



# Introduction to AI for postgraduate students

Lecture Note 2-1  
Linear Algebra

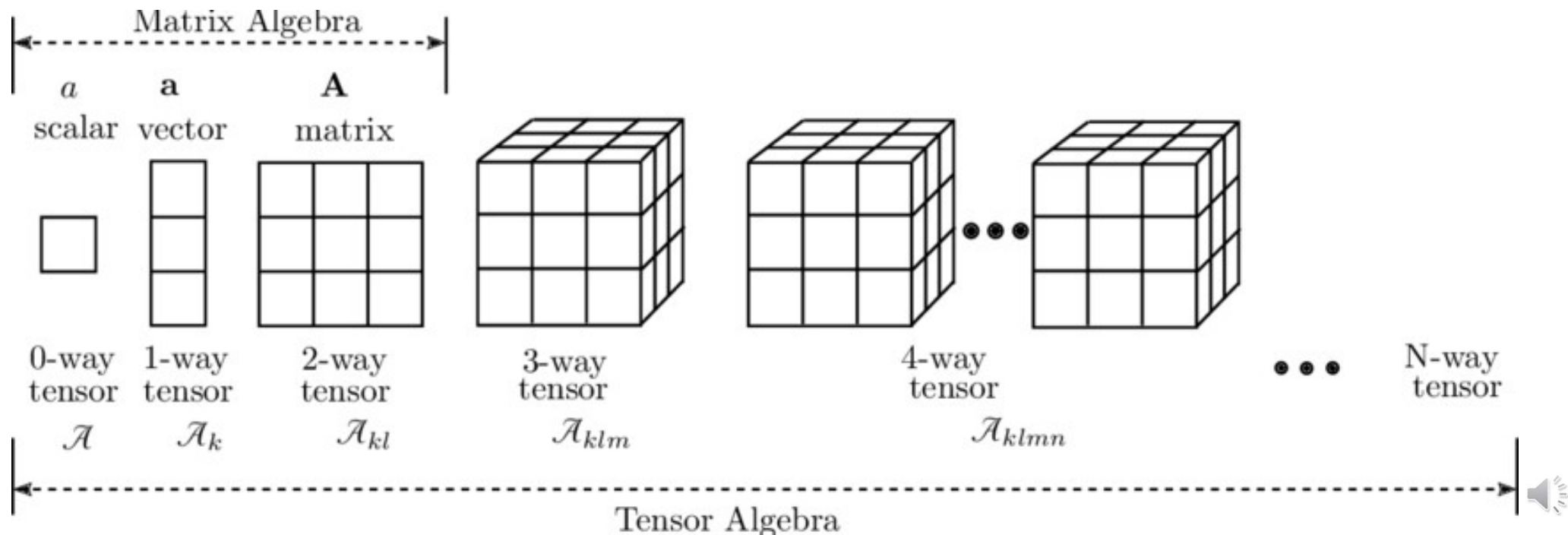
Hyun Jong Yang ([hyunyang@postech.ac.kr](mailto:hyunyang@postech.ac.kr))

**POSTECH**



# Definition

- Scalar: single number
- Vector: array of numbers
- Matrix: 2-D array of numbers
- Tensors: array with more than two axes



# Operations

## Transpose

$$A = \begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \\ A_{3,1} & A_{3,2} \end{bmatrix} \Rightarrow A^T = \begin{bmatrix} A_{1,1} & A_{2,1} & A_{3,1} \\ A_{1,2} & A_{2,2} & A_{3,2} \end{bmatrix}$$

## Broadcasting

$$C = A + b$$

Matrix

vector



$$C_{i,j} = A_{i,j} + b_j$$

`np.arange(3)+5`

$$\begin{bmatrix} 0 & 1 & 2 \end{bmatrix} + \begin{bmatrix} 5 & 5 & 5 \end{bmatrix} = \begin{bmatrix} 5 & 6 & 7 \end{bmatrix}$$

`np.ones((3,3))+np.arange(3)`

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} + \begin{bmatrix} 0 & 1 & 2 \\ 0 & 1 & 2 \\ 0 & 1 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \\ 1 & 2 & 3 \end{bmatrix}$$

