

Goals:

- Review of some linear algebra
- Deriving Principal Component Analysis.

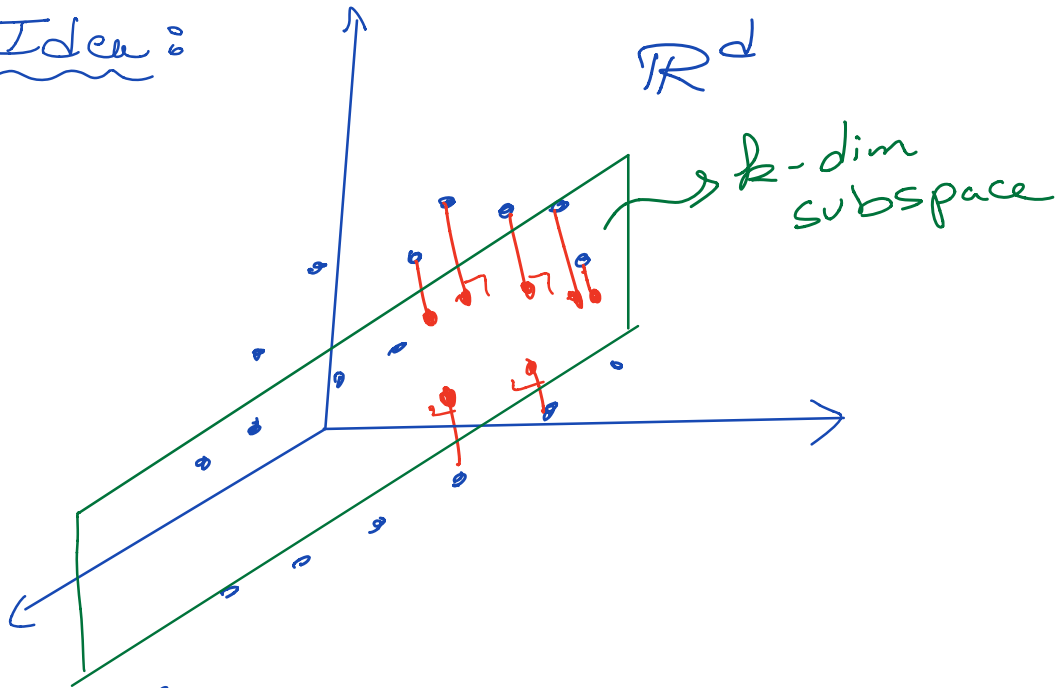
↙

Idea: Given some data set $\{x_1, x_2, \dots, x_p\} \subseteq \mathbb{R}^d$, want to find a linear projection to $k < p$ dimensions while preserving key properties of the data set.

Why?

- Visualization ($k = 2$ or 3)
- "Interpretation"

Idea:



Let $X = \{x_1, \dots, x_p\}$
and define the sample mean

$$\mu_p = \frac{1}{p} \sum_{i=1}^p x_i$$

and sample covariance

$$\Sigma_p = \frac{1}{p-1} \sum_{i=1}^p (x_i - \mu_p)(x_i - \mu_p)^T$$

Homework: Prove that μ_p & Σ_p are unbiased estimators for the mean & variance of the distribution from which each x_i is drawn independently.

Now, assuming v_1, v_2, \dots, v_k form an orthonormal basis for a k -dim. subspace, we want to minimize

$$\min_{\mu, \beta} \sum_{i=1}^P \|x_i - (\mu + V\beta_i)\|_2^2$$

$V^T V = I$

$\in \mathbb{R}^{d \times k}$ $\in \mathbb{R}^k$, vector of coefficients
 \downarrow \downarrow
 $\left[v_1 | v_2 | \dots | v_k \right]$
 \downarrow
 $\in \mathbb{R}^d$

offline shift
 $\in \mathbb{R}^d$

(I) optimal μ :

