

Introduction to ML

Anastasia Makarova and Mohammad Reza Karimi

ETH Zürich

22.04.2020

Tutorial plan

- Questions for Homework 3 with Anastasia
- Dimensionality Reduction via PCA with Mohammad

Note 1: Please open EduApp on your mobile phone or web browser (link in the chat)

Note 2: Please bring paper and pen for the second part

HW3 Q4: Valid kernel. Problem

Problem 2 (Kernels):

Use the basic rules for kernel decomposition discussed in class or otherwise and assuming that $k(x, y)$ kernel, letting $f : \mathbb{R} \rightarrow \mathbb{R}$ in a) and b), $g : \mathcal{X} \rightarrow \mathbb{R}_+$ for d), $f : \mathcal{X} \rightarrow \mathbb{R}$ for e) and f), and $\phi : \mathcal{X} \rightarrow \mathcal{X}'$.

4. Mark the following statements as True or False. Try to justify your answers to yourself.

- (a) $k_a(x, y) = f(k(x, y))$ is a valid kernel, if f is a polynomial with non-negative coefficients.
- (b) $k_b(x, y) = f(k(x, y))$ is a valid kernel, if f is any polynomial.
- (c) $k_c(x, y) = \exp(k(x, y))$ is a valid kernel.
- (d) $k_d(x, y) = g(x)k(x, y)g(y)$ is a valid kernel.
- (e) $k_e(x, y) = f(x)k(x, y)f(y)$ is a valid kernel.
- (f) $k_f(x, y) = k(\phi(x), \phi(y))$ is a valid kernel.

HW3 Q4: Valid kernel. Recap

Definition 1: Kernel

A kernel k is a two-argument real-valued function over $\mathcal{X} \times \mathcal{X}$, such that for any $\mathbf{x}, \mathbf{y} \in \mathcal{X}$:

$$k(x, y) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle_{\mathcal{F}}$$

for some inner-product space \mathcal{F} such that $\forall \mathbf{x} \in \mathcal{X} \quad \phi(\mathbf{x}) \in \mathcal{F}$.

How to show that k is a *valid* kernel?

- It is sufficient to present a mapping ϕ ,
- Or show that \mathbf{K} is PSD for finite subset of \mathcal{X}
 \mathbf{K} denotes kernel matrix of corresponding kernel k .

HW3 Q4: Valid kernel. Recap

Definition 2: Kernel

A kernel k is a two-argument real-valued function over $\mathcal{X} \times \mathcal{X}$, such that for any $\mathbf{x}, \mathbf{y} \in \mathcal{X}$:

$$k(x, y) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle_{\mathcal{F}}$$

for some inner-product space \mathcal{F} such that $\forall \mathbf{x} \in \mathcal{X} \quad \phi(\mathbf{x}) \in \mathcal{F}$.

How to show that k is a *valid* kernel?

- It is sufficient to present a mapping ϕ ,
- Or show that \mathbf{K} is PSD for finite subset of \mathcal{X}
 \mathbf{K} denotes kernel matrix of corresponding kernel k .

HW3 Q4: Valid kernel. Recap

Definition 3: Kernel

A kernel k is a two-argument real-valued function over $\mathcal{X} \times \mathcal{X}$, such that for any $\mathbf{x}, \mathbf{y} \in \mathcal{X}$:

$$k(x, y) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle_{\mathcal{F}}$$

for some inner-product space \mathcal{F} such that $\forall \mathbf{x} \in \mathcal{X} \quad \phi(\mathbf{x}) \in \mathcal{F}$.

How to show that k is a *valid* kernel?

