

Question

सफेदी की चमकर... ज्यादा सफेद

The whitening transform

This is given by $\tilde{x}_i = \Delta^{-1/2} U^T p_i$. Here, we deal with

$k \times 1$ Type I normalised pattern vectors (i.e., normalised with respect to the mean) p_i , $0 \leq i \leq n-1$.

Δ is a diagonal matrix of their eigenvalues, and U is their corresponding eigenvector matrix. What is the covariance matrix of the transformed patterns \tilde{x}_i ? (The following will not be evaluated: why do you think this transformation is called so?)

$$\tilde{x}_i = \Delta^{-1/2} U^T p_i$$

stacking all n of them together

$$\tilde{R} = \Delta^{-1/2} U^T P$$

$$\begin{aligned} \text{Cov} &\triangleq \frac{1}{n} \tilde{R} \tilde{R}^T = \frac{1}{n} \Delta^{-1/2} U^T P (\Delta^{-1/2} U^T P)^T \\ &= \left(\frac{1}{n} \right) \Delta^{-1/2} U^T P P^T U \Delta^{-1/2} \\ &\quad (\because U^T U = I, \Delta^{-1/2 T} = \Delta^{-1/2} : \text{diag}) \\ &= \Delta^{-1/2} U^T \underbrace{A U}_{U^T A U = \Lambda} \Delta^{-1/2} \\ &= \Delta^{-1/2} \Lambda \Delta^{-1/2} \\ &= I \quad (\text{Multiplication of diagonal} \\ &\quad \text{matrices} = \text{multiplication of} \\ &\quad \text{individual terms}) \end{aligned}$$

Why is this called so? No correlation, same variance along each direction \sim white light / white noise.

**Question: Orthopaedic Orthonormality:
Bone-Breaking Normal work**

Show that for non-repeated eigenvalues of a symmetric matrix, the eigenvectors will be orthonormal. Consider the case of repeated eigenvalues of a symmetric matrix A . Specifically use the Gram Schmidt orthogonalisation process to show that it is still possible to obtain a set of orthonormal vectors.

Let the symmetric matrix be $A \in \mathbb{R}^n \times \mathbb{R}^n$ with n eigenvalues λ_k and corresponding eigenvectors \underline{u}_k i.e., $A \underline{u}_k = \lambda_k \underline{u}_k$

Further, given an orthogonal set, it is always possible to convert it to an orthonormal one by dividing the vector by its norm e.g., Euclidean.

Non-repeated eigenvalues

$$\begin{aligned} \text{Consider } \lambda_i \underline{u}_i \cdot \underline{u}_j &= \lambda_i \underline{u}_i^T \underline{u}_j \\ &= (\lambda_i \underline{u}_i)^T \underline{u}_j \quad (\text{transpose of a scalar}) \\ &= (A \underline{u}_i)^T \underline{u}_j \\ &= \underline{u}_i^T A^T \underline{u}_j = \underline{u}_i^T A \underline{u}_j \quad (\text{transpose of a symmetric matrix}) \\ &= (\underline{u}_i^T) (A \underline{u}_j) = \underline{u}_i^T \lambda_j \underline{u}_j \\ &= \lambda_j \underline{u}_i^T \underline{u}_j \end{aligned}$$

$$\Rightarrow \lambda_i \underline{u}_i^T \underline{u}_j = \lambda_j \underline{u}_i^T \underline{u}_j \Rightarrow \underbrace{(\lambda_i - \lambda_j)}_{\neq 0 \text{ as non-repeated eigenvalues}} \underbrace{\underline{u}_i^T \underline{u}_j}_{\Rightarrow \underline{u}_i \perp \underline{u}_j} = 0$$

are orthogonal.

Repeated eigenvalues of degree k (say)

FIRST, [consider the Gram-Schmidt orthogonalisation construction]

Given a set of k basis vectors $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_k$

(Basis vectors \Rightarrow they are linearly independent)

To construct an orthogonal set $\underline{u}_1, \underline{u}_2, \dots, \underline{u}_k$ from it.

