





Advanced Multimodal Machine Learning

Lecture 10.1: Probabilistic

Graphical Models

Louis-Philippe Morency Tadas Baltrusaitis

Lecture Objectives

- Probabilistic Graphical Models
- Markov Random Fields
 - Boltzmann/Gibbs distribution
 - Factor graphs
- Conditional Random Fields
 - Multi-View Conditional Random Fields
- CRFs and Deep Learning
 - DeepConditional Neural Fields
 - CRF and Bilinear LSTM
- Continuous and Fully-Connected CRFs

Administrative Stuff

Lecture Schedule

Classes	Lectures
Spring break 3/13 – 3/17	
Week 9	Multimodal alignment
3/21 & 3/23	 Attention models and multi-instance learning
	 Multimodal synchrony and prediction
Week 10	Probabilistic Graphical Models
3/28 & 3/30	 Markov random field and Boltzmann machines
	 Latent, continuous and fully-connected CRFs
Week 11	Mid-term project assignment - Presentations
4/4 & 4/6	

Lecture Schedule

Classes	Lectures
Week 12	Multimodal fusion
4/11 & 4/13	 Multi-kernel learning and fusion
	 Sample-based late fusion
Week 13	Advanced multimodal representations
4/18 & 4/20	 Image and video description
	 Guest lecture: Ruslan Salakhutdinov
Week 14	Multilingual computational models
4/25 & 4/27	 Neural Machine Translation
	 Guest lecture: Graham Neubig
Week 15 5/2 & 5/4	Final project assignment - Presentations

Upcoming Schedule

- Pre-proposal (tomorrow Wednesday 2/5 at 9am)
- First project assignment:
 - Proposal presentation (2/21 and 2/23)
 - First project report (Sunday 3/5)
- Second project assignment
 - Midterm presentations (4/4 and 4/6)
 - Midterm report (Sunday 4/9)
- Final project assignment
 - Final presentation (5/2 & 5/4)
 - Final report (Sunday 5/7)

Midterm Presentation Instructions

- 8-9 minute presentations (12-18 slides)
 - +2-3 minutes for written feedback and notes
- All team members should be involved during the presentation.
- The ordering of the presentations (Tuesday vs. Thursday) will be inverted based on proposal presentations.
- The presentations will be from 4:30pm 6:15 to give everyone more time
 - Let us know if you need to leave before then

Midterm Project Report Instructions

PART 1

- Research problem: describe and motivate the research problem you are planning to work on. Explain why this problem is important for the research community and, if possible, the society in general. Define in generic terms the main computational challenges involved in this research problem.
- Related Work: Present an overview of the work happening in this research area. This section should include about 12-15 citations of prior work, grouped in similar topics. Also, you should present in more details the 3-4 research papers most related to your proposed work. The related work section should end emphasizing how your proposed approach differ from previous work.
- Dataset and Input Modalities: Describe the dataset(s) you are planning to use for this project. If many options exist, please motivate your choice of dataset for this research problem. Describe the input modalities and annotations available in this dataset. Specify which subset of these modalities and annotations you are planning to use.

Read carefully all the comments we gave you in the proposal reports. We will be stricter for the midterm reports.

Midterm Project Report Instructions

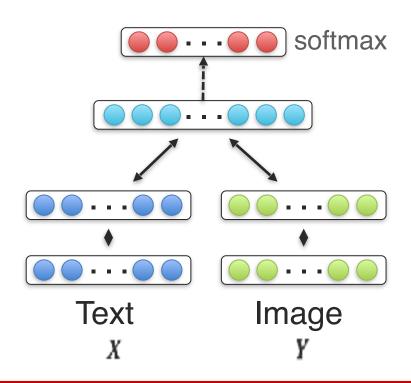
PART 2

- Problem statement: formalize mathematically your research problem. This should include the mathematical definition of the variables involved in your problem.
- Multimodal baseline models: Describe mathematically at least one multimodal baseline model for your research problem.
- Experimental methodology: Describe your experimental methodology for evaluating the multimodal baseline model(s).
- Results and Discussion: Present in tables and/or figures your experimental results. This section should include more than re-running existing baseline models.
- Proposed approach: Describe what models you are planning to test for the final report experiments. Whenever possible, you should write down the loss function of these models, following the same mathematical formulation previously used.

Quick Recap

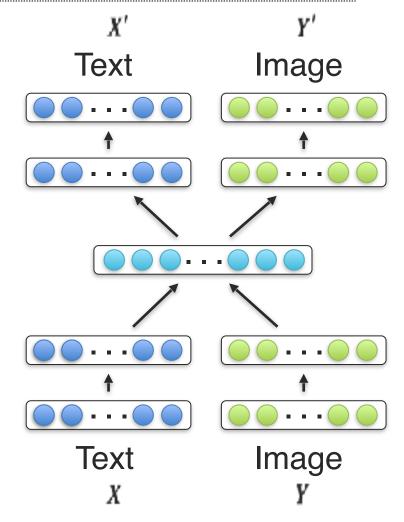
Learn (unsupervised) a joint representation between multiple modalities where similar unimodal concepts are closely projected.

Deep MultimodalBoltzmann machines



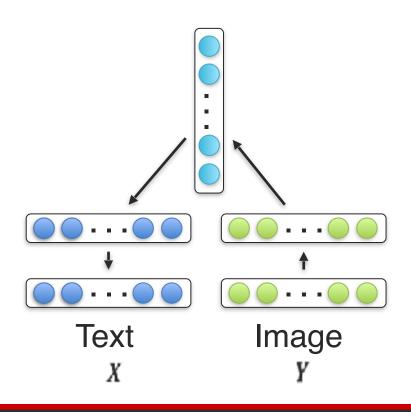
Learn (unsupervised) a joint representation between multiple modalities where similar unimodal concepts are closely projected.

- Deep Multimodal Boltzmann machines
- Stacked Autoencoder



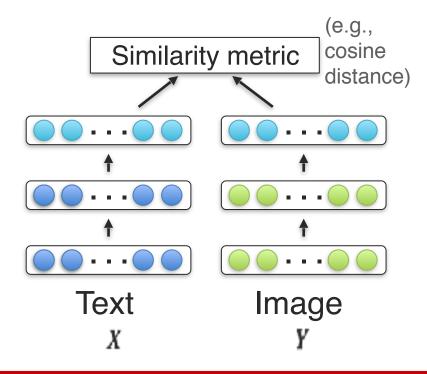
Learn (unsupervised) a joint representation between multiple modalities where similar unimodal concepts are closely projected.

- Deep Multimodal Boltzmann machines
- Stacked Autoencoder
- Encoder-Decoder

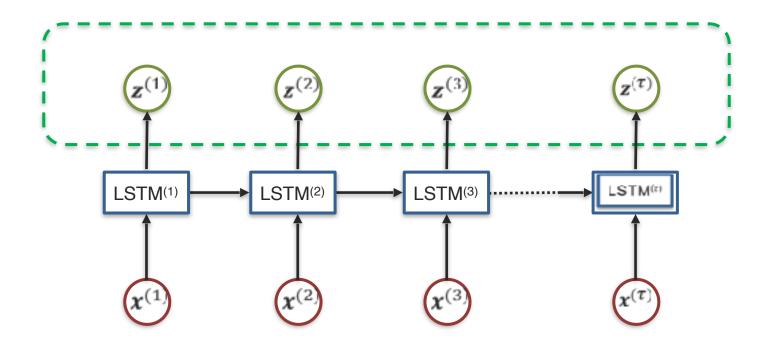


Learn (unsupervised) a joint representation between multiple modalities where similar unimodal concepts are closely projected.

- Deep MultimodalBoltzmann machines
- Stacked Autoencoder
- Encoder-Decoder
- "Minimum-distance"Multimodal Embedding



Recurrent Neural Network using LSTM Units



How can we improve reasoning by including prior domain knowledge?



