

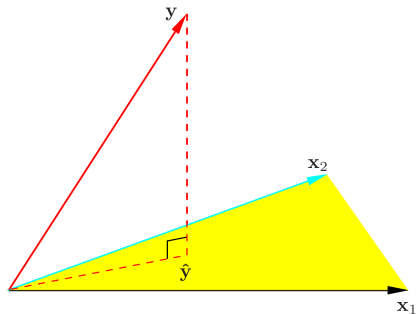
Dimensionality Reduction

Zoltán Szabó – CMAP, École Polytechnique

Data Science @ HEC Paris
May 10, 2019

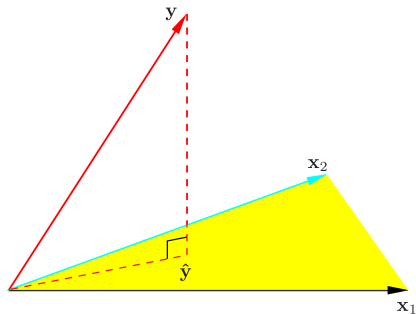
Recall from Tuesday

- We projected to a fixed subspace, $\text{span}(\{\mathbf{x}_i\}_{i=1}^n)$:



Recall from Tuesday

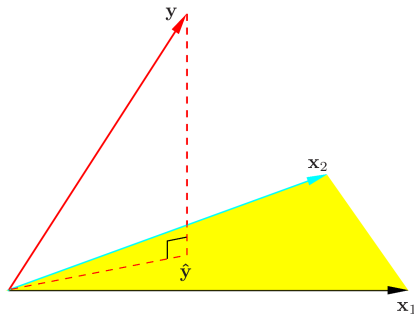
- We projected to a fixed subspace, $\text{span}(\{\mathbf{x}_i\}_{i=1}^n)$:



- Non-linear extensions:
 - $\varphi(x)$: explicit,

Recall from Tuesday

- We projected to a fixed subspace, $\text{span}(\{\mathbf{x}_i\}_{i=1}^n)$:



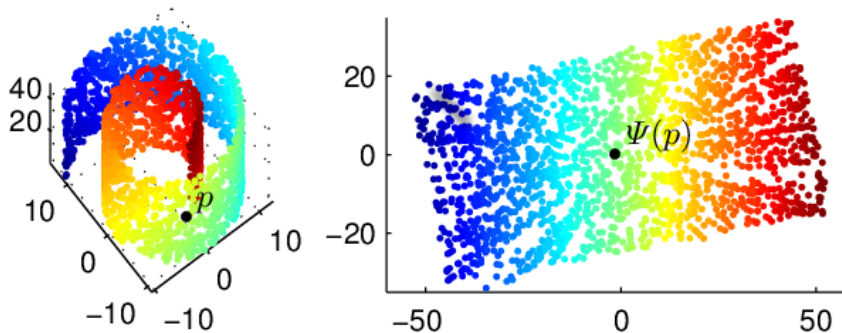
- Non-linear extensions:
 - $\varphi(x)$: explicit,
 - $k(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{H}_k}$
 - implicit usage of features,
 - $\mathcal{H}_k = \overline{\{\sum_{i=1}^n \alpha_i \varphi(x_i)\}}$.

Today: dimensionality reduction

- Given: a set of observations $X = \{\mathbf{x}_i\}_{i=1}^n \subset \mathbb{R}^D$.
- Goal: find $X' = \{\mathbf{x}'_i\}_{i=1}^n \subset \mathbb{R}^d$ 'preserving' the geometry of X .
- $d \ll D$: compression (images, music, ...).



Dimensionality reduction = manifold learning



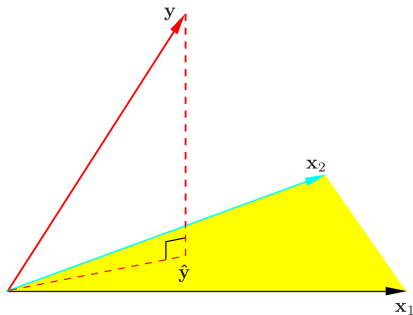
Why?

- Visualization, computational reason, noise reduction.

Why?

- Visualization, computational reason, noise reduction.
- Simplest example:

We optimize the subspace of projection (PCA).



Principal Component Analysis (PCA)

PCA example: 100%

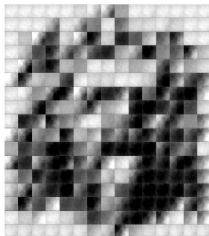


(A)

PCA example: 100% \rightarrow 1%



(A)

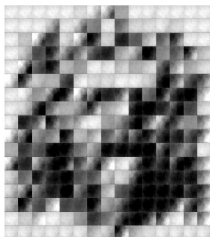


(B)

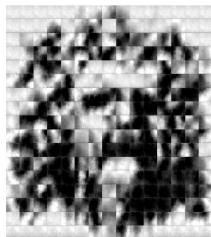
PCA example: 100% \rightarrow 2%



(A)



(B)

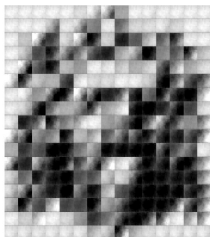


(C)

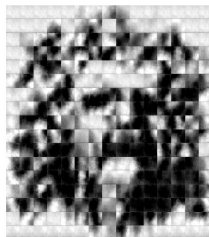
PCA example: 100% \rightarrow 5%



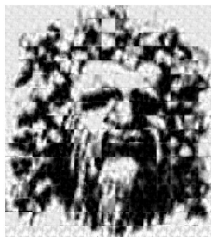
(A)



(B)



(C)

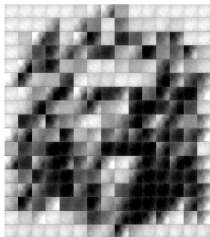


(D)

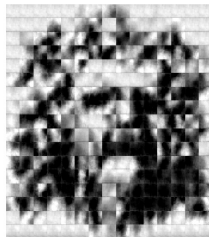
PCA example: 100% \rightarrow 10%



(A)



(B)



(C)



(D)

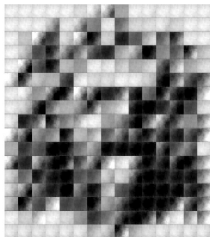


(E)

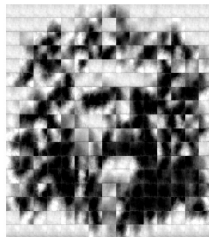
PCA example: 100% \rightarrow 20%



(A)



(B)



(C)



(D)

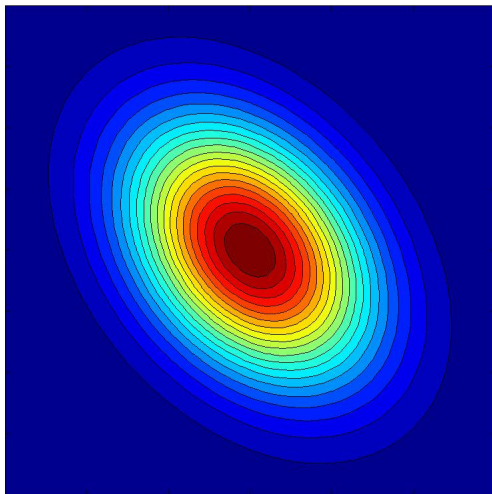


(E)



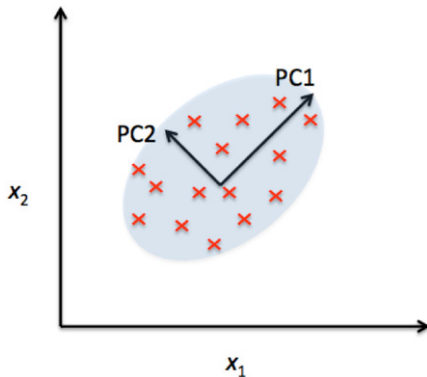
(F)

Conjecture? Most important direction?



PCA: intuition

Task: find the best d -dimensional subspace approximating $\{\mathbf{x}_i\}_{i=1}^n \subset \mathbb{R}^D$.



Cov, var, corr: properties – recall

- Covariance:

$$\text{cov}(x, y) = \mathbb{E}_{xy}[(x - \mathbb{E}x)(y - \mathbb{E}y)].$$

Cov, var, corr: properties – recall

- Covariance: \rightarrow values? $\text{cov}(ax, by) = ?$

$$\text{cov}(x, y) = \mathbb{E}_{xy}[(x - \mathbb{E}x)(y - \mathbb{E}y)].$$

Cov, var, corr: properties – recall

- Covariance:

$$\text{cov}(x, y) = \mathbb{E}_{xy}[(x - \mathbb{E}x)(y - \mathbb{E}y)].$$

- Variance:

$$\text{var}(x) = \text{cov}(x, x) = \mathbb{E}(x^2) - \mathbb{E}^2(x)$$

Cov, var, corr: properties – recall

- Covariance:

$$\text{cov}(x, y) = \mathbb{E}_{xy}[(x - \mathbb{E}x)(y - \mathbb{E}y)].$$

- Variance, std: \rightarrow values? min?

$$\text{var}(x) = \text{cov}(x, x) = \mathbb{E}(x^2) - \mathbb{E}^2(x), \quad \sigma(x) = \sqrt{\text{var}(x)}.$$

Cov, var, corr: properties – recall

- Covariance:

$$\text{cov}(x, y) = \mathbb{E}_{xy}[(x - \mathbb{E}x)(y - \mathbb{E}y)].$$

- Variance, std:

$$\text{var}(x) = \text{cov}(x, x) = \mathbb{E}(x^2) - \mathbb{E}^2(x), \quad \sigma(x) = \sqrt{\text{var}(x)}.$$

- Correlation:

$$\text{corr}(x, y) = \frac{\text{cov}(x, y)}{\sigma(x)\sigma(y)}.$$

Cov, var, corr: properties – recall

- Covariance:

$$\text{cov}(x, y) = \mathbb{E}_{xy}[(x - \mathbb{E}x)(y - \mathbb{E}y)].$$

- Variance, std:

$$\text{var}(x) = \text{cov}(x, x) = \mathbb{E}(x^2) - \mathbb{E}^2(x), \quad \sigma(x) = \sqrt{\text{var}(x)}.$$

- Correlation: \rightarrow intuition? values? max? zero?

$$\text{corr}(x, y) = \frac{\text{cov}(x, y)}{\sigma(x)\sigma(y)}.$$

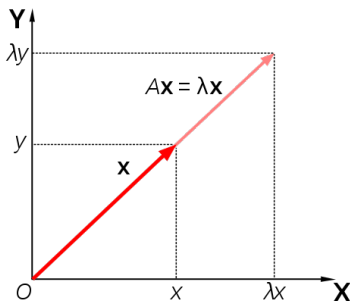
Eigenvectors, eigenvalues – recall

- Simplest transformation: scaling.

Eigenvectors, eigenvalues – recall

- Simplest transformation: scaling.
- $\mathbf{x} \neq \mathbf{0}$ is an eigenvector of \mathbf{A} with eigenvalue $\lambda \in \mathbb{R}$ if

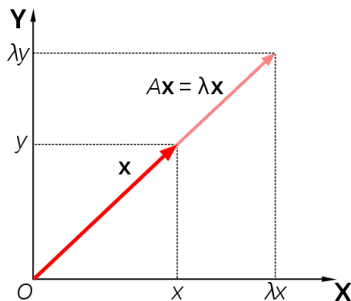
$$\mathbf{Ax} = \lambda\mathbf{x}.$$



Eigenvectors, eigenvalues – recall

- Simplest transformation: scaling.
- $\mathbf{x} \neq \mathbf{0}$ is an eigenvector of \mathbf{A} with eigenvalue $\lambda \in \mathbb{R}$ if

$$\mathbf{Ax} = \lambda\mathbf{x}.$$



- Size of \mathbf{A} ?

Eigensystems: continued

Examples:

- Identity: $\mathbf{A} = \mathbf{I}$.

Eigensystems: continued

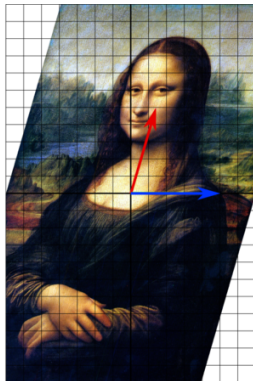
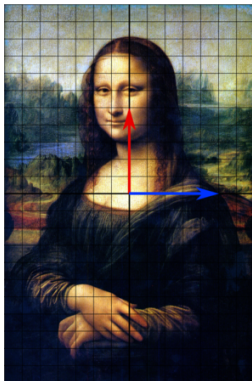
Examples:

- Identity: $\mathbf{A} = \mathbf{I}$.
- Diagonal matrix: $\mathbf{A} = \text{diag}(a_i)$, spec: reflection.

Eigensystems: continued

Examples:

- Identity: $\mathbf{A} = \mathbf{I}$.
- Diagonal matrix: $\mathbf{A} = \text{diag}(a_i)$, spec: reflection.
- Shear mapping on Mona Lisa:



- Diagonal matrix: we saw that the eigensystem is orthogonal.

- Diagonal matrix: we saw that the eigensystem is orthogonal.
- A symmetric \mathbf{A} ($\mathbf{A} = \mathbf{A}^T$) behaves similarly:

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{U}^T,$$

where $\mathbf{\Sigma} = \text{diag}(\lambda_i)$, \mathbf{U} : orthogonal.

Let us apply these observations in PCA!

- We are looking for the best **one-dimensional projection**.



- $\mathbb{E} :=$ empirical/population expectation: $\mathbb{E}\mathbf{x} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$.
- Assumption: $\mathbb{E}\mathbf{x} = \mathbf{0}$.

- We are looking for the best **one-dimensional projection**.



- $\mathbb{E} :=$ empirical/population expectation: $\mathbb{E}\mathbf{x} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$.
- Assumption: $\mathbb{E}\mathbf{x} = \mathbf{0}$.
 - centering: $\mathbf{x} \rightarrow \mathbf{x} - \mathbb{E}\mathbf{x}$.

