# PRINCIPAL COMPONENT ANALYSIS

## CS6140

Predrag Radivojac

KHOURY COLLEGE OF COMPUTER SCIENCES
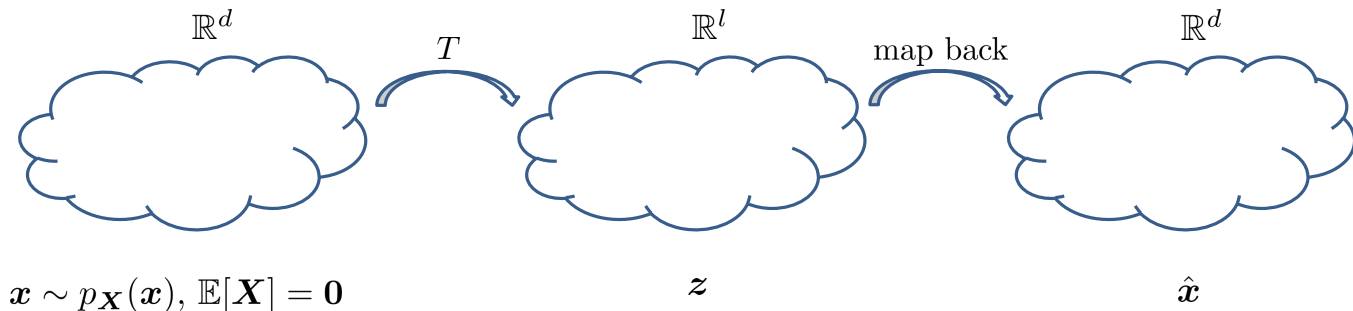
NORTHEASTERN UNIVERSITY

Fall 2024

# PROBLEM FORMULATION

**Given:** a set of vectors $\{\boldsymbol{x}_i\}_{i=1}^n$, where $\boldsymbol{x}_i \in \mathbb{R}^d$, sampled from $p_{\boldsymbol{X}}(\boldsymbol{x})$

**Objective:** find a linear mapping $T : \mathbb{R}^d \to \mathbb{R}^l$, where $l \leq d$, such that the reconstruction of projections back to $\mathbb{R}^d$ is optimal in the mean-squared-error sense.



$\boldsymbol{x} \sim p_{\boldsymbol{X}}(\boldsymbol{x}), \mathbb{E}[\boldsymbol{X}] = \boldsymbol{0}$

$\uparrow$
A minor additional constraint: data is centered.

# LINEAR MAPPING

A function $T : \mathbb{R}^d \longrightarrow \mathbb{R}^l$ is a linear mapping if for $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and $\forall c \in \mathbb{R}$

$$T(\mathbf{x} + \mathbf{y}) = T(\mathbf{x}) + T(\mathbf{y})$$

and

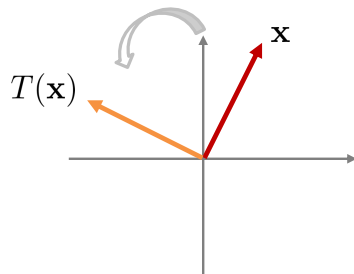$$T(c\mathbf{x}) = cT(\mathbf{x})$$

**Claim:** every linear map $T$ can be represented by an $l \times d$ matrix $\mathbf{T}$ as $T(\mathbf{x}) = \mathbf{T}\mathbf{x}$

**Example:** rotation by $90°$ in 2D space.

$$\mathbf{T} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \qquad \mathbf{x} = (2, 4)$$

$$T(\mathbf{x}) = \mathbf{T}\mathbf{x} = (-4, 2)$$

# PROBLEM FORMULATION

**Matrix view:** $\mathbf{x} \in \mathbb{R}^{d\times 1}$, $\mathbf{T} \in \mathbb{R}^{l\times d}$. The goal is to find $\mathbf{T}$, $\mathbf{z}$.

$$\mathbf{Tx} = \mathbf{z}$$



Minimize: $\mathbb{E}[||\boldsymbol{X} - \hat{\boldsymbol{X}}||^2]$

$$\tilde{\mathbf{T}}\mathbf{z} = \hat{\mathbf{x}}$$

It will turn out later that $\tilde{\mathbf{T}}$ is in fact $\mathbf{T}^T$

# IDEA



Haykin. Neural networks. 1999.

