

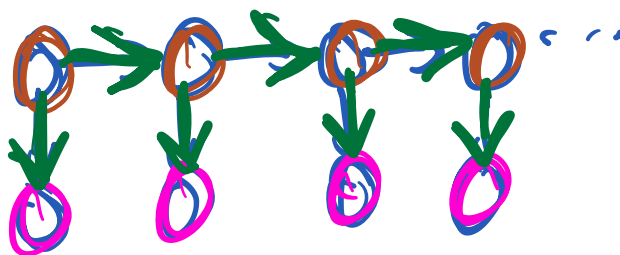
$P(x, y) \rightarrow$ generative model

$$P(y|x) = \frac{P(x, y)}{P(x)}$$

discriminative model

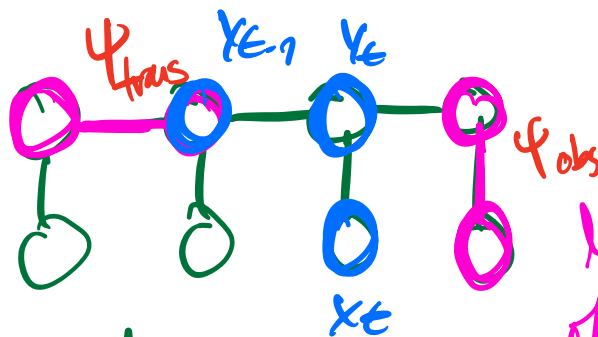
$P(x|y)$

Hint:
Bayes' network



idea:
 \rightarrow MRF

Markov Random Field



clustering
of size 2.

Instead of conditional distributions, we'll use
potential functions $\psi \rightarrow$ defined on cliques

$$P(x|y) = \frac{1}{Z} \prod_t \psi_{\text{trans}}(y_t, y_{t-1}) \cdot \prod \psi_{\text{obs}}(x_t, y_t)$$



↑
potential functions

⇒ non-negative !!!

$$Z(x|y) = \sum_y \prod_t \psi_{\text{trans}}(y_t, y_{t-1}) \cdot \prod \psi_{\text{obs}}(x_t, y_t)$$

compare Bayes Theorem:

$$P(y|x) = \frac{P(x|y)P(y)}{\sum_y P(x|y)P(y)}$$

How to implement the potential functions ψ ?

⇒ log-linear combinations of basis functions

$$\psi_{\text{obs}}(x_t, y_t) = \exp \left\{ \sum_{d=1}^{D_0} w_d \phi_d^0(x_t, y_t) \right\}$$

$$\begin{aligned} \psi_{\text{trans}}(y_t, y_{t-1}) &= \exp \left\{ \sum_{d=1}^{D_1} w_d \phi_d^1(y_t, y_{t-1}) \right\} \\ &= \exp \left\{ \mathbf{w}^T \boldsymbol{\phi}(y_t, y_{t-1}) \right\} \end{aligned}$$

inner product

$$\begin{pmatrix} w_1 \\ \vdots \\ w_d \end{pmatrix}$$

$$\begin{pmatrix} \phi_1^t(\dots) \\ \phi_2^t(\dots) \\ \vdots \\ \phi_{j^*}^t(\dots) \end{pmatrix}$$

\Rightarrow example :

$$\phi_1^t(\underline{x_t}, \underline{y_{t-1}}) = \begin{cases} 1: & y_t = N \wedge y_{t-1} = A \\ 0: & \text{otherwise} \end{cases}$$

trans. \rightarrow

$$\phi_2^t(\underline{x_t}, \underline{y_{t-1}}) = \begin{cases} 1: & y_t = N \wedge y_{t-1} = N \\ 0: & \text{otherwise} \end{cases}$$

also $\rightarrow \phi_{1234}^0(x_t, y_t) = \# \text{ class in } x_t$

$$\phi_{245}^0(x_t, y_t) = \begin{cases} 1: & x_t \text{ contains a number} \\ 0: & \text{otherwise} \end{cases}$$

$$\phi_{1111}^0(x_t, y_t) = \begin{cases} 1: & x_t \text{ is capitalized and } y_t \text{ is a N} \\ 0: & \text{otherwise} \end{cases}$$

