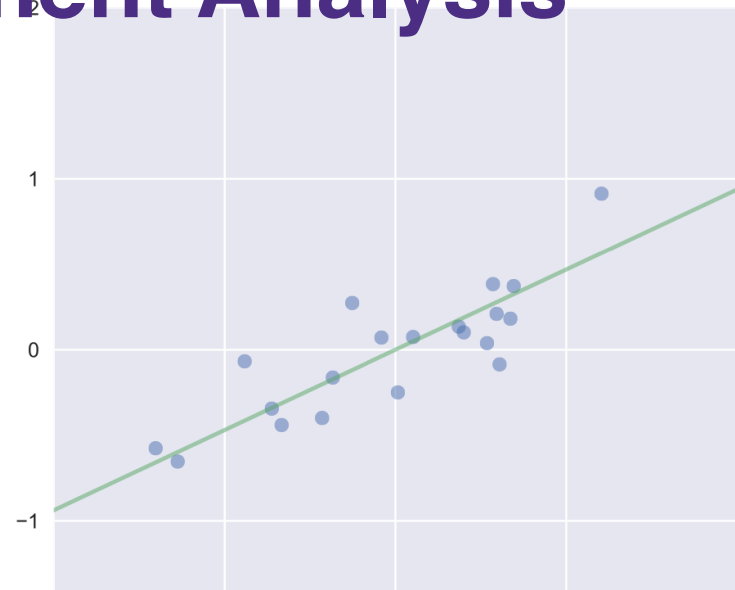


Principal Component Analysis



Principal components is the subspace that minimizes the reconstruction error

$$\underset{u_1, \dots, u_r}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \|x_i - p_i\|_2^2$$

$$p_i = \sum_{j=1}^r (u_j^T x_i) u_j = \mathbf{U} \mathbf{U}^T x_i$$

$$\text{where } \mathbf{U} = [u_1 \ u_2 \ \cdots \ u_r] \in \mathbb{R}^{d \times r}$$

$$\begin{aligned} &\underset{\mathbf{U}}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \|x_i - \mathbf{U} \mathbf{U}^T x_i\|_2^2 \\ &\text{subject to} \quad \mathbf{U}^T \mathbf{U} = \mathbf{I}_{r \times r} \end{aligned}$$

Q. How do we solve this optimization?

Minimizing reconstruction error to find principal components

$$\frac{1}{n} \sum_{i=1}^n x_i^T x_i - 2 x_i^T U U^T x_i + \underbrace{x_i^T U U^T U U^T x_i}_{II}$$

$$= \frac{1}{n} \sum_{i=1}^n x_i^T x_i - x_i^T U U^T x_i$$

$$= \underbrace{\frac{1}{n} \sum_{i=1}^n x_i^T x_i}_{\text{Const.}} - \frac{1}{n} \sum_{i=1}^n \underbrace{x_i^T U U^T x_i}_{\sum_{j=1}^r (x_i^T u_j)^2}$$

$$U = \begin{bmatrix} u_1 & u_2 & \dots & u_r \end{bmatrix}$$

$$\text{minimize}_{u_1, \dots, u_r} - \sum_{j=1}^r \frac{1}{n} \sum_{i=1}^n (x_i^T u_j)^2$$

$$\text{s.t.} \quad U^T U = I$$

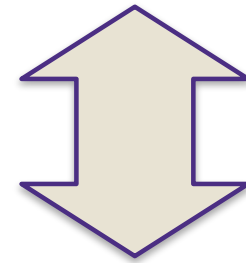
$$\underset{U}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \|x_i - U U^T x_i\|_2^2$$

$$\text{subject to} \quad U^T U = I_{r \times r}$$

Minimizing reconstruction error to find principal components

$$\begin{aligned}
 & \frac{1}{n} \sum_{i=1}^n \|x_i - UU^T x_i\|_2^2 \\
 &= \frac{1}{n} \sum_{i=1}^n \left\{ \|x_i\|_2^2 - 2x_i^T UU^T x_i + x_i^T U \underbrace{U^T U}_{=I} U^T x_i \right\} \\
 &= \underbrace{\frac{1}{n} \sum_{i=1}^n \|x_i\|_2^2}_{\text{does not depend on } U} - \frac{1}{n} \sum_{i=1}^n x_i^T UU^T x_i \\
 &= C - \sum_{j=1}^r \underbrace{\frac{1}{n} \sum_{i=1}^n (u_j^T x_i)^2}_{\text{Variance in direction } u_j}
 \end{aligned}$$

$$\begin{aligned}
 & \underset{U}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \|x_i - UU^T x_i\|_2^2 \\
 & \text{subject to} \quad U^T U = \mathbf{I}_{r \times r}
 \end{aligned}$$



Variance along U_j

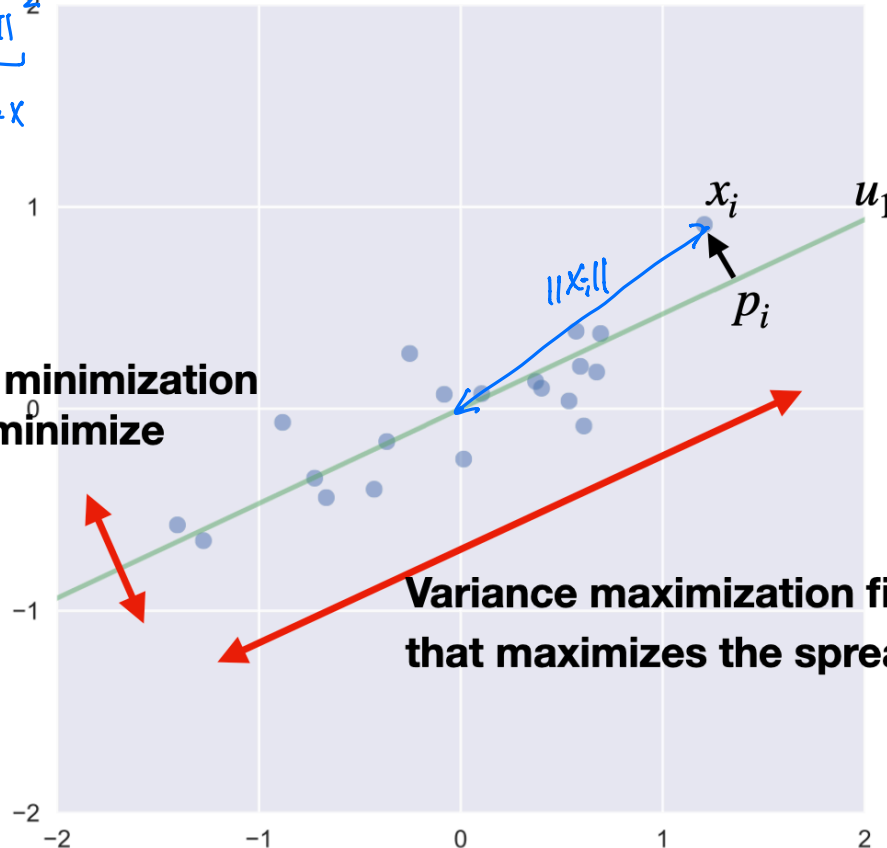
$$\begin{aligned}
 & \underset{U}{\text{maximize}} \quad \sum_{j=1}^r \underbrace{\frac{1}{n} \sum_{i=1}^n (u_j^T x_i)^2}_{\text{Variance along } U_j} \\
 & \text{subject to} \quad U^T U = \mathbf{I}_{r \times r}
 \end{aligned}$$