Probabilistic Graphical Models

Probabilistic Graphical Model

Definition: A probabilistic graphical model (PGM) is a graph formalism for compactly modeling joint probability distributions and dependence structures over a set of random variables.

- Random variables: X₁,...,X_n
- P is a joint distribution over X₁,...,X_n

Can we represent P more compactly?

Key: Exploit independence properties

Independent Random Variables

- Two variables X and Y are independent if
 - P(X=x|Y=y) = P(X=x) for all values x,y
 - Equivalently, knowing Y does not change predictions of X
- If X and Y are independent then:

$$P(X, Y) = P(X|Y)P(Y) = P(X)P(Y)$$





- If $X_1,...,X_n$ are independent then:
 - $P(X_1,...,X_n) = P(X_1)...P(X_n)$

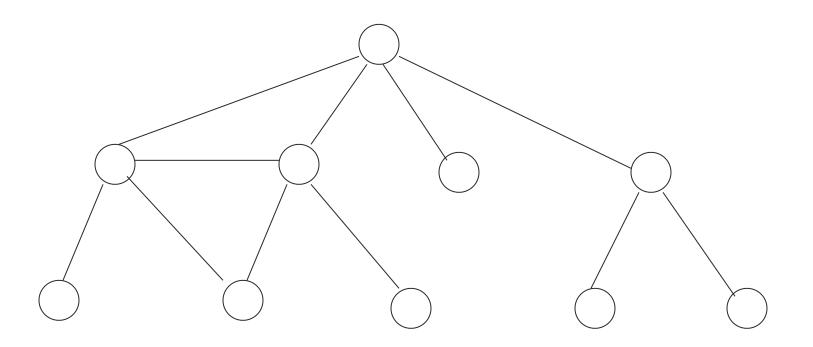
Conditional Independence

- X and Y are conditionally independent given Z if
 - P(X=x|Y=y, Z=z) = P(X=x|Z=z) for all values x, y, z
 - Equivalently, if we know Z, then knowing Y does not change predictions of X

Graphical Model

- A tool that visually illustrate <u>conditional</u> <u>dependence</u> among variables in a given problem.
- Consisting of nodes (Random variables or States) and edges (Connecting two nodes, directed or undirected).
- The lack of edge represents conditional independence between variables.

Graphical Model



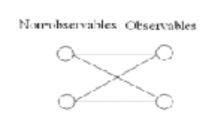
 Chain, Path, Cycle, Directed Acyclic Graph (DAG), Parents and Children

Reasoning

 The activity of guessing the state of the domain from prior knowledge and observations.

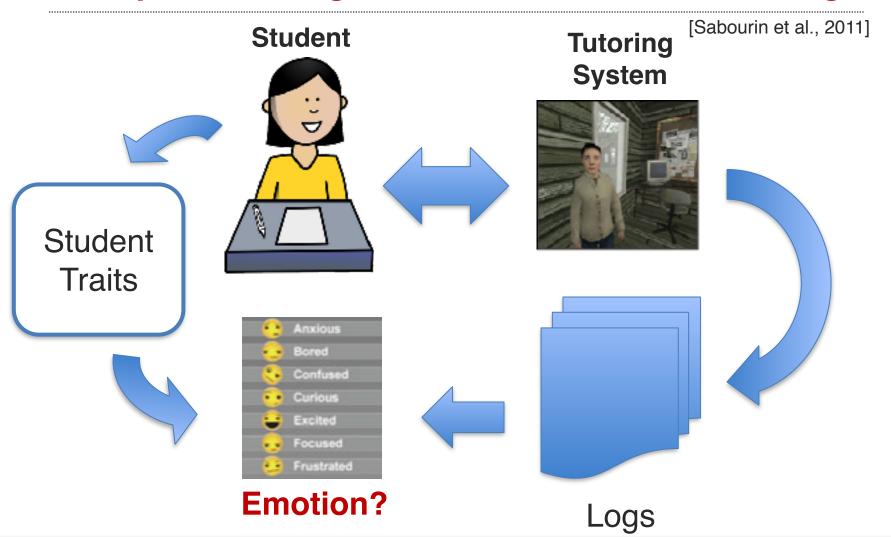
Uncertain Reasoning (Guessing)

- Some aspects of the domain are often unobservable and must be estimated indirectly through other observations.
- The relationships among domain events are often uncertain, particularly the relationship between the observables and non-observables.



Developing a Graphical Model

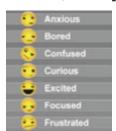
Example: Inferring Emotion from Interaction Logs



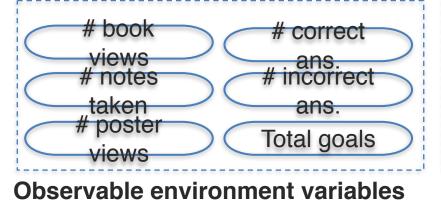
Example: Graphical Model Representation

Emotion

[Sabourin et al., 2011]



Evidences (observable)



