

Feb 11, 2020

Last time: PCA, robust PCA

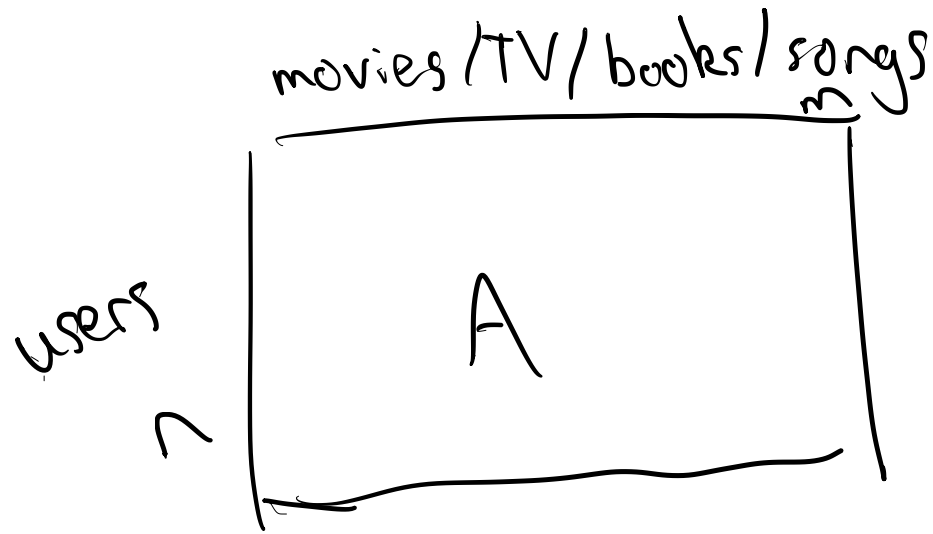
truncated SVD



Today: latent factor models thru matrix completion

HW1 due Thurs 11:59pm ET CMS

Recommender systems



A_{ij} = how much user i likes song j (# plays)

song $j \in \mathbb{R}^k \Rightarrow y_j$
user $i \in \mathbb{R}^k \Rightarrow x_i$

$(y_j)_1 \approx$ how much "rap"

$(y_j)_2 \approx$ length

$$A_{ij} \approx x_i^T y_j$$

$(x_i)_1 \approx$ proclivity for rap

$(x_i)_2 \approx$ short songs

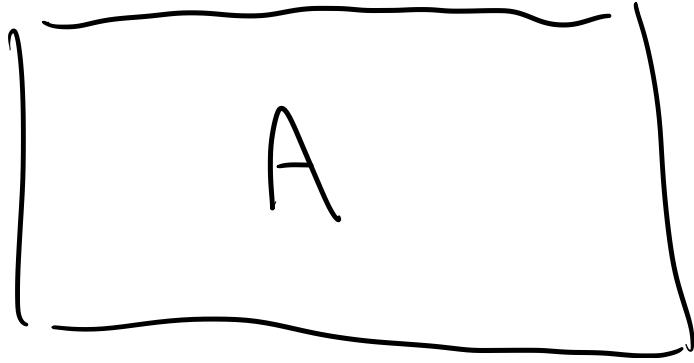
should user r like song s ?

$$A \approx XY^T \quad X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \quad Y^T = [y_1 \dots y_m]$$

Word embeddings

"contexts"

words



context $j \in \mathbb{R}^k \Rightarrow y_j$
word $i \in \mathbb{R}^k \Rightarrow x_i$

want: $\Pr(\text{choose } i \mid \text{context } j) \approx \frac{\exp(x_i^T y_j)}{\sum_r \exp(x_r^T y_j)}$

$$A_{ij} = \log(\text{Prob}(\text{choose } i \mid \text{context } j))$$

$$= x_i^T y_j - \log\left(\sum_r \exp(x_r^T y_j)\right)$$

constant
per column
force to zero

$$A \approx XY^T$$

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$$

$$Y^T = [y_1 \dots y_m]$$

$$A \approx XY^T \quad \min_{X, Y} \|A - XY^T\|_F^2$$

$X \in \mathbb{R}^{n \times k} \quad Y \in \mathbb{R}^{m \times k}$

Idea: truncated SVD

One problem: X, Y^T not unique

$$X = U_k \sqrt{\Sigma_k} \quad Y = \sqrt{\Sigma_k} V_k^T$$

$$(XB)(B^{-1}Y^T) = XY^T$$

Add some regularization:

$$\min_{X, Y} \frac{1}{2} \|A - XY^T\|_F^2 + \lambda (\|X\|_F^2 + \|Y\|_F^2)$$

Major problem: only observe a few entries of A

$$\min_{X, Y} \frac{1}{2} \|P_{\Omega}(A - XY^T)\|_F^2 + \lambda (\|X\|_F^2 + \|Y\|_F^2)$$

