

For Personal Use Only -bkwk2

Probability and random variables

Probability

sample space, probability and probability space.

- A random expt is any situation which has a set of possible outcomes, each of which occurs w/ a particular probability
- To mathematically describe a random expt, we must specify:

↳ 1) The sample space Ω : The set of all possible outcomes of the random expt.

We can call any subset $A \subseteq \Omega$ an event

↳ 2) The probability, P : A mapping/function from events to a no. in the interval $[0, 1]$

- We call (Ω, P) the probability space.

Axioms of probability

- Probability theory is based on the Kolmogorov axioms:

↳ Axiom I: The probability of an event is a non-negative real no.

$$P(A) \in \mathbb{R}, \quad P(A) \geq 0 \quad \forall A \subseteq \Omega$$

↳ Axiom II: The sample space / certain event has unit probability

$$P(\Omega) = 1.$$

↳ Axiom III: Additivity for incompatible events (disjoint sets)

$$P(A \cup B) = P(A) + P(B) \quad \text{if } A \cap B = \emptyset.$$

- These three axioms are sufficient to prove the following:

↳ Monotonicity:

$$\text{If } A \subseteq B, \text{ then } P(A) \leq P(B)$$

↳ Probability of the empty set:

$$P(\emptyset) = 0$$

↳ Complement rule:

$$P(A') = 1 - P(A)$$

↳ Numeric bound:

$$0 \leq P(A) \leq 1 \quad \forall A \subseteq \Omega$$

↳ Addition law:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

↳ Sum rule:

$$P(A) = P(A|B) + P(A|B')$$

↳ Additivity extended: $P(A) = \sum_{\omega \in A} P(\omega)$ where $\{\omega\}$ are mutually-disjoint individual outcomes

For Personal Use Only -bkwk2

Defining probability

- Define the indicator function for a set or event E

$$I_E(t) = \begin{cases} 0 & t \in E \\ 1 & t \notin E \end{cases}$$

- For a finite discrete set $\Omega = \{\omega_1, \omega_2, \dots, \omega_n\}$, let p_1, p_2, \dots, p_n be non-negative no. that sum to 1. For any event E , set

$$P(E) = \sum_{i=1}^n I_E(\omega_i) p_i$$

then P satisfies the axioms.

- For an infinite discrete set $\Omega = \{\omega_1, \omega_2, \dots\}$, let p_1, p_2, \dots be a non-negative sequence that sum to 1. For any event E , set

$$P(E) = \sum_{i=1}^{\infty} I_E(\omega_i) p_i$$

then P satisfies the axioms

- For the set of real no $\Omega = \mathbb{R}$, we define a pdf $f(t)$ that satisfies

$$f(t) \geq 0 \quad \text{and} \quad \int_{-\infty}^{\infty} f(t) dt = 1$$

For any event E , set

$$P(E) = \int_{-\infty}^{\infty} I_E(t) f(t) dt$$

then P satisfies the axioms.

Conditional probability

- The conditional probability of event A occurring given that event B has occurred is defined to be.

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \text{if } P(B) > 0$$

- From this, we can get the multiplication rule

$$P(A \cap B) = P(A|B) \cdot P(B)$$

In general, we can define the probability chain rule

$$P(A_1, \dots, A_n | A_{n+1}, \dots, A_m) = P(A_1) \prod_{i=2}^m P(A_i | A_1, \dots, A_{i-1})$$

- Applying the sum rule, we get the law of total probability

$$\begin{aligned} P(A) &= P(A \cap B) + P(A \cap B') \\ &= P(A|B) P(B) + P(A|B') P(B') \end{aligned}$$

In general, for a set of pairwise incompatible events B_i , s.t. $\bigcup B_i = \Omega$,

$$P(A) = \sum_i P(A|B_i) P(B_i)$$

For Personal Use Only -bkwk2

Bayes' theorem

- We can write $P(A \cap B)$ in two ways using the multiplication rule.

$$P(A \cap B) = P(A|B)P(B) = P(B|A)P(A)$$

Rearranging, we get Bayes's thm. (relates $P(A|B)$ w/ $P(B|A)$)

$$P(B|A) = \frac{P(A|B)P(B)}{P(A)}$$

- Bayes's thm. is often used w/ the law of total probability.

$$P(B|A) = \frac{P(A|B)P(B)}{P(A|B)P(B) + P(A|B')P(B')}$$

Independence

- Two events are independent if the knowledge that one occurred does not change the probability that the other one occurs.
- Mathematically, A and B are independent if.

$$P(A) = P(A|B) = P(A|B')$$

Using the multiplication rule, we get the formal defn of independence

$$P(A \cap B) = P(A)P(B)$$

Random variables

Random variables (RVs)

- Given a probability space (Ω, P) , a RV is a function $X(\omega)$ which maps each element ω of the sample space Ω onto a pt. on the real line.

- We can define a RV by listing the range of $X(\omega)$ along w/ the probability the RV takes those values.