$$\nabla_{\mu}(\) = 0 \Rightarrow$$

$$\nabla_{\mu} \sum_{i=1}^P (x_i - \mu_{\text{opt}} - V\beta_i)^T (x_i - \mu_{\text{opt}} - V\beta_i) = 0$$

so
$$\sum_{i=1}^P (x_i - \mu_{\text{opt}} - V\beta_i) = 0$$

&
$$\mu_{\text{opt}} = \frac{1}{n} \sum_{i=1}^P x_i + \frac{1}{n} V \underbrace{\sum_{i=1}^P \beta_i}_{=0}$$

assume, w2OG that $\sum \beta_i = 0$

\Rightarrow $\mu_{\text{opt}} = \frac{1}{n} \sum_{i=1}^P x_i$ ← sample mean

(II) Optimal β_i , $i=1, \dots, P$

$$\min_{\beta_i} \sum_{i=1}^P \|x_i - \mu_{\text{opt}} - V\beta_i\|_2^2$$

We can solve for each β_i separately
(Obj. Function is separable)

$$\Rightarrow \text{want to minimize } \|x_i - \mu_{\text{opt}} - V\beta_i\|_2^2$$

$$= (x_i - \mu_{\text{opt}} - V\beta_i)^T (x_i - \mu_{\text{opt}} - V\beta_i)$$

$$\nabla_{\beta_i} = 0 \Rightarrow V^T (x_i - \mu_{\text{opt}} - V\beta_i) = 0 \quad (\text{chain rule})$$

$$\Rightarrow \beta_i = V^T (x_i - \mu_{\text{opt}})$$

(III) Optimal V :

$$\text{want } \min_{V, V^T V = I} \sum_{i=1}^P \|x_i - \mu_{\text{opt}} - V\beta_i^{\text{opt}}\|_2^2$$

$$\Leftrightarrow \min_{V, V^T V = I} \sum_{i=1}^P \|x_i - \mu_{\text{opt}} - V V^T (x_i - \mu_{\text{opt}})\|_2^2$$

$$\| (I - V V^T) (x_i - \mu_{\text{opt}}) \|_2^2$$

$$= \sum_{i=1}^P \underbrace{\|x_i - \mu_{\text{opt}}\|_2^2}_{\text{ind. of } V} - (x_i - \mu_{\text{opt}})^T V V^T (x_i - \mu_{\text{opt}})$$

so we want to maximize (over $V, V^T V = I$)

$$\sum_{i=1}^P (x_i - \mu_{opt})^T V V^T (x_i - \mu_{opt})$$

$$= \sum_{i=1}^P \text{Tr} \left((x_i - \mu_{opt})^T V V^T (x_i - \mu_{opt}) \right)$$

↳ trace: $\text{trace}(X) = \sum_j X_{jj}$

$$= \sum_{i=1}^P \text{Tr} \left(V^T (x_i - \mu_{opt}) (x_i - \mu_{opt})^T V \right)$$

↳ cyclicity of trace: $\text{Tr}(A X B) = \text{Tr}(B A X)$

$$= \text{Tr} \sum_{i=1}^P \left(V^T (x_i - \mu_{opt}) (x_i - \mu_{opt})^T V \right)$$

↳ linearity of trace: $(\text{Tr}(A+B) = \text{Tr}(A) + \text{Tr}(B))$

$$= \text{Tr} V^T \underbrace{\sum_{i=1}^P (x_i - \mu_{opt}) (x_i - \mu_{opt})^T}_{} V$$

$= (P-1) \Sigma_P$

$$= (P-1) \text{Tr} (V^T \Sigma_P V)$$

So we want to maximize

$$\max \text{Tr}(V^T \Sigma_P V)$$

$V^T V = I$
 $V \in \mathbb{R}^{d \times k}$

$\Rightarrow V =$ matrix of k leading eigenvectors
of Σ_P

\rightarrow Why? Linear algebra

Alternative Interpretation: PCA finds

the directions that "preserve the most variance."

