

Nonnegative matrix factorization (NMF) (review)

• How to approximate $X \approx VY$ for $X, V, Y \geq 0$?

• Approximation error:

$$E_{KL} = \sum_{in} \left[X_{in} \log \frac{X_{in}}{(VY)_{in}} - X_{in} + (VY)_{in} \right]$$

• Auxiliary function:

$$G(V, W) = \sum_{in} \left[X_{in} \log \frac{X_{in}}{(VY)_{in}} - X_{in} + (VY)_{in} - X_{in} \sum_{\alpha} \left(\frac{W_{\alpha} Y_{\alpha in}}{(VY)_{in}} \log \frac{V_{i\alpha}}{W_{i\alpha}} \right) \right]$$

↖ p_{α} for Jensen's inequality

• Properties:

(i) $G(V, V) = E_{KL}(V)$

(ii) $G(V, W) \geq E_{KL}(V)$

• Update:

$$V_{t+1} = \underset{V}{\operatorname{argmin}} G(V, V_t)$$

• Jensen's inequality

- interchanges log-sum and sum-log operations

- replaces $\log(\sum v_{i\alpha} y_{\alpha in})$ in cost function by

$$\sum p_{\alpha} \log v_{i\alpha} \text{ in auxiliary function}$$

- simplifies optimization

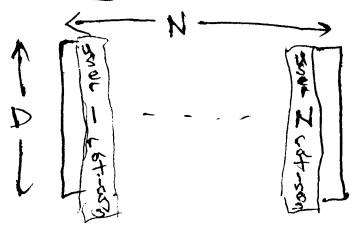
$\partial E_{KL} / \partial v_{i\alpha} = 0$ cannot be solved in closed form.

$\partial G / \partial v_{i\alpha} = 0$ can be solved in closed form

- also used to derive EM in discrete BNs

Binary matrix factorization

Ex: movie-user data
 $X_{ij} \in \{0,1\}$ binary rating
 $i = 1 \dots D$ movies
 $j = 1 \dots N$ users



How to fill in missing matrix elements?

Netflix challenge: \$1 million

Matrix has some missing elements:

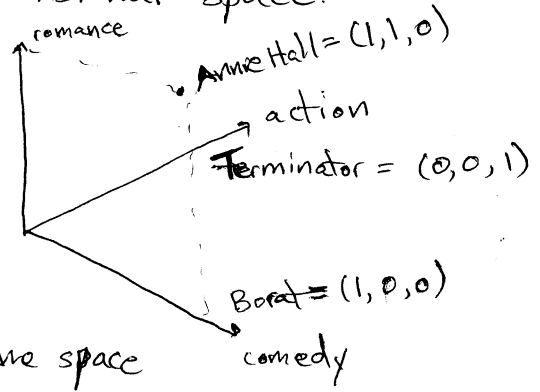
$X_{ij} \in \{0,1\}$ rating is absent / present

Logistic PCA

• Intuition

Represent movies as points in $d \ll D$ dimensional space.

$\vec{v}_i \in \mathbb{R}^d$ (i th movie) Ex: $d=3$



Represent user preferences as points in same space

$\vec{y}_j \in \mathbb{R}^d$ (j th user)

$\vec{y}_j \cdot \hat{e}_\alpha =$ user's preference for movie genre α ($\alpha = 1$ comedy, $= 2$ romance, $= 3$ action)

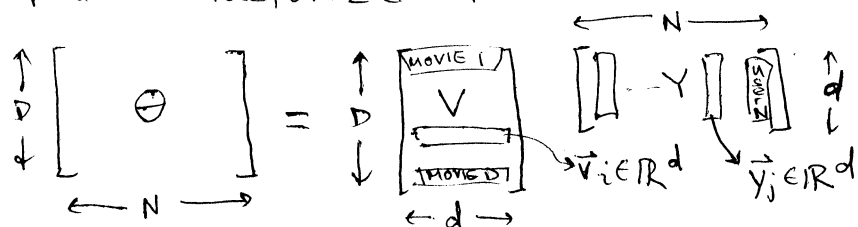
Model ratings as binary random variables: $X_{ij} \in \{0,1\}$

~~Fitting~~ Fit ratings by logistic model:

$$P(X_{ij} = 1) = \sigma(\vec{v}_i \cdot \vec{y}_j) \text{ where } \sigma(z) = \frac{1}{1 + e^{-z}}$$

Use predictions to fill in matrix elements.

• Matrix factorization



$$\Theta = VY \text{ (exact factorization)}$$

$$P(X_{ij}=1) = \sigma(\Theta_{ij}) \text{ (indirection)}$$

X is generated from Θ

• Log-likelihood

$$\mathcal{L} = \sum_{ij} r_{ij} \left[X_{ij} \log \sigma(\Theta_{ij}) + (1 - X_{ij}) \log [1 - \sigma(\Theta_{ij})] \right]$$

non-missing ratings

$$= \sum_{ij} r_{ij} \left[X_{ij} \log \sigma(\vec{v}_i \cdot \vec{y}_j) + (1 - X_{ij}) \log \sigma(-\vec{v}_i \cdot \vec{y}_j) \right]$$

How to maximize w.r.t. \vec{v}_i and \vec{y}_j ?