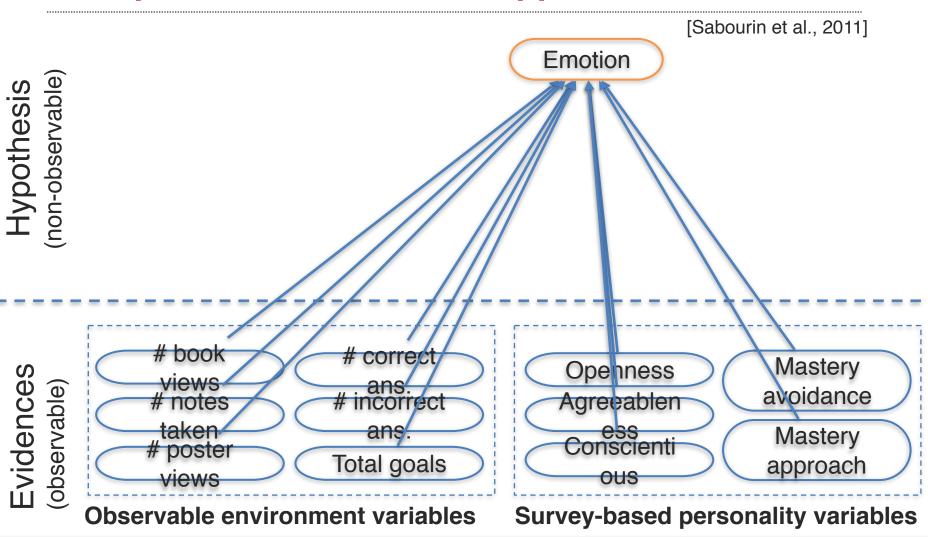
Openness
Agreeablen
ess
Conscientious

Mastery
avoidance

Mastery
approach

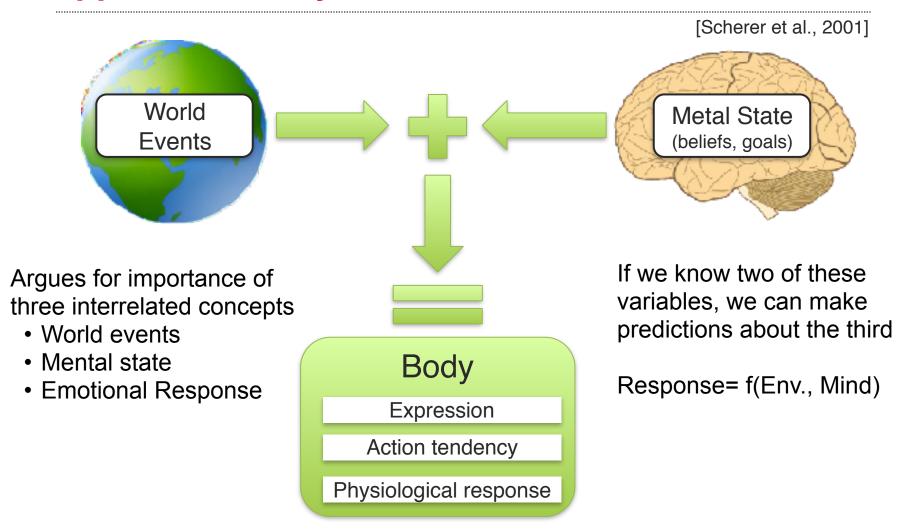
Survey-based personality variables

Example: Direct Prediction Approach

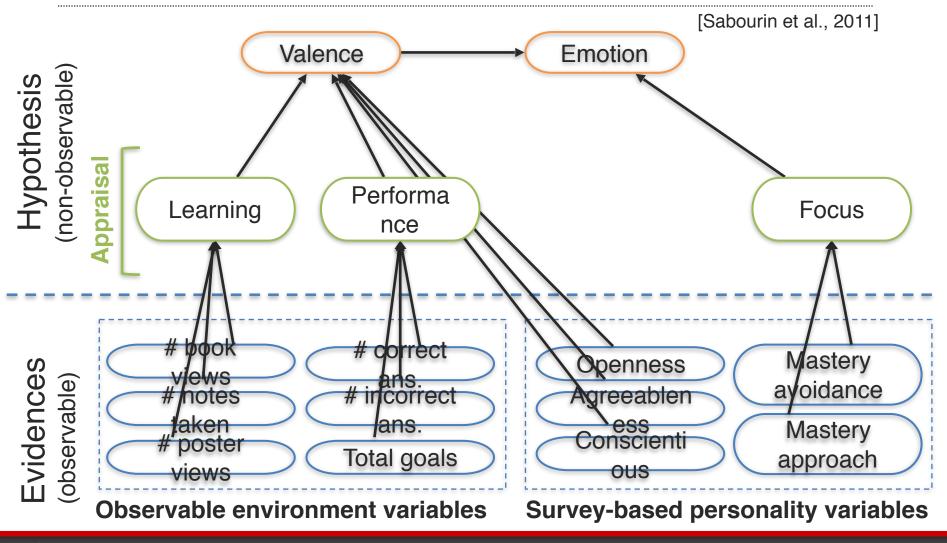




Appraisal Theory of Emotion

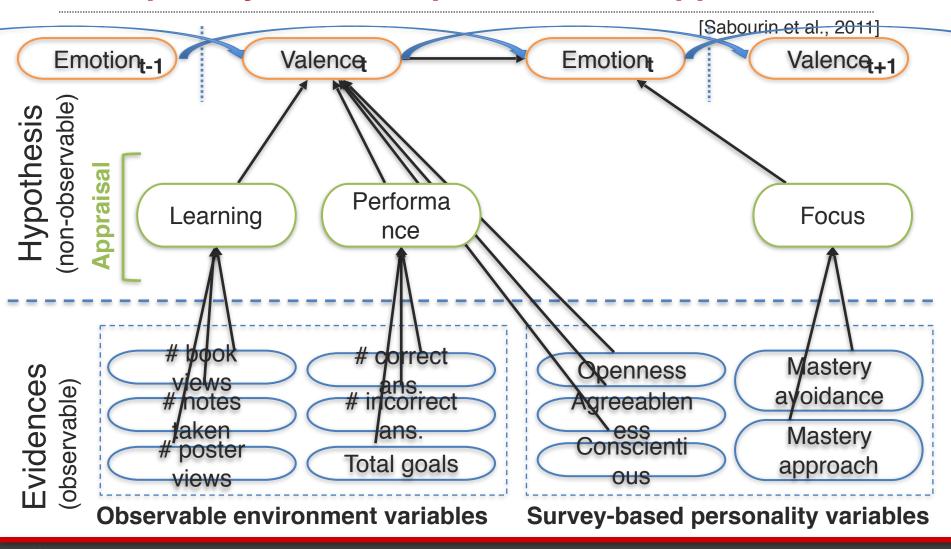


Example: Graphical Model Approach





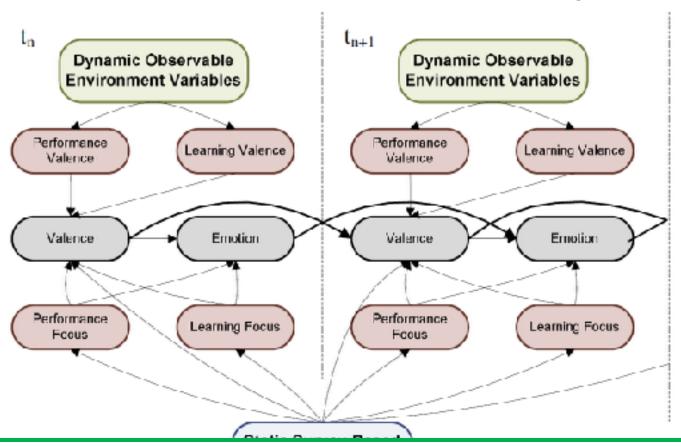
Example: Dynamic Graphical Model Approach





Example: Dynamic Bayesian Network Approach

[Sabourin et al., 2011]



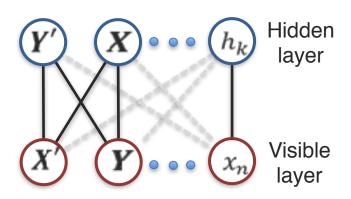
What if the "evidences" require neural network architectures to perform automatic perception?



Markov Random Fields

Restricted Boltzmann Machine (RBM)

- Undirected Graphical Model
- A generative rather than discriminative model
- Connections from every hidden unit to every visible one
- No connections across units (hence Restricted), makes it easier to train and do inference



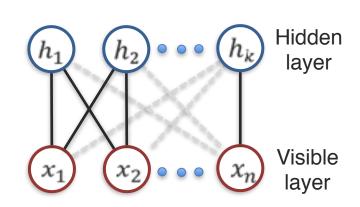
Restricted Boltzmann Machine (RBM)

$$p(\mathbf{x}, \mathbf{h}; \theta) = \frac{\exp(-\mathbf{E}(\mathbf{x}, \mathbf{h}; \theta))}{\sum_{\mathbf{x}'} \sum_{\mathbf{h}'} \exp(-\mathbf{E}(\mathbf{x}', \mathbf{h}'; \theta))} \leftarrow \frac{\text{Partition}}{\text{function } \mathbf{z}}$$

- Hidden and visible layers are binary (e.g. x = {0, ..., 1,0,1})
- Model parameters $\theta = \{W, \boldsymbol{b}, \boldsymbol{a}\}$

$$E = -xWh - bx - ah$$

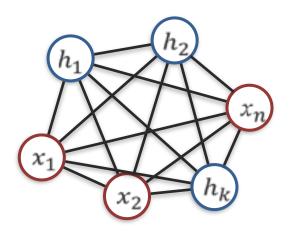
$$E = -\sum_{i} \sum_{j} w_{i,j} x_{i} h_{j} - \sum_{i} b_{i} x_{i} - \sum_{j} a_{j} h_{j}$$
Interaction Bias terms term

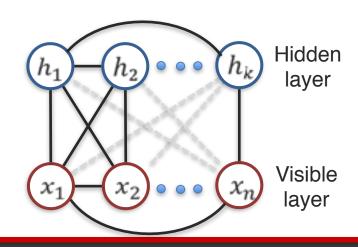


Boltzmann Machine

$$p(\mathbf{x}, \mathbf{h}; \theta) = \frac{\exp(-E(\mathbf{x}, \mathbf{h}; \theta))}{\sum_{\mathbf{x}'} \sum_{\mathbf{h}'} \exp(-E(\mathbf{x}', \mathbf{h}'; \theta))}$$

• Hidden and visible layers are binary (e.g. $x = \{0, ..., 1, 0, 1\}$)



