

Learning Graphical Models from Data

- Learning Graph *structure*
- Learning Potentials given fixed graph

Graph Structure



- 1 Manual: Designed by domain expert
 - ▶ Used in applications where dependency structure is well-understood
 - ▶ Example: QMR systems, Kalman filters, Vision (Grids), HMM for speech recognition and IE.
- 2 Learned from examples
 - ▶ NP hard to find the optimal structure.
 - ▶ Widely researched, mostly posed as a branch and bound search problem.
 - ▶ Useful in dynamic situations

Parameters in Potentials

- 1 Manual: Provided by domain expert
 - ▶ Used in infrequently constructed graphs, example QMR systems
 - ▶ Also where potentials are an easy function of the attributes of connected graphs, example: vision networks.
- 2 Learned: from examples
 - ▶ More popular since difficult for humans to assign numeric values
 - ▶ Many variants of parameterizing potentials.
 - 1 Table potentials: each entry a parameter, example, HMMs
 - 2 Potentials: combination of shared parameters and data attributes: example, CRFs.

Learning potentials

Given sample of data ^D generated from a distribution represented by a graphical model with known structure G , learn potentials.

Two settings:

- ① All variables observed or not.
 - ① Fully observed: each training sample has all n variables observed.
 - ② Partially observed: a subset of the variables are observed.
- ② Potentials coupled with a (log-partition function) or not.
 - ① No: Closed form solutions ^{$\log 2$}
 - ② Yes: Potentials attached to arbitrary overlapping subset of variables in a UDGM. Example = edge potentials in a grid graph. Solve using gradient descent kind of iterative algorithms.

DOUBT

Learning potentials: Two settings

1 Generative:

- 1 $P(\mathbf{x}) = P(x_1, \dots, x_n)$ represented as a graphical model G .
- 2 Samples are $D = \{\mathbf{x}^1, \dots, \mathbf{x}^N\}$ $\underline{x^i} = \{x_1^i, x_2^i, \dots, x_n^i\}$
- 3 Potentials to be learned are $\psi_{\underline{c}}(\underline{\mathbf{x}}_c)$

2 Conditional:

- 1 $P(\underline{\mathbf{y}}|\underline{\mathbf{x}}) = P(y_1, \dots, y_n | \underline{\mathbf{x}})$ represented as a graphical model G over the \mathbf{y} variables.
- 2 Training samples are $D = \{(\underline{\mathbf{x}}^1, \mathbf{y}^1), (\underline{\mathbf{x}}^2, \mathbf{y}^2), \dots, (\underline{\mathbf{x}}^N, \mathbf{y}^N)\}$
- 3 Potentials to be learned are $\psi_{\underline{c}}(\underline{\mathbf{y}}_c, \underline{\mathbf{x}}_c)$ $\psi_c(y_c, \underline{\mathbf{x}}_c)$

Most of the topics under learning we will discuss apply equally well to the two settings.

General framework for Parameter learning in graphical models

- Conditional distribution $\Pr(\mathbf{y}|\mathbf{x}, \theta)$, potentials are function of \mathbf{x} and parameters θ to be learned.
- $\mathbf{y} = y_1, \dots, y_n$ forms a graphical model: directed or undirected.
- Undirected:

$$\Pr(y_1, \dots, y_n | \mathbf{x}, \theta) = \frac{\prod_c \psi_c(\mathbf{y}_c, \mathbf{x}, \theta)}{Z_\theta(\mathbf{x})}$$

$F_\theta(\mathbf{y}_c, \mathbf{x}, \theta) = \log \psi_c(\mathbf{y}_c, \mathbf{x}, \theta) \rightarrow = \frac{1}{Z_\theta(\mathbf{x})} \exp(\sum_c F_\theta(\mathbf{y}_c, \mathbf{x}, \theta))$

where $Z_\theta(\mathbf{x}) = \sum_{\mathbf{y}'} \exp(\sum_c F_\theta(\mathbf{y}'_c, \mathbf{x}, \theta))$
clique potential $\psi_c(\mathbf{y}_c, \mathbf{x}) = \exp(F_\theta(\mathbf{y}_c, \mathbf{x}, \theta))$

Forms of $F_{\theta}(\mathbf{y}_c, c, \mathbf{x})$

- Log-linear model over user-defined features. E.g. CRFs, Maxent models, etc.