Step 0: Take $\underline{u}_1 = \underline{v}_1$

Step I: To construct $\underline{u}_2 \perp \underline{u}_1$ such that

the span of $\underline{u}_1, \underline{u}_2 = \text{span of } \underline{v}_1, \underline{v}_2$

Take $\underline{u}_2 = a_1 \underline{u}_1 + \underline{v}_2$ (linear combination)

To find a_1 : take a dot product with \underline{u}_1

$$\underbrace{\underline{u}_2 \cdot \underline{u}_1}_{=0} = a_1 (\underline{u}_1 \cdot \underline{u}_1) + \underline{v}_2 \cdot \underline{u}_1$$
$$\Rightarrow a_1 = \frac{-\underline{v}_2 \cdot \underline{u}_1}{\underline{u}_1 \cdot \underline{u}_1} = -\frac{\langle \underline{v}_2, \underline{u}_1 \rangle}{\langle \underline{u}_1, \underline{u}_1 \rangle}$$

$$\Rightarrow \underline{u}_2 = \underline{v}_2 - \frac{\langle \underline{v}_2, \underline{u}_1 \rangle}{\langle \underline{u}_1, \underline{u}_1 \rangle} \underline{u}_1$$

Step II: To construct $\underline{u}_3 \perp (\underline{u}_1, \underline{u}_2)$ such that

the span of $\underline{u}_1, \underline{u}_2, \underline{u}_3 = \text{span of } \underline{v}_1, \underline{v}_2, \underline{v}_3$

Take $\underline{u}_3 = a_1 \underline{u}_1 + a_2 \underline{u}_2 + \underline{v}_3$

To find a_1, a_2 : take dot products with $\underline{u}_1, \underline{u}_2$

$$\rightarrow \underbrace{\underline{u}_1 \cdot \underline{u}_3}_{=0} = a_1 \underline{u}_1 \cdot \underline{u}_1 + a_2 \underbrace{\underline{u}_2 \cdot \underline{u}_1}_{=0} + \underline{v}_3 \cdot \underline{u}_1$$
$$\Rightarrow a_1 = \frac{-\underline{v}_3 \cdot \underline{u}_1}{\underline{u}_1 \cdot \underline{u}_1} = -\frac{\langle \underline{v}_3, \underline{u}_1 \rangle}{\langle \underline{u}_1, \underline{u}_1 \rangle}$$

$$\rightarrow \underbrace{\underline{u}_2 \cdot \underline{u}_3}_{=0} = a_1 \underbrace{\underline{u}_1 \cdot \underline{u}_2}_{=0} + a_2 \underline{u}_2 \cdot \underline{u}_2 + \underline{v}_3 \cdot \underline{u}_2$$
$$\Rightarrow a_2 = \frac{-\underline{v}_3 \cdot \underline{u}_2}{\underline{u}_2 \cdot \underline{u}_2} = -\frac{\langle \underline{v}_3, \underline{u}_2 \rangle}{\langle \underline{u}_2, \underline{u}_2 \rangle}$$

$$\Rightarrow \underline{u}_3 = \underline{v}_3 - \frac{\langle \underline{v}_3, \underline{u}_1 \rangle}{\langle \underline{u}_1, \underline{u}_1 \rangle} \underline{u}_1 - \frac{\langle \underline{v}_3, \underline{u}_2 \rangle}{\langle \underline{u}_2, \underline{u}_2 \rangle} \underline{u}_2$$

Step (R):

$$\underline{u}_k = \underline{v}_k - \sum_{j=1}^k \frac{\langle \underline{v}_k, \underline{u}_j \rangle}{\langle \underline{u}_j, \underline{u}_j \rangle} \underline{u}_j$$

THE MOST IMPORTANT STEP: To show that

$\underline{u}_1, \underline{u}_2, \dots, \underline{u}_k$ are also eigenvectors of A , like $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_k$

Step (O): $\underline{u}_1 = \underline{v}_1 \Rightarrow \underline{u}_1$ is an eigenvector of A

Step (I) $\underline{u}_2 = \underline{v}_2 - \frac{\langle \underline{v}_2, \underline{u}_1 \rangle}{\langle \underline{u}_1, \underline{u}_1 \rangle} \underline{u}_1$

$$\begin{aligned} \Rightarrow A \underline{u}_2 &= A \underline{v}_2 - \frac{\langle \underline{v}_2, \underline{u}_1 \rangle}{\langle \underline{u}_1, \underline{u}_1 \rangle} A \underline{u}_1 \\ &= \lambda_2 \underline{v}_2 - \frac{\langle \underline{v}_2, \underline{u}_1 \rangle}{\langle \underline{u}_1, \underline{u}_1 \rangle} \lambda_1 \underline{u}_1 \end{aligned}$$

