

Lecture 12

High Dimensional Geometry, SVD and Best-Fit Subspace

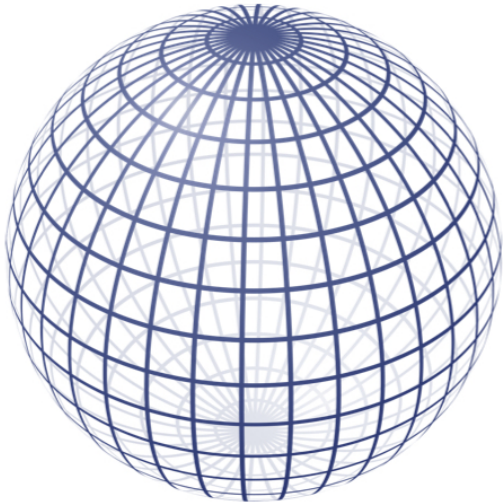
Bao-Jun Cai, 5/20/2026

Introduction to Algorithms for Data Science and Physics IMP@Fudan, 2026

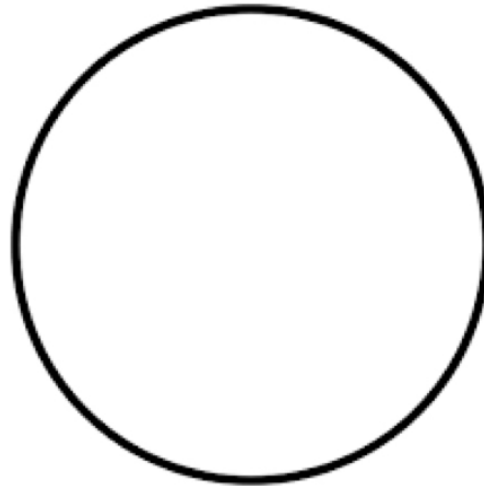
Topics of this lecture:

- volume of ball in high dimensions $\text{Vol}(B((1 - \epsilon)R)) / \text{Vol}(B(R)) = (1 - \epsilon)^d$
- law of large numbers $P(|m^{-1} \sum_{i=1}^m x^{(i)} - E[x]| \geq \epsilon) \leq \text{var}[x] / m\epsilon^2$
- random projection theorem, Johnson-Lindenstrauss lemma
- singular-value decomposition $A = U\vec{\Sigma}V^T$
- largest variance, best-fit subspace, effective approximation $A \approx \sum_{j=1}^{d_{\text{eff}}} \sigma_j \mathbf{u}_j \mathbf{u}_j^T$
- MNIST dataset “small” $y \approx$ “large” x

Circle, sphere, ...



volume: $V = \frac{4\pi R^3}{3}$
surface: $S = 4\pi R^2$



volume: $V = \pi R^2$
surface: $S = 2\pi R$



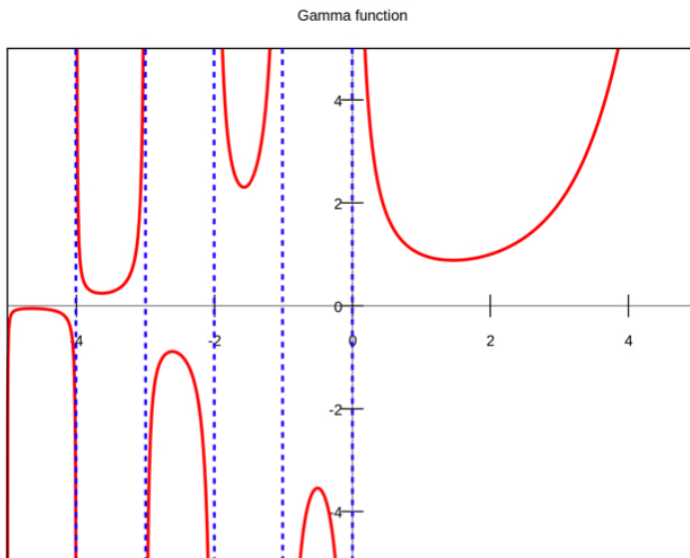
volume: $V = 2R$
surface: $S = 2$

“Sphere” in d dimensions

$$x_1^2 + x_2^2 + \dots + x_d^2 = R^2$$

$$\rightarrow \text{Vol}(d) = \frac{2\pi^{d/2}R^d}{d\Gamma(d/2)} \sim R^d \rightarrow \text{Suf}(d) = \frac{\partial \text{Vol}(d)}{\partial R} = \frac{d}{R} \text{Vol}(d)$$

Ex.: what is the expression for volume of the d-dimensional sphere?



$$\lim_{d \rightarrow \infty} \frac{\text{Suf}(d)}{\text{Vol}(d)} = \frac{d}{R} \rightarrow \infty$$

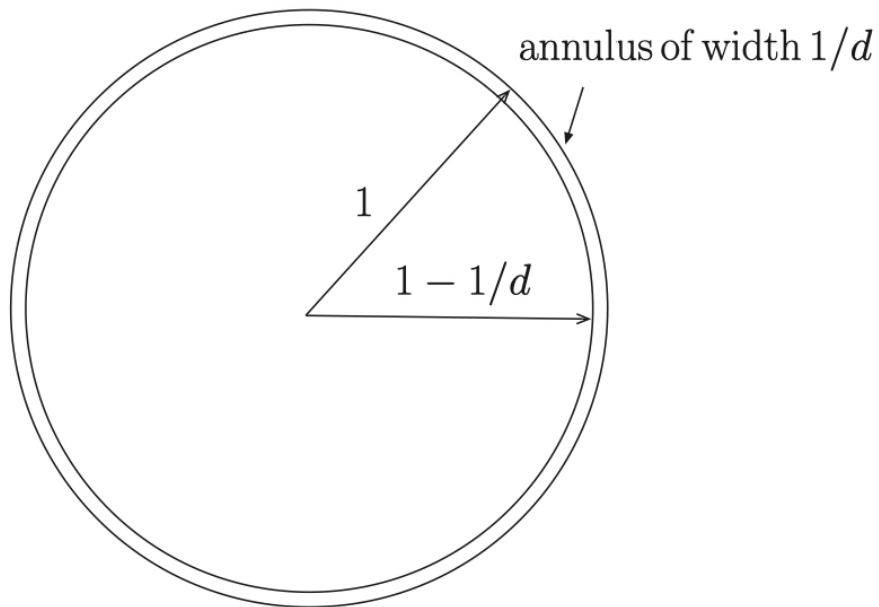
curse of dimensionality!