$$||\mathbf{c}||^2 = (||\mathbf{b}|| - ||\mathbf{a}|| \cos \alpha)^2 + (||\mathbf{a}|| \sin \alpha)^2$$
$$= ||\mathbf{b}||^2 - 2||\mathbf{a}|| \cdot ||\mathbf{b}|| \cos \alpha + ||\mathbf{a}||^2 \cos^2 \alpha + ||\mathbf{a}||^2 \sin^2 \alpha$$
$$= ||\mathbf{a}||^2 + ||\mathbf{b}||^2 - 2||\mathbf{a}|| \cdot ||\mathbf{b}|| \cos \alpha$$

$$||\mathbf{c}||^2 = \mathbf{c}^T \mathbf{c}$$
$$= (\mathbf{a} - \mathbf{b})^T (\mathbf{a} - \mathbf{b})$$
$$= \mathbf{a}^T \mathbf{a} - 2\mathbf{a}^T \mathbf{b} + \mathbf{b}^T \mathbf{b}$$
$$= ||\mathbf{a}||^2 - 2\mathbf{a}^T \mathbf{b} + ||\mathbf{b}||^2$$

Combine the two:

$$||\mathbf{a}||^2 - 2\mathbf{a}^T \mathbf{b} + ||\mathbf{b}||^2 = ||\mathbf{a}||^2 + ||\mathbf{b}||^2 - 2||\mathbf{a}|| \cdot ||\mathbf{b}|| \cos \alpha$$



$$\mathbf{a} = \mathbf{b} + \mathbf{c}$$

$$\cos(\alpha) = \frac{\mathbf{a}^T \mathbf{b}}{||\mathbf{a}|| \cdot ||\mathbf{b}||}$$

# PROJECTION TO ONE DIMENSION

Let us project a vector $\mathbf{x}$ to a unit vector $\mathbf{v}$. Note: $\mathbf{v}^T \mathbf{v} = 1$ or $||\mathbf{v}|| = 1$.



$$\cos(\alpha) = \frac{z}{||\mathbf{x}||} = \frac{\mathbf{x}^T \mathbf{v}}{||\mathbf{x}|| \cdot ||\mathbf{v}||} \qquad \Rightarrow \qquad z = \mathbf{x}^T \mathbf{v}$$

Let us project a random vector $\boldsymbol{X} \overset{d \times 1}{\sim} p(\boldsymbol{x})$ to some unit vector $\mathbf{v}$.

$$Z = \boldsymbol{X}^T \mathbf{v} = \mathbf{v}^T \boldsymbol{X}$$

$$\mathbb{E}[Z] = \mathbf{v}^T \mathbb{E}[\boldsymbol{X}] = 0$$

$$\mathbb{E}[Z^2] = \mathbb{E}[\mathbf{v}^T \boldsymbol{X} \boldsymbol{X}^T \mathbf{v}] = \mathbf{v}^T \mathbb{E}[\underset{\underset{d \times d}{\downarrow}}{\boldsymbol{X} \boldsymbol{X}^T}] \mathbf{v} = \mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v} \quad \Rightarrow \quad \mathbb{V}[Z] = \mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v}$$

# PROJECTION TO ONE DIMENSION

For a set of vectors, let us find a unit vector $\mathbf{v}$ so that the projection has maximum variance $\mathbb{V}[Z] = \mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v}$.

**Objective:** Given $\boldsymbol{\Sigma}$, find $\mathbf{v}$ to maximize variance of the projection.

$$\max \mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v} \quad \text{s.t.} \quad \mathbf{v}^T \mathbf{v} = 1$$

$$L(\mathbf{v}, \lambda) = \mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v} + \lambda(1 - \mathbf{v}^T \mathbf{v}) \quad \overset{\text{Solve}}{\Rightarrow} \quad \boldsymbol{\Sigma} \mathbf{v} = \lambda \mathbf{v} \qquad \text{The eigenvalue problem}$$

# PROJECTION TO *d* DIMENSIONS

Consider now projecting to $d$ orthogonal vectors:

$$\mathbf{\Sigma V} = \mathbf{V \Lambda} \qquad\qquad \text{← matrix version}$$

where $\mathbf{V} = [\mathbf{v}_1\ \mathbf{v}_2\ \ldots\ \mathbf{v}_d]$, with $\mathbf{V}^T\mathbf{V} = \mathbf{I}$   ← because $\mathbf{V}$ is orthogonal

$\mathbf{\Lambda} = \text{diag}\{\lambda_1, \lambda_2, \ldots, \lambda_d\}$, with $\lambda_1 \geq \lambda_2 \ldots \geq \lambda_d$

Let us re-write: $\mathbf{V}^T\mathbf{\Sigma V} = \mathbf{\Lambda}$

$$\mathbf{v}_i^T\mathbf{\Sigma v}_j = \begin{cases} \lambda_i & \text{for } i = j \\ 0 & \text{for } i \neq j \end{cases} \qquad \text{← variance of projection } Z_i$$