- It is sufficient to present a mapping ϕ ,
- Or show that \mathbf{K} is PSD for finite subset of \mathcal{X}

\mathbf{K} denotes kernel matrix of corresponding kernel k . $x_i \in \mathcal{X}$

$$\mathbf{K} = \begin{bmatrix} k(x_1, x_1) & & \\ & \dots & \\ & & k(x_n, x_n) \end{bmatrix}$$

\mathbf{K} is PSD \Leftrightarrow eigenvalues are non negative

HW3 Q4: Valid kernel. Basics

① $k(x, y) = k_1(x, y) + k_2(x, y)$

$K = K_1 + K_2 \rightarrow$

② $k(x, y) = \alpha k_1(x, y), \alpha \geq 0 \rightarrow$

③ $k(x, y) = k_1(x, y)k_2(x, y)$

④ $k_4(x, y) = f(x)f(y), f : \mathcal{X} \rightarrow R$

By defining $\phi : x \rightarrow f(x)$ we show

$\exists \phi : \mathcal{X} \rightarrow \mathcal{F} : k(x, y) = \langle \phi(x), \phi(y) \rangle$

HW3 Q4: Valid kernel. Basics

① $k(x, y) = k_1(x, y) + k_2(x, y)$

$$K = K_1 + K_2 \rightarrow$$

$$\beta \in \mathbb{R}^n, K \in \mathbb{R}^{n \times n}$$

$\beta K \beta^T \stackrel{?}{\geq} 0$ why? \rightarrow
 $\beta K \beta^T = \beta K_1 \beta^T + \beta K_2 \beta^T \geq 0$

② $k(x, y) = \alpha k_1(x, y), \alpha \geq 0 \rightarrow$

③ $k(x, y) = k_1(x, y)k_2(x, y)$

④ $k_4(x, y) = f(x)f(y), f: \mathcal{X} \rightarrow \mathbb{R}$

By defining $\phi: x \rightarrow f(x)$ we show

$$\exists \phi: \mathcal{X} \rightarrow \mathcal{F}: k(x, y) = \langle \phi(x), \phi(y) \rangle$$

HW3 Q4: Valid kernel. Basics

① $k(x, y) = k_1(x, y) + k_2(x, y)$
 $K = K_1 + K_2 \rightarrow$

② $k(x, y) = \alpha k_1(x, y), \alpha \geq 0 \rightarrow$

$\alpha \beta k_1 \beta^T \geq 0$

③ $k(x, y) = k_1(x, y)k_2(x, y)$

④ $k_4(x, y) = f(x)f(y), f : \mathcal{X} \rightarrow R$

By defining $\phi : x \rightarrow f(x)$ we show

$\exists \phi : \mathcal{X} \rightarrow \mathcal{F} : k(x, y) = \langle \phi(x), \phi(y) \rangle$

HW3 Q4: Valid kernel. Basics

① $k(x, y) = k_1(x, y) + k_2(x, y)$
 $K = K_1 + K_2 \rightarrow$

② $k(x, y) = \alpha k_1(x, y), \alpha \geq 0 \rightarrow$

③ $k(x, y) = k_1(x, y)k_2(x, y)$

eigenvalue decompositions
 $K_1 = \sum_{i=1}^n \lambda_i v_i v_i^\top; K_2 = \sum_{i=1}^n \mu_i u_i u_i^\top$

$K = K_1 \circ K_2 = \sum_{i,j} \lambda_i \mu_j (v_i \otimes u_j)(v_i \otimes u_j)^\top \succeq 0$ rank-1 matrices

④ $k_4(x, y) = f(x)f(y), f: \mathcal{X} \rightarrow \mathbb{R}$
By defining $\phi: x \rightarrow f(x)$ we show
 $\exists \phi: \mathcal{X} \rightarrow \mathcal{F}: k(x, y) = \langle \phi(x), \phi(y) \rangle$

HW3 Q4: Valid kernel. Basics

① $k(x, y) = k_1(x, y) + k_2(x, y)$
 $K = K_1 + K_2 \rightarrow$

② $k(x, y) = \alpha k_1(x, y), \alpha \geq 0 \rightarrow$

③ $k(x, y) = k_1(x, y)k_2(x, y)$

④ $k_4(x, y) = f(x)f(y), f : \mathcal{X} \rightarrow R$

By defining $\phi : x \rightarrow f(x)$ we show

$\exists \phi : \mathcal{X} \rightarrow \mathcal{F} : k(x, y) = \langle \phi(x), \phi(y) \rangle$

