

Discussion 5

EE599: Deep Learning

Keith M. Chugg

Spring 2020



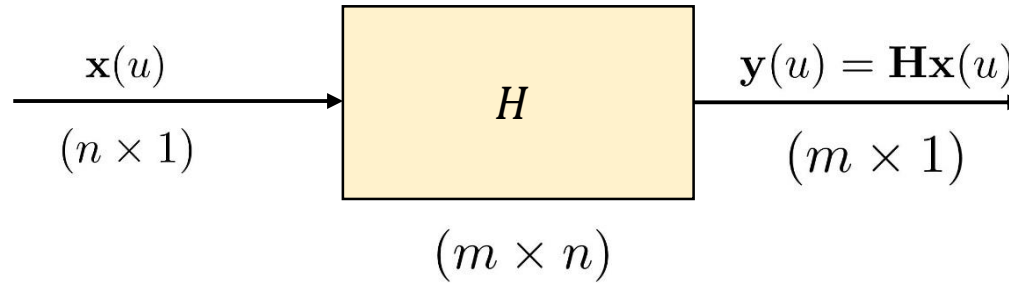
USC University of
Southern California

Outline



- Principal Component Analysis
- Filter Design

Random Vectors



$$\mathbf{m}_y = \mathbf{H}\mathbf{m}_x$$

$$\mathbf{R}_y = \mathbf{H}\mathbf{R}_x\mathbf{H}^T$$

$$\mathbf{K}_y = \mathbf{H}\mathbf{K}_x\mathbf{H}^T$$

Correlation

$$\begin{aligned}\mathbf{R}_y &= \mathbb{E}[\mathbf{y}(u)\mathbf{y}^T(u)] \\ &= \mathbb{E}[\mathbf{H}(\mathbf{x}(u))(\mathbf{H}\mathbf{x}(u))^T] \\ &= \mathbb{E}[\mathbf{H}\mathbf{x}(u)\mathbf{x}(u)^T\mathbf{H}^T] \\ &= \mathbf{H}\mathbb{E}[\mathbf{x}(u)\mathbf{x}(u)^T]\mathbf{H}^T \\ &= \mathbf{H}\mathbf{R}_x\mathbf{H}^T\end{aligned}$$

Special case

$$y(u) = \mathbf{b}^T \mathbf{x}(u) \quad (1 \times 1)$$

$$m_y = \mathbf{b}^T \mathbf{m}_x$$

$$\mathbb{E}[y(u)^2] = \mathbf{b}^T \mathbf{R}_x \mathbf{b}$$

$$\sigma_y^2 = \mathbf{b}^T \mathbf{K}_x \mathbf{b}$$

Note that covariance / correlation matrices are symmetric, non-negative definite.

Karhunen-Loève (KL) Expansion

Can always find orthonormal set of eigen-vectors of \mathbf{K} .

These are an alternate coordinate systems (rotations, reflections)

in this eigen-coordinate system, the components are uncorrelated

(principle components)

Eigen-values are the variant (energy) in each of these principle directions

(can be used to reduce dimensions by throwing out components with low energy)

KL-Expansion

$$\mathbf{K}_x \mathbf{e}_k = \lambda_k \mathbf{e}_k \quad k = 0, 1, \dots, N - 1$$

(Eigen equation)

$$\mathbf{e}_k^T \mathbf{e}_l = \delta[k - l] \quad \lambda_k \geq 0$$

(Orthonormal eigen-vectors)

$$\mathbf{x}(u) = \sum_{k=0}^{N-1} X_k(u) \mathbf{e}_k$$

(Change of coordinates)

$$X_k(u) = \mathbf{e}_k^T \mathbf{x}(u)$$

$$\mathbb{E}[X_k(u)X_l(u)] = \mathbf{e}_k^T \mathbf{K}_x \mathbf{e}_l = \lambda_k \delta(k - l)$$