# TRANSFORMATION

Let us express the $i$-th projection as $z_i = \mathbf{v}_i^T \mathbf{x} = \mathbf{x}^T \mathbf{v}_i$

Thus,

$$\mathbf{z} = (z_1, z_2, \ldots, z_d) = (\mathbf{v}_1^T \mathbf{x}, \mathbf{v}_2^T \mathbf{x}, \ldots, \mathbf{v}_d^T \mathbf{x}) = \mathbf{V}^T \mathbf{x} = \sum_{i=1}^{d} x_i \mathbf{v}_i^T$$

Let us reconstruct $\mathbf{x}$ now. Remember, $\mathbf{V}^{-1} = \mathbf{V}^T$.

$$\mathbf{x} = \mathbf{V}\mathbf{z} = \sum_{i=1}^{d} z_i \mathbf{v}_i$$

# DIMENSIONALITY REDUCTION

Let us now keep the first $l$ components of $\mathbf{z}$.

$$d$$

$$d \quad \mathbf{V}^T \quad \cdot \mathbf{x} = \mathbf{z} \qquad \rightarrow \qquad l \quad \boxed{\mathbf{T}} \quad \cdot \mathbf{x} = \mathbf{z}$$

**Matrix view:**

$$\overset{n \times d}{\underset{\downarrow}{}}$$
$$\mathbf{Z} = \mathbf{X}\mathbf{V} \qquad \rightarrow$$

$\ddot{\mathbf{V}}$ is $\mathbf{V}$ reduced to $l$ columns
$\ddot{\mathbf{V}}_{d \times l}$ reminds us of $\ddot{\mathbf{V}}$'s dimensions

$$\overset{n \times l}{\underset{\downarrow}{}} \quad \overset{}{\underset{\downarrow}{}}$$
$$\mathbf{Z} = \mathbf{X}\ddot{\mathbf{V}}_{d \times l}$$
$$= \mathbf{X}\mathbf{T}^T$$

# RECONSTRUCTION

Let us reconstruct $\mathbf{x}$ now:

$$\mathbf{V} \cdot \mathbf{z} = \mathbf{x} \qquad \rightarrow \qquad \mathbf{T}^T \cdot \mathbf{z} = \hat{\mathbf{x}}$$

**Matrix view:**

$$\overset{n \times d}{\underset{\downarrow}{\mathbf{X} = \mathbf{Z}\mathbf{V}^T}} \qquad \rightarrow \qquad \begin{aligned} \hat{\mathbf{X}} &= \mathbf{Z}\ddot{\mathbf{V}}_{d \times l}^T \\ &= \mathbf{Z}\mathbf{T} \end{aligned}$$

# RECONSTRUCTION ERROR

Let us reconstruct $\mathbf{x}$ now:

$$\hat{\mathbf{x}} = \sum_{i=1}^{l} z_i \mathbf{v}_i$$

The error vector $\mathbf{e} = \mathbf{x} - \hat{\mathbf{x}}$ is now

$$\mathbf{e} = \sum_{i=l+1}^{d} z_i \mathbf{v}_i$$

because

$$\mathbf{x} - \hat{\mathbf{x}} = \sum_{i=1}^{d} z_i \mathbf{v}_i - \sum_{i=1}^{l} z_i \mathbf{v}_i$$

We now have

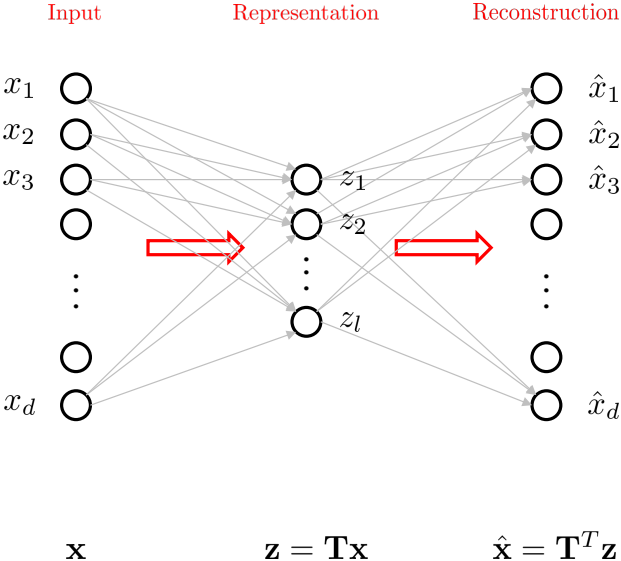$$\mathbb{E}[\boldsymbol{X} - \hat{\boldsymbol{X}}] = \mathbf{0} - \sum_{i=1}^{l} \mathbb{E}[Z_i]\mathbf{v}_i = \mathbf{0}$$

$$\mathbb{E}[||\boldsymbol{X} - \hat{\boldsymbol{X}}||^2] = \sum_{i=l+1}^{d} \mathbf{v}_i^T \boldsymbol{\Sigma} \mathbf{v}_i = \sum_{i=l+1}^{d} \lambda_i \qquad \leftarrow \text{proved later}$$