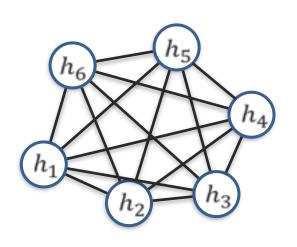


Statistical Mechanics: Boltzmann Distribution

[also called Gibbs measure]

$$p(\mathbf{h}; \theta) = \frac{\exp(-E(\mathbf{h}; \theta)/kT)}{\sum_{\mathbf{h}'} \exp(-E(\mathbf{h}'; \theta)/kT)}$$

probability distribution that gives the probability that a system will be in a certain state h



 $E(h; \theta)$: Energy of state h

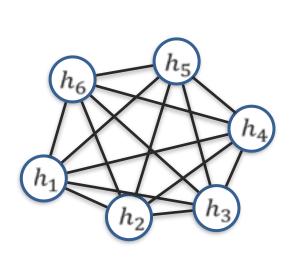
k: Boltzmann constant

T: Thermodynamic temperature

Markov Random Fields

$$p(H = \mathbf{h}; \theta) = \frac{\exp(-E(\mathbf{h}; \theta))}{\sum_{\mathbf{h}'} \exp(-E(\mathbf{h}'; \theta))} = \frac{\Phi(\mathbf{h}; \theta)}{\sum_{\mathbf{h}'} \Phi(\mathbf{h}'; \theta)}$$

Set of random variables H having a Markov property described by undirected graph



$$\Phi(\mathbf{h}; \theta) = \prod_{k} \phi_{k}(\mathbf{h}; \theta_{k})$$

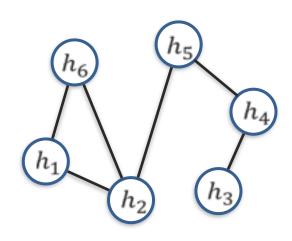
$$= \exp \left(-\sum_{k} E_{k}(\mathbf{h}; \theta_{k}) \right)$$

$$= \exp \left(-\sum_{k} E_{k}(\mathbf{h}; \theta_{k}) \right)$$

Markov Random Fields

$$p(H = \boldsymbol{h}; \theta) = \frac{\Phi(\boldsymbol{h}; \theta)}{\sum_{\boldsymbol{h}'} \Phi(\boldsymbol{h}'; \theta)} = \frac{\sum_{\boldsymbol{k}} \phi_{\boldsymbol{k}}(\boldsymbol{y}, \boldsymbol{x}; \theta)}{\sum_{\boldsymbol{y}'} \sum_{\boldsymbol{k}} \phi_{\boldsymbol{k}}(\boldsymbol{y}', \boldsymbol{x}; \theta)}$$

$$\Phi(h;\theta) = \phi_{12}(h_1, h_2; \theta_{12}) \times \\ \phi_{16}(h_1, h_6; \theta_{16}) \times \\ \phi_{26}(h_2, h_6; \theta_{26}) \times \\ \phi_{25}(h_2, h_5; \theta_{25}) \times \\ \phi_{45}(h_4, h_5; \theta_{45}) \times \\ \phi_{34}(h_3, h_4; \theta_{34})$$



Markov Random Fields: Factor Graphs

$$p(H = \boldsymbol{h}; \theta) = \frac{\Phi(\boldsymbol{h}; \theta)}{\sum_{\boldsymbol{h}'} \Phi(\boldsymbol{h}'; \theta)} = \frac{\sum_{\boldsymbol{k}} \phi_{\boldsymbol{k}}(\boldsymbol{y}, \boldsymbol{x}; \theta)}{\sum_{\boldsymbol{y}'} \sum_{\boldsymbol{k}} \phi_{\boldsymbol{k}}(\boldsymbol{y}', \boldsymbol{x}; \theta)}$$

$$\Phi(\boldsymbol{h};\theta) = \phi_{12}(h_1,h_2;\theta_{12}) \times$$

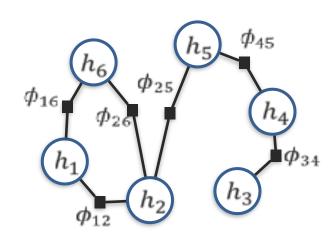
$$\phi_{16}(h_1, h_6; \theta_{16}) \times$$

$$\phi_{26}(h_2, h_6; \theta_{26}) \times$$

$$\phi_{25}(h_2, h_5; \theta_{25}) \times$$

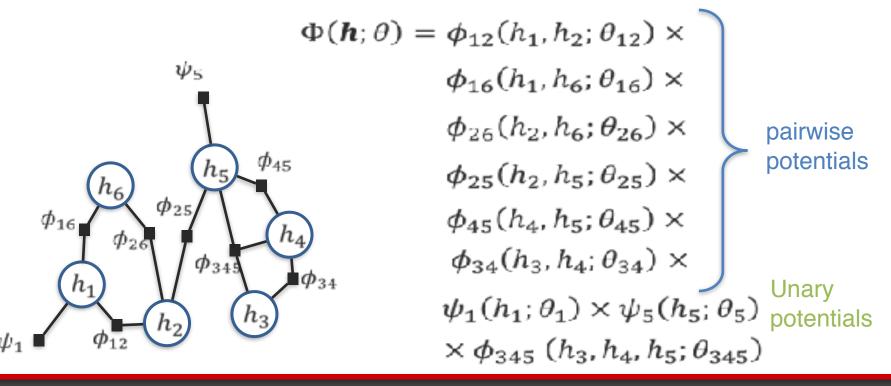
$$\phi_{45}(h_4,h_5;\theta_{45}) \times$$

$$\phi_{34}(h_3, h_4; \theta_{34})$$



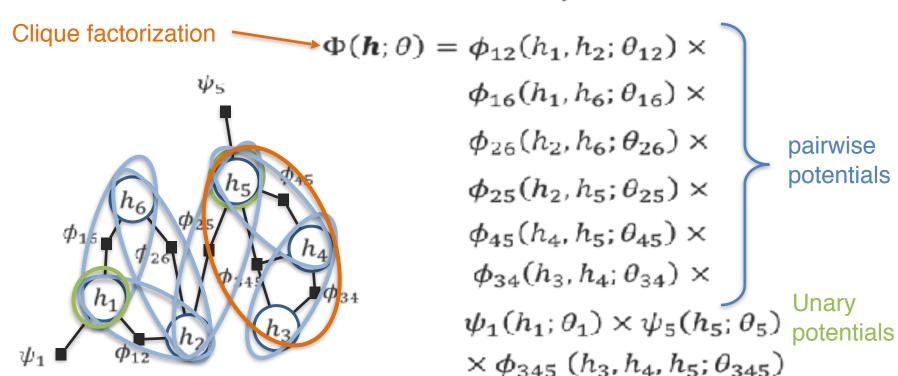
Markov Random Fields (Factor Graphs)

$$p(H = \mathbf{h}; \theta) = \frac{\Phi(\mathbf{h}; \theta)}{\sum_{\mathbf{h}'} \Phi(\mathbf{h}'; \theta)} = \frac{\sum_{k} \phi_{k}(\mathbf{y}, \mathbf{x}; \theta)}{\sum_{\mathbf{y}'} \sum_{k} \phi_{k}(\mathbf{y}', \mathbf{x}; \theta)}$$



Markov Random Fields – Clique Factorization

$$p(H = \mathbf{h}; \theta) = \frac{\Phi(\mathbf{h}; \theta)}{\sum_{\mathbf{h}'} \Phi(\mathbf{h}'; \theta)} = \frac{\sum_{k} \phi_{k}(\mathbf{y}, \mathbf{x}; \theta)}{\sum_{\mathbf{y}'} \sum_{k} \phi_{k}(\mathbf{y}', \mathbf{x}; \theta)}$$



