecture 11: Structured rediction and Probabilistic Graphical Models Stephen Scott and Vinod Variyam troduction befinitions opplications Graphical foodels raining	Object localization $f: \{images\} \rightarrow \{bounding box coordinates\}$ input: input: image output: object position (left, top right, bottom)	
10/74		QQ

Applications

- Natural language processing (e.g., translation; output is sentences)
- Bioinformatics (e.g., structure prediction; output is graphs)
- Speech processing (e.g., recognition; output is sentences)
- Robotics (e.g., planning; output is action plan)
- Image denoising (output is "clean" version of image)

Nebraska Lincoln	Graphical Models Probabilistic Modeling
CSCE 496/096 Lecture 11: Ecture 11: Prodabilistic Graphical Models Stephen Scott and Vinod Variyam ntroduction Definitions Graphical Models Stephen Scott university Stephen Scott Carphical Stephen Scott Carphical Definitions Stephens Ste	 To represent structured outputs, we will often employ probabilistic modeling Joint distributions (e.g., P(A, B, C)) Conditional distributions (e.g., P(A B, C)) Can estimate joint and conditional probabilities by counting and normalizing, but have to be careful about representation

きょうかん 回 ふかく 山下 ふゆう ふしゃ

efinitions

Applications

raphical lodels

aininc

Graphical Models Probabilistic Modeling (2)

troduction Definitions pplication Graphical Aodels

Nebraska

F.a.	I have a	coin	with	unknown	probability	n	of h	nead	s
Ľ.g.,	i nave a	COIII	****		probability	P		icau.	-

- I want to estimate the probability of flipping it ten times and getting the sequence HHTTHHTTTT
- One way of representing this joint distribution is a single, big lookup table:
- Each experiment consists of ten coin flips
- For each outcome, incre its counter

HHTTHHTTTT's counter get the estimate Will this work?

	ten coin flips	Outcome	Count
•	For each outcome, increment	TTHHTTHHTH	1
	its counter	HHHTHTTTHH	0
•	After n experiments divide	HTTTTTHHHT	0
	HHTTHHTTTT'S counter by n to	TTHTHTHHTT	1
	get the estimate	÷	÷
_			

- **Graphical Models** Nebraska Probabilistic Modeling (3)
 - Problem: Number of possible outcomes grows exponentially with number of variables (flips)
 - \Rightarrow Most outcomes will have count = 0, a few with 1, probably none with more

- ⇒ Lousy probability estimates
- Ten flips is bad enough, but consider 100 ~
- How would you solve this problem?

Graphical Models **Graphical Models** Nebraska Nebraska Factoring a Distribution Factoring a Distribution (2) • Of course, we recognize that all flips are independent, Another example: Relay racing team so $\Pr[\texttt{HHTTHHTTTT}] = p^4 (1-p)^6$ • Alice, then Bob, then Carol • Let t_A = Alice's finish time (in seconds), t_B = Bob's, • So we can count *n* coin flips to estimate *p* and use the Variya $t_C = Carol's$ formula above ntroduction ntroduction • Want to model the joint distribution $\Pr[t_A, t_B, t_C]$ I.e., we factor the joint distribution into independent Definitions components and multiply the results: oplication Application • Let $t_C, t_B, t_A \in \{1, \ldots, 1000\}$ Graphical Models Graphical Models $\Pr[\texttt{HHTTHHTTTT}] = \Pr[f_1 = \texttt{H}] \Pr[f_2 = \texttt{H}] \Pr[f_3 = \texttt{T}] \cdots \Pr[f_{10} = \texttt{T}]$ • How large would the table be for $\Pr[t_A, t_B, t_C]$? • How many races must they run to populate the table? We greatly reduce the number of parameters to nergy eparatio estimate Training Training

Definitions

Graphical Models







pplications

Graphical Vodels

Graphical Models Factoring a Distribution (5)

 $\Pr[t_A, t_B, t_C] = \Pr[t_A] \Pr[t_B \mid t_A] \Pr[t_C \mid t_B]$

- Table for $\Pr[t_A]$ requires¹ 1000 entries, while $\Pr[t_B | t_A]$ requires 10^6 , as does $\Pr[t_C \mid t_B]$ \Rightarrow Total 2.001 \times 10⁶, versus 10⁹
- Idea easily extends to continuous distributions by changing discrete probability $\Pr[\cdot]$ to pdf $p(\cdot)$

