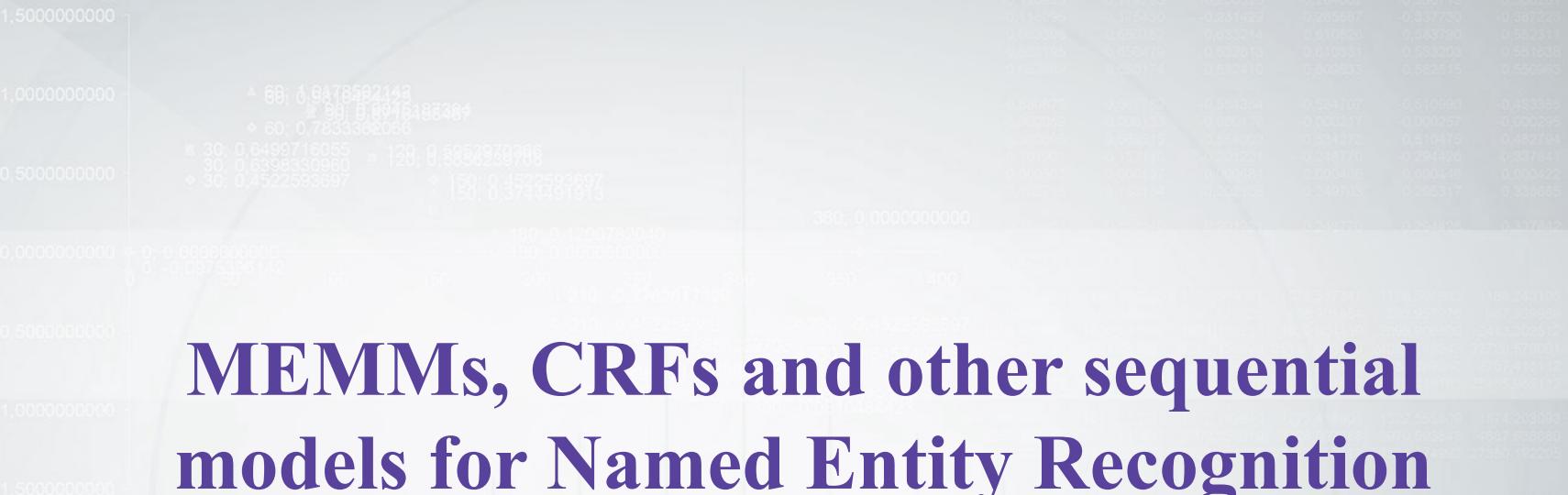


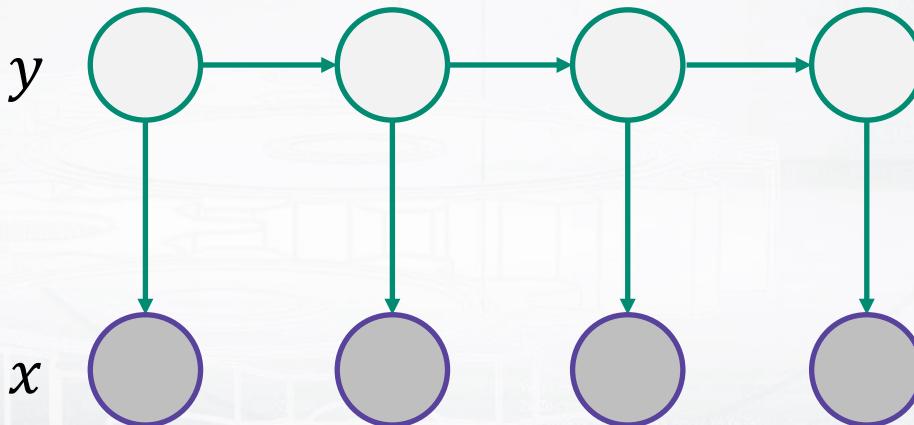
# MEMMs, CRFs and other sequential models for Named Entity Recognition