Chain Markov Random Fields (Factor Graphs)

$$p(H = \mathbf{h}; \theta) = \frac{\Phi(\mathbf{h}; \theta)}{\sum_{\mathbf{h}'} \Phi(\mathbf{h}'; \theta)} = \frac{\sum_{k} \phi_{k}(\mathbf{y}, \mathbf{x}; \theta)}{\sum_{\mathbf{y}'} \sum_{k} \phi_{k}(\mathbf{y}', \mathbf{x}; \theta)}$$

$$\Phi(\mathbf{h}; \theta) = \phi_{12}(h_{1}, h_{2}; \theta_{12}) \times \phi_{23}(h_{2}, h_{3}; \theta_{23}) \times \phi_{34}(h_{3}, h_{4}; \theta_{34}) \times \phi_{34}(h_{3}, h_{4}; \theta_{34}) \times \phi_{34}(h_{1}; \theta_{1}) \times \phi_{23}(h_{2}; \theta_{2}) \times \phi_{34}(h_{3}; \theta_{3}) \times \phi_{34}(h_$$

Conditional Random Fields

Conditional Random Fields (Factor Graphs)

$$p(\mathbf{y}|\mathbf{x};\theta) = \frac{\Phi(\mathbf{y},\mathbf{x};\theta)}{\sum_{\mathbf{y}'} \Phi(\mathbf{y}',\mathbf{x};\theta)} = \frac{\sum_{k} \phi_{k}(\mathbf{y},\mathbf{x};\theta)}{\sum_{\mathbf{y}'} \sum_{k} \phi_{k}(\mathbf{y}',\mathbf{x};\theta)}$$

$$\Phi(\mathbf{y},\mathbf{x};\theta) = \phi_{12}(y_{1},y_{2},\mathbf{x};\theta_{12}) \times \phi_{23}(y_{2},y_{3},\mathbf{x};\theta_{23}) \times \phi_{34}(y_{3},y_{4},\mathbf{x};\theta_{34}) \times \phi_{34}(y_{3},y_{4},\mathbf{x};\theta_{34}) \times \phi_{12}(y_{2},\mathbf{x};\theta_{1}) \times \phi_{12}(y_{2},\mathbf{x};\theta_{1}) \times \phi_{12}(y_{2},\mathbf{x};\theta_{2}) \times \phi_{13}(y_{3},\mathbf{x};\theta_{3}) \times \phi_{13$$

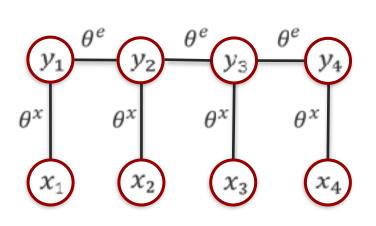
Conditional Random Fields (Factor Graphs)

$$p(\mathbf{y}|\mathbf{x};\theta) = \frac{\Phi(\mathbf{y},\mathbf{x};\theta)}{\sum_{\mathbf{y}'} \Phi(\mathbf{y}',\mathbf{x};\theta)} = \frac{\sum_{k} \phi_{k}(\mathbf{y},\mathbf{x};\theta)}{\sum_{\mathbf{y}'} \sum_{k} \phi_{k}(\mathbf{y}',\mathbf{x};\theta)}$$

$$\Phi(\mathbf{y},\mathbf{x};\theta) = \phi_{12}(y_{1},y_{2},\mathbf{x};\theta_{12}) \times \phi_{23}(y_{2},y_{3},\mathbf{x};\theta_{23}) \times \phi_{34}(y_{3},y_{4},\mathbf{x};\theta_{34}) \times \phi_{34}(y_{3},y_{4},\mathbf{x};\theta_{34}) \times \phi_{12}(y_{2},y_{3},\mathbf{x};\theta_{12}) \times \phi_{12}(y_{2},y_{3},\mathbf{x};\theta_{12}) \times \phi_{12}(y_{2},\mathbf{x}_{2};\theta_{2}) \times \phi_{13}(y_{3},\mathbf{x}_{3};\theta_{3}) \times \phi_{13}(y_{3},\mathbf{x}_{3};\theta_{3}) \times \phi_{13}(y_{3},\mathbf{x}_{3};\theta_{3}) \times \phi_{13}(y_{3},\mathbf{x}_{3};\theta_{3}) \times \phi_{14}(y_{4},\mathbf{x}_{4};\theta_{4})$$
Unary potentials

Conditional Random Fields (Log-linear Model)

$$p(\mathbf{y}|\mathbf{x};\theta) = \frac{\Phi(\mathbf{y},\mathbf{x};\theta)}{\sum_{\mathbf{y}'} \Phi(\mathbf{y}',\mathbf{x};\theta)} = \frac{\sum_{k} \phi_{k}(\mathbf{y},\mathbf{x};\theta)}{\sum_{\mathbf{y}'} \sum_{k} \phi_{k}(\mathbf{y}',\mathbf{x};\theta)}$$
$$= \frac{\exp(\sum_{k} \theta_{k} f_{k}(\mathbf{y},\mathbf{x}))}{\sum_{\mathbf{y}'} \exp(\sum_{k} \theta_{k} f_{k}(\mathbf{y}',\mathbf{x}))}$$



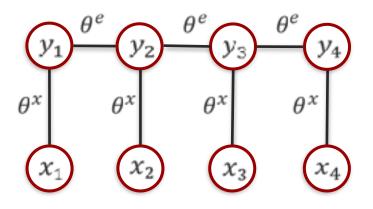
 $f_k(y,x)$: feature function

- Pairwise feature function $f_k(y_i, y_i, x; \theta^e)$
- Unary feature function $f_k(y_i, x; \theta^x)$

Learning Parameters of a CRF Model

$$\underset{\hat{y}}{\operatorname{argmax}} \log (p(\mathbf{y}|\mathbf{x};\theta))$$

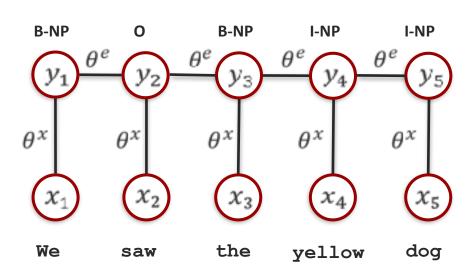
- Gradient can be computed analytically
 - Inference of marginal probabilities using belief propagation (or loopy belief propagation for cyclic graphs)
- Optimized with stochastic or batch approaches



CRFs for Shallow Parsing

$$p(\mathbf{y}|\mathbf{x};\theta) = \frac{\Phi(\mathbf{y},\mathbf{x};\theta)}{\sum_{\mathbf{y}'} \Phi(\mathbf{y}',\mathbf{x};\theta)} = \frac{\exp(\sum_{k} \theta_{k} f_{k}(\mathbf{y},\mathbf{x}))}{\sum_{\mathbf{y}'} \exp(\sum_{k} \theta_{k} f_{k}(\mathbf{y}',\mathbf{x}))}$$

- How many θ^x parameters?
- ➤ What did θ^x learn?



\triangleright	What	did	θ^e	learn?

	B-NP	I-NP	0	
B-NP	$egin{array}{c cccc} \theta_{11} & \theta_{21} & \theta_{31} \\ \theta_{12} & \theta_{22} & \theta_{32} \\ \theta_{13} & \theta_{23} & \theta_{33} \\ \end{array}$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$egin{array}{c cccc} g_{11} & g_{21} & g_{31} \\ g_{12} & g_{22} & g_{32} \\ g_{13} & g_{23} & g_{33} \\ \end{array}$	
I-NP	$ \begin{array}{c cccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$egin{array}{c cccc} g_{11} & g_{21} & g_{31} \\ g_{12} & g_{22} & g_{32} \\ \end{array}$	
0	$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	

Labels:

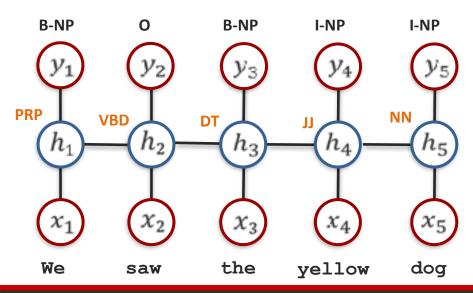
B-NP: Beginning of a noun phrase I-NP: Continuation of a noun phrase

