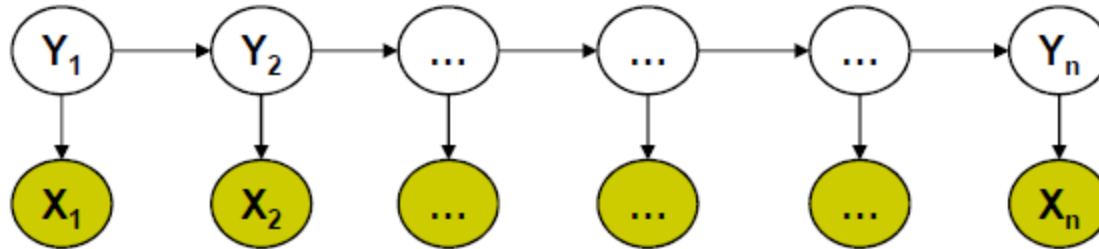# Conditional Random Field

**Dr. Jianlin Cheng**

**Department of Computer Science**
**University of Missouri, Columbia**

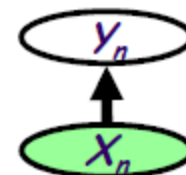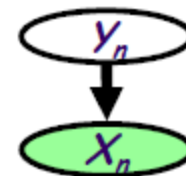**Slides Adapted from Book and CMU, MU, Stanford Machine Learning Courses**
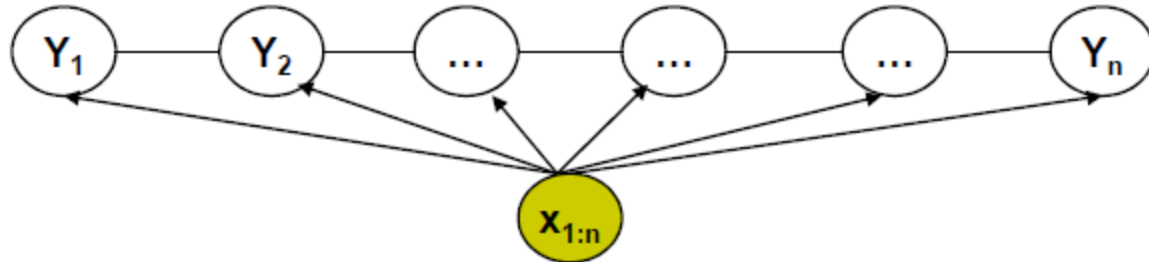**Fall, 2011**

# Shortcomings of HMMs



- 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 $P(Y, X)$, but in a prediction task, we need the conditional probability $P(Y|X)$

# Generative vs. Discriminative Classifiers

- Goal: Wish to learn f: X → Y, e.g., P(Y|X)

- Generative classifiers (e.g., Naïve Bayes):
  - Assume some functional form for P(X|Y), P(Y)
    This is a '*generative*' model of the data!
  - Estimate parameters of P(X|Y), P(Y) directly from training data
  - Use Bayes rule to calculate P(Y|X= x)

- Discriminative classifiers (e.g., logistic regression)
  - Directly assume some functional form for P(Y|X)
    This is a '*discriminative*' model of the data!
  - Estimate parameters of P(Y|X) directly from training data
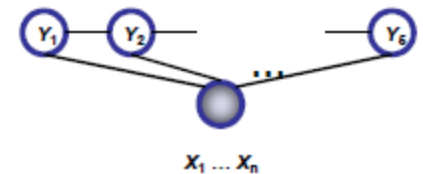
# Structured Conditional Models



- Conditional probability $P$(label sequence **y** | observation sequence **x**) rather than joint probability $P(\mathbf{y}, \mathbf{x})$
  - Specify the probability of possible label sequences given an observation sequence

- Allow arbitrary, non-independent features on the observation sequence $\mathbf{X}$

- The probability of a transition between labels may depend on past and future observations