# Hidden Markov Model

$$p(\mathbf{y}, \mathbf{x}) = \prod_{t=1}^T p(y_t | y_{t-1}) p(x_t | y_t)$$

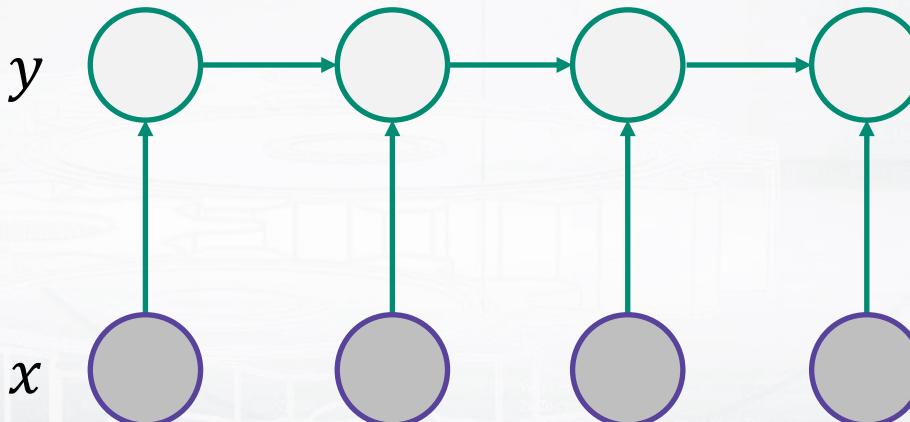
Generative  
model



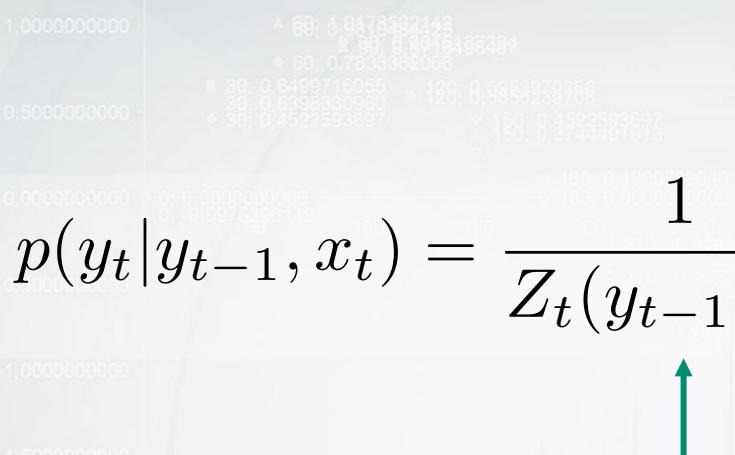
# Maximum Entropy Markov Model

$$p(\mathbf{y}|\mathbf{x}) = \prod_{t=1}^T p(y_t|y_{t-1}, x_t)$$

↑  
Discriminative  
model



# Maximum Entropy Markov Model

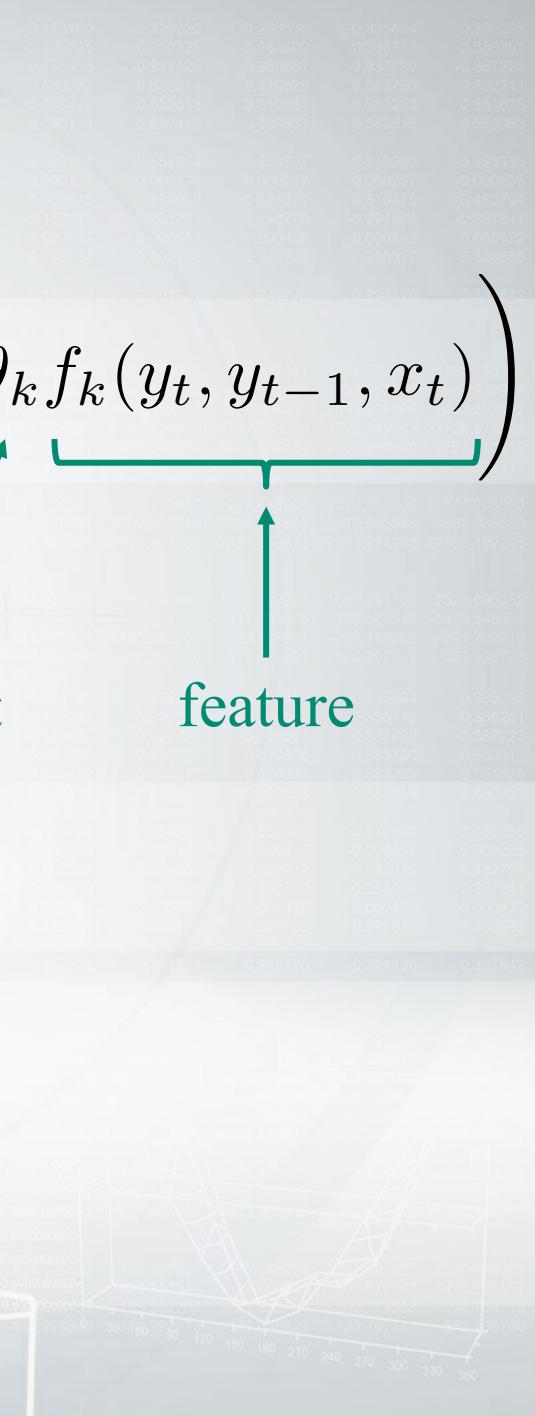


$$p(y_t|y_{t-1}, x_t) = \frac{1}{Z_t(y_{t-1}, x_t)} \exp \left( \sum_{k=1}^K \theta_k f_k(y_t, y_{t-1}, x_t) \right)$$