Projection ($\|\mathbf{w}\|_2 = 1$):

- $\hat{\mathbf{x}} = \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}$.
- zero mean: $\mathbf{0} \stackrel{?}{=} \mathbb{E} \hat{\mathbf{x}} = \mathbb{E} [\langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}]$

Projection ($\|\mathbf{w}\|_2 = 1$):

- $\hat{\mathbf{x}} = \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}$.
- zero mean: $\mathbf{0} \stackrel{?}{=} \mathbb{E} \hat{\mathbf{x}} = \mathbb{E} [\langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}] = \langle \mathbf{w}, \underbrace{\mathbb{E} \mathbf{x}}_{=\mathbf{0}} \rangle \mathbf{w}$.

PCA: min residual \Leftrightarrow max squared projection

- Goal: $\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \rightarrow \min_{\mathbf{w}}$.

PCA: min residual \Leftrightarrow max squared projection

- Goal: $\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \rightarrow \min_{\mathbf{w}}$.
- Residual \Rightarrow objective:

$$\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \|\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}\|_2^2$$

PCA: min residual \Leftrightarrow max squared projection

- Goal: $\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \rightarrow \min_{\mathbf{w}}$.
- Residual \Rightarrow objective:

$$\begin{aligned}\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \|\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}\|_2^2 \\ &\stackrel{\|\mathbf{w}\|_2^2=1}{=} \|\mathbf{x}\|_2^2 - \langle \mathbf{w}, \mathbf{x} \rangle^2 \Rightarrow\end{aligned}$$

PCA: min residual \Leftrightarrow max squared projection

- Goal: $\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \rightarrow \min_{\mathbf{w}}$.
- Residual \Rightarrow objective:

$$\begin{aligned}\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \|\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}\|_2^2 \\ &\stackrel{\|\mathbf{w}\|_2^2=1}{=} \|\mathbf{x}\|_2^2 - \langle \mathbf{w}, \mathbf{x} \rangle^2 \Rightarrow \\ \mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \mathbb{E} \left[\|\mathbf{x}\|_2^2 - \langle \mathbf{w}, \mathbf{x} \rangle^2 \right] = \underbrace{\mathbb{E} \|\mathbf{x}\|_2^2}_{\text{independent of } \mathbf{w}} - \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 \Leftrightarrow\end{aligned}$$

PCA: min residual \Leftrightarrow max squared projection

- Goal: $\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \rightarrow \min_{\mathbf{w}}$.
- Residual \Rightarrow objective:

$$\begin{aligned}\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \|\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}\|_2^2 \\ &\stackrel{\|\mathbf{w}\|_2^2=1}{=} \|\mathbf{x}\|_2^2 - \langle \mathbf{w}, \mathbf{x} \rangle^2 \Rightarrow \\ \mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \mathbb{E} \left[\|\mathbf{x}\|_2^2 - \langle \mathbf{w}, \mathbf{x} \rangle^2 \right] = \underbrace{\mathbb{E} \|\mathbf{x}\|_2^2}_{\text{independent of } \mathbf{w}} - \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 \Leftrightarrow\end{aligned}$$

Solution

maximizes the mean squared projection.

PCA: max squared projection \Leftrightarrow max variance of projection

By using $\mathbb{E}y^2 = (\mathbb{E}y)^2 + \text{var}(y)$:

$$\max_{\mathbf{w}} \leftarrow \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 = \underbrace{(\mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle)^2}_{=0} + \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle).$$

PCA: max squared projection \Leftrightarrow max variance of projection

By using $\mathbb{E}y^2 = (\mathbb{E}y)^2 + \text{var}(y)$:

$$\max_{\mathbf{w}} \leftarrow \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 = \underbrace{\left(\mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle \right)^2}_{=0} + \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle).$$

To sum up:

Minimize MSE of the residual : $\min_{\mathbf{w}} \mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \Leftrightarrow$

PCA: max squared projection \Leftrightarrow max variance of projection

By using $\mathbb{E}y^2 = (\mathbb{E}y)^2 + \text{var}(y)$:

$$\max_{\mathbf{w}} \leftarrow \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 = \underbrace{\left(\mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle \right)^2}_{=0} + \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle).$$

To sum up:

Minimize MSE of the residual : $\min_{\mathbf{w}} \mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \Leftrightarrow$

Maximize mean squared projection : $\max_{\mathbf{w}} \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 \Leftrightarrow$

PCA: max squared projection \Leftrightarrow max variance of projection

By using $\mathbb{E}y^2 = (\mathbb{E}y)^2 + \text{var}(y)$:

$$\max_{\mathbf{w}} \leftarrow \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 = \underbrace{\left(\mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle \right)^2}_{=0} + \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle).$$

To sum up:

Minimize MSE of the residual : $\min_{\mathbf{w}} \mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \Leftrightarrow$

Maximize mean squared projection : $\max_{\mathbf{w}} \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 \Leftrightarrow$

Maximize variance of the projection : $\max_{\mathbf{w}} \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle).$

By the bilinearity of cov:

$$\text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) = \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x})$$

PCA: optimization

By the bilinearity of cov:

$$\text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) = \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x}) = \mathbf{w}^T \text{cov}(\mathbf{x}) \mathbf{w} =: \mathbf{w}^T \Sigma \mathbf{w} \rightarrow \max_{\|\mathbf{w}\|_2=1} .$$

PCA: optimization

By the bilinearity of cov:

$$\text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) = \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x}) = \mathbf{w}^T \text{cov}(\mathbf{x}) \mathbf{w} =: \mathbf{w}^T \Sigma \mathbf{w} \rightarrow \max_{\|\mathbf{w}\|_2=1}.$$

Lagrange function, solving for 'derivatives = 0':

PCA: optimization

By the bilinearity of cov:

$$\text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) = \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x}) = \mathbf{w}^T \text{cov}(\mathbf{x}) \mathbf{w} =: \mathbf{w}^T \Sigma \mathbf{w} \rightarrow \max_{\|\mathbf{w}\|_2=1}.$$

Lagrange function, solving for 'derivatives = 0':

$$L(\mathbf{w}, \lambda) = \underbrace{\mathbf{w}^T \Sigma \mathbf{w}}_{=\text{objective}} - \lambda(\underbrace{\mathbf{w}^T \mathbf{w} - 1}_{=\text{condition}}) \Rightarrow$$

PCA: optimization

By the bilinearity of cov:

$$\text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) = \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x}) = \mathbf{w}^T \text{cov}(\mathbf{x}) \mathbf{w} =: \mathbf{w}^T \Sigma \mathbf{w} \rightarrow \max_{\|\mathbf{w}\|_2=1}.$$

Lagrange function, solving for 'derivatives = 0':

$$\begin{aligned} L(\mathbf{w}, \lambda) &= \underbrace{\mathbf{w}^T \Sigma \mathbf{w}}_{=\text{objective}} - \lambda (\underbrace{\mathbf{w}^T \mathbf{w} - 1}_{=\text{condition}}) \Rightarrow \\ 0 &= \frac{\partial L(\mathbf{w}, \lambda)}{\partial \lambda} = -(\mathbf{w}^T \mathbf{w} - 1), \end{aligned}$$

PCA: optimization

By the bilinearity of cov:

$$\text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) = \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x}) = \mathbf{w}^T \text{cov}(\mathbf{x}) \mathbf{w} =: \mathbf{w}^T \mathbf{\Sigma} \mathbf{w} \rightarrow \max_{\|\mathbf{w}\|_2=1}.$$

Lagrange function, solving for 'derivatives = 0':

$$\begin{aligned} L(\mathbf{w}, \lambda) &= \underbrace{\mathbf{w}^T \mathbf{\Sigma} \mathbf{w}}_{=\text{objective}} - \lambda (\underbrace{\mathbf{w}^T \mathbf{w} - 1}_{=\text{condition}}) \Rightarrow \\ 0 &= \frac{\partial L(\mathbf{w}, \lambda)}{\partial \lambda} = -(\mathbf{w}^T \mathbf{w} - 1), \\ 0 &= \frac{\partial L(\mathbf{w}, \lambda)}{\partial \mathbf{w}} = 2\mathbf{\Sigma} \mathbf{w} - 2\lambda \mathbf{w} \Rightarrow \end{aligned}$$

PCA: optimization

By the bilinearity of cov:

$$\text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) = \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x}) = \mathbf{w}^T \text{cov}(\mathbf{x}) \mathbf{w} =: \mathbf{w}^T \mathbf{\Sigma} \mathbf{w} \rightarrow \max_{\|\mathbf{w}\|_2=1}.$$

Lagrange function, solving for 'derivatives = 0':

$$\begin{aligned} L(\mathbf{w}, \lambda) &= \underbrace{\mathbf{w}^T \mathbf{\Sigma} \mathbf{w}}_{=\text{objective}} - \lambda(\underbrace{\mathbf{w}^T \mathbf{w} - 1}_{=\text{condition}}) \Rightarrow \\ 0 &= \frac{\partial L(\mathbf{w}, \lambda)}{\partial \lambda} = -(\mathbf{w}^T \mathbf{w} - 1), \\ 0 &= \frac{\partial L(\mathbf{w}, \lambda)}{\partial \mathbf{w}} = 2\mathbf{\Sigma} \mathbf{w} - 2\lambda \mathbf{w} \Rightarrow \end{aligned}$$

Solution

\mathbf{w}^* : eigenvector associated to $\lambda_{\max}(\mathbf{\Sigma})$.

- Goal: approximate with a d -dimensional subspace.
- ONB in the subspace ($\mathbf{W}^T \mathbf{W} = \mathbf{I}$):

$$\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_d] \in \mathbb{R}^{D \times d},$$

- Approximation:

$$\hat{\mathbf{x}} = \sum_{i=1}^d \langle \mathbf{w}_i, \mathbf{x} \rangle \mathbf{w}_i = \mathbf{W} \mathbf{W}^T \mathbf{x}.$$

After similar calculation than for $d = 1 \dots$

- The d principal components:

$$\{\mathbf{w}_i\}_{i=1}^d = \text{top } d \text{ eigenvectors of } \text{cov}(\mathbf{x}).$$

- The d principal components:

$$\{\mathbf{w}_i\}_{i=1}^d = \text{top } d \text{ eigenvectors of } \text{cov}(\mathbf{x}).$$

- $\mathbf{\Sigma} := \text{cov}(\mathbf{x})$: symmetric, positive semi-definite $\Rightarrow \{\mathbf{w}_i\}$: ONS,
 $\lambda_i \geq 0$.

- The d principal components:

$$\{\mathbf{w}_i\}_{i=1}^d = \text{top } d \text{ eigenvectors of } \text{cov}(\mathbf{x}).$$

- $\mathbf{\Sigma} := \text{cov}(\mathbf{x})$: symmetric, positive semi-definite $\Rightarrow \{\mathbf{w}_i\}$: ONS, $\lambda_i \geq 0$.
- Variance decomposition: $\text{cov}(\mathbf{x}) = \sum_{i=1}^D \lambda_i \mathbf{w}_i \mathbf{w}_i^T$.

- The d principal components:

$$\{\mathbf{w}_i\}_{i=1}^d = \text{top } d \text{ eigenvectors of } \text{cov}(\mathbf{x}).$$

- $\Sigma := \text{cov}(\mathbf{x})$: symmetric, positive semi-definite $\Rightarrow \{\mathbf{w}_i\}$: ONS, $\lambda_i \geq 0$.
- Variance decomposition: $\text{cov}(\mathbf{x}) = \sum_{i=1}^D \lambda_i \mathbf{w}_i \mathbf{w}_i^T$.
- Energy preserved using d components: $\sum_{i=1}^d \lambda_i \Rightarrow$

$$R = R(d) := \frac{\sum_{i=1}^d \lambda_i}{\sum_{i=1}^D \lambda_i} \in [0, 1].$$

- The d principal components:

$$\{\mathbf{w}_i\}_{i=1}^d = \text{top } d \text{ eigenvectors of } \text{cov}(\mathbf{x}).$$

- $\mathbf{\Sigma} := \text{cov}(\mathbf{x})$: symmetric, positive semi-definite $\Rightarrow \{\mathbf{w}_i\}$: ONS, $\lambda_i \geq 0$.
- Variance decomposition: $\text{cov}(\mathbf{x}) = \sum_{i=1}^D \lambda_i \mathbf{w}_i \mathbf{w}_i^T$.
- Energy preserved using d components: $\sum_{i=1}^d \lambda_i \Rightarrow$

$$R = R(d) := \frac{\sum_{i=1}^d \lambda_i}{\sum_{i=1}^D \lambda_i} \in [0, 1].$$

- In practice: choose d such that $R \approx 0.8 - 0.9$.

Non-linear PCA

- PCA:
 - objective: maximize the variance of the projection.
 - solution: leading eigenvectors of $\mathbf{\Sigma} = \text{cov}(\mathbf{x})$.

Non-linear PCA: idea

- PCA:
 - objective: maximize the variance of the projection.
 - solution: leading eigenvectors of $\mathbf{\Sigma} = \text{cov}(\mathbf{x})$.
- Non-linear PCA:
 - Take $\varphi(\mathbf{x})$.

Non-linear PCA: idea

- PCA:
 - objective: maximize the variance of the projection.
 - solution: leading eigenvectors of $\Sigma = \text{cov}(\mathbf{x})$.
- Non-linear PCA:
 - Take $\varphi(\mathbf{x})$.
 - What is $\Sigma := \text{cov}(\varphi(\mathbf{x}))$?













Non-linear PCA: idea

- PCA:
 - objective: maximize the variance of the projection.
 - solution: leading eigenvectors of $\Sigma = \text{cov}(\mathbf{x})$.
- Non-linear PCA:
 - Take $\varphi(\mathbf{x})$.
 - What is $\Sigma := \text{cov}(\varphi(\mathbf{x}))$?
 - Eigenvectors of an operator?

Non-linear PCA: idea

- PCA:
 - objective: maximize the variance of the projection.
 - solution: leading eigenvectors of $\Sigma = \text{cov}(\mathbf{x})$.
- Non-linear PCA:
 - Take $\varphi(\mathbf{x})$.
 - What is $\Sigma := \text{cov}(\varphi(\mathbf{x}))$?
 - Eigenvectors of an operator?
 - Computational tractability?

In denoising application: PCA vs non-linear PCA

		Gaussian noise									
orig.											
noisy											
$n = 1$											
4											
16											
64											
256											
$n = 1$											
4											
16											
64											
256											

Kernel PCA: idea for ' $d = 1$ ' $\leftrightarrow f$

Let $\mathcal{H} = \mathcal{H}_k$.