$$\Gamma(x) = \int_0^{\infty} t^{x-1} e^{-t} dt$$

$$\Gamma(1/2) = \sqrt{\pi}, \Gamma(1) = 1, \Gamma(3/2) = \sqrt{\pi}/2, \Gamma(2) = 1, \Gamma(5/2) = 3\sqrt{\pi}/4, \dots$$

Theorem of annulus

B : radius R ; B' : radius $R' = (1 - \epsilon)R$



$$\frac{\text{Vol}(B')}{\text{Vol}(B)} = (1 - \epsilon)^d \leq e^{-\epsilon d}$$

Ex.: prove that $(1 - \epsilon)^d \leq e^{-\epsilon d}$

Law of large numbers

$$\text{Prob} \left(\left| \frac{x^{(1)} + x^{(2)} + \dots + x^{(m)}}{m} - \mu \right| \geq \epsilon \right) \leq \frac{\sigma^2}{m\epsilon^2}$$

$$x^{(i)} \sim p(x|\mu, \sigma^2)$$

proof:

Markov inequality

+Chebyshev inequality

Markov: $P(x \geq a) \leq E[x]/a, x \geq 0$

Chebyshev: $P(|x - E[x]| \geq c) \leq \text{var}[x]/c^2$

Two random vectors in dimensions d

$$\mathbf{x}, \mathbf{y} \in \mathbf{R}^d, \quad x_i, y_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

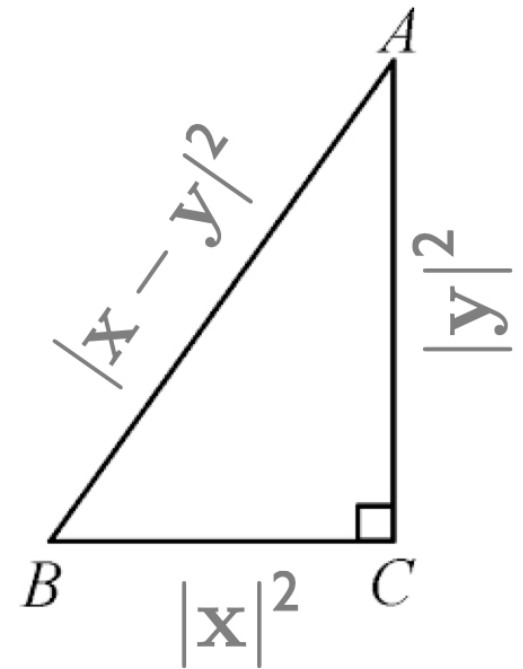
$$\text{distance} = |\mathbf{x} - \mathbf{y}|^2 = \sum_{i=1}^d (x_i - y_i)^2$$

$$E[\mathbf{x}^2] = \sum_{i=1}^d E[x_i^2] = d \text{var}[x_i] = d$$

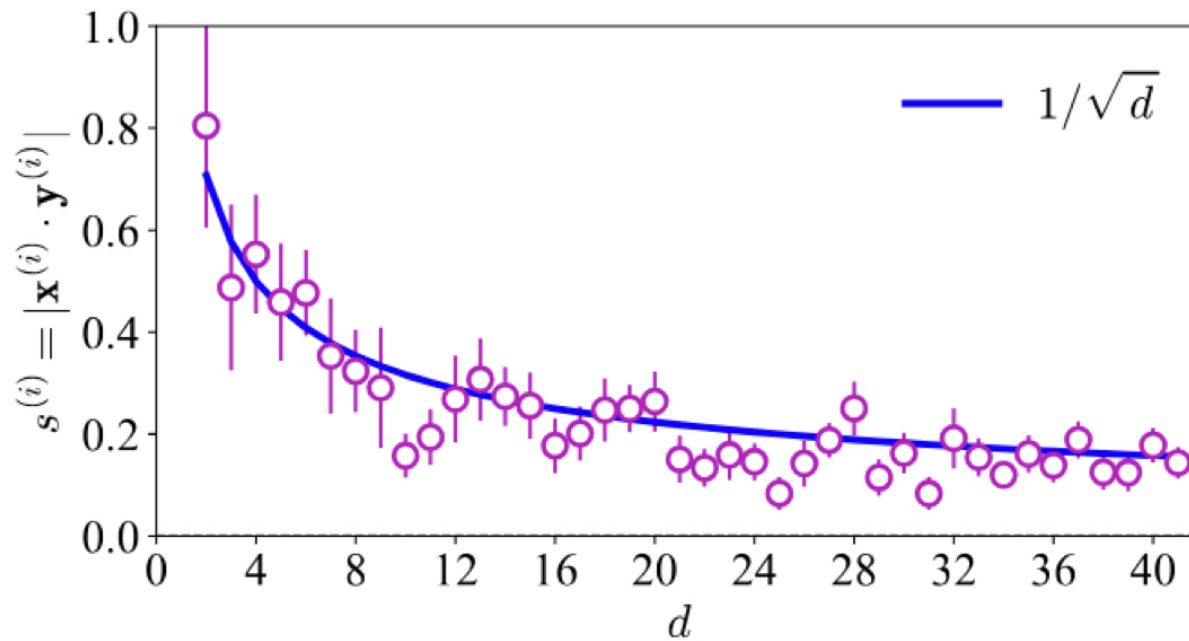
$$= \sum_{i=1}^d (E[x_i^2] + E[y_i^2] - 2E[x_i]E[y_i])$$