Let K be number of features. Denote a feature as $f_k(\mathbf{y}_c, c, \mathbf{x})$.
Then,

$$F_{\theta}(\mathbf{y}_c, c, \mathbf{x}) = \sum_{k=1}^K \theta_k f_k(\mathbf{y}_c, c, \mathbf{x})$$

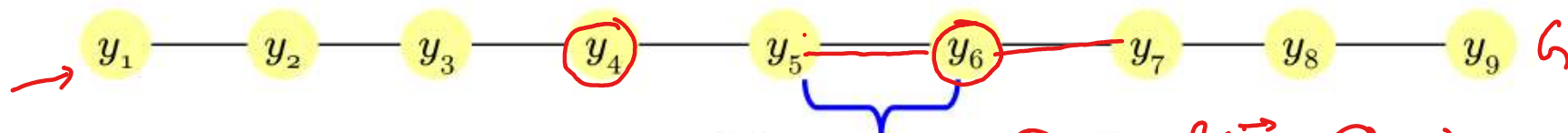
- Arbitrary function, e.g. a neural network that takes as input $\mathbf{y}_c, c, \mathbf{x}$ and transforms them possibly non-linearly into a real value. θ are the parameters of the network.

Example: Named Entity Recognition

My review of Fermat's last theorem by S. Singh

$y_i \in \{ \text{Title}, \text{Author}, \text{Other} \}$

t	1	2	3	4	5	6	7	8	9
x	My	review	of	Fermat's	last	theorem	by	S.	Singh
y	Other	Other	Other	Title	Title	Title	other	Author	Author



$f(y_i, y_{i-1}, i, x)$ $f(\vec{y}_c, c, x)$
 $c = (i-1, i)$
 (y_i, y_{i-1})

Features decompose over adjacent labels.

$$f(x, y) = \sum_{i=1}^{|x|} f(y_i, y_{i-1}, i, x)$$

Named Entity Recognition: Features

- Feature vector for each position

$$\mathbf{f}(y_i, \mathbf{x}, i, y_{i-1})$$

User provided

i-th label

Word i & neighbors

previous label

- Examples

$$f_2(y_i, \mathbf{x}, i, y_{i-1}) = 1 \text{ if } y_i \text{ is Person \& } x_i \text{ is Douglas}$$

$$f_3(y_i, \mathbf{x}, i, y_{i-1}) = 1 \text{ if } y_i \text{ is Person \& } y_{i-1} \text{ is Other}$$

$$\rightarrow f_{k+1}(y_i, y_{i-1}, i, \mathbf{x}) = 1 \text{ if } y_i = y_{i-1} \text{ and } 0 \text{ otherwise.}$$

$$\vec{f}_k(y_i, y_{i-1}, i, \mathbf{x}) = \text{embedding of the } i^{\text{th}} \text{ word from BERT.}$$



Training

Given

- N input output pairs $D = \{(\mathbf{x}^1, \mathbf{y}^1), (\mathbf{x}^2, \mathbf{y}^2), \dots, (\mathbf{x}^N, \mathbf{y}^N)\}$
- Form of F_θ
- Learn parameters θ by maximum likelihood.

$$\max_{\theta} LL(\theta, D) = \max_{\theta} \sum_{i=1}^N \log \Pr(\mathbf{y}^i | \mathbf{x}^i, \theta)$$

Training undirected graphical model

$$\begin{aligned}\underline{LL(\theta, D)} &= \sum_{i=1}^N \log \Pr(\mathbf{y}^i | \mathbf{x}^i, \theta) \\ &= \sum_{i=1}^N \log \frac{1}{Z_{\theta}(\mathbf{x}^i)} \exp\left(\sum_c F_{\theta}(\mathbf{y}_c^i, c, \mathbf{x}^i)\right) \\ &= \sum_i \left[\sum_c \underline{F_{\theta}(\mathbf{y}_c^i, c, \mathbf{x}^i)} - \log Z_{\theta}(\mathbf{x}^i) \right]\end{aligned}$$

The first part is easy to compute but the second term requires to invoke an inference algorithm to compute $Z_{\theta}(\mathbf{x}^i)$ for each i .