Variance maximization vs. reconstruction error minimization

- both give the same principal components as optimal solution

$$\|x_i\|_2^2 = \underbrace{\|x_i - p_i\|^2}_{\min} + \underbrace{\|p_i\|^2}_{\max}$$

Reconstruction error minimization
finds directions that minimize
the distances to p_i 's



Variance maximization finds directions
that maximizes the spread of p_i 's

Maximizing variance to find principal components

$$\begin{aligned} &\underset{U}{\text{maximize}} \quad \sum_{j=1}^r \frac{1}{n} \sum_{i=1}^n (u_j^T x_i)^2 \\ &\text{subject to} \quad U^T U = \mathbf{I}_{r \times r} \end{aligned}$$

We will solve it for $r = 1$ case,
and the general case follows similarly

$$\underset{u: \|u\|_2=1}{\text{maximize}} \quad \frac{1}{n} \sum_{i=1}^n (u^T x_i)^2$$

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n u^T x_i x_i^T u$$

$$C = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$

$$\underset{u: \|u\|_2=1}{\text{maximize}} \quad u^T C u$$

Maximizing variance to find principal components

$$C = \frac{1}{n} \sum_{i=1}^n \boxed{x_i x_i^T} \quad \text{PSD}$$

$$\text{maximize}_u u^T C u \quad (a)$$

$$\text{subject to } \|u\|_2^2 = 1$$

- we first claim that this optimization problem has the same optimal solution as the following **inequality constrained** problem

$$\text{maximize}_u u^T C u \quad (b)$$

$$\text{subject to } \|u\|_2^2 \leq 1$$

- the reason is that, because $u^T C u \geq 0$ for all $u \in \mathbb{R}^d$, the optimal solution of (b) has to have $\|u\|_2^2 = 1$
- if it did not have $\|u\|_2^2 = 1$, say $\|u\|_2^2 = 0.9$, then we can just multiply this u by a constant factor of $\sqrt{10/9}$ and increase the objective by a factor of $10/9$ while still satisfying the constraints

$$\hat{u} = \sqrt{\frac{10}{9}} u^* \rightarrow \frac{10}{9} \text{ larger}$$

$$\text{maximize}_u u^T \mathbf{C} u \quad (b)$$

$$\text{subject to } \|u\|_2^2 \leq 1$$

- we are maximizing the variance, while **keeping u small**
- this can be reformulated as an unconstrained problem, with Lagrangian encoding, to move the constraint into the objective

$$\text{maximize}_u \underbrace{u^T \mathbf{C} u - \lambda \|u\|_2^2}_{F_\lambda(u)} \quad (c)$$

- this encourages small u as we want, and we can make this connection precise: there exists a (unknown) choice of λ such that the optimal solution of (c) is the same as the optimal solution of (b)
- further, for this choice of λ , the optimal u has $\|u\|_2 = 1$

Solving the unconstrained optimization

$$\text{maximize}_u \underbrace{u^T \mathbf{C} u - \lambda \|u\|_2^2}_{F_\lambda(u)}$$

- to find such λ and the corresponding u , we solve the unconstrained optimization, by setting the gradient to zero

$$\nabla_u F_\lambda(u) = 2\mathbf{C}u - 2\lambda u = 0$$

- the candidate solution satisfies: $\mathbf{C}u^* = \lambda u^*$,
i.e. an eigenvector of \mathbf{C}

$$\text{maximize}_u u^T \mathbf{C} u$$

$$\text{subject to } \|u\|_2^2 = 1$$

- let $(\lambda^{(1)}, u^{(1)})$ denote the largest eigenvalue and corresponding eigenvector of \mathbf{C} , with norm one, i.e. $\|u^{(1)}\|_2^2 = 1$
- The maximum is achieved when $u = u^{(1)}$

The principal component analysis

- so far we considered finding ONE principal component $u \in \mathbb{R}^d$
- it is the eigenvector corresponding to the maximum eigenvalue of the covariance matrix

$$\mathbf{C} = \frac{1}{n} \mathbf{X}^T \mathbf{X} \in \mathbb{R}^{d \times d}$$

- We can use Singular Value Decomposition (SVD) to find such eigen vector
- note that is the data is not centered at the origin, we should re-center the data before applying SVD
- in general we define and use multiple principal components
- if we need r principal components, we take r eigenvectors corresponding to the largest r eigenvalues of \mathbf{C}

Algorithm: Principal Component Analysis

- **input:** data points $\{x_i\}_{i=1}^n$, target dimension $r \ll d$

- **output:** r -dimensional subspace U

- **algorithm:**

- compute mean $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

- compute covariance matrix

$$\mathbf{C} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$

- let (u_1, \dots, u_r) be the set of (normalized) eigenvectors with corresponding to the largest r eigenvalues of \mathbf{C}

