

$$P(X) = \frac{1}{Z(\theta)} e^{\frac{F(X, \theta)}{Z(\theta)}}$$

$$\frac{\partial}{\partial \theta_j} \ln \ell(\theta) = \frac{\partial}{\partial \theta_j} \ln \left( \frac{1}{Z(\theta)} e^{\frac{F(X, \theta)}{Z(\theta)}} \right)$$

$$\Rightarrow \nabla \ln \ell(\theta) = \frac{\partial}{\partial \theta} \ln \ell(\theta)$$

$$\left[ \begin{array}{c} \frac{\partial}{\partial \theta_1} \ln \ell(\theta) \\ \frac{\partial}{\partial \theta_2} \ln \ell(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_p} \ln \ell(\theta) \end{array} \right]$$

$$= \frac{\partial}{\partial \theta} \ln \ell(\theta)$$

$$\theta_1 \quad \theta_2$$

gradient ascent

$$\frac{\partial}{\partial \theta_j} \ln \ell(\theta) = m E_P \left\{ \frac{\partial}{\partial \theta_j} F(X, \theta) \right\} - m E_{P_\theta(X)} \left\{ \frac{\partial}{\partial \theta_j} F(X, \theta) \right\}$$

$t \leftarrow 0$   
 $\theta_0$   
 while not converged

$$\frac{\partial}{\partial \theta} \ln \ell(\theta) \Big|_{\theta = \theta_t}$$

$$\theta_{t+1} = \theta_t + \lambda \frac{\partial}{\partial \theta} \ln \ell(\theta) \Big|_{\theta = \theta_t}$$

$t \leftarrow t+1$





problem: How to compute  $E_{P_{\theta}(X)} \left\{ \frac{\partial}{\partial \theta_j} F(X, \theta) \right\}$

take samples  $X^1, X^2, \dots, X^m$  from  $P_{\theta}(X) = \frac{1}{Z(\theta)} e^{F(X, \theta)}$

$$E_{P_{\theta}(X)} \left\{ \frac{\partial}{\partial \theta_j} F(X, \theta) \right\} \approx \frac{1}{M} \sum_{i=1}^m \frac{\partial}{\partial \theta_j} F(X^i, \theta)$$

$$\frac{\partial}{\partial \theta_j} \ell(\theta) = \sum_{i=1}^m \frac{\partial}{\partial \theta_j} F(X^i, \theta) - \frac{m}{M} \sum_{i=1}^m \frac{\partial}{\partial \theta_j} F(X^i, \theta)$$

training data  $X^1, X^2, \dots, X^m$

samples from  $P_{\theta}(X)$   $X^1, X^2, \dots, X^m$

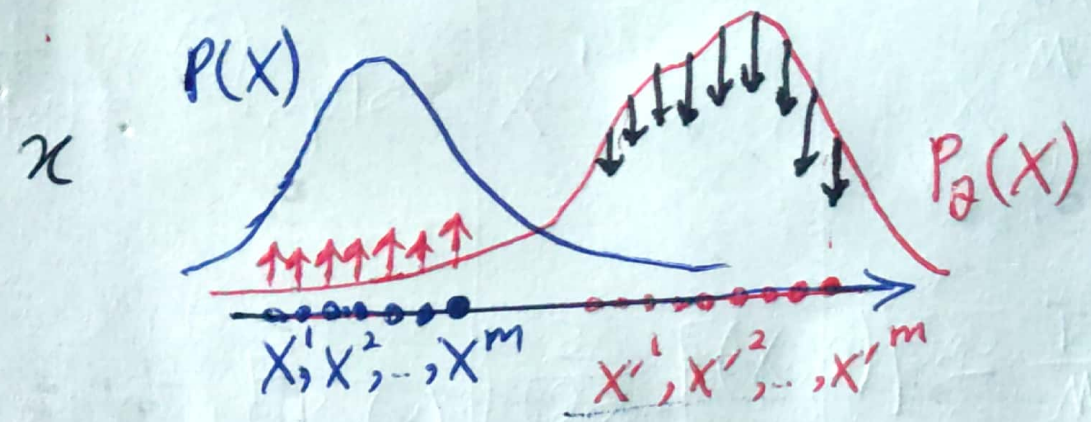
usually  
 $M=m$

$$\frac{\partial}{\partial \theta_j} \ell(\theta) = \underbrace{\sum_{i=1}^m \frac{\partial F}{\partial \theta_j} F(X^i, \theta)}_{\text{positive force}} - \underbrace{\sum_{i=1}^m \frac{\partial F}{\partial \theta_j} F(X^i, \theta)}_{\text{negative force}}$$