- Objective function:

$$J(f) = \frac{1}{n} \sum_{i=1}^n \left\langle f, \underbrace{\varphi(x_i) - \frac{1}{n} \sum_{j=1}^n \varphi(x_j)}_{=:\tilde{\varphi}(x_i)} \right\rangle^2 = \text{var}(f) \rightarrow \max_{f: \|f\|_{\mathcal{H}} \leq 1}.$$

Kernel PCA: idea for ' $d = 1$ ' $\leftrightarrow f$

Let $\mathcal{H} = \mathcal{H}_k$.

- Objective function:

$$J(f) = \frac{1}{n} \sum_{i=1}^n \left\langle f, \underbrace{\varphi(x_i) - \frac{1}{n} \sum_{j=1}^n \varphi(x_j)}_{=: \tilde{\varphi}(x_i)} \right\rangle^2 = \text{var}(f) \rightarrow \max_{f: \|f\|_{\mathcal{H}} \leq 1}.$$

- The solution can be searched in the form ($\mathcal{H} \ni f \leftrightarrow \mathbf{a} \in \mathbb{R}^n$):

$$f = \sum_{i=1}^n a_i \tilde{\varphi}(x_i)$$

since component $\perp \text{span}(\{\tilde{\varphi}(x_i)\}_{i=1}^n)$ has no contribution.

Kernel PCA: idea for ' $d = 1$ ' $\leftrightarrow f$

Let $\mathcal{H} = \mathcal{H}_k$.

- Objective function:

$$J(f) = \frac{1}{n} \sum_{i=1}^n \left\langle f, \underbrace{\varphi(x_i) - \frac{1}{n} \sum_{j=1}^n \varphi(x_j)}_{=: \tilde{\varphi}(x_i)} \right\rangle^2 = \text{var}(f) \rightarrow \max_{f: \|f\|_{\mathcal{H}} \leq 1}.$$

- The solution can be searched in the form ($\mathcal{H} \ni f \leftrightarrow \mathbf{a} \in \mathbb{R}^n$):

$$f = \sum_{i=1}^n a_i \tilde{\varphi}(x_i)$$

since component $\perp \text{span}(\{\tilde{\varphi}(x_i)\}_{i=1}^n)$ has no contribution.

- We will get an **eigenvalue problem for \mathbf{a}** .

(Empirical) covariance operator

$$C := \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \otimes \tilde{\varphi}(x_i).$$

$c \otimes d$ is the analogue of cd^T :

$$(c \otimes d)(e) = c \langle d, e \rangle_{\mathcal{H}}.$$

(Empirical) covariance operator

$$C := \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \otimes \tilde{\varphi}(x_i).$$

$c \otimes d$ is the analogue of cd^T :

$$(c \otimes d)(e) = c \langle d, e \rangle_{\mathcal{H}}.$$

Similarly to the finite-dimensional case:

$$Cf_j = \lambda_j f_j.$$

Challenge

How do we solve this **eigenvalue problem**?

Assume j is fixed ($Cf = \lambda f$):

$$Cf = \left[\frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \otimes \tilde{\varphi}(x_i) \right] f$$

Assume j is fixed ($Cf = \lambda f$):

$$\begin{aligned} Cf &= \left[\frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \otimes \tilde{\varphi}(x_i) \right] f \\ &\stackrel{\otimes}{=} \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \left\langle \tilde{\varphi}(x_i), \sum_{j=1}^n a_j \tilde{\varphi}(x_j) \right\rangle_{\mathcal{H}} \end{aligned}$$

Computation of Cf_j

Assume j is fixed ($Cf = \lambda f$):

$$Cf = \left[\frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \otimes \tilde{\varphi}(x_i) \right] f$$
$$\stackrel{\otimes \text{def}}{=} \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \left\langle \tilde{\varphi}(x_i), \sum_{j=1}^n a_j \tilde{\varphi}(x_j) \right\rangle_{\mathcal{H}} = \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \sum_{j=1}^n a_j \tilde{k}(x_i, x_j)$$

with $\tilde{\mathbf{G}} = \mathbf{H}\mathbf{G}\mathbf{H} = \left[\tilde{k}(x_i, x_j) \right]_{i,j=1}^n$, $\mathbf{H} = \mathbf{I}_n - \frac{\mathbf{E}_n}{n}$, $\mathbf{E}_n = [1] \in \mathbb{R}^{n \times n}$.

Computation of Cf_j

Assume j is fixed ($Cf = \lambda f$):

$$Cf = \left[\frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \otimes \tilde{\varphi}(x_i) \right] f$$
$$\stackrel{\otimes \text{def}}{=} \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \left\langle \tilde{\varphi}(x_i), \sum_{j=1}^n a_j \tilde{\varphi}(x_j) \right\rangle_{\mathcal{H}} = \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \sum_{j=1}^n a_j \tilde{k}(x_i, x_j)$$

with $\tilde{\mathbf{G}} = \mathbf{H}\mathbf{G}\mathbf{H} = \left[\tilde{k}(x_i, x_j) \right]_{i,j=1}^n$, $\mathbf{H} = \mathbf{I}_n - \frac{\mathbf{E}_n}{n}$, $\mathbf{E}_n = [1] \in \mathbb{R}^{n \times n}$.

Since $f = \sum_{j=1}^n a_j \tilde{\varphi}(x_j)$

multiplying by $\tilde{\varphi}(x_r)$ [$r = 1, \dots, n$] gives expressions in terms of $\tilde{\mathbf{G}}$.

Eigenvalue problem

- We want to solve $Cf = \lambda f$; Cf and f : functions of $\tilde{\varphi}(x_i)$.
- By multiplying with $\tilde{\varphi}(x_r)$:

$$\langle \tilde{\varphi}(x_r), \lambda f \rangle_{\mathcal{H}} = \lambda (\tilde{\mathbf{G}}\mathbf{a})_r,$$

$$\langle \tilde{\varphi}(x_r), Cf \rangle_{\mathcal{H}} = \frac{1}{n} (\tilde{\mathbf{G}}^2 \mathbf{a})_r.$$

- Eigenvalue problem: $\tilde{\mathbf{G}}^2 \mathbf{a} = n\lambda \tilde{\mathbf{G}}\mathbf{a}$, i.e. $\tilde{\mathbf{G}}\mathbf{a} = (n\lambda)\mathbf{a}$.

Orthogonal eigenvectors in kernel PCA

Taking two eigenvectors:

$$\mathbf{f}_1 = \sum_{i=1}^n a_{1i} \tilde{\varphi}(x_i),$$

$$\tilde{\mathbf{G}} \mathbf{a}_1 = \lambda_1 \mathbf{a}_1,$$

$$\mathbf{f}_2 = \sum_{j=1}^n a_{2j} \tilde{\varphi}(x_j),$$

$$\tilde{\mathbf{G}} \mathbf{a}_2 = \lambda_2 \mathbf{a}_2.$$

one has

$$0 \stackrel{?}{=} \langle \mathbf{f}_1, \mathbf{f}_2 \rangle_{\mathcal{H}}$$

Orthogonal eigenvectors in kernel PCA

Taking two eigenvectors:

$$\begin{aligned} f_1 &= \sum_{i=1}^n a_{1i} \tilde{\varphi}(x_i), & \tilde{\mathbf{G}} \mathbf{a}_1 &= \lambda_1 \mathbf{a}_1, \\ f_2 &= \sum_{j=1}^n a_{2j} \tilde{\varphi}(x_j), & \tilde{\mathbf{G}} \mathbf{a}_2 &= \lambda_2 \mathbf{a}_2. \end{aligned}$$

one has

$$0 \stackrel{?}{=} \langle f_1, f_2 \rangle_{\mathcal{H}} = \left\langle \sum_{i=1}^n a_{1i} \tilde{\varphi}(x_i), \sum_{j=1}^n a_{2j} \tilde{\varphi}(x_j) \right\rangle_{\mathcal{H}}$$

Orthogonal eigenvectors in kernel PCA

Taking two eigenvectors:

$$\begin{aligned} \mathbf{f}_1 &= \sum_{i=1}^n a_{1i} \tilde{\varphi}(x_i), & \tilde{\mathbf{G}} \mathbf{a}_1 &= \lambda_1 \mathbf{a}_1, \\ \mathbf{f}_2 &= \sum_{j=1}^n a_{2j} \tilde{\varphi}(x_j), & \tilde{\mathbf{G}} \mathbf{a}_2 &= \lambda_2 \mathbf{a}_2. \end{aligned}$$

one has

$$0 \stackrel{?}{=} \langle f_1, f_2 \rangle_{\mathcal{H}} = \left\langle \sum_{i=1}^n a_{1i} \tilde{\varphi}(x_i), \sum_{j=1}^n a_{2j} \tilde{\varphi}(x_j) \right\rangle_{\mathcal{H}} = \mathbf{a}_1^T \tilde{\mathbf{G}} \mathbf{a}_2$$

Orthogonal eigenvectors in kernel PCA

Taking two eigenvectors:

$$\mathbf{f}_1 = \sum_{i=1}^n a_{1i} \tilde{\varphi}(x_i),$$

$$\tilde{\mathbf{G}}\mathbf{a}_1 = \lambda_1 \mathbf{a}_1,$$

$$\mathbf{f}_2 = \sum_{j=1}^n a_{2j} \tilde{\varphi}(x_j),$$

$$\tilde{\mathbf{G}}\mathbf{a}_2 = \lambda_2 \mathbf{a}_2.$$

one has

$$0 \stackrel{?}{=} \langle f_1, f_2 \rangle_{\mathcal{H}} = \left\langle \sum_{i=1}^n a_{1i} \tilde{\varphi}(x_i), \sum_{j=1}^n a_{2j} \tilde{\varphi}(x_j) \right\rangle_{\mathcal{H}} = \mathbf{a}_1^T \tilde{\mathbf{G}} \mathbf{a}_2 = \mathbf{a}_1^T \lambda_2 \mathbf{a}_2.$$

Orthogonality \Rightarrow projection is easy

- Projection of a new x^* to the first d -PCs:

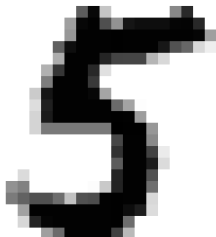
$$\Pi [\tilde{\varphi}(x^*)] = \sum_{j=1}^d \langle \tilde{\varphi}(x^*), f_j \rangle_{\mathcal{H}} f_j.$$

Orthogonality \Rightarrow projection is easy

- Projection of a new x^* to the first d -PCs:

$$\Pi [\tilde{\varphi}(x^*)] = \sum_{j=1}^d \langle \tilde{\varphi}(x^*), f_j \rangle_{\mathcal{H}} f_j.$$

- The pre-image problem we solved in denoising:



$$\widehat{x^*} = \arg \min_{x \in \mathcal{X}} \|\tilde{\varphi}(x) - \Pi [\tilde{\varphi}(x^*)]\|_{\mathcal{H}}^2.$$

Canonical Correlation Analysis (CCA)

- Given a pair of random variables: $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{2d}$.
- Find the directions $(\mathbf{a} \in \mathbb{R}^d, \mathbf{b} \in \mathbb{R}^d)$ in which \mathbf{x} and \mathbf{y} are maximally correlated:

$$\text{CCA}(\mathbf{x}, \mathbf{y}) = \max_{\mathbf{a}, \mathbf{b}} \text{corr}_{\mathbf{x}, \mathbf{y}}(\mathbf{a}^T \mathbf{x}, \mathbf{b}^T \mathbf{y}) .$$

Examples

follow where dependence measures are useful!

Outlier-robust image registration

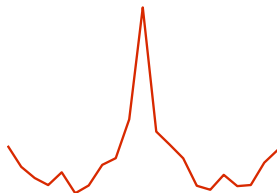
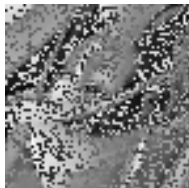
Given two images:



Goal: find the transformation which takes the right one to the left.

Outlier-robust image registration

Given two images:



Goal: find the transformation which takes the right one to the left.

Outlier-robust image registration: equations

- Reference image: \mathbf{y}_{ref} ,
- test image: \mathbf{y}_{test} ,
- possible transformations: Θ .

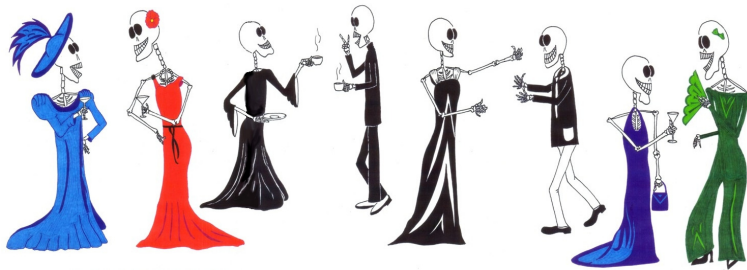
Objective:

$$J(\theta) = \underbrace{l(\mathbf{y}_{\text{ref}}, \mathbf{y}_{\text{test}}(\theta))}_{\text{similarity}} \rightarrow \max_{\theta \in \Theta}.$$

In the example: $l = \text{Non-linear CCA}$.

Cocktail party problem:

- independent groups of people / music bands,
- observation = mixed sources.



Observation:

$$\mathbf{x}_t = \mathbf{A}\mathbf{s}_t, \quad \mathbf{s} = \begin{bmatrix} \mathbf{s}^1; \dots; \mathbf{s}^M \end{bmatrix}.$$

Goal: $\hat{\mathbf{s}}$ from $\{\mathbf{x}_1, \dots, \mathbf{x}_T\}$. Assumptions:

- independent groups: $l(\mathbf{s}^1, \dots, \mathbf{s}^M) = 0$,
- \mathbf{s}^m -s: non-Gaussian,
- \mathbf{A} : invertible.

Find \mathbf{W} which makes the estimated components independent:

$$\mathbf{y} = \mathbf{W}\mathbf{x} = \begin{bmatrix} \mathbf{y}^1; \dots; \mathbf{y}^M \end{bmatrix},$$
$$J(\mathbf{W}) = I(\mathbf{y}^1, \dots, \mathbf{y}^M) \rightarrow \min_{\mathbf{W}}.$$

Recall: feature selection

- **Goal:** find
 - the feature subset ($\#$ of rooms, criminal rate, local taxes)
 - most relevant for house price prediction (y).



Here we consider a non-linear alternative of Lasso .