# PRINCIPAL COMPONENT ANALYSIS AS REPRESENTATION LEARNING

# RELATIONSHIP WITH SINGULAR VALUE DECOMPOSITION (SVD)

$n \times d$
$\downarrow$

Every matrix $\mathbf{X}$ has a SVD: $\mathbf{X} = \mathbf{USV}^T$.

$\mathbf{U}$ = orthogonal, $n \times n$
$\mathbf{S}$ = diagonal, $n \times d$
$\mathbf{V}^T$ = orthogonal, $d \times d$

In MATLAB: $[\mathrm{U, S, V}] = \mathrm{svd(X)}$

Let's look at $\mathbf{X}^T\mathbf{X}$

$$\mathbf{X}^T\mathbf{X} = (\mathbf{USV}^T)^T(\mathbf{USV}^T) = \mathbf{VS}^T\mathbf{SV}^T.$$

Recall, $\frac{1}{n-1}\mathbf{X}^T\mathbf{X}$ is the estimated covariance matrix when $\mathbf{X}$ is normalized

$$\mathbf{\Sigma} = \frac{1}{n-1}\mathbf{X}^T\mathbf{X} = \frac{1}{n-1}\mathbf{VS}^T\mathbf{SV}^T = \mathbf{V\Lambda V}^T.$$

$$\mathbf{\Lambda} = \frac{1}{n-1}\mathbf{S}^T\mathbf{S}. \qquad \leftarrow \text{eigenvalue matrix}$$

# EIGENDECOMPOSITION VS. SINGULAR VALUE DECOMPOSITION

Eigendecomposition: $\quad \frac{1}{n-1}\mathbf{X}^T\mathbf{X} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$

Singular value decomposition: $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^T$

In MATLAB: $[\mathrm{V}, \Lambda] = \mathrm{eig}(\Sigma)$

$[\mathrm{U}, \mathrm{S}, \mathrm{V}] = \mathrm{svd}(\mathrm{X})$

**Q:** Is matrix $\mathbf{V}$ exactly the same in both?

**A:** Should be but not necessarily. Vectors in $\mathbf{V}$ can have opposite directions.

Depends on the software we use.

We were solving the following system:

$$\mathbf{\Sigma V} = \mathbf{V\Lambda}$$

where $\mathbf{V} = [\mathbf{v}_1 \ \mathbf{v}_2 \ \ldots \ \mathbf{v}_d]$, with $\mathbf{V}^T\mathbf{V} = \mathbf{I}$

$\mathbf{\Lambda} = \text{diag}\{\lambda_1, \lambda_2, \ldots, \lambda_d\}$, with $\lambda_1 \geq \lambda_2 \ldots \geq \lambda_d$

**Total complexity:** $O(d^3 + nd^2)$
- ○ computing the covariance matrix ($\mathbf{\Sigma}$): $O(nd^2)$
- ○ computing eigenvectors ($\mathbf{V}$) and eigenvalues ($\mathbf{\Lambda}$): $O(d^3)$

Singular value decomposition takes $O(\min\{nd^2, dn^2\})$

Strang. Introduction to linear algebra. Wellesley-Cambridge Press, 2021.

# HANDLING HIGH-DIMENSIONAL DATA

$n \times d$
$\downarrow$

Consider a centered data matrix $\mathbf{X}$, where $d \gg n$.

$d \times d$
$\downarrow$

$\boldsymbol{\Sigma}$ cannot fit in memory!

