Language Technologies Institute

Carnegie Mellon University

# Advanced Multimodal Machine Learning

## Lecture 9.1: Probabilistic Graphical Models

**Louis-Philippe Morency**

* Original version co-developed with 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

# Upcoming Schedule

- First project assignment:
    - Proposal presentations (10/2 and 10/4)
    - First project reports (10/7)

- **Midterm project assignment**
    - **Midterm presentations**
        - Tuesday 11/6 & Thursday 11/8 (DH A302)
    - **Midterm reports (Sunday 11/11)**

- Final project assignment
    - Final presentations (12/3 & 12/4)
    - Final reports (12/11)

# Lecture Schedule

| Classes | Lectures |
| --- | --- |
| **Week 9**<br>10/23 – 10/25 | **Probabilistic graphical models**<br>• Boltzmann distribution and CRFs<br>• Continuous and fully-connected CRFs |
| **Week 10**<br>10/30 & 11/1 | **Multimodal optimization**<br>• Variational Auto-encoder<br>• Generative-Adversarial Networks |
| **Week 11**<br>11/6 & 11/8 | ***Mid-term project assignment - Pre*** |
| **Week 12**<br>11/13 & 11/15 | **Multimodal fusion and new directions**<br>• Multi-kernel learning and fusion<br>• New directions in multimodal machine learning |

Thursday in DH A302.
Midterm due on 11/11.

Language Technologies Institute

Carnegie Mellon University

# Lecture Schedule

| Classes | Lectures |
|---|---|
| **Week 13** <br> 11/20 & 11/22 | ***Thanksgiving week (+ Project preparation)*** |
| **Week 14** <br> 11/27 & 11/29 | **Multi-lingual representations and** <br> ● Neural machine translation <br> ● Guest lecture: Graham Neubig |
| **Week 15** <br> 12/3 & 12/4 <br> * Final * | ***Final project assignment - Present*** |

Lecture on Thursday.
Discussions on Tuesday.

Poster presentations in GHC 6121.
Final project due: 12/11.

Language Technologies Institute

Carnegie Mellon University

# Quick Recap

# Multimodal Representation Learning

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

❑ Deep Multimodal Boltzmann machines

softmax

Text
$X$

Image
$Y$

# Multimodal Representation Learning

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

- ❑ Deep Multimodal Boltzmann machines
- ❑ Stacked Autoencoder

Language Technologies Institute

Carnegie Mellon University

# Multimodal Representation Learning

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

❑ Deep Multimodal Boltzmann machines

❑ Stacked Autoencoder

❑ Encoder-Decoder



Text
*X*

Image
*Y*

# Multimodal Representation Learning

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

- ❑ Deep Multimodal Boltzmann machines

- ❑ Stacked Autoencoder

- ❑ Encoder-Decoder

- ❑ "Minimum-distance" Multimodal Embedding

Language Technologies Institute

Carnegie Mellon University

# Recurrent Neural Network using LSTM Units



How can we improve reasoning by including prior domain knowledge?

Language Technologies Institute

**Carnegie Mellon University**

# 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: $X_1,\ldots,X_n$
- P is a joint distribution over $X_1,\ldots,X_n$

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 $X_1,\ldots,X_n$ are independent then:
  - $P(X_1,\ldots,X_n) = P(X_1)\ldots P(X_n)$

$$X \qquad Y$$

# 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 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.

Non-observables   Observables

# Developing a Graphical Model

Carnegie Mellon University

# Example: Inferring Emotion from Interaction Logs

[Sabourin et al., 2011]

**Student**

**Tutoring System**

Student Traits

Anxious
Bored
Confused
Curious
Excited
Focused
Frustrated

**Emotion?**

Logs

# Example: Graphical Model Representation

[Sabourin et al., 2011]

**Hypothesis (non-observable)**

Emotion

- Anxious
- Bored
- Confused
- Curious
- Excited
- Focused
- Frustrated

**Evidences (observable)**

# book views
# correct ans.
# notes taken
# incorrect ans.
# poster views
Total goals

**Observable environment variables**

Openness
Mastery avoidance
Agreeableness
Conscientious
Mastery approach

**Survey-based personality variables**

# Example: Direct Prediction Approach

[Sabourin et al., 2011]