O: Outside a noun phrase

Latent-Dynamic CRF

$$p(\mathbf{y}|\mathbf{x};\theta) = \sum_{\mathbf{h}} p(\mathbf{y}|\mathbf{h};\theta) p(\mathbf{h}|\mathbf{x};\theta)$$

$$= \sum_{\mathbf{h}:\forall h_t \in \mathcal{H}_{\mathbf{y}_t}} p(\mathbf{h}|\mathbf{x};\theta) = \sum_{\mathbf{h}:\forall h_t \in \mathcal{H}_{\mathbf{y}_t}} \frac{\Phi(\mathbf{h},\mathbf{x};\theta)}{\sum_{\mathbf{h}'} \Phi(\mathbf{h}',\mathbf{x};\theta)}$$



Latent variables (e.g., POS tags)

$$h = \{h_1, h_2, h_3, \dots, h_t\}$$
 where $h_t \in \{\mathcal{H}_{y_t}\}$

For example:

$$\mathcal{H} = \{\mathcal{H}_{R-NP}, \mathcal{H}_{I-NP}, \mathcal{H}_{O}\}$$

$$\mathcal{H} = \{B_1, B_2, B_3, B_4, I_1, I_2, I_3, I_4, O_1, O_2, O_3, O_4\}$$

Latent-Dynamic CRF

$$\mathbf{p}(\mathbf{y}|\mathbf{x};\boldsymbol{\theta}) = \sum_{h \in \mathcal{H}_{k} \in \mathcal{H}_{k}} \frac{\exp(\sum_{k} \theta_{k} f_{k}(\mathbf{h}, \mathbf{x}))}{\sum_{\mathbf{h}'} \exp(\sum_{k} \theta_{k} f_{k}(\mathbf{h}', \mathbf{x}))}$$

- How many θ^x parameters?
- ➤ What did θ^x learn?
- **B-NP** 0 **B-NP** I-NP I-NP θ^e θ^e θ^e h_5 θ^{x} θ^x θ^{x} θ^{x} θ^{x} We the yellow dog saw

- How many θ^e parameters?
- What did θ^e learn?
 - Intrinsic dynamics
 - Extrinsic dynamics

Latent variables (e.g., POS tags)

$$h = \{h_1, h_2, h_3, ..., h_t\}$$
 where $h_t \in \{\mathcal{H}_{y_t}\}$

For example:

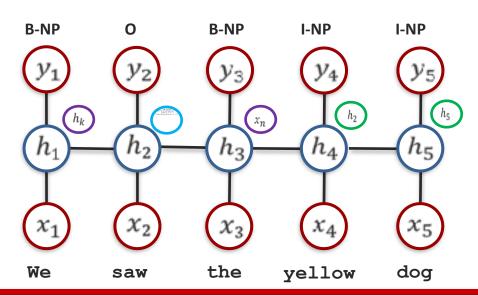
$$\mathcal{H} = \{\mathcal{H}_{B-NP}, \mathcal{H}_{I-NP}, \mathcal{H}_{O}\}$$

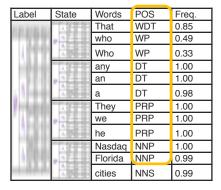
$$\mathcal{H} = \{B_1, B_2, B_3, B_4, I_1, I_2, I_3, I_4, O_1, O_2, O_3, O_4\}$$

Latent-Dynamic CRF for Shallow Parsing

Experiment – Analyzing latent variables

- Task: Shallow parsing with CoNLL 2000 dataset
- Input features: word feature only
- Output labels: Noun phrase labels
 - 1) Select hidden state a^* with highest marginal: $a^* = \arg\max_{a} p(h_t = a|x; \theta)$
- Compute relative frequency for each word





Label	State	Words	POS	Freq.
OR SHEET	0 7	but	CC	0.88
		by	IN	0.73
1,799	3711	or	IN	0.67
1444	4 1	4.6	CD	1.00
300	0.7	1	CD	1.00
1944	3711	11	CD	0.62
14666	B 2 2 2	were	VBD	0.94
1000	0.7	rose	VBD	0.93
3333	2	have	VBP	0.92
733	4 1	been	VBN	0.97
1999	0 %	be	VB	0.94
		to	TO	0.92

Latent variables (e.g., POS tags)

$$h = \{h_1, h_2, h_3, \dots, h_t\}$$
 where $h_t \in \{\mathcal{H}_{y_t}\}$

For example:

$$\mathcal{H} = \{\mathcal{H}_{B-NP}, \mathcal{H}_{I-NP}, \mathcal{H}_{O}\}$$

$$\mathcal{H} = \{B_1, B_2, B_3, B_4, I_1, I_2, I_3, I_4, O_1, O_2, O_3, O_4\}$$

Hidden Conditional Random Field

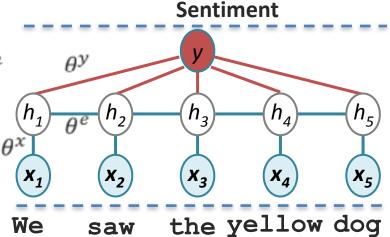
Sequence label:

k Boltzmann constant $\mathbf{E}(h; \theta)$: Energy of state h

Latent variables with shared hidden states:

$$\mathbf{h} = \{h_1, h_2, h_3, \dots, h_t\}$$

$$\Phi(\boldsymbol{h};\theta) = \prod_{k} \phi_{k}(\boldsymbol{h};\theta_{k})$$



 $p(H=h;\theta) = \frac{\exp(-\mathbb{E}(h;\theta))}{\sum_{h'} \exp(-\mathbb{E}(h';\theta))}$

T: Thermodynamic to

- Inference is tractable: O(YH²T)
 - Linear in sequence length T!
- Parameter learning $(\theta^x, \theta^e, \theta^y)$:
 - Gradient descent or L-BFGS

Shared hidden states



Learning Multimodal Structure

Modality-*private* structure

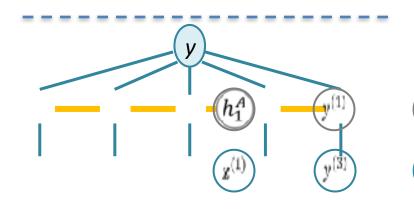
Internal grouping of observations

Modality-*shared* structure

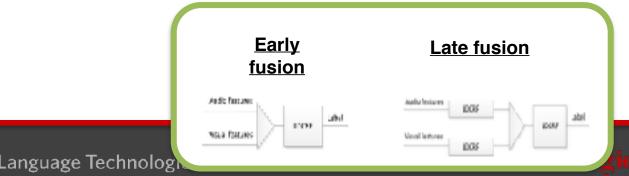
Interaction and synchrony

Early / Late fusion is inappropriate

- Early strong modality can dominate
- Late cross-modality dependency is discarded







Multi-view Latent Variable Discriminative Models

Modality-*private* structure

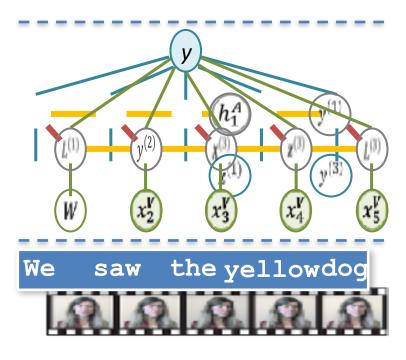
Internal grouping of observations

Modality-*shared* structure

Interaction and synchrony

Early / Late fusion is inappropriate

- Early strong modality can dominate
- Late cross-modality dependency is discarded



$$p(y|\mathbf{x}^A, \mathbf{x}^V; \boldsymbol{\theta}) = \sum_{\mathbf{h}^A, \mathbf{h}^V} p(y, \mathbf{h}^A, \mathbf{h}^V | \mathbf{x}^A, \mathbf{x}^V; \boldsymbol{\theta})$$

