

EECS225A Course Notes

Anmol Parande

Fall 2021 - Professor Jiantao Jiao

Disclaimer: These notes reflect EECS225A when I took the course (Fall 2021). They may not accurately reflect current course content, so use at your own risk. If you find any typos, errors, etc, please raise an issue on the [GitHub repository](#).

Contents

1 Hilbert Space Theory	2
2 Linear Estimation	3
2.1 Affine Estimation	3
2.2 Least Squares	4
3 Discrete Time Random Processes	4
3.1 Spectral Analysis	4
3.2 LTI Filtering	7
3.2.1 Wiener Filtering	8

1 Hilbert Space Theory

Complex random vectors form a Hilbert space with inner product $\langle X, Y \rangle = \mathbb{E}[XY^*]$. If we have a random complex vector, then we can use Hilbert Theory in a more efficient manner by looking at the matrix of inner products. For simplicity, we will call this the “inner product” of two complex vectors.

Definition 1 Let the inner product between two random, complex vectors $\mathbf{Z}_1, \mathbf{Z}_2$

$$\langle \mathbf{Z}_1, \mathbf{Z}_2 \rangle = \mathbb{E}[\mathbf{Z}_1 \mathbf{Z}_2^*]$$

The ij -th entry of the matrix is simply the scalar inner product $\mathbb{E}[\mathbf{X}_i \mathbf{Y}_j^*]$ where \mathbf{X}_i and \mathbf{Y}_j are the i th and j th entries of \mathbf{X} and \mathbf{Y} respectively. This means the matrix is equivalent to the cross correlation \mathbf{R}_{XY} between the two vectors. We can also specify the auto-correlation $\mathbf{R}_X = \langle \mathbf{X}, \mathbf{X} \rangle$ and auto-covariance $\Sigma_X = \langle \mathbf{X} - \mathbb{E}[\mathbf{X}], \mathbf{X} - \mathbb{E}[\mathbf{X}] \rangle$. One reason why we can think of this matrix as the inner product is because it also satisfies the properties of inner products. In particular, it is

1. Linear: $\langle \alpha_1 \mathbf{V}_1 + \alpha_2 \mathbf{V}_2, \mathbf{u} \rangle = \alpha_1 \langle \mathbf{V}_1, \mathbf{u} \rangle + \alpha_2 \langle \mathbf{V}_2, \mathbf{u} \rangle$.
2. Reflexive: $\langle \mathbf{U}, \mathbf{V} \rangle = \langle \mathbf{V}, \mathbf{U} \rangle^*$.
3. Non-degeneracy: $\langle \mathbf{V}, \mathbf{V} \rangle = 0 \Leftrightarrow \mathbf{V} = \mathbf{0}$.

Since we are thinking of the matrix as an inner product, we can also think of the norm as a matrix.

Definition 2 The norm of a complex random vector is given by $\|\mathbf{Z}\|^2 = \langle \mathbf{Z}, \mathbf{Z} \rangle$.

Since we are thinking of inner products as matrices instead of scalars, we can rewrite the Hilbert Projection Theorem to use matrices instead.

Theorem 1 (Hilbert Projection Theorem) The minimization problem $\min_{\hat{\mathbf{X}}(\mathbf{Y})} \|\hat{\mathbf{X}}(\mathbf{Y}) - \mathbf{X}\|^2$ has a unique solution which is a linear function of \mathbf{Y} . The error is orthogonal to the linear subspace of \mathbf{Y} (i.e. $\langle \mathbf{X} - \hat{\mathbf{X}}, \mathbf{Y} \rangle = 0$)