- Relax strong independence assumptions in generative models
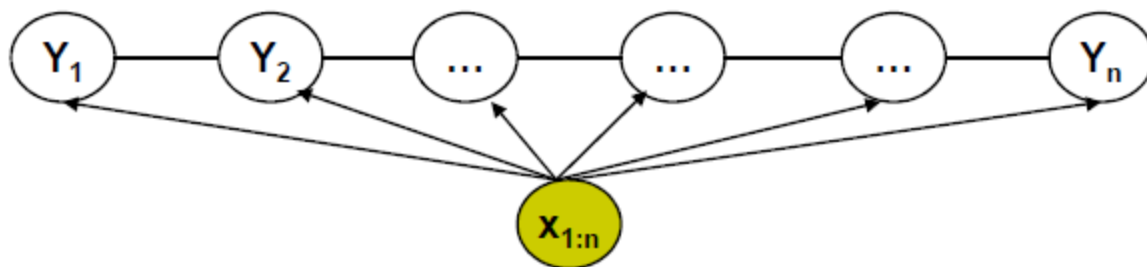
# Conditional Distribution

- If the graph $G = (V, E)$ of $\mathbf{Y}$ is a tree, the conditional distribution over the label sequence $\mathbf{Y} = \mathbf{y}$, given $\mathbf{X} = \mathbf{x}$, by the Hammersley Clifford theorem of random fields is:

$$p_\theta(\mathbf{y} \mid \mathbf{x}) \propto \exp\left( \sum_{e \in E, k} \lambda_k f_k(e, \mathbf{y}|_e, \mathbf{x}) + \sum_{v \in V, k} \mu_k g_k(v, \mathbf{y}|_v, \mathbf{x}) \right)$$

- x is a data sequence
- y is a label sequence
- v is a vertex from vertex set V = set of label random variables
- e is an edge from edge set E over V
- $f_k$ and $g_k$ are given and fixed. $g_k$ is a Boolean vertex feature; $f_k$ is a Boolean edge feature
- k is the number of features
- $\theta = (\lambda_1, \lambda_2, \cdots, \lambda_n; \mu_1, \mu_2, \cdots, \mu_n); \lambda_k$ and $\mu_k$ are parameters to be estimated
- y|$_e$ is the set of components of y defined by edge e
- y|$_v$ is the set of components of y defined by vertex v
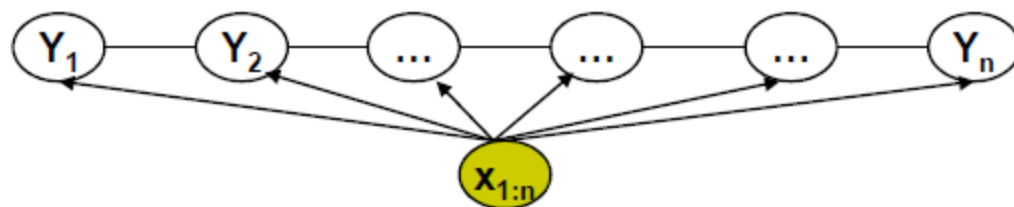
# Conditional Random Fields



$$P(\mathbf{y}_{1:n}|\mathbf{x}_{1:n}) = \frac{1}{Z(\mathbf{x}_{1:n})} \prod_{i=1}^{n} \phi(y_i, y_{i-1}, \mathbf{x}_{1:n}) = \frac{1}{Z(\mathbf{x}_{1:n}, \mathbf{w})} \prod_{i=1}^{n} \exp(\mathbf{w}^T \mathbf{f}(y_i, y_{i-1}, \mathbf{x}_{1:n}))$$

