



# 10-708 Probabilistic Graphical Models

Machine Learning Department  
School of Computer Science  
Carnegie Mellon University



## Learning MRFs & CRFs + Neural Potential Functions

Matt Gormley  
Lecture 7  
Feb. 22, 2021

# Q&A

**Q:** What are the subsets in Question 4 of HW1?  
Is the question even asking what you meant to ask?

**A:** The subset notation in Q4 follows I-map notation (e.g. page 60, of Koller & Friedman textbook). Alas, we never introduced this notation in class and, as such, the whole question was not asking what we wanted.

Since this was our mistake, we'll figure out someday to even the playing field on Q4. I'll post more detail on Piazza after class...

# Reminders

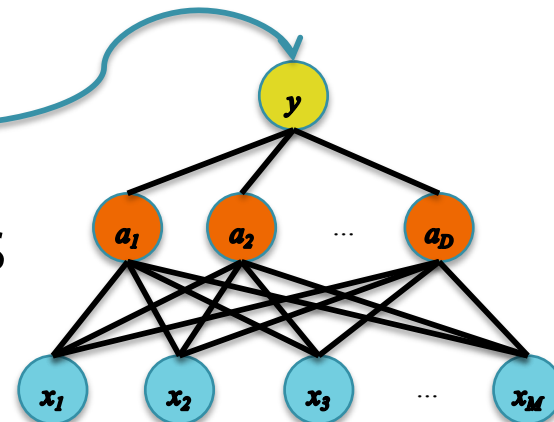
- **Homework 1: PGM Representation**
  - Out: Mon, Feb. 15
  - Due: Mon, Feb. 22 at 11:59pm
- **Homework 2: Exact inference and supervised learning (CRF+RNN)**
  - Out: Mon, Feb. 22
  - Due: Mon, Mar. 08 at 11:59pm

# **MRF LEARNING (LOG-LINEAR CASE)**

# Options for MLE of MRFs

- **Setting I:**  $\psi_C(\mathbf{x}_C) = \theta_{C, \mathbf{x}_C}$ 
  - A. MLE by inspection (Decomposable Models)
  - B. Iterative Proportional Fitting (IPF)
- **Setting II:**  $\psi_C(\mathbf{x}_C) = \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}_C))$ 
  - C. Generalized Iterative Scaling
  - D. Gradient-based Methods

- **Setting III:**  $\psi_C(\mathbf{x}_C) =$ 
  - E. Gradient-based Methods



# Title Here

## ***Whiteboard***

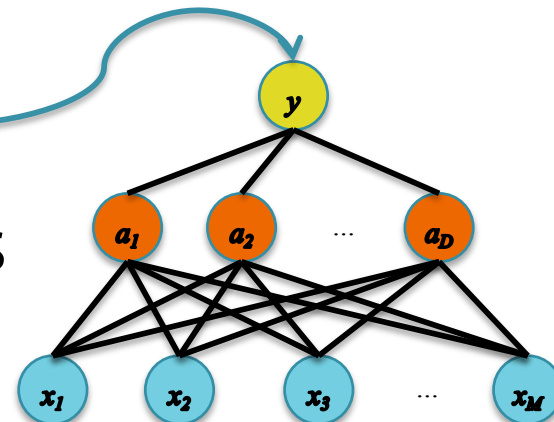
- log-linear MRF model (i.e. with feature based potentials)
- log-linear MRF derivatives
- log-linear MRF training with SGD

# **LOG-LINEAR PARAMETERIZATION OF CONDITIONAL RANDOM FIELD**

# Options for MLE of MRFs

- **Setting I:**  $\psi_C(\mathbf{x}_C) = \theta_{C, \mathbf{x}_C}$ 
  - A. MLE by inspection (Decomposable Models)
  - B. Iterative Proportional Fitting (IPF)
- **Setting II:**  $\psi_C(\mathbf{x}_C) = \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}_C))$ 
  - C. Generalized Iterative Scaling
  - D. Gradient-based Methods

- **Setting III:**  $\psi_C(\mathbf{x}_C) =$ 
  - E. Gradient-based Methods

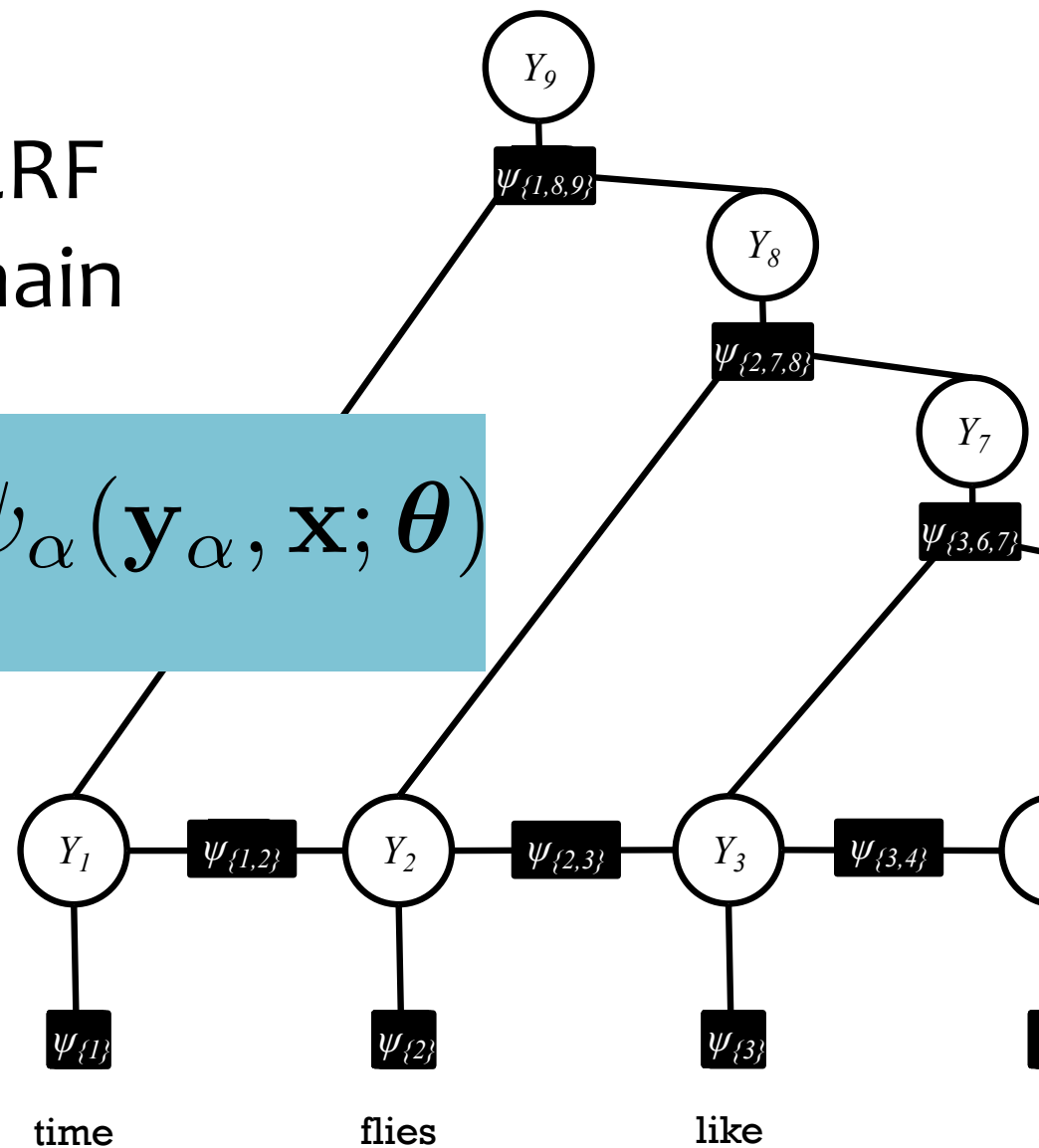




# General CRF

The topology of the graphical model for a CRF doesn't have to be a chain

$$p_{\theta}(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{\alpha} \psi_{\alpha}(\mathbf{y}_{\alpha}, \mathbf{x}; \theta)$$