$$= \sum_{i=1}^d (\text{var}[x_i] + \text{var}[y_i] - 2E[x_i]E[y_i]) = 2d$$

$$|\mathbf{x} - \mathbf{y}|^2 \approx |\mathbf{x}|^2 + |\mathbf{y}|^2$$



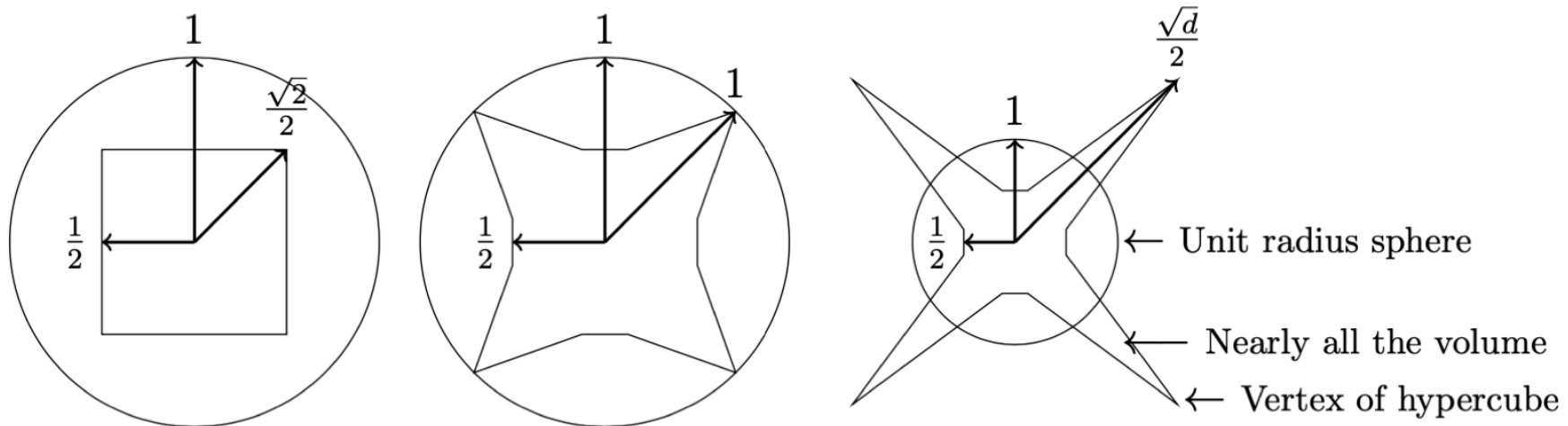
Simulation for vectors in high-D



$$\|\mathbf{x}^{(i)}\| \geq 1 - \frac{2 \ln m}{d}, \quad \mathbf{x}^{(i), \top} \mathbf{x}^{(j)} = \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)} \leq \sqrt{\frac{6 \ln m}{d-1}}$$

How to understand the result?

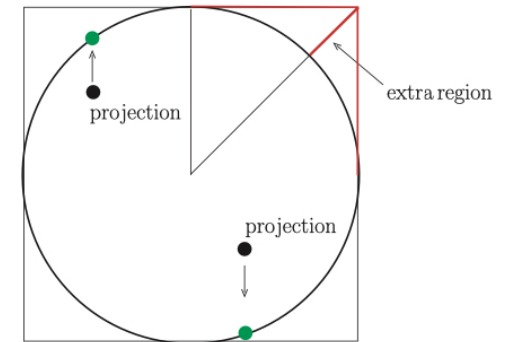
A. Blum, J. Hopcroft, and R. Kannan, *Foundations of Data Science*, Cambridge, 2020.



one can show that most of the volume of the unit ball in high dimensions is concentrated near its “equator”, for for any unit-length vector \mathbf{a} defining the “north pole”, most of the volume of the unit ball lies in the thin slab of points whose inner product with \mathbf{a} has magnitude $O(d^{-1/2})$

Generate uniform random samples on surface/volume

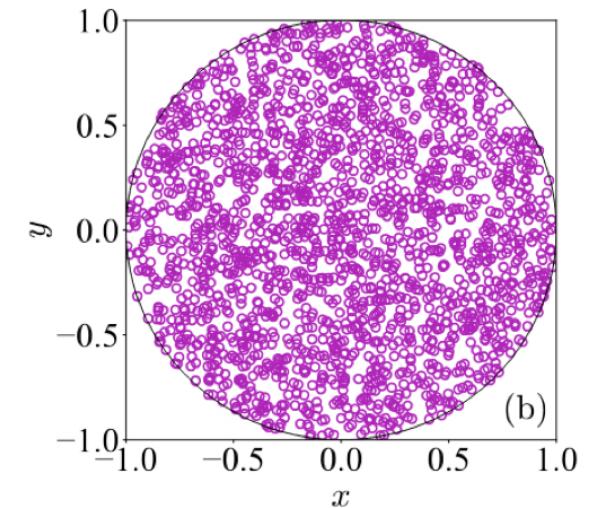
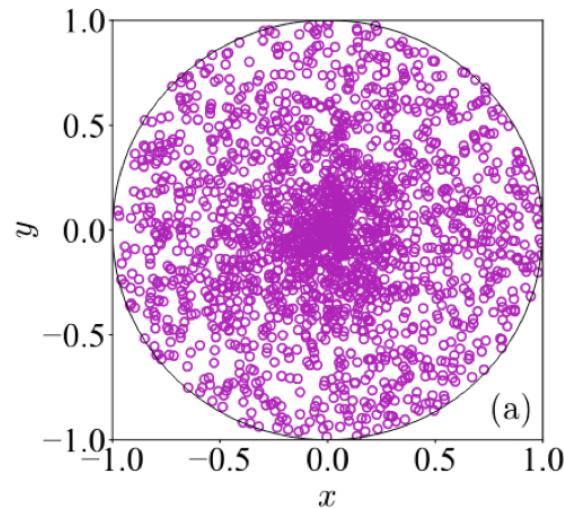
$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{x_1^2 + x_2^2 + \dots + x_d^2}{2}\right), \quad x_i \sim \text{Box-Muller}$$