- CRF is a partially directed model
  - Discriminative model
  - Usage of global normalizer $Z(\mathbf{x})$
  - Models the dependence between each state and the entire observation sequence
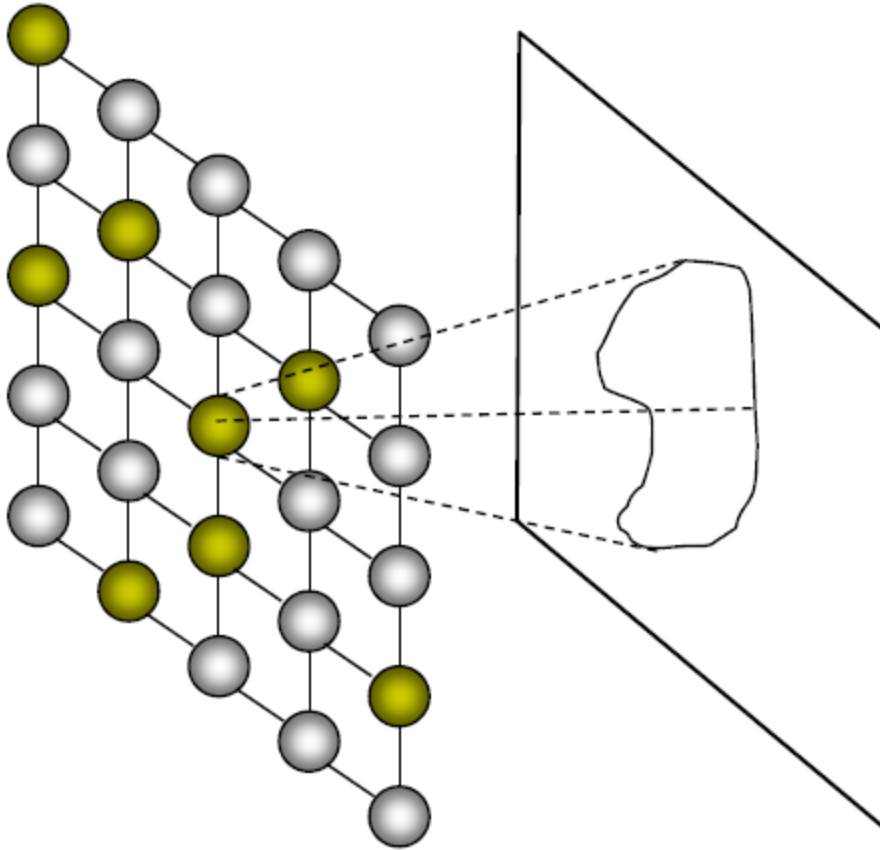
# Conditional Random Field

- General parametric form:



$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x}, \lambda, \mu)} \exp\left(\sum_{i=1}^{n}\left(\sum_{k} \lambda_k f_k(y_i, y_{i-1}, \mathbf{x}) + \sum_{l} \mu_l g_l(y_i, \mathbf{x})\right)\right)$$

$$= \frac{1}{Z(\mathbf{x}, \lambda, \mu)} \exp\left(\sum_{i=1}^{n}\left(\lambda^T \mathbf{f}(y_i, y_{i-1}, \mathbf{x}) + \mu^T \mathbf{g}(y_i, \mathbf{x})\right)\right)$$

$$\text{where } Z(\mathbf{x}, \lambda, \mu) = \sum_{\mathbf{y}} \exp\left(\sum_{i=1}^{n}\left(\lambda^T \mathbf{f}(y_i, y_{i-1}, \mathbf{x}) + \mu^T \mathbf{g}(y_i, \mathbf{x})\right)\right)$$