Approximate inference using loopy-belief

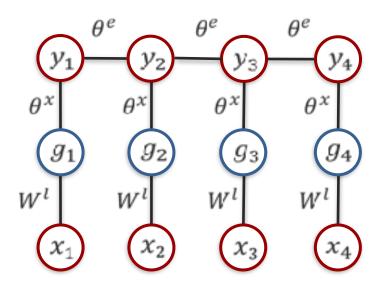
CRFs and Deep Learning

Conditional Neural Fields

$$\boldsymbol{\mathcal{G}^l}\big(\boldsymbol{x_i}, \boldsymbol{W^l}\big) = \big[g_1^l\big(\boldsymbol{x_i} \cdot \boldsymbol{W_1^l}\big), g_2^l\big(\boldsymbol{x_t} \cdot \boldsymbol{W_i^l}\big), \dots, g_n^l(\boldsymbol{x_i} \cdot \boldsymbol{W_n^l})\big]$$

$$p(\mathbf{y} \mid \mathbf{x}; \boldsymbol{\theta}) \propto \exp \left\{ \sum_{i} \boldsymbol{\theta}^{x} \cdot \boldsymbol{f}^{x}(y_{i}, \mathbf{x}_{i}) + \sum_{i} \boldsymbol{\theta}^{e} \cdot \boldsymbol{f}^{e}(y_{i}, y_{i-1}) \right\}$$

$$f^{x}(y_{i}, \mathbf{x}_{i}) = \mathbb{I}[y_{i} = y'] \cdot \mathcal{G}(\mathbf{x}_{i}, \mathbf{W}^{l})$$



Deep Conditional Neural Fields

$$\begin{aligned}
\mathbf{g}^{l}(\mathbf{x}_{i}, \mathbf{W}^{l}) &= \left[g_{1}^{l}(\mathbf{x}_{i} \cdot \mathbf{W}_{1}^{l}), g_{2}^{l}(\mathbf{x}_{i} \cdot \mathbf{W}_{i}^{l}), \dots, g_{n}^{l}(\mathbf{x}_{i} \cdot \mathbf{W}_{n}^{l})\right] \\
p(\mathbf{y} \mid \mathbf{x}; \boldsymbol{\theta}) &\propto \exp\left\{\sum_{i} \boldsymbol{\theta}^{x} \cdot f^{x}(y_{i}, \mathbf{x}_{i}) + \sum_{i} \boldsymbol{\theta}^{e} \cdot f^{e}(y_{i}, y_{i-1})\right\} \\
\underbrace{\begin{pmatrix} y_{1} & \theta^{e} & \theta^{e} & \theta^{e} \\ y_{2} & y_{3} & \theta^{x} \\ \theta^{x} & \theta^{x} & \theta^{x} \\ \theta^{x} & \theta^{x} & \theta^{x} \\ g_{1}^{2} & g_{2}^{2} & g_{3}^{2} & g_{4}^{2} \\ y_{2} & y_{3} & g_{4}^{2} \\ y_{3} & \theta^{x} & \theta^{x} \\ g_{1}^{2} & g_{2}^{2} & g_{3}^{2} & g_{4}^{2} \\ y_{4} & \alpha^{l} &= \mathcal{G}(\boldsymbol{a}_{i}^{l-1}, \boldsymbol{\theta}^{g}) & \text{for } l = 2 \dots m-1 \end{aligned} \right\}$$

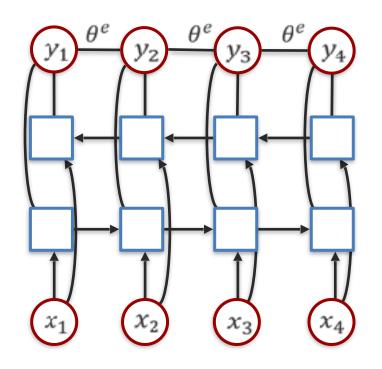
$$\mathbf{g}^{l}(\mathbf{x}_{i} \cdot \mathbf{W}^{l}) = \mathbf{g}^{l}(\mathbf{y}_{i}, \mathbf{x}_{i}) = \mathbf{g}^{l}(\mathbf{y}_{i}, \mathbf{y}_{i}) = \mathbf{g}^{l}(\mathbf{y}_{i},$$

CRF and Bilinear LSTM

[Dyer, 2016]

Learning:

- 1. Feedforward
- 2. Gradient
 - a) Belief propagation
- 3.Backpropagation



Output labels:

Name entities

Input features:

Word embedding

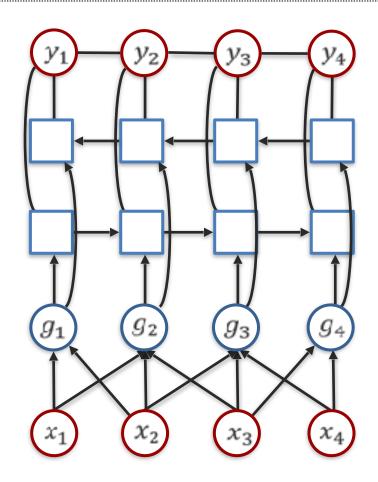
- What did θ^e paramters learn?
- ➤ What does LSTM parameters learns?

CNN and CRF and Bilinear LSTM

[Hovy, 2016]

Learning:

- 1. Feedforward
- 2. Gradient
 - a) Belief propagation
- 3.Backpropagation



Output labels:

Name entities

Input features:

 Character embedding

Continuous and Fully-Connected CRFs

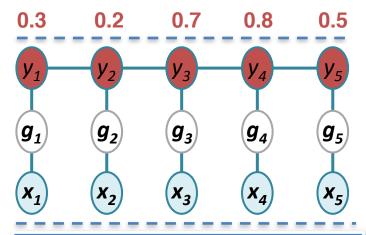
Continuous Conditional Neural Field

[Baltrusaitis 2014]

Continuous output variables: (e.g., continuous emotional label)

$$\psi_1(y_1, \mathbf{x}; \theta_1) \times \mathbf{y}_1$$

$$p(\mathbf{y}|\ \mathbf{x};\boldsymbol{\theta}) = \frac{1}{\mathcal{Z}(\mathbf{x};\boldsymbol{\theta})} \exp \left\{ \sum_t \boldsymbol{\theta} \cdot F(y_t, y_{t-1}, \mathbf{x}_t, \boldsymbol{\theta}^g) \right\}$$



We saw the yellowdog

> How to solve

Multivariate Gaussian integral:

$$\int_{-\infty}^{\infty} \exp\left\{\frac{1}{2} \mathbf{y}^T \Sigma^{-1} \mathbf{y} + \mathbf{y} \Sigma^{-1} \boldsymbol{\mu}\right\} d\mathbf{y}$$
$$= \frac{(2\pi)^{n/2}}{|\Sigma^{-1}|^{1/2}} \exp\left(\frac{1}{2} \boldsymbol{\mu} \Sigma^{-1} \boldsymbol{\mu}\right)$$

[Radosavljevic et al., 2010]

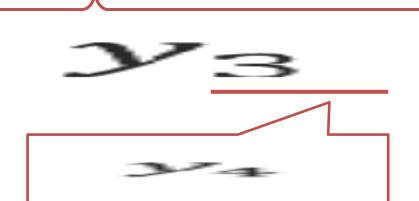
negie Mellon University

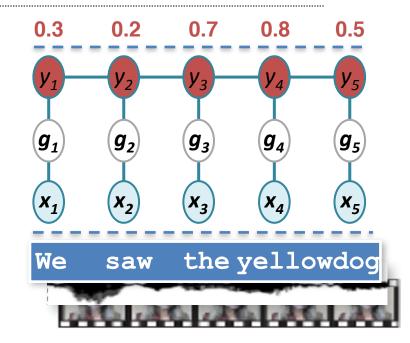
Continuous Conditional Neural Field

Continuous output variables: (e.g., continuous emotional label)

$$\psi_1(y_1, \mathbf{x}; \theta_1) \times \mathbf{y}_1$$

$$p(\mathbf{y}|\mathbf{x};\boldsymbol{\theta}) = \frac{1}{Z(\mathbf{x};\boldsymbol{\theta})} \exp \left\{ \sum_{t} \boldsymbol{\theta} \cdot F(y_{t}, y_{t-1}, \mathbf{x}_{t}, \boldsymbol{\theta}^{g}) \right\}$$