- For any set $A \subset (-\infty, \infty)$, we define

$$P(X \in A) = P(\{\omega : X(\omega) \in A\})$$

- For a discrete RV X , we define the probability mass function (pmf) P_X to be.

$$P_X(x_i) = P(X=x_i), \text{ where } \sum_i P_X(x_i) = 1.$$

$$\text{For any set } A, P(X \in A) = \sum_i I_A(x_i) P_X(x_i)$$

- For a continuous RV X , we define the probability density function (pdf) f_X to be

$$f_X(x) = P(X \leq x), \text{ where } \int_{-\infty}^{\infty} f_X(x) dx = 1.$$

$$\text{For any set } A, P(X \in A) = \int_{-\infty}^{\infty} I_A(x) f_X(x) dx.$$

For Personal Use Only -bkwk2

cumulative distribution function (cdf)

- For discrete/continuous RV, the cdf F_X is defined as

$$F_X(x) = P(X \leq x)$$

note $P(X > x) = 1 - F_X(x)$

- The cdf $F_X(x)$ has the following properties:

$$\hookrightarrow 0 \leq F_X(x) \leq 1$$

$$\hookrightarrow P(X_1 < X \leq X_2) = F_{X_2}(x_2) - F_{X_1}(x_1)$$

$\hookrightarrow F_X(x)$ is non-decreasing as x increases

$$\hookrightarrow \lim_{x \rightarrow -\infty} F_X(x) = 0 \text{ and } \lim_{x \rightarrow \infty} F_X(x) = 1.$$

- For continuous RV X , $F_X(x) = \int_{-\infty}^x f_X(t) dt \rightarrow$ area under f_X is a continuous function $\rightarrow F_X$ is continuous.

- For continuous RV X , the cdf F_X and pdf f_X are related by.

$$F_X(x) = \int_{-\infty}^x f_X(t) dt \xrightarrow{\text{FTC}} f_X(x) = \frac{dF_X(x)}{dx} \quad \text{useful for finding pdf of transformed variables!}$$

- For discrete RV X w/ range $\{x_1, \dots, x_M\}$, $F_X(x) = \sum_{k=1}^M p(x_k) I_{\{x_k \leq x\}}(x) \rightarrow F_X(x)$ is a step function that

jumps w/ each $x_k \rightarrow F_X$ is right-continuous, i.e. $F_X(x) = \lim_{x \rightarrow x^+} F_X(x) \quad \forall x$.

Expectation

- The expected value of X , $E[X]$ is defined as.

$$E[X] = \begin{cases} \sum x_i p(x_i) & \text{[discrete]} \\ \int x f_X(x) dx & \text{[continuous]} \end{cases}$$

- The law of large nos (LLN) states that.

$$\frac{1}{n} \sum_{i=1}^n x_i \rightarrow E[X] \quad \text{as } n \rightarrow \infty$$

so we can think of $E[X]$ as the empirical average of independent samples of X , x_i .

- For any function $g(x)$, its expectation $E[g(x)]$ is computed using.

$$E[g(x)] = \begin{cases} \sum g(x_i) p(x_i) & \text{[discrete]} \\ \int g(x) f_X(x) dx & \text{[continuous]} \end{cases}$$

- For $g(x) = I_A(x)$, its expectation is the probability of event A .

$$E[I_A(x)] = \left[\frac{\sum I_A(x_i) p(x_i)}{\int I_A(x) f_X(x) dx} \right] = P(X \in A)$$

The LLN formalises the freq. interpretation of probability.

$$\frac{1}{n} \sum_{i=1}^n I_A(x_i) \rightarrow E[I_A(x)] = P(A \text{ fix}) \quad \text{as } n \rightarrow \infty.$$

Transformation of random variable.

- Consider a continuous RV X w/ pdf f_X and cdf F_X . Define a new RV $Y = g(X)$. For strictly increasing/decreasing g , there exist an inverse g^{-1} , and the pdf of the transformed RV, f_Y , is given by.

$$f_Y(y) = f_X(g^{-1}(y)) \left| \frac{d(g^{-1}(y))}{dy} \right|$$

\hookrightarrow increasing $g [g(x) \leq y \Leftrightarrow x \leq g^{-1}(y)]$: $F_Y(y) = P(g(X) \leq y) = P(X \leq g^{-1}(y)) = F_X(g^{-1}(y))$

$$\therefore f_Y(y) = \frac{dF_Y(y)}{dy} = \frac{d}{dy} F_X(g^{-1}(y)) = \frac{dF_X(g^{-1}(y))}{dx} \frac{dx}{dy} = f_X(g^{-1}(y)) \left| \frac{d(g^{-1}(y))}{dy} \right|$$

decreasing $g [g(x) \leq y \Leftrightarrow x \geq g^{-1}(y)]$: $F_Y(y) = P(g(X) \leq y) = P(X \geq g^{-1}(y)) = 1 - F_X(g^{-1}(y))$

$$\therefore f_Y(y) = \frac{dF_Y(y)}{dy} = \frac{d}{dy} (1 - F_X(g^{-1}(y))) = - \frac{dF_X(g^{-1}(y))}{dx} \frac{dx}{dy} = f_X(g^{-1}(y)) \left| \frac{d(g^{-1}(y))}{dy} \right|$$

* In general, express $F_Y(y)$ in terms of F_X , then differentiate w/ respect to y to get $f_Y(y)$.

For Personal Use Only -bkwk2

Bivariate random variables.

Bivariate discrete random variables.

- Consider two discrete RVs X, Y where $X \in \{x_1, \dots, x_n\}$ and $Y \in \{y_1, \dots, y_m\}$. Define the joint pmf P_{XY}

$$P_{XY}(x_i, y_j) = P(X=x_i, Y=y_j)$$

- We can derive the marginal pmfs P_X, P_Y given the joint pmf P_{XY} .

$$P_X(x_k) = \sum_{j=1}^m P_{XY}(x_k, y_j)$$

$$P_Y(y_l) = \sum_{k=1}^n P_{XY}(x_k, y_l)$$

- Two discrete RVs X and Y are independent if

$$P_{XY}(x, y) = P_X(x) P_Y(y) \quad \forall (x, y).$$

- For the discrete RVs X, Y , the conditional pmf of X given $Y=y$, $P_{X|Y}$ is defined as

$$P_{X|Y}(x|y) = \frac{P_{XY}(x,y)}{P_Y(y)}.$$

Bivariate continuous random variables.