# Log-linear CRF Parameterization

$$p_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{\alpha} \psi_{\alpha}(\mathbf{y}_{\alpha}, \mathbf{x}; \boldsymbol{\theta})$$

Define each potential function in terms of a fixed set of feature functions:

$$\psi_{\alpha}(\mathbf{y}_{\alpha}, \mathbf{x}; \boldsymbol{\theta}) = \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\alpha}(\mathbf{y}_{\alpha}, \mathbf{x}))$$

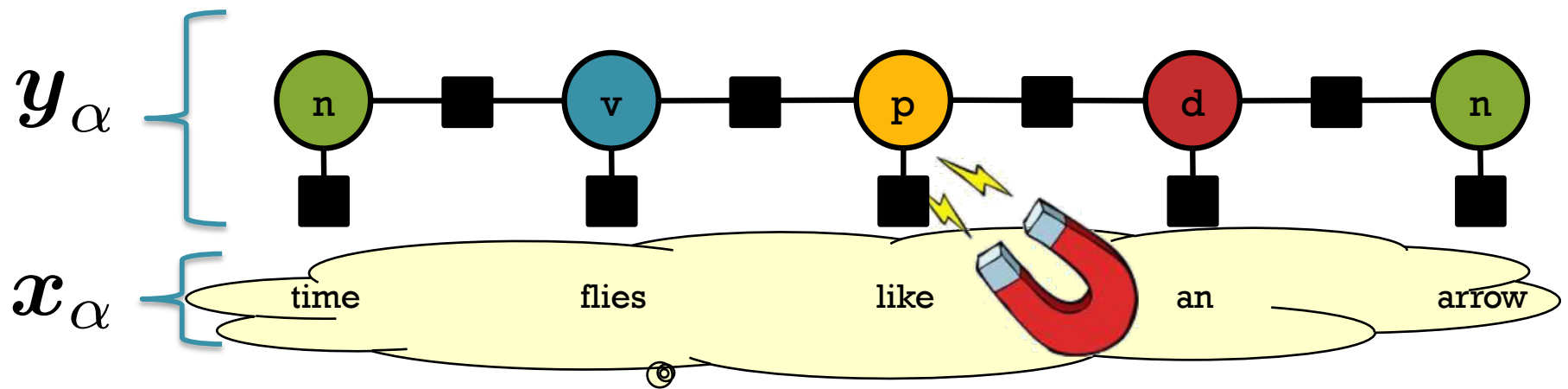
Predicted  
variables

Observed  
variables

# Log-linear CRF Parameterization

Define each potential function in terms of a fixed set of feature functions:

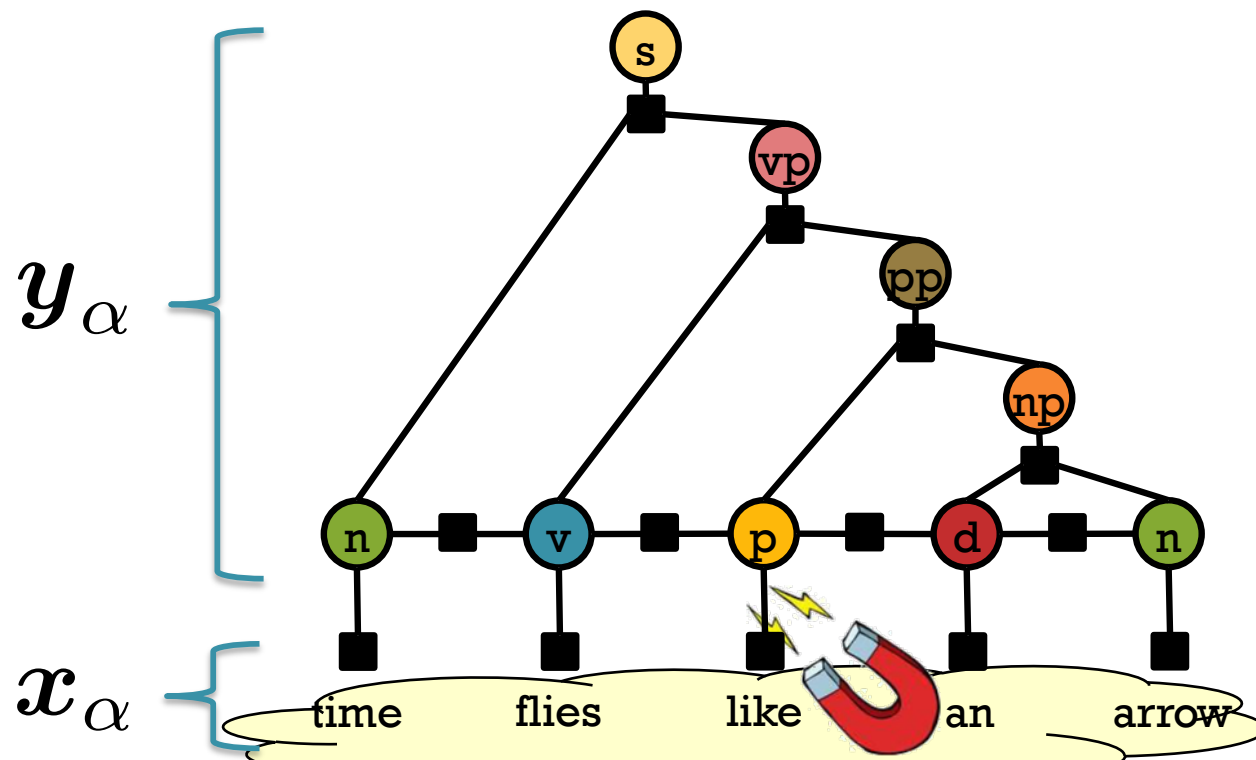
$$\psi_{\alpha}(\mathbf{y}_{\alpha}, \mathbf{x}; \boldsymbol{\theta}) = \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\alpha}(\mathbf{y}_{\alpha}, \mathbf{x}))$$



# Log-linear CRF Parameterization

Define each potential function in terms of a fixed set of feature functions:

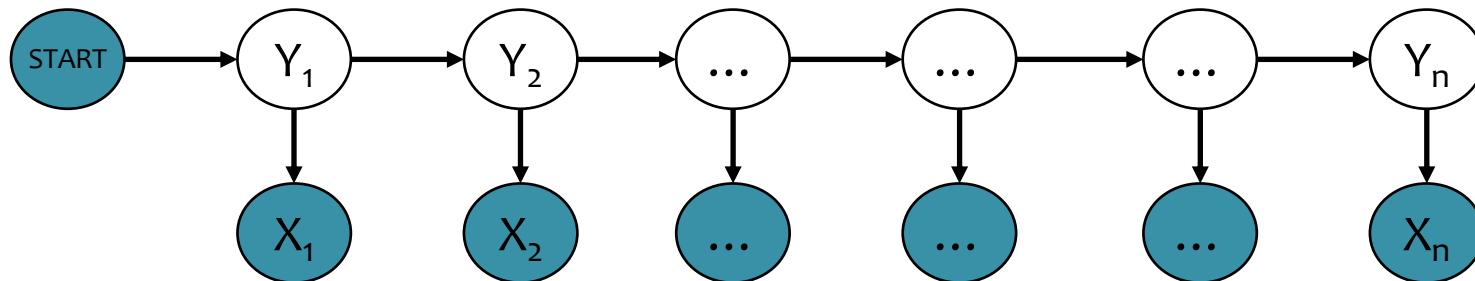
$$\psi_{\alpha}(\mathbf{y}_{\alpha}, \mathbf{x}; \boldsymbol{\theta}) = \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\alpha}(\mathbf{y}_{\alpha}, \mathbf{x}))$$



Conditional Random Fields (CRFs) for time series data

# **LINEAR-CHAIN CRFS**

# Shortcomings of Hidden Markov Models

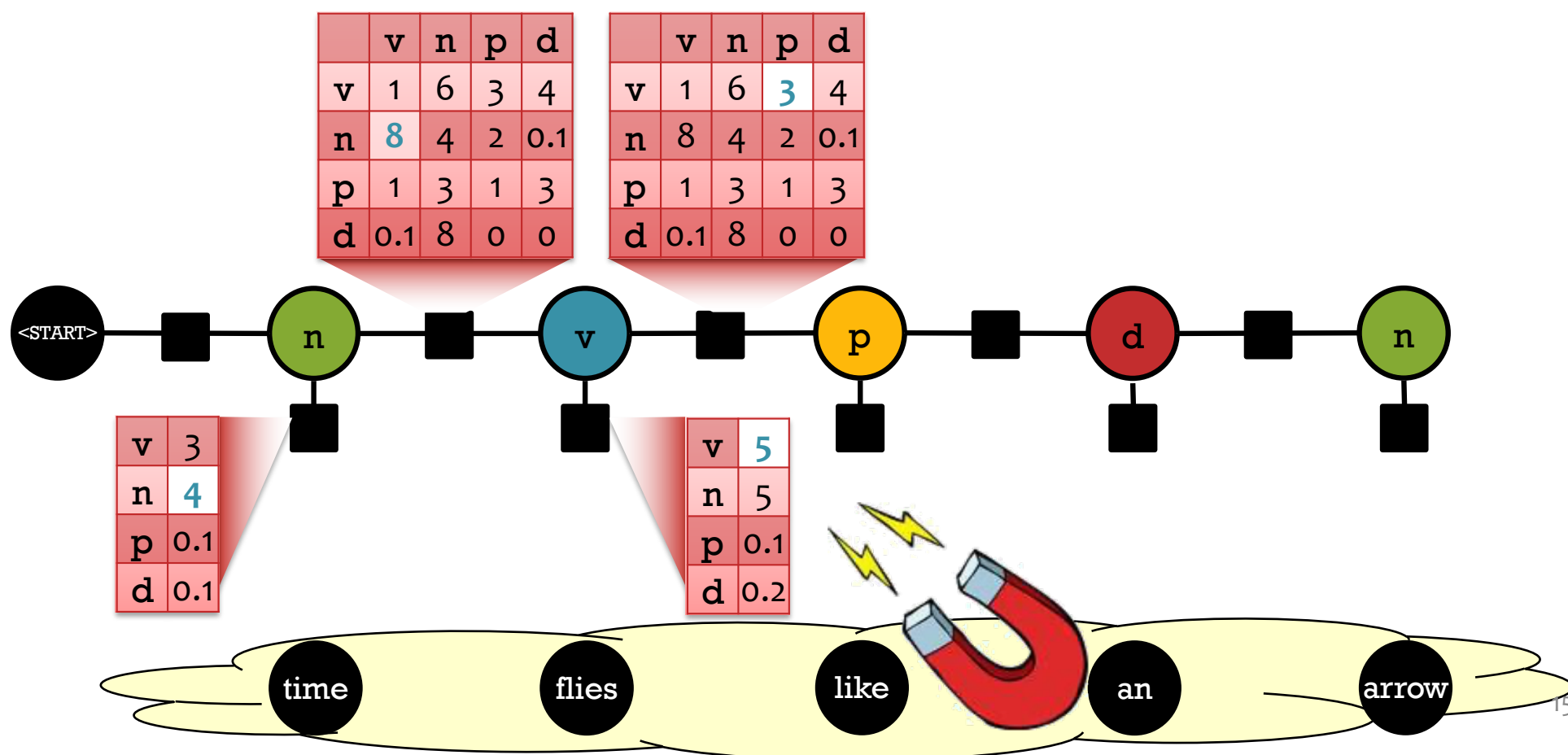


- HMM models capture dependences between each state and **only** its corresponding observation
  - NLP example: In a sentence segmentation task, each segmental state may depend not just on a single word (and the adjacent segmental stages), but also on the (non-local) features of the whole line such as line length, indentation, amount of white space, etc.
- Mismatch between learning objective function and prediction objective function
  - HMM learns a joint distribution of states and observations  $P(\mathbf{Y}, \mathbf{X})$ , but in a prediction task, we need the conditional probability  $P(\mathbf{Y}|\mathbf{X})$