`np.arange(3).reshape((3,1))+np.arange(3)`

$$\begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix} + \begin{bmatrix} 0 & 1 & 2 \\ 0 & 1 & 2 \\ 0 & 1 & 2 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 2 \\ 1 & 2 & 3 \\ 2 & 3 & 4 \end{bmatrix}$$



# Multiplying Matrices and Vectors

- Matrix multiplication:

- Dimensions of  $\mathbf{A}, \mathbf{B}, \mathbf{C}$ :  $(m, n), (n, p), (m, p)$

$$\mathbf{C} = \mathbf{AB}$$

$$C_{i,j} = \sum_k A_{i,k} B_{k,j}$$

- Dot product:  $\mathbf{x} \cdot \mathbf{y} = \mathbf{x}^T \mathbf{y}$

- Hadamard product

$$\begin{matrix} & G & & & H & & & N \\ \begin{bmatrix} 3 & 5 & 7 \\ 4 & 9 & 8 \end{bmatrix} & \circ & \begin{bmatrix} 1 & 6 & 3 \\ 0 & 2 & 9 \end{bmatrix} & = & \begin{bmatrix} 3 \times 1 & 5 \times 6 & 7 \times 3 \\ 4 \times 0 & 9 \times 2 & 8 \times 9 \end{bmatrix} \end{matrix}$$



# Multiplying Matrices and Vectors

- Properties of matrix multiplication

- Distributive

$$A(B + C) = AB + AC$$

- Associative

$$A(BC) = (AB)C$$

- Not commutative for matrix multiplication

$$AB \neq BA$$

- But commutative for dot product

$$x^\top y = y^\top x$$



# Identity and Inverse Matrices

- Identity matrix

$$\forall \mathbf{x} \in \mathbb{R}^n, \mathbf{I}_n \mathbf{x} = \mathbf{x}, \quad \mathbf{I}_n \in \mathbb{R}^{n \times n}$$

$$\mathbf{I}_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

- Inverse matrix

$$\mathbf{A}^{-1} \mathbf{A} = \mathbf{I}_n$$

- Solving a linear equation

$$\mathbf{A} \mathbf{x} = \mathbf{b}$$

$$\mathbf{A}^{-1} \mathbf{A} \mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$$

$$\mathbf{I}_n \mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$$

$$\mathbf{x} = \mathbf{A}^{-1} \mathbf{b}.$$



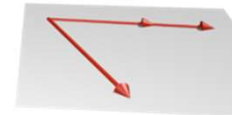
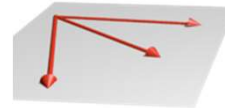
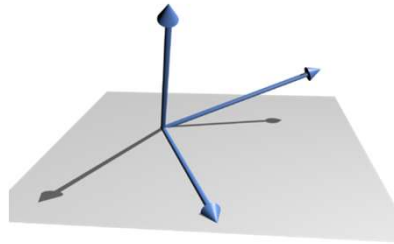
# Linear Dependence

- Linearly dependent

A sequence of vectors  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$  from a vector space  $V$  is said to be *linearly dependent*, if there exist scalars  $a_1, a_2, \dots, a_k$ , not all zero, such that

$$a_1 \mathbf{v}_1 + a_2 \mathbf{v}_2 + \dots + a_k \mathbf{v}_k = \mathbf{0},$$

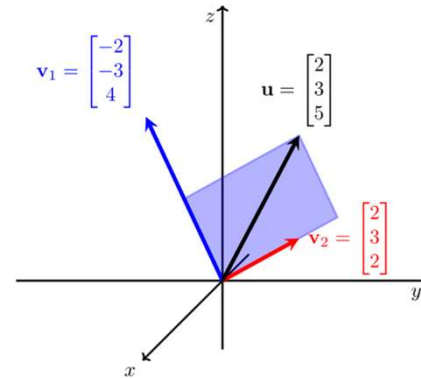
where  $\mathbf{0}$  denotes the zero vector.