Normalization  
constant

$$Z_t(y_{t-1}, x_t) = \sum_{y_t} \exp \left( \sum_{k=1}^K \theta_k f_k(y_t, y_{t-1}, x_t) \right)$$

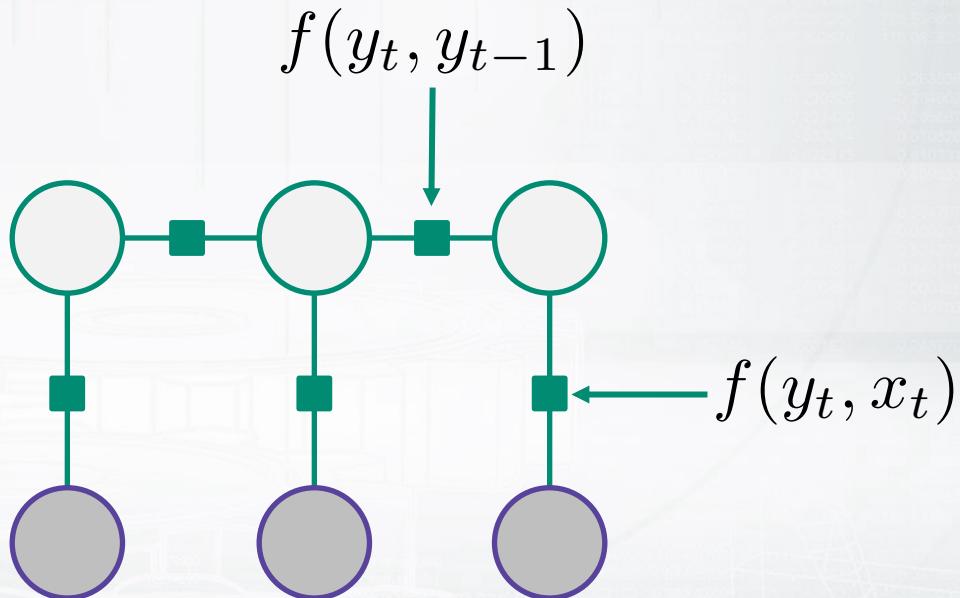
↑  
weight      feature



# Conditional Random Field (linear chain)

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(x)} \prod_{t=1}^T \exp \left( \sum_{k=1}^K \theta_k f_k(y_t, y_{t-1}, x_t) \right)$$

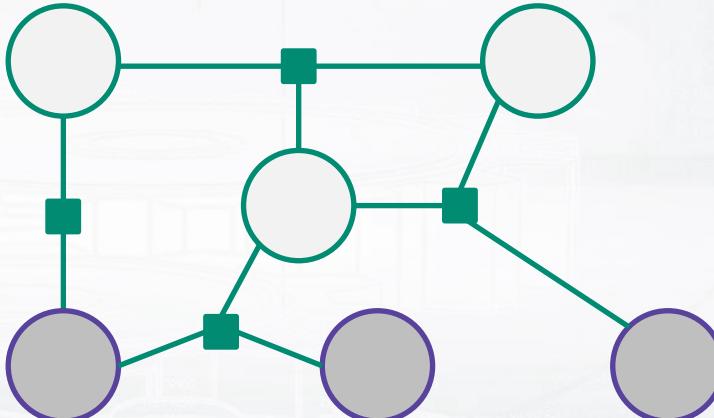
Undirected graph:



# Conditional Random Field (general form)

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(x)} \prod_{a=1}^A \Psi_a(y_a, x_a)$$

Arbitrary factors



# Black-box implementations

CRF++	<a href="https://sourceforge.net/projects/crfpp/">https://sourceforge.net/projects/crfpp/</a>
MALLET	<a href="http://mallet.cs.umass.edu/">http://mallet.cs.umass.edu/</a>
GRMM	<a href="http://mallet.cs.umass.edu/grmm/">http://mallet.cs.umass.edu/grmm/</a>
CRFSuite	<a href="http://www.chokkan.org/software/crfsuite/">http://www.chokkan.org/software/crfsuite/</a>
FACTORIE	<a href="http://www.factorie.cc">http://www.factorie.cc</a>

<http://homepages.inf.ed.ac.uk/csutton/publications/crftut-fnt.pdf>

# Features engineering

**Label-observation** features:

$$\bullet \quad f(y_t, y_{t-1}, x_t) = [y_t = y] g_m(x_t)$$

$$\bullet \quad f(y_t, y_{t-1}, x_t) = [y_t = y][y_{t-1} = y']$$

$$\bullet \quad f(y_t, y_{t-1}, x_t) = [y_t = y][y_{t-1} = y'] g_m(x_t)$$

# Observation function examples

-0.115751	0.173160	0.229220	0.263536	-0.335684	0.385298
-0.115939	-0.174295	-0.230325	-0.264602	-0.336713	0.386261
-0.116093	-0.175430	-0.231423	-0.265687	-0.337730	0.387223
-0.116246	0.157767	0.155744	0.161626	0.153780	0.152311
-0.116398	0.158904	0.156735	0.162615	0.154915	0.153168

$w_t = v$   $\forall v \in$

part-of-speech tag for  $w_t$  is  $j$   $\forall$  tags  $j$

$w_t$  is in a phrase of syntactic type  $j$   $\forall$  tags  $j$

Capitalized  $w_t$  matches [A-Z][a-z]+

AllCaps  $w_t$  matches [A-Z]+

EndsInDot  $w_t$  matches [^\.\.]+\.\*\.

$w_t$  matches a dash

$w_t$  appears in a list of stop words

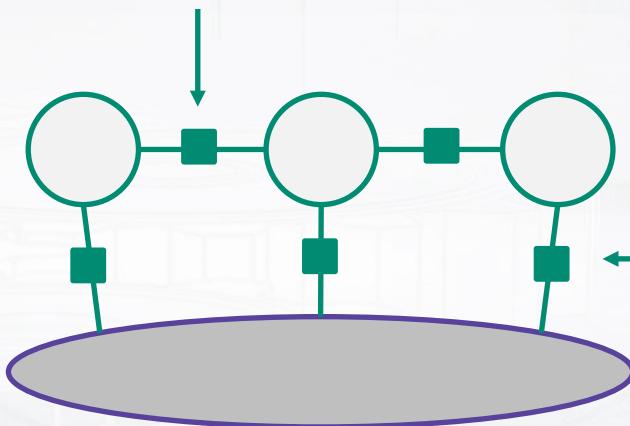
$w_t$  appears in list of capitals

# Dependencies on input

**Trick:** Pretend the current input  $x_t$  contains not only the current word  $w_t$ , but also  $w_{t-1}$  and  $w_{t+1}$  and build observation functions for them as well.

The model is discriminative, so we can use the whole input:

$$[y_t = y][y_{t-1} = y']$$



$$[y_t = y] g_m(x_t)$$

# Resume for the lesson

## Probabilistic graphical models:

- Hidden Markov Models (generative, directed)
- Maximum Entropy Markov Models (discriminative, directed)
- Conditional Random Field (discriminative, undirected)

## Tasks:

- Training – fit parameters (Baum-Welch for HMM)
- Decoding – find the most probable tags (Viterbi for HMM)

## Practice:

- Features engineering + black-box implementation