Feature selection: equations

- Features: x^1, \dots, x^F . Subset: $S \subseteq \{1, \dots, F\}$.
- MaxRelevance - MinRedundancy principle:

$$J(S) = \frac{1}{|S|} \sum_{i \in S} I(x^i, y) - \frac{1}{|S|^2} \sum_{i, j \in S} I(x^i, x^j) \rightarrow \max_{S \subseteq \{1, \dots, F\}} .$$

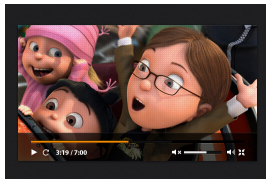
Example: independence testing-1

- We are given **paired samples**. Task: test **independence**.
- Examples:
 - (song, year of release) pairs



Example: independence testing-1

- We are given **paired samples**. Task: test **independence**.
- Examples:
 - (song, year of release) pairs
 - (video, caption) pairs

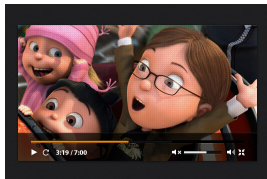


Example: independence testing-1

- We are given **paired samples**. Task: test **independence**.
- Examples:
 - (song, year of release) pairs



- (video, caption) pairs



- $\{(x_i, y_i)\}_{i=1}^n \xrightarrow{?} \mathbb{P}_{xy} = \mathbb{P}_x \mathbb{P}_y.$

Example: independence testing-2

- How do we detect dependency? (**paired** samples)

x_1 : Honourable senators, I have a question for the Leader of the Government in the Senate with regard to the support funding to farmers that has been announced. Most farmers have not received any money yet.

x_2 : No doubt there is great pressure on provincial and municipal governments in relation to the issue of child care, but the reality is that there have been no cuts to child care funding from the federal government to the provinces. In fact, we have increased federal investments for early childhood development.

...

y_1 : Honorables sénateurs, ma question s'adresse au leader du gouvernement au Sénat et concerne l'aide financière qu'on a annoncée pour les agriculteurs. La plupart des agriculteurs n'ont encore rien reçu de cet argent.

y_2 : Il est évident que les ordres de gouvernements provinciaux et municipaux subissent de fortes pressions en ce qui concerne les services de garde, mais le gouvernement n'a pas réduit le financement qu'il verse aux provinces pour les services de garde. Au contraire, nous avons augmenté le financement fédéral pour le développement des jeunes enfants.

...

Example: independence testing-2

- How do we detect dependency? (paired samples)

x_1 : Honourable senators, I have a question for the Leader of the Government in the Senate with regard to the support funding to farmers that has been announced. Most farmers have not received any money yet.

x_2 : No doubt there is great pressure on provincial and municipal governments in relation to the issue of child care, but the reality is that there have been no cuts to child care funding from the federal government to the provinces. In fact, we have increased federal investments for early childhood development.

...

y_1 : Honorables sénateurs, ma question s'adresse au leader du gouvernement au Sénat et concerne l'aide financière qu'on a annoncée pour les agriculteurs. La plupart des agriculteurs n'ont encore rien reçu de cet argent.

y_2 : Il est évident que les ordres de gouvernements provinciaux et municipaux subissent de fortes pressions en ce qui concerne les services de garde, mais le gouvernement n'a pas réduit le financement qu'il verse aux provinces pour les services de garde. Au contraire, nous avons augmenté le financement fédéral pour le développement des jeunes enfants.

...

Are the French paragraphs translations of the English ones, or have nothing to do with it, i.e. $\mathbb{P}_{xy} = \mathbb{P}_x \mathbb{P}_y$?

Towards non-linear CCA – History

- Given: random variable $(x, y) \in \mathcal{X} \times \mathcal{Y}$, $(x, y) \sim \mathbb{P}_{xy}$.
- **Goal**: measure the dependence of x and y .



Towards non-linear CCA – History

- Given: random variable $(x, y) \in \mathcal{X} \times \mathcal{Y}$, $(x, y) \sim \mathbb{P}_{xy}$.
- **Goal:** measure the dependence of x and y .
- **Desiderata** for a $Q(\mathbb{P}_{xy})$ independence measure:
 1. $Q(\mathbb{P}_{xy})$ is well-defined,
 2. $Q(\mathbb{P}_{xy}) \in [0, 1]$,
 3. $Q(\mathbb{P}_{xy}) = 0$ iff. $x \perp y$.
 4. $Q(\mathbb{P}_{xy}) = 1$ iff. $y = f(x)$ or $x = g(y)$.



- $Q(\mathbb{P}_{xy}) = \sup_{f,g} \text{corr}(f(x), g(y))$ satisfies 1-4.

- $Q(\mathbb{P}_{xy}) = \sup_{f,g} \text{corr}(f(x), g(y))$ satisfies 1-4.
- Too ambitious:
 - computationally intractable.
 - many functions.

Independence measures: restriction to continuous functions

- $C_b(\mathcal{X}) = \{f : \mathcal{X} \rightarrow \mathbb{R}, \text{ bounded continuous}\}$ would also work.
- Still too large!

Independence measures: restriction to continuous functions

- $C_b(\mathcal{X}) = \{f : \mathcal{X} \rightarrow \mathbb{R}, \text{ bounded continuous}\}$ would also work.
- Still too large!
- Idea:
 - certain \mathcal{H}_k function classes are dense in $C_b(\mathcal{X})$.
 - computationally tractable.

KCCA: definition

- Given: $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$.
- Associated:
 - feature maps $\varphi(x) = k(\cdot, x)$, $\psi(y) = \ell(\cdot, y)$,
 - RKHS-s \mathcal{H}_k , \mathcal{H}_ℓ .

- Given: $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$.
- Associated:
 - feature maps $\varphi(x) = k(\cdot, x)$, $\psi(y) = \ell(\cdot, y)$,
 - RKHS-s \mathcal{H}_k , \mathcal{H}_ℓ .
- KCCA measure of $(x, y) \in \mathcal{X} \times \mathcal{Y}$

$$\rho_{\text{KCCA}}(x, y; \mathcal{H}_k, \mathcal{H}_\ell) = \sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \text{corr}(f(x), g(y)),$$
$$\text{corr}(f(x), g(y)) = \frac{\text{cov}_{xy}(f(x), g(y))}{\sqrt{\text{var}_x f(x) \text{var}_y g(y)}}.$$

- Optimization domain: $\mathcal{H}_k \times \mathcal{H}_\ell \ni (f, g)$.
- By **reproducing property**: we will get a **finite-D task**.
- k, ℓ linear: traditional CCA.

- Optimization domain: $\mathcal{H}_k \times \mathcal{H}_\ell \ni (f, g)$.
- By reproducing property: we will get a finite-D task.
- k, ℓ linear: traditional CCA.
- In practice:
 - we have $\{(x_n, y_n)\}_{n=1}^N$ samples from (x, y) ,

- Optimization domain: $\mathcal{H}_k \times \mathcal{H}_\ell \ni (f, g)$.
- By **reproducing property**: we will get a **finite-D task**.
- k, ℓ linear: traditional CCA.
- In **practice**:
 - we have $\{(x_n, y_n)\}_{n=1}^N$ **samples** from (x, y) ,
 - it is worth applying **regularization**

$$\hat{\rho}_{\text{KCCA}}(x, y; \mathcal{H}_k, \mathcal{H}_\ell, \kappa) = \sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \widehat{\text{corr}}(f(x), g(y); \kappa),$$

$$\widehat{\text{corr}}(f(x), g(y); \kappa) = \frac{\widehat{\text{cov}}_{xy}(f(x), g(y))}{\sqrt{\widehat{\text{var}}_x f(x) + \kappa \|f\|_{\mathcal{H}_k}^2} \sqrt{\widehat{\text{var}}_y g(y) + \kappa \|g\|_{\mathcal{H}_\ell}^2}}.$$

KCCA solution: one-page summary

- Representer theorem $\Rightarrow f = \sum_{i=1}^N c_i \tilde{\varphi}(x_i)$, $g = \sum_{i=1}^N d_i \tilde{\psi}(y_i)$.

KCCA solution: one-page summary

- Representer theorem $\Rightarrow \mathbf{f} = \sum_{i=1}^N \mathbf{c}_i \tilde{\varphi}(x_i)$, $\mathbf{g} = \sum_{i=1}^N \mathbf{d}_i \tilde{\psi}(y_i)$.
- Objective in terms of \mathbf{c} and \mathbf{d} :

$$\widehat{\rho_{\text{KCCA}}}(x, y) := \sup_{\mathbf{c} \in \mathbb{R}^N, \mathbf{d} \in \mathbb{R}^N} \frac{\mathbf{c}^T \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \mathbf{d}}{\sqrt{\mathbf{c}^T (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 \mathbf{c}} \sqrt{\mathbf{d}^T (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \mathbf{d}}}.$$

KCCA solution: one-page summary

- Representer theorem $\Rightarrow \mathbf{f} = \sum_{i=1}^N \mathbf{c}_i \tilde{\varphi}(x_i)$, $\mathbf{g} = \sum_{i=1}^N \mathbf{d}_i \tilde{\psi}(y_i)$.
- Objective in terms of \mathbf{c} and \mathbf{d} :

$$\widehat{\rho_{\text{KCCA}}}(x, y) := \sup_{\mathbf{c} \in \mathbb{R}^N, \mathbf{d} \in \mathbb{R}^N} \frac{\mathbf{c}^T \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \mathbf{d}}{\sqrt{\mathbf{c}^T (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 \mathbf{c}} \sqrt{\mathbf{d}^T (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \mathbf{d}}}.$$

- Stationary points of $\widehat{\rho_{\text{KCCA}}}(x, y)$:

$$\mathbf{0} = \frac{\partial \widehat{\rho_{\text{KCCA}}}(x, y)}{\partial \mathbf{c}}, \quad \mathbf{0} = \frac{\partial \widehat{\rho_{\text{KCCA}}}(x, y)}{\partial \mathbf{d}}.$$

KCCA solution: one-page summary

- Representer theorem $\Rightarrow f = \sum_{i=1}^N c_i \tilde{\varphi}(x_i)$, $g = \sum_{i=1}^N d_i \tilde{\psi}(y_i)$.
- Objective in terms of \mathbf{c} and \mathbf{d} :

$$\widehat{\rho_{\text{KCCA}}}(x, y) := \sup_{\mathbf{c} \in \mathbb{R}^N, \mathbf{d} \in \mathbb{R}^N} \frac{\mathbf{c}^T \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \mathbf{d}}{\sqrt{\mathbf{c}^T (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 \mathbf{c}} \sqrt{\mathbf{d}^T (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \mathbf{d}}}.$$

- Stationary points of $\widehat{\rho_{\text{KCCA}}}(x, y)$:

$$\mathbf{0} = \frac{\partial \widehat{\rho_{\text{KCCA}}}(x, y)}{\partial \mathbf{c}}, \quad \mathbf{0} = \frac{\partial \widehat{\rho_{\text{KCCA}}}(x, y)}{\partial \mathbf{d}}.$$

- We just need the maximal eigenvalues ($\mathbf{A}\mathbf{z} = \lambda \mathbf{B}\mathbf{z}$) of

$$\begin{bmatrix} (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 & \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \\ \tilde{\mathbf{G}}_y \tilde{\mathbf{G}}_x & (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix} = \lambda \begin{bmatrix} (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 & \mathbf{0} \\ \mathbf{0} & (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix}.$$

2-variables $[(x, y)]$:

$$\begin{bmatrix} (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 & \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \\ \tilde{\mathbf{G}}_y \tilde{\mathbf{G}}_x & (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix} = \lambda \begin{bmatrix} (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 & \mathbf{0} \\ \mathbf{0} & (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix}$$

KCCA: M -variables

2-variables $[(x, y)]:$

$$\begin{bmatrix} (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 & \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \\ \tilde{\mathbf{G}}_y \tilde{\mathbf{G}}_x & (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix} = \lambda \begin{bmatrix} (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 & \mathbf{0} \\ \mathbf{0} & (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix}$$

For M -variables (pairwise dependence):

$$\begin{bmatrix} (\tilde{\mathbf{G}}_1 + \kappa \mathbf{I}_N)^2 & \tilde{\mathbf{G}}_1 \tilde{\mathbf{G}}_2 & \dots & \tilde{\mathbf{G}}_1 \tilde{\mathbf{G}}_M \\ \tilde{\mathbf{G}}_2 \tilde{\mathbf{G}}_1 & (\tilde{\mathbf{G}}_2 + \kappa \mathbf{I}_N)^2 & \dots & \tilde{\mathbf{G}}_2 \tilde{\mathbf{G}}_M \\ \vdots & \vdots & \ddots & \vdots \\ \tilde{\mathbf{G}}_M \tilde{\mathbf{G}}_1 & \tilde{\mathbf{G}}_M \tilde{\mathbf{G}}_2 & \dots & (\tilde{\mathbf{G}}_M + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c}_1 \\ \mathbf{c}_2 \\ \vdots \\ \mathbf{c}_M \end{bmatrix} = \lambda \begin{bmatrix} (\tilde{\mathbf{G}}_1 + \kappa \mathbf{I}_N)^2 & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & (\tilde{\mathbf{G}}_2 + \kappa \mathbf{I}_N)^2 & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & (\tilde{\mathbf{G}}_M + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c}_1 \\ \mathbf{c}_2 \\ \vdots \\ \mathbf{c}_M \end{bmatrix}.$$

If $x \perp y$, then $\rho_{\text{KCCA}}(x, y; \mathcal{H}_k, \mathcal{H}_\ell, \kappa) = 0$. Opposite direction:

- For 'rich' $\mathcal{H}_k, \mathcal{H}_\ell$.

KCCA as an independence measure

If $x \perp y$, then $\rho_{\text{KCCA}}(x, y; \mathcal{H}_k, \mathcal{H}_\ell, \kappa) = 0$. Opposite direction:

- For 'rich' $\mathcal{H}_k, \mathcal{H}_\ell$.
- Enough: **universal kernel**.

If $x \perp y$, then $\rho_{\text{KCCA}}(x, y; \mathcal{H}_k, \mathcal{H}_\ell, \kappa) = 0$. Opposite direction:

- For 'rich' $\mathcal{H}_k, \mathcal{H}_\ell$.
- Enough: **universal kernel**.
- **Example** ($\gamma > 0$):
 - Gaussian: $k(x, x') = e^{-\gamma \|x - x'\|_2^2}$.
 - Laplacian kernel: $k(x, x') = e^{-\gamma \|x - x'\|_2}$.