Now, $\lambda_1 = \lambda_2 = \dots = \lambda_k$ for all, $= \lambda$ (say)

$$= \lambda \left[\underline{v}_2 - \frac{\langle \underline{v}_2, \underline{u}_1 \rangle}{\langle \underline{u}_1, \underline{u}_1 \rangle} \underline{u}_1 \right]$$

$\Rightarrow A \underline{u}_2 = \lambda \underline{u}_2 \Rightarrow \underline{u}_2$ is also an eigenvector of A

Step (II) $\underline{u}_3 = \underline{v}_3 - \frac{\langle \underline{v}_3, \underline{u}_2 \rangle}{\langle \underline{u}_2, \underline{u}_2 \rangle} \underline{u}_2 - \frac{\langle \underline{v}_3, \underline{u}_1 \rangle}{\langle \underline{u}_1, \underline{u}_1 \rangle} \underline{u}_1$

$$\begin{aligned} \Rightarrow A \underline{u}_3 &= A \underline{v}_3 - \frac{\langle \underline{v}_3, \underline{u}_2 \rangle}{\langle \underline{u}_2, \underline{u}_2 \rangle} A \underline{u}_2 - \frac{\langle \underline{v}_3, \underline{u}_1 \rangle}{\langle \underline{u}_1, \underline{u}_1 \rangle} A \underline{u}_1 \\ &= \lambda_3 \underline{v}_3 - \frac{\langle \underline{v}_3, \underline{u}_2 \rangle}{\langle \underline{u}_2, \underline{u}_2 \rangle} \lambda_2 \underline{u}_2 - \frac{\langle \underline{v}_3, \underline{u}_1 \rangle}{\langle \underline{u}_1, \underline{u}_1 \rangle} \lambda_1 \underline{u}_1 \end{aligned}$$

Now, $\lambda_1 = \lambda_2 = \dots = \lambda_k$ for all, $= \lambda$ (say)

$$= \lambda \left[\underline{v}_3 - \frac{\langle \underline{v}_3, \underline{u}_2 \rangle}{\langle \underline{u}_2, \underline{u}_2 \rangle} \underline{u}_2 - \frac{\langle \underline{v}_3, \underline{u}_1 \rangle}{\langle \underline{u}_1, \underline{u}_1 \rangle} \underline{u}_1 \right]$$

$= \lambda \underline{u}_3 \Rightarrow \underline{u}_3$ is also an eigenvector of A

step (K) follows

Question

Projection Interjection: Nothing to lose?

What is the physical significance of projecting a new vector onto an eigenspace? Explain using mathematical expressions, what the above implies, for both the KL-Transform, as well as the SVD

Now, explain how the KL-Transform and the SVD respectively perform lossy compression.

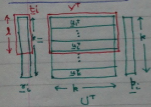
Projecting onto an eigenspace:

KL-Transform: Taking a dot product with the eigenvectors of the cov matrix $\frac{1}{n} P P^T$

SVD: Taking a dot product with the orthonormal basis vectors u_i , NOT the eigenvectors of $P P^T$ (see the figures below, for more detail on the notation)

Lossy Compression

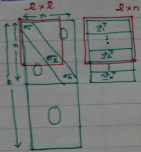
KL-Transform



$t_i = U^T p_i$

$t_i = V^T p_i$

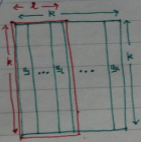
SVD



choosing l :

$\frac{\sum_{i=1}^l \lambda_i}{\sum_{i=1}^n \lambda_i} \geq 0.95$

min l such that this is valid



choosing l :

$\frac{\sum_{i=1}^l \sigma_i}{\sum_{i=1}^n \sigma_i} \geq 0.95$

min l such that this is valid.