surface: $\mathbf{x}/|\mathbf{x}|$

$$\rho(r) = cr^{d-1}$$

volume: $\rho(r)\mathbf{x}/|\mathbf{x}|$



Random projection theorem

$$\mathbf{a}^{(i)} \in \mathbb{R}^d, i = 1 \sim s, \mathbf{a}_j^{(i)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\mathbb{R}^d \rightarrow \mathbb{R}^s : \mathbf{f}(\mathbf{b}) = (\mathbf{a}^{(1)} \cdot \mathbf{b}, \dots, \mathbf{a}^{(s)} \cdot \mathbf{b})^\top$$

$$|\mathbf{f}(\mathbf{b})| \approx \sqrt{s}|\mathbf{b}| : \text{Prob} (||\mathbf{f}(\mathbf{b})| - \sqrt{s}|\mathbf{b}|| \geq \epsilon\sqrt{s}|\mathbf{b}|) \leq 3e^{-\phi s \epsilon^2}, \phi > 0$$



$$\text{Prob} \sim 1 - 2 \binom{m}{2} e^{-s\epsilon^2/8}, s \gtrsim \frac{16}{\epsilon^2} \ln \frac{m}{\epsilon}$$

$$(1 - \epsilon)\sqrt{s} |\mathbf{a}^{(i)} - \mathbf{a}^{(j)}| \leq |\mathbf{f}(\mathbf{a}^{(i)}) - \mathbf{f}(\mathbf{a}^{(j)})| \leq (1 + \epsilon)\sqrt{s} |\mathbf{a}^{(i)} - \mathbf{a}^{(j)}|$$

Johnson-Lindenstrauss lemma

Singular-value-decomposition (SVD)

d unknowns, m equations

$$\mathbf{Ax} = \mathbf{b}, \mathbf{A} \in \mathbf{R}^{m \times d}$$

$$a_1^{(1)}x_1 + a_2^{(1)}x_2 + \cdots + a_d^{(1)}x_d = b_1$$

$$a_1^{(2)}x_1 + a_2^{(2)}x_2 + \cdots + a_d^{(2)}x_d = b_2$$

\vdots

$$a_1^{(m)}x_1 + a_2^{(m)}x_2 + \cdots + a_d^{(m)}x_d = b_m$$

SVD of \mathbf{A}

$$\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_d$$

$$\mathbf{A} = \mathbf{U}\vec{\Sigma}\mathbf{V}^\top$$

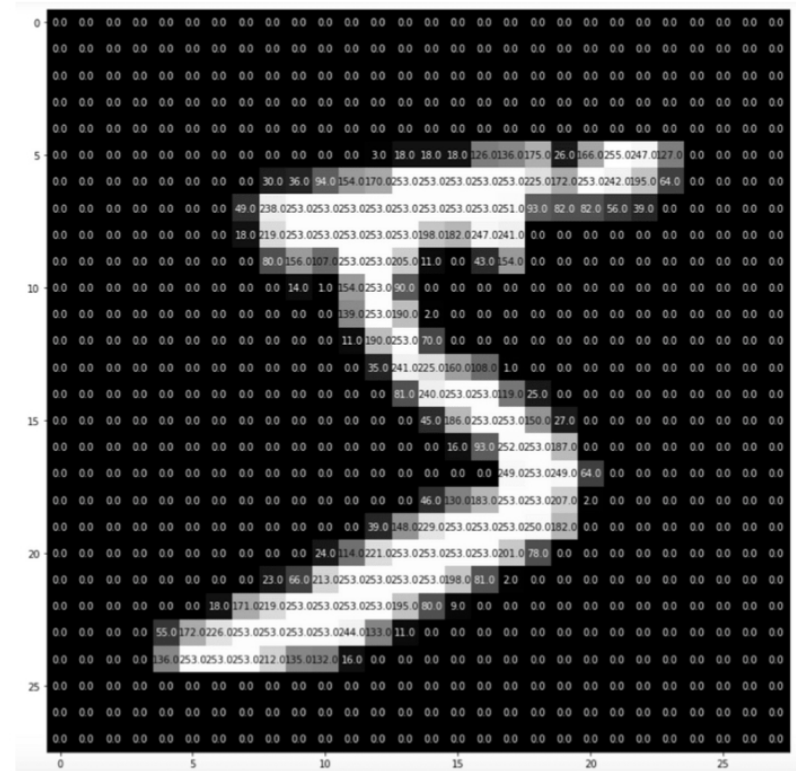
$$\mathbf{U} \in \mathbf{R}^{m \times d}, \vec{\Sigma} \in \mathbf{R}^{d \times d}, \mathbf{V} \in \mathbf{R}^{d \times d}$$

$$a_j^{(i)} = \sum_{k=1}^d \sigma_k U_{ik} V_{jk}$$

Another motivation for SVD

3 6 8 1 7 9 6 6 9 1
 6 7 5 7 8 6 3 4 8 5
 2 1 7 9 7 1 2 8 4 5
 4 8 1 9 0 1 8 8 9 4
 7 6 1 8 6 4 1 5 6 0
 7 5 9 2 6 5 8 1 9 7
 2 2 2 2 2 3 4 4 8 0
 0 2 3 8 0 7 3 8 5 7
 0 1 4 6 4 6 0 2 4 3
 7 1 2 8 7 6 9 8 6 1

$$\mathbf{x} \in \mathbb{R}^{28 \times 28} \sim \mathbb{R}^{784} \sim \mathbb{R}^b$$




“small” $\mathbf{y} \approx \mathbf{x}$?