Definition

Assume:

- \mathcal{X} : compact metric space.
- k : continuous kernel on \mathcal{X} .

k is called **universal** if \mathcal{H}_k is dense in $(C_b(\mathcal{X}), \|\cdot\|_\infty)$.

Properties of universal kernels

If k is universal, then

- $k(x, x) > 0$ for all $x \in \mathcal{X}$.

Properties of universal kernels

If k is universal, then

- $k(x, x) > 0$ for all $x \in \mathcal{X}$.
- Every restriction of k to an $\mathcal{X}' \subseteq \mathcal{X}$ compact set is universal.

Properties of universal kernels

If k is universal, then

- $k(x, x) > 0$ for all $x \in \mathcal{X}$.
- Every restriction of k to an $\mathcal{X}' \subseteq \mathcal{X}$ compact set is universal.
- $\varphi(x) = k(\cdot, x)$ is injective, i.e.

$$\rho_k(x, y) = \|\varphi(x) - \varphi(y)\|_{\mathcal{H}_k}$$

is a metric.

Properties of universal kernels

If k is universal, then

- $k(x, x) > 0$ for all $x \in \mathcal{X}$.
- Every restriction of k to an $\mathcal{X}' \subseteq \mathcal{X}$ compact set is universal.
- $\varphi(x) = k(\cdot, x)$ is injective, i.e.

$$\rho_k(x, y) = \|\varphi(x) - \varphi(y)\|_{\mathcal{H}_k}$$

is a metric.

- The normalized kernel

$$\tilde{k}(x, y) := \frac{k(x, y)}{\sqrt{k(x, x)k(y, y)}}$$

is universal.

- For an $C^\infty \ni f : (-r, r) \rightarrow \mathbb{R}$

$$f(t) = \sum_{n=0}^{\infty} a_n t^n \quad t \in (-r, r), r \in (0, \infty].$$

- For an $C^\infty \ni f : (-r, r) \rightarrow \mathbb{R}$

$$f(t) = \sum_{n=0}^{\infty} a_n t^n \quad t \in (-r, r), r \in (0, \infty].$$

- If $a_n > 0 \ \forall n$, then

$$k(x, y) = f(\langle x, y \rangle)$$

is universal on $\mathcal{X} := \{x \in \mathbb{R}^d : \|x\|_2 \leq \sqrt{r}\}$.

- $k(\mathbf{x}, \mathbf{y}) = e^{\alpha \langle \mathbf{x}, \mathbf{y} \rangle}$: previous result with $a_n = \frac{\alpha^n}{n!}$.

Universal kernels, $\alpha > 0$

- $k(\mathbf{x}, \mathbf{y}) = e^{\alpha \langle \mathbf{x}, \mathbf{y} \rangle}$: previous result with $a_n = \frac{\alpha^n}{n!}$.
- $k(\mathbf{x}, \mathbf{y}) = e^{-\alpha \|\mathbf{x} - \mathbf{y}\|_2^2}$: exp. kernel & normalization.

Universal kernels, $\alpha > 0$

- $k(\mathbf{x}, \mathbf{y}) = (1 - \langle \mathbf{x}, \mathbf{y} \rangle)^{-\alpha}$ binomial kernel
 - on \mathcal{X} compact $\subset \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 < 1\}$.
 - $f(t) = (1 - t)^{-\alpha} = \sum_{n=0}^{\infty} \underbrace{\binom{-\alpha}{n}}_{>0} (-1)^n t^n \quad (|t| < 1),$

where $\binom{b}{n} = \sum_{i=1}^n \frac{b-i+1}{i}$.

Artifacts of too much free time

<https://bitbucket.org/szzoli/ite-in-python/>

Artifacts of too much free time

<https://bitbucket.org/szzoli/ite-in-python/>

Import ITE, generate observations:

```
>>> import ite
>>> from numpy.random import randn
>>> from numpy import array
>>> ds = array([2, 3, 4])
>>> t = 1000
>>> y = randn(t, sum(ds))
```

Estimate KCCA:

```
>>> co = ite.cost.BIKCCA()  
>>> kcca = co.estimate(y, ds)
```


Estimate KCCA:

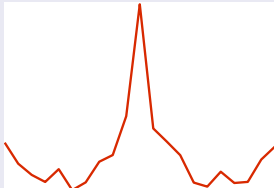
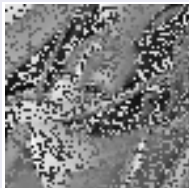
```
>>> co = ite.cost.BIKCCA()  
>>> kcca = co.estimated(y, ds)
```

Alternative initialization:

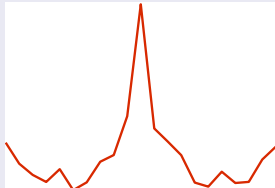
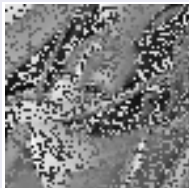
```
>>> co2 = ite.cost.BIKCCA(eta=1e-4, kappa=0.02)  
>>> kcca2 = co2.estimated(y, ds)
```

where η : low-rank approximation, κ : regularization constant.

Recall: outlier-robust image registration (it was KCCA)



Recall: outlier-robust image registration (it was KCCA)



Can solving eigenvalue problems be avoided? Analytical solution?

CCA Alternative: HSIC

HSIC: intuition. \mathcal{X} : images, \mathcal{Y} : descriptions



Their noses guide them through life, and they're never happier than when following an interesting scent. They need plenty of exercise, about an hour a day if possible.



A large animal who slings slobber, exudes a distinctive houndy odor, and wants nothing more than to follow his nose. They need a significant amount of exercise and mental stimulation.



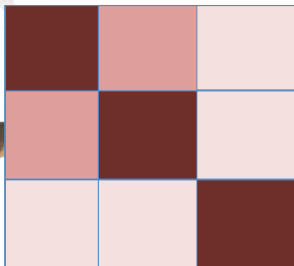
Known for their curiosity, intelligence, and excellent communication skills, the Javanese breed is perfect if you want a responsive, interactive pet, one that will blow in your ear and follow you everywhere.

Text from dogtime.com and petfinder.com

HSIC intuition: Gram matrices

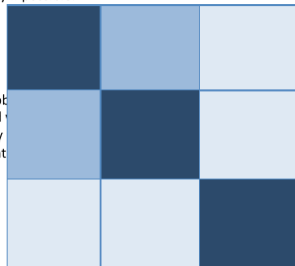


$\tilde{\mathbf{G}}_x$



Their noses guide them through life, and they're never happier than when following an interesting scent. They need plenty of exercise, about an hour a day if possible.

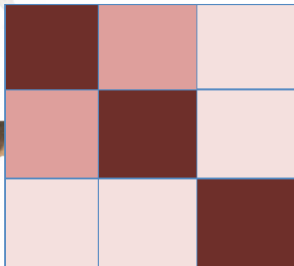
$\tilde{\mathbf{G}}_y$



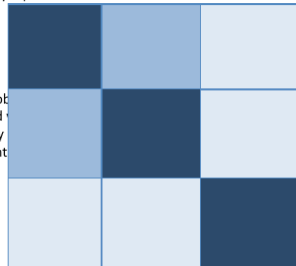
A large animal who slings slobbery drool, has a distinctive houndy odor, and is more interested in following his nose. They need a lot of exercise and mental stimulation.

Known for their curiosity, intelligence, and excellent communication skills, the Javanese breed is perfect if you want a responsive, interactive pet, one that will blow in your ear and follow you everywhere.

HSIC intuition: Gram matrices

 $\tilde{\mathbf{G}}_x$ 

Their noses guide them through life, and they're never happier than when following an interesting scent. They need plenty of exercise, about an hour a day if possible.

 $\tilde{\mathbf{G}}_y$ 

A large animal who slings slobbery, distinctive houndy odor, and who is more interested in sniffing than to follow his nose. They need a lot of exercise and mental stimulation.

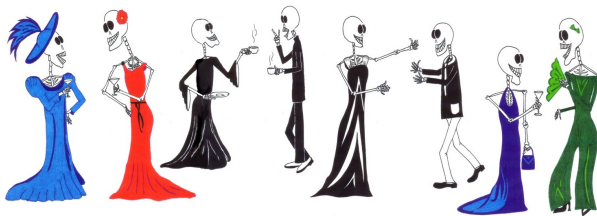
Known for their curiosity, intelligence, and excellent communication skills, the Javanese breed is perfect if you want a responsive, interactive pet, one that will blow in your ear and follow you everywhere.

Empirical estimate[†]:

$$\widehat{\text{HSIC}}^2 = \frac{1}{n^2} \left\langle \tilde{\mathbf{G}}_x, \tilde{\mathbf{G}}_y \right\rangle_F. \quad \leftarrow \text{analytical!}$$

[†] Visual illustration credit: Arthur Gretton

Cocktail party: HSIC demo



$$\mathbf{x} = \mathbf{A}\mathbf{s}, \quad \mathbf{s} = \begin{bmatrix} \mathbf{s}^1; \dots; \mathbf{s}^M \end{bmatrix},$$

where \mathbf{s}^m -s are non-Gaussian & independent.

- Goal: $\{\mathbf{x}_t\}_{t=1}^T \rightarrow \mathbf{W} = \mathbf{A}^{-1}, \{\mathbf{s}_t\}_{t=1}^T,$

$$\mathbf{x} = \mathbf{A}\mathbf{s}, \quad \mathbf{s} = \begin{bmatrix} \mathbf{s}^1; \dots; \mathbf{s}^M \end{bmatrix},$$

where \mathbf{s}^m -s are non-Gaussian & independent.

- Goal: $\{\mathbf{x}_t\}_{t=1}^T \rightarrow \mathbf{W} = \mathbf{A}^{-1}, \{\mathbf{s}_t\}_{t=1}^T,$
- Objective function:

$$\hat{\mathbf{s}} = \mathbf{W}\mathbf{x},$$
$$J(\mathbf{W}) = I(\hat{\mathbf{s}}^1, \dots, \hat{\mathbf{s}}^M) \rightarrow \min_{\mathbf{W}}.$$

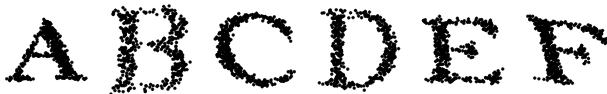
- Hidden sources (s):



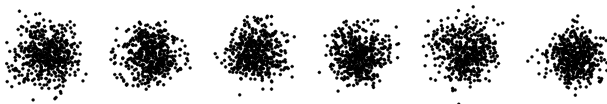
A B C D E F

ISA: source, observation

- Hidden sources (\mathbf{s}):



- Observation (\mathbf{x}):



ISA: estimated sources using HSIC, ambiguity

- Estimated sources (\hat{s}):

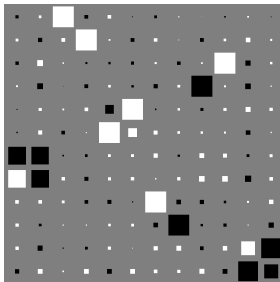
BROADV

ISA: estimated sources using HSIC, ambiguity

- Estimated sources ($\hat{\mathbf{s}}$):



- Performance ($\hat{\mathbf{W}}\mathbf{A}$), ambiguity:

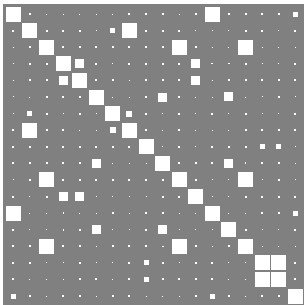


Conjecture: ISA separation theorem

- $\text{ISA} = \text{ICA} + \text{permutation}$.

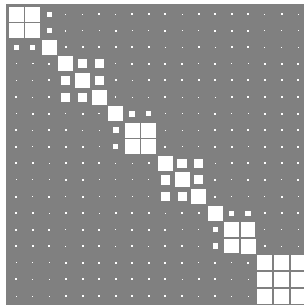
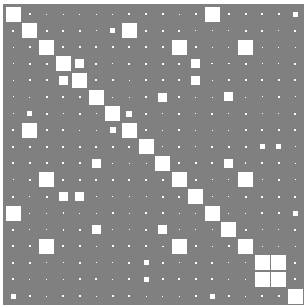
Conjecture: ISA separation theorem

- ISA = ICA + permutation. $\widehat{\text{HSIC}}(\hat{s}_i, \hat{s}_j)$. Here: $\dim(\mathbf{s}^m) = 3$.



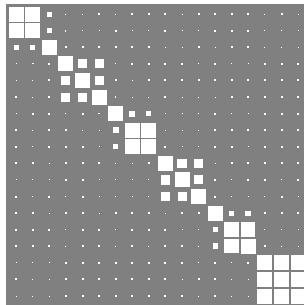
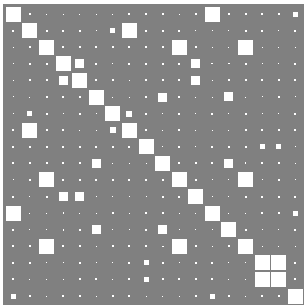
Conjecture: ISA separation theorem

- ISA = ICA + permutation. $\widehat{\text{HSIC}}(\hat{s}_i, \hat{s}_j)$. Here: $\dim(\mathbf{s}^m) = 3$.



Conjecture: ISA separation theorem

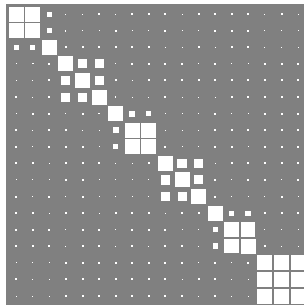
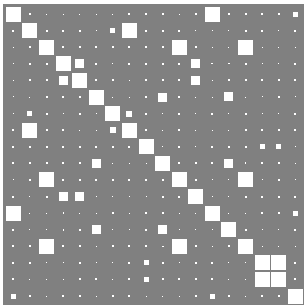
- ISA = ICA + permutation. $\widehat{\text{HSIC}}(\hat{s}_i, \hat{s}_j)$. Here: $\dim(\mathbf{s}^m) = 3$.



- Basis of the state-of-the-art ISA solvers.

Conjecture: ISA separation theorem

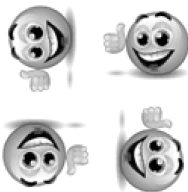
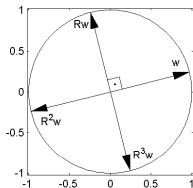
- ISA = ICA + permutation. $\widehat{\text{HSIC}}(\hat{s}_i, \hat{s}_j)$. Here: $\dim(\mathbf{s}^m) = 3$.