- Consider two continuous RVs X, Y , we define the joint pdf f_{XY} to be a non-negative function s.t.

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I_A(x) I_B(y) f_{XY}(x, y) dx dy = P(X \in A, Y \in B), \text{ where } \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{XY}(x, y) dx dy = 1.$$

- We can derive the marginal pdfs f_X, f_Y given the joint pdf f_{XY} .

$$f_X(x) = \int_{-\infty}^{\infty} f_{XY}(x, y) dy$$

$$f_Y(y) = \int_{-\infty}^{\infty} f_{XY}(x, y) dx$$

- Two continuous RVs X and Y are independent iff

$$f_{XY}(x, y) = f_X(x) f_Y(y)$$

- For the continuous RVs X, Y , the conditional pdf of X given $Y=y$, $f_{X|Y}$ is defined as

$$f_{X|Y}(x|y) = \frac{f_{XY}(x, y)}{f_Y(y)}$$

Sum of random variables.

- Let X, Y be two independent RVs and let $Z = X+Y$.

$$\text{For discrete } X, Y, P_Z(z) = \sum_{x \in X} P_X(x) P_Y(z-x) = P_X * P_Y$$

$$\hookrightarrow P_Z(z) = \sum_{x \in X} P_{XY}(x, z) = \sum_{x \in X} P_{Z|X}(z|x) P_X(x) = \sum_{x \in X} P(X+Y=z | X=x) P_X(x) = \sum_{x \in X} P(Y=z-x | X=x) P_X(x) \\ = \sum_{x \in X} P_{Y|X}(z-x) P_X(x) = \sum_{x \in X} P_Y(z-x) P_X(x) = P_X * P_Y.$$

$$\text{For continuous } X, Y, f_Z(z) = \int_{-\infty}^{\infty} f_X(x) f_Y(z-x) dx = f_X * f_Y$$

$$\hookrightarrow F_Z(z) = P(Z \leq z) = P(X+Y \leq z) = \iint_{X+Y \leq z} f_{XY}(x, y) dx dy = \int_{-\infty}^{\infty} \int_{-\infty}^{z-y} f_{XY}(x, y) dx dy$$

$$\text{Sub } u=x+y: F_Z(z) = \int_{-\infty}^{\infty} \int_{-\infty}^z f_{XY}(u-y, y) du dy = \int_{-\infty}^{\infty} \int_{-\infty}^z f_X(u-y) f_Y(y) du dy$$

$$\text{Diff. w.r.t } z: f_Z(z) = \frac{d}{dz} F_Z(z) = \int_{-\infty}^{\infty} f_X(z-y) f_Y(y) dy = f_X * f_Y.$$

For Personal Use Only -bkwk2

Expectation of bivariate

- The expectation of a function $g(X,Y)$ of the bivariate (X,Y) is given by.

$$E[g(x,y)] = \begin{cases} \sum_{x,y} g(x,y) P_{X,Y}(x,y) & [\text{discrete}] \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x,y) f_{X,Y}(x,y) dx dy & [\text{continuous}] \end{cases}$$

The conditional expectation of $g(x,y)$ given $Y=y$ is then

$$E[g(x,y) | Y=y] = \begin{cases} \sum_x g(x,y) P_{X|Y}(x|y) & [\text{discrete}] \\ \int_{-\infty}^{\infty} g(x,y) f_{X|Y}(x|y) dx & [\text{continuous}] \end{cases}$$

- We can compute $E[g(x,y)]$ by considering $E[g(x,y) | Y=Y]$. For the continuous case,

$$E[g(x,y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x,y) f_{X,Y}(x,y) dx dy = \int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} g(x,y) f_{X|Y}(x|y) dx \right) f_Y(y) dy$$

$$E[g(x,y)] = \int_{-\infty}^{\infty} E[g(x,y) | Y=y] f_Y(y) dy = E_Y [E[g(x,y) | Y]]$$

Note that for a function $h(Y)$, $E[h(Y)] = \int_{-\infty}^{\infty} h(Y=y) f_Y(y) dy$

- * The above result also holds for the discrete case

Law of iterated expectation (Adam's law)

- The law of iterated expectation states that

$$E[X] = E_Y [E_X[X|Y]]$$

- Consider the continuous case (result also holds for discrete case)

$$E[X] = \int_{-\infty}^{\infty} x f_X(x) dx = \int_{-\infty}^{\infty} x \left(\int_{-\infty}^{\infty} f_{X|Y}(x|y) dy \right) dx = \int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} x f_{X|Y}(x|y) dy \right) f_Y(y) dy = \int_{-\infty}^{\infty} E[X|Y=y] f_Y(y) dy = E_Y [E_X[X|Y]]$$

Multivariate random variables

Random vectors

- Let X_1, X_2, \dots, X_n be n continuous RVs. We call $\mathbf{X} = (X_1, \dots, X_n) \in \mathbb{R}^n$ a continuous random vector

(similarly we can define a discrete random vector, but we will consider the continuous case in this section)

- For a continuous random vector $\mathbf{X} = (X_1, \dots, X_n) \in \mathbb{R}^n$, we define the pdf f to be a non-negative function that is integrated to unity which satisfies the following for all events A_1, \dots, A_n .