# Span

- Span

$$\text{span}(S) = \left\{ \sum_{i=1}^k \lambda_i v_i \mid k \in \mathbb{N}, v_i \in S, \lambda_i \in K \right\}$$





# Norms

- $L^p$  norm for  $p \in \mathbb{R}, p \geq 1$ :

$$\|\mathbf{x}\|_p = \left( \sum_i |x_i|^p \right)^{\frac{1}{p}}$$

- $L^2$  norm: Euclidean norm

- $L^1$  norm:  $\|\mathbf{x}\|_1 = \sum_i |x_i|$

- $L^\infty$  norm:  $\|\mathbf{x}\|_\infty = \max_i |x_i|$

- Frobenius norm:  $\|A\|_F = \sqrt{\sum_{i,j} A_{i,j}^2}$



# Special Kinds of Matrices and Vectors

- Diagonal matrix

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & -2 \end{bmatrix}$$

- Symmetric matrix

$$\mathbf{A} = \mathbf{A}^\top$$

- Unit vector

$$\|\mathbf{x}\|_2 = 1$$

- Orthogonality

$$\mathbf{x}^\top \mathbf{y} = 0$$

- Orthogonal matrix

$$\mathbf{A}^\top \mathbf{A} = \mathbf{A} \mathbf{A}^\top = \mathbf{I}$$

$$\mathbf{A}^{-1} = \mathbf{A}^\top$$



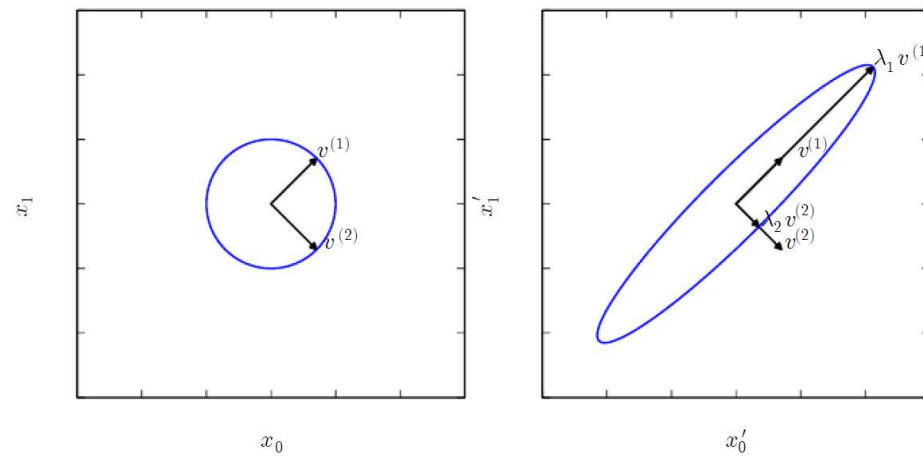
# Eigen Decomposition

- Square matrix:  $A$ , Eigen vector:  $v$ , Eigen value:  $\lambda$

$$Av = \lambda v$$

$$A[v_1 \ v_2] = [v_1 \ v_2] \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$$

$$A = V \text{diag}(\lambda) V^{-1}$$



# Singular Value Decomposition

$$A = U D V^T$$

Singular values  
(diagonal matrix)

Left singular matrix  
(orthogonal matrix)

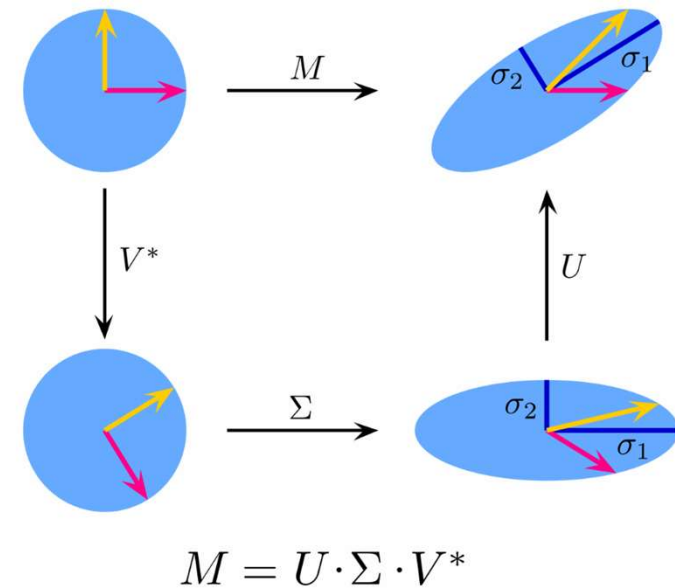
Right singular matrix  
(orthogonal matrix)