(Uncorrelated components)

$$\mathbf{K}_x = \sum_{k=0}^{N-1} \lambda_k \mathbf{e}_k \mathbf{e}_k^T$$

(Mercer's Theorem)

$$\mathbb{E}[||\mathbf{x}(u)||^2] = \text{tr}(\mathbf{K}_x) = \sum_{k=0}^{N-1} \lambda_k$$

(Total Energy)

Always exists because \mathbf{K} is positive semi-definite (PSD) and symmetric

KL-Expansion

$$d_k(u) = \mathbf{e}_k^T \mathbf{x}(u) \quad k = 0, 1, \dots, D - 1$$

$$\mathbf{d}(u) = \mathbf{E}^T \mathbf{x}(u)$$

$$\begin{aligned} \mathbf{K}_d &= \mathbf{E}^T \mathbf{K}_x \mathbf{E} \\ &= \mathbf{E}^T (\mathbf{E} \mathbf{\Lambda} \mathbf{E}^T) \mathbf{E} \\ &= \mathbf{\Lambda} = \mathbf{diag}(\lambda_k) \end{aligned}$$

Multiplying by \mathbf{E}^T makes the component uncorrelated

$$\mathbf{E} = \left[\mathbf{e}_0 \mid \mathbf{e}_1 \mid \mathbf{e}_2 \mid \dots \mid \mathbf{e}_{D-1} \right]$$

KL-Expansion: Relation to Whitening

$$w_k(u) = \frac{X_k(u)}{\sqrt{\lambda_k}} = \frac{\mathbf{e}_k^T \mathbf{x}(u)}{\sqrt{\lambda_k}} \quad k = 0, 1, \dots, D - 1$$

$$\begin{aligned} \mathbf{w}(u) &= \mathbf{\Lambda}^{-1/2} \mathbf{E}^T \mathbf{x}(u) \\ \mathbf{K}_w &= \mathbf{\Lambda}^{-1/2} \mathbf{E}^T \mathbf{K}_x \mathbf{E} \mathbf{\Lambda}^{-1/2} \\ &= \mathbf{\Lambda}^{-1/2} \mathbf{\Lambda} \mathbf{\Lambda}^{-1/2} \\ &= \mathbf{I} \end{aligned}$$

For any orthogonal matrix \mathbf{U} , this whitening matrix also works:

$$\mathbf{G} = \mathbf{U} \mathbf{\Lambda}^{-1/2} \mathbf{E}^T$$

KL- Expansion: Relation to PCA

$$\tilde{x}_k(u) = \mathbf{e}_k^K \mathbf{x}(u)$$

$$k = 0, 1, \dots, K - 1$$

$$\tilde{\mathbf{x}}(u) = \mathbf{E}_{[:K]} \mathbf{x}(u)$$

First K components

$$\mathbf{K}_{\tilde{\mathbf{x}}} = \mathbf{\Lambda}_{[:K]}$$

Assumes ordered eigen-values: $\lambda_0, \lambda_1, \dots, \lambda_{D-1}$

$$\mathbb{E}[\|\tilde{\mathbf{x}}(u)\|^2] = \sum_{k=0}^{K-1} \lambda_k$$

$$\mathbb{E}[\|\mathbf{x}(u) - \tilde{\mathbf{x}}(u)\|^2] = \sum_{k=K}^{D-1} \lambda_k$$

minimizes approximation error (lossy compression)

PCA is simply taking only the K most important Eigen-directions or principal components

$$\mathbf{E}_{[:K]} = \left[\mathbf{e}_0 \mid \mathbf{e}_1 \mid \mathbf{e}_2 \mid \dots \mid \mathbf{e}_{K-1} \right]$$