$$P(X_1 \in A_1, \dots, X_n \in A_n) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} I_{A_1}(x_1) I_{A_2}(x_2) \dots I_{A_n}(x_n) f(x_1, \dots, x_n) dx_1 \dots dx_n$$

If event $A_i = [a_i, b_i]$, then we can replace the range w/

$$P(X_1 \in A_1, \dots, X_n \in A_n) = \int_{a_1}^{b_1} \int_{a_2}^{b_2} \dots \int_{a_n}^{b_n} f(x_1, \dots, x_n) dx_1 \dots dx_n$$

- The marginal pdf of X_i can be obtained by integrating $f(x_1, \dots, x_n)$ over the full range of all variables except x_i :

$$f_{X_i}(x_i) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f(x_1, \dots, x_n) dx_1 \dots dx_{i-1} dx_{i+1} \dots dx_n$$

- The n RVs X_1, \dots, X_n are independent iff for every A_1, \dots, A_n

$$P(X_1 \in A_1, \dots, X_n \in A_n) = P(X_1 \in A_1) \dots P(X_n \in A_n)$$

This is equivalent to checking that the joint pdf reduces to the product of marginals.

$$f(x_1, \dots, x_n) = f_{X_1}(x_1) \dots f_{X_n}(x_n)$$

For Personal Use Only -bkwk2

Transformation of random vectors

- Consider the transformation $\mathbf{Y} = G(\mathbf{X})$, i.e.

$$\begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} g_1(x_1, \dots, x_n) \\ \vdots \\ g_n(x_1, \dots, x_n) \end{bmatrix}$$

If G is invertible then $\mathbf{X} = G^{-1}(\mathbf{Y})$. Let $H(t) = G^{-1}(t)$, then

$$\begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix} = \begin{bmatrix} h_1(x_1, \dots, x_n) \\ \vdots \\ h_n(x_1, \dots, x_n) \end{bmatrix}$$

Performing a change of variable during integration, we get

$$f_Y(y) = f_X(H(y)) | \det J(y) |$$

where the Jacobian $|\det J(y)|$ is the magnitude of the determinant of the matrix of partial derivatives of $H(t)$

$$J(y) = \begin{bmatrix} \frac{\partial h_1}{\partial y_1} & \cdots & \frac{\partial h_1}{\partial y_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial h_n}{\partial y_1} & \cdots & \frac{\partial h_n}{\partial y_n} \end{bmatrix}$$

Characteristic function.

- The characteristic function of a RV X is defined as

$$\phi_X(t) = E[\exp(itX)], \quad t \in \mathbb{R}.$$

The characteristic function of a random vector $\mathbf{x} = (x_1, \dots, x_n)$ is defined as

$$\phi_{\mathbf{x}}(t) = E[\exp(it^T \mathbf{x})], \quad t \in \mathbb{R}^n$$

- Suppose that \mathbf{X} and \mathbf{Y} are RVs/ random vectors w/ $\phi_{\mathbf{X}}(t) = \phi_{\mathbf{Y}}(t)$ for all $t \in \mathbb{R}^n$, then \mathbf{X} and \mathbf{Y} have the same probability distribution. (since the characteristic function uniquely describes a pdf).

- consider the RV $\mathbf{Y} = \sum_{i=1}^n \mathbf{X}_i$, where \mathbf{X}_i are independent RVs. Its characteristic function $\phi_{\mathbf{Y}}(t)$ is

$$\phi_{\mathbf{Y}}(t) = E[\exp(it^T \mathbf{Y})] = E[\exp(it\mathbf{X}_1) \dots \exp(it\mathbf{X}_n)] = E[\exp(it\mathbf{X}_1)] \dots E[\exp(it\mathbf{X}_n)] = \phi_{\mathbf{X}_1}(t) \dots \phi_{\mathbf{X}_n}(t)$$

i.e. the characteristic function of the sum of independent RVs is the product of the individual characteristic functions.

- we can use the characteristic function to compute moments, $E[X^n]$. consider $\frac{d^n}{dt^n} \phi_{\mathbf{x}}(t)$.

$$\frac{d^n}{dt^n} \phi_{\mathbf{x}}(t) = E\left[\frac{d^n}{dt^n} \exp(it\mathbf{x})\right] = t^n E[i^n \mathbf{x}^n \exp(it\mathbf{x})]$$

$$\text{If we set } t=0, \text{ we get } \frac{d^n}{dt^n} \phi_{\mathbf{x}}(t) \Big|_{t=0} = i^n E[\mathbf{x}^n]. \rightarrow E[\mathbf{x}^n] = \frac{1}{i^n} \frac{d^n}{dt^n} \phi_{\mathbf{x}}(t) \Big|_{t=0}$$

The Gaussian vector.

- For a Gaussian vector $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$, the pdf $f(x_1, \dots, x_n)$ is given by

$$f(x_1, \dots, x_n) = \frac{1}{(2\pi)^{n/2} |\mathbf{S}|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{x} - \mathbf{m})^T \mathbf{S}^{-1} (\mathbf{x} - \mathbf{m})\right)$$

where \mathbf{m} is the mean vector and \mathbf{S} is the covariance matrix.

$$m_i = E[x_i]$$

$$[S]_{ij} = E[(x_i - m_i)(x_j - m_j)]$$

* For independent x_i , the covariance matrix \mathbf{S} is diagonal, i.e. $[S]_{ij} = 0$ for $i \neq j$.