# Conditional Random Field



$$p_\theta(y \mid x) = \frac{1}{Z(\theta, x)} \exp\left\{\sum_c \theta_c f_c(x, y_c)\right\}$$
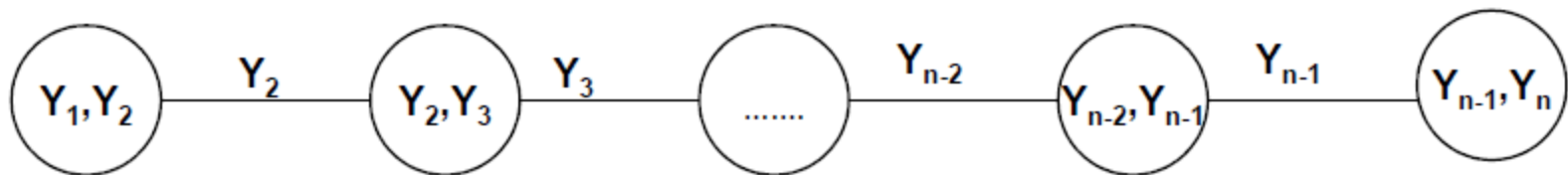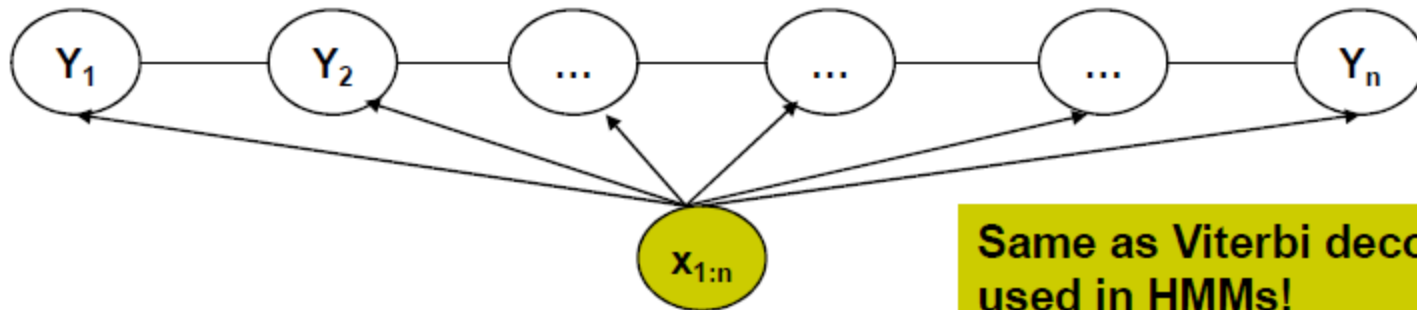
- Allow arbitrary dependencies on input

- Clique dependencies on labels

- Use approximate inference for general graphs

# CRFs Inference

- Given CRF parameters $\lambda$ and $\mu$, find the $\mathbf{y}^*$ that maximizes $P(\mathbf{y}|\mathbf{x})$

$$\mathbf{y}^* = \arg\max_{\mathbf{y}} \exp\left(\sum_{i=1}^{n}(\lambda^T \mathbf{f}(y_i, y_{i-1}, \mathbf{x}) + \mu^T \mathbf{g}(y_i, \mathbf{x}))\right)$$

  - Can ignore $Z(\mathbf{x})$ because it is not a function of y

- Run the max-product algorithm on the junction-tree of CRF:



**Same as Viterbi decoding used in HMMs!**

# CRF Learning

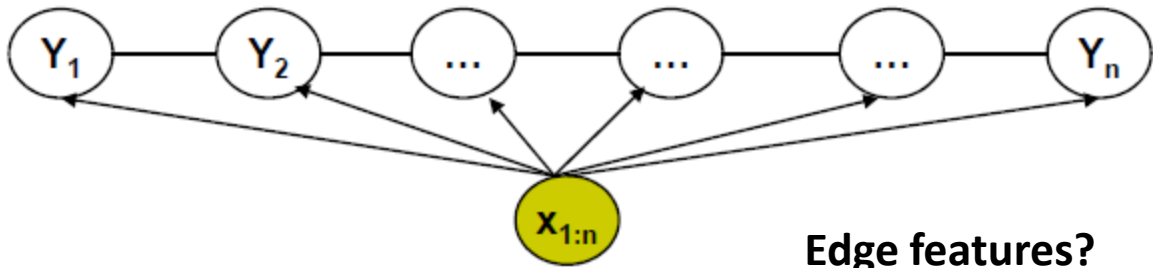- Given $\{(\mathbf{x}_d, \mathbf{y}_d)\}_{d=1}^N$, find $\lambda*, \mu*$ such that