Best-fit subspace

$$\text{square: } \mathbf{A} = \sum_{j=1}^d \sigma_j \mathbf{u}_j \mathbf{u}_j^T$$

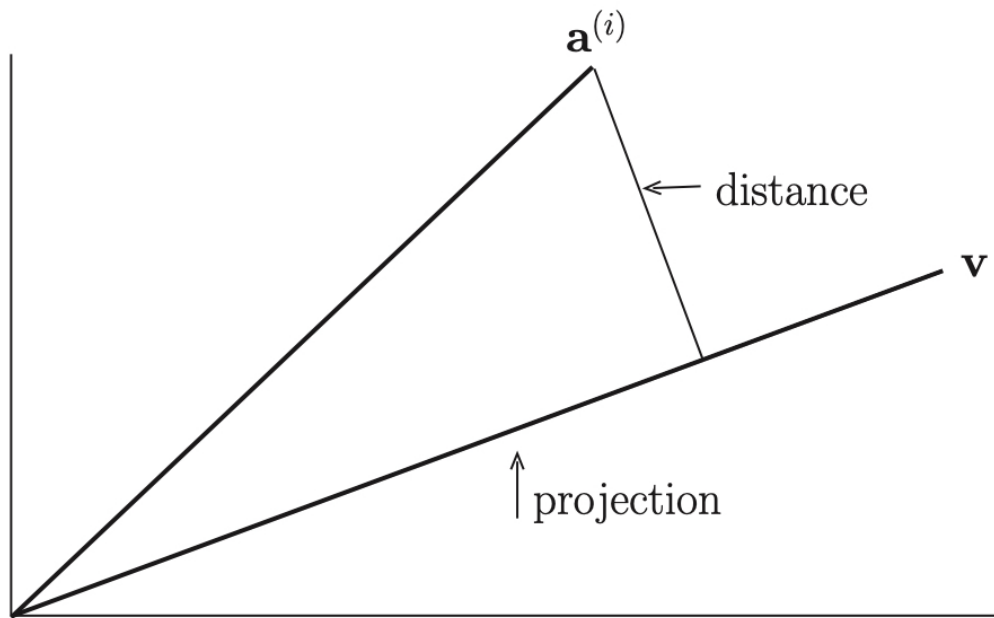
separable of SVs

$$\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_{d_{\text{eff}}} \gg \sigma_{d_{\text{eff}}+1} \geq \cdots \geq \sigma_d$$


$$\mathbf{A} \approx \sum_{j=1}^{d_{\text{eff}}} \sigma_j \mathbf{u}_j \mathbf{u}_j^T$$

d_{eff} could be far smaller than d

Geometric meaning of best-fit subspace



$$\mathbf{v}_1 = \operatorname{argmax}_{|\mathbf{v}|=1} |\mathbf{A}\mathbf{v}|^2$$

$$\mathbf{v}_2 = \operatorname{argmax}_{\mathbf{v} \perp \mathbf{v}_1, |\mathbf{v}|=1} |\mathbf{A}\mathbf{v}|^2$$

\vdots

$$[\text{distance of point to line}]^2 = \mathbf{a}_1^{(i),2} + \mathbf{a}_2^{(i),2} + \dots + \mathbf{a}_d^{(i),2} - \overbrace{[\text{length of projection}]^2}^{|\mathbf{A}\mathbf{v}|^2}$$

Directions with largest variance

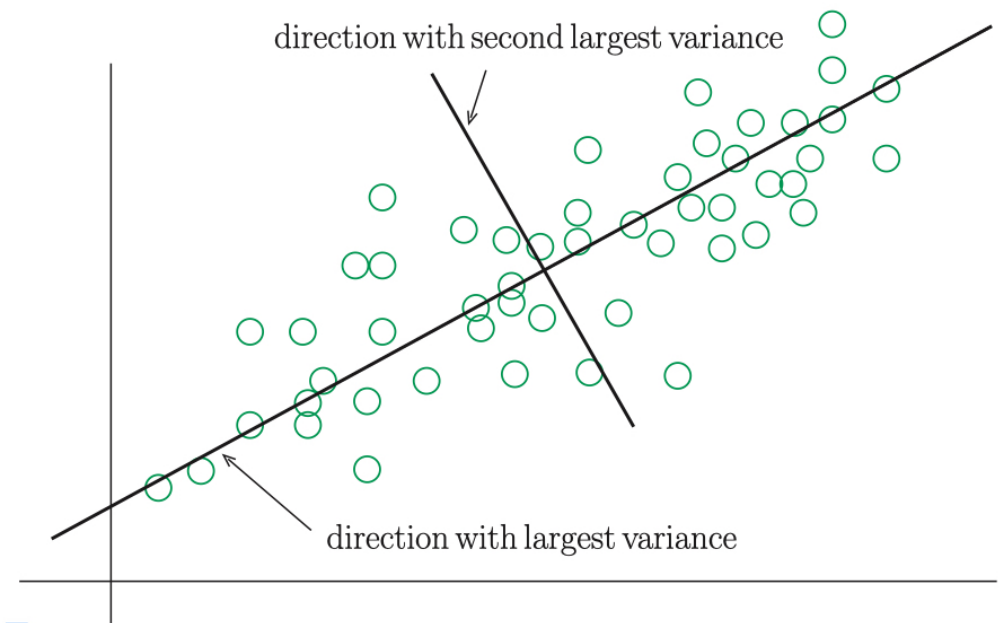
definition:

for $\mathbf{x} \in \mathbb{R}^b$, define the first s principle components y_j 's such that $\text{var}[y_1] \geq \text{var}[y_2] \cdots \geq \text{var}[y_s]$

$$y_1 = \mathbf{u}_1^\top \mathbf{x}, \quad \mathbf{u}_1^\top \mathbf{u}_1 = 1, \quad \mathbf{u}_1 \in \mathbb{R}^b$$

$$\mathbf{u}_1^* = \operatorname{argmax}_{|\mathbf{u}|=1} \text{var}[\mathbf{u}^\top \mathbf{x}]$$