# Conditional Random Field (CRF)

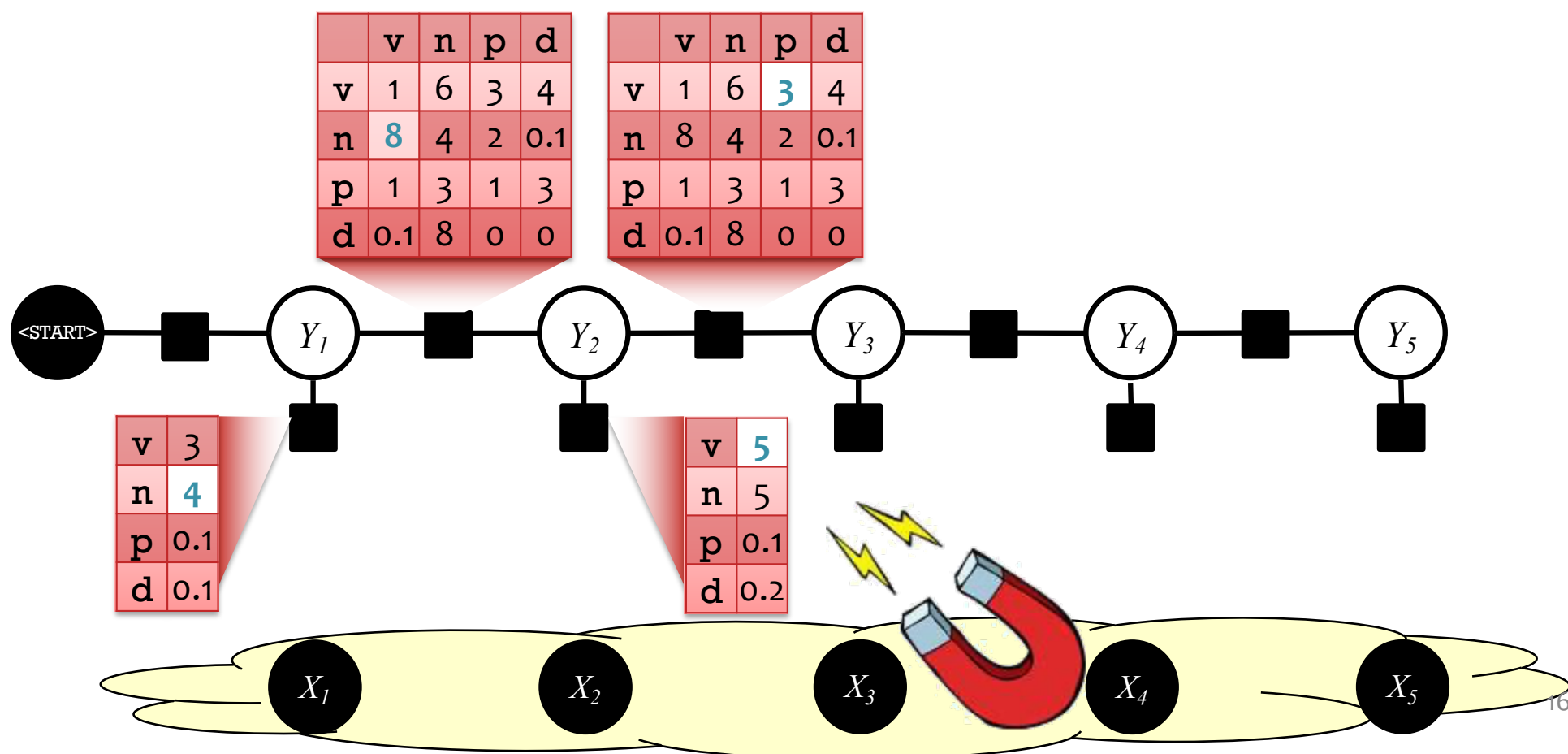
Conditional distribution over tags  $X_i$  given words  $w_i$ .  
The factors and  $Z$  are now specific to the sentence  $w$ .

$$p(n, v, p, d, n \mid \text{time, flies, like, an, arrow}) = \frac{1}{Z} (4 * 8 * 5 * 3 * \dots)$$



# Conditional Random Field (CRF)

Recall: Shaded nodes in a graphical model are **observed**

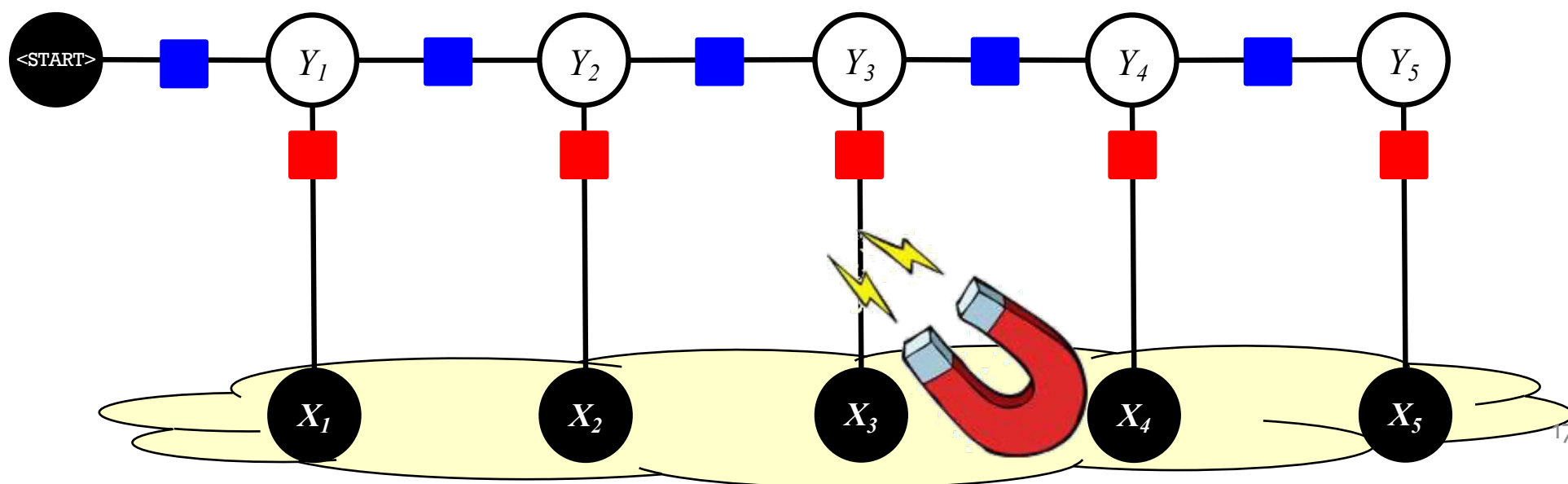




# Conditional Random Field (CRF)

This **linear-chain CRF** is just **like an HMM**, except that its factors are **not** necessarily probability distributions

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \psi_{\text{em}}(y_k, x_k) \psi_{\text{tr}}(y_k, y_{k-1})$$
$$= \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\text{em}}(y_k, x_k)) \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\text{tr}}(y_k, y_{k-1}))$$



# Exercise

$$\begin{aligned} p(\mathbf{y}|\mathbf{x}) &= \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \psi_{\text{em}}(y_k, x_k) \psi_{\text{tr}}(y_k, y_{k-1}) \\ &= \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\text{em}}(y_k, x_k)) \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\text{tr}}(y_k, y_{k-1})) \end{aligned}$$