- Basis of the state-of-the-art ISA solvers.
- Sufficient conditions:
 - \mathbf{s}^m : spherical.

ISA separation theorem

For $\dim(\mathbf{s}^m) = 2$: less is sufficient.

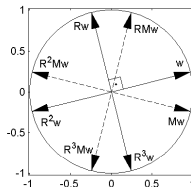
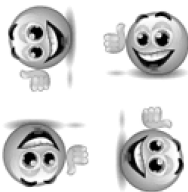
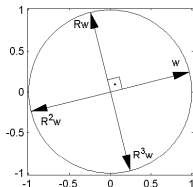


Invariance to

- 90° rotation: $f(u_1, u_2) = f(-u_2, u_1) = f(-u_1, -u_2) = f(u_2, -u_1)$.

ISA separation theorem

For $\dim(\mathbf{s}^m) = 2$: less is sufficient.

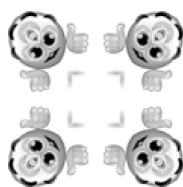
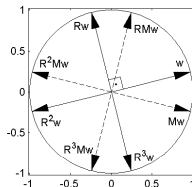
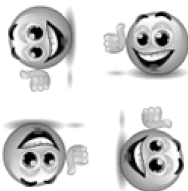
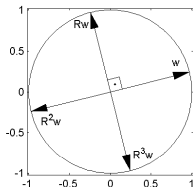


Invariance to

- 90° rotation: $f(u_1, u_2) = f(-u_2, u_1) = f(-u_1, -u_2) = f(u_2, -u_1)$.
- permutation and sign: $f(\pm u_1, \pm u_2) = f(\pm u_2, \pm u_1)$.

ISA separation theorem

For $\dim(\mathbf{s}^m) = 2$: less is sufficient.



Invariance to

- 90° rotation: $f(u_1, u_2) = f(-u_2, u_1) = f(-u_1, -u_2) = f(u_2, -u_1)$.
- permutation and sign: $f(\pm u_1, \pm u_2) = f(\pm u_2, \pm u_1)$.
- L^p -spherical: $f(u_1, u_2) = h(\sum_i |u_i|^p)$ ($p > 0$).

Idea: $\mathbb{P}_{xy} \mapsto C_{xy}$.

- Covariance matrix

$$C_{xy} = \mathbb{E}_{xy} \left[(x - \mathbb{E}x) (y - \mathbb{E}y)^T \right]$$

Idea: $\mathbb{P}_{xy} \mapsto C_{xy}$.

- Covariance matrix

$$C_{xy} = \mathbb{E}_{xy} \left[(x - \mathbb{E}x) (y - \mathbb{E}y)^T \right],$$

$$S = \|C_{xy}\|_F$$

Idea: $\mathbb{P}_{xy} \mapsto C_{xy}$.

- Covariance matrix

$$C_{xy} = \mathbb{E}_{xy} \left[(x - \mathbb{E}x) (y - \mathbb{E}y)^T \right],$$

$$S = \|C_{xy}\|_F \stackrel{?}{=} 0$$

Idea: $\mathbb{P}_{xy} \mapsto C_{xy}$.

- Covariance matrix

$$C_{xy} = \mathbb{E}_{xy} \left[(x - \mathbb{E}x) (y - \mathbb{E}y)^T \right],$$

$$S = \|C_{xy}\|_F \stackrel{?}{=} 0 \leftrightarrow \text{linear dependence.}$$

Idea: $\mathbb{P}_{xy} \mapsto C_{xy}$.

- Covariance matrix

$$C_{xy} = \mathbb{E}_{xy} \left[(x - \mathbb{E}x) (y - \mathbb{E}y)^T \right],$$

$$S = \|C_{xy}\|_F \stackrel{?}{=} 0 \leftrightarrow \text{linear dependence.}$$

- Covariance operator: take features of x and y

$$C_{xy} = \mathbb{E}_{xy} \left[\underbrace{(\varphi(x) - \mathbb{E}_x \varphi(x))}_{\text{centering in feature space}} \otimes (\psi(y) - \mathbb{E}_y \psi(y)) \right]$$

Idea: $\mathbb{P}_{xy} \mapsto C_{xy}$.

- Covariance matrix

$$C_{xy} = \mathbb{E}_{xy} \left[(x - \mathbb{E}x) (y - \mathbb{E}y)^T \right],$$
$$S = \|C_{xy}\|_F \stackrel{?}{=} 0 \leftrightarrow \text{linear dependence.}$$

- Covariance operator: take features of x and y

$$C_{xy} = \mathbb{E}_{xy} \left[\underbrace{(\varphi(x) - \mathbb{E}_x \varphi(x))}_{\text{centering in feature space}} \otimes (\psi(y) - \mathbb{E}_y \psi(y)) \right],$$
$$S = \|C_{xy}\|_{HS} =: \text{HSIC}(\mathbb{P}_{xy}).$$

We capture **non-linear dependencies** via φ, ψ !

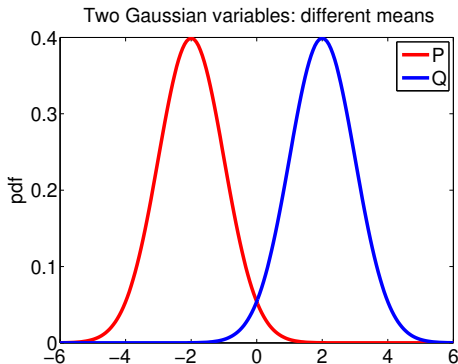
- Independence: $\mathbb{P}_{xy} = \mathbb{P}_x \otimes \mathbb{P}_y$.

Questions

- How do we check this equality?
- How can distributions be represented?

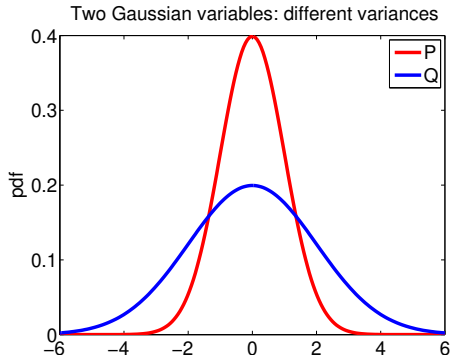
Representations of distributions: EX

- Given: 2 Gaussians with different means.
- Solution: *t*-test.



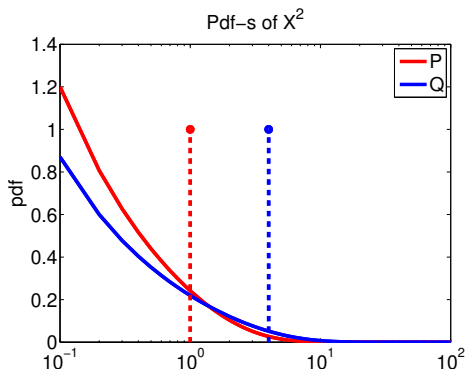
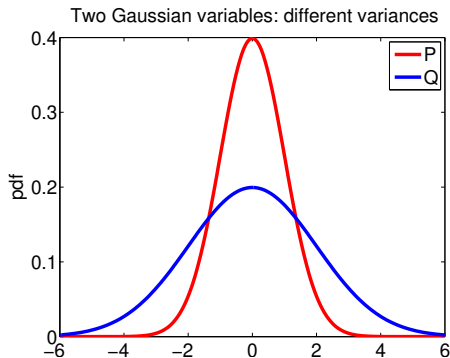
Representations of distributions: $\mathbb{E}X^2$

- Setup: 2 Gaussians; same means, different variances.
- Idea: look at the 2nd-order features of RVs.



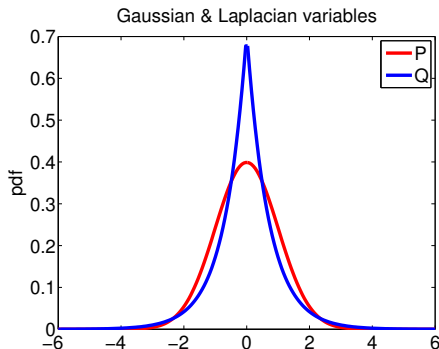
Representations of distributions: $\mathbb{E}X^2$

- Setup: 2 Gaussians; same means, different variances.
- Idea: look at the 2nd-order features of RVs.
- $\varphi(x) = x^2 \Rightarrow$ difference in $\mathbb{E}X^2$.



Representations of distributions: further moments

- Setup: a Gaussian and a Laplacian distribution.
- Challenge: their means *and* variances are the same.
- Idea: look at **higher-order features**.



$$\varphi(\mathbf{x}) = e^{i\langle \cdot, \mathbf{x} \rangle}: \text{characteristic function, } \mathcal{X} = \mathbb{R}^d.$$

Distribution representation via functions

$$\mathbb{P} \mapsto \mu_{\mathbb{P}} = \int_{\mathcal{X}} \varphi(x) d\mathbb{P}(x).$$

Distribution representation via functions

$$\mathbb{P} \mapsto \mu_{\mathbb{P}} = \int_{\mathcal{X}} \varphi(x) d\mathbb{P}(x).$$

- Cdf:

$$\mathbb{P} \mapsto F_{\mathbb{P}}(z) = \mathbb{E}_{x \sim \mathbb{P}} \chi_{(-\infty, z)}(x).$$

Distribution representation via functions

$$\mathbb{P} \mapsto \mu_{\mathbb{P}} = \int_{\mathcal{X}} \varphi(x) d\mathbb{P}(x).$$

- Cdf:

$$\mathbb{P} \mapsto F_{\mathbb{P}}(z) = \mathbb{E}_{x \sim \mathbb{P}} \chi_{(-\infty, z)}(x).$$

- Characteristic function:

$$\mathbb{P} \mapsto c_{\mathbb{P}}(z) = \int e^{i\langle z, x \rangle} d\mathbb{P}(x).$$

Distribution representation via functions

$$\mathbb{P} \mapsto \mu_{\mathbb{P}} = \int_{\mathcal{X}} \varphi(x) d\mathbb{P}(x).$$

- Cdf:

$$\mathbb{P} \mapsto F_{\mathbb{P}}(z) = \mathbb{E}_{x \sim \mathbb{P}} \chi_{(-\infty, z]}(x).$$

- Characteristic function:

$$\mathbb{P} \mapsto c_{\mathbb{P}}(z) = \int e^{i\langle z, x \rangle} d\mathbb{P}(x).$$

- Moment generating function:

$$\mathbb{P} \mapsto M_{\mathbb{P}}(z) = \int e^{\langle z, x \rangle} d\mathbb{P}(x).$$

Distribution representation via functions

$$\mathbb{P} \mapsto \mu_{\mathbb{P}} = \int_{\mathcal{X}} \varphi(x) d\mathbb{P}(x).$$

- Cdf:

$$\mathbb{P} \mapsto F_{\mathbb{P}}(z) = \mathbb{E}_{x \sim \mathbb{P}} \chi_{(-\infty, z]}(x).$$

- Characteristic function:

$$\mathbb{P} \mapsto c_{\mathbb{P}}(z) = \int e^{i\langle z, x \rangle} d\mathbb{P}(x).$$

- Moment generating function:

$$\mathbb{P} \mapsto M_{\mathbb{P}}(z) = \int e^{\langle z, x \rangle} d\mathbb{P}(x).$$

Trick

φ : on any kernel-endowed domain! $\varphi(x) := k(\cdot, x)$, $\mu_{\mathbb{P}} \in \mathcal{H}_k$.

- Mean embedding:

$$\mu_{\mathbb{P}} := \int_{\mathcal{X}} \varphi(\mathbf{x}) \, \mathrm{d}\mathbb{P}(\mathbf{x})$$

We got

- Mean embedding:

$$\mu_{\mathbb{P}} := \int_{\mathcal{X}} \underbrace{\varphi(x)}_{k(\cdot, x)} d\mathbb{P}(x) \in \mathcal{H}_k.$$

- Maximum mean discrepancy:

$$\text{MMD}_k(\mathbb{P}, \mathbb{Q}) := \|\mu_k(\mathbb{P}) - \mu_k(\mathbb{Q})\|_{\mathcal{H}_k}.$$

Recall: $\langle \mu_k(\hat{\mathbb{P}}), \mu_k(\hat{\mathbb{Q}}) \rangle_{\mathcal{H}_k}$



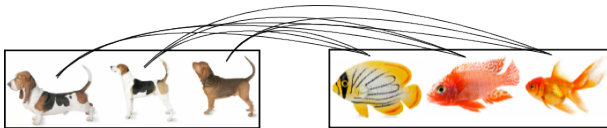
- Mean embedding:

$$\mu_{\mathbb{P}} := \int_{\mathcal{X}} \underbrace{\varphi(x)}_{k(\cdot, x)} d\mathbb{P}(x) \in \mathcal{H}_k.$$

- Maximum mean discrepancy:

$$\text{MMD}_k(\mathbb{P}, \mathbb{Q}) := \|\mu_k(\mathbb{P}) - \mu_k(\mathbb{Q})\|_{\mathcal{H}_k}.$$

Recall: $\langle \mu_k(\hat{\mathbb{P}}), \mu_k(\hat{\mathbb{Q}}) \rangle_{\mathcal{H}_k}$



- Hilbert-Schmidt independence criterion, $k = \bigotimes_{m=1}^M k_m$:

$$\text{HSIC}_k(\mathbb{P}) := \text{MMD}_k\left(\mathbb{P}, \bigotimes_{m=1}^M \mathbb{P}_m\right).$$

MMD with $k = \bigotimes_{m=1}^M k_m$:

$$k(x, x') := \prod_{m=1}^M k_m(x_m, x'_m),$$

$$\text{HSIC}_k(\mathbb{P}) = \text{MMD}_k\left(\mathbb{P}, \bigotimes_{m=1}^M \mathbb{P}_m\right).$$

MMD with $k = \bigotimes_{m=1}^M k_m$:

$$k(x, x') := \prod_{m=1}^M k_m(x_m, x'_m),$$

$$\text{HSIC}_k(\mathbb{P}) = \text{MMD}_k\left(\mathbb{P}, \bigotimes_{m=1}^M \mathbb{P}_m\right).$$

Applications :

- blind source separation,
- feature selection, post selection inference,
- independence testing, causal inference.