$$M = U \Sigma V^*$$

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

$$U U^* = I_m$$

$$V V^* = I_n$$



# Moore-Penrose Pseudoinverse

- Moore-Penrose pseudoinverse matrix

$$\mathbf{A}^+ = \lim_{\alpha \searrow 0} (\mathbf{A}^\top \mathbf{A} + \alpha \mathbf{I})^{-1} \mathbf{A}^\top$$

$$\mathbf{A}^+ = \mathbf{V} \mathbf{D}^+ \mathbf{U}^\top$$

- $\mathbf{D}^+$  of a diagonal matrix  $\mathbf{D}$  is obtained by taking the reciprocal of its nonzero elements then taking the transpose of the resulting matrix.



# Trace Operator

- Definition

$$\text{Tr}(\mathbf{A}) = \sum_i \mathbf{A}_{i,i}$$

- Frobenius norm

$$\|\mathbf{A}\|_F = \sqrt{\text{Tr}(\mathbf{A}\mathbf{A}^\top)}$$

- Properties

$$\text{Tr}(\mathbf{A}) = \text{Tr}(\mathbf{A}^\top)$$

$$\text{Tr}(\mathbf{A}\mathbf{B}\mathbf{C}) = \text{Tr}(\mathbf{C}\mathbf{A}\mathbf{B}) = \text{Tr}(\mathbf{B}\mathbf{C}\mathbf{A})$$



# Determinant

- Definition

$$\det(A) = \prod_{i=1}^n \lambda_i = \lambda_1 \lambda_2 \cdots \lambda_n$$

- Examples

$$|A| = \begin{vmatrix} a & b \\ c & d \end{vmatrix} = ad - bc.$$

$$\begin{aligned} |A| &= \begin{vmatrix} a & b & c \\ d & e & f \\ g & h & i \end{vmatrix} = a \begin{vmatrix} e & f \\ h & i \end{vmatrix} - b \begin{vmatrix} d & f \\ g & i \end{vmatrix} + c \begin{vmatrix} d & e \\ g & h \end{vmatrix} \\ &= aei + bfg + cdh - ceg - bdi - afh. \end{aligned}$$



# Vector Calculus

- Gradient of  $f$  (with respect to  $A \in \mathbb{R}^{m \times n}$ )

$$\nabla_A f(A) \in \mathbb{R}^{m \times n} = \begin{bmatrix} \frac{\partial f(A)}{\partial A_{11}} & \frac{\partial f(A)}{\partial A_{12}} & \cdots & \frac{\partial f(A)}{\partial A_{1n}} \\ \frac{\partial f(A)}{\partial A_{21}} & \frac{\partial f(A)}{\partial A_{22}} & \cdots & \frac{\partial f(A)}{\partial A_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f(A)}{\partial A_{m1}} & \frac{\partial f(A)}{\partial A_{m2}} & \cdots & \frac{\partial f(A)}{\partial A_{mn}} \end{bmatrix}$$

- Linearity of Gradient

$$\nabla_x (f(x) + g(x)) = \nabla_x f(x) + \nabla_x g(x).$$

$$\text{For } t \in \mathbb{R}, \nabla_x (t f(x)) = t \nabla_x f(x).$$





# Vector Calculus

- Examples

$$f(x) = b^T x \quad \Rightarrow \quad f(x) = \sum_{i=1}^n b_i x_i \quad \Rightarrow \quad \frac{\partial f(x)}{\partial x_k} = \frac{\partial}{\partial x_k} \sum_{i=1}^n b_i x_i = b_k. \quad \Rightarrow \quad b$$

$$f(x) = x^T A x \quad \Rightarrow \quad f(x) = \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j$$

$$\begin{aligned} \Rightarrow \quad \frac{\partial f(x)}{\partial x_k} &= \frac{\partial}{\partial x_k} \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j \\ &= \frac{\partial}{\partial x_k} \left[ \sum_{i \neq k} \sum_{j \neq k} A_{ij} x_i x_j + \sum_{i \neq k} A_{ik} x_i x_k + \sum_{j \neq k} A_{kj} x_k x_j + A_{kk} x_k^2 \right] \\ &= \sum_{i \neq k} A_{ik} x_i + \sum_{j \neq k} A_{kj} x_j + 2A_{kk} x_k \\ &= \sum_{i=1}^n A_{ik} x_i + \sum_{j=1}^n A_{kj} x_j = 2 \sum_{i=1}^n A_{ki} x_i, \quad \Rightarrow \quad 2Ax \end{aligned}$$