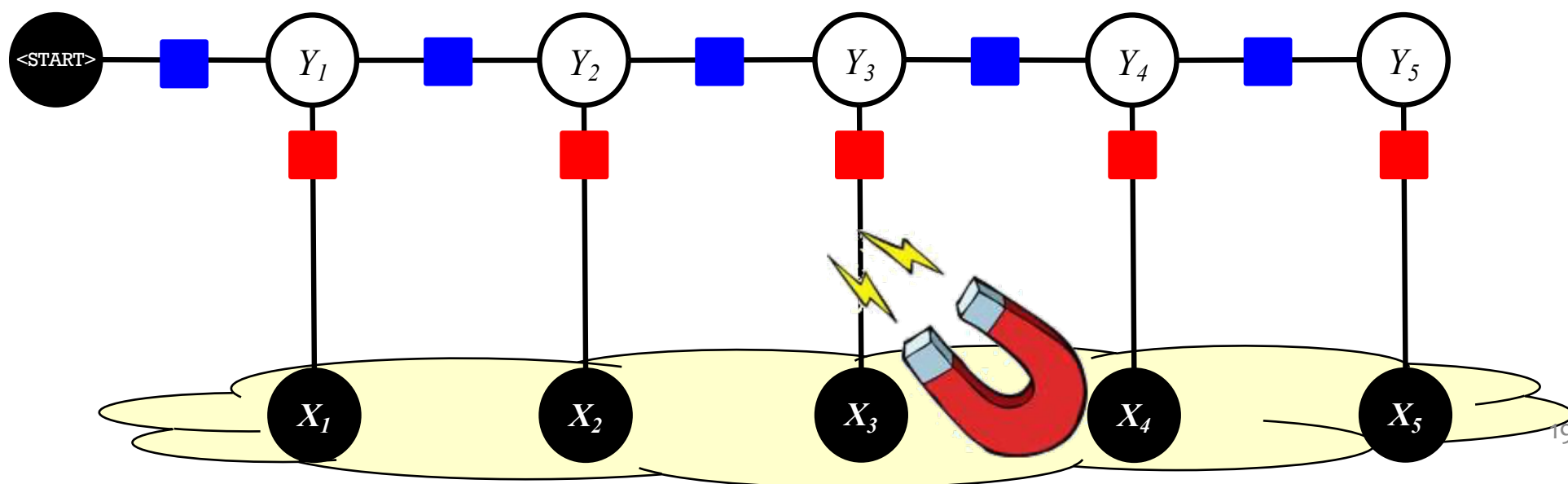
**Select All That Apply:** Which model does the above distribution share the most in common with?

- A. Hidden Markov Model
- B. Bernoulli Naïve Bayes
- C. Gaussian Naïve Bayes
- D. Logistic Regression

# Conditional Random Field (CRF)

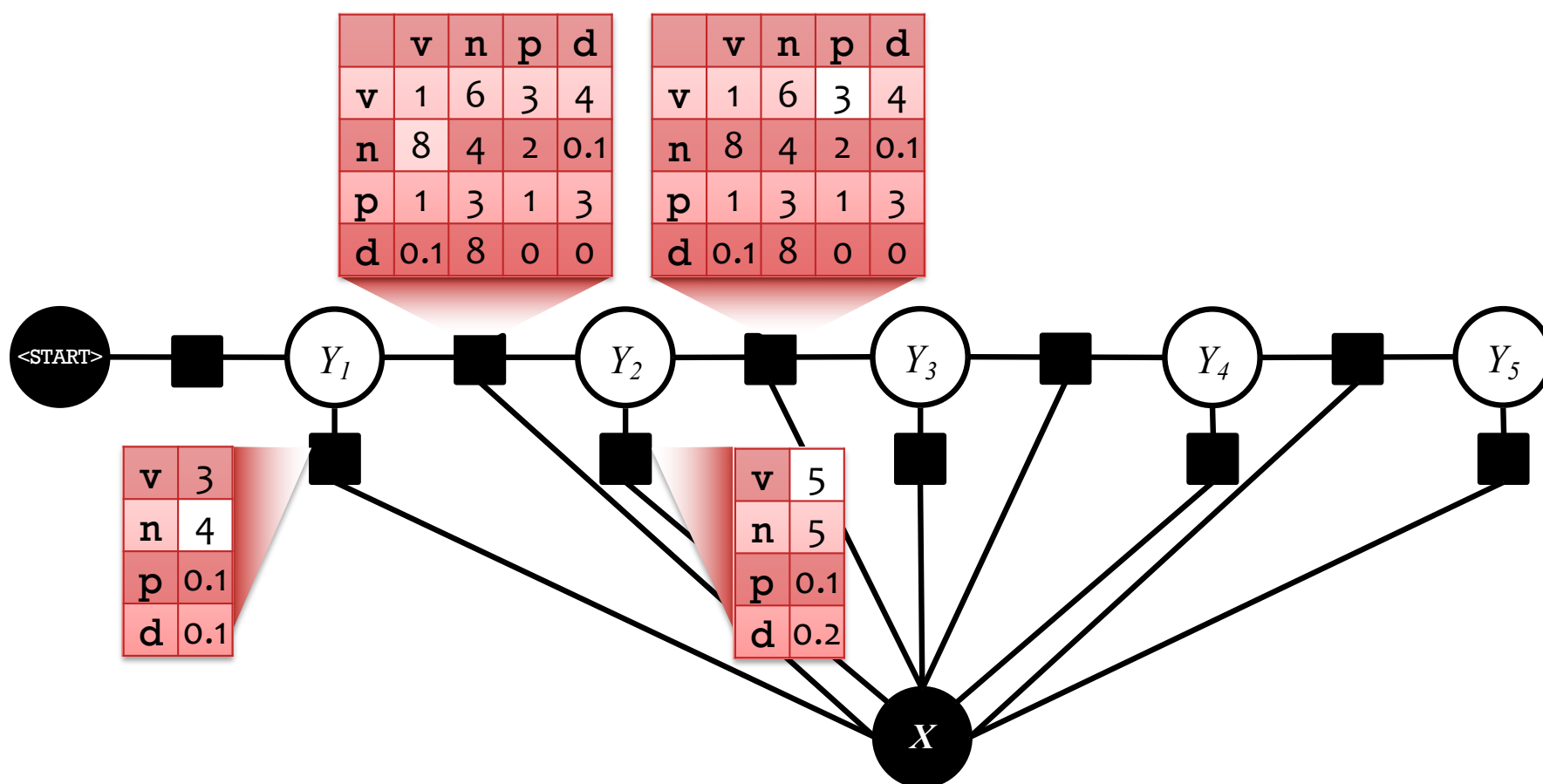
This **linear-chain CRF** is just **like an HMM**, except that its factors are **not** necessarily probability distributions

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \psi_{\text{em}}(y_k, x_k) \psi_{\text{tr}}(y_k, y_{k-1})$$
$$= \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\text{em}}(y_k, x_k)) \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\text{tr}}(y_k, y_{k-1}))$$