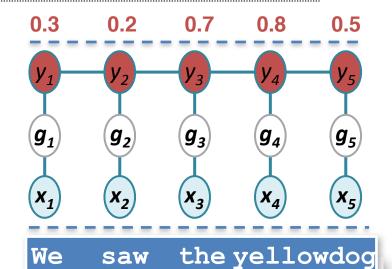
Continuous Conditional Neural Field

Continuous output variables: (e.g., continuous emotional label)

$$\psi_1(y_1, \mathbf{x}; \theta_1) \times \mathbf{y}_1$$

Multivariate Gaussian distribution:

$$\psi_3(y_3, x; \theta_3) \times$$



 $\psi_4(\gamma$

 $\phi_{23}(y_2, y_3)$

Since CCNF can be viewed as a multivariate Gaussian, the prediction of y' is simply the mean value of distribution:

$$y' = \arg\max_{y} (P(y|x)) = \mu$$

Optimized using gradient ascent or BFGS.

High-Order Continuous Conditional Neural Field

Continuous output variables: (e.g., continuous emotional label)

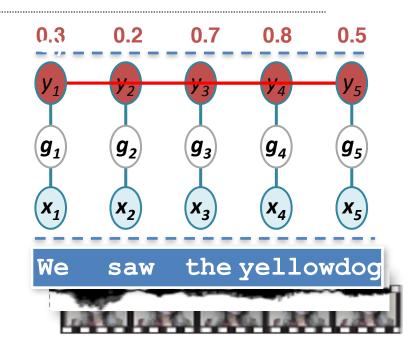
$$\psi_1(y_1, \mathbf{x}; \theta_1) \times \mathbf{y}_1$$

Multivariate Gaussian distribution:

$$\psi_3(y_3, x; \theta_3) \times$$

k-order potential functions:

$$\psi_1(y_1,x_1;\theta_1) \times$$



Fully-Connected Continuous Conditional Neural Field

Continuous output variables: (e.g., continuous emotional label)

$$\psi_1(y_1, \mathbf{x}; \theta_1) \times \mathbf{y}_1$$

Multivariate Gaussian distribution:

$$\psi_3(y_3, x; \theta_3) \times$$

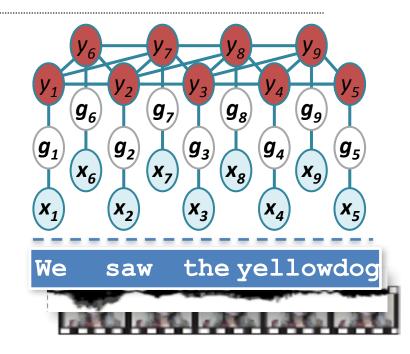
k-order potential functions:

$$\psi_1(y_1,x_1;\theta_1) \times$$

Grid potential functions:

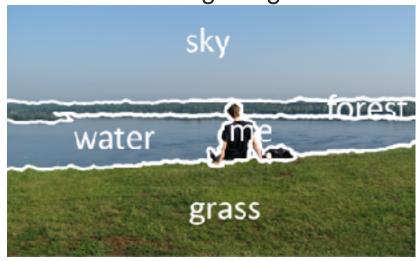
$$f^{2D}(y_i, y_j) = -\frac{1}{2} S_{ij} (y_i - y_j)^2$$

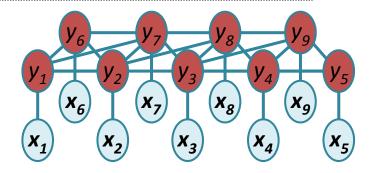
where $S_{t,t}$ specifies which nodes are connected.



Fully-Connected CRF_[Krahenbuhl and Koltun, 2013]

"Semantic" image segmentation





yi: object class label

x_i: local pixel features

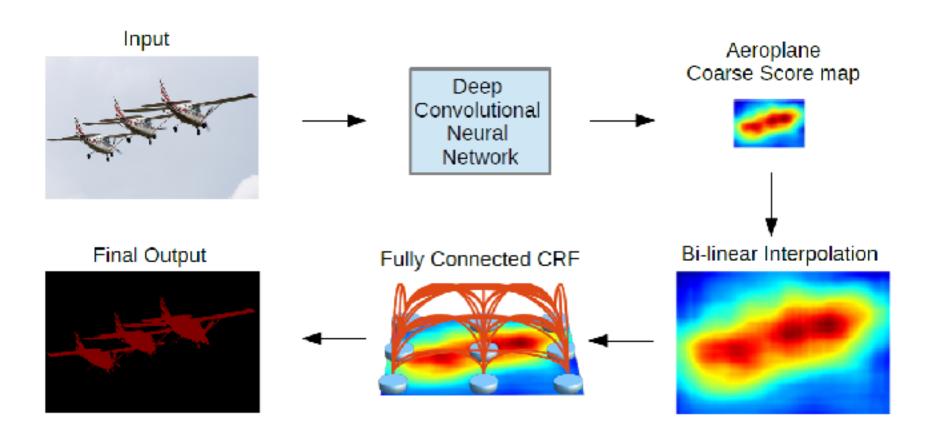
$$p(\mathbf{y}|\mathbf{x};\theta) = \frac{\Phi(\mathbf{y},;\theta)}{\sum_{\mathbf{y}'} \Phi(\mathbf{y}',\mathbf{x};\theta)}$$

where
$$\Phi_{ij}(y_i, y_j; \boldsymbol{\theta}) = \sum_{m=1}^{m}$$

$$p(\mathbf{y}|\mathbf{x};\theta) = \frac{\Phi(\mathbf{y},;\theta)}{\sum_{\mathbf{y}'} \Phi(\mathbf{y}',\mathbf{x};\theta)}$$
 Mixture of kernels where
$$\Phi_{ij}(y_i,y_j;\theta) = \sum_{m=1}^{C} u^{(m)}(y_i,y_j|\theta)k^{(m)}(\mathbf{x}_i,\mathbf{x}_j)$$

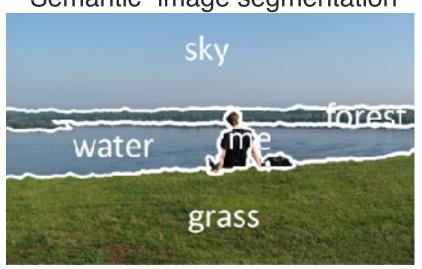
CNN and Fully-Connected CRF

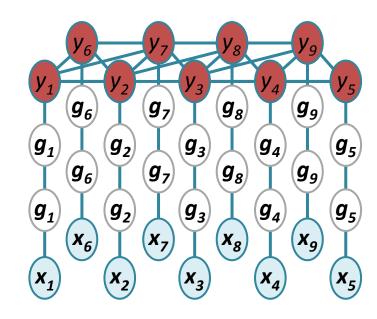
[Chen et al., 2014]



Fully Connected Deep Structured Networks [Zheng et al., 2015; Schwing and Urtasun, 2015]

"Semantic" image segmentation





Algorithm: Learning Fully Connected Deep Structured Models Repeat until stopping criteria

- 1. Forward pass to compute $f_r(x, \hat{y}_r; w) \ \forall r \in \mathcal{R}, y_r \in \mathcal{Y}_r$
- 2. Computation of marginals $q_{(x,n),i}^t(\hat{y}_i)$ via filtering for $t \in \{1,\ldots,T\}$
- approximation 3. Backtracking through the marginals $q_{(x,y),i}^t(\tilde{y}_i)$ from t=T-1 down to t=1
- 4. Backward pass through definition of function via chain rule
- Parameter update



Using mean field