KL- Expansion / PCA for Data

Everything is the same, except that we use the data averaging instead of the expectation $\mathbb{E}[\cdot]$

$$\begin{aligned}\hat{\mathbf{R}}_x &= \langle \mathbf{x}\mathbf{x}^T \rangle_{\mathcal{D}} \\ &= \frac{1}{N} \sum_{n=0}^{N-1} \mathbf{x}_n \mathbf{x}_n^T \\ &= \frac{1}{N} \mathbf{X}^T \mathbf{X}\end{aligned}$$

Both KL and PCA can be applied to **R** or **K**. Center x if you want to use **K**

$$\mathbf{x} \leftarrow \mathbf{x} - \mathbf{m}$$

(same if mean is zero)

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_0^T \\ \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_{N-1}^T \end{bmatrix}$$

$$\mathbf{X}^T = [\mathbf{x}_0 \quad \mathbf{x}_1 \dots \mathbf{x}_{N-1}]$$

“stacked” data matrix

KL- Expansion / PCA for Data

PCA for data

$$\tilde{\mathbf{x}}_n = \mathbf{E}_{[:K]}^T \mathbf{x}_n$$

First K components

$$\mathbf{E}_{[:K]} = \left[\mathbf{e}_0 \mid \mathbf{e}_1 \mid \mathbf{e}_2 \mid \dots \mid \mathbf{e}_{K-1} \right]$$

Apply to the “stacked” data matrix

$$\tilde{\mathbf{X}} = \begin{bmatrix} (\mathbf{E}_{[:K]}^T \mathbf{x}_0)^T \\ (\mathbf{E}_{[:K]}^T \mathbf{x}_1)^T \\ \vdots \\ (\mathbf{E}_{[:K]}^T \mathbf{x}_{N-1})^T \end{bmatrix} = \mathbf{X} \mathbf{E}_{[:K]}$$

$$\underset{N \times K}{\tilde{\mathbf{X}}} = \underset{N \times D}{\mathbf{X}} \underset{D \times K}{\mathbf{E}_{[:K]}}$$

Dimension reduced from D to K

$$\tilde{\mathbf{X}}^T = [\mathbf{E}_{[:K]}^T \mathbf{x}_0 \quad \mathbf{E}_{[:K]}^T \mathbf{x}_1 \quad \dots \quad \mathbf{E}_{[:K]}^T \mathbf{x}_{N-1}] = \mathbf{E}_{[:K]}^T \mathbf{X}^T$$

PCA reduces the number of features from D to K

KL / PCA for Data: Singular Value Decomposition (SVD)

SVD for an arbitrary matrix **A**

$$\underset{m \times n}{\mathbf{A}} = \underset{m \times m}{\mathbf{U}} \underset{m \times n}{\mathbf{\Sigma}} \underset{n \times n}{\mathbf{V}^T}$$

The **SVD** for **A** provides the KL factorization for the non-negative definite, symmetric $\mathbf{A}^T \mathbf{A}$

Use **SVD** to expand matrix $\mathbf{A}^T \mathbf{A}$

$$\begin{aligned} \mathbf{A}^T \mathbf{A} &= (\mathbf{U} \mathbf{\Sigma} \mathbf{V})^T \mathbf{U} \mathbf{\Sigma} \mathbf{V} \\ &= \mathbf{V} \mathbf{\Sigma}^T \mathbf{U}^T \mathbf{U} \mathbf{\Sigma} \mathbf{V} \\ &= \underset{n \times n}{\mathbf{V}} \underset{n \times n}{\mathbf{\Sigma} \mathbf{\Sigma}^T} \underset{n \times n}{\mathbf{V}^T} \\ &= \mathbf{E} \mathbf{\Lambda} \mathbf{E}^T \end{aligned}$$

U and **V** orthogonal matrices, **Σ** is the diagonal with singular values on diagonal

Note that this is also the **SVD** for $\mathbf{A}^T \mathbf{A}$

KL-Expansion / PCA for Data: Singular Value Decomposition (SVD)

