

**Mayank Kejriwal**  
University of Southern California

Conditional Random Fields

# Outline

- Modeling
- Inference
- Training
- Applications

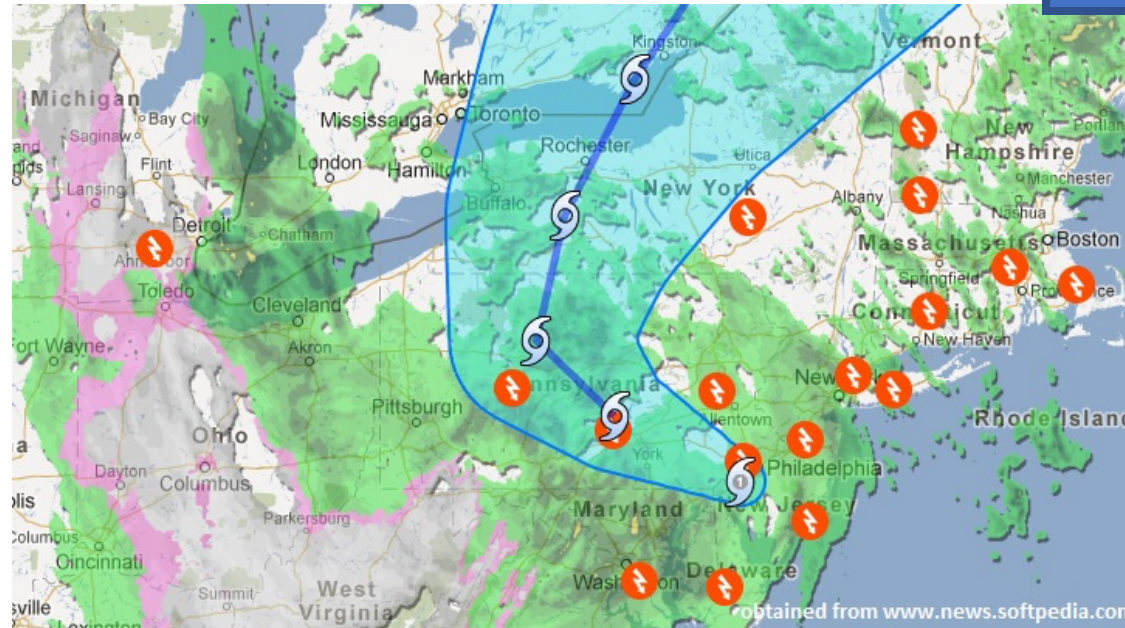
# Problem Description

- Given  $X$  (observations), find  $Y$  (predictions)
- For example,

$$\begin{cases} X = \{temperature, moisture, pressure, \dots\} \\ Y = \{Sunny, Rainy, Stormy, \dots\} \end{cases}$$

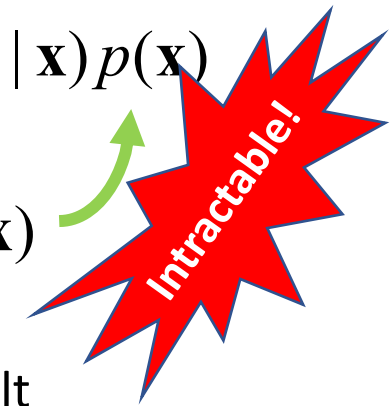
Might depend on  
previous days and  
each other

Might depend on  
previous days and  
each other



# Problem Description

- The relational connection occurs in many applications, NLP, Computer Vision, Signal Processing, ....
- Traditionally in graphical models,  $\longrightarrow p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y} | \mathbf{x}) p(\mathbf{x})$ 
  - Modeling the joint distribution can lead to difficulties
  - rich local features occur in relational data,  $\longrightarrow p(\mathbf{x})$
  - features may have complex dependencies,
    - constructing probability distribution over them is difficult
- **Solution:** directly model the conditional,  $p(\mathbf{y} | \mathbf{x})$ 
  - is sufficient for classification!
- **CRF** is simply a **conditional distribution**  $p(\mathbf{y} | \mathbf{x})$  with an **associated graphical structure**