Pick now any eigenvalue $\lambda$ and the corresponding eigenvector $\mathbf{v}$

$$\hat{\boldsymbol{\Sigma}}\mathbf{v} = \lambda\mathbf{v}$$

$$\frac{1}{n-1}\mathbf{X}^T\mathbf{X}\mathbf{v} = \lambda\mathbf{v}$$

$$\frac{1}{n-1}\mathbf{X}\mathbf{X}^T\mathbf{X}\mathbf{v} = \lambda\mathbf{X}\mathbf{v}$$

$\Longrightarrow$

$n \times n$
$\downarrow$

$$\frac{1}{n-1}\mathbf{X}\mathbf{X}^T\underbrace{\mathbf{X}\mathbf{v}}_{\underset{n \times 1}{\uparrow}} = \lambda\underbrace{\mathbf{X}\mathbf{v}}_{\underset{n \times 1}{\uparrow}}$$

Note: $\mathbf{X}$ is still column-normalized

# Handling High-Dimensional Data

$$\frac{1}{n-1}\mathbf{X}\mathbf{X}^T \underbrace{\overset{n \times d}{\overset{\downarrow}{\mathbf{X}\ddot{\mathbf{V}}}}}_{\mathbf{Z}} = \underbrace{\mathbf{X}\ddot{\mathbf{V}}}_{\mathbf{Z}}\mathbf{\Lambda}$$

Note: $d \gg n$, so we reduce $\mathbf{V}$ to $\ddot{\mathbf{V}}_{d \times l}$.

Eigenvalues of $\frac{1}{n-1}\mathbf{X}^T\mathbf{X}$ are the same as eigenvalues of $\frac{1}{n-1}\mathbf{X}\mathbf{X}^T$

There are at most $n$ nonzero eigenvalues, for both $\frac{1}{n-1}\mathbf{X}^T\mathbf{X}$ and $\frac{1}{n-1}\mathbf{X}\mathbf{X}^T$

**Solution:**

$$\frac{1}{n-1}\mathbf{X}\mathbf{X}^T\overset{n \times n}{\overset{\downarrow}{\mathbf{W}}} = \mathbf{W}\mathbf{\Lambda}$$

Note: we can reduce $\mathbf{W}$ to $\ddot{\mathbf{W}}_{n \times l}$.

$$\frac{1}{n-1}\mathbf{X}^T\mathbf{X}\underbrace{\mathbf{X}^T\mathbf{W}}_{\mathbf{V}'} = \underbrace{\mathbf{X}^T\mathbf{W}}_{\mathbf{V}'}\mathbf{\Lambda}$$

The norm of each column of $\mathbf{W}$ is 1, but not for $\mathbf{X}^T\mathbf{W}$.

<span style="color:red">$\leftarrow$ we centered $\mathbf{X}$ not $\mathbf{X}^T$</span>

$\mathbf{V} \leftarrow$ normalize($\mathbf{V}'$) so that column norms are 1.

# HANDLING HIGH-DIMENSIONAL DATA

Normalizing $\mathbf{V}'$ has a closed-form formula: $\mathbf{V} = \overset{\underset{n \times d}{\downarrow}}{\mathbf{X}^T} \underset{\underset{n \times n}{\uparrow}}{\mathbf{W}} \cdot \mathrm{diag}\left\{ \sqrt{\overset{\underset{n \times n}{\downarrow}}{\mathbf{W}^T} \mathbf{X} \mathbf{X}^T \mathbf{W}} \right\}$

$$\frac{1}{n-1} \mathbf{X}\mathbf{X}^T \underbrace{\mathbf{X}\ddot{\mathbf{V}}}_{\mathbf{Z}} = \underbrace{\mathbf{X}\ddot{\mathbf{V}}}_{\mathbf{Z}} \mathbf{\Lambda}$$

**Algorithm:**

Solve $\frac{1}{n-1}\mathbf{X}\mathbf{X}^T\mathbf{W} = \mathbf{W}\mathbf{\Lambda}$ to find $\mathbf{\Lambda}$ and $\mathbf{W}$

Keep $l$ columns of $\mathbf{W}$ to obtain $\ddot{\mathbf{W}}_{n \times l}$

$\ddot{\mathbf{V}}_{d \times l} = \mathbf{X}^T \ddot{\mathbf{W}}_{n \times l} \cdot \mathrm{diag}\left\{ \sqrt{\ddot{\mathbf{W}}_{n \times l}^T \mathbf{X} \mathbf{X}^T \ddot{\mathbf{W}}_{n \times l}} \right\}$

$\mathbf{Z} = \mathbf{X}\ddot{\mathbf{V}}_{d \times l}$

**Additional considerations:**

What if $\mathbf{X}$ is sparse with huge $d$ and we cannot center it?

What if some columns of $\mathbf{X}$ are constant?