# Conditional Random Field (CRF)

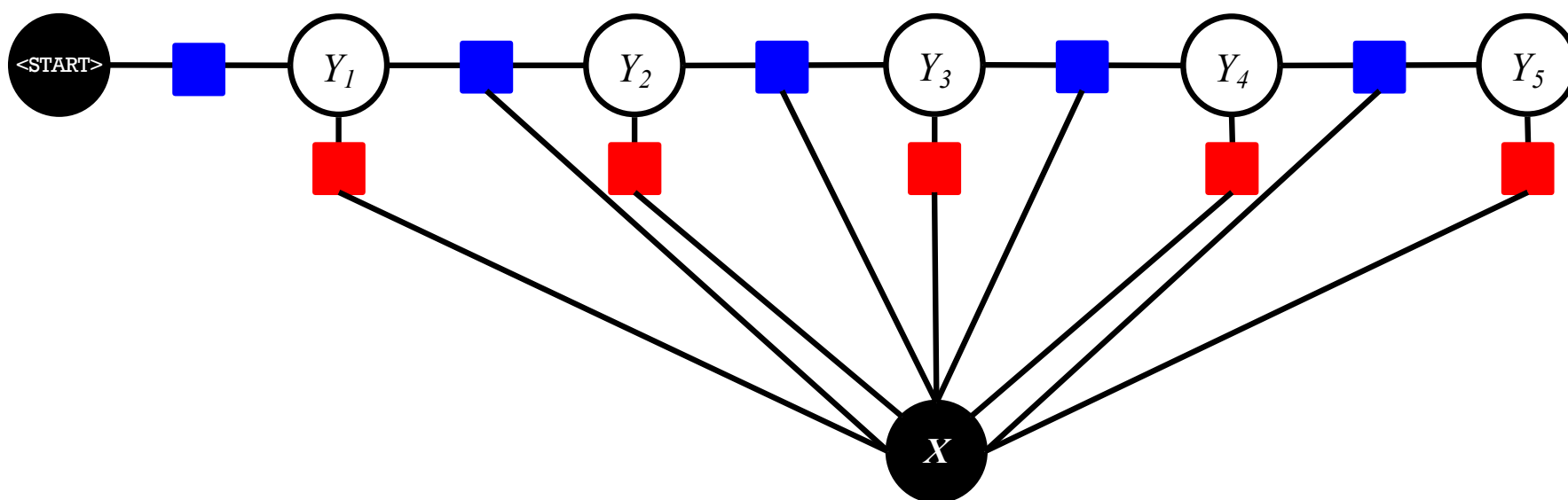
- That is the **vector**  $X$
- Because it's observed, we can condition on it for free
- Conditioning is how we converted from the MRF to the CRF (i.e. when taking a slice of the emission factors)



# Conditional Random Field (CRF)

- This is the **standard** linear-chain CRF definition
- It permits rich, overlapping features of the vector  $\mathbf{X}$

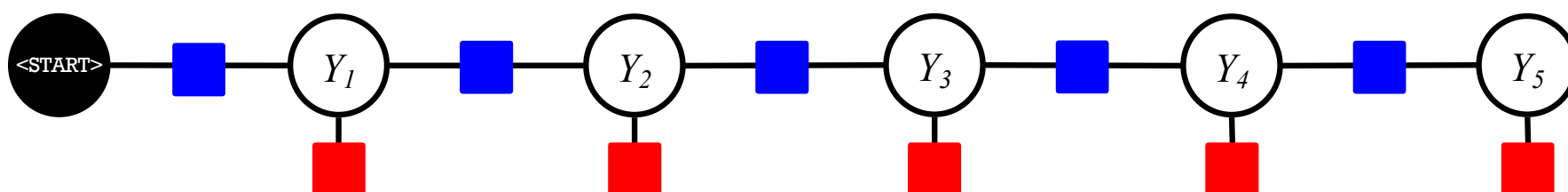
$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \psi_{\text{em}}(y_k, \mathbf{x}) \psi_{\text{tr}}(y_k, y_{k-1}, \mathbf{x})$$
$$= \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\text{em}}(y_k, \mathbf{x})) \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\text{tr}}(y_k, y_{k-1}, \mathbf{x}))$$



# Conditional Random Field (CRF)

- This is the **standard** linear-chain CRF definition
- It permits rich, overlapping features of the vector  $\mathbf{X}$

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \psi_{\text{em}}(y_k, \mathbf{x}) \psi_{\text{tr}}(y_k, y_{k-1}, \mathbf{x})$$
$$= \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\text{em}}(y_k, \mathbf{x})) \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\text{tr}}(y_k, y_{k-1}, \mathbf{x}))$$



**Visual Notation:** Usually we draw a CRF **without** showing the variable corresponding to  $\mathbf{X}$

# **CRF LEARNING (LOG-LINEAR CASE)**

# Title Here

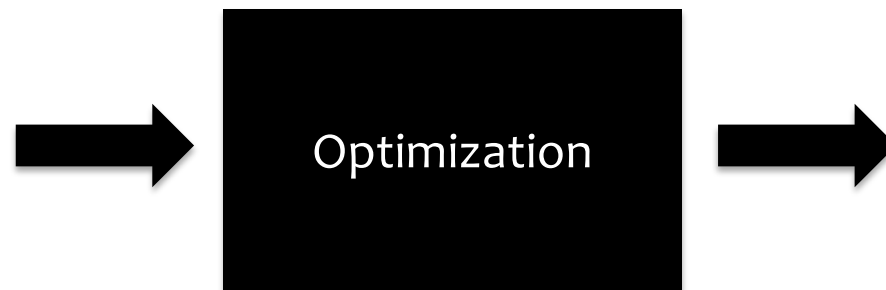
## ***Whiteboard***

- log-linear CRF model (i.e. with feature based potentials)
- log-linear CRF derivatives
- log-linear CRF training with SGD



# Recipe for Gradient-based Learning

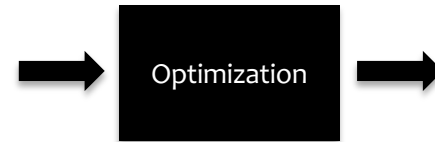
1. Write down the objective function
2. Compute the partial derivatives of the objective (i.e. gradient, and maybe Hessian)
3. Feed objective function and derivatives into black box



4. Retrieve optimal parameters from black box

# Optimization Algorithms

**What is the black box?**



- Newton's method
- Hessian-free / Quasi-Newton methods
  - Conjugate gradient
  - L-BFGS
- Stochastic gradient methods
  - Stochastic gradient descent (SGD)
  - SGD with momentum
  - Adam

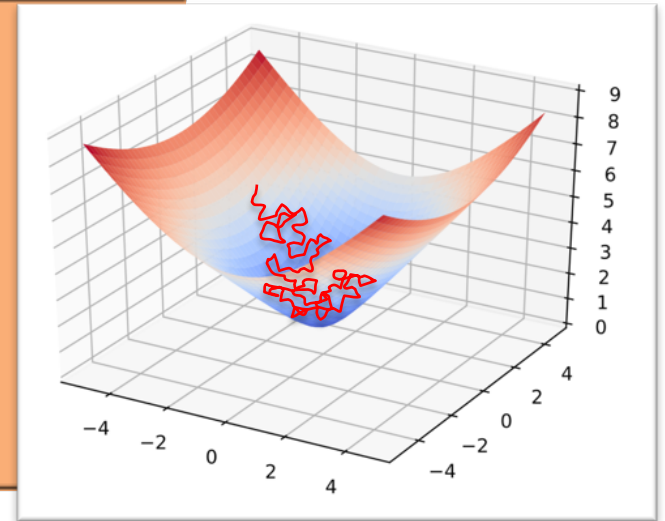
# Stochastic Gradient Descent (SGD)

Assume we have an objective that decomposes additively:

$$\text{Let } J(\boldsymbol{\theta}) = \sum_{i=1}^N J^{(i)}(\boldsymbol{\theta})$$