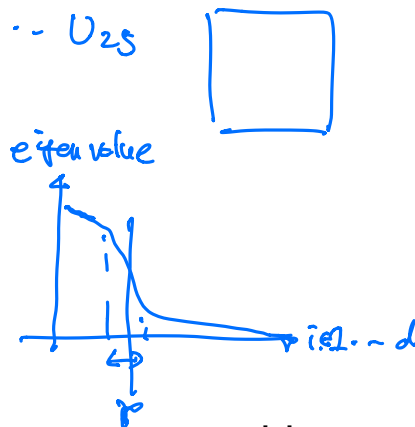
- return $\mathbf{U} = [u_1 \ u_2 \ \dots \ u_r]$

- further the data points can be represented compactly via

$$a_i = \mathbf{U}^T(x_i - \bar{x}) \in \mathbb{R}^r$$

$x_i \sim \text{face}$

$u_1 \dots u_{25}$



Singular Value Decomposition (SVD)

$$V = \begin{bmatrix} | & | & \dots & | \\ v_1 & v_2 & \dots & v_r \\ | & | & \dots & | \end{bmatrix}$$

Theorem (SVD): Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ with rank $r \leq \min\{m, n\}$. Then $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^T$ where $\mathbf{S} \in \mathbb{R}^{r \times r}$ is diagonal with positive entries, $\mathbf{U}^T\mathbf{U} = I$, $\mathbf{V}^T\mathbf{V} = I$.

$$\underbrace{V S U^T U S}_{\parallel} \underbrace{\begin{bmatrix} e_i \\ \vdots \\ \tilde{v_i^T \cdot v_i} \end{bmatrix}}_{\parallel} = V \cdot S^2 \cdot e_i = \delta_i^2 \cdot v_i$$

What is $\underline{A^T A v_i = \delta_i^2 v_i}$

What is $AA^T u_i = \delta_i^2 u_i$

$$\begin{aligned} AA^T &= U S^2 U^T \\ A^T A &= V S^2 V^T \end{aligned}$$

$$A = \begin{bmatrix} | & | & \dots & | \\ \vdots & \delta_1 & \vdots & \vdots \\ | & | & \dots & | \end{bmatrix} \begin{bmatrix} \vdots & \vdots & \dots & \vdots \\ \delta_1 & \delta_2 & \dots & \delta_r \\ \vdots & \vdots & \dots & \vdots \end{bmatrix} \begin{bmatrix} | & | & \dots & | \\ \vdots & \vdots & \dots & \vdots \\ | & | & \dots & | \end{bmatrix}$$

- v_i 's are the r eigen vectors of $A^T A$ with corresponding eigen values S_{jj}^2 's
- u_i 's are the r eigen vectors of AA^T with corresponding eigen values S_{jj}^2 's
- Computing SVD takes $O(mnr)$ operations

Singular Value Decomposition (SVD)

$$u_k = v_k$$

- Consider a full rank matrix $A \in \mathbb{R}^{m \times n}$ whose SVD is $A = USV^T$, and we want to find the best rank- r approximation of A that minimizes the error



$$\text{minimize}_{L \in \mathbb{R}^{m \times n}} \sum_{i=1}^m \sum_{j=1}^n (A_{i,j} - L_{i,j})^2$$

$$\text{subject to } \text{rank}(L) = r$$

- The optimal rank- r approximation is $U_{1:r} S_{1:r,1:r} V_{1:r}^T$

How do we compute singular vectors?

- In practice: Lanczos method
- We will learn: power iteration
- Let $C = USU^T \in \mathbb{R}^{d \times d}$ be SVD of the matrix we want to compute the top one singular vector
 - $U = [u_1, u_2, \dots, u_d]$ are the singular vectors
(ordered in the decreasing order of the corresponding singular values)
 - We also assume $\lambda_1 > \lambda_2$ in order to ensure uniqueness of u_1

$v_0 \sim \text{anything random}$

$$\begin{aligned}\tilde{v}_{t+1} &\leftarrow C v_t \\ v_{t+1} &\leftarrow \frac{\tilde{v}_{t+1}}{\|\tilde{v}_{t+1}\|_2}\end{aligned}$$

Claim: $v_t \xrightarrow[t \rightarrow \infty]{} u_1$

Power iteration

$$u_1^T C = \lambda_1 \cdot u_1^T$$

1st singular value
↓

$$\tilde{v}_{t+1} \leftarrow C v_t$$

$$v_{t+1} \leftarrow \frac{\tilde{v}_{t+1}}{\|\tilde{v}_{t+1}\|_2}$$

$$\begin{aligned} \|u_1 - v_{t+1}\|^2 &= 1 - 2 u_1^T v_{t+1} + 1 \\ &= 2 (1 - u_1^T v_{t+1}) = 2 (1 - u_1^T \frac{C v_t}{\|C v_t\|_2}) \end{aligned}$$

$$= 2 (1 - \frac{\lambda_1}{\|C v_t\|_2} \cdot u_1^T v_t)$$

$$= 2 (1 - u_1^T v_t + u_1^T v_t - \frac{\lambda_1}{\|C v_t\|_2} \cdot u_1^T v_t)$$

$$= \underbrace{\|u_1 - v_t\|^2}_{\text{optimal dist.}} - 2 \underbrace{\left(\frac{\lambda_1}{\|C v_t\|_2} - 1 \right)}_{> 0} \cdot u_1^T v_t \leq \|u_1 - v_t\|^2_t$$

equal only if
 $v_t = u_1$

$$\|C \cdot v_t\|_2 < \lambda_1$$

Power iteration for general rank- r

$$\tilde{v}_{t+1} \leftarrow C v_t$$

$$v_{t+1} \leftarrow \frac{\tilde{v}_{t+1}}{\|\tilde{v}_{t+1}\|_2}$$

$$C - \lambda_i \cdot u_i u_i^T \rightarrow C$$

- First approach:

- Repeat r times

- Run rank-1 power iteration

- Subtract $C - (v_t^T C v_t) v_t v_t^T \rightarrow C$
 $u_1, u_2, u_3 \dots, u_r$

- Second approach:

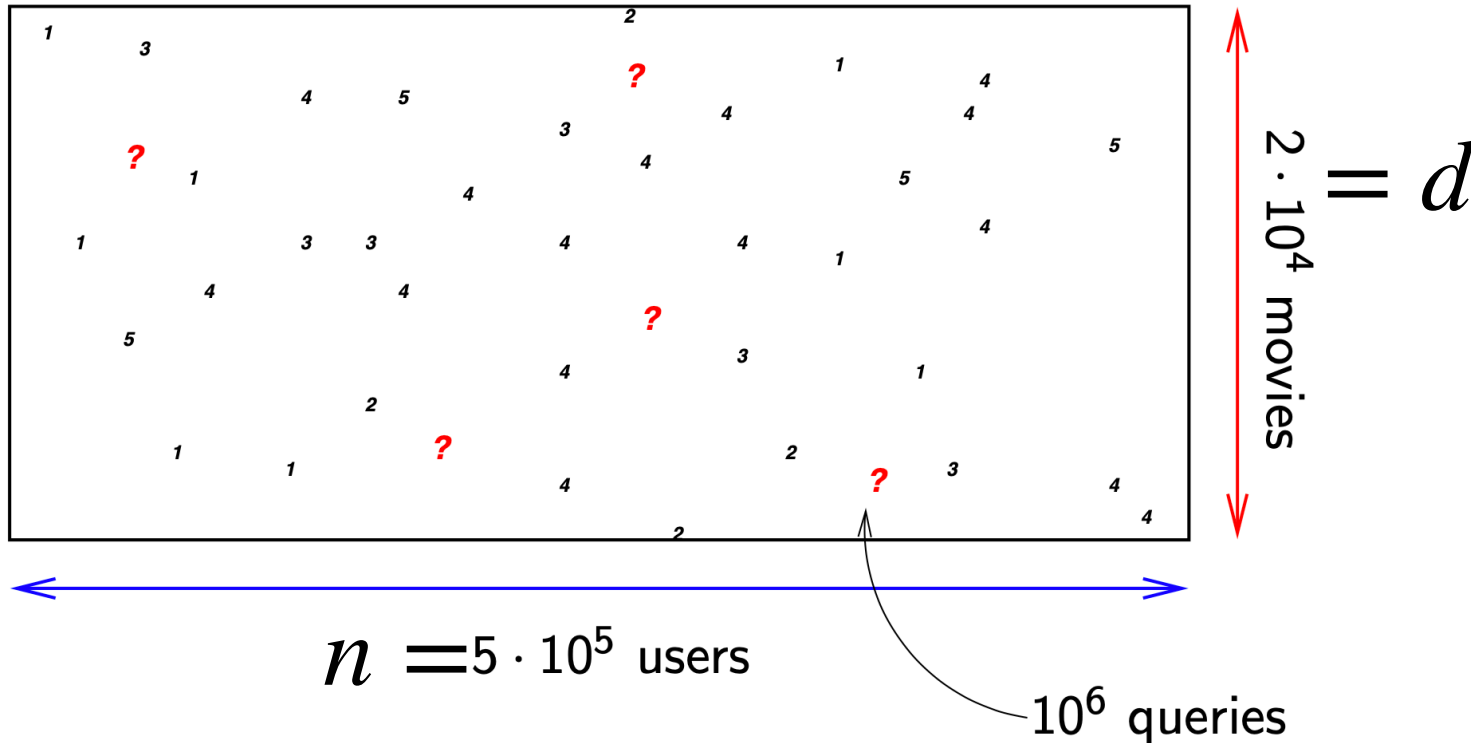
$$V_0 = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}^r, \text{ s.t. } V_0^T V_0 = I$$

$$\tilde{V}_{t+1} \leftarrow C \cdot V_t$$

$$V_{t+1} \leftarrow \underbrace{Q, R}_{\text{QR decomposition}}(\tilde{V}_{t+1})$$

Matrix completion for recommendation systems

Netflix challenge dataset



- users provide ratings on a few movies, and we want to predict the missing entries in this ratings matrix, so that we can make recommendations
- without any assumptions, the missing entries can be anything, and no prediction is possible

Matrix completion problem

- however, the ratings are not arbitrary, but people with similar tastes rate similarly
- such structure can be modeled using low dimensional representation of the data as follows
- we will find a set of principal component vectors

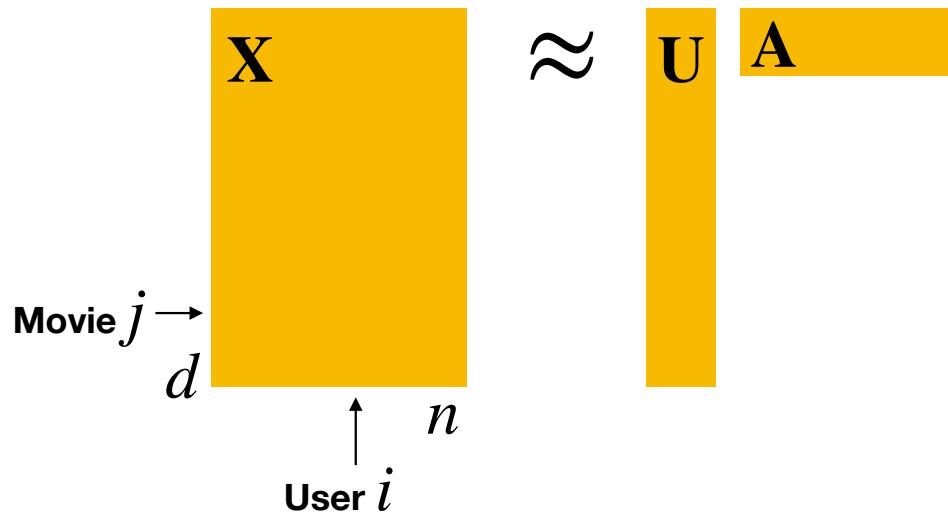
$$\mathbf{U} = [u_1 \quad u_2 \quad \cdots \quad u_r] \in \mathbb{R}^{d \times r} = d \left\{ \begin{array}{c} \text{ } \end{array} \right\} \# \text{ of movies}$$
- such that that ratings $x_i \in \mathbb{R}^d$ of user i , can be represented as

$$\begin{aligned} x_i &= a_i[1]u_1 + \cdots a_i[r]u_r \\ &= \mathbf{U}a_i \end{aligned}$$

for some lower-dimensional $a_i \in \mathbb{R}^r$ for i -th user and some $r \ll d$
- for example, $u_1 \in \mathbb{R}^d$ means how horror movie fans like each of the d movies,
- and $a_i[1]$ means how much user i is fan of horror movies