Ex.: find out \mathbf{u}_1 and generally \mathbf{u}_j



Stacking principal components

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_s \end{pmatrix} \in \mathbb{R}^s, \mathbf{U} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_s) \in \mathbb{R}^{b \times s}$$

$$\mathbf{y} = \mathbf{U}^\top \mathbf{x}$$



$$\vec{\Sigma}_{\mathbf{y}} = \mathbf{U}^\top \vec{\Sigma}_{\mathbf{x}} \mathbf{U}$$

$$\widehat{\vec{\Sigma}}_m = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)} \mathbf{x}^{(i)\top} = \frac{1}{m} \mathbf{X} \mathbf{X}^\top \approx \vec{\Sigma}_{\mathbf{x}}$$

An optimization problem for reconstruction

$$\min_{\mathbf{U}^T \mathbf{U} = \mathbf{I}_s, \{\mathbf{y}^{(i)}\}} \sum_{i=1}^m \|\mathbf{x}^{(i)} - \mathbf{U} \mathbf{y}^{(i)}\|^2, \quad \sum_{i=1}^m \mathbf{y}^{(i)} = \mathbf{0}$$

original data

reconstruction of $\mathbf{x}^{(i)}$

error $\sum_{i=s+1}^b \sigma_i^2$

$$\min_{\mathbf{U}^T \mathbf{U} = \mathbf{I}_s} \sum_{i=1}^m \|\mathbf{x}^{(i)} - \mathbf{U} \mathbf{U}^T \mathbf{x}^{(i)}\|^2$$

Ex.: SVD of $\mathbf{X} \mathbf{X}^T$ is necessary, why?

Notations for image application(s)



$$\mathbf{X}_{(k)} = \left(\mathbf{x}_{(k)}^{(1)}, \dots, \mathbf{x}_{(k)}^{(m)} \right) \in \mathbf{R}^{b \times m}$$

$$\vec{\mu}_{(k)} = m^{-1} \sum_{i=1}^m \mathbf{x}_{(k)}^{(i)}$$

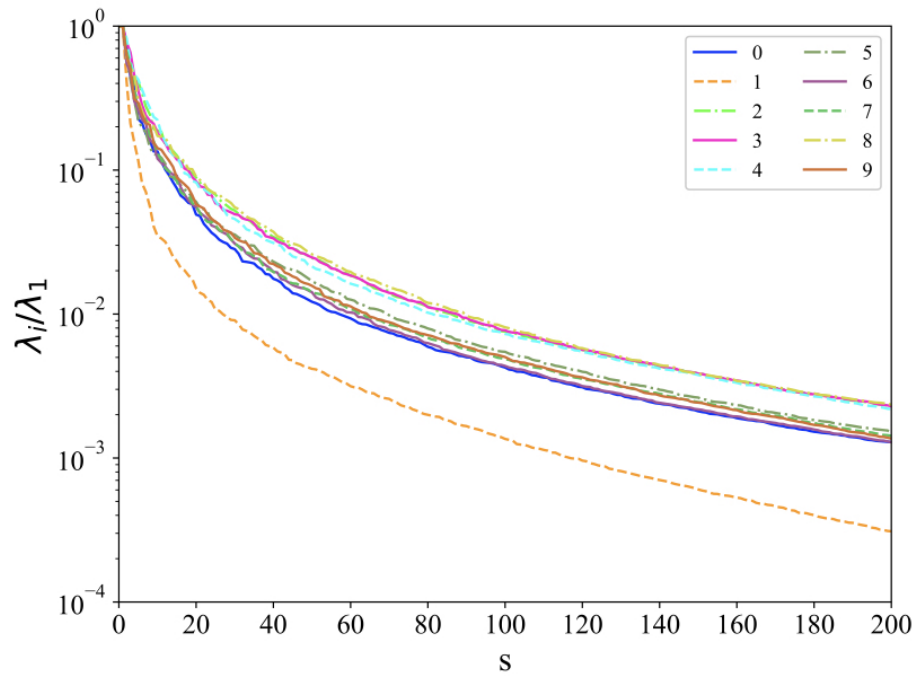
$$\hat{\mathbf{X}}_{(k)} = \mathbf{X}_{(k)} - \vec{\mu}_{(k)}$$

$$\hat{\mathbf{X}}_{(k)} \hat{\mathbf{X}}_{(k)}^{\top} = \mathbf{U}_{(k)} \text{diag}(\lambda_{1,(k)}, \dots, \lambda_{784,(k)}) \mathbf{U}_{(k)}^{\top}$$

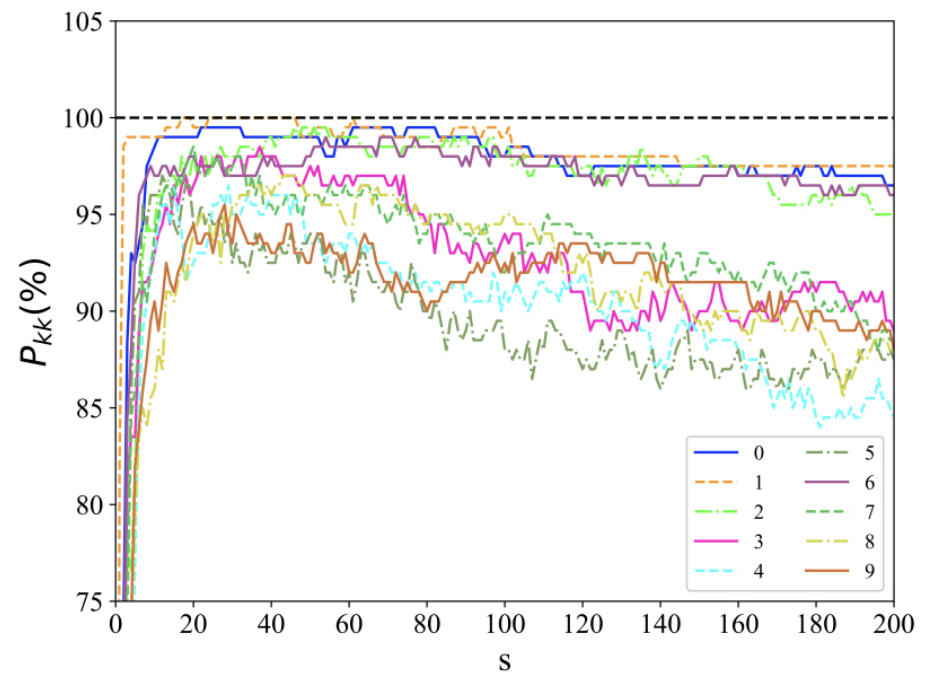
$$\Delta_{(k)} = \text{tr} \left(\mathbf{U}_{(k)}^{(s),\top} \hat{\mathbf{f}}_{(k)} \hat{\mathbf{f}}_{(k)}^{\top} \mathbf{U}_{(k)}^{(s)} \right)$$