Hypothesis (non-observable)

Evidences (observable)

Emotion

# book views
# notes taken
# poster views

# correct ans.
# incorrect ans.
Total goals

Openness
Agreeableness
Conscientious

Mastery avoidance
Mastery approach

**Observable environment variables**

**Survey-based personality variables**

Language Technologies Institute

Carnegie Mellon University

# Appraisal Theory of Emotion

[Scherer et al., 2001]

World Events

+

Metal State
(beliefs, goals)

=

Body

Expression

Action tendency

Physiological response

Argues for importance of three interrelated concepts
- World events
- Mental state
- Emotional Response

If we know two of these variables, we can make predictions about the third

Response= f(Env., Mind)

Language Technologies Institute

Carnegie Mellon University

# Example: Graphical Model Approach

[Sabourin et al., 2011]



**Observable environment variables**     **Survey-based personality variables**

# Example: Dynamic Graphical Model Approach



[Sabourin et al., 2011]

# Example: Dynamic Bayesian Network Approach

[Sabourin et al., 2011]



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

# Markov Random Fields

Language Technologies Institute

Carnegie Mellon University

# 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

$h_1$ $h_2$ $\bullet \bullet \bullet$ $h_k$    Hidden layer

$x_1$ $x_2$ $\bullet \bullet \bullet$ $x_n$    Visible layer

# Restricted Boltzmann Machine (RBM)

$$p(\boldsymbol{x}, \boldsymbol{h}; \theta) = \frac{\exp(-\mathrm{E}(\boldsymbol{x}, \boldsymbol{h}; \theta))}{\sum_{\boldsymbol{x'}} \sum_{\boldsymbol{h'}} \exp(-\mathrm{E}(\boldsymbol{x'}, \boldsymbol{h'}; \theta))}$$ ← **Partition function $Z$**

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

- Model parameters $\theta = \{W, \boldsymbol{b}, \boldsymbol{a}\}$

$$\mathrm{E} = -\boldsymbol{x}W\boldsymbol{h} - \boldsymbol{b}\boldsymbol{x} - \boldsymbol{a}\boldsymbol{h}$$

$$\mathrm{E} = -\underbrace{\sum_i \sum_j w_{i,j} x_i h_j}_{\text{Interaction term}} - \underbrace{\sum_i b_i x_i - \sum_j a_j h_j}_{\text{Bias terms}}$$

Hidden layer: $h_1$, $h_2$, $\dots$, $h_k$

Visible layer: $x_1$, $x_2$, $\dots$, $x_n$

Language Technologies Institute

Carnegie Mellon University

# Boltzmann Machine

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

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

Carnegie Mellon University

# Statistical Mechanics: Boltzmann Distribution

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

➢ probability distribution that gives the probability that a system will be in a certain state $\boldsymbol{h}$

$\mathrm{E}(\boldsymbol{h}; \theta)$: Energy of state $\boldsymbol{h}$

$k$: Boltzmann constant

$T$: Thermodynamic temperature

Language Technologies Institute

Carnegie Mellon University

# Markov Random Fields

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

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



**Potential functions** $\phi_k(\boldsymbol{h}; \theta) > 0$

$$\Phi(\boldsymbol{h}; \theta) = \prod_k \phi_k(\boldsymbol{h}; \theta_k)$$

$$= \exp\left(-\sum_k \mathrm{E}_k(\boldsymbol{h}; \theta_k)\right)$$

Language Technologies Institute

Carnegie Mellon University

# Markov Random Fields

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

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

Language Technologies Institute

Carnegie Mellon University

# Markov Random Fields: Factor Graphs

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

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

Language Technologies Institute

Carnegie Mellon University

# Markov Random Fields (Factor Graphs)

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

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

pairwise potentials

$$\psi_1(h_1; \theta_1) \times \psi_5(h_5; \theta_5)$$

Unary potentials

$$\times \phi_{345}(h_3, h_4, h_5; \theta_{345})$$

Language Technologies Institute

Carnegie Mellon University

# Markov Random Fields – Clique Factorization

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

Clique factorization

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

pairwise potentials

$$\psi_1(h_1; \theta_1) \times \psi_5(h_5; \theta_5)$$

Unary potentials

$$\times \phi_{345}(h_3, h_4, h_5; \theta_{345})$$

$\psi_5$ $\phi_{45}$ $h_5$ $h_6$ $\phi_{25}$ $\phi_{16}$ $h_4$ $\phi_{26}$ $\phi_{345}$ $\phi_{34}$ $h_1$ $h_2$ $h_3$ $\psi_1$ $\phi_{12}$