Matrix completion

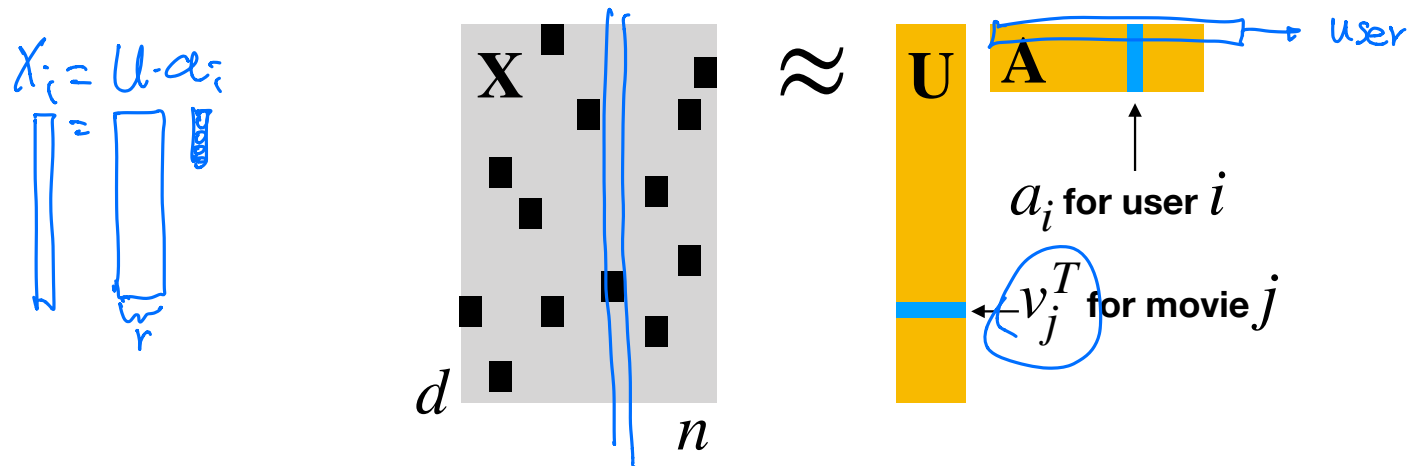
- let $\mathbf{X} = [x_1 \ x_2 \ \cdots \ x_n] \in \mathbb{R}^{d \times n}$ be the ratings matrix, and assume it is fully observed, i.e. we know all the entries
- then we want to find $\mathbf{U} \in \mathbb{R}^{d \times r}$ and $\mathbf{A} = [a_1 \ a_2 \ \cdots \ a_n] \in \mathbb{R}^{r \times n}$ that approximates \mathbf{X}



- if we **observe all entries** of \mathbf{X} , then we can find the best rank- r approximation with SVD

Matrix completion

- in practice, we only observe \mathbf{X} partially
- let $S_{\text{train}} = \{(i_\ell, j_\ell)\}_{\ell=1}^N$ denote N observed ratings for user i_ℓ on movie j_ℓ



- let v_j^T denote the j -th row of \mathbf{U} and a_i denote i -th column of \mathbf{A}
- then user i 's rating on movie j , i.e. X_{ji} is approximated by $v_j^T a_i$, which is the inner product of v_j (a column vector) and a column vector a_i
- we can also write it as $\langle v_j, a_i \rangle = v_j^T a_i$

Matrix completion

- a natural approach to fit v_j 's and a_i 's to given training data is to solve

$$\text{minimize}_{\mathbf{U}, \mathbf{A}} \sum_{(i,j) \in \mathcal{S}_{\text{train}}} (\mathbf{X}_{ji} - v_j^T a_i)^2$$

Handwritten annotations: $d \times r$ (pointing to \mathbf{U}), $r \times n$ (pointing to \mathbf{A}), and $r(d+n)$ (boxed).

- this can be solved, for example via gradient descent or alternating minimization
- this can be quite accurate, with small number of samples

- Theorem [Keshavan, Montanari, Oh 2009]

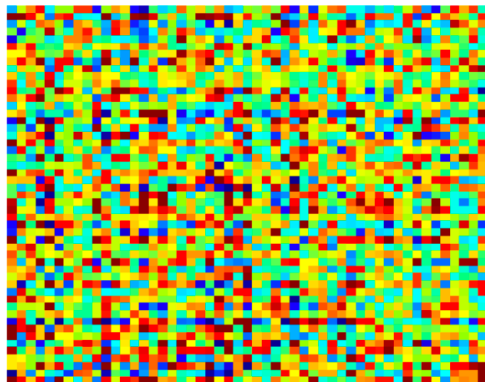
Assume the ground truths \mathbf{X} has rank r , then

(a variant of) gradient descent finds the optimal solution if we

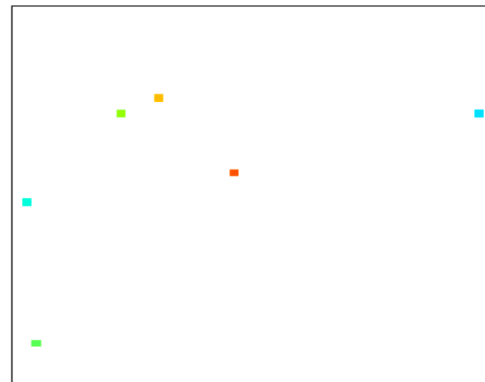
observe more than $c r (d + n) \log(dn)$ entries at random positions