# APPLICATION: EIGENFACES

**Given:** a set of $n$ images $\mathbf{X}_{n \times d}$, where each row is a flattened matrix.

**Find:** transformation matrix $\mathbf{T}_{l \times d}$.

$\longleftarrow$ a sample from Yale Faces B set, with $n = 5000+$ images of 28 subjects

$\leftarrow$ each row is a sample of 5 images for the same subject

$\leftarrow$ each image is processed to a $48 \times 42$ matrix, so $d = 48 \cdot 42 = 2016$

https://www.face-rec.org

$\mathbf{X}$

Mean image:



First 20 eigenvectors, shown as scaled matrices:

$\mathbf{T}$



Sirovich & Kirby. Low-dimensional procedure for the characterization of human faces. *J Opt Soc Am A*, 1987. Turk & Pentland. Eigenfaces for recognition. *J Cogn Neurosci*, 1991.

Yale Faces B data set



Original      Reconstructed

$l = $   1   2   4   8   16   32   64   128   256

$$R^2 = 1 - \frac{\sum_{i=1}^{n} ||\mathbf{x}_i - \hat{\mathbf{x}}_i||^2}{\sum_{i=1}^{n} ||\mathbf{x}_i - \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_i||^2}$$

$R^2 = 0.1878$     $R^2 = 0.8906$     $R^2 = 0.9949$

Note: reconstruction error is measured on the "training" set

# HOW MANY COMPONENTS TO KEEP?

Yale Faces B data set



95%, 49 components

← 71%, 2 components

← 37%, 1 component ($\frac{\lambda_1}{\sum_{i=1}^{n} \lambda_i} = 0.37$)

99%, 181 components

99.9%, 556 components

It is often better to specify the percent of ratained variance, and not $l$.

# APPLICATION: LATENT SEMANTIC ANALYSIS FOR DOCUMENT RETRIEVAL

**Given:** an $n \times d$ text document matrix $\mathbf{X}$     $n$ = number of documents
$d$ = dictionary size

**Find:** latent semantic spaces for document retrieval and term similarity.

Semantic space = space where "terms and documents that are closely associated are placed near one another" (Deerwester et al., 1990).

| | access | document | retrieval | information | theory | database | indexing | computer | REL | MATCH |
|---|---|---|---|---|---|---|---|---|---|---|
| Doc 1 | x | x | x | | | x | x | | R | |
| Doc 2 | | | | x* | x | | | x* | | M |
| Doc 3 | | | x | x* | | | | x* | R | M |

R = relevant
M = matched

Query:   "IDF in *computer*-based *information* look-up"

$$\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^T \approx \ddot{\mathbf{U}}_{n \times l}\ddot{\mathbf{S}}_{l \times l}\ddot{\mathbf{V}}_{d \times l}^T$$

$$\mathbf{X}^T = \mathbf{V}\mathbf{S}^T\mathbf{U}^T \approx \ddot{\mathbf{V}}_{d \times l}\ddot{\mathbf{S}}_{l \times l}^T\ddot{\mathbf{U}}_{n \times l}^T$$

Term similarities: $\mathbf{X}^T\mathbf{X} \approx \ddot{\mathbf{V}}\ddot{\mathbf{S}}^T\ddot{\mathbf{S}}\ddot{\mathbf{V}}^T$

Document similarities: $\mathbf{X}\mathbf{X}^T \approx \ddot{\mathbf{U}}\ddot{\mathbf{S}}\ddot{\mathbf{S}}^T\ddot{\mathbf{U}}^T$

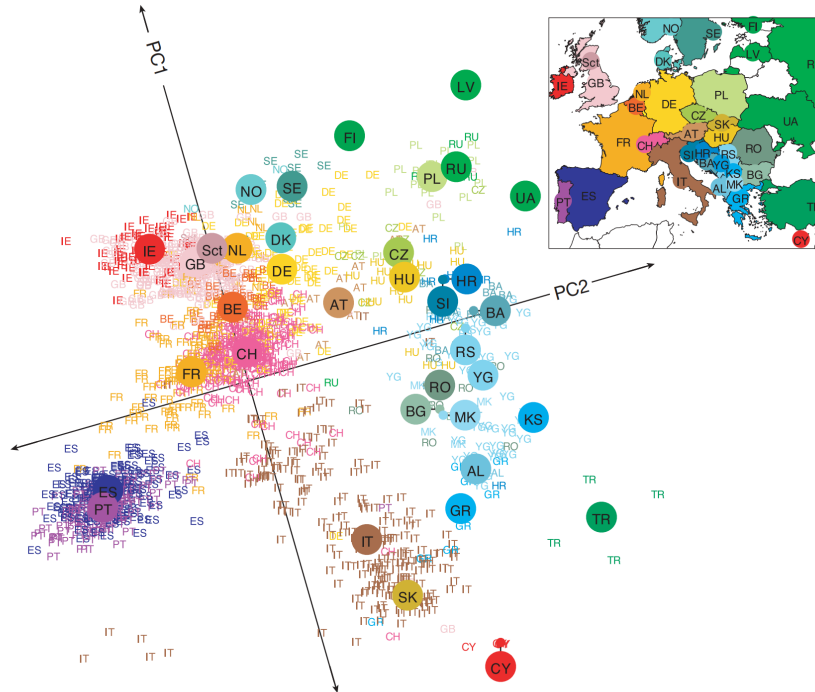Deerwester et al. Indexing by latent semantic analysis. *J Am Soc Inf Sci*, 1990.