# Discriminative Vs. Generative

$$p(\mathbf{y}, \mathbf{x})$$

- **Generative Model:** A model that generate observed data randomly
- **Naïve Bayes:** once the class label is known, all the features are independent

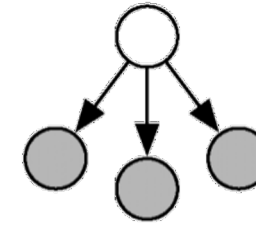
$$p(y, \mathbf{x}) = p(y) \prod_{k=1}^K p(x_k | y)$$

- **Discriminative:** Directly estimate the posterior probability; Aim at modeling the “discrimination” between different outputs

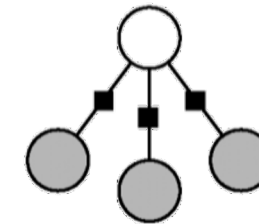
$$p(\mathbf{y} | \mathbf{x})$$

- **MaxEnt** classifier: linear combination of feature function in the exponent,

$$p(y|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp \left\{ \sum_{k=1}^K \theta_k f_k(y, \mathbf{x}) \right\}$$



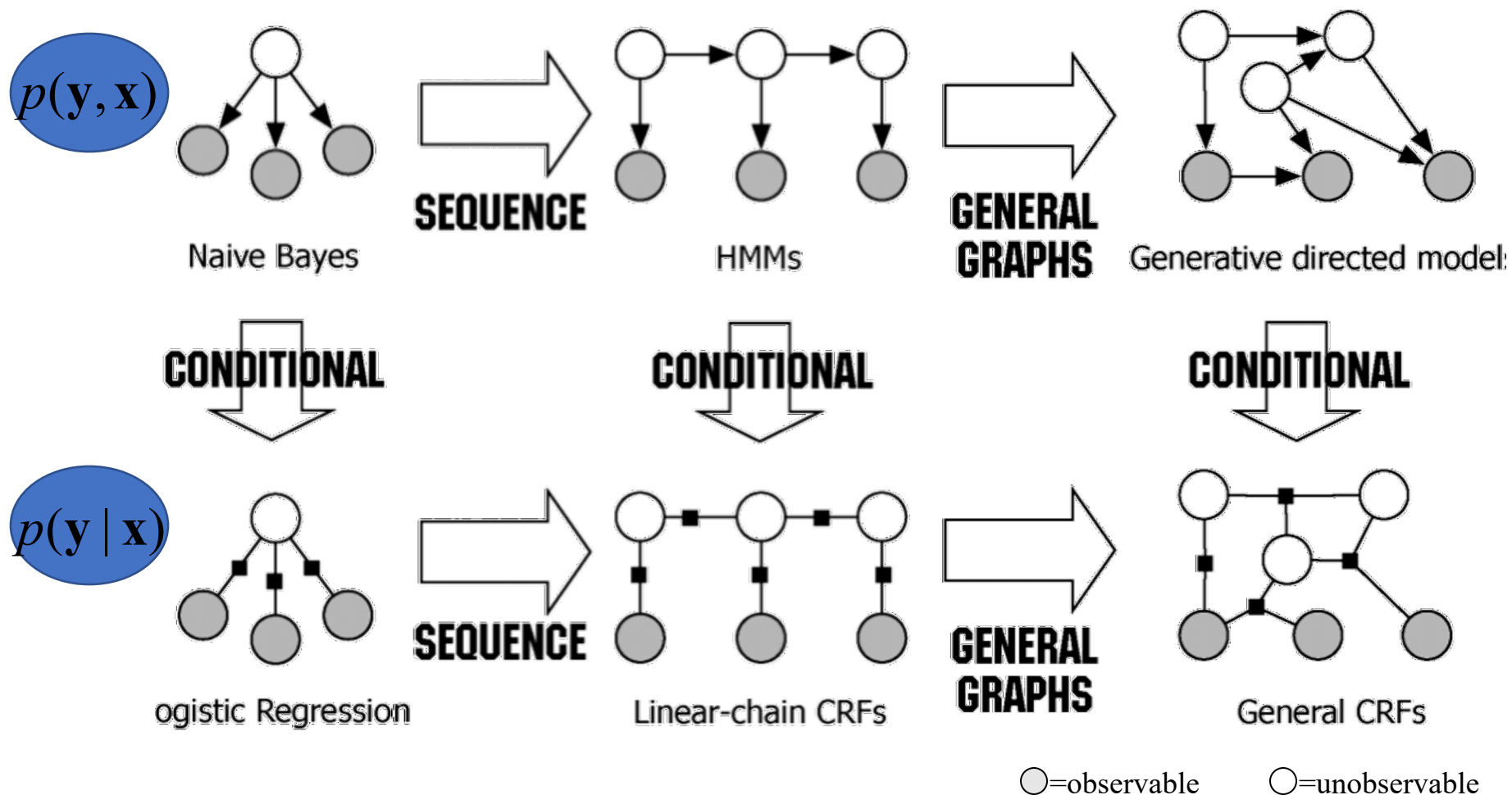
Naive Bayes



Logistic Regression