SVD for a stacked data matrix **X**

$$\underset{N \times D}{\mathbf{X}} = \underset{N \times N}{\mathbf{U}} \underset{N \times D}{\mathbf{\Sigma}} \underset{D \times D}{\mathbf{V}^T}$$

$$\underset{D \times D}{\mathbf{X}^T \mathbf{X}} = \underset{D \times D}{\mathbf{V}} \underset{D \times D}{\mathbf{\Sigma \Sigma}^T} \underset{D \times D}{\mathbf{V}^T}$$

$$= \mathbf{E} \mathbf{\Lambda} \mathbf{E}^T$$

Equivalent approaches:

1. Find the **SVD** of **X**. Take **V**
2. Find the Eigen decomposition of **X^TX**. Take **E = V**
3. Find the **SVD** of **X^TX**. Take **V = U = E**

$$\underset{N \times K}{\tilde{\mathbf{X}}} = \underset{N \times D}{\mathbf{X}} \underset{D \times K}{\mathbf{V}_{[:K]}}$$

KL-Expansion / PCA for Data: Singular Value Decomposition (SVD)

Equivalent approaches:

1. Find the **SVD** of \mathbf{X} . Take \mathbf{V}
2. Find the Eigen decomposition of $\mathbf{X}^T\mathbf{X}$. Take $\mathbf{E} = \mathbf{V}$
3. Find the **SVD** of $\mathbf{X}^T\mathbf{X}$. Take $\mathbf{V} = \mathbf{U} = \mathbf{E}$

May want to use method 3, with `numpy.linalg.svd`, instead of method 2, with `numpy.linalg.eig`, since the **SVD** returns the Eigen-vectors in sorted order and `eig` does not.

$$\begin{matrix} \tilde{\mathbf{X}} & = & \mathbf{X} & \mathbf{V}_{[:K]} \\ N \times K & & N \times D & D \times K \end{matrix}$$

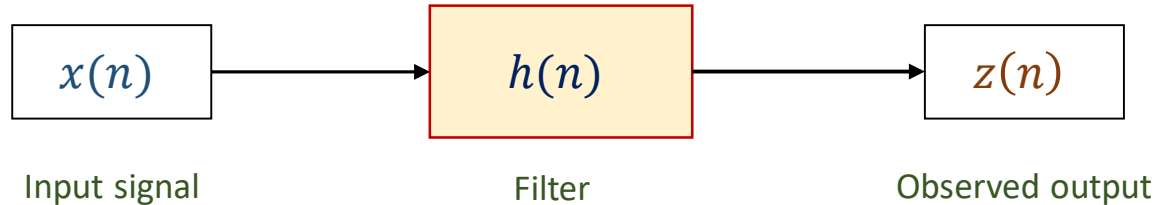
Outline

- Principal Component Analysis

-  • Filter Design

Filter Design

This involves creating a signal processing system that satisfies a pre-defined set of requirements.