$$
\begin{aligned}
\lambda*, \mu* &= \arg\max_{\lambda,\mu} L(\lambda,\mu) = \arg\max_{\lambda,\mu} \prod_{d=1}^N P(\mathbf{y}_d|\mathbf{x}_d, \lambda, \mu) \\
&= \arg\max_{\lambda,\mu} \prod_{d=1}^N \frac{1}{Z(\mathbf{x}_d, \lambda, \mu)} \exp(\sum_{i=1}^u (\lambda^T \mathbf{f}(y_{d,i}, y_{d,i-1}, \mathbf{x}_d) + \mu^T \mathbf{g}(y_{d,i}, \mathbf{x}_d))) \\
&= \arg\max_{\lambda,\mu} \sum_{d=1}^N (\sum_{i=1}^n (\lambda^T \mathbf{f}(y_{d,i}, y_{d,i-1}, \mathbf{x}_d) + \mu^T \mathbf{g}(y_{d,i}, \mathbf{x}_d)) - \log Z(\mathbf{x}_d, \lambda, \mu))
\end{aligned}
$$

- Computing the gradient w.r.t $\lambda$:

> **Gradient of the log-partition function in an exponential family is the expectation of the sufficient statistics.**
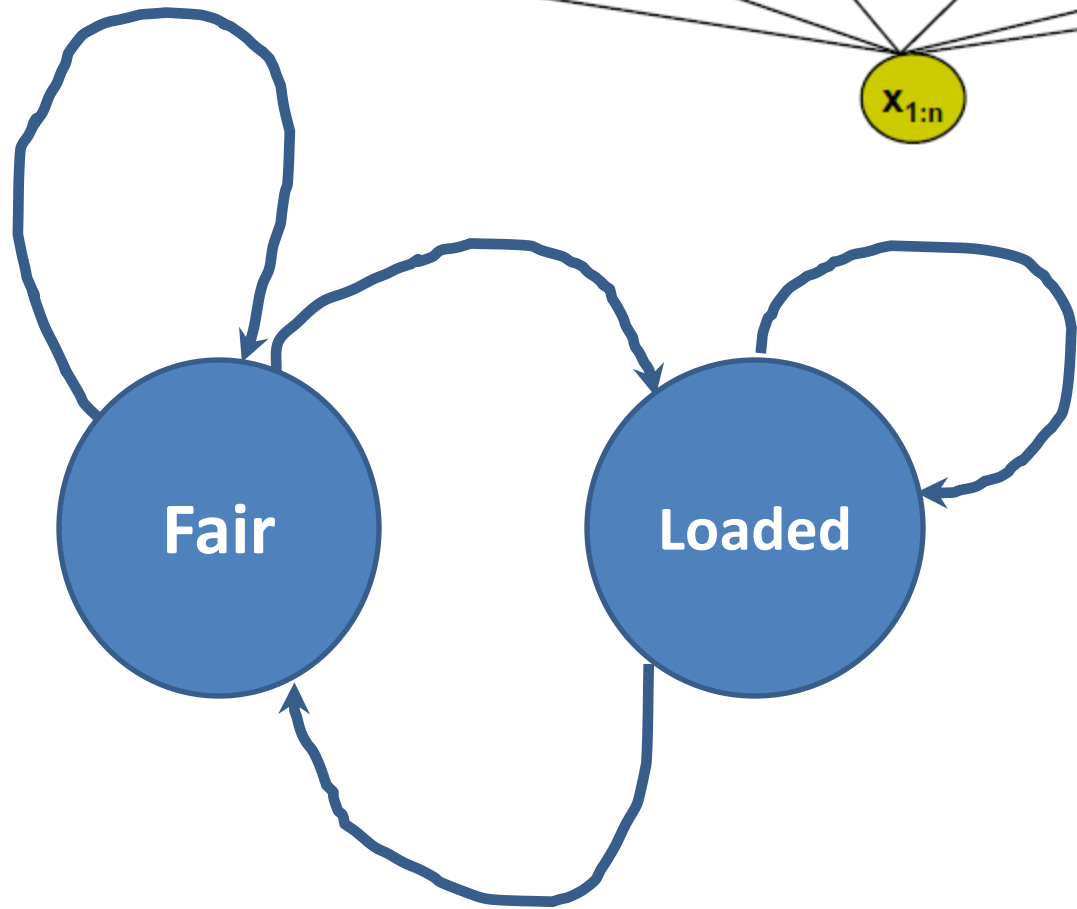
$$
\nabla_\lambda L(\lambda, \mu) = \sum_{d=1}^N (\sum_{i=1}^n \mathbf{f}(y_{d,i}, y_{d,i-1}, \mathbf{x}_d) - \sum_{\mathbf{y}} (P(\mathbf{y}|\mathbf{x_d}) \sum_{i=1}^n \mathbf{f}(y_{d,i}, y_{d,i-1}, \mathbf{x}_d)))
$$