Example: 2000×2000 rank-8 random matrix

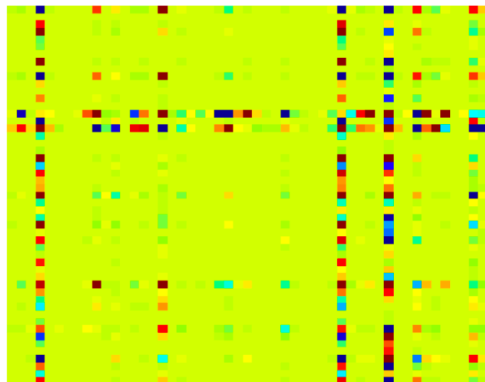
low-rank matrix \mathbf{X}



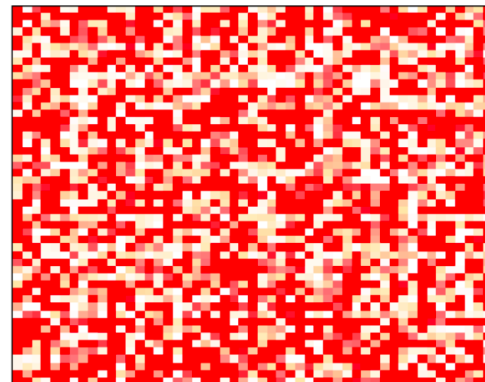
sampled matrix



Gradient descent output \mathbf{UA}



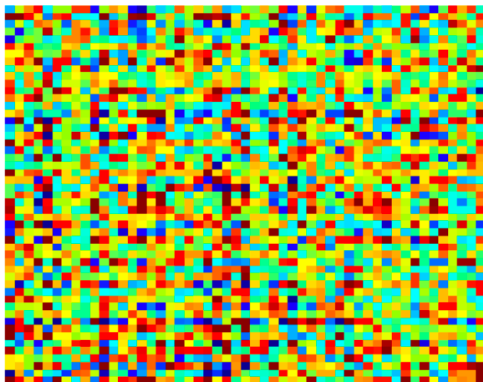
squared error $(\mathbf{X}_{ji} - (\mathbf{UA})_{ji})^2$



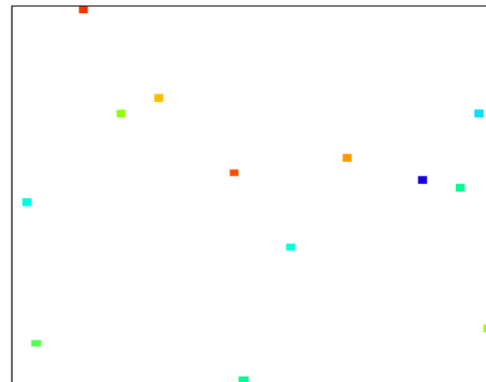
0.25% sampled

Example: 2000×2000 rank-8 random matrix

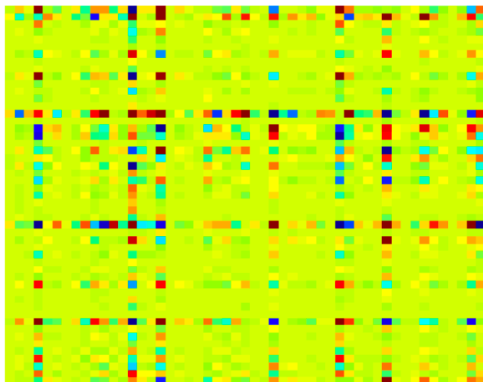
low-rank matrix \mathbf{X}



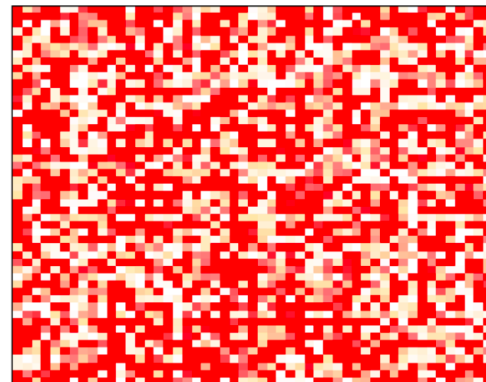
sampled matrix



Gradient descent output $\mathbf{U}\mathbf{A}$



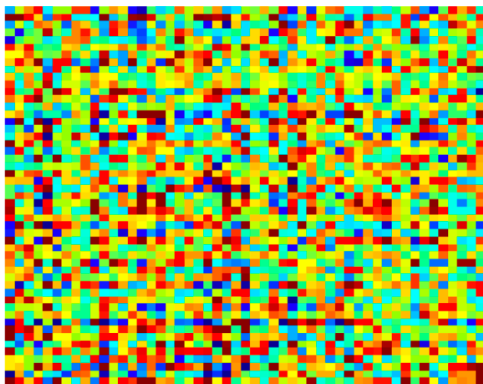
squared error $(\mathbf{X}_{ji} - (\mathbf{U}\mathbf{A})_{ji})^2$