• Alternating maximization

— Optimize $\{\vec{v}_i\}_{i=1}^D$ for fixed $\{\vec{y}_j\}_{j=1}^N$

— Optimize $\{\vec{y}_j\}_{j=1}^N$ for fixed $\{\vec{v}_i\}_{i=1}^D$

— Converge to local maximum of \mathcal{L}

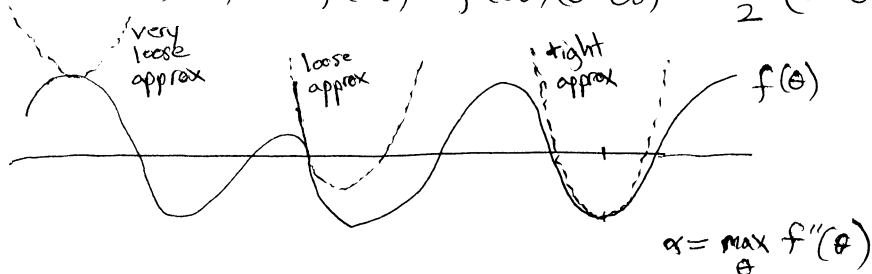
— subproblems are logistic regression

— simpler update rules? Look for auxiliary function.

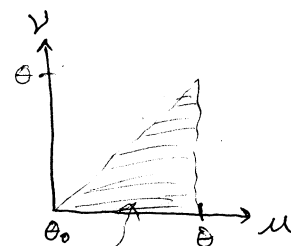
• simple inequality

Lemma: for any twice differentiable function $f(\theta)$,
let $\alpha \triangleq \max_{\theta} f''(\theta)$. Then:

$$f(\theta) \leq f(\theta_0) + f'(\theta_0)(\theta - \theta_0) + \frac{\alpha}{2} (\theta - \theta_0)^2$$



Proof: $f(\theta) = f(\theta_0) + \int_{\theta_0}^{\theta} f'(u) du$
 $= f(\theta_0) + \int_{\theta_0}^{\theta} du \left[f'(\theta_0) + \int_{\theta_0}^u dv f''(v) \right]$
 $= f(\theta_0) + f'(\theta_0)(\theta - \theta_0) + \int_{\theta_0}^{\theta} du \int_{\theta_0}^u dv f''(v)$
 $\leq f(\theta_0) + f'(\theta_0)(\theta - \theta_0) + [\max_v f''(v)] \cdot \frac{1}{2} (\theta - \theta_0)^2$



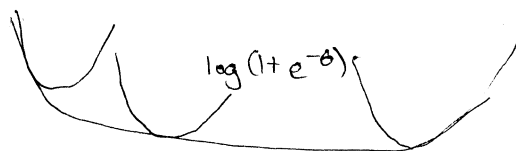
area of triangle

- Apply to negative log-sigmoid function:

$$-\log \sigma(\theta) = \log [1 + e^{-\theta}]$$

$$-\frac{d}{d\theta} \log \sigma(\theta) = -\sigma(-\theta)$$

$$-\frac{d^2}{d\theta^2} \log \sigma(\theta) = \sigma(\theta) \sigma(-\theta) \leq [\sigma(0)]^2 = \frac{1}{4}$$



$$-\log \sigma(\theta) \leq -\log \sigma(\theta_0) - \sigma(-\theta_0)(\theta - \theta_0) + \frac{1}{8} (\theta - \theta_0)^2$$

• simple auxiliary function for logistic PCA

$$F(v) = -\sum_{ij} n_{ij} [x_{ij} \log \sigma(\vec{v}_i \cdot \vec{y}_j) + (1 - x_{ij}) \log \sigma(-\vec{v}_i \cdot \vec{y}_j)]$$

$$G(v, w) = F(w) + \sum_i \frac{\partial F}{\partial w_i} (\vec{v}_i - \vec{w}_i) + \frac{1}{8} \sum_{ij} n_{ij} [(\vec{v}_i - \vec{w}_i) \cdot \vec{y}_j]^2$$

Property (i): $G(v, v) = F(v)$ b/c last terms vanish at $v=w$

Property (ii): $G(v, v) \geq F(v)$ from lemma (and algebra)

- Simpler update rule:

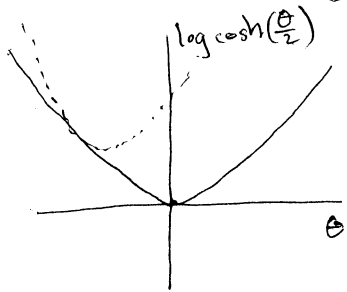
$$v_{t+1} = \underset{v}{\operatorname{argmin}} G(v, v_t)$$

Now $G(v, w)$ is quadratic in first argument.