\Leftrightarrow want $V = [v_1 | v_2 | \dots | v_k]$, $V^T V = I$

such that $\left\{ \begin{pmatrix} v_1^T x_i \\ \vdots \\ v_k^T x_i \end{pmatrix} \right\}_{i=1}^P$ has maximal variance

$\left\{ \begin{pmatrix} v_1^T x_i \\ \vdots \\ v_k^T x_i \end{pmatrix} \right\}_{i=1}^P \in \mathbb{R}^k$

$$\begin{aligned}
 \Leftrightarrow \text{want } \max_{\substack{V: V^T V = I \\ V \in \mathbb{R}^{d \times k}}} & \underbrace{\sum_{i=1}^P \left\| V^T x_i - \frac{1}{P} \sum_{j=1}^P V^T x_j \right\|_2^2}_{=} \\
 & = \sum_{i=1}^P \left\| V^T (x_i - \mu_P) \right\|_2^2 \\
 & = \text{Tr}(V^T \Sigma_P V)
 \end{aligned}$$

as before

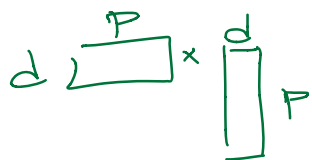
~ x ~

Computational Complexity:

Cost of Computing Σ_P via SVD $\stackrel{\text{PCA via}}{=} \uparrow$ leading \uparrow eigenvectors

cost of Computing Σ_P + cost of SVD

$$= O(d^2 + kdP)$$



↳ can be reduced via randomized algorithms

Choice of k is for dimensionality

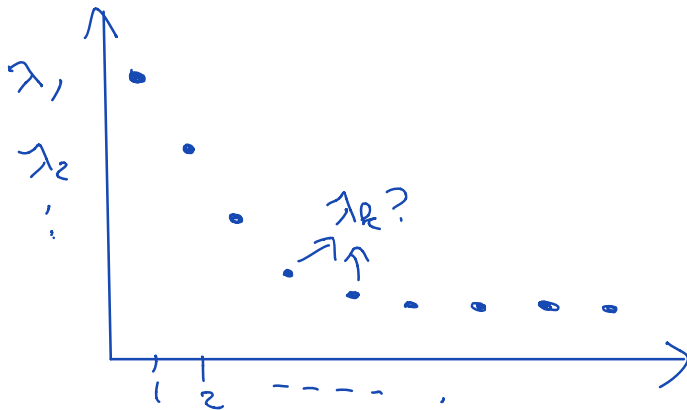
reduction is usually done heuristically.

• Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k \geq \dots \geq 0$

be the eigenvalues of Σ_P ,
select k so that $\sum_{i=1}^k \lambda_i / \text{Tr}(\Sigma_P)$
is "large enough".

smallest $= \sum_{i=1}^d \lambda_i$

• Or pick k at the "elbow"
of a scree plot.



Extensions to Classical PCA:

- Nonlinear dimensionality reduction
e.g.: • kernel PCA (may discuss towards end of course)
 - principal manifold embeddings

- Non-negative matrix factorization
want

$$\min_{L, R} \|X - LR\|_F$$

$L_{ij} \geq 0$
 $R_{ij} \geq 0$

- Multilinear PCA (X is a tensor)

- Robust PCA ($X = L + S$)
Low-rank sparse

$$\min \|L\|_* + \lambda \|S\|_1$$

s.t. $L + S = M$

nuclear norm l_1 -norm of vectorized S

- many others

Properties / Drawbacks:

PCA is data dependent, so if an encoder performs dimensionality reduction via PCA, i.e., computes the map

$$\mathbb{R}^d \ni X = [x_1 | x_2 | \dots | x_p] \mapsto \begin{matrix} V^T X \\ \downarrow \\ \begin{matrix} k \times d \\ \text{PCA matrix} \end{matrix} \end{matrix}$$

a decoder would need to know V to recover an approximation of any x_i .

But V depends on X , hence x_i .

A solution to this is to use a fixed transformation that takes

↳ data-independent

into account general properties of the signal model, but not the data itself.

⇒ Discrete cosine transform, wavelets, etc....