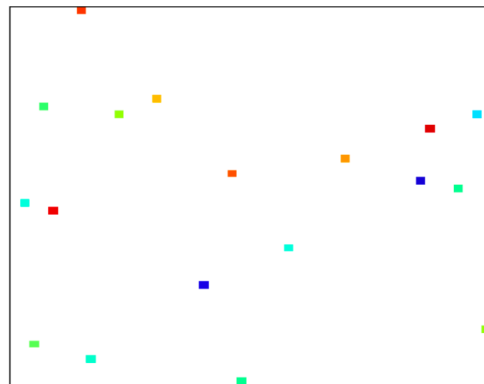
0.50% sampled

Example: 2000×2000 rank-8 random matrix

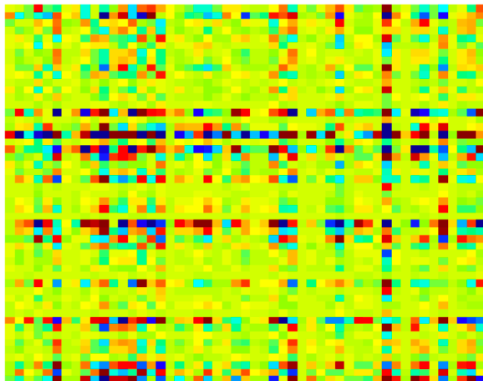
low-rank matrix \mathbf{X}



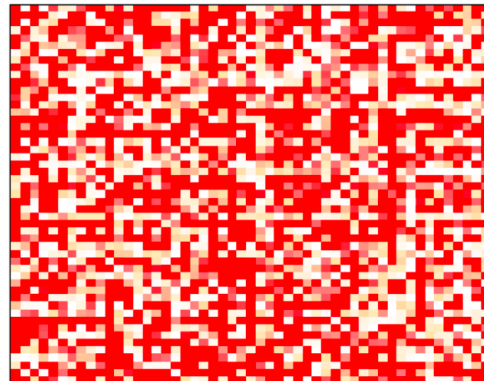
sampled matrix



Gradient descent output $\mathbf{U}\mathbf{A}$



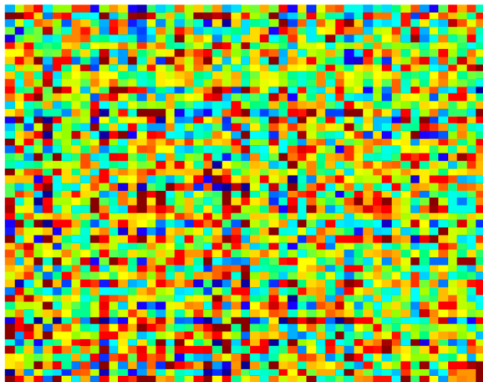
squared error $(\mathbf{X}_{ji} - (\mathbf{U}\mathbf{A})_{ji})^2$



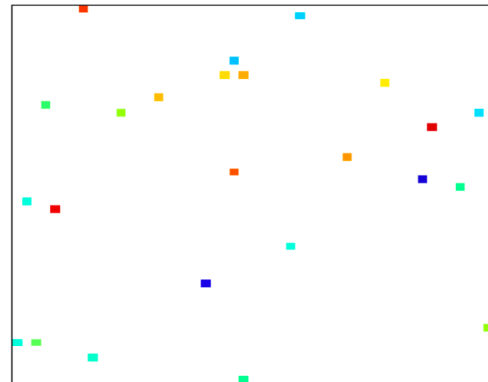
0.75% sampled

Example: 2000×2000 rank-8 random matrix

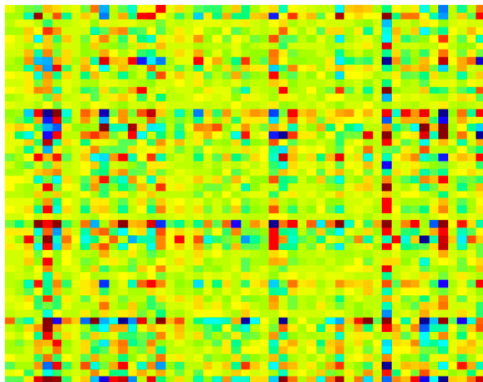
low-rank matrix \mathbf{X}



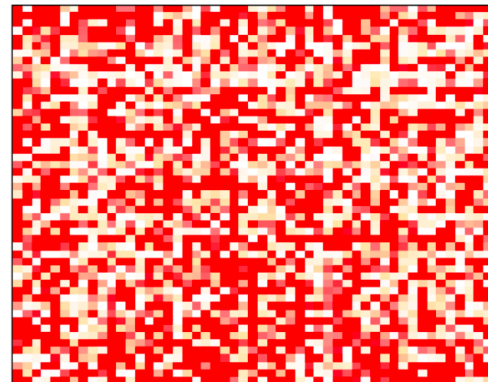
sampled matrix



Gradient descent output $\mathbf{U}\mathbf{A}$



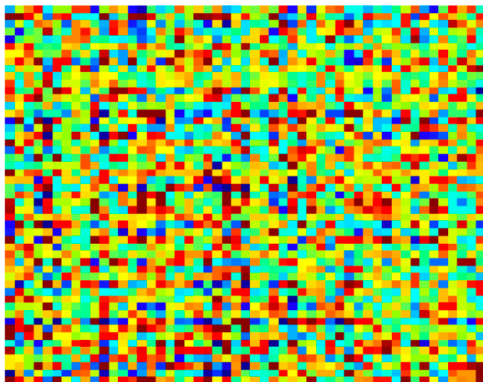
squared error $(\mathbf{X}_{ji} - (\mathbf{U}\mathbf{A})_{ji})^2$