# APPLICATION: GENOMIC DATA VISUALIZATION

$\mathbf{X}_{n \times d} =$ sparse matrix

$n = 3192$ subjects
$d = 500568$ genomic loci
$l = 2$ principal directions



Novembre et al. Genes mirror geography within Europe. *Nature*, 2008.

# KERNEL PCA

Consider high-dimensional data. We had $\mathbf{\Gamma} = \frac{1}{n-1}\mathbf{X}\mathbf{X}^T$, where $\Gamma_{ij} = \frac{1}{n-1}\mathbf{x}_i^T\mathbf{x}_j$.

$$\overset{\overset{\textcolor{red}{n \times n}}{\textcolor{red}{\downarrow}}}{}$$

A mapping $\boldsymbol{\varphi} : \mathcal{X} \to \mathcal{F}$ allows us to use any domain for inputs; i.e., $x \in \mathcal{X}$.

$$K_{ij} = \boldsymbol{\varphi}(x_i)^T\boldsymbol{\varphi}(x_j) \qquad \textcolor{red}{\leftarrow \mathbf{K} \text{ is positive semi-definite}}$$

**Problem:** $\frac{1}{n}\sum \mathbf{x}_i = \mathbf{0}$, but $\frac{1}{n}\sum \boldsymbol{\varphi}(x_i)$ is arbitrary.

$$\mathbf{K} \leftarrow \mathbf{K} - \mathbf{1}_n\mathbf{K} - \mathbf{K}\mathbf{1}_n + \mathbf{1}_n\mathbf{K}\mathbf{1}_n \qquad \textcolor{red}{\leftarrow \text{proved later}}$$

$$\mathbf{1}_n = \text{an } n \times n \text{ matrix where each element is } \frac{1}{n}$$

Schölkopf et al. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Comput*, 1998.

# Handling Test Data with Kernel PCA

**Given:** training set $\{x_i\}_{i=1}^n$ and test set $\{t_i\}_{i=1}^m$;

i.e., an $n \times n$ training kernel matrix $\mathbf{K}$ and $m \times n$ test matrix $\mathbf{K}^{\text{test}}$.

$$K_{ij} = \boldsymbol{\varphi}(x_i)^T \boldsymbol{\varphi}(x_j)$$

$$K_{ij}^{\text{test}} = \boldsymbol{\varphi}(t_i)^T \boldsymbol{\varphi}(x_j)$$

**Centering:**

$$\overset{\overset{m \times n}{\downarrow}}{\mathbf{K}^{\text{test}}} \leftarrow \mathbf{K}^{\text{test}} - \mathbf{1}_n' \overset{\overset{n \times n}{\downarrow}}{\mathbf{K}} - \mathbf{K}^{\text{test}} \mathbf{1}_n + \mathbf{1}_n' \mathbf{K} \mathbf{1}_n$$

$\mathbf{1}_n' =$ an $m \times n$ matrix where each element is $\frac{1}{n}$

$\mathbf{1}_n =$ an $n \times n$ matrix where each element is $\frac{1}{n}$

Schölkopf et al. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Comput*, 1998.

# DIFFERENCES BETWEEN KERNEL PCA AND PCA

**Kernel PCA vs. PCA**

- ○ KPCA offers many choices of similarity functions
  - ◇ $\mathbf{K}$ must be symmetric positive semi-definite
- ○ input space for KPCA need not be $\mathbb{R}^d$
  - ◇ KPCA can directly operate on seqeunces, strings, graphs
- ○ classification accuracy often improved over PCA, given $l$
- ○ KPCA allows $l > d$, PCA does not
- ○ loss of interpretability with KPCA
  - ◇ cannot easily visualize eigenvectors for images
  - ◇ requires separate optimization
- ○ computing time a problem for KPCA when $n$ is large
- ○ additional numerical problems with KPCA
  - ◇ centering may cause that $K_{ij} \neq K_{ji}$ which gives complex $\mathbf{\Lambda}$

Schölkopf et al. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Comput*, 1998.