Both generative models and discriminative models describe distributions over  $(y, x)$ , but they work in different directions.

# Discriminative Vs. Generative



# Sequence prediction

- Like NER: identifying and classifying proper names in text, e.g. China as location; George Bush as people; United Nations as organizations

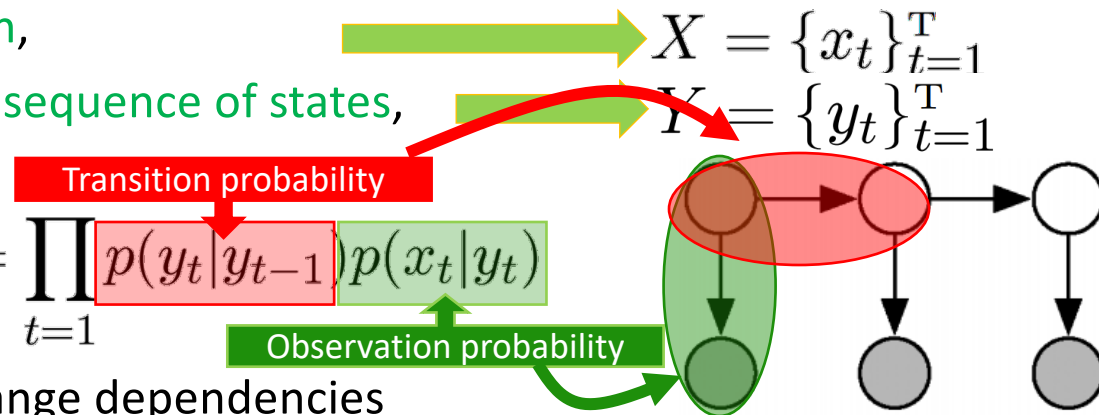
- Set of **observation**,

- Set of **underlying sequence of states**,

- HMM is generative:

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

- Doesn't model long-range dependencies
- Not practical to represent multiple interacting features (hard to model  $p(\mathbf{x})$ )
- The primary advantage of CRFs over hidden Markov models is their conditional nature, resulting in the relaxation of the independence assumptions
- And it can handle overlapping features



