

Convergence of CG and polynomial approximation:

$$\frac{\|x_n - x_*\|_Q}{\|x_0 - x_*\|_Q} \leq \min_{r(t)} \max_{\lambda_1, \dots, \lambda_m} |r(\lambda_i)|$$

$r(t)$ varies over polynomials of degree $\leq n$ s.t. $r(0) = 1$

Proof:

$$\|x_0 - x_*\|_Q^2 = \|x_* - x_*\|_Q^2 = \|x_*\|_Q^2 = x_*^T Q x_*$$

If $Q = UDU^T$ is an eigenvalue decomp., and $x_* = Uc$ are the coordinates of x_* in the eigenvector basis, then

$$x_*^T Q x_* = c^T \cancel{U^T} U D U^T U c = [c_1 \dots c_m] \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_m \end{bmatrix} \begin{bmatrix} c_1 \\ \vdots \\ c_m \end{bmatrix}$$

$$= \lambda_1 c_1^2 + \lambda_2 c_2^2 + \dots + \lambda_m c_m^2.$$

$$\|x_n - x_*\|$$

$$x_n = \arg \min_{x \in K_n(Q, v)} \|x - x_*\|_Q$$

$$K_n(Q, v) = \{ p(Q)v : p \text{ polynomial of degree } \leq n \}$$

$$x_n = \arg \min_{p(t)} \|p(Q)v - x_*\|_Q = \arg \min_{p(t)} \|x_* - p(Q)Qx_*\|_Q$$

$$= \arg \min_{r(t)} \|r(Q)x_*\|_Q$$

$$\text{where } r(t) = 1 - p(t)t$$

varies over all polynomials of degree $\leq n$ such that $r(0) = 1$

$$\text{If } Q = UDU^T, \quad r(Q) = U \begin{bmatrix} r(\lambda_1) & & \\ & r(\lambda_2) & \\ & & \ddots \\ & & & r(\lambda_m) \end{bmatrix} U^T$$

(from earlier results)

$$= \arg \min_{r(\cdot)} \left\| U \begin{bmatrix} r(\lambda_1) \\ \vdots \\ r(\lambda_m) \end{bmatrix} U^T c \right\|_Q = \arg \min_{r(\cdot)} \left\| U \begin{bmatrix} r(\lambda_1) c_1 \\ \vdots \\ r(\lambda_m) c_m \end{bmatrix} \right\|_Q$$

$$= \arg \min_r \sqrt{c_1^2 |r(\lambda_1)|^2 \lambda_1 + c_2^2 |r(\lambda_2)|^2 \lambda_2 + \dots + c_m^2 |r(\lambda_m)|^2 \lambda_m}$$

$$\frac{\|X_u - X_*\|_Q^2}{\|X_o - X_*\|_Q^2} = \min_{r(\cdot)} \frac{c_1^2 |r(\lambda_1)|^2 \lambda_1 + \dots + c_m^2 |r(\lambda_m)|^2 \lambda_m}{c_1^2 \lambda_1 + \dots + c_m^2 \lambda_m} \leq$$

$$\leq \min_{r(\cdot)} \frac{\max |r(\lambda_i)|^2 (c_1^2 \lambda_1 + \dots + c_m^2 \lambda_m)}{c_1^2 \lambda_1 + \dots + c_m^2 \lambda_m}$$

Singular value decomposition

$$A \in \mathbb{R}^{m \times n} \quad m > n \quad A = USV^T, \text{ with } U, V \text{ orthogonal,}$$

$m \times m \quad m \times n \quad n \times n$

$$S = \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_n \\ & & & & 0 \end{bmatrix} \quad \sigma_1 > \sigma_2 > \dots > \sigma_n$$

$$A = U_1 \sigma_1 V_1^T + \dots + U_{\min(m,n)} \sigma_{\min(m,n)} V_{\min(m,n)}^T$$

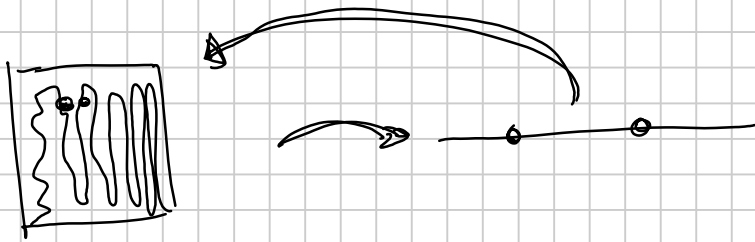
Theorem (Eckart-Young): the solution of

$$X_k = \arg \min_{\text{rk } X \leq k} \|A - X\|$$

is

$$X_k = U_1 \sigma_1 V_1^T + \dots + U_k \sigma_k V_k^T$$

$\boxed{\text{rk } X = k} \rightarrow m+n+1$



If we compress one $m \times n$ matrix to rank k ,
 we save space

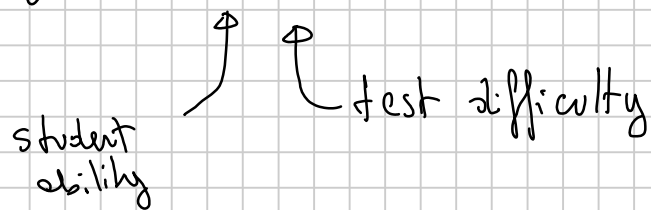
$$mn \rightarrow mk + nk + k$$

	Test A	Test B	Test C
Student 1	90	95	85
Student 2	32	35	31
Student 3	70	75	71
Student 4	70	31	72