# Appendix: Proof #1

**Proof for the squared norm of the error vector:**

$$\mathbb{E}[||(\boldsymbol{X} - \hat{\boldsymbol{X}})||^2] = \mathbb{E}[(\boldsymbol{X} - \hat{\boldsymbol{X}})^T(\boldsymbol{X} - \hat{\boldsymbol{X}})]$$
$$= \mathbb{E}[\boldsymbol{X}^T\boldsymbol{X}] - 2\mathbb{E}[\boldsymbol{X}^T\hat{\boldsymbol{X}}] + \mathbb{E}[\hat{\boldsymbol{X}}^T\hat{\boldsymbol{X}}]$$

We investigate one of these terms

$$\mathbb{E}[\hat{\boldsymbol{X}}^T\hat{\boldsymbol{X}}] = \mathbb{E}[\sum_{i=1}^{l} Z_i\mathbf{v}_i^T \cdot \sum_{j=1}^{l} Z_j\mathbf{v}_j] = \mathbb{E}[\sum_{i=1}^{l} Z_i^2\mathbf{v}_i^T\mathbf{v}_i] = \mathbb{E}[\sum_{i=1}^{l} Z_i^2]$$

because $\mathbf{v}_i^T\mathbf{v}_j = 0$ when $i \neq j$ and $\mathbf{v}_i^T\mathbf{v}_j = 1$ when $i = j$. This makes a double sum above a single sum.

# APPENDIX: PROOF #1

We similarly have

$$\mathbb{E}[\boldsymbol{X}^T \boldsymbol{X}] = \mathbb{E}[\sum_{i=1}^{d} Z_i \mathbf{v}_i^T \cdot \sum_{j=1}^{d} Z_j \mathbf{v}_j] = \mathbb{E}[\sum_{i=1}^{d} Z_i^2]$$

$$\mathbb{E}[\boldsymbol{X}^T \hat{\boldsymbol{X}}] = \mathbb{E}[\sum_{i=1}^{d} Z_i \mathbf{v}_i^T \cdot \sum_{j=1}^{l} Z_j \mathbf{v}_j] = \mathbb{E}[\sum_{i=1}^{l} Z_i^2]$$

Finally, we have

$$\mathbb{E}[(\boldsymbol{X} - \hat{\boldsymbol{X}})^T (\boldsymbol{X} - \hat{\boldsymbol{X}})] = \sum_{i=1}^{d} \lambda_i - 2\sum_{i=1}^{l} \lambda_i + \sum_{i=1}^{l} \lambda_i = \sum_{i=l+1}^{d} \lambda_i$$
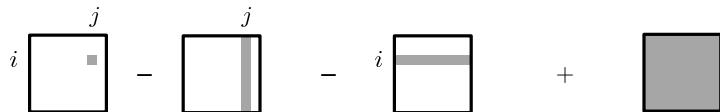
Q.E.D.

# APPENDIX: PROOF #2

**Proof for the kernel normalized in the feature space:**

$$K_{ij} \leftarrow (\boldsymbol{\varphi}(x_i) - \tfrac{1}{n}\sum_{k=1}^{n}\boldsymbol{\varphi}(x_k))^T(\boldsymbol{\varphi}(x_j) - \tfrac{1}{n}\sum_{l=1}^{n}\boldsymbol{\varphi}(x_l)) \qquad \textcolor{red}{\leftarrow \text{centering in feature space}}$$

$$= \boldsymbol{\varphi}(x_i)^T\boldsymbol{\varphi}(x_j) - \tfrac{1}{n}\sum_{l=1}^{n}\boldsymbol{\varphi}(x_i)^T\boldsymbol{\varphi}(x_l) - \tfrac{1}{n}\sum_{l=1}^{n}\boldsymbol{\varphi}(x_k)^T\boldsymbol{\varphi}(x_j) + \tfrac{1}{n^2}\sum_{k=1}^{n}\sum_{l=1}^{n}\boldsymbol{\varphi}(x_k)^T\boldsymbol{\varphi}(x_l)$$

$$= K_{ij} - \tfrac{1}{n}\sum_{k=1}^{n}K_{kj} - \tfrac{1}{n}\sum_{l=1}^{n}K_{il} + \tfrac{1}{n^2}\sum_{k=1}^{n}\sum_{l=1}^{n}K_{kl}$$



**Matrix form:**

$$\mathbf{K} \leftarrow \mathbf{K} - \mathbf{1}_n\mathbf{K} - \mathbf{K}\mathbf{1}_n + \mathbf{1}_n\mathbf{K}\mathbf{1}_n$$

$$\mathbf{1}_n = \text{an } n \times n \text{ matrix where each element is } \tfrac{1}{n}$$

Q.E.D.