HW3 Q4: Valid kernel. Basics

① $k(x, y) = k_1(x, y) + k_2(x, y)$
 $K = K_1 + K_2 \rightarrow$

② $k(x, y) = \alpha k_1(x, y), \alpha \geq 0 \rightarrow$

③ $k(x, y) = k_1(x, y)k_2(x, y)$

④ $k_4(x, y) = f(x)f(y), f : \mathcal{X} \rightarrow R$
By defining $\phi : x \rightarrow f(x)$ we show
 $\exists \phi : \mathcal{X} \rightarrow \mathcal{F} : k(x, y) = \langle \phi(x), \phi(y) \rangle$

HW3 Q4: Valid kernel. Solution

- a $k_a(x, y) = f(k(x, y))$ polynomial with non-negative coefficients

follows from (1), (2) and (3)

- b $k_b(x, y) = f(k(x, y))$ polynomial. Counter example

- c $k_c(x, y) = \exp(k(x, y))$

HW3 Q4: Valid kernel. Solution

a $k_a(x, y) = f(k(x, y))$ polynomial with non-negative coefficients

b $k_b(x, y) = f(k(x, y))$ polynomial. Counter example

$$\begin{aligned} f(x) &= -x \\ k(x, y) &= x^T y \end{aligned} \rightarrow k_b(x, y) = -x^T y \text{ that will} \\ \text{have negative} \\ \text{eigenvalues}$$

c $k_c(x, y) = \exp(k(x, y))$

negation of a valid kernel

HW3 Q4: Valid kernel. Solution

- a $k_a(x, y) = f(k(x, y))$ polynomial with non-negative coefficients

- b $k_b(x, y) = f(k(x, y))$ polynomial. Counter example

- c $k_c(x, y) = \exp(k(x, y))$

$$e^x = 1 + x + \frac{x^2}{2!} + \dots$$

HW3 Q4: Valid kernel. Solution

$g: \mathcal{X} \rightarrow \mathbb{R}_+$ \rightarrow 1. let's prove for a general case
 $f: \mathcal{X} \rightarrow \mathbb{R}$ of f

2. $f(x)f(y)$ is a kernel (4)
3. product of kernel (3)
- d $k_d(x, y) = g(x)k(x, y)g(y)$
- e $k_e(x, y) = f(x)k(x, y)f(y) = k_4(x, y)k(x, y)$

- f $k_f = k(\phi(x), \phi(y))$ We want to show that K_f is PSD.
As $k(\cdot)$ is a kernel, when applied to any set of vectors $\{\phi(x_i)\}_{i=1}^N$ will result if PSD matrix.

HW3 Q4: Valid kernel. Solution

- d $k_d(x, y) = g(x)k(x, y)g(y)$
- e $k_e(x, y) = f(x)k(x, y)f(y) = k_4(x, y)k(x, y)$
- f $k_f = k(\phi(x), \phi(y))$ We want to show that \mathbf{K}_f is PSD.
As $k(\cdot)$ is a kernel, when applied to any set of vectors $\{\phi(x_i)\}_{i=1}^N$ will result if PSD matrix.

HW3 Q4: Valid kernel. Solution

d $k_d(x, y) = g(x)k(x, y)g(y)$

e $k_e(x, y) = f(x)k(x, y)f(y) = k_4(x, y)k(x, y)$

f $k_f = k(\phi(x), \phi(y))$ We want to show that \mathbf{K}_f is PSD.

As $k(\cdot)$ is a kernel, when applied to any set of vectors $\{\phi(x_i)\}_{i=1}^N$ will result if PSD matrix.

HW3 Q4: Valid kernel. Solution

- d $k_d(x, y) = g(x)k(x, y)g(y)$
- e $k_e(x, y) = f(x)k(x, y)f(y) = k_4(x, y)k(x, y)$
- f $k_f = k(\phi(x), \phi(y))$ We want to show that \mathbf{K}_f is PSD.
As $k(\cdot)$ is a kernel, when applied to any set of vectors $\{\phi(x_i)\}_{i=1}^N$ will result if PSD matrix.