¹Technically, we only need 999 entries, since the value of the last one is implied since probabilities must sum to one. However, then the analysis requires the use of a lot of "9"s, and that's not something I'm willing to take on at this point in my life. э



496/896 ecture 1 Structure

Definitions

pplication



 $(\forall x_i, y_i, z_k) \Pr[X = x_i \mid Y = y_i, Z = z_k] = \Pr[X = x_i \mid Z = z_k]$

more compactly, we write

$$\Pr[X \mid Y, Z] = \Pr[X \mid Z]$$

Example: Thunder is conditionally independent of Rain, given Lightning

Pr[*Thunder* | *Rain*, *Lightning*] = Pr[*Thunder* | *Lightning*]



- immediate predecessors • E.g., Given Storm and BusTourGroup, Campfire is CI of
- Lightning and Thunder ・ロト・1日・1日・1日・1日・1日・10人の



Nebraska	Directed Models	Nebraska	Directed Models
Lincoln	Generative Models	Lincoln	Predicting Most Likely Label
CSCE 496/896 Lecture 11: Structured Prediction and Probabilistic Graphical Models Stephen Scott and Vinod Variyam Introduction Definitions Applications Graphical Models <u>Directed</u> Energy Separation Training	Represents joint probability dis- tribution over $\langle Y_1,, Y_n \rangle$, e.g., $\Pr[Storm, BusTourGroup,, ForestFire]$ • In general, for y_i = value of Y_i $\Pr[y_1,, y_n] = \prod_{i=1}^{n} \Pr[y_i Parents(Y_i)]$ (Parents(Y_i) denotes immediate predecessors of Y_i) • E.g., $\Pr[S, B, C, \neg L, \neg T, \neg F] =$ $\Pr[S] \cdot \Pr[B] \cdot \Pr[C \mid B, S] \cdot \Pr[\neg L \mid S] \cdot \Pr[\neg T \mid \neg L] \cdot \Pr[\neg F \mid S, \neg L, \neg C]$ • If variables continuous, use pdf $p(\cdot)$ instead of $\Pr[\cdot]$	CSCE 496/896 Lecture 11: Structured Prediction and Probabilistic Graphical Models Stephen Scott and Vinod Variyam Introduction Definitions Applications Graphical Models Directed Emergy Emergy Emergy Training	 We sometimes call graphical models generative (vs discriminative) models since they can be used to generate instances ⟨Y₁,,Y_n⟩ according to joint distribution Can use for classification Label <i>r</i> to predict is one of the variables, represented by a node If we can determine the most likely value of <i>r</i> given the rest of the nodes, can predict label One idea: Go through all possible values of <i>r</i>, and compute joint distribution (previous slide) with that value and other attribute values, then return one that maximizes

Definitions

pplications

aphical

Directed Models Predicting Most Likely Label (cont'd)



E.g., if *Storm* (*S*) is the label to predict, and we are given values of *B*, *C*, $\neg L$, $\neg T$, and $\neg F$, can use formula to compute $\Pr[S, B, C, \neg L, \neg T, \neg F]$ and $\Pr[\neg S, B, C, \neg L, \neg T, \neg F]$, then predict more likely one

Easily handles unspecified attribute values

Issue: Takes time exponential in number of values of unspecified attributes

More efficient approach: **Pearl's message passing** algorithm for chains and trees and polytrees (at most one path between any pair of nodes)

Nebraska Undirected Models

Definitions

- Since directed edges imply causal relationships, might want to use undirected edges if causality not modeled
 - E.g., let h_y = 1 if you are healthy, 0 if sick
 h_r same but for your roommate, h_c for coworker
 - h_y and h_r directly influence each other, but causality unknown and irrelevant
 - *h_y* and *h_c* also directly influence each other
 - h_r and h_c only indirect influence, via h_y
- Can model $\Pr[h_r, h_y, h_c]$ with undirected model, aka Markov random field (MRF), aka Markov network