When we do a minimization over a matrix, we are minimizing it in a PSD sense, so for any other linear function \mathbf{X}' ,

$$\|\mathbf{X} - \hat{\mathbf{X}}\|^2 \preceq \|\mathbf{X} - \mathbf{X}'\|^2.$$

2 Linear Estimation

In Linear Estimation, we are trying to estimate a random variable \mathbf{X} using an observation \mathbf{Y} with a linear function of \mathbf{Y} . If \mathbf{Y} is finite dimensional, then we can say $\hat{\mathbf{X}}(\mathbf{Y}) = \mathbf{W}\mathbf{Y}$ where \mathbf{W} is some matrix. Using theorem 1 and the orthogonality principle, we know that

$$\langle \mathbf{X} - \mathbf{W}\mathbf{Y}, \mathbf{Y} \rangle = \mathbf{0} \Leftrightarrow \mathbf{R}_{XY} = \mathbf{W}\mathbf{R}_Y$$

This is known as the **Normal Equation**. If \mathbf{R}_Y is invertible, then we can apply the inverse to find \mathbf{W} . Otherwise, we can apply the pseudoinverse \mathbf{R}_Y^\dagger to find \mathbf{W} , which may not be unique. If we want to measure the quality of the estimation, since $\mathbf{X} = \hat{\mathbf{X}} + (\mathbf{X} - \hat{\mathbf{X}})$,

$$\begin{aligned} \|\mathbf{X}\|^2 &= \|\hat{\mathbf{X}}\|^2 + \|\mathbf{X} - \hat{\mathbf{X}}\|^2 \implies \\ \|\mathbf{X} - \hat{\mathbf{X}}\|^2 &= \|\mathbf{X}\|^2 - \|\hat{\mathbf{X}}\|^2 = \mathbf{R}_X - \mathbf{R}_{XY}\mathbf{R}_Y^{-1}\mathbf{R}_{YX} \end{aligned}$$

2.1 Affine Estimation

If we allow ourselves to consider an affine function for estimation $\hat{\mathbf{X}}(\mathbf{Y}) = \mathbf{W}\mathbf{Y} + b$, then this is equivalent to instead finding an estimator

$$\hat{\mathbf{X}}(\mathbf{Y}') = \mathbf{W}\mathbf{Y}' \quad \text{where } \mathbf{Y}' = \begin{bmatrix} \mathbf{Y} \\ 1 \end{bmatrix}$$

This is equivalent to the following orthogonality conditions:

1. $\langle \mathbf{X} - \hat{\mathbf{X}}, \mathbf{Y} \rangle$
2. $\langle \mathbf{X} - \hat{\mathbf{X}}, 1 \rangle$

Solving gives us

$$\hat{\mathbf{X}}(\mathbf{Y}) = \mathbf{W}(\mathbf{Y} - \boldsymbol{\mu}_Y) + \mu_x \quad \text{where } \mathbf{W}\boldsymbol{\Sigma}_Y = \boldsymbol{\Sigma}_{XY}.$$

$\boldsymbol{\Sigma}_Y$ and $\boldsymbol{\Sigma}_{XY}$ are the auto-covariance and cross-covariance respectively. Recall that if

$$\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_Y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_X & \boldsymbol{\Sigma}_{XY} \\ \boldsymbol{\Sigma}_{YX} & \boldsymbol{\Sigma}_Y \end{bmatrix} \right)$$

then

$$\mathbf{X}|\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu}_X + \boldsymbol{\Sigma}_{XY}\boldsymbol{\Sigma}_Y^{-1}(\mathbf{Y} - \boldsymbol{\mu}_Y), \boldsymbol{\Sigma}_X - \boldsymbol{\Sigma}_{XY}\boldsymbol{\Sigma}_Y^{-1}\boldsymbol{\Sigma}_{YX})$$

The mean is the best affine estimator of \mathbf{X} , and the covariance is the estimation error. This has two interpretations.

1. Under the Gaussian assumption, the best nonlinear estimator $\mathbb{E}[\mathbf{X}|\mathbf{Y}]$ is affine
2. Gaussian random variables are the hardest to predict because nonlinearity should improve our error, but it does not in the Gaussian case. This means if affine estimation works well, we shouldn't try and find better non-linear estimators.