$$k^* = \text{argmax}_k \Delta_{(k)}$$

Sis and correctness

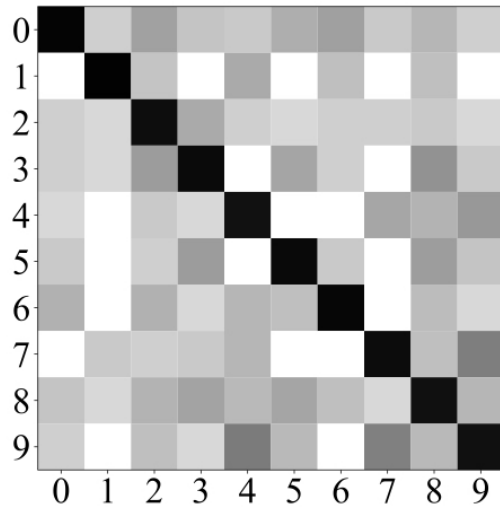


$$\lambda_i = \sigma_i^2$$

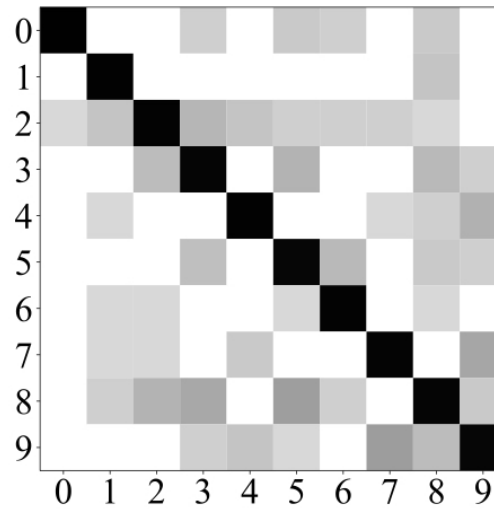


Correlation(s)

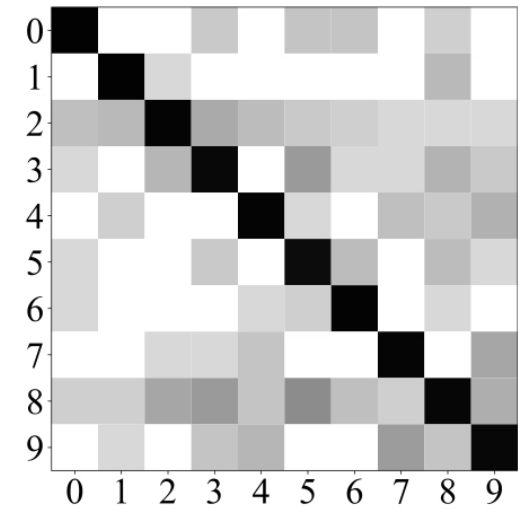
sum rule for each row $\sum_{k'=0}^9 P_{kk'} = 1$