# A CRF Example



Edge features?
Node features?

Output the maximum number of three dice throws

# Edge Features

| $Y_i$ | $Y_{i+1}$ | $X_{i-1}$ | $X_i$ | $X_{i+1}$ |
|-------|-----------|-----------|-------|-----------|
| F | F | 0 | 1 | 1 |
| F | F | 1 | 1 | 1 |
| ... | ... | ... | ... | ... |
| L | F | 1 | 1 | 1 |
| L | F | 2 | 1 | 1 |
| ... | ... | ... | ... | ... |

# Node Features

| $Y_i$ | $X_{i-1}$ | $X_i$ | $X_{i+1}$ |
|-------|-----------|-------|-----------|
| F | 0 | 1 | 1 |
| F | 1 | 1 | 1 |
| ... | ... | ... | ... |
| L | 1 | 1 | 1 |
| L | 2 | 1 | 1 |
| ... | ... | ... | ... |

# Evaluation & Decoding

- **Model parameters**: weights for features

- **Calculate** P(Path | Obs). Y = FFLF. X = 1216.

- **Edge features**: FF (012), FL (121), LF (216)

- **Node features**: F(012), F(121), L(216), F(160)

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x}, \lambda, \mu)} \exp\left(\sum_{i=1}^{n}\left(\sum_{k} \lambda_k f_k(y_i, y_{i-1}, \mathbf{x}) + \sum_{l} \mu_l g_l(y_i, \mathbf{x})\right)\right)$$

$$= \frac{1}{Z(\mathbf{x}, \lambda, \mu)} \exp\left(\sum_{i=1}^{n}\left(\lambda^T \mathbf{f}(y_i, y_{i-1}, \mathbf{x}) + \mu^T \mathbf{g}(y_i, \mathbf{x})\right)\right)$$