# Vector Calculus

- Hessian

$$\nabla_x^2 f(x) \in \mathbb{R}^{n \times n} = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_n \partial x_1} & \frac{\partial^2 f(x)}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_n^2} \end{bmatrix}$$

- Hessian is symmetric since  $\frac{\partial^2 f(x)}{\partial x_i \partial x_j} = \frac{\partial^2 f(x)}{\partial x_j \partial x_i}$ .

- Example

$$f(x) = x^T A x \quad \Rightarrow \quad f(x) = \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j$$

$$\Rightarrow \frac{\partial^2 f(x)}{\partial x_k \partial x_\ell} = \frac{\partial}{\partial x_k} \left[ \frac{\partial f(x)}{\partial x_\ell} \right] = \frac{\partial}{\partial x_k} \left[ 2 \sum_{i=1}^n A_{li} x_i \right] = 2A_{\ell k} = 2A_{k\ell}. \quad \Rightarrow \quad \nabla_x^2 x^T A x = 2A$$



# Principal Components Analysis

- Original Data:  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  in  $\mathbb{R}^n$
- Goal: apply lossy compression to these points, which requires less memory losing as little precision as possible to get the mapped code vectors  $\mathbf{c}^{(i)} \in \mathbb{R}^l$  ( $l < n$ )
- Assumption: *Linear* encoding/decoding

Encoding:  $f(\mathbf{x}) = \mathbf{c}$

Encoding function

Decoding:  $\mathbf{x} \approx g(f(\mathbf{x})) = D\mathbf{c}$

Decoding function

- Part1: Optimizing the code  $\mathbf{c}$  vector for given  $\mathbf{x}$

$$\begin{aligned}\mathbf{c}^* &= \arg \min_{\mathbf{c}} \|\mathbf{x} - g(\mathbf{c})\|_2 \\ &= \arg \min_{\mathbf{c}} \|\mathbf{x} - g(\mathbf{c})\|_2^2.\end{aligned}$$



# Principal Components Analysis

- Part1: Optimizing the code  $c$  vector for given  $x$

$$\begin{aligned} (x - g(c))^T (x - g(c)) &= x^T x - x^T g(c) - g(c)^T x + g(c)^T g(c) \\ &= x^T x - 2x^T g(c) + g(c)^T g(c) \end{aligned}$$

➔

$$\begin{aligned} c^* &= \arg \min_c -2x^T g(c) + g(c)^T g(c) \\ &= \arg \min_c -2x^T Dc + c^T D^T Dc \\ &= \arg \min_c -2x^T Dc + c^T I_l c \\ &= \arg \min_c -2x^T Dc + c^T c \end{aligned}$$



# Principal Components Analysis

$$\nabla_{\mathbf{c}}(-2\mathbf{x}^\top \mathbf{D}\mathbf{c} + \mathbf{c}^\top \mathbf{c}) = \mathbf{0}$$

$$-2\mathbf{D}^\top \mathbf{x} + 2\mathbf{c} = \mathbf{0} \quad \Rightarrow \quad f(\mathbf{x}) = \mathbf{D}^\top \mathbf{x} \quad \Rightarrow \quad r(\mathbf{x}) = g(f(\mathbf{x})) = \mathbf{D}\mathbf{D}^\top \mathbf{x}$$

$$\mathbf{c} = \mathbf{D}^\top \mathbf{x}.$$

- Part 2: Choosing the encoding matrix  $\mathbf{D}$

$$\mathbf{D}^* = \arg \min_{\mathbf{D}} \sqrt{\sum_{i,j} \left( x_j^{(i)} - r(\mathbf{x}^{(i)})_j \right)^2} \text{ subject to } \mathbf{D}^\top \mathbf{D} = \mathbf{I}_l$$