MMD with $k = \bigotimes_{m=1}^M k_m$:

$$k(x, x') := \prod_{m=1}^M k_m(x_m, x'_m),$$

$$\text{HSIC}_k(\mathbb{P}) = \text{MMD}_k\left(\mathbb{P}, \bigotimes_{m=1}^M \mathbb{P}_m\right).$$

Applications :

- blind source separation,
- feature selection, post selection inference,
- independence testing, causal inference.

The 2 views are equivalent; we estimated HSIC empirically.

Applications:

- two-sample testing,
- domain adaptation, -generalization,
- kernel Bayesian inference,
- approximate Bayesian computation, probabilistic programming,
- model criticism, goodness-of-fit,
- distribution classification, distribution regression,
- topological data analysis.

When is

- $\text{MMD}_k(\mathbb{P}, \mathbb{Q}) := \|\mu_k(\mathbb{P}) - \mu_k(\mathbb{Q})\|_{\mathcal{H}_k}$ a metric? In this case k is called characteristic.
- $\text{HSIC}_k(\mathbb{P})$ an independence measure?

When is

- $\text{MMD}_k(\mathbb{P}, \mathbb{Q}) := \|\mu_k(\mathbb{P}) - \mu_k(\mathbb{Q})\|_{\mathcal{H}_k}$ a metric? In this case k is called characteristic.
- $\text{HSIC}_k(\mathbb{P})$ an independence measure?

MMD: for continuous, bounded, shift-invariant k

- By the Bochner's theorem:

$$k(\mathbf{x}, \mathbf{x}') = k_0(\mathbf{x} - \mathbf{x}') = \int_{\mathbb{R}^d} e^{i\langle \mathbf{x} - \mathbf{x}', \boldsymbol{\omega} \rangle} d\Lambda(\boldsymbol{\omega}).$$

When is

- $\text{MMD}_k(\mathbb{P}, \mathbb{Q}) := \|\mu_k(\mathbb{P}) - \mu_k(\mathbb{Q})\|_{\mathcal{H}_k}$ a metric? In this case k is called characteristic.
- $\text{HSIC}_k(\mathbb{P})$ an independence measure?

MMD: for continuous, bounded, shift-invariant k

- By the Bochner's theorem:

$$k(\mathbf{x}, \mathbf{x}') = k_0(\mathbf{x} - \mathbf{x}') = \int_{\mathbb{R}^d} e^{i\langle \mathbf{x} - \mathbf{x}', \boldsymbol{\omega} \rangle} d\Lambda(\boldsymbol{\omega}).$$

- \Rightarrow MMD in terms of characteristic functions:

$$\text{MMD}_k^2(\mathbb{P}, \mathbb{Q}) = \|c_{\mathbb{P}} - c_{\mathbb{Q}}\|_{L^2(\Lambda)}^2.$$

Theorem

k is characteristic iff. $\text{supp}(\Lambda) = \mathbb{R}^d$.

Theorem

k is characteristic iff. $\text{supp}(\Lambda) = \mathbb{R}^d$.

Example on \mathbb{R} :

kernel name	k_0	$\hat{k}_0(\omega)$	$\text{supp}(\hat{k}_0)$
Gaussian	$e^{-\frac{x^2}{2\sigma^2}}$	$\sigma e^{-\frac{\sigma^2\omega^2}{2}}$	\mathbb{R}
Laplacian	$e^{-\sigma x }$	$\sqrt{\frac{2}{\pi}} \frac{\sigma}{\sigma^2 + \omega^2}$	\mathbb{R}
Sinc	$\frac{\sin(\sigma x)}{x}$	$\sqrt{\frac{\pi}{2}} \chi_{[-\sigma, \sigma]}(\omega)$	$[-\sigma, \sigma]$

Theorem

k is characteristic iff. $\text{supp}(\Lambda) = \mathbb{R}^d$.

Example on \mathbb{R} :

kernel name	k_0	$\hat{k}_0(\omega)$	$\text{supp}(\hat{k}_0)$
Gaussian	$e^{-\frac{x^2}{2\sigma^2}}$	$\sigma e^{-\frac{\sigma^2\omega^2}{2}}$	\mathbb{R}
Laplacian	$e^{-\sigma x }$	$\sqrt{\frac{2}{\pi}} \frac{\sigma}{\sigma^2 + \omega^2}$	\mathbb{R}
Sinc	$\frac{\sin(\sigma x)}{x}$	$\sqrt{\frac{\pi}{2}} \chi_{[-\sigma, \sigma]}(\omega)$	$[-\sigma, \sigma]$

Note:

- universality \Rightarrow characteristic.
- $k = \otimes_m k_m$: characteristic \Rightarrow HSIC: \checkmark . How about in terms of k_m -s?

Description when HSIC is 'valid'

Proposition (characteristic property)

- $\otimes_{m=1}^M k_m$: *characteristic* $\Rightarrow (k_m)_{m=1}^M$ *are characteristic*.
- $\Leftrightarrow [|\mathcal{X}_m| = 2, k_m(x, x') = 2\delta_{x, x'} - 1]$

Description when HSIC is 'valid'

Proposition (characteristic property)

- $\bigotimes_{m=1}^M k_m$: *characteristic* $\Rightarrow (k_m)_{m=1}^M$ are *characteristic*.
- $\Leftrightarrow [|\mathcal{X}_m| = 2, k_m(x, x') = 2\delta_{x, x'} - 1]$

Proposition (\mathcal{I} -characteristic property)

- k_1, k_2 : *characteristic* $\Rightarrow k_1 \otimes k_2$: \mathcal{I} -*characteristic*.
- \Leftarrow : for $\forall M \geq 2$.
- k_1, k_2, k_3 : *characteristic* $\Rightarrow \bigotimes_{m=1}^3 k_m$: \mathcal{I} -*characteristic* [Ex].

Description when HSIC is 'valid'

Proposition (characteristic property)

- $\bigotimes_{m=1}^M k_m$: characteristic $\Rightarrow (k_m)_{m=1}^M$ are characteristic.
- $\Leftrightarrow [|\mathcal{X}_m| = 2, k_m(x, x') = 2\delta_{x, x'} - 1]$

Proposition (\mathcal{I} -characteristic property)

- k_1, k_2 : characteristic $\Rightarrow k_1 \otimes k_2$: \mathcal{I} -characteristic.
- \Leftarrow : for $\forall M \geq 2$.
- k_1, k_2, k_3 : characteristic $\Rightarrow \bigotimes_{m=1}^3 k_m$: \mathcal{I} -characteristic [Ex].

Proposition ($\mathcal{X}_m = \mathbb{R}^{d_m}$, k_m : continuous, shift-invariant, bounded)

$(k_m)_{m=1}^M$ -s are characteristic $\Leftrightarrow \bigotimes_{m=1}^M k_m$: \mathcal{I} -characteristic \Leftrightarrow
 $\bigotimes_{m=1}^M k_m$: characteristic.

Description when HSIC is 'valid'

Proposition (characteristic property)

- $\bigotimes_{m=1}^M k_m$: *characteristic* $\Rightarrow (k_m)_{m=1}^M$ are *characteristic*.
- $\Leftrightarrow [|\mathcal{X}_m| = 2, k_m(x, x') = 2\delta_{x, x'} - 1]$

Proposition (\mathcal{I} -characteristic property)

- k_1, k_2 : *characteristic* $\Rightarrow k_1 \otimes k_2$: \mathcal{I} -*characteristic*.
- \Leftarrow : for $\forall M \geq 2$.
- k_1, k_2, k_3 : *characteristic* $\Rightarrow \bigotimes_{m=1}^3 k_m$: \mathcal{I} -*characteristic* [Ex].

Proposition ($\mathcal{X}_m = \mathbb{R}^{d_m}$, k_m : continuous, shift-invariant, bounded)

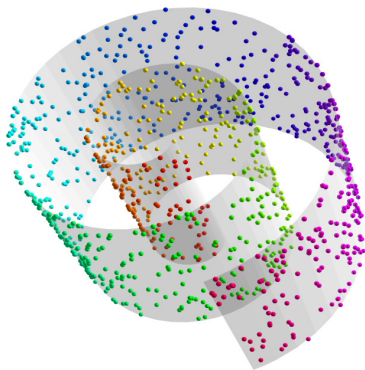
$(k_m)_{m=1}^M$ -s are *characteristic* $\Leftrightarrow \bigotimes_{m=1}^M k_m$: \mathcal{I} -*characteristic* \Leftrightarrow
 $\bigotimes_{m=1}^M k_m$: *characteristic*.

Proposition (universality)

$\bigotimes_{m=1}^M k_m$: *universal* $\Leftrightarrow (k_m)_{m=1}^M$ are *universal*.

Other dimensionality reduction techniques

Other non-linear methods



Goal: $\{\mathbf{x}_i\}_{i=1}^n \subset \mathbb{R}^D \xrightarrow{?} \{\mathbf{x}'_i\}_{i=1}^n \subset \mathbb{R}^d$, retaining the geometry of $\{\mathbf{x}_i\}_{i=1}^n$.

Multidimensional scaling (MDS)

- Given: $\mathbf{D} = [d_{ij}]_{i,j=1}^n$ distance matrix, $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$.

Multidimensional scaling (MDS)

- Given: $\mathbf{D} = [d_{ij}]_{i,j=1}^n$ distance matrix, $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$.
- Objective function:

$$\min_{\mathbf{x}'} \sum_{i,j} \underbrace{\left(d_{ij}^2 - \|\mathbf{x}'_i - \mathbf{x}'_j\|_2^2 \right)}_{\text{preserve (large) distances}}, \text{ s.t. } \mathbf{x}'_i = \mathbf{W}\mathbf{x}_i, \|\mathbf{w}_i\|_2^2 = 1, \forall i.$$

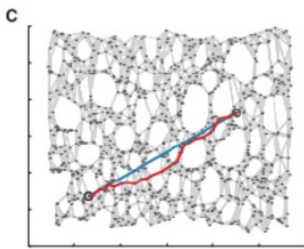
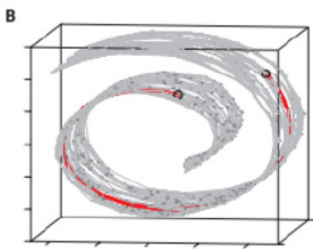
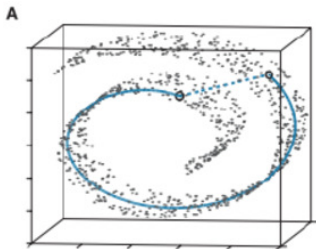
Multidimensional scaling (MDS)

- Given: $\mathbf{D} = [d_{ij}]_{i,j=1}^n$ distance matrix, $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$.
- Objective function:

$$\min_{\mathbf{X}'} \sum_{i,j} \underbrace{\left(d_{ij}^2 - \|\mathbf{x}'_i - \mathbf{x}'_j\|_2^2 \right)}_{\text{preserve (large) distances}}, \text{ s.t. } \mathbf{x}'_i = \mathbf{W}\mathbf{x}_i, \|\mathbf{w}_i\|_2^2 = 1, \forall i.$$

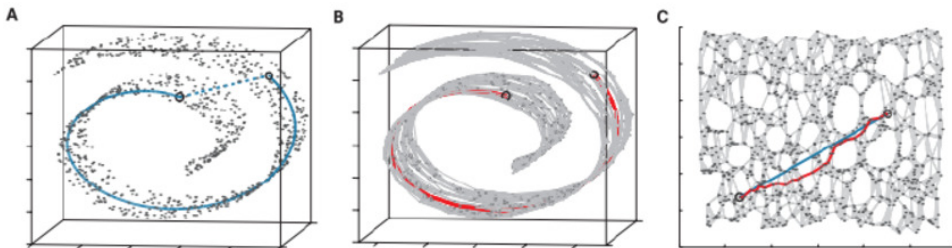
- Solution: $\mathbf{G} = \mathbf{X}^T \mathbf{X} = [\langle \mathbf{x}_i, \mathbf{x}_j \rangle]_{i,j=1}^n$ Gram matrix.
 - Top d eigenvalues, eigenvectors of \mathbf{G} : λ_i, \mathbf{v}_i ($i = 1, \dots, d$).
 - $\mathbf{x}'_i = \sqrt{\lambda_i} \mathbf{v}_i$.

ISOMAP \Leftarrow MDS



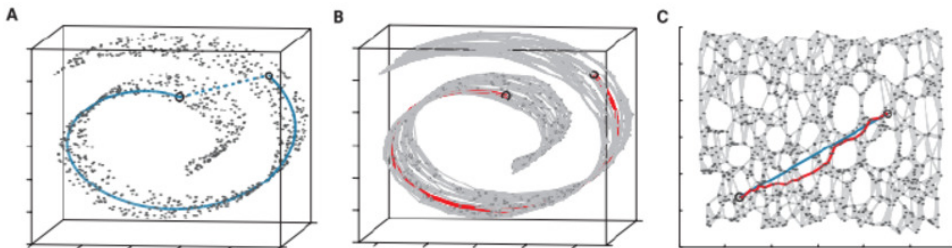
- Idea: For curved manifold let us rely on neighborhoods.

ISOMAP \Leftarrow MDS



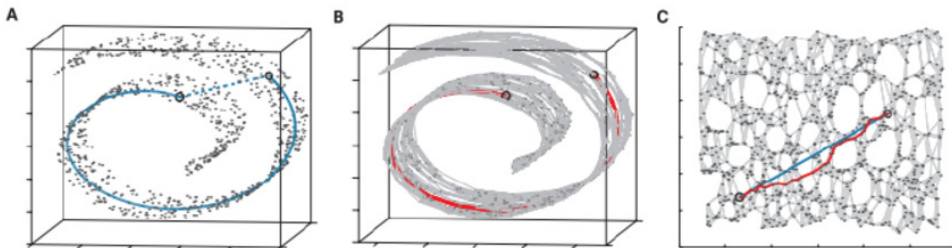
- Idea: For curved manifold let us rely on **neighborhoods**.
- Steps:
 - 1 $\hat{d}_{\text{geodesic}}(\mathbf{x}_i, \mathbf{x}_j)$ = shortest path of \mathbf{x}_i and \mathbf{x}_j on kNN graph.
(Dijkstra/Floyd's alg.)