2.2 Least Squares

The theory of linear estimation is very closely connected with the theory behind least squares in linear algebra. In least squares, we have a deterministic \mathbf{x} and assume nothing else about it, meaning we are looking for an unbiased estimator. theorem 2 tells us how to find the best linear unbiased estimator in a linear setting.

Theorem 2 (Gauss Markov Theorem) Suppose that $\mathbf{Y} = \mathbf{H}\mathbf{x} + \mathbf{Z}$ and \mathbf{Z} is zero-mean with $\langle \mathbf{Z}, \mathbf{Z} \rangle = \mathbf{I}$, \mathbf{H} is full-column rank, then $\hat{\mathbf{x}}_b = (\mathbf{H}^* \mathbf{H})^{-1} \mathbf{H}^* \mathbf{Y}$ is the best linear unbiased estimator.

3 Discrete Time Random Processes

Definition 3 A Discrete-Time Random Process is a countably infinite collection of random variables on the same probability space $\{X_n : n \in \mathbb{Z}\}$.

Discrete Time Random Processes have a mean function $\mu_n = \mathbb{E}[X_n]$ and an autocorrelation function $R_X(n_1, n_2) = \mathbb{E}[X_{n_1} X_{n_2}^*]$

Definition 4 A Wide-Sense Stationary Random Process is a discrete-time random process with constant mean, finite variance, and an autocorrelation function that can be re-written to only depend on $n_1 - n_2$.

We call this wide-sense stationary because the mean and covariance do not change as the process evolves. In a strict-sense stationary process, the distribution of each random variable in the process would not change.

3.1 Spectral Analysis

Recall that the Discrete Time Fourier Transform is given by

$$X(e^{j\omega}) = \sum_{n=-\infty}^{\infty} x[n]e^{-j\omega n}.$$

The Inverse Discrete Time Fourier Transform is given by

$$x[n] = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(e^{j\omega}) e^{j\omega n} d\omega.$$

Since the DTFT is an infinite summation, it may or may not converge.

Definition 5 A signal $x[n]$ belongs to the l^1 class of signals if the series converges absolutely. In other words,

$$\sum_{k=-\infty}^{\infty} |x[k]| < \infty.$$

This class covers most real-world signals.

Theorem 3 If $x[n]$ is a l^1 signal, then the DTFT $X(e^{j\omega})$ converges uniformly and is well-defined for every ω . $X(e^{j\omega})$ is also a continuous function.

Definition 6 A signal $x[n]$ belongs to the l^2 class of signals if it is square summable. In other words,

$$\sum_{k=-\infty}^{\infty} |x[k]|^2 < \infty.$$

The l^2 class contains important functions such as sinc.

Theorem 4 If $x[n]$ is a l^2 signal, then the DTFT $X(e^{j\omega})$ is defined almost everywhere and only converges in the mean-squared sense:

$$\lim_{N \rightarrow \infty} \int_{-\pi}^{\pi} \left| \left(\sum_{k=-N}^N x[k] e^{-j\omega n} \right) - X(\omega) \right|^2 d\omega = 0$$

Tempered distributions like the Dirac Delta function are other functions which are important for computing the DTFT, and they arise from the theory of generalized functions.

Suppose we want to characterize the signal using its DTFT.

Definition 7 The energy of a deterministic, discrete-time signal $x[n]$ is given by

$$\sum_{n \in \mathbb{Z}} |x[n]|^2.$$

The autocorrelation of $x[n]$, given by $a[n] = x[n] * x^*[-n]$, is closely related to the energy of the signal since $a[0] = \sum_{n \in \mathbb{Z}} |x(n)|^2$.

Definition 8 *The Energy Spectral Density $x[n]$ with auto-correlation $a[n]$ is given by*

$$A(\omega) = \sum_{n \in \mathbb{Z}} a[n] e^{j\omega n}$$

We call the DTFT of the autocorrelation the energy spectral density because by the Inverse DTFT,

$$a[0] = \frac{1}{2\pi} \int_{-\pi}^{\pi} A(\omega) d\omega.$$

Since summing over each frequency gives us the energy spectral density, we can think of $A(\omega)$ as storing the energy density of each spectral component of the signal. We can apply this same idea to wide-sense stationary stochastic processes.