## Algorithm 2 Stochastic Gradient Descent (SGD)

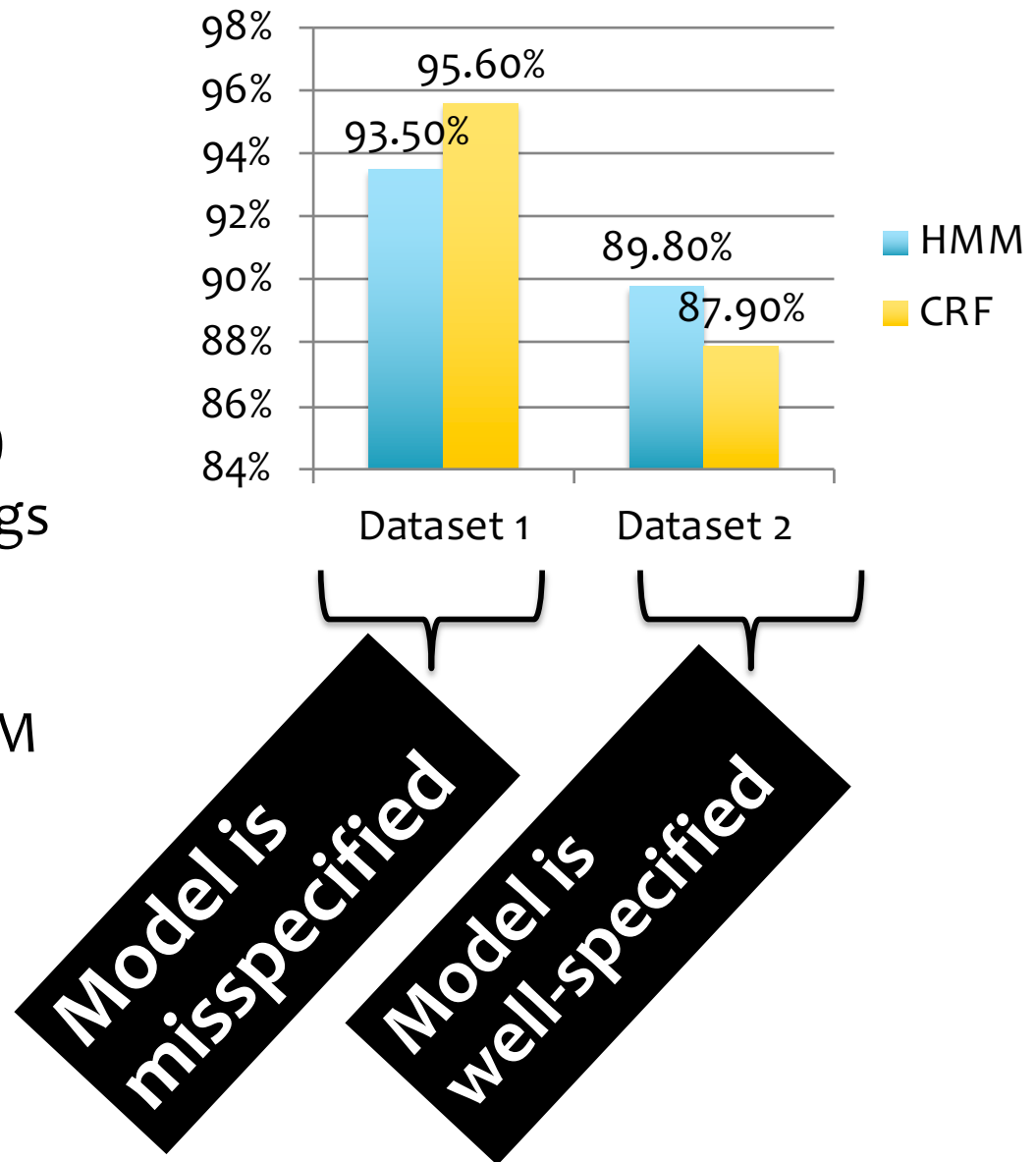
```
1: procedure SGD( $\mathcal{D}, \boldsymbol{\theta}^{(0)}$ )  
2:    $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}^{(0)}$   
3:   while not converged do  
4:      $i \sim \text{Uniform}(\{1, 2, \dots, N\})$   
5:      $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \gamma \nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta})$   
6:   return  $\boldsymbol{\theta}$ 
```



# Generative vs. Discriminative

Liang & Jordan (ICML 2008) compares **HMM** and **CRF** with **identical features**

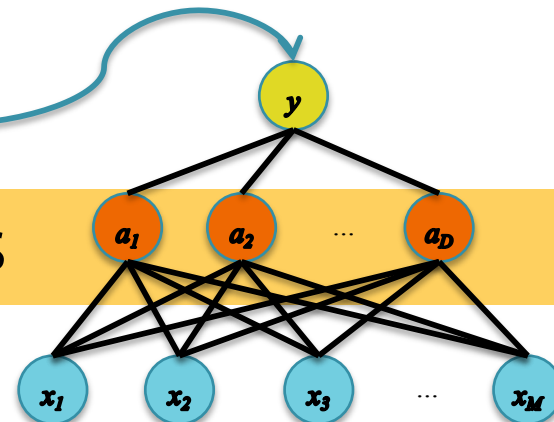
- Dataset 1: (Real)
  - WSJ Penn Treebank (38K train, 5.5K test)
  - 45 part-of-speech tags
- Dataset 2: (Artificial)
  - Synthetic data generated from HMM learned on Dataset 1 (1K train, 1K test)
- Evaluation Metric: Accuracy



# Options for MLE of MRFs

- **Setting I:**  $\psi_C(\mathbf{x}_C) = \theta_{C, \mathbf{x}_C}$ 
  - A. MLE by inspection (Decomposable Models)
  - B. Iterative Proportional Fitting (IPF)
- **Setting II:**  $\psi_C(\mathbf{x}_C) = \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}_C))$ 
  - C. Generalized Iterative Scaling
  - D. Gradient-based Methods

- **Setting III:**  $\psi_C(\mathbf{x}_C) =$ 
  - E. Gradient-based Methods



# **NEURAL POTENTIAL FUNCTIONS**

# Hybrids of Graphical Models and Neural Networks

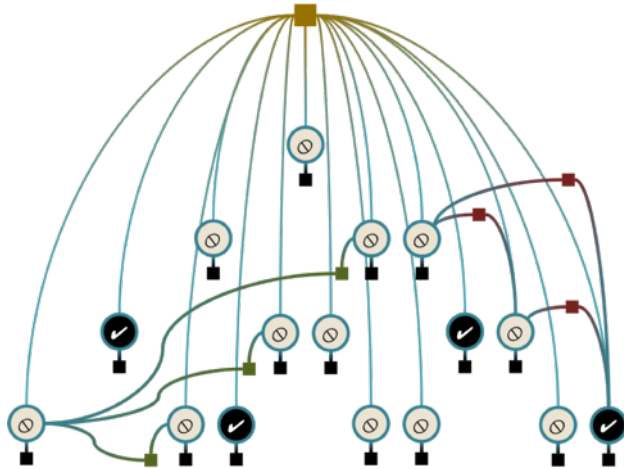
This lecture is not about a  
convergence of the two fields.