Nebraska Lincoln	Undirected Models Factors (3)								
CSCE 496/896 Lecture 11: Structured Prediction and Probabilistic Graphical Models	Model: $\phi(\mathcal{C}_{ry}) \mid h_y$ $h_r = 0$ $h_r = 1$	$\frac{=0}{2}$	<i>h</i> _y	= 1 1 0	. (h.)(h _y)(h _t	$\underbrace{\frac{\phi(\mathcal{C}_{yc})}{h_c = 0}}_{h_c = 1}$	$\frac{h_y = 0}{5}$	$\frac{h_y = 1}{15}$
Stephen Scott and Vinod Variyam	Distribution:								
		h.	h.,	h.	$\phi(C_{m})$	$\phi(C_{m})$	$\tilde{P}(\mathbf{v})$	Pr[v]	
Introduction	-	0	0	0	2	5	10	0.051	
Definitions		Ō	ō	1	2	2	4	0.020	
		0	1	0	1	1	1	0.005	
Applications		0	1	1	1	15	15	0.076	
Graphical		1	0	0	1	5	5	0.025	
Directed		1	0	1	1	2	2	0.010	
Undirected		1	1	0	10	1	10	0.051	
Energy		1	1	1	10	15	150	0.762	
- · ·	-						Z = 197	1.0	
Training	What is time complexity of brute-force approach?								
00/74									



onhen S

and Vino Variyam

Applications araphical Iodels Directed Undirected

Graphic Models

and Vinc Variyarr

ntroduction

Applications Graphical Models Directed Undirected

raining

Undirected Models Factor Graphs (2)

- Formally, a factor graph is a bipartite graph $(V, \mathcal{F}, \mathcal{E})$, where V = variable nodes, $\mathcal{F} =$ factor nodes and edges $\mathcal{E} \subseteq V \times \mathcal{F}$ with one endpoint *V* and one in \mathcal{F}
- The **scope** $N : \mathcal{F} \to 2^V$ of factor $f \in \mathcal{F}$ is the set of neighboring variables:

$$N(f) = \{i \in V : (i,f) \in \mathcal{E}\}$$

• Now compute distribution similar to before:

$$\Pr[\mathbf{y}] = \frac{1}{Z} \prod_{f \in \mathcal{F}} \phi_f(\mathbf{y}_{N(f)})$$

EVALUATE: For the set of
$$f(y_{N(f)})$$

 $f(y_{N(f)})$
Undirected Models
Definitional Random Fields
Conditional Random
field (CRF) is a factor
graph Used to directly
model a conditional
distribution
 $Pr[Y = y | X = x]$
Conditional Random
field (CRF) is a factor
graph Used to directly
model a conditional
distribution
 $Pr[Y = y | X = x]$
Conditional Random
Field
Conditional Random
field (CRF) is a factor
graph Used to directly
for

ected Models Based Functions	Nebraska Lincoln	Undirected Models Energy-Based Functions (2)
We now know how to factor the distribution graphically, ut what form will $\phi(\cdot)$ take? Want to learn them to infer a distribution eed $\tilde{p}(\mathbf{x}) > 0$ for all \mathbf{x} in order to get a distribution efine an energy function $E_f : \mathcal{Y}_{N(f)} \to \mathbb{R}$ for factor f hen define $\phi_f = \exp(-E_f(y_f)) > 0$ and get $p(Y = \mathbf{y}) = \frac{1}{Z} \prod_{f \in \mathcal{F}} \phi_f(y_f) = \frac{1}{Z} \prod_{f \in \mathcal{F}} \exp\left(-E_f(y_f)\right)$ $= \frac{1}{Z} \exp\left(-\sum_{f \in \mathcal{F}} E_f(y_f)\right)$	CSCE 499/896 Lecture 11: Structured Prediction and Probabilistic Graphical Models Stephen Scott and Vinod Variyam Introduction Definitions Applications Graphical Definitions Graphical Definitions Applications Graphical Definitions Training	Using this form of ϕ allows us to fa as well! a b d c $E(a,b,c,d,e,f) = E_{a,b}(a,b) + E_{b,c}(b,c)$