Language Technologies Institute

Carnegie Mellon University

# Chain Markov Random Fields (Factor Graphs)

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

$$\Phi(\boldsymbol{h}; \theta) = \phi_{12}(h_1, h_2; \theta_{12}) \times$$
$$\phi_{23}(h_2, h_3; \theta_{23}) \times$$
$$\phi_{34}(h_3, h_4; \theta_{34}) \times$$

pairwise potentials

$$\psi_1(h_1; \theta_1) \times$$
$$\psi_2(h_2; \theta_2) \times$$
$$\psi_3(h_3; \theta_3) \times$$
$$\psi_4(h_4; \theta_4)$$

Unary potentials

Language Technologies Institute

Carnegie Mellon University

# Conditional Random Fields

# Conditional Random Fields (Factor Graphs)

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

$$\Phi(\boldsymbol{y},\boldsymbol{x};\theta) = \phi_{12}(y_1,y_2,\boldsymbol{x};\theta_{12}) \times$$
$$\phi_{23}(y_2,y_3,\boldsymbol{x};\theta_{23}) \times$$
$$\phi_{34}(y_3,y_4,\boldsymbol{x};\theta_{34}) \times$$

pairwise potentials

$$\psi_1(y_1,\boldsymbol{x};\theta_1) \times$$
$$\psi_2(y_2,\boldsymbol{x};\theta_2) \times$$
$$\psi_3(y_3,\boldsymbol{x};\theta_3) \times$$
$$\psi_4(y_4,\boldsymbol{x};\theta_4)$$

Unary potentials

# Conditional Random Fields (Factor Graphs)

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

$$\Phi(\boldsymbol{y},\boldsymbol{x};\theta) = \phi_{12}(y_1,y_2,\boldsymbol{x};\theta_{12}) \times$$
$$\phi_{23}(y_2,y_3,\boldsymbol{x};\theta_{23}) \times$$
$$\phi_{34}(y_3,y_4,\boldsymbol{x};\theta_{34}) \times$$

pairwise potentials

$$\psi_1(y_1,x_1;\theta_1) \times$$
$$\psi_2(y_2,x_2;\theta_2) \times$$
$$\psi_3(y_3,x_3;\theta_3) \times$$
$$\psi_4(y_4,x_4;\theta_4)$$

Unary potentials

Language Technologies Institute

Carnegie Mellon University

# Conditional Random Fields (Log-linear Model)

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

$$= \frac{\exp(\sum_k \theta_k f_k(\boldsymbol{y},\boldsymbol{x}))}{\sum_{\boldsymbol{y}'}\exp(\sum_k \theta_k f_k(\boldsymbol{y}',\boldsymbol{x}))}$$

$f_k(\boldsymbol{y},\boldsymbol{x})$: feature function

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

Carnegie Mellon University

# Learning Parameters of a CRF Model

$$\underset{\hat{y}}{\mathrm{argmax}} \log\big(p(\boldsymbol{y}|\boldsymbol{x};\theta)\big)$$

- 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

$$p(\boldsymbol{y}|\boldsymbol{x}; \theta) = \frac{\Phi(\boldsymbol{y}, \boldsymbol{x}; \theta)}{\sum_{\boldsymbol{y'}} \Phi(\boldsymbol{y'}, \boldsymbol{x}; \theta)} = \frac{\exp(\sum_k \theta_k f_k(\boldsymbol{y}, \boldsymbol{x}))}{\sum_{\boldsymbol{y'}} \exp(\sum_k \theta_k f_k(\boldsymbol{y'}, \boldsymbol{x}))}$$

➤ How many $\theta^x$ parameters?

➤ What did $\theta^x$ learn?