For Personal Use Only -bkwk2

Transformation of the Gaussian vector

- consider a Gaussian vector $\underline{x} = (x_1 \dots x_n)$, where each x_i is independent and $x_i \sim N(0, 1)$.
- let $\underline{\Sigma}$ be an invertible matrix and \underline{m} be a column vector, and define $\underline{y} = \underline{m} + \underline{\Sigma} \underline{x}$,
- Using the change of variable result, where $\underline{y} = g(\underline{x}) = \underline{m} + \underline{\Sigma} \underline{x} \rightarrow h(\underline{y}) = g(\underline{y}) = \underline{\Sigma}^{-1}(\underline{y} - \underline{m}) \rightarrow J(\underline{y}) = \underline{\Sigma}^{-1}$.

$$f_{\underline{y}}(\underline{y}) = f_{\underline{x}}(\underline{\Sigma}^{-1}(\underline{y} - \underline{m})) / |\det \underline{\Sigma}|^{-1}$$

$$\text{where } f_{\underline{x}}(x_1 \dots x_n) = \frac{1}{(2\pi)^{n/2}} \exp(-\frac{1}{2} \underline{x}^T \underline{x}).$$

- we can therefore rewrite $f_{\underline{y}}(\underline{y})$ as

$$f_{\underline{y}}(\underline{y}) = \frac{|\det \underline{\Sigma}|^{-1}}{(2\pi)^{n/2}} \exp\left[-\frac{1}{2}(\underline{y} - \underline{m})^T (\underline{\Sigma}^{-1})^T \underline{\Sigma}^{-1} (\underline{y} - \underline{m})\right]$$

i.e., \underline{y} is a Gaussian vector w/ mean \underline{m} and covariance matrix $\underline{\Sigma}^{-1}$

(note that $|\det \underline{\Sigma}^{-1}| = \frac{1}{|\det \underline{\Sigma}|}$ and $|\det \underline{\Sigma}^{-1}|^{1/2} = \sqrt{|\det \underline{\Sigma}|^{-1}} = \frac{1}{\sqrt{|\det \underline{\Sigma}|}}$).

- In general, an affine transformation (shift + linear stretch) of a Gaussian vector is a Gaussian vector.

We can use this result to generate any Gaussian vector $N(\underline{m}, \underline{\Sigma})$.

↳ Decompose the symmetric matrix $\underline{\Sigma} = \underline{\Sigma}^T$ (use Cholesky decomposition)

↳ Output $\underline{m} + \underline{\Sigma} \underline{x}$, where $\underline{x} = (x_1 \dots x_n)$ and each x_i is independent w/ $x_i \sim N(0, 1)$.

Characteristic function of a Gaussian

- consider a Gaussian RV X , where $X \sim N(\mu, \sigma^2)$. Its characteristic function $\phi_X(t)$ is

$$E[\exp(itX)] = \int_{-\infty}^{\infty} e^{itx} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) dx \stackrel{s=x-\mu}{=} e^{itm} \int_{-\infty}^{\infty} e^{its} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{s^2}{2\sigma^2}\right) ds = e^{itm} e^{-\frac{1}{2}\sigma^2 t^2}$$

$$\text{since } FT[e^{-\frac{1}{2}\sigma^2 t^2}] = \frac{\sqrt{\pi}}{\sigma \sqrt{2}} e^{-\frac{t^2}{2\sigma^2}} = \frac{\sqrt{\pi}}{\sigma} e^{-\frac{t^2}{2\sigma^2}} \rightarrow FT\left[\frac{\sqrt{\pi}}{\sigma} e^{\frac{its}{\sigma^2}}\right] = \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{\sqrt{\pi}}{\sigma} e^{-\frac{s^2}{2\sigma^2}} e^{its} ds = e^{-\frac{t^2}{2\sigma^2}}$$

- consider the transformation of a Gaussian vector $\underline{x} = (x_1 \dots x_n)$, x_i iid $\sim N(0, 1)$ and $\underline{y} = \underline{m} + \underline{\Sigma} \underline{x}$.

Noting that $i \underline{t}^T \underline{y} = i \underline{t}^T \underline{m} + i \underline{t}^T \underline{\Sigma} \underline{x} = i \underline{t}^T \underline{m} + i \underline{t}^T \underline{x}$, where we define $\underline{t} = \underline{\Sigma}^{-1}$

$$\phi_Y(t) = E[\exp(i \underline{t}^T \underline{y})] = E[\exp(i \underline{t}^T \underline{x}) \exp(i \underline{t}^T \underline{m})] = E[\exp(i \underline{t}^T \underline{x}) \exp(i t_1 x_1) \dots \exp(i t_n x_n)]$$

$$= \exp(i \underline{t}^T \underline{m}) E[\exp(i t_1 x_1)] \dots E[\exp(i t_n x_n)] = \exp(i \underline{t}^T \underline{m}) \exp(-\frac{1}{2} t_1^2) \dots \exp(-\frac{1}{2} t_n^2)$$

$$= \exp(i \underline{t}^T \underline{m}) \exp(-\frac{1}{2} \sum_{i=1}^n t_i^2) = \exp(i \underline{t}^T \underline{m}) \exp(-\frac{1}{2} |\underline{t}^T \underline{\Sigma}^{-1}|^2) = \exp(i \underline{t}^T \underline{m}) \exp(-\frac{1}{2} \underline{t}^T \underline{\Sigma} \underline{\Sigma}^{-1} \underline{t})$$

i.e., the characteristic function of a multivariate Gaussian w/ mean \underline{m} and covariance $\underline{\Sigma} = \underline{\Sigma}^T$.