→□>→■>→E>→E> E→ 9000

Undir Nebraska Energy-I

- W bı
- W
- N
- De
- TI

$$p(Y = \mathbf{y}) = \frac{1}{Z} \prod_{f \in \mathcal{F}} \phi_f(y_f) = \frac{1}{Z} \prod_{f \in \mathcal{F}} \exp\left(-E_f(y_f)\right)$$
$$= \frac{1}{Z} \exp\left(-\sum_{f \in \mathcal{F}} E_f(y_f)\right)$$

actor our energy function



 $+E_{a,d}(a,d)+E_{b,e}(b,e)+E_{e,f}(e,f)$

+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
+□ >
<p

Nebraska Lincoln	Undirected Models Energy-Based Functions (3)	Nebraska Lincoln	Separation and D-Separation
CSCE 496/896 Lecture 11: Structured Probabilistic Graphical Models Stephen Scott and Vinod Variyam Introduction Definitions Applications Graphical Models Undirected Undirected Steparation Training	 Still need a form for E(·) to parameterize and learn Define E_f(y_f; w) to depend on weight vector w ∈ ℝ^d: E_f : 𝔅_{N(f)} × ℝ^d → ℝ E.g., say we are doing binary image segmentation Want adjacent pixes to try to take same value, so define E_f : {0, 1} × {0, 1} × ℝ² → ℝ as E_f(0, 0; w) = E_f(1, 1; w) = w₁ E_f(0, 1; w) = E_f(0, 1; w) = w₂ We learn w₁ and w₂ from training data, expecting w₁ > w₂ More on this later 	CSCE 496/896 Lecture 11: Structured Prediction and Probabilistic Graphical Models Stephen Scott and Vinod Vanyam Introduction Definitions Applications Graphical Models Directed Understed Energy Separation	 An edge between two nodes indicates a direct interaction between the variables Paths between nodes indicate indirect interactions Observing (instantiating) some variables change the interactions between others Useful to know which subsets of variables are conditionally independent from each other, given values of other variables Say that set of variables A is separated (if undirected model) or d-separated (if directed) from set B given set S if the graph implies that A and B are conditionally independent given S
35/74	< □ > < 图 > < 差 > < 差 > のへの	36/74	< □ > < 個 > < 置 > < 置 > < 置 > のQ(の



Separation and D-Separation Nebraska Separation in Undirected Models

- If a variable is observed. it blocks all paths through it
- In an undirected model, two nodes are separated if all paths between them are blocked



・ロト・日下・日下・日、 日、 のへの

• E.g., a and c are blocked, as are d and c, but not a and d (even though one of their paths is blocked)

Nebraska

troduction

efinitions

pplication

Separation and D-Separation D-Separation in Directed Models

- In directed models, d-separation is more complicated
- Depends on the direction of the edges involved
- When considering nodes a and b connected via c, can classify connection as tail-to-tail, head-to-tail, and head-to-head
- For each case, assuming no other path exists (ignoring) edge direction) between a and b, we will determine if aand b are independent, or conditionally independent given c





+□> <畳> < Ξ> < Ξ> < Ξ> < Ξ</p>









Nebraska Markov Blankets



Learning Graphical Models Nebraska Conditional Random Fields

Definitions

pplication

Graphical Aodels

Training

 Learning a CRF with input x, parameterized by weight vector w:

$$\Pr[\mathbf{y} \mid \mathbf{x}, \mathbf{w}] = \frac{1}{Z(\mathbf{x}, \mathbf{w})} \exp\left(-E(\mathbf{y}, \mathbf{x}, \mathbf{w})\right)$$

where $Z(\mathbf{x}, \mathbf{w}) = \sum_{\mathbf{y} \in \mathcal{Y}} \exp\left(-E(\mathbf{y}, \mathbf{x}, \mathbf{w})\right)$

- Let energy function $E(\mathbf{y}, \mathbf{x}, \mathbf{w}) = \langle \mathbf{w}, \varphi(\mathbf{x}, \mathbf{y}) \rangle$ • I.e., a weighted sum of features produced by feature function $\varphi(x, y)$
 - $\varphi(\mathbf{x}, \mathbf{y})$ could be a deep network, possibly trained earlier • w is trained to get $\Pr_P[y \mid x, w]$ "close" to the true distribution $\Pr_D[y \mid x]$