ISOMAP \Leftarrow MDS



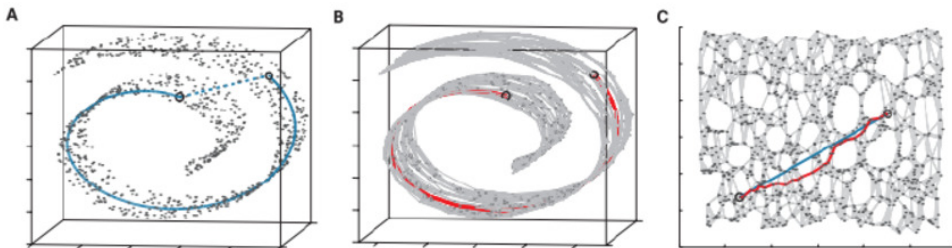
- Idea: For curved manifold let us rely on neighborhoods.
- Steps:
 - 1 $\hat{d}_{\text{geodesic}}(\mathbf{x}_i, \mathbf{x}_j)$ = shortest path of \mathbf{x}_i and \mathbf{x}_j on kNN graph.
(Dijkstra/Floyd's alg.)
 - 2 $\mathbf{D} := \left[\hat{d}_{\text{geodesic}}(\mathbf{x}_i, \mathbf{x}_j) \right]$.

ISOMAP \Leftarrow MDS



- Idea: For curved manifold let us rely on **neighborhoods**.
- Steps:
 - 1 $\hat{d}_{\text{geodesic}}(\mathbf{x}_i, \mathbf{x}_j)$ = shortest path of \mathbf{x}_i and \mathbf{x}_j on kNN graph.
(Dijkstra/Floyd's alg.)
 - 2 $\mathbf{D} := \left[\hat{d}_{\text{geodesic}}(\mathbf{x}_i, \mathbf{x}_j) \right]$.
 - 3 Call **MDS** on \mathbf{D} .

ISOMAP \Leftarrow MDS



- Idea: For curved manifold let us rely on **neighborhoods**.
- Steps:
 - 1 $\hat{d}_{\text{geodesic}}(\mathbf{x}_i, \mathbf{x}_j)$ = shortest path of \mathbf{x}_i and \mathbf{x}_j on kNN graph. (Dijkstra/Floyd's alg.)
 - 2 $\mathbf{D} := \left[\hat{d}_{\text{geodesic}}(\mathbf{x}_i, \mathbf{x}_j) \right]$.
 - 3 Call **MDS** on \mathbf{D} .
- It can be **slow**.

Sammon mapping = MDS & local distance preservation

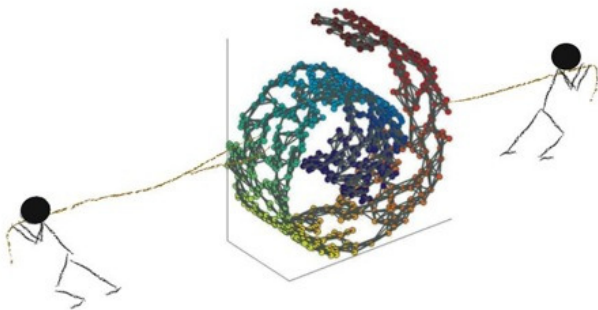
- Recall (MDS):

$$\min_{\mathbf{x}'} \sum_{i,j} \underbrace{\left(d_{ij}^2 - \|\mathbf{x}'_i - \mathbf{x}'_j\|_2^2 \right)}_{\text{preserve (large) distances}}, \text{ s.t. } \mathbf{x}'_i = \mathbf{W}\mathbf{x}_i, \|\mathbf{w}_i\|_2^2 = 1, \forall i.$$

- MDS cares mostly about **large** distances.
- Sammon mapping: weights := $\frac{1}{d_{ij}}$.

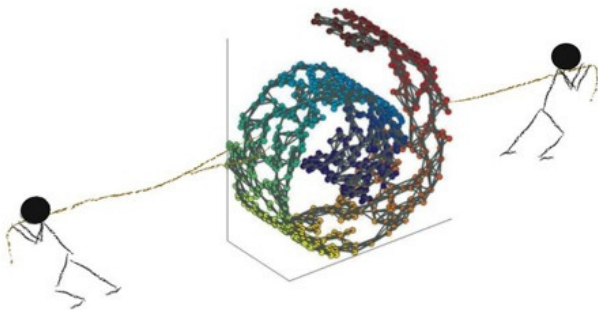
$$\min_{\mathbf{x}'} \frac{1}{\sum_{i \neq j} d_{ij}} \sum_{i \neq j} \frac{\left(d_{ij} - \|\mathbf{x}'_i - \mathbf{x}'_j\|_2 \right)^2}{d_{ij}}.$$

MVU = MDS & explicit unfolding



$G := \text{kNN graph of } \{\mathbf{x}_i\}_{i=1}^n.$

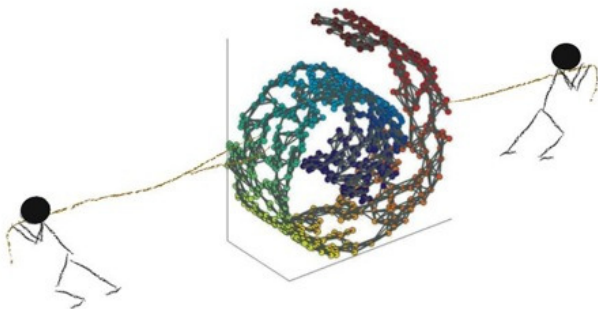
MVU = MDS & explicit unfolding



$G := \text{kNN graph of } \{\mathbf{x}_i\}_{i=1}^n$. Objective:

$$\max_{\mathbf{x}'} \sum_{ij} \|\mathbf{x}'_i - \mathbf{x}'_j\|_2^2 \quad \text{s.t.} \quad \|\mathbf{x}'_i - \mathbf{x}'_j\|_2^2 = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \quad (i, j) \in G$$

MVU = MDS & explicit unfolding



$G := \text{kNN graph of } \{\mathbf{x}_i\}_{i=1}^n$. Objective:

$$\max_{\mathbf{x}'} \sum_{ij} \|\mathbf{x}'_i - \mathbf{x}'_j\|_2^2 \quad \text{s.t.} \quad \|\mathbf{x}'_i - \mathbf{x}'_j\|_2^2 = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \quad (i, j) \in G$$

Leads to SDP: linear objective on positive-semidefinite matrices.

Locally linear embedding (LLE)

- Assumption: **local linearity**.
- Steps:
 - 1 $G := k\text{NN graph} \Rightarrow \mathbf{x}_{i_j} := j^{\text{th}} \text{ NN of } \mathbf{x}_i.$

Locally linear embedding (LLE)

- Assumption: **local linearity**.
- Steps:
 - $G := k\text{NN graph} \Rightarrow \mathbf{x}_{i_j} := j^{\text{th}} \text{ NN of } \mathbf{x}_i$.
 - $\mathbf{w}_j := \arg \min_{\mathbf{w}} \left\| \mathbf{x}_i - \sum_{j=1}^k w_{ij} \mathbf{x}_{i_j} \right\|_2$. Objective:

$$\min_{\mathbf{X}'} \sum_i \underbrace{\left\| \mathbf{x}'_i - \sum_j \mathbf{w}_{ij} \mathbf{x}'_{i_j} \right\|_2^2}_{\text{local linearity preserving}} \quad \text{s.t.} \quad \underbrace{\left\| \mathbf{x}'^{(k)} \right\|_2^2 = 1, \forall k}_{\text{to avoid } \mathbf{X}' = \mathbf{0}}.$$

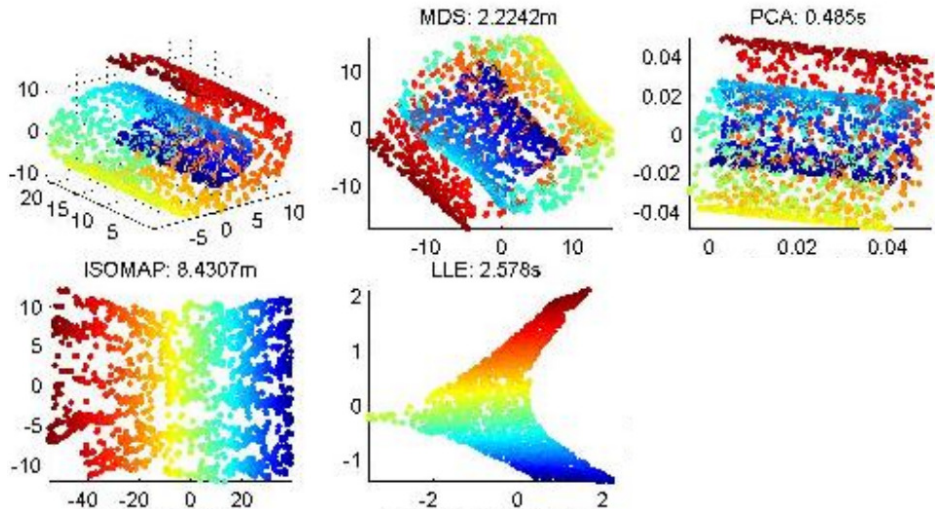
Locally linear embedding (LLE)

- Assumption: local linearity.
- Steps:
 - $G := k\text{NN graph} \Rightarrow \mathbf{x}_{ij} := j^{\text{th}} \text{ NN of } \mathbf{x}_i$.
 - $\mathbf{w}_j := \arg \min_{\mathbf{w}} \left\| \mathbf{x}_i - \sum_{j=1}^k w_{ij} \mathbf{x}_{ij} \right\|_2$. Objective:

$$\min_{\mathbf{X}'} \sum_i \underbrace{\left\| \mathbf{x}'_i - \sum_j \mathbf{w}_{ij} \mathbf{x}'_{ij} \right\|_2^2}_{\text{local linearity preserving}} \quad \text{s.t.} \quad \underbrace{\left\| \mathbf{x}'^{(k)} \right\|_2^2 = 1, \forall k}_{\text{to avoid } \mathbf{X}' = \mathbf{0}}.$$

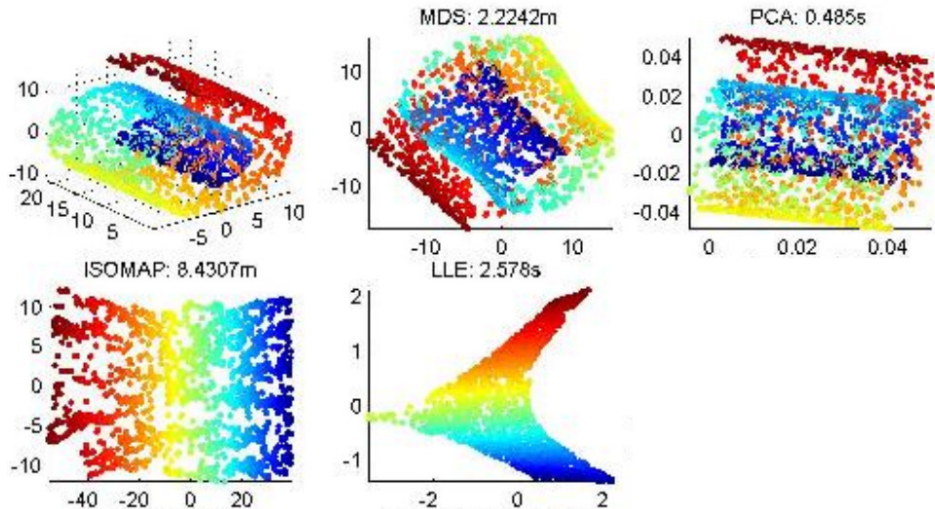
- Solution: from eigensystem of $(\mathbf{I} - \mathbf{W})^T (\mathbf{I} - \mathbf{W})$, $\mathbf{W} = 1 - \chi_G$.

Manifold embedding: demo[†]



[†]Todd Wittman

Manifold embedding: demo[†]



MDS, ISOMAP: slow. MDS, PCA: fail to unroll (no manifold info).

[†]Todd Wittman

Techniques:

- PCA, KPCA: maximum variance projection.
- CCA, KCCA: maximally dependent projection.
- HSIC:
 - analytical KCCA alternative,
 - norm of covariance operator.
- MDS: (large) distance retaining.
- ISOMAP: geodesic distance preserving.
- Sammon mapping: distance retaining (including small ones).
- MVU: kNN distance preserving & explicit unrolling.
- LLE: local linearity preserving.

Applications:

- image compression & registration,
- non-linear feature selection,
- media annotation, translation testing,
- cocktail party (ISA).

Thank you for the attention!



Why do we get eigenvalue problems?

- $\mathbf{A} \in \mathbb{R}^{n \times n}$: symmetric matrix.
- Objective:

$$\max_{\mathbf{V} \in \mathbb{R}^{n \times d}: \mathbf{V}^T \mathbf{V} = \mathbf{I}} \text{Tr}(\mathbf{V}^T \mathbf{A} \mathbf{V}).$$

Why do we get eigenvalue problems?

- $\mathbf{A} \in \mathbb{R}^{n \times n}$: symmetric matrix.
- Objective:

$$\max_{V \in \mathbb{R}^{n \times d}: V^T V = I} \text{Tr}(\mathbf{V}^T \mathbf{A} \mathbf{V}).$$

- Optimal solution:
 - $\mathbf{V}^* = d$ leading eigenvectors of \mathbf{A} .
 - uniqueness up to subspace.

Why do we get generalized eigenvalue problems?

- $\mathbf{A} \in \mathbb{R}^{n \times n}$: symmetric matrix. $\mathbf{B} \in \mathbb{R}^{n \times n}$: positive definite.
- Objective:

$$\max_{V \in \mathbb{R}^{n \times d}: \mathbf{V}^T \mathbf{B} \mathbf{V} = \mathbf{I}} \text{Tr} \left(\mathbf{V}^T \mathbf{A} \mathbf{V} \right).$$

Why do we get generalized eigenvalue problems?

- $\mathbf{A} \in \mathbb{R}^{n \times n}$: symmetric matrix. $\mathbf{B} \in \mathbb{R}^{n \times n}$: positive definite.
- Objective:

$$\max_{V \in \mathbb{R}^{n \times d}: \mathbf{V}^T \mathbf{B} \mathbf{V} = \mathbf{I}} \text{Tr} \left(\mathbf{V}^T \mathbf{A} \mathbf{V} \right).$$

- Solution: $\mathbf{V}^* = d$ leading (\mathbf{B} -orthogonal) eigenvectors of the generalized eigenvalue problem

$$\mathbf{A} \mathbf{x} = \lambda \mathbf{B} \mathbf{x}.$$