positive force

negative force





Contrastive divergence (CD)

$$CRF = P_{\theta}(Y|X) = \frac{1}{Z(\theta, X)} e^{F_{\theta}(X, Y)}$$

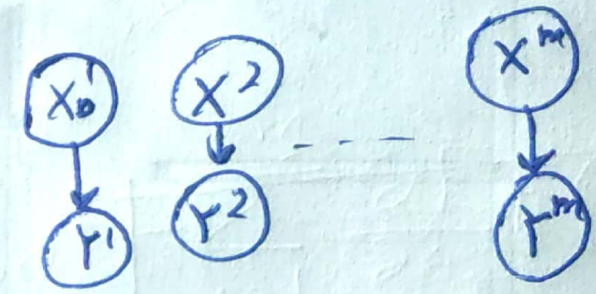
Conditional likelihood

Data  
 $(X^1, Y^1), (X^2, Y^2), \dots, (X^m, Y^m)$

$$cLL(\theta) = Pr(Y^1, Y^2, \dots, Y^m | X^1, X^2, \dots, X^m)$$

$$\prod_{i=1}^m Pr(Y^i | X^1, \dots, X^m)$$

$$\prod_{i=1}^m Pr(Y^i | X^i) \Rightarrow cLL(\theta) = \prod_{i=1}^m P_{\theta}(Y^i | X^i)$$





$$\begin{aligned}
 \text{c}ll(\theta) = \log \text{cl}(\theta) &= \sum_{i=1}^m \log P_{\theta}(Y^i | X^i) \\
 &= \sum_{i=1}^m \log \frac{1}{Z_{\theta}(X^i)} e^{F_{\theta}(X^i, Y^i)} \\
 &= - \sum_{i=1}^m \log Z_{\theta}(X^i) + \sum_{i=1}^m F_{\theta}(X^i, Y^i)
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial}{\partial \theta_j} \text{c}ll(\theta) &= - \sum_{i=1}^m \sum_Y P_{\theta}(Y | X^i) F_{\theta}(X^i, Y) + \sum_{i=1}^m F_{\theta}(X^i, Y^i) \\
 &= - \sum_{i=1}^m \underbrace{E_{P_{\theta}(Y|X^i)} \{ F_{\theta}(X^i, Y) \}}_{\text{Needs inference per } X^i} + \sum_{i=1}^m E_D \{ F_{\theta}(X, Y) \}
 \end{aligned}$$

$$\begin{aligned}
 \text{c}ll(\theta) = \log \text{cl}(\theta) &= \sum_{i=1}^m \log P_{\theta}(Y^i | X^i) \\
 &= \sum_{i=1}^m \log \frac{1}{Z_{\theta}(X^i)} e^{F_{\theta}(X^i, Y^i)} \\
 &= -\sum_{i=1}^m \log Z_{\theta}(X^i) + \sum_{i=1}^m F_{\theta}(X^i, Y^i)
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial}{\partial \theta_j} \text{c}ll(\theta) &= -\sum_{i=1}^m \sum_Y P_{\theta}(Y | X^i) F_{\theta}(X^i, Y) + \sum_{i=1}^m F_{\theta}(X^i, Y^i) \\
 &= -\sum_{i=1}^m \underbrace{E_{P_{\theta}(Y|X^i)} \left\{ F_{\theta}(X^i, Y) \right\}}_{\text{Needs inference per } X^i} + \sum_{i=1}^m E_D \left\{ F_{\theta}(X, Y) \right\}
 \end{aligned}$$



$$F(\theta) = \theta \left[ f_1(X_1, X_2) + f_2(X_2, X_3) + f_3(X_3, X_4) \right]$$

$$E \left\{ \frac{\partial F}{\partial \theta} \right\} = E \left\{ f_1(X_1, X_2) + f_2(X_2, X_3) + f_3(X_3, X_4) \right\}$$
$$= E \left\{ f_1(X_1, X_2) \right\} + E \left\{ f_2(X_2, X_3) \right\} + E \left\{ f_3(X_3, X_4) \right\}$$