s=5



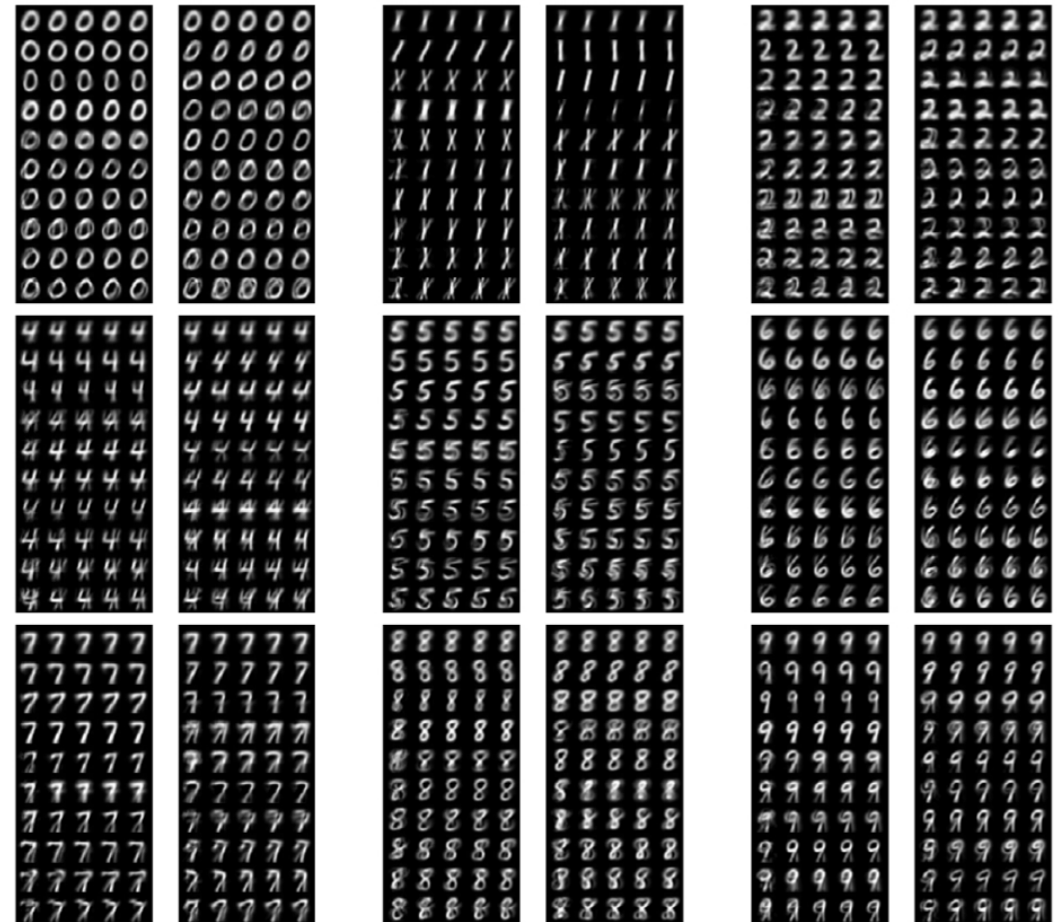
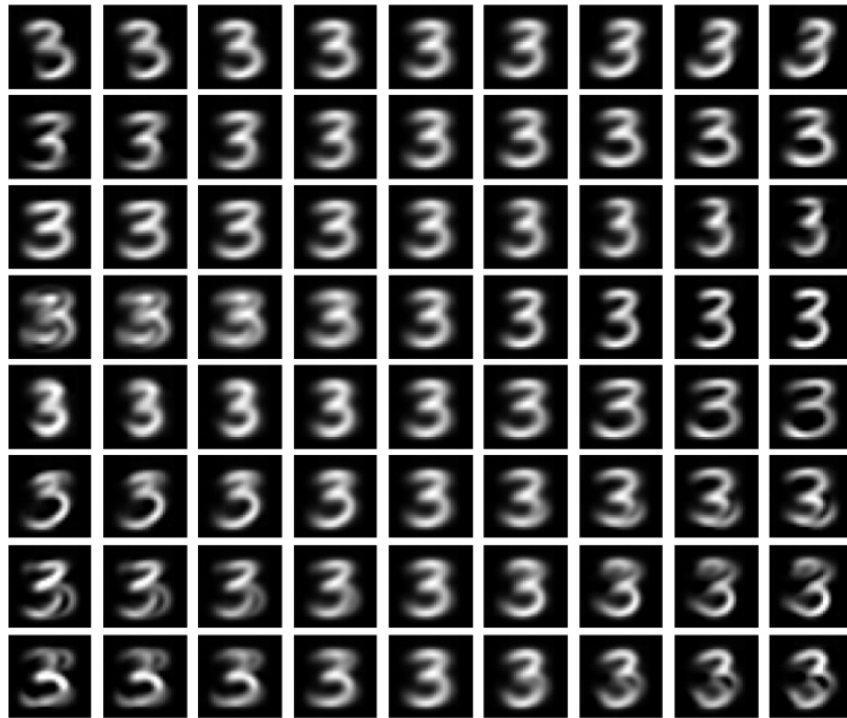
s=40



s=100

20

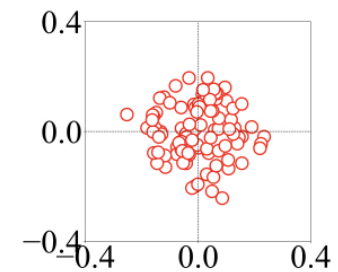
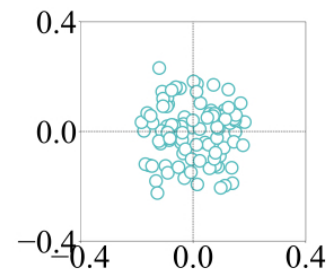
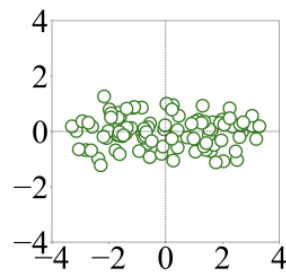
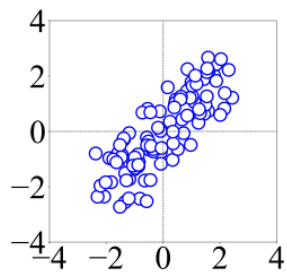
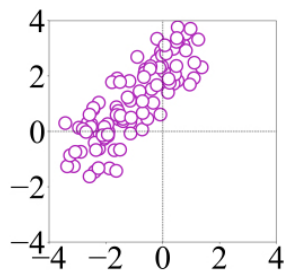
Feature(s) discovered by PCA



$$\hat{\mu}_m + \chi \mathbf{u}_i, i = 1 \sim 8$$

Application: whitening & SPCA*

$$\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} - \vec{\mu}, \quad \mathbf{x}^{(i)} \leftarrow \mathbf{V}^\top \mathbf{x}^{(i)}, \quad \mathbf{x}_j^{(i)} \leftarrow \mathbf{x}^{(i)} / \sigma_j, \quad \mathbf{x}^{(i)} \leftarrow \mathbf{V} \mathbf{x}^{(i)}$$



sparse PCA (SPCA):

$$\hat{\mathbf{B}} = \operatorname{argmin}_{\mathbf{A}, \mathbf{B}} \left[\sum_{i=1}^m \left\| \mathbf{x}^{(i)} - \mathbf{A} \mathbf{B}^\top \mathbf{x}^{(i)} \right\|^2 + \lambda \sum_{j=1}^s \left\| \vec{\beta}_j \right\|^2 + \sum_{j=1}^s \lambda'_j \left\| \vec{\beta}_j \right\|_1 \right]$$