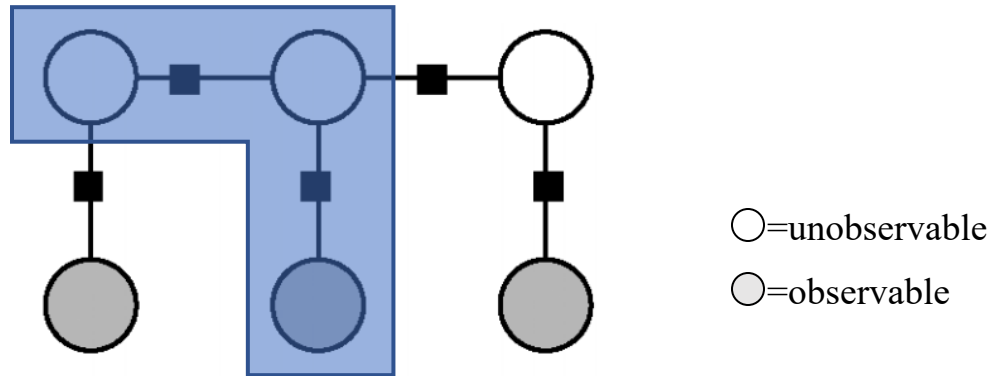
# Chain CRFs

- Each potential function will operate on pairs of adjacent label variables

$$p(\mathbf{y}|\mathbf{x}, \boldsymbol{\lambda}) = \frac{1}{Z(\mathbf{x})} \exp \left( \sum_j \lambda_j F_j(\mathbf{y}, \mathbf{x}) \right)$$

$$F_j(\mathbf{y}, \mathbf{x}) = \sum_{i=1} f_j(y_{i-1}, y_i, \mathbf{x}, i), \quad \text{Feature functions}$$

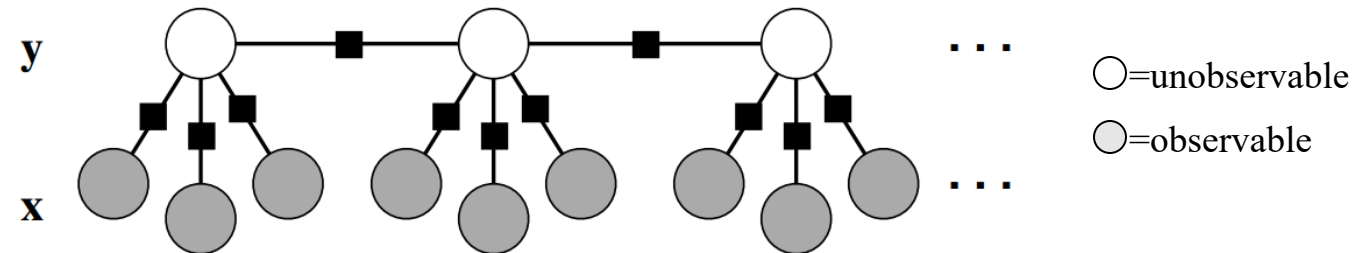
- Parameters to be estimated,  $\lambda_j$