= A

If A were rank 1, $A = xy^T = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} [y_1 \ y_2 \ y_3] = \begin{bmatrix} x_1 y_1 & x_1 y_2 & x_1 y_3 \\ x_2 y_1 & x_2 y_2 & x_2 y_3 \\ \vdots & \vdots & \vdots \\ x_4 y_1 & x_4 y_2 & x_4 y_3 \end{bmatrix}$

result of student i on test j is $x_i y_j$



rank 1 \Rightarrow perfect structure in your data.

If A has rank 2,

$$A = xy^T + Wz^T = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} [y_1 \ y_2 + y_3] + \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix} [z_1 \ z_2 \ z_3]$$

$$A_{ij} = x_i y_j + w_i z_j$$

$$x, y, z, w \geq 0$$

Idea:

2 different topics, e.g. programming and mathematics

Each student i gets a score that reflect their abilities in the two subjects (x_i, w_i) and the amount of math difficulty / programming difficulty in each exercise j
 (y_j, z_j)

svd approximation: $A \approx U, \Sigma, V, ^T$ gives the best possible rank-1 approximation of true student scores

Entries of $u_i \sim$ estimated ranking of students
 $v_i \sim$ est. ranking of problems.

$$\text{"Best"} = \text{minimum} \sum_{i=1}^m \sum_{j=1}^n (A_{ij} - u_i v_j)^2$$

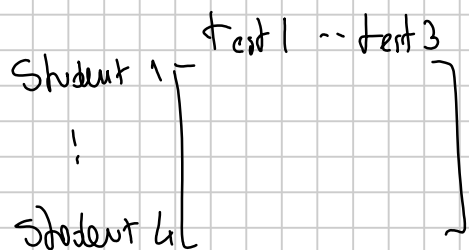
In statistics terms:

$$\text{model: } \text{Score}_{ij} = x_i y_j + \epsilon_{ij}$$

$$\epsilon_{ij} = \text{error} \sim \text{Gaussian}(0, \sigma)$$

maximum likelihood estimator \Leftrightarrow most likely choice of x, y that could have generated this data.

$$ML \Leftrightarrow \min \sum_{i=1}^m \sum_{j=1}^n (A_{ij} - x_i y_j)^2 = \min \|A - xy^T\|_F^2$$



$$\text{rank 2: } A_{ij} = x_i y_j + w_i z_j$$

two different topics, each with ability / difficulty

With an SVD:

$$A_{ij} = U_1 \sigma_1 U_1^T + U_2 \sigma_2 U_2^T + U_3 \sigma_3 U_3^T$$

↑
overall ability

higher rank: further corrections for further topics

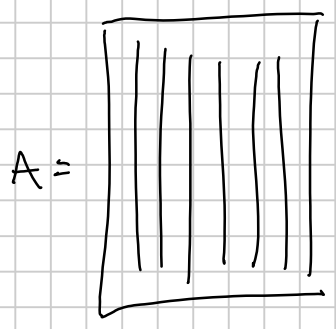
Since U_i, V_j have norm 1, the σ_i tell me how large the corrections are.

$\sigma_1 \gg \sigma_2$: the rank-2 approximation is only slightly changed w.r.t. the rank-1 approximation
→ the classification with only 1 topic was already pretty good

$\sigma_1 \approx \sigma_2$: classification into 2 topics makes a big change.

Related: Latent semantic analysis principal component analysis (PCA)

PCA: given a data matrix



x_j = columns of A = observations

$$\mu = \text{mean} = \frac{1}{n} \sum_{j=1}^n x_j$$

$$\hat{x}_j = x_j - \bar{x}$$

$$\hat{A} = \begin{bmatrix} \hat{x}_1 & \dots & \hat{x}_n \\ \vdots & & \vdots \end{bmatrix}$$

Variances and covariances of each component of the data:

$$\text{Covariance matrix } C = \frac{1}{n-1} \left(\hat{x}_1 \hat{x}_1^T + \dots + \hat{x}_n \hat{x}_n^T \right)$$

principal components = eigenvectors of the covariance matrix

This can also be seen in terms of the SVD:

$$\hat{A} = USV^T$$

$$C = \frac{1}{n-1} \hat{A} \hat{A}^T = \frac{1}{n-1} USV^T \cancel{VS^T} U^T = \frac{1}{n-1} U \begin{bmatrix} \sigma_1^2 & & & \\ & \sigma_2^2 & & \\ & & \dots & \\ & & & \sigma_n^2 \\ & & & & \dots \\ & & & & & 0 \end{bmatrix} U^T$$