$$\Rightarrow P(x|y) = \frac{1}{Z(x|y)} \prod_t \exp(w_0^T \phi_0(x_t, y_t))$$

$$\prod_{\underline{b}} \exp(w_{\underline{t}}^T \phi_{\underline{t}}(y_{\underline{b}}, y_{\underline{b}-1}))$$

$$= \frac{1}{Z(x|y)} \exp\left(\sum_{\underline{b}} w_0^T \phi_0(x_{\underline{b}}, y_{\underline{b}})\right) \times$$

ok

$$\exp \left(\sum_{\mathbf{b}} \mathbf{w}_{\mathbf{b}}^T \phi_{\mathbf{b}}(y_{\mathbf{b}}, y_{\mathbf{b}-1}) \right) \quad \uparrow$$

Example: Det — N — V — Det — N
 | | | | |
 The cat closes the dog

$$\phi_{\mathbf{b}}(N, \text{cat}) = \begin{pmatrix} x = \text{dog} ? \\ x = \text{cat} ? \\ x = \text{cat} \wedge y = V ? \\ x = \text{cat} \wedge y = N ? \\ \vdots \end{pmatrix} \Rightarrow \begin{pmatrix} 0 \\ 1 \\ 0 \\ 1 \\ \vdots \end{pmatrix}$$

$$P(x, y) = \frac{1}{Z(x, y)} \cdot \exp \left(\sum_{\mathbf{b}} \mathbf{w}_{\mathbf{b}}^T \phi_{\mathbf{b}}(x_{\mathbf{b}}, y_{\mathbf{b}}) \right)$$

$$\in \mathbb{R}^{D_0 \times D_e}$$

$$\Rightarrow P(x, y) = \frac{1}{Z(x, y)} \exp \left[\mathbf{w}^T \begin{pmatrix} \sum_{\mathbf{b}} \phi_{\mathbf{b}}(x_{\mathbf{b}}, y_{\mathbf{b}}) \\ \sum_{\mathbf{b}} \psi_{\mathbf{b}}(y_{\mathbf{b}}, y_{\mathbf{b}-1}) \end{pmatrix} \right]$$

$$\mathbf{w} = \begin{pmatrix} w_0 \\ w_e \end{pmatrix}$$

$$\Phi(x, y)$$

$$\begin{pmatrix} \phi_{01}(x_e, y_e) \\ \vdots \\ \phi_{0D_e}(x_e, y_e) \\ \phi_{e1}(y_e, y_{e-1}) \\ \vdots \\ \phi_{eD_e}(y_e, y_{e-1}) \end{pmatrix}$$

joint input
output space!!

↓
feature representation of
input X and output Y

→ move to conditional model: discard function
in z over inputs

$$\Rightarrow z = \sum_y \exp(\mathbf{w}^T \Phi(x, y))$$

$$\Rightarrow \text{model: } P(Y|X) = \frac{1}{z(X)} \exp(\mathbf{w}^T \Phi(x, y))$$

⇒ following this approach leads to
conditional random fields (CRFs)
Lafferty et al

⇒ or simplify $P(Y|X)$ by

– MAP (maximum a posteriori) approach

$$\Rightarrow \operatorname{argmax}_y P(Y|X)$$

$$= \operatorname{argmax}_y \frac{1}{Z(x)} \sup \{ w^T \Phi(x|y) \}$$

$$= \operatorname{argmax}_y \underline{w^T \Phi(x|y)}$$

perceptron \leftrightarrow primal / dual
probably Collins & Duffy (2002)

$$f(x) = w^T x$$

$$w \leq 0$$

$$w \leq w + \underline{\gamma x}$$

$$w \leq 0$$

$$w \leq 0 + \gamma_1 x_1$$

$$w \leq 0 + \gamma_1 x_1 + \gamma_2 x_2$$

$$\underline{w \leq 0 + \gamma_1 x_1 + \gamma_2 x_2 + \gamma_3 x_3}$$

dual perceptron:

$$\alpha_1, \alpha_2, \dots, \alpha_N \leq 0$$

$$\text{update: } \alpha_n \leftarrow \alpha_n + 1$$

$$\text{model: } w = \sum \alpha_n \gamma_n x_n$$

putting things together: $f(x) = w^T \underline{x}$

$$= \sum d_n \gamma_n \underline{x_n^T} \underline{x}$$

$$= \sum d_n \gamma_n K(x_n, x)$$