➤ What did $\theta^e$ learn?

|  | B-NP | I-NP | O |
|---|---|---|---|
| B-NP | $\theta_{11}$ | $\theta_{21}$ | $\theta_{31}$ |
| I-NP | $\theta_{12}$ | $\theta_{22}$ | $\theta_{32}$ |
| O | $\theta_{13}$ | $\theta_{23}$ | $\theta_{33}$ |

**Labels:**

**B-NP: Beginning of a noun phrase**

**I-NP: Continuation of a noun phrase**

**O: Outside a noun phrase**

**Dictionary size: 10,000 words**

# Latent-Dynamic CRF

$$p(\boldsymbol{y}|\boldsymbol{x};\theta) = \sum_{\boldsymbol{h}} p(\boldsymbol{y}|\boldsymbol{h};\boldsymbol{\theta})p(\boldsymbol{h}|\boldsymbol{x};\boldsymbol{\theta}) \quad \text{where} \quad p(\boldsymbol{y}|\boldsymbol{h};\boldsymbol{\theta}) = \begin{cases} 1 & \text{if } \forall h_t \in \mathcal{H}_{y_t} \\ 0 & \text{otherwise} \end{cases}$$

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



**Latent variables**   (e.g., POS tags)

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

**For example:**

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

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

**Dictionary size: 10,000 words**

Language Technologies Institute

Carnegie Mellon University

# Latent-Dynamic CRF

$$p(\boldsymbol{y}|\boldsymbol{x};\theta) = \sum_{\boldsymbol{h}:\forall h_t \in \mathcal{H}_{y_t}} \frac{\exp(\sum_k \theta_k f_k(\boldsymbol{h}, \boldsymbol{x}))}{\sum_{\boldsymbol{h}'} \exp(\sum_k \theta_k f_k(\boldsymbol{h}', \boldsymbol{x}))}$$

➢ How many $\theta^x$ parameters?

➢ What did $\theta^x$ learn?

➢ How many $\theta^e$ parameters?

➢ What did $\theta^e$ learn?

- Intrinsic dynamics
- Extrinsic dynamics

**Latent variables**  (e.g., POS tags)

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

**For example:**

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

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

**Dictionary size: 10,000 words**



B-NP — $y_1$ — $\theta^e$ — $h_1$ — $\theta^x$ — $x_1$ — **We**

O — $y_2$ — $\theta^e$ — $h_2$ — $\theta^x$ — $x_2$ — **saw**

B-NP — $y_3$ — $\theta^e$ — $h_3$ — $\theta^x$ — $x_3$ — **the**

I-NP — $y_4$ — $\theta^e$ — $h_4$ — $\theta^x$ — $x_4$ — **yellow**

I-NP — $y_5$ — $h_5$ — $\theta^x$ — $x_5$ — **dog**

Language Technologies Institute

Carnegie Mellon University

# 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:
$$a^* = \arg \max_a p(\boldsymbol{h_t} = a | \boldsymbol{x}; \boldsymbol{\theta})$$

2) Compute relative frequency for each word

| Label | State | Words | POS | Freq. |
|---|---|---|---|---|
| *B* | $B_1$ | That | WDT | 0.85 |
| | | who | WP | 0.49 |
| | | Who | WP | 0.33 |
| | $B_2$ | any | DT | 1.00 |
| | | an | DT | 1.00 |
| | | a | DT | 0.98 |
| | $B_3$ | They | PRP | 1.00 |
| | | we | PRP | 1.00 |
| | | he | PRP | 1.00 |
| | $B_4$ | Nasdaq | NNP | 1.00 |
| | | Florida | NNP | 0.99 |
| | | cities | NNS | 0.99 |

| Label | State | Words | POS | Freq. |
|---|---|---|---|---|
| *O* | $O_1$ | but | CC | 0.88 |
| | | by | IN | 0.73 |
| | | or | IN | 0.67 |
| | $O_2$ | 4.6 | CD | 1.00 |
| | | 1 | CD | 1.00 |
| | | 1 1 | CD | 0.62 |
| | $O_3$ | were | VBD | 0.94 |
| | | rose | VBD | 0.93 |
| | | have | VBP | 0.92 |
| | $O_4$ | been | VBN | 0.97 |
| | | be | VB | 0.94 |
| | | to | TO | 0.92 |