HW3 Q4: Valid kernel. Eduapp question

Please open eduapp

HW3 Q4: Valid kernel. Eduapp question

Which of the following functions are kernels?

X $k_1(x, y) = -\frac{\langle x, y \rangle}{\|x\|_2 \|y\|_2}$ negation of a valid kernel

X $k_2(x, y) = \sqrt{(\|x - y\|_2^2 + 1)}$ counter example unit vectors $x = (1, 0)$ $y = (0, 1)$
 $k_2 = \begin{pmatrix} 1 & \sqrt{3} \\ \sqrt{3} & 1 \end{pmatrix} \rightarrow$
 $\det \begin{pmatrix} 1-\lambda & \sqrt{3} \\ \sqrt{3} & 1-\lambda \end{pmatrix} = 0 \rightarrow (1-\lambda)^2 - 3 = 0$
 $(1-\lambda-\sqrt{3})(1-\lambda+\sqrt{3}) = 0 \rightarrow k_2$ has negative eigenvalue

Valid $k_3(x, y) = \prod_{i=1}^D \underbrace{\exp\left(\frac{x_i^2 - c}{2}\right)}_{f(x_i)} \underbrace{\exp\left(\frac{y_i^2 - c}{2}\right)}_{f(y_i)}$

can be rewritten as $k_3(x, y) = f(x)f(y)$
 \rightarrow according to (4) is valid

HW4 Q3: SVM update rule via SGD

The exercise suggests to train an SVM, where we penalise the margin violation given by $(1 - y\mathbf{w}^T \mathbf{x})_+ = \max(1 - y\mathbf{w}^T \mathbf{x}, 0)$, not linearly but with the square root instead. Correspondingly, our modified SVM seeks to optimise the following objective

$$L(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \sqrt{(1 - y\mathbf{w}^T \mathbf{x})_+} + \lambda \|\mathbf{w}\|_2$$

Dimensionality Reduction via PCA

1) Review of Some Linear algebra

Spectral Decomp $X \in \mathbb{R}^{n \times n}$ symmetric $X = V \Lambda V^T$
 V is orthogonal $= \sum \lambda_i v_i v_i^T$
(eigenvectors)

if all $\lambda_i \geq 0 \rightarrow$ PSD Matrix \rightarrow Cholskey Dcomp.

$$X = V \Lambda^{1/2} \Lambda^{1/2} V^T = (V \Lambda^{1/2}) (V \Lambda^{1/2})^T = Q Q^T.$$

Trace For square matrices $X \in \mathbb{R}^{n \times n}$,
 $\text{tr}(X) = \sum_{i=1}^n X_{ii} =$ sum of diagonal entries.

* Properties of Trace (Exercise):

$$- \text{tr}(AB) = \text{tr}(BA) = \sum_{i,j=1}^n A_{ij} B_{ij}.$$

$$- \text{tr}(ABC) = \text{tr}(BCA) = \text{tr}(CAB).$$

$$- \text{tr}(P A P^{-1}) = \text{tr}(A).$$

\Rightarrow trace of a "linear transformation" is well defined as it does not depend on the matrix repr. of the linear operator.

Frobenius norm of a matrix A is

$$\|A\|_F = \sqrt{\sum_{i,j} A_{ij}^2} = \sqrt{\text{tr}(A A^T)}.$$

Also, if A 's columns are a_1, \dots, a_n ,

$$\|A\|_F^2 = \|a_1\|_2^2 + \dots + \|a_n\|_2^2$$

An Optimization Problem

We are interested in
($X \in \mathbb{R}^{n \times n}$ symmetric).

$$\begin{aligned} \max_{A \in \mathbb{R}^{n \times d}} \operatorname{tr}(A^T X A) \\ A^T A = I \end{aligned}$$

We claim that the optimal A is indeed the set of eigenvectors of X corresponding to the d largest eigenvalues.