For Personal Use Only -bkwk2

Random process

Random process

Random process

- A discrete time random process is defined as an ensemble of functions

$$\{X_n(\omega)\}, \quad n \in (-\infty, \infty)$$

where ω is a RV w/ pdf $f_w(\omega)$.

- We can consider a generative model for the waveforms. We first draw a random value $\tilde{\omega}$ from the density $f_w(\omega)$. Then, the waveform for this value $\omega = \tilde{\omega}$ is given by $X_n(\tilde{\omega}), \quad n \in (-\infty, \infty)$
- The ensemble is built up by considering all possible values $\tilde{\omega}$ and their corresponding waveforms $X_n(\tilde{\omega})$. The rel. freq. of each waveform $X_n(\omega)$ is determined by $f_w(\omega)$.
- Usually, for simplicity, we can consider the random process as an infinite collection of RVs, e.g.

$$\{X_n\}_{n=-\infty}^{\infty} = \{..., X_{-1}, X_0, X_1, ...\} \quad \text{or} \quad \{X_n\}_{n=0}^{\infty} = \{X_0, X_1, ...\}$$

To completely specify the random process, we must specify

↳ For continuous RVS X_0, X_1, \dots , its n th order joint pdf

$$f_{X_0, X_1, \dots, X_n}(x_0, x_1, \dots, x_n) \quad \forall n \in \mathbb{Z}^+$$

↳ For discrete RVS X_0, X_1, \dots , its n th order joint pmf.

$$P_{X_0, X_1, \dots, X_n}(x_0, x_1, \dots, x_n) \quad \forall n \in \mathbb{Z}^+$$

Correlation and covariance functions

- The autocorrelation function of a random process $\{X_n\}$ is

$$r_{XX}[n, m] = E[X_n X_m]$$

The autocovariance function of a random process $\{X_n\}$ is

$$c_{XX}[n, m] = E[(X_n - M_n)(X_m - M_m)]$$

- The crosscorrelation function b/w random processes $\{X_n\}$ and $\{Y_n\}$ is

$$r_{XY}[n, m] = E[X_n Y_m]$$

The crosscovariance function b/w random processes $\{X_n\}$ and $\{Y_n\}$ is

$$c_{XY}[n, m] = E[(X_n - M_{X,n})(Y_m - M_{Y,m})]$$

strict-sense stationarity (SSS)

- A discrete time random process $\{X_n\}$ is SSS if, for any finite k, m and $(\alpha_0, \alpha_1, \dots, \alpha_k)$,

$$f_{X_0 X_1 \dots X_k}(d_0, d_1, \dots, d_k) = f_{X_0 X_1 \dots X_k}(d_0 + \alpha_0, d_1 + \alpha_1, \dots, d_k + \alpha_k)$$

i.e. it has the same statistical characteristics irrespective of shifts along the time axis.

- + If X_0, X_1, \dots are discrete, we would use the joint pmf instead of the joint pdf to define SSS.

For Personal Use Only -bkwk2

Wide sense stationarity (WSS)

- A discrete-time random process $\{X_n\}$ is WSS if all the following are satisfied
 - (i) $E[X_n] = M_x \quad \forall n$ (const. mean)
 - (ii) $\text{Var}[X_n] < \infty / E[X_n^2] < \infty \quad \forall n$ (finite variance/power)
 - (iii) $r_{xx}[n, m] \rightarrow r_{xx}[m-n] \quad (r_{xx} = r_{xx}(m-n))$
- Two discrete-time random processes $\{X_n\}, \{Y_n\}$ are jointly WSS if the following are satisfied:
 - (i) $E[X_n] = M_x, E[Y_n] = M_y, \forall n$ (const. mean)
 - (ii) $\text{Var}[X_n] < \infty / E[X_n^2] < \infty, \text{Var}[Y_n] < \infty / E[Y_n^2] < \infty, \forall n$ (finite variance/power)
 - (iii) $r_{xx}[n, m] \rightarrow r_{xx}[m-n], r_{yy}[n, m] \rightarrow r_{yy}[m-n], r_{xy}[n, m] \rightarrow r_{xy}[m-n] \quad (r = r(m-n))$
- For WSS $\{X_n\}$, $r_{xx}[k] = r_{xx}[-k]$, i.e. $r_{xx}[k]$ is even

$$r_{xx}[-k] = E[X_n X_{n+k}] = E[X_{n+k} X_n] = r_{xx}[k].$$

and $\max|r_{xx}[k]| = r_{xx}[0]$,

$$0 \leq E[(X_{n+k} - aX_n)^2] = E[X_{n+k}^2] - 2aE[X_{n+k} X_n] + a^2 E[X_n^2] = r_{xx}[0](1+a^2) - r_{xx}[k](2a).$$

$$\therefore r_{xx}[0](1+a^2) \geq 2a r_{xx}[k]$$

setting $a=1 \rightarrow r_{xx}[0] \geq r_{xx}[k]$; setting $a=-1 \rightarrow r_{xx}[0] \geq -r_{xx}[k] \Rightarrow r_{xx}[0] \geq |r_{xx}[k]|$
- For jointly WSS $\{X_n\}, \{Y_n\}$, $r_{xy}[k] = r_{yx}[-k]$

$$r_{xy}[k] = E[X_n Y_{n+k}] = E[Y_{n+k} X_n] = r_{yx}[-k].$$

Power spectral density.