B-NP    O    B-NP    I-NP    I-NP

$y_1$   $y_2$   $y_3$   $y_4$   $y_5$

$B_3$   $O_3$   $B_2$   $I_2$   $I_3$

$h_1$   $h_2$   $h_3$   $h_4$   $h_5$

$x_1$   $x_2$   $x_3$   $x_4$   $x_5$

We    saw    the    yellow    dog

**Latent variables**   (e.g., POS tags)

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

**For example:**

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

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

**Dictionary size: 10,000 words**

Language Technologies Institute

Carnegie Mellon University

# Hidden Conditional Random Field

Sequence label:

$y \in \mathcal{Y}$ for example, $\mathcal{Y}$: {**positive, negative**}

Latent variables with shared hidden states:

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

**Sentiment**



We saw the yellow dog

$$p(y, \boldsymbol{h} \mid \mathbf{x}; \boldsymbol{\theta}) = \frac{1}{\mathcal{Z}(\boldsymbol{x}; \boldsymbol{\theta})} \exp\left\{ \sum_t \boldsymbol{\theta}^x \cdot f^x(h_t, \mathbf{x}_t) + \sum_t \boldsymbol{\theta}^e \cdot f^e(h_t, h_{t-1}, y) + \sum_t \boldsymbol{\theta}^y \cdot f^y(y, \boldsymbol{h}_t) \right\}$$

$$p(y \mid x; \boldsymbol{\theta}) = \sum_h p($$

**Shared hidden states**

- **Inference is tractable: $O(YH^2T)$**
  - **Linear in sequence length T !**
- **Parameter learning ($\theta^x, \theta^e, \theta^y$):**
  - **Gradient descent or L-BFGS**

# Learning Multimodal Structure

Modality-*private* structure

- Internal grouping of observations

Modality-*shared* structure

- Interaction and synchrony



We saw the yellow dog

# Multi-view Latent Variable Discriminative Models

Modality-***private*** structure

- Internal grouping of observations

Modality-***shared*** structure

- Interaction and synchrony



We saw the yellow dog

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

➢ Approximate inference using loopy-belief

# CRFs and Deep Learning

# Conditional Neural Fields

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

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

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

Language Technologies Institute

Carnegie Mellon University

# Deep Conditional Neural Fields

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

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

$$f^x(y_i, \mathbf{x}_i) = \mathbb{I}[y_i = y'] \cdot \mathcal{G}\big(\boldsymbol{a_i^{m-1}}, \boldsymbol{W^l}\big)$$

$$a^l = \mathcal{G}\big(\boldsymbol{a_i^{l-1}}, \boldsymbol{\theta^g}\big) \quad \text{for } l = 2 \dots m - 1$$

**Iterate**

Language Technologies Institute

Carnegie Mellon University

# CRF and Bilinear LSTM [Dyer, 2016]



Output labels:
- Name entities

**Learning:**

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

Input features:
- Word embedding

➤ What did $\theta^e$ paramters learn?

➤ What does LSTM parameters learns?

Language Technologies Institute

Carnegie Mellon University

# 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)

$$\boldsymbol{y} = \{y_1, y_2, y_3, \ldots, y_t\} \text{ where } y_t \in \mathbb{R}$$

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

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

| 0.3 | 0.2 | 0.7 | 0.8 | 0.5 |

$y_1$ — $y_2$ — $y_3$ — $y_4$ — $y_5$

$g_1$ $g_2$ $g_3$ $g_4$ $g_5$

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$

We saw the yellowdog

➤ How to solve

**Multivariate Gaussian integral:**

$$\int_{-\infty}^{\infty} \exp\left\{ \tfrac{1}{2} \boldsymbol{y}^T \Sigma^{-1} \boldsymbol{y} + \boldsymbol{y} \Sigma^{-1} \boldsymbol{\mu} \right\} d\boldsymbol{y}$$

$$= \frac{(2\pi)^{n/2}}{|\Sigma^{-1}|^{1/2}} \exp\left( \tfrac{1}{2} \boldsymbol{\mu} \, \Sigma^{-1} \boldsymbol{\mu} \right)$$

[Radosavljevic et al., 2010]

Language Technologies                    Carnegie Mellon University

# Continuous Conditional Neural Field

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

$$\boldsymbol{y} = \{y_1, y_2, y_3, \dots, y_t\} \text{ where } y_t \in \mathbb{R}$$