2. Likelihood:

$$p(\underline{t} | \underline{X}, \underline{w}) = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1 - t_n}$$

where $y_n = y(\underline{x}_n) = \sigma(\underline{w}^T \underline{x}_n)$

If the data set is separable, any decision boundary will have the property:

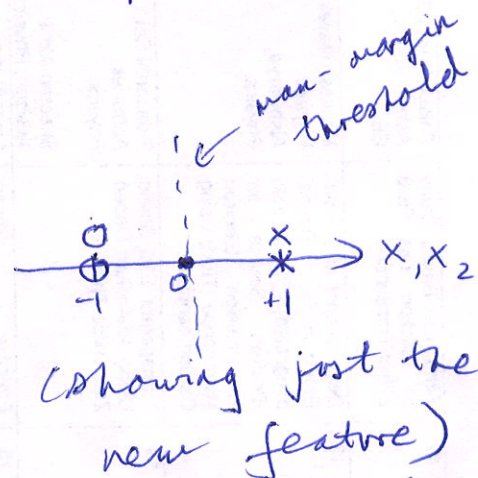
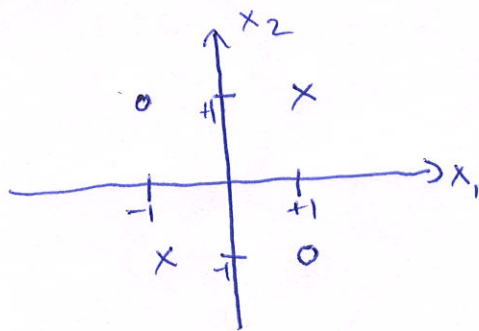
$$y(\underline{x}_n) > 0.5 \Rightarrow \underline{w}^T \underline{x}_n > 0 \text{ [if } t_n = 1\text{]}$$

$$y(\underline{x}_n) < 0.5 \Rightarrow \underline{w}^T \underline{x}_n < 0 \text{ [if } t_n = 0\text{]}$$

Clearly, the likelihood will be maximised when $y_n = t_n \forall n$. For y_n to be 1/0, $\underline{w}^T \underline{x}_n$ should go to $+\infty / -\infty$. Thus, once \underline{w} specifies a separating hyperplane as above, then taking $\|\underline{w}\| \rightarrow \infty$ will maximise the likelihood.

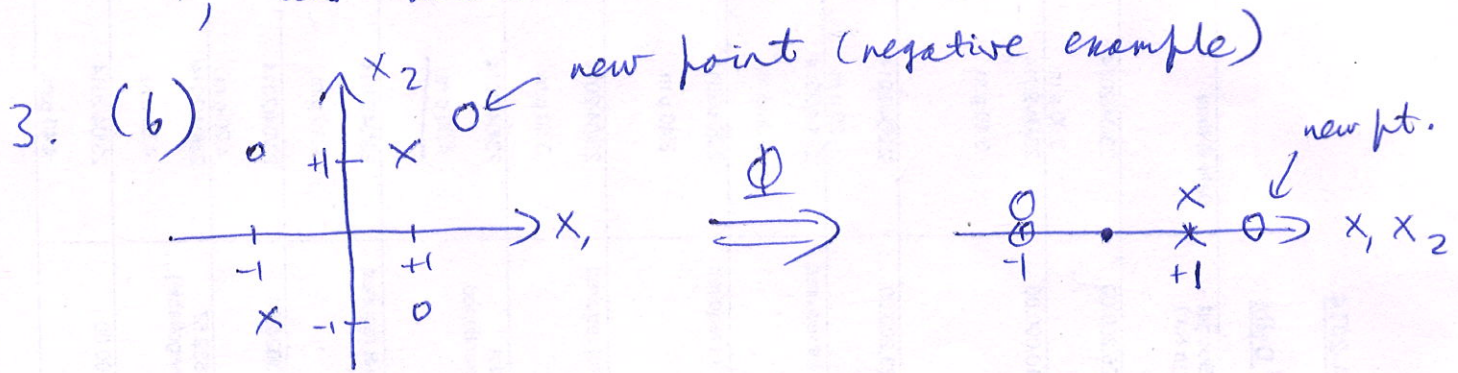
3. (a) No, not separable. (XOR problem)

Max-margin:



only the new feature is needed, with a threshold at 0; thus, the simplest $\underline{w} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ -1 \end{pmatrix}$.

Note that this w gives functional margin
 $= 1$, as discussed in class.



Clearly non-separable, even with the new feature.

(Positive examples lie in-between negative examples on all three dimensions.)

$$(c) \quad K(\underline{x}, \underline{x}') = \underline{\Phi}(\underline{x}) \cdot \underline{\Phi}(\underline{x}')$$

$$= \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ x_1 x_2 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ x_1' \\ x_2' \\ x_1' x_2' \end{pmatrix}$$

$$= 1 + x_1 x_1' + x_2 x_2' + x_1 x_1' x_2 x_2'$$

It is a dot product, and thus by definition a Mercer kernel.