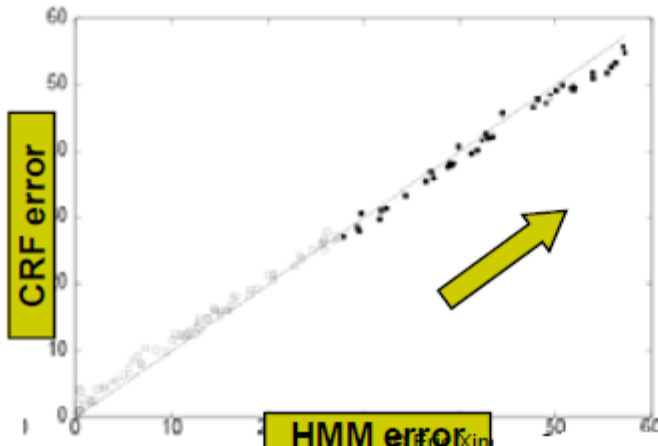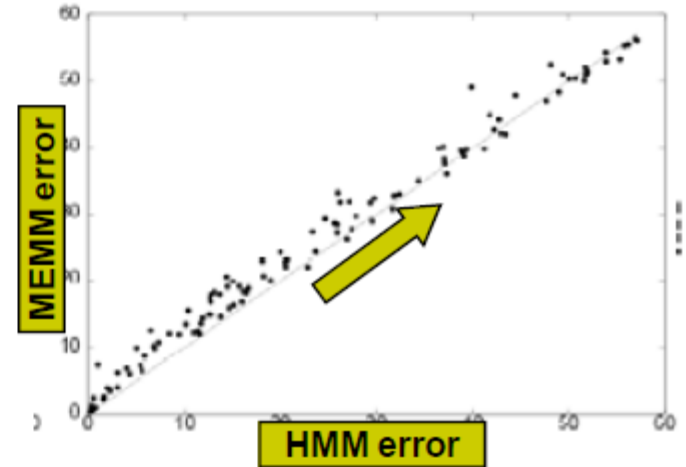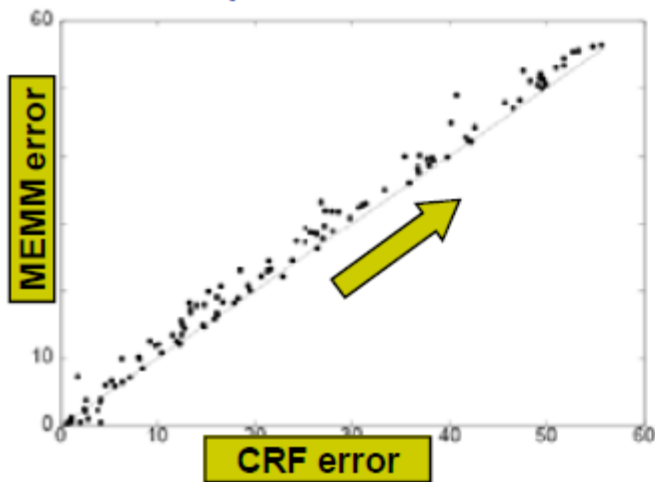# Learning – Gradient Ascend

**Fit the weights of the features**

| $Y_i$ | $Y_{i+1}$ | $X_{i-1}$ | $X_i$ | $X_{i+1}$ |
|---|---|---|---|---|
| F | F | 0 | 1 | 1 |
| F | F | 1 | 1 | 1 |
| … | … | … | … | … |
| L | F | 1 | 1 | 1 |
| L | F | 2 | 1 | 1 |
| … | … | … | … | … |

| $Y_i$ | $X_{i-1}$ | $X_i$ | $X_{i+1}$ |
|---|---|---|---|
| F | 0 | 1 | 1 |
| F | 1 | 1 | 1 |
| … | … | … | … |
| L | 1 | 1 | 1 |
| L | 2 | 1 | 1 |
| … | … | … | … |

$$\nabla_\lambda L(\lambda, \mu) \ = \ \sum_{d=1}^{N} \left( \sum_{i=1}^{n} \mathbf{f}(y_{d,i}, y_{d,i-1}, \mathbf{x}_d) \right) - \sum_{\mathbf{y}} \left( P(\mathbf{y}|\mathbf{x_d}) \sum_{i=1}^{n} \mathbf{f}(y_{d,i}, y_{d,i-1}, \mathbf{x}_d) \right)$$

# Comparison on Synthetic Data

- Comparison of error rates on synthetic data



Data is increasingly higher order in the direction of arrow

**MEMM: maximum entropy Markov models**

# CRFs: Some Empirical Results on Speech Tagging

- Parts of Speech tagging

| model | error | oov error |
|:---:|:---:|:---:|
| HMM | 5.69% | 45.99% |
| MEMM | 6.37% | 54.61% |
| CRF | 5.55% | 48.05% |
| MEMM$^+$ | 4.81% | 26.99% |
| CRF$^+$ | 4.27% | 23.76% |

$^+$Using spelling features

- Using same set of features: HMM >=< CRF > MEMM
- Using additional overlapping features: CRF$^+$ > MEMM$^+$ >> HMM

# More References

- **Collection of papers and tools**: [http://www.inference.phy.cam.ac.uk/hmw26/crf/](http://www.inference.phy.cam.ac.uk/hmw26/crf/)

- **Tutorial**:  H.M. Wallach. Conditional Random Fields: An Introduction

- **Paper**: J. Lafferty, A. McCallum, F. Perreira. Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data