Nebraska Lincoln	Learning Graphical Models Conditional Random Fields: RMCL				
CSCE 496/896 Lecture 11: Structuredd Probabilistic Graphical Models Stephen Scott and Vinod Variyam Introduction Definitions Applications Graphical Models Training	 I.e., we choose a model (w*) that maximizes the conditional log likelihood of the data If all (x,y) instances are drawn iid, then w* maximizes the probability of seeing all the ys given all the xs Throw in a regularizer for good measure Definition: Let Pr[y x, w] = 1/(Z(x,w)) exp(-⟨w, φ(x,y)⟩) be a probability distribution parameterized by w ∈ ℝ^d and let D = {(xⁿ, yⁿ)}_{n=1,,N} be a set of training examples. For any λ > 0, regularized maximum conditional likelihood (RMCL) training chooses 				
53/74	$w \in \mathbb{R}^d \qquad n=1 \qquad \qquad$				

Learning Graphical Models Conditional Random Fields: RMCL (2) Nebraska Goal: find w minimizing $\mathcal{L}(\boldsymbol{w}) = \lambda \|\boldsymbol{w}\|^2 + \sum_{n=1}^N \langle \boldsymbol{w}, \varphi(\boldsymbol{x}^n, \boldsymbol{y}^n) \rangle + \sum_{n=1}^N \log Z(\boldsymbol{x}^n, \boldsymbol{w})$

Compute the gradient:

Definitions pplicatior Graphical Nodels Training

$$\begin{aligned} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) &= 2\lambda \mathbf{w} + \sum_{n=1}^{N} \left[\varphi(\mathbf{x}^{n}, \mathbf{y}^{n}) - \sum_{\mathbf{y} \in \mathcal{Y}} \left(\frac{\exp(-\langle \mathbf{w}, \varphi(\mathbf{x}^{n}, \mathbf{y}) \rangle)}{\sum_{\mathbf{y}' \in \mathcal{Y}} \exp(-\langle \mathbf{w}, \varphi(\mathbf{x}^{n}, \mathbf{y}') \rangle)} \right) \varphi(\mathbf{x}^{n}, \mathbf{y}) \right] \\ &= 2\lambda \mathbf{w} + \sum_{n=1}^{N} \left[\varphi(\mathbf{x}^{n}, \mathbf{y}^{n}) - \sum_{\mathbf{y} \in \mathcal{Y}} \Pr[\mathbf{y} \mid \mathbf{x}^{n}, \mathbf{w}] \varphi(\mathbf{x}^{n}, \mathbf{y}) \right] \\ &= \left[2\lambda \mathbf{w} + \sum_{n=1}^{N} \left[\varphi(\mathbf{x}^{n}, \mathbf{y}^{n}) - \mathsf{E}_{\mathbf{y} \sim P(\mathbf{y} \mid \mathbf{x}^{n}, \mathbf{w})} \left[\varphi(\mathbf{x}^{n}, \mathbf{y}) \right] \right] \end{aligned}$$

Nebraska

Learning Graphical Models Conditional Random Fields: RMCL (3)

- Definitions polication
- The gradient has a nice, compact form, and is convex \Rightarrow Any local optimum is a global one
- Problem: Computing expectation requires summing over exponentially many combinations of values of y
- We can factor energy function, and therefore its derivative, and therefore the expectation of its derivative
- Let's focus on an individual factor *f*:

$$\mathsf{E}_{\mathbf{y}_{f} \sim P(\mathbf{y}_{f} \mid \mathbf{x}^{n}, \mathbf{w})} \left[\varphi_{f}(\mathbf{x}^{n}, \mathbf{y}_{f}) \right] = \sum_{\mathbf{y}_{f} \in \mathcal{Y}_{f}} \Pr_{P}(\mathbf{y}_{f} \mid \mathbf{x}, \mathbf{w}) \varphi_{f}(\mathbf{x}^{n}, \mathbf{y}_{f})$$

- Summation still has exponentially many terms, but instead of $K^{|V|}$ now it's $K^{|N(f)|}$ (more manageable)
- Still need to compute each factor's marginal probability