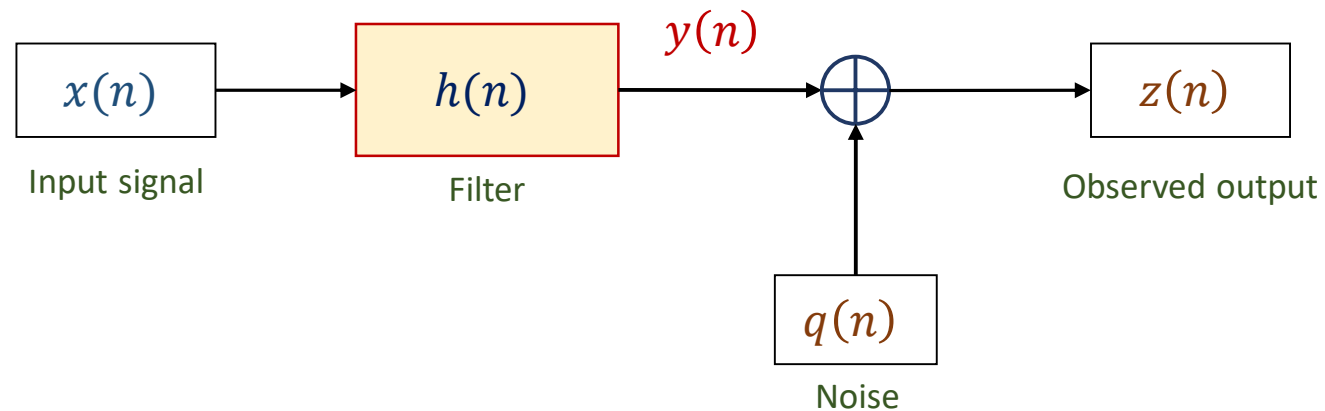
Example of Design requirements

- Frequency Response : *Low-pass, high-pass, band-pass, or band-stop*
- Impulse Response : *Finite (FIR) or infinite (IIR)*
- Causality : *Causal or non-causal*
- Phase and group delay

The goal is to find $h(n)$ that minimizes the error between the observed output $z(n)$ and the expected output $y(n)$

Homework 2: Filter Design for Additive Noise

Input signal $x(n)$ goes through a linear system and is observed in additive Gaussian noise.



$$z(n) = y(n) + q(n)$$

$$y(n) = h_0 x(n) + h_1 x(n-1) + h_2 x(n-2)$$

We need a system that minimize the mean squared error between $y(n)$ and $z(n)$

Wiener Filter

This is a *linear* and *time invariant* filter that separates original signal $y_n(u)$ from a corrupted signal.



Assumptions:

- Original signal and (additive) noise are stationary linear stochastic processes.
- Known spectral characteristics or known autocorrelation and cross-correlation.
- It is a minimum mean-square error (MMSE) estimator.

$$\hat{z}_n(u) = w_{lmmse}^T x_n(u)$$

$$E(w) = \mathbb{E}[\|y_n(u) - w^T x_n(u)\|^2]$$

$$w_{lmmse} = \arg \max_w E(w) = R_x^{-1} r_{xz}$$

Homework 2: Wiener Filter



$$w_{lmmse} = \arg \max_w E(w) = R_v^{-1} r_{vz} \quad (1)$$

$$y(n) = h_0 x(n) + h_1 x(n-1) + h_2 x(n-2)$$

$$E(w) = \mathbb{E}[\|y_n(u) - \mathbf{w}^T x_n(u)\|^2]$$

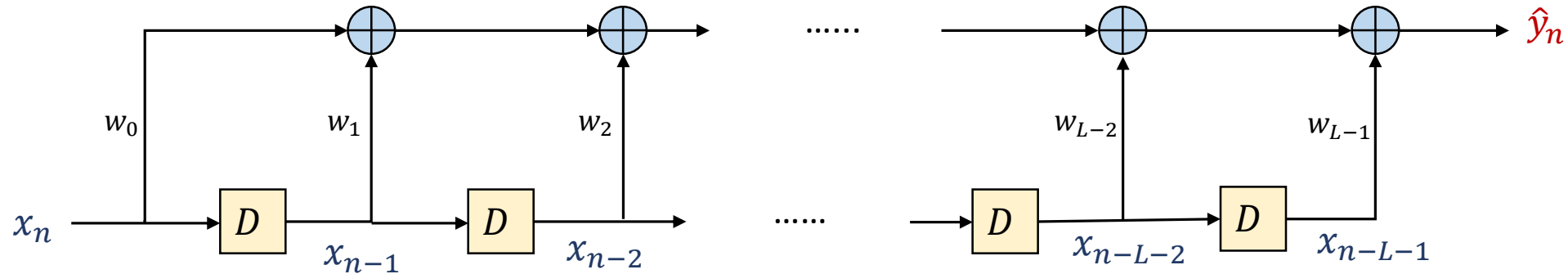
$$v_n(u) = [x_n(u), x_{n-1}(u), x_{n-1}(u)]^T$$