Case $d=1$: $\max_{\substack{a \in \mathbb{R}^n \\ \|a\|=1}} a^T X a$

$$\begin{aligned} X &= V \Lambda V^T \\ \Rightarrow a^T X a &= a^T V \Lambda V^T a = (V^T a)^T \Lambda (V^T a) \\ &= u^T \Lambda u \end{aligned}$$

and V ortho $\Rightarrow \|u\| = \|a\| + V$ one-to-one

$$\Rightarrow \text{OPT} = \max_{\|u\|=1} u^T \Lambda u$$

But $u^T \Lambda u = \sum_{i=1}^n u_i^2 \lambda_i$, $\|u\|^2 = \sum_{i=1}^n u_i^2 = 1$

max is when

$$u_1 = 1, u_i = 0 \quad (i > 1). \quad \blacksquare$$

2) Variance point-of-view to PCA

Goal: find orthonormal basis v_1, \dots, v_d such that the projections of x_1, \dots, x_n on span of $\{v_1, \dots, v_d\}$ has the **most variance**.

(Variance above should be read "sample variance"...)

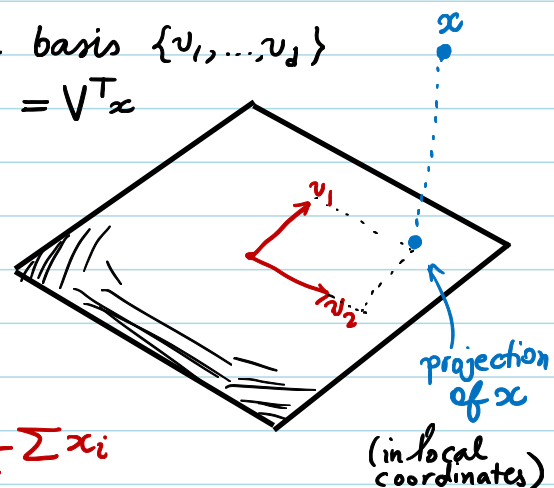
The projection of x_i , expressed in the basis $\{v_1, \dots, v_d\}$ has coordinates $(v_1^T x_i, \dots, v_d^T x_i) = V^T x_i$

Hence, the variance of the projection is

$$\sum_{i=1}^n \|V^T x_i - V^T \bar{x}\|^2$$

$$V = [v_1 | \dots | v_d]$$

$$\text{mean} = \frac{1}{n} \sum x_i$$



The goal of PCA is then to find

$$\max_{V^T V = I} \sum_{i=1}^n \|V^T x_i - V^T \bar{x}\|^2.$$

Now some algebraic tricks:

$$\|V^T x_i - V^T \bar{x}\|^2 = \|V^T (x_i - \bar{x})\|^2$$

$$= (x_i - \bar{x})^T V V^T (x_i - \bar{x})$$

$$= \text{tr} \left\{ (x_i - \bar{x})^T V V^T (x_i - \bar{x}) \right\}$$

$$= \text{tr} \left\{ V^T (x_i - \bar{x}) (x_i - \bar{x})^T V \right\}$$

$$\Rightarrow \sum_{i=1}^n \|V^T x_i - V^T \bar{x}\|^2$$

$$= \sum \text{tr} \left\{ V^T (x_i - \bar{x}) (x_i - \bar{x})^T V \right\}$$

$$= \text{tr} \left\{ \sum V^T (x_i - \bar{x}) (x_i - \bar{x})^T V \right\}$$

$$= \text{tr} \left\{ V^T \left[\sum (x_i - \bar{x}) (x_i - \bar{x})^T \right] V \right\}$$

$$= (n-1) \text{tr} V^T \Sigma_n V$$

↖ the covariance matrix!

... and we know the answer!

3) PCA in high dimensions

Assume that the data points x_1, \dots, x_n are drawn i.i.d from some distribution.

What PCA computes $\xrightarrow[\text{on}]{\text{depends}}$ Spectrum of Σ_n ,
the sample covariance matrix

What we "Hopefully" want $\xrightarrow[\text{on}]{\text{depends}}$ Spectrum of Σ ,
the covariance matrix
of the distribution

if the data is "really" having a
low-dimensional structure,
does PCA see it?