- For WSS $\{X_n\}$, the PSD function is the DTFT of the autocorrelation function $r_{xx}[k]$.

$$S_x(f) = \sum_{k=-\infty}^{\infty} r_{xx}[k] e^{-j2\pi fk}$$

$$r_{xx}[k] = \int_{-\frac{1}{2}}^{\frac{1}{2}} S_x(f) e^{j2\pi fk} df$$
- For jointly WSS $\{X_n\}, \{Y_n\}$, the cross PSD function is the DTFT of the cross correlation function $r_{xy}[k]$

$$S_{xy}(f) = \sum_{k=-\infty}^{\infty} r_{xy}[k] e^{-j2\pi fk}$$

$$r_{xy}[k] = \int_{-\frac{1}{2}}^{\frac{1}{2}} S_{xy}(f) e^{j2\pi fk} df$$
- The PSD $S_x(f)$ is periodic w/ period 1

$$S_x(f+1) = \sum_{k=-\infty}^{\infty} r_{xx}[k] e^{-j2\pi(f+1)k} = \sum_{k=-\infty}^{\infty} r_{xx}[k] e^{j2\pi k} e^{-j2\pi fk} = S_x(f)$$
- The PSD $S_x(f)$ is realvalued and even. It is also continuous if $\sum_{k=-\infty}^{\infty} |r_{xx}[k]| < \infty$.

$$S_x(f) = \sum_{k=-\infty}^{\infty} r_{xx}[k] e^{-j2\pi fk} = \sum_{k=-\infty}^{\infty} r_{xx}[k] \cos(2\pi fk) - j \sum_{k=-\infty}^{\infty} r_{xx}[k] \sin(2\pi fk) = \sum_{k=-\infty}^{\infty} \underbrace{r_{xx}[k]}_{\text{even}} \underbrace{\cos(2\pi fk)}_{\text{even}}$$

$\therefore S_x(f)$ is the sum of realvalued, even and continuous functions $\rightarrow S_x(f)$ is real, even and continuous.
- The PSD $S_x(f)$ can be regarded as the density of power — integrating gives power,
so it must be nonnegative.
- The cross PSD $S_{xy}(f)$ has the symmetry $S_{xy}(f) = S_{yx}^*(f)$

$$S_{xy}(f) = \sum_{k=-\infty}^{\infty} r_{xy}[k] e^{-j2\pi fk} = \sum_{k=-\infty}^{\infty} r_{yx}[-k] e^{-j2\pi fk} = \sum_{k=-\infty}^{\infty} r_{yx}[k] e^{j2\pi fk} = \left[\sum_{k=-\infty}^{\infty} r_{yx}[k] e^{j2\pi fk} \right]^* = S_{yx}^*(f)$$
- The cross PSD $S_{xy}(f)$ can be regarded as a measure of the coherence between the two processes for a certain frequency.

For Personal Use Only -bkwk2

Linear systems and random processes.

- If the i/p $\{x_n\}$ of a LTI system w/ impulse resp. $\{h_n\}$ is WSS, and $\sum_{k=-\infty}^{\infty} |h_k| < \infty$, then the o/p $\{y_n\}$ is also WSS.

(i) Consider the mean $E[y_n]$.

$$E[y_n] = E\left[\sum_{k=-\infty}^{\infty} h_k x_{n-k}\right] = \sum_{k=-\infty}^{\infty} h_k E[x_{n-k}] = E[x_0] \sum_{k=-\infty}^{\infty} h_k \rightarrow \text{const. over all } n.$$

(ii) Consider the autocorrelation $r_{yy}[k]$.

$$\begin{aligned} r_{yy}[k] &= E\left[\left(\sum_{j=-\infty}^{\infty} h_j x_{n-j}\right)\left(\sum_{j=-\infty}^{\infty} h_j x_{n+k-j}\right)\right] = E\left[\sum_{j=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} h_i h_j x_{n-i} x_{n+k-j}\right] \\ &= \sum_{j=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} h_i h_j E[x_{n-i} x_{n+k-j}] = \sum_{j=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} h_i h_j R_{xx}[k+i-j] = \sum_{j=-\infty}^{\infty} h_i \left(\sum_{j=-\infty}^{\infty} h_j R_{xx}[k+j] \right) \\ &= \sum_{j=-\infty}^{\infty} h_i (h_{k+i} R_{xx}[k+i]) = \tilde{h}_k * h_k * R_{xx}[k] \rightarrow \text{only a function of } k. \end{aligned}$$

where \tilde{h}_k is the filter-reversed sequence of h_k , i.e. $\tilde{h}_k = h_{-k}$.

$$(\text{Note } a_n * b_n = \sum_{k=-\infty}^{\infty} a_k b_{n-k}, \quad \tilde{a} * b_n = \sum_{k=-\infty}^{\infty} \tilde{a}_k b_{n-k} = \sum_{k=-\infty}^{\infty} a_{-k} b_{n-k} = \sum_{k=-\infty}^{\infty} a_k b_{n+k}).$$

(iii) Consider the power $E[y_n^2] = r_{yy}[0]$.

$$r_{yy}[0] = \sum_{j=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} h_i h_j R_{xx}[i-j] \leq \sum_{j=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} h_i h_j R_{xx}[0] = R_{xx}[0] \sum_{j=-\infty}^{\infty} h_i h_j \leq R_{xx}[0] \sum_{j=-\infty}^{\infty} |h_i| |\tilde{h}_j|$$

This is finite if $\sum_{j=-\infty}^{\infty} |h_i| < \infty$.