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

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

$$f^x(y_t, x_t, \theta^g) = -\left(y_t - g_k(x_t, \theta_k^g)\right)^2$$

$$f^e(y_t, y_{t-1}) = -\frac{1}{2}(y_t - y_{t-1})^2$$



| 0.3 | 0.2 | 0.7 | 0.8 | 0.5 |

We    saw    the yellow dog

# Continuous Conditional Neural Field

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

$$\boldsymbol{y} = \{y_1, y_2, y_3, \ldots, y_t\} \text{ where } y_t \in \mathbb{R}$$

**Multivariate Gaussian distribution:**

$$p(\boldsymbol{y} \mid \boldsymbol{x}; \boldsymbol{\theta}) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\boldsymbol{y} - \boldsymbol{\mu})^T \Sigma^{-1}(\boldsymbol{y} - \boldsymbol{\mu})\right)$$

where $\Sigma$

matrix ve

and $\boldsymbol{\mu} = \Sigma$

| 0.3 | 0.2 | 0.7 | 0.8 | 0.5 |

$y_1$ — $y_2$ — $y_3$ — $y_4$ — $y_5$

$g_1$  $g_2$  $g_3$  $g_4$  $g_5$

$x_1$  $x_2$  $x_3$  $x_4$  $x_5$

We    saw    the yellowdog

**Since CCNF can be viewed as a multivariate Gaussian, the prediction of $\boldsymbol{y}'$ is simply the mean value of distribution:**

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

➢ Optimized using gradient ascent or BFGS.

# High-Order Continuous Conditional Neural Field

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

$$\boldsymbol{y} = \{y_1, y_2, y_3, \ldots, y_t\} \text{ where } y_t \in \mathbb{R}$$

**Multivariate Gaussian distribution:**

$$p(\boldsymbol{y}\,|\,\boldsymbol{x}; \boldsymbol{\theta}) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\boldsymbol{y} - \boldsymbol{\mu})^T \Sigma^{-1}(\boldsymbol{y} - \boldsymbol{\mu})\right)$$

*k*-order potential functions:

$$f^{e_k}(y_t, y_{t-k}) = -\frac{1}{2}(y_t - y_{t-k})^2$$

0.3  0.2  0.7  0.8  0.5

$y_1$  $y_2$  $y_3$  $y_4$  $y_5$

$g_1$  $g_2$  $g_3$  $g_4$  $g_5$

$x_1$  $x_2$  $x_3$  $x_4$  $x_5$

We    saw    the yellow dog

# Fully-Connected Continuous Conditional Neural Field

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

$$\boldsymbol{y} = \{y_1, y_2, y_3, \ldots, y_t\} \text{ where } y_t \in \mathbb{R}$$

**Multivariate Gaussian distribution:**

$$p(\boldsymbol{y}|\,\boldsymbol{x}; \boldsymbol{\theta}) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\boldsymbol{y} - \boldsymbol{\mu})^T \Sigma^{-1}(\boldsymbol{y} - \boldsymbol{\mu})\right)$$

$k$-order potential functions:

$$f^{e_k}(y_t, y_{t-k}) = -\frac{1}{2}(y_t - y_{t-k})^2$$

*Grid* potential functions:

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

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



We    saw    the yellow dog

# Fully-Connected CRF [Krahenbuhl and Koltun, 2013]

"Semantic" image segmentation





$y_i$: object class label

$x_i$: local pixel features

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

**Mixture of kernels**

where $\quad \Phi_{ij}(y_i, y_j; \boldsymbol{\theta}) = \sum_{m=1}^{C} u^{(m)}(y_i, y_j|\boldsymbol{\theta}) k^{(m)}(\boldsymbol{x_i}, \boldsymbol{x_j})$

# CNN and Fully-Connected CRF [Chen et al., 2014]

Language Technologies Institute

Carnegie Mellon University

# 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 $q^t_{(x,y),i}(\hat{y}_i)$ via filtering for $t \in \{1, \ldots, T\}$ ← Using mean field approximation
3. Backtracking through the marginals $q^t_{(x,y),i}(\hat{y}_i)$ from $t = T - 1$ down to $t = 1$
4. Backward pass through definition of function via chain rule
5. Parameter update