Can PCA trick us to "see"
something which is not true?

MAIN
QUESTION

\rightarrow requires that spectra of Σ_n is "close" to $\Sigma \dots$

ANSWER: maybe! yes! so, be careful!

BABY example: Gaussian data

Suppose that $x_1, \dots, x_n \sim \mathcal{N}(0, \Sigma)$ multivariate
in p dimensions

$$\text{Construct } \Sigma_n = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$

law of large numbers says for fixed p , $\Sigma_n \rightarrow \Sigma$
as $n \rightarrow \infty \dots$

In practice, p often is of order of $n \dots$ Cannot use
LLN!

Assume that $\Sigma = I$.

[Show the notebook ...]

Lots of eigenvalues > 1 !

PCA with small d gets fooled ... but the original distribution is rotation-invariant (isotropic)!

BABY example 2 : Spike model

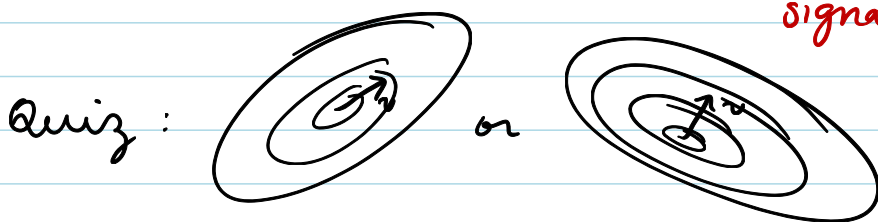
The data distribution really has low-dimensional structure. Does PCA find it?

$$\Sigma = I + \beta v v^T \quad \beta > 0, \|v\| = 1.$$

if $x \sim N(0, \Sigma)$, $x = \underbrace{\sqrt{\beta} g_0 v}_{\text{signal}} + \underbrace{\epsilon}_{\text{noise}}$

$\epsilon \sim N(0, I), g_0 \sim N(0, 1)$

β : How powerful is the signal (SNR).



Reformulated Question: How large β should be so that PCA "sees" the direction v ?

Intuition: something should be visible in spectra of Σ_n ...

[show notebook ...]

BBP Transition Suppose $\frac{p}{n} = \gamma$ and let $n, p \rightarrow \infty$.

then if $\beta > \sqrt{\gamma}$, $\lambda_{\max}(\Sigma_n)$ gets out of MP curve
if $\beta \leq \sqrt{\gamma}$, $\lambda_{\max}(\Sigma_n)$ stays inside ...

Beautiful theory, but little time ...

kernel k-means & Spectral Clustering

Setup a set of points x_1, \dots, x_n
 k-means wish to find clusters with means $\{\mu_1, \dots, \mu_k\} = \mu$

st. $C(z, \mu) = \sum_{i=1}^n \|x_i - \mu_{z_i}\|^2$ is minimized.

the cluster assigned to point x_i

Now, assume that before clustering, we first apply a feature map ϕ to our data, and then perform clustering.

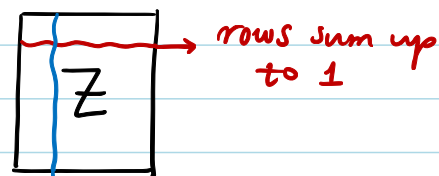
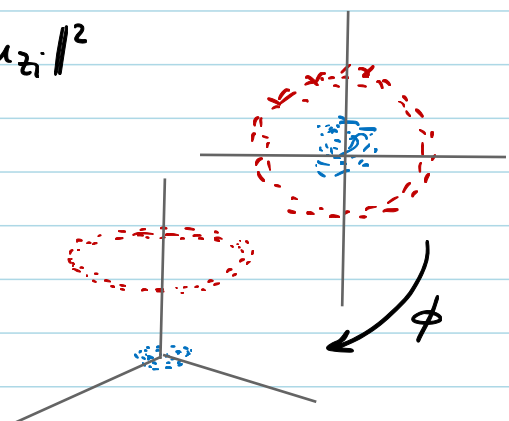
New obj: $C(z, \mu) = \sum_{i=1}^n \|\phi(x_i) - \mu_{z_i}\|^2$

To write the objective nicer, we define a "cluster assignment matrix" $Z_{n \times K}$ as

$$z_{ik} = \begin{cases} 1 & \text{if } z_i = k \\ 0 & \text{otherwise} \end{cases}$$

Define $L = \text{diag}(\frac{1}{N_1}, \dots, \frac{1}{N_K})$.

and $\Phi = [\phi(x_1) | \dots | \phi(x_n)]$.