$\min_v G(v, w)$ reduces to least squares optimization

- Tighter inequality for $\log \sigma(\theta)$

Note: $\log \sigma(\theta) = -\log [1 + e^{-\theta}]$

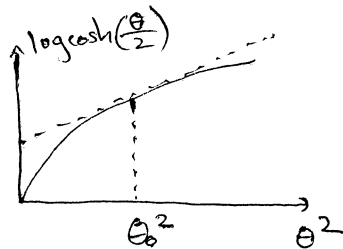


$$= -\log \{ e^{-\theta/2} e^{\theta/2} [1 + e^{-\theta}] \}$$

$$= \frac{\theta}{2} - \log (e^{\theta/2} + e^{-\theta/2})$$

$$= \frac{\theta}{2} - \log \cosh\left(\frac{\theta}{2}\right) + \log 2$$

"like a smooth absolute value function"



- One can show that $\log \cosh\left(\frac{\theta}{2}\right)$ is concave function of θ^2

Bound from concavity:

$$\log \cosh\left(\frac{\theta}{2}\right) \leq \log \cosh\left(\frac{\theta_0}{2}\right) + (\theta^2 - \theta_0^2) \left. \frac{\partial}{\partial [\theta^2]} \log \cosh\left(\frac{\theta}{2}\right) \right|_{\theta = \theta_0}$$

$$= \log \cosh\left(\frac{\theta_0}{2}\right) + (\theta^2 - \theta_0^2) \frac{\tanh(\theta_0/2)}{4\theta_0}$$

- Assembling the above:

$$-\log \sigma(\theta) \leq -\frac{\theta}{2} + \log \cosh\left(\frac{\theta_0}{2}\right) + (\theta^2 - \theta_0^2) \frac{\tanh(\theta_0/2)}{4\theta_0} - \log 2$$

∴ (after some algebra)

$$\boxed{-\log \sigma(\theta) \leq -\log \sigma(\theta_0) - \sigma(\theta_0)(\theta - \theta_0) + \frac{\tanh(\theta_0/2)}{4\theta_0} (\theta - \theta_0)^2}$$

Compare to previous bound:

Same except last term $\frac{1}{8}(\theta - \theta_0)^2$

Adaptive quadratic bound
 \Rightarrow tighter auxiliary function
 \Rightarrow faster convergence

Exponential family PCA

- Exponential family
 Random variable X
 Natural parameter θ

$$\log P(X|\theta) = \log P_0(X) + X\theta - G(\theta)$$

- Examples

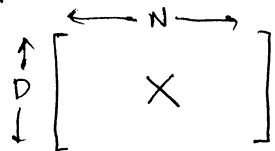
Normal $X \in \mathbb{R}$ $P_0(X) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$ $G(\theta) = \frac{\theta^2}{2}$ $P(X) = \frac{1}{\sqrt{2\pi}} e^{-(x-\theta)^2/2}$

Bernoulli $X \in \{0,1\}$ $P_0(X) = 1$ $G(\theta) = \log(1+e^\theta)$ $P(X) = \sigma(\theta)^X \sigma(-\theta)^{1-X}$

Poisson $X \in \{0,1,2,\dots,\infty\}$ $P_0(X) = \frac{1}{X!}$ $G(\theta) = e^\theta$ $P(X) = e^{-\lambda} \frac{\lambda^X}{X!}$
 with $\lambda = e^\theta$

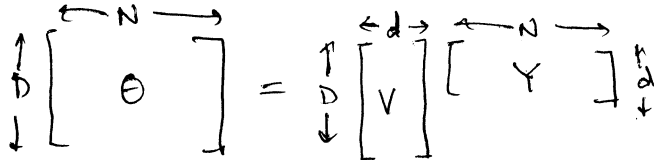
logistic PCA

- Matrix data



Assume X_{ij} has exponential family distribution with natural parameter $\theta_{ij} \in \mathbb{R}$

- Matrix factorization



$$\theta_{ij} = \vec{v}_i \cdot \vec{y}_j$$

- Log-likelihood

$$\mathcal{L} = \sum_{ij} \tau_{ij} \log P(X_{ij})$$

$$= \sum_{ij} \tau_{ij} [\log P_0(X_{ij}) + X_{ij}\theta_{ij} - G(\theta_{ij})]$$

$$= \sum_{ij} \tau_{ij} [\log P_0(X_{ij}) + X_{ij}(\vec{v}_i \cdot \vec{y}_j) - G(\vec{v}_i \cdot \vec{y}_j)]$$

- Alternating maximization
 - optimize \vec{v}_i for fixed \vec{y}_j
 - optimize \vec{y}_j for fixed \vec{v}_i
 - repeat until convergence