Learning Graphical Models Nebraska Inference

Definitions

pplication

raphical Training

- Efficient inference of marginal probabilities and Z in a graphical model is itself a major research area
- Depends on the structural model we're using
- Start with belief propagation in acyclic models
- Then approximate loopy belief propagation for cyclic models

Learning Graphical Models Inference: Sum-Product Algorithm Learning Graphical Models Nebraska Nebraska Inference: Sum-Product Algorithm (2) • Belief propagation is a general approach to inference in directed and undirected graphical models Generally, some node i sends a message to another $M(i) = \{ f \in \mathcal{F} : (i, f) \in \mathcal{E} \}$ node *j* regarding *i*'s belief about variable *y* • *i* informs *j* its belief about marginal probability Pr[y]• E.g., message value high \Rightarrow belief is $\Pr[y]$ also high adjacent to i · Each node messages each of its neighbors about its belief for each value of the random variable roduction • Sum-Product Algorithm uses belief propagation to efinitions oplication message is find marginal probabilities and Z in tree-structured factor graphs (connected and acyclic) • Each edge $(i, f) \in \mathcal{E} \subset V \times \mathcal{F}$ has 6 Training $\mathbf{0} \quad q_{Y_i \to f} \in \mathbb{R}^{|\mathcal{Y}_i|} \text{ is a variable-to-factor message}$ 2 $r_{f \to Y_i} \in \mathbb{R}^{|\mathcal{Y}_i|}$ is a factor-to-variable message Note they are vector quantities, one component per

value of Y_i

Variable-to-Factor Message • For variable $i \in V$, let

be the set of factors

 $r_{B \to Y}$

For each value y_i of variable i, variable-to-factor

$$q_{Y_i \to f}(y_i) = \sum_{f' \in \mathcal{M}(i) \setminus \{f\}} r_{f' \to Y_i}(y_i)$$

• Variable node *i* sums up all factor-to-variable messages from all factors except f and transmits result to f



braska Lincoln	Learning Graphical Models Inference: Sum-Product Algorithm (4)
CSCE 96/896 cture 11: ructured liction and babilistic raphical Models	 Since we have a tree structure, there is always at least one variable adjacent to only one factor or one factor adjacent to one variable
hen Scott d Vinod	 These messages depend on nothing, so start there
anyam	 Then order the other message computations via
duction	precedence graph

- Designate an arbitrary variable node to be the root
- Two phases of algorithm:
 - Leaf-to-root phase: start at leaves and compute messages toward root
 - 8 Root-to-leaf phase: start at root and compute messages toward leaves



Learning Graphical Models Inference: Sum-Product Algorithm (7)

To compute factor marginals:

Nebraska

Nebraska

efinitions

raphical odels

Training

Learning Graphical Models Nebraska Inference: Sum-Product Algorithm (6)

To **compute** Z, sum over factor-to-variable messages directed to root Y_r :

$$\log Z = \log \sum_{y_r \in \mathcal{Y}_r} \exp \left(\sum_{f \in M(r)} r_{f \to Y_r}(y_r) \right)$$

Learning Graphical Models Inference: Sum-Product Algorithm (8) Nebraska







 $q_{Y_i \to F}$

Learning Graphical Models Nebraska Inference: Loopy Belief Propagation

- When graph has a cycle, can still perform message passing to approximate Z and marginal probabilities
- Initialize messages to fixed value
- Perform updates in random order until convergence
- Factor-to-variable messages $r_{f \rightarrow Y_i}$ computed as before
- Variable-to-factor messages computed differently



496/896 Lecture 11 Structured

Graphica Models

Stephen So and Vino Variyam ntroductior Definitions

pplications

raphical odels

Training

Learning Graphical Models Inference: Loopy Belief Propagation (2)

Variable-to-factor messages:

$$\begin{array}{lll} \bar{q}_{Y_i \to f}(y_i) &=& \displaystyle \sum_{f' \in \mathcal{M}(i) \setminus \{f\}} r_{f' \to Y_i}(y_i) \\ \delta &=& \displaystyle \log \sum_{y_i \in \mathcal{Y}_i} \exp \left(\bar{q}_{Y_i \to f}(y_i) \right) \\ q_{Y_i \to f}(y_i) &=& \bar{q}_{Y_i \to f}(y_i) - \delta \end{array}$$