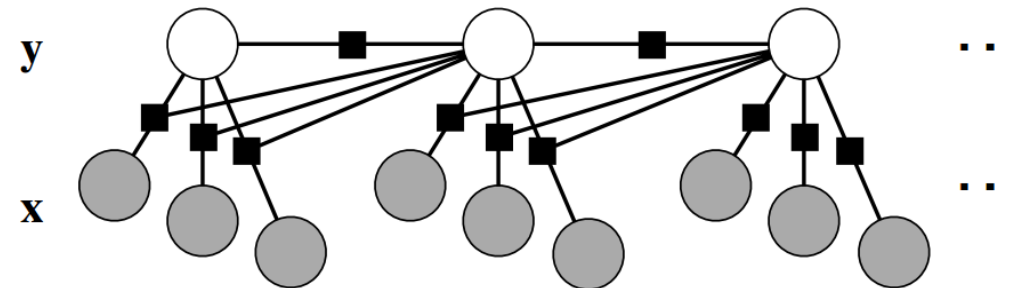


# Chain CRF

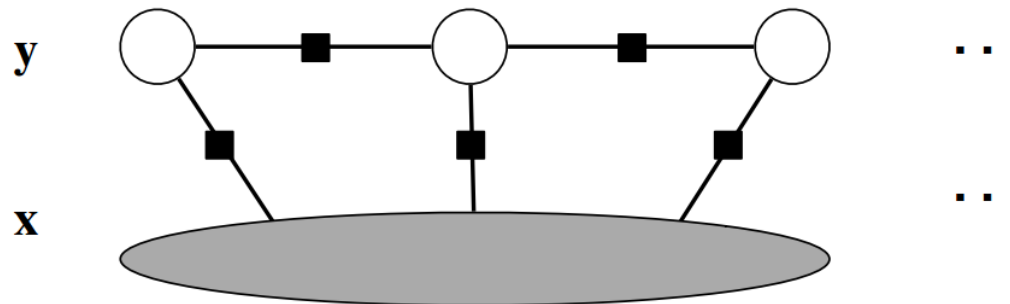
- We can change it so that each state depends on more observations



- Or inputs at previous steps



- Or all inputs

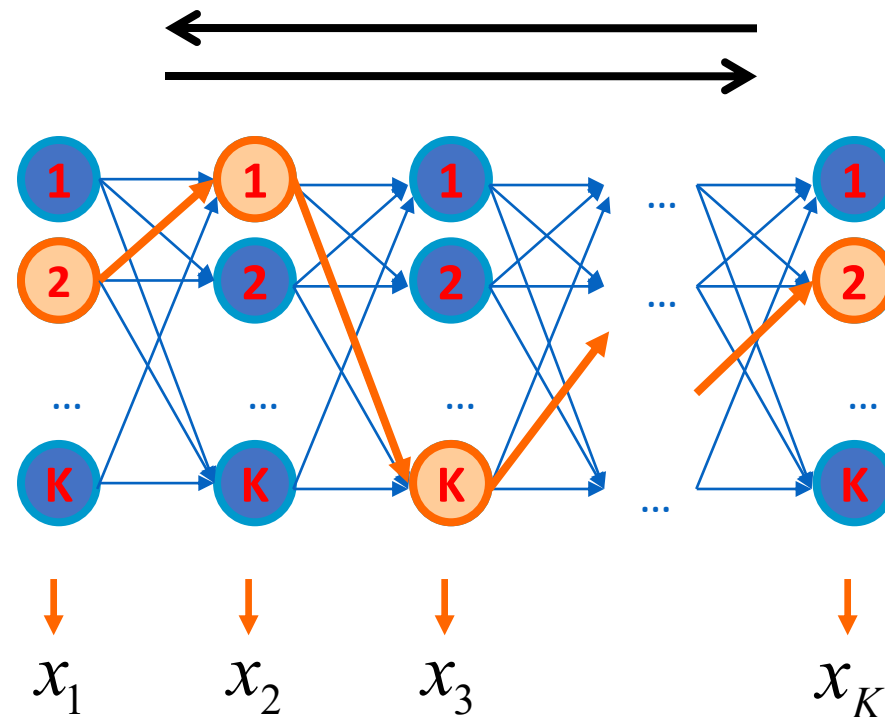


# Outline

- Modeling
- Inference
- Training
- Applications

# Inference in HMM

- Dynamic Programming:
  - Forward
  - Backward
  - Viterbi



# Inference: Chain-CRF

- The inference of linear-chain CRF is very similar to that of HMM
- We can write the marginal distribution:

$$p(Y_{i-1} = y', Y_i = y | \mathbf{x}^{(k)}, \boldsymbol{\lambda}) = \frac{\alpha_{i-1}(y' | \mathbf{x}) M_i(y', y | \mathbf{x}) \beta_i(y | \mathbf{x})}{Z(\mathbf{x})}$$

- Solve **Chain-CRF** using **Dynamic Programming** (Similar to Viterbi)!
- 1. First computing  $\alpha$  for all  $t$  (forward), then compute  $\beta$  for all  $t$  (backward).
- 2. Return the marginal distributions computed.
- 3. Run viterbi to find the optimal sequence