Now check that $M = \Phi Z L Z^T$ has n columns, each column is the cluster center associated to that point via z .

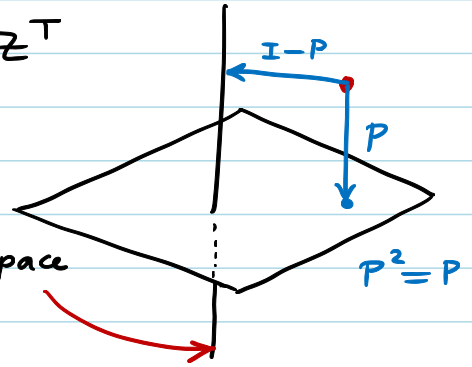
copy of

\Rightarrow K-means cost $C = \|\Phi - M\|_F^2$

$= \text{tr} \{ (\Phi - M)(\Phi - M)^T \}$.

- Some facts:
- $Z^T Z = L^{-1}$
 - $(Z L Z^T)^2 = Z L Z^T$

$\Rightarrow Z L Z^T$ is a projection
 $\Rightarrow I - Z L Z^T$ is projection on the complement space



$$\Phi - M = \Phi(I - Z L Z^T)$$

Now we simplify the objective:

$$\begin{aligned} C &= \text{tr} \{ (\Phi - M)(\Phi - M)^T \} \\ &= \text{tr} \{ \Phi(I - Z L Z^T)^2 \Phi^T \} \\ &= \text{tr} \{ \Phi(I - Z L Z^T) \Phi^T \} \\ &= \text{tr} \{ \Phi \Phi^T \} - \text{tr} \{ \Phi Z L Z^T \Phi^T \} \\ &= \text{tr} \{ K \} - \text{tr} \{ L^{1/2} Z^T K Z L^{1/2} \} \quad (K = \Phi^T \Phi) \end{aligned}$$

does not depend on clustering Z

\Rightarrow new objective: $\max \text{tr} \{ L^{1/2} Z^T K Z L^{1/2} \}$
 s.t. Z is a binary clustering matrix.

"Binary" makes it hard \rightarrow relax!

$$Z^T Z = L^{-1} \Rightarrow L^{1/2} Z^T Z L^{1/2} = I \Rightarrow H^T H = I$$

relaxed problem $\max \text{tr} \{ H^T K H \}$
 s.t. $H^T H = I$

... our friend in PCA!

Hence, optimal H is obtained via taking top k eigenvectors of K .

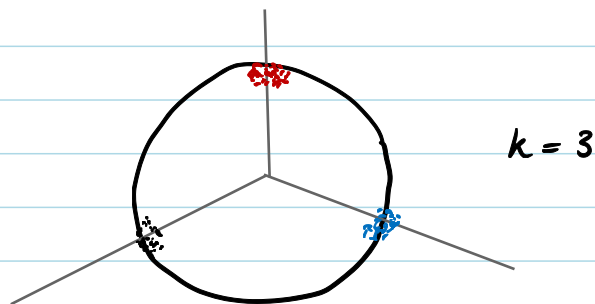
Last part: How to find the clusters? (Z)

proposal: H approximates Z ...
each row of Z summed up to 1 ...
 \Rightarrow for each row of H , set max to 1
and the rest to 0 (rounding)

Better: Normalize each row of H by its norm.
 \rightarrow rows of H are now points on the unit sphere.

Run k -means and extract clusters!

dream:



Called "Spectral Clustering".

There are other approaches as well, but we don't discuss them here.