$$z(n) = y(n) + q(n)$$

- Show that : $r_{xz}[m] = \mathbb{E}[x_n(u)x_{n+m}(u)] = h_m$
- Compute auto-correlation R_{v_n} : $R_{v_n} = \mathbb{E}[v_n(u)v_n^T(u)]$
- Compute cross-correlation r_{vz} : $r_{vz} = \mathbb{E}[v_n(u)z_n(u)]$
- Compute MMSE : $E(w) = \mathbb{E}[\|y_n(u) - w^T x_n(u)\|^2]$

Least Mean Squares (LMS) Algorithm

This is a method for designing adaptive filters (finding the coefficients) that minimizes the mean square error.



$$\mathbf{v}_n = [x_n, x_{n-1}, \dots, x_{n-L-2}, x_{n-L-1}]^T$$

$$\hat{y}_n = \sum_{l=0}^{L-1} w_l x_{n-l} = \mathbf{w}^T \mathbf{v}_n$$

$$E(w) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{y}^{(i)} - \mathbf{w}^T \mathbf{v}^{(i)}\|^2$$

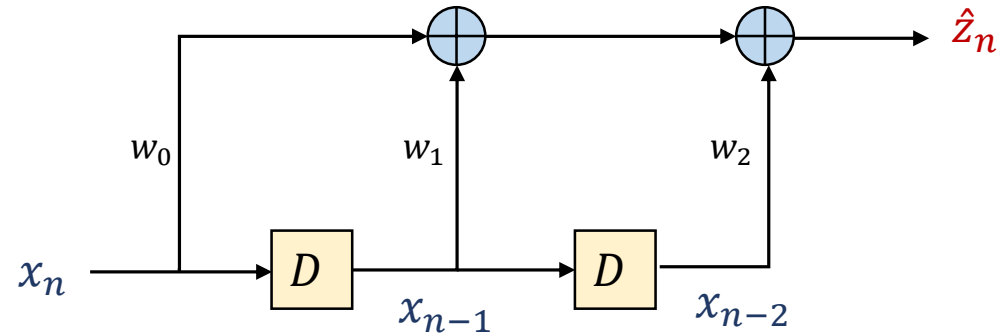
- Filter is adaptive
- It uses stochastic gradient descent (update rule)
- We need real samples $\{\mathbf{x}^{(i)}, \mathbf{y}^{(i)}\}_{i=1}^N$ to use the LMS algorithm.

Update rule with single point stochastic gradient descent

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \eta \nabla_{\mathbf{w}} E(w) = \mathbf{w}_n + \eta (\mathbf{y}^{(i)} - \mathbf{w}_n^T \mathbf{v}^{(i)}) \mathbf{v}^{(i)}$$

Homework 2 : LMS Algorithm

- Dataset : $\{\mathbf{x}^{(i)}, \mathbf{z}^{(i)}\}_{i=1}^N$
- Number of sequence N : 600
- Length of sequence : $n \in \{0,1,2, \dots, 500\}$
- Vary the learning rate η .



$$\mathbf{v}_n = [x_n, x_{n-1}, x_{n-2}]^T$$

$$\hat{z}_n = \sum_{l=0}^2 w_l x_{n-l} = \mathbf{w}^T \mathbf{v}_n$$

$$E(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{z}^{(i)} - \mathbf{w}^T \mathbf{v}^{(i)}\|^2$$

Update rule with single point stochastic gradient descent

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \eta \nabla_{\mathbf{w}} E(\mathbf{w}) = \mathbf{w}_n + \eta (\mathbf{y}^{(i)} - \mathbf{w}_n^T \mathbf{v}^{(i)}) \mathbf{v}^{(i)}$$

Code for the LMS algorithm and PCA !!!