Definition 9 *The Power Spectral Density of a Wide-Sense Stationary random process is given by*

$$S_X(\omega) = \sum_{k \in \mathbb{Z}} R_X(k) e^{-j\omega k}.$$

Note that when considering stochastic signals, the metric changes from energy to power. This is because stochastic processes are not generally l_2 , and if X_n is Wide-Sense Stationary, then

$$\mathbb{E} \left[\sum_{n \in \mathbb{Z}} |x(n)|^2 \right] = \infty,$$

so energy doesn't even make sense. To build our notion of power, let $A_T(\omega)$ be a truncated DTFT of the auto-correlation of a wide-sense stationary process, then

$$\begin{aligned} \lim_{T \rightarrow \infty} \frac{\mathbb{E}[A_T(\omega)]}{2T+1} &= \lim_{T \rightarrow \infty} \frac{1}{2T+1} \left(\sum_{n=-T}^T x[n] e^{-j\omega n} \right) \left(\sum_{m=-T}^T x^*[m] e^{j\omega m} \right) \\ &= \lim_{T \rightarrow \infty} \frac{1}{2\pi} \sum_{n,m \in [-T,T]} \mathbb{E}[x[n] x^*[m]] e^{-j\omega(n-m)} \\ &= \lim_{T \rightarrow \infty} \frac{1}{2\pi} \sum_{n,m \in [-T,T]} R_x(n-m) e^{-j\omega(n-m)} \\ &= \lim_{T \rightarrow \infty} \sum_{k=-2T}^{2T} R_X(k) e^{-j\omega k} \left(1 - \frac{|k|}{2T+1} \right) \\ &= \sum_{k=-\infty}^{\infty} R_X(k) e^{-j\omega k} \end{aligned}$$

The DTFT of the auto-correlation function naturally arises out of taking the energy spectral density and normalizing it by time (the truncated sequence is made of $2T + 1$ points). In practice, this means to measure the PSD, we need to either use the distribution of the signal to compute R_X , or estimate the *PSD* by averaging multiple realizations of the signal.

The inverse DTFT formula tells us that we can represent a deterministic, discrete-time signal $x[n]$ as a sum of complex exponentials weighted by $\frac{X(\omega)d\omega}{2\pi}$. This representation has an analog for stochastic signals as well.

Theorem 5 (Cramer-Khinchin) *For a complex-valued WSS stochastic process X_n with power spectral density $S_X(\omega)$, there exists a unique right-continuous stochastic process $F(\omega), \omega \in (-\pi, \pi]$ with square-integrable, orthogonal increments such that*

$$X_n = \int_{-\pi}^{\pi} e^{j\omega n} dF(\omega)$$

where for any interval $[\omega_1, \omega_2], [\omega_3, \omega_4] \subset [-\pi, \pi]$,

$$\mathbb{E} [(F(\omega_2) - F(\omega_1))(F(\omega_4) - F(\omega_3))^*] = f((\omega_1, \omega_2] \cap (\omega_3, \omega_4])$$

where f is the structural measure of the stochastic process and has Radon-Nikodym derivative $\frac{S_X(\omega)}{2\pi}$.

Besides giving us a decomposition of a WSS random process, theorem 5 tells a few important facts.

1. $\omega_1 \neq \omega_2 \implies \langle dF(\omega_1), dF(\omega_2) \rangle = 0$ (i.e different frequencies are uncorrelated).
2. $\mathbb{E} [|dF(\omega)|^2] = \frac{S_X(\omega)}{2\pi}$

3.2 LTI Filtering

Recall that the Z-transform converts a discrete-time signal into a complex representation. It is given by

$$X(z) = \sum_{n=-\infty}^{\infty} x[n]z^{-n}.$$

It is a special type of series called a **Laurent Series**.