$$\|P_{\Omega}(M)\|_F^2 = \sum_{(i,j) \in \Omega} M_{ij}^2$$

Alternating scheme

Fix Y , opt X

$$X^{k+1} = \min_X \frac{1}{2} \|P_{\Omega}(A - XY^T)\|_F^2 + \lambda \|X\|_F^2$$

i th row of A
 $J_i = \{j \mid (i,j) \in \Omega\}$

$$\frac{1}{2} \sum_{j \in J_i} (A_{ij} - x_i^T y_j)^2$$

$$\|Y_{J_i} x_i - A_{i, J_i}\|_2^2 + \lambda \|x_i\|_2^2$$

l_2 -reg LLS on each row

$$\left\| \begin{pmatrix} Y_{j_1} & Y_{j_2} & \dots & Y_{j_n} \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} - \begin{pmatrix} A_{1j_1} \\ \vdots \\ A_{nj_n} \end{pmatrix} \right\|_2^2 + \lambda \left\| \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \right\|_2^2$$

Fix X , opt Y

$$Y^{k+1} = \underset{Y}{\text{min}} \frac{1}{2} \left\| P_{\Omega} (A - X^{k+1} Y^T) \right\|_F^2 + \lambda \|Y\|_F^2$$

$$\Omega^T = \left\{ \begin{pmatrix} \omega_{1i} \\ \vdots \\ \omega_{ni} \end{pmatrix} \mid \omega_{ij} \in \Omega \right\}$$

$$\frac{1}{2} \left\| P_{\Omega^T} (A^T - Y X^T) \right\|_F^2 + \lambda \|Y\|_F^2$$

same!

Nuclear norm approach

$$\min_Z \frac{1}{2} \|P_\Omega(A-Z)\|_F^2 + \lambda \|Z\|_*$$

$\sum \sigma_i$
↓

Last lecture: $\arg \min_Z \frac{1}{2} \|A-Z\|_F^2 + \lambda \|Z\|_*$

$$\begin{aligned} S_\lambda(A) &= U S_\lambda(\Sigma) V^T = \sum_{i=1}^n u_i v_i^T \max(\sigma_i - \lambda, 0) \\ &= U_k \Sigma_k V_k^T \end{aligned}$$

Idea: "fill in" A with Z ,
use prox,
iterate

$$\bar{A}_{ij} = \begin{cases} A_{ij} & (i,j) \in \Omega \\ z^k_{ij} & \text{o/w} \end{cases} \quad \bar{A} = z^k + Q_{\Omega}(A - z^k)$$

puts zero if $(i,j) \notin \Omega$

$$z^{k+1} = \min_z \frac{1}{2} \|\bar{A} - z\|_F^2 + \lambda \|z\|_*$$

$$s_{\lambda}(\bar{A}) = S_{\lambda}(z^k + Q_{\Omega}(A - z^k))$$

SGD

$$\min_{X, Y} \frac{1}{2} \| P_{\Omega}(A - XY^T) \|_F^2 + \lambda (\|X\|_F^2 + \|Y\|_F^2)$$

$$\sum_{(i,j) \in \Omega} (a_{ij} - x_i^T y_j)^2 + \frac{\lambda}{|J_i|} \|x_i\|_2^2 + \frac{\lambda}{|I_j|} \|y_j\|_2^2$$

$J_i = \{(i,j) \in \Omega\}$ $I_j = \{(i,j) \in \Omega\}$

$$\min \frac{1}{|\Omega|} \sum_{(i,j) \in \Omega} f_{ij}(x, y)$$

$\nabla_x f_{ij}$	non-zeroes on x_i
$\nabla_y f_{ij}$	" " y_j

sparse gradients
 $O(k)$ parameters
to update

Before: $A_{ij} \approx x_i^T y_j$ $A \approx XY^T$

Maybe: $A_{ij} \approx x_i^T y_j + b_i + c_j + \mu$

user i bias *item j bias* *overall bias*

$$A \approx XY^T + \mathbf{1}b^T + c\mathbf{1}^T + \mu\mathbf{1}\mathbf{1}^T$$

$$\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} (b_1 \dots b_n)$$

Should this actually work?

$$n \begin{array}{|c|} \hline \begin{array}{c} m \\ \hline A \\ \hline \end{array} \\ \hline \end{array}$$

$$m \geq n$$

$$A = e_i e_j^T = \begin{bmatrix} 0 & \dots & 1 & \dots & 0 \\ \vdots & & \vdots & & \vdots \\ 0 & \dots & 0 & \dots & 0 \end{bmatrix}$$

$$A = \sum_{i=1}^k \sigma_i u_i v_i^T \quad u_i, v_i \text{ sampled uniformly}$$

$$\mathcal{O}(m^{5/4} \log(m) k) \text{ UAR samples}$$

$$\text{Thm: } \min \|z\|_*$$

$$\text{s.t. } z_\Omega = A_\Omega$$

recover A w.h.p. (Candès & Recht)

(incoherence

$$\max_i \|(UU^T)(i, :)\|_2$$

$$\|UU^T\|_2 \rightarrow \infty$$