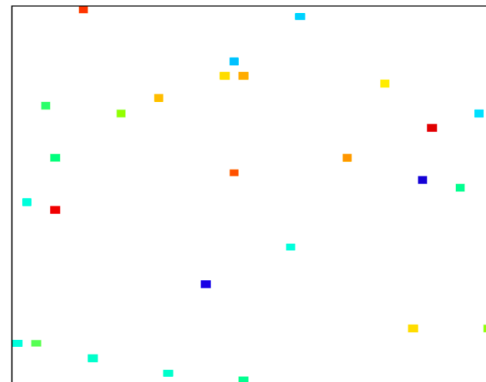
1.00% sampled

Example: 2000×2000 rank-8 random matrix

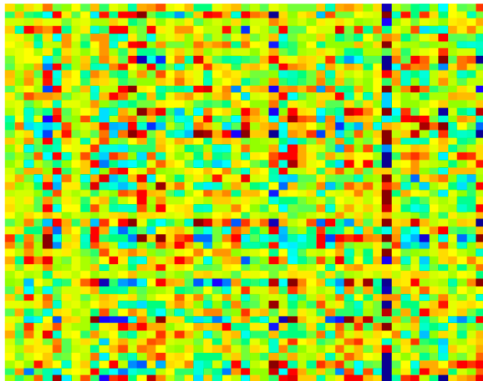
low-rank matrix \mathbf{X}



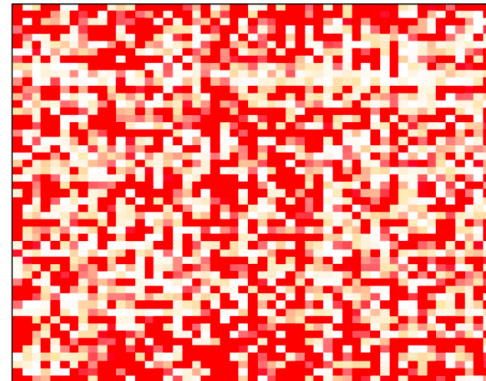
sampled matrix



Gradient descent output $\mathbf{U}\mathbf{A}$



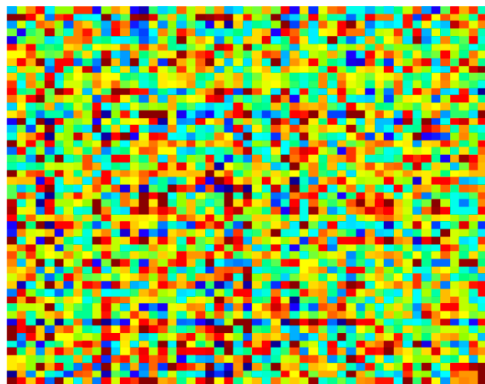
squared error $(\mathbf{X}_{ji} - (\mathbf{U}\mathbf{A})_{ji})^2$



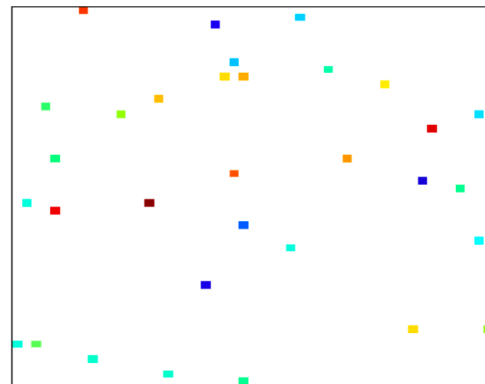
1.25% sampled

Example: 2000×2000 rank-8 random matrix

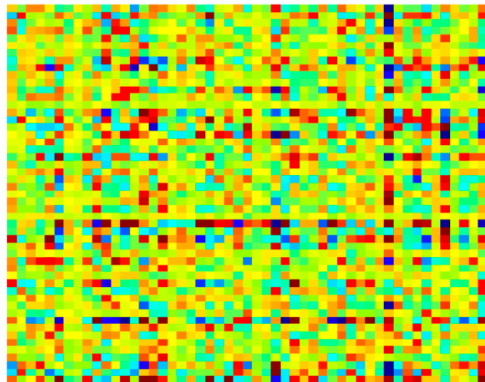
low-rank matrix \mathbf{X}



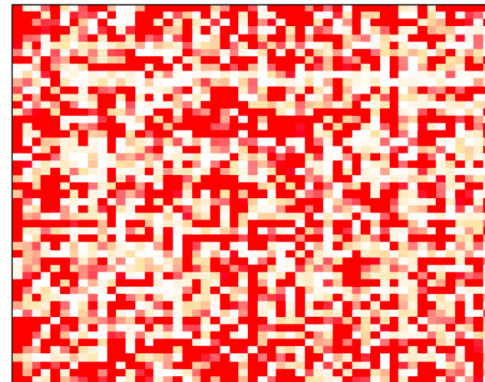
sampled matrix



Gradient descent output \mathbf{UA}



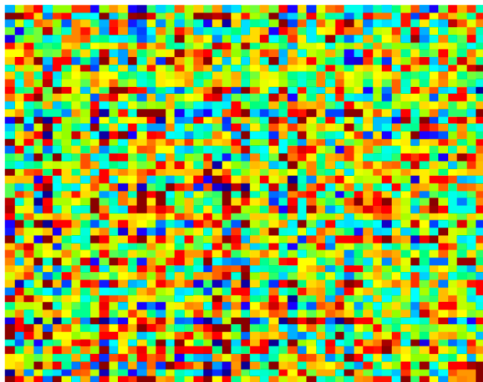
squared error $(\mathbf{X}_{ji} - (\mathbf{UA})_{ji})^2$



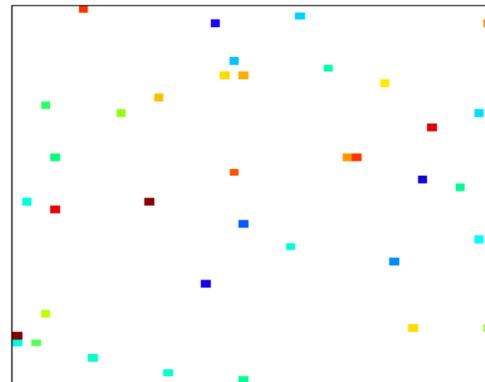
1.50% sampled

Example: 2000×2000 rank-8 random matrix

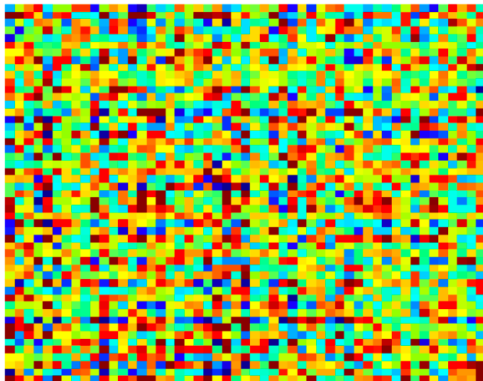
low-rank matrix \mathbf{X}



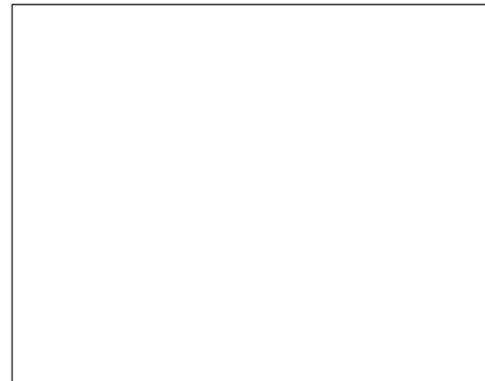
sampled matrix



Gradient descent output $\mathbf{U}\mathbf{A}$



squared error $(\mathbf{X}_{ji} - (\mathbf{U}\mathbf{A})_{ji})^2$



1.75% sampled