Theorem 6 *A Laurent Series will converge absolutely on an open annulus*

$$A = \{z | r < |z| < R\}$$

for some r and R .

We can compute r and R using the signal $x[n]$.

$$r = \limsup_{n \rightarrow \infty} |x[n]|^{\frac{1}{n}}, \quad \frac{1}{R} = \limsup_{n \rightarrow \infty} |x[-n]|^{\frac{1}{n}}.$$

Definition 10 For two jointly WSS processes X_n, Y_n , the z -cross spectrum is the Z -Transform of the correlation function $R_{YX}(k) = \mathbb{E}[Y_n X_{n-k}^*]$.

$$S_{YX}(z) = \sum_{k \in \mathbb{Z}} R_{YX}(k) z^{-k}$$

Using this definition, we can see that

$$S_{XY}(z) = S_{YZ}^*(z^{-*}).$$

Theorem 7 When $Y(n)$ is formed by passing a WSS process X_n through a stable LTI system with impulse response $h[n]$ and transfer function $H(z)$, then $S_Y(z) = H(z)S_X(z)H^*(z^{-*})$ and $S_{YX}(z) = H(z)S_X(z)$. If we have a third process Z_n that is jointly WSS with (Y_n, X_n) , then $S_{ZY}(z) = S_{ZX}(z)H^*(z^{-*})$.

3.2.1 Wiener Filtering

Suppose we have a stochastic WSS process Y_n that is jointly WSS with X_n . The best linear estimator of X_n given the observations Y_n can be written as

$$\hat{X}_n = \sum_{m \in \mathbb{Z}} h(m) Y_{n-m} = h[n] * Y_n.$$

This is identical to passing Y_n through an LTI filter. Starting with the orthogonality principle,

$$\begin{aligned} \mathbb{E}[(X_n - \hat{X}_n)Y_{n-k}^*] &= 0 \implies \mathbb{E}[X_n Y_{n-k}^*] = \sum_{m \in \mathbb{Z}} h(m) \mathbb{E}[Y_{n-m} Y_{n-k}^*] \\ \therefore R_{XY}(k) &= \sum_{m \in \mathbb{Z}} h(m) R_Y(k-m) \implies S_{XY}(e^{j\omega}) = H(e^{j\omega}) S_Y(e^{j\omega}) \\ \therefore H(e^{j\omega}) &= \frac{S_{XY}(e^{j\omega})}{S_Y(e^{j\omega})} \end{aligned}$$

Definition 11 *The best linear estimator of X_n using Y_n where (X_n, Y_n) is jointly WSS is given by the non-causal Wiener filter.*

$$H(e^{j\omega}) = \frac{S_{XY}(e^{j\omega})}{S_Y(e^{j\omega})}$$

For a specific ω , we can understand $H(e^{j\omega})$ as an optimal linear estimator for $F_X(\omega)$ where $F_X(\omega)$ is the stochastic process given by the Cramer-Khinchin decomposition (theorem 5). More specifically, we can use the Cramer-Khinchin decomposition of Y_n .

$$\begin{aligned}\hat{X}_n &= \sum_{i \in \mathbb{Z}} h[i] \int_{-\pi}^{\pi} e^{j\omega(n-i)} dF_Y(\omega) \\ &= \int_{-\pi}^{\pi} \left(\sum_{i \in \mathbb{Z}} h[i] e^{-j\omega i} \right) e^{j\omega n} dF_Y(\omega) \\ &= \int_{-\pi}^{\pi} H(e^{j\omega}) e^{j\omega n} dF_Y(\omega)\end{aligned}$$

Since F_X and F_Y have jointly orthogonal increments, this tells us that $H(e^{j\omega})$ is just the optimal linear estimator of $dF_X(\omega)$ using $dF_Y(\omega)$. $dF_X(\omega)$ and $dF_Y(\omega)$ exist on a Hilbert space, meaning we are essentially projecting each frequency component of X_n onto the corresponding frequency component of Y_n .