Computing the gradient of the above objective with respect to θ also requires inference.

Training via gradient descent

Assume log-linear models like in CRFs where $F_{\theta}(\mathbf{y}_c^i, c, \mathbf{x}^i) = \theta \cdot \mathbf{f}(\mathbf{x}^i, \mathbf{y}_c^i, c)$ Also, for brevity write $\mathbf{f}(\mathbf{x}^i, \mathbf{y}^i) = \sum_c \mathbf{f}(\mathbf{x}^i, \mathbf{y}_c^i, c)$

Handwritten note: $\rightarrow = \sum_{k=1}^K \theta_k f_k(\mathbf{y}_c^i, c, \mathbf{x}^i)$

$$LL(\theta) = \sum_i \log \Pr(\mathbf{y}^i | \mathbf{x}^i, \theta) = \sum_i (\theta \cdot \mathbf{f}(\mathbf{x}^i, \mathbf{y}^i) - \log Z_{\theta}(\mathbf{x}^i))$$

Add a regularizer to prevent over-fitting.

$$\max_{\theta} \sum_i (\theta \cdot \mathbf{f}(\mathbf{x}^i, \mathbf{y}^i) - \log Z_{\theta}(\mathbf{x}^i)) - \|\theta\|^2 / C$$

Concave in θ \Rightarrow gradient descent methods will work.

Handwritten notes:

$$\log Z_{\theta}(\mathbf{x}^i) = \log \sum_{\mathbf{y}} \exp(\theta \cdot \mathbf{f}(\mathbf{x}^i, \mathbf{y}))$$

$$\nabla_{\theta} \log Z_{\theta}(\mathbf{x}^i) = \frac{\sum_{\mathbf{y}} \exp(\theta \cdot \mathbf{f}(\mathbf{x}^i, \mathbf{y})) \mathbf{f}(\mathbf{x}^i, \mathbf{y})}{Z_{\theta}(\mathbf{x}^i)}$$

Gradient of the training objective

$\#(y') = 3$
 $\exp(\sum_c \theta \cdot f(x^i, y'_c))$

$$\nabla L(\theta) = \sum_i \mathbf{f}(\mathbf{x}^i, \mathbf{y}^i) - \frac{\sum_{\mathbf{y}'} \mathbf{f}(\mathbf{x}^i, \mathbf{y}') \exp(\theta \cdot \mathbf{f}(\mathbf{x}^i, \mathbf{y}'))}{Z_\theta(\mathbf{x}^i)} - 2\theta / C$$

$$= \sum_i \mathbf{f}(\mathbf{x}^i, \mathbf{y}^i) - \sum_{\mathbf{y}'} \mathbf{f}(\mathbf{x}^i, \mathbf{y}') \Pr(\mathbf{y}' | \theta, \mathbf{x}^i) - 2\theta / C$$

$$= \sum_i \mathbf{f}(\mathbf{x}^i, \mathbf{y}^i) - E_{\Pr(\mathbf{y}' | \theta, \mathbf{x}^i)} \mathbf{f}(\mathbf{x}^i, \mathbf{y}') - 2\theta / C$$

$$E_{\Pr(\mathbf{y}' | \theta, \mathbf{x}^i)} f_k(\mathbf{x}^i, \mathbf{y}') = \sum_{\mathbf{y}'} f_k(\mathbf{x}^i, \mathbf{y}') \Pr(\mathbf{y}' | \theta, \mathbf{x}^i)$$

$$= \sum_{\mathbf{y}'} \sum_c f_k(\mathbf{x}^i, \mathbf{y}'_c, c) \Pr(\mathbf{y}' | \theta, \mathbf{x}^i)$$

$$= \sum_c \sum_{\mathbf{y}'_c} f_k(\mathbf{x}^i, \mathbf{y}'_c, c) \Pr(\mathbf{y}'_c | \theta, \mathbf{x}^i)$$

$$= \sum_c \sum_{\mathbf{y}'_c} f_k(\mathbf{x}^i, \mathbf{y}'_c, c) \left\{ \sum_{\mathbf{y}'_c} \Pr(\mathbf{y}'_c | \theta, \mathbf{x}^i) \right\}$$