Nebraska Inference: Loopy Belief Propagation (3)

CSCE 496/896 Lecture 11: Structured

Definitions

pplication

raphical

Training

68/74

To compute factor marginals:

$$\begin{split} \bar{\mu}_f(\mathbf{y}_f) &= -E_f(\mathbf{y}_f) + \sum_{j \in N(f)} q_{Y_j \to f}(y_j) \\ z_f &= \log \sum_{\mathbf{y}_f \in \mathcal{Y}_f} \exp(\bar{\mu}_f(\mathbf{y}_f)) \\ \mu_f(\mathbf{y}_f) &= \exp\left(\bar{\mu}_f(\mathbf{y}_f) - z_f\right) \end{split}$$

Nebraska

Graphical Models

Variyam Introduction Definitions Applications Graphical Models Training

Learning Graphical Models
Inference: Loopy Belief Propagation (4)

To compute variable marginals:

$$\begin{split} \bar{\mu}_i(y_i) &= \sum_{f' \in \mathcal{M}(i)} r_{f' \to Y_i}(y_i) \\ z_i &= \log \sum_{y_i \in \mathcal{Y}_i} \exp(\bar{\mu}_i(y_i)) \\ \mu_i(y_i) &= \exp(\bar{\mu}_i(y_i) - z_i) \end{split}$$

Nebraska Lincoln	Learning Graphical Models Inference: Loopy Belief Propagation (5)				
CSCE 496/896 Lecture 11: Structured Prediction and Probabilistic Graphical Models	To compute Z:				
Stephen Scott and Vinod Variyam	$\log Z = \sum_{i} \left(M(i) - 1 \right) \left[\sum_{i} \mu_i(y_i) \log \mu_i(y_i) \right]$				
Introduction	$i \in V$ $y_i \in \mathcal{Y}_i$				
Definitions					
Applications	$-\sum \sum \mu_f(\mathbf{y}_f)(E_f(\mathbf{y}_f) + \log \mu_f(\mathbf{y}_f))$				
Graphical Models	$f \in \mathcal{F} \mathbf{y}_f \in \mathcal{Y}_f$				
Training					



70/74

```
(ロ)・(問)・(目)・(目)、目、の()
```

(ロ)

CSCE 496/896 Lecture 11: Structured Prediction and Probabilistic Graphical Models Stephen Scott and Vinod Variyam

Definitions Applications

Graphical Models

Training

Learning Graphical Models Conditional Random Fields: Case Study (2): CRF

• Energy function:

$$E(\mathbf{y}) = \sum_{i} \theta_{i}(y_{i}) + \sum_{i,j} \theta_{ij}(y_{i}, y_{j})$$

where $y_i \in \{0, 1\}$ is label assignment for pixel *i* • Use $\theta_i(y_i) = -\log P(y_i)$ and

г

$$\theta_{ij}(y_i, y_j) = \mu(y_i, y_j) \left[w_1 \exp\left(-\frac{\|p_i - p_j\|^2}{2\sigma_{\alpha}^2} - \frac{\|I_i - I_j\|^2}{2\sigma_{\beta}^2} \right) + w_2 \exp\left(-\frac{\|p_i - p_j\|^2}{2\sigma_{\gamma}^2} \right) \right]$$

where

- $\mu(y_i, y_j) = 1$ iff $y_i \neq y_j$ (different labels) $p_i = \text{position of pixel } i$ $I_i = \text{RGB color of pixel } i$
- $\sigma = parameters$
- Inference via specialized algorithms for
- Gaussian-based functions (日)(個)(目)(目)(日)(0)(0)

Nebraska Lincoln	Learning Graphical Models Conditional Random Fields: Case Study (3): CRF Training Example							
CSCE 496/896 Lecture 11: Structured Prediction and Probabilistic Graphical Models								
Stephen Scott and Vinod Variyam	E T	Å	÷	+	+			
Introduction	Image/G T	DCNN output	CRF Iteration 1	CRF Iteration 2	CRF Iteration 10			
Applications								
Graphical Models								
Training				() <) < ≥ > ≥ √Q (0			