Rather, it is about state-of-the-art  
collaboration between two  
complementary techniques.

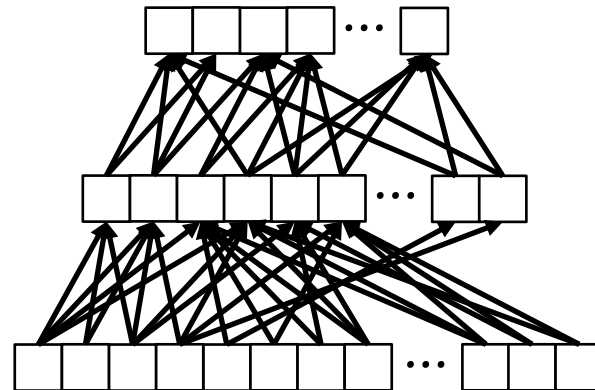
# Motivation:

## Hybrid Models

**Graphical models** let you encode domain knowledge



**Neural nets** are really good at fitting the data discriminatively to make good predictions



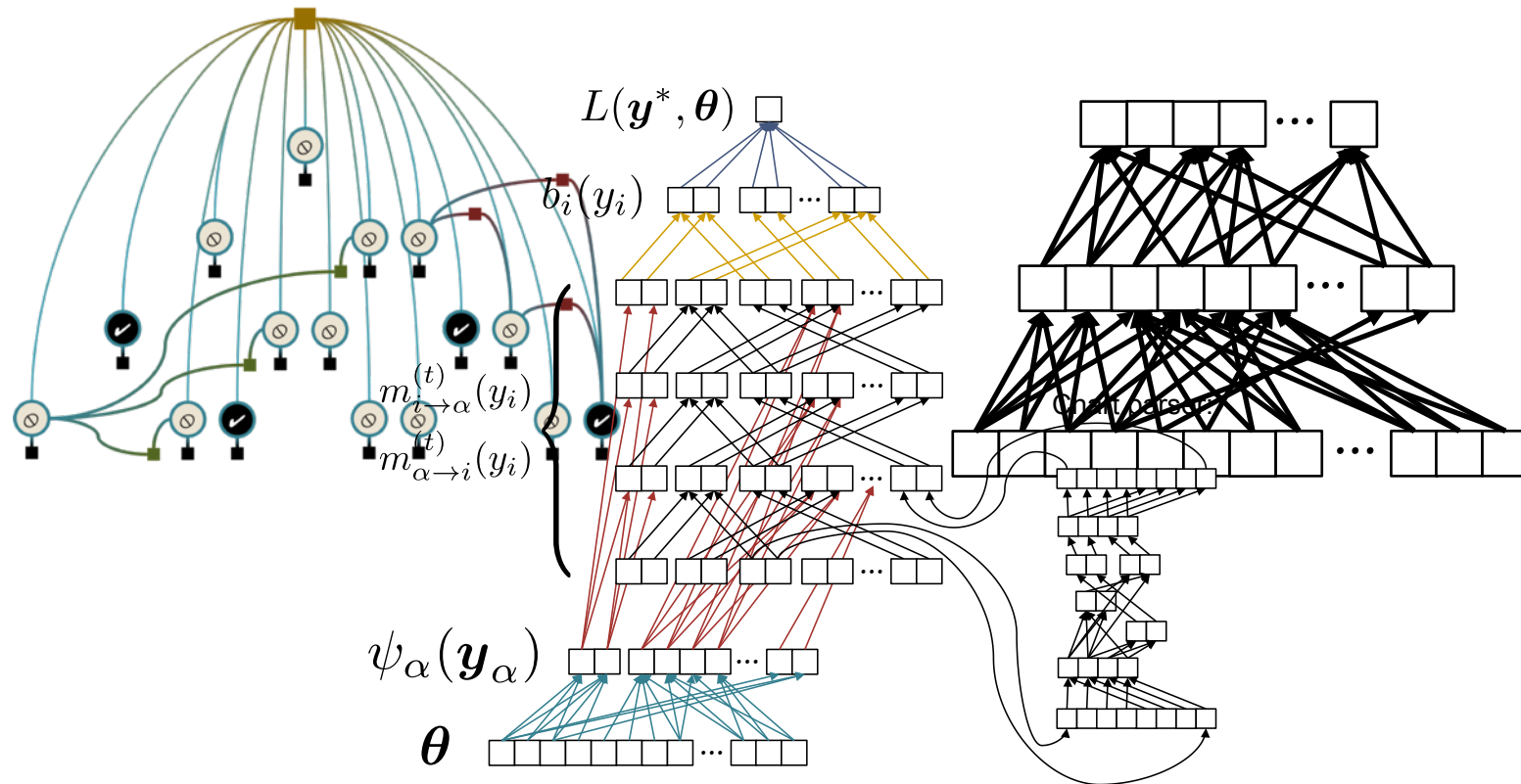
**Could we define a neural net  
that incorporates  
domain knowledge?**



# Motivation:

## Hybrid Models

*Key idea:* Use a NN to learn features for a GM, then train the entire model by backprop



# A Recipe for Neural Networks

1. Given training data:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

2. Choose each of these:

– Decision function

$$\hat{\mathbf{y}} = f_{\theta}(\mathbf{x}_i)$$

– Loss function

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

*Face*



*Face*



*Not a face*



**Examples:** Linear regression,  
Logistic regression, Neural Network

**Examples:** Mean-squared error,  
Cross Entropy

# A Recipe for Neural Networks

1. Given training data:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

2. Choose each of these:

– Decision function

$$\hat{\mathbf{y}} = f_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

– Loss function

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

3. Define goal:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^N \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

4. Train with SGD:

(take small steps  
opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

# **BACKPROPAGATION AND BELIEF PROPAGATION**

## ***Whiteboard:***

- CRF w/LSTM potentials
- Gradient of MRF/CRF log-likelihood with respect to log potentials
- Gradient of MRF/CRF log-likelihood with respect to potentials
- Backprop with MRF/CRF log-likelihood as a loss function

# Factor Derivatives

Log-probability:

$$\log p(\mathbf{y}) = \left[ \sum_{\alpha} \log \psi_{\alpha}(\mathbf{y}_{\alpha}) \right] - \log \sum_{\mathbf{y}' \in \mathcal{Y}} \prod_{\alpha} \psi_{\alpha}(\mathbf{y}'_{\alpha}) \quad (1)$$

Derivatives:

$$\frac{\partial \log p(\mathbf{y})}{\partial \log \psi_{\alpha}(\mathbf{y}'_{\alpha})} = \mathbb{1}(\mathbf{y}_{\alpha} = \mathbf{y}'_{\alpha}) - p(\mathbf{y}'_{\alpha}) \quad (2)$$

$$\frac{\partial \log p(\mathbf{y})}{\partial \psi_{\alpha}(\mathbf{y}'_{\alpha})} = \frac{\mathbb{1}(\mathbf{y}_{\alpha} = \mathbf{y}'_{\alpha}) - p(\mathbf{y}'_{\alpha})}{\psi_{\alpha}(\mathbf{y}'_{\alpha})} \quad (3)$$