Assumption:  $\ell=1$   $\mathbf{D} \rightarrow \mathbf{d}$



$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \sum_i \|\mathbf{x}^{(i)} - \mathbf{d}\mathbf{d}^\top \mathbf{x}^{(i)}\|_2^2 \text{ subject to } \|\mathbf{d}\|_2 = 1$$

$$= \arg \min_{\mathbf{d}} \sum_i \|\mathbf{x}^{(i)} - \mathbf{d}^\top \mathbf{x}^{(i)} \mathbf{d}\|_2^2 \text{ subject to } \|\mathbf{d}\|_2 = 1$$

$$= \arg \min_{\mathbf{d}} \sum_i \|\mathbf{x}^{(i)} - \mathbf{x}^{(i)\top} \mathbf{d}\mathbf{d}\|_2^2 \text{ subject to } \|\mathbf{d}\|_2 = 1$$



# Principal Components Analysis

Defining the augmented matrix  $\mathbf{X}_{i,:} = \mathbf{x}^{(i)\top}$

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \sum_i \|\mathbf{x}^{(i)} - \mathbf{x}^{(i)\top} \mathbf{d} \mathbf{d}^\top\|_2^2 \text{ subject to } \|\mathbf{d}\|_2 = 1$$

From

$$= \arg \min_{\mathbf{d}} \|\mathbf{X} - \mathbf{X} \mathbf{d} \mathbf{d}^\top\|_F^2 \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1$$

$$\arg \min_{\mathbf{d}} \|\mathbf{X} - \mathbf{X} \mathbf{d} \mathbf{d}^\top\|_F^2 = \arg \min_{\mathbf{d}} \text{Tr} \left( \left( \mathbf{X} - \mathbf{X} \mathbf{d} \mathbf{d}^\top \right)^\top \left( \mathbf{X} - \mathbf{X} \mathbf{d} \mathbf{d}^\top \right) \right)$$

$$= \arg \min_{\mathbf{d}} \text{Tr}(\mathbf{X}^\top \mathbf{X} - \mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top - \mathbf{d} \mathbf{d}^\top \mathbf{X}^\top \mathbf{X} + \mathbf{d} \mathbf{d}^\top \mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top)$$

$$= \arg \min_{\mathbf{d}} \text{Tr}(\mathbf{X}^\top \mathbf{X}) - \text{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) - \text{Tr}(\mathbf{d} \mathbf{d}^\top \mathbf{X}^\top \mathbf{X}) + \text{Tr}(\mathbf{d} \mathbf{d}^\top \mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top)$$

$$= \arg \min_{\mathbf{d}} - \text{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) - \text{Tr}(\mathbf{d} \mathbf{d}^\top \mathbf{X}^\top \mathbf{X}) + \text{Tr}(\mathbf{d} \mathbf{d}^\top \mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top)$$

$$= \arg \min_{\mathbf{d}} -2 \text{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) + \text{Tr}(\mathbf{d} \mathbf{d}^\top \mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top)$$

$$= \arg \min_{\mathbf{d}} -2 \text{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) + \text{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top \mathbf{d} \mathbf{d}^\top)$$



# Principal Components Analysis

Thus, we have

$$\begin{aligned} & \arg \min_d -2 \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) + \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top \mathbf{d} \mathbf{d}^\top) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1 \\ &= \arg \min_d -2 \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) + \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1 \\ &= \arg \min_d -\operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1 \\ &= \arg \max_d \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1 \\ &= \arg \max_d \operatorname{Tr}(\mathbf{d}^\top \mathbf{X}^\top \mathbf{X} \mathbf{d}) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1. \\ & \quad \quad \quad \uparrow \\ & \quad \quad \quad ||\mathbf{X}\mathbf{d}||^2 \end{aligned}$$





# Principal Components Analysis

Therefore, the optimal

$$\begin{aligned} X &= USV^T \\ \mathbf{d}^* &= \arg \max ||X\mathbf{d}||^2 = \arg \max ||USV^T \mathbf{d}||^2 \\ &= \arg \max (USV^T \mathbf{d})^T (USV^T \mathbf{d}) \\ &= \arg \max \mathbf{d}^T VS^2V^T \mathbf{d} \end{aligned}$$

→ the optimal  $\mathbf{d}$  is the right singular vector corresponding to the largest singular value of  $X$