- The cross-correlation $r_{xy}[k]$ and P/P autocorrelation $r_{yy}[k]$ can be expressed as convolutions.

$$r_{xy}[k] = h_k * r_{xx}[k]$$

$$r_{yy}[k] = \tilde{h}_k * h_k * R_{xx}[k].$$

$$\hookrightarrow r_{xy}[k] = E[X_n Y_{n+k}] = E[X_n \sum_{i=-\infty}^{\infty} h_i X_{n+k-i}] = \sum_{i=-\infty}^{\infty} h_i E[X_n X_{n+k-i}] = \sum_{i=-\infty}^{\infty} h_i R_{xx}[k-i] = h_k * R_{xx}[k].$$

\hookrightarrow Derivation for $r_{yy}[k] = \tilde{h}_k * h_k * R_{xx}[k]$ found above.

$$\xrightarrow{\text{Part}(k)} |h_k| \xrightarrow{\text{Part}(k)} |h_{-k}| \xrightarrow{\text{Part}(k)}$$

- Taking DTFT of the convolution form of the o/p autocorrelation $r_{yy}[k]$,

$$\text{DTFT}[r_{yy}[k] = \tilde{h}_k * h_k * R_{xx}[k]] \rightarrow S_y(f) = H^*(f) H(f) S_x(f) = |H(f)|^2 S_x(f)$$

$$\text{where DTFT}[h_k] = \sum_{k=-\infty}^{\infty} \tilde{h}_k e^{j2\pi f k} = \sum_{k=-\infty}^{\infty} h_k e^{j2\pi f k} = \sum_{k=-\infty}^{\infty} h_k e^{j2\pi f k} = \left(\sum_{k=-\infty}^{\infty} h_k e^{j2\pi f k} \right)^* = H^*(f).$$

Physical interpretation of PSD

- For a deterministic signal x_n , the instantaneous power is x_n^2 and the average power is $\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N x_n^2$.

- For a random process $\{x_n\}$, the expected instantaneous power is $E[x_n^2]$ and the expected average power is $E\left[\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N x_n^2\right] = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N E[x_n^2]$.

- If $\{x_n\}$ is WSS, then $E[x_n^2] = R_{xx}[0]$ for all n , so the expected average power is $R_{xx}[0]$.

$$R_{xx}[0] = \int_{-\infty}^{\infty} S_x(f) e^{j2\pi f(0)} df = \int_{-\infty}^{\infty} S_x(f) df$$

i.e. the area under $S_x(f)$ b/w $[f_1/2, f_2/2]$ gives the total power of $\{x_n\}$.

- To find the power of a WSS process $\{x_n\}$ in a particular freq. band, say $[f_1, f_2]$ s.t. $0 < f_1 < f_2 < 1/2$, consider passing $\{x_n\}$ through an ideal bandpass filter w/ TF $H(f) = \begin{cases} 1 & \text{if } f_1 \leq f \leq f_2 \\ 0 & \text{else} \end{cases}$, then find the o/p power $r_{yy}[0]$.

$$r_{yy}[0] = \int_{-1/2}^{1/2} S_x(f) df = \int_{-1/2}^{1/2} |H(f)|^2 S_x(f) df = \int_{f_1}^{f_2} S_x(f) df + \int_{f_2}^{1/2} S_x(f) df = 2 \int_{f_1}^{f_2} S_x(f) df.$$

\rightarrow The power of $\{x_n\}$ in a particular freq. band $f_1 \leq f \leq f_2$ is

$$2 \int_{f_1}^{f_2} S_x(f) df$$

For Personal Use Only -bkwk2

Ergodic random process.

- For an ergodic random process, we can estimate expectations by performing time-averaging on a single sample function (rather than averaging over the whole ensemble).

$$(i) \text{ mean ergodic} \quad M = E[X_N] = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=0}^{N-1} X_n \xrightarrow{\substack{\text{ensemble average} \\ \text{time average}}} \sum_{n=0}^{N-1} X_n$$

$$(ii) \text{ correlation ergodic} \quad C_{xx}[k] = E[X_N X_{N+k}] = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=0}^{N-1} X_n X_{n+k} \approx \frac{1}{N} \sum_{n=0}^{N-1} X_n X_{n+k}$$

- A necessary and sufficient condition for mean ergodicity is given by

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-1} C_{xx}[k] = 0.$$

A simpler sufficient condition for mean ergodicity is that

$$C_{xx}[0] < \infty$$

and

$$\lim_{N \rightarrow \infty} C_{xx}[N] = 0$$

(NTS condition true \rightarrow statement true, NTS condition false \rightarrow statement false)

S condition true \rightarrow statement true, S condition false \rightarrow statement indeterminate)

- Unless otherwise stated, we can always assume that all signals encountered are ergodic.

Markov chains (MC)

Markov chains (MC)

limited memory \rightarrow Markov property.

- A random process $\{X_n\}$ is a MC if for all times $n \geq 0$,

$$P_{X_0 X_1 \dots X_{n-1}}(i_0 i_1 \dots i_{n-1}) = P_{X_n | X_{n-1}}(i_n | i_{n-1})$$

where the RVS X_0, X_1, \dots can take values in the state space $S = \{1, \dots, L\}$.

- The MC is completely defined by the transitional probability matrix Q and the initial distribution of the chain λ .

\hookrightarrow The transitional probability matrix Q has non-negative entries and each row sums to 1.

$$Q = \begin{bmatrix} Q_{1,1} & Q_{1,2} & \dots & Q_{1,L} \\ Q_{2,1} & Q_{2,2} & \dots & Q_{2,L} \\ \vdots & \vdots & \ddots & \vdots \\ Q_{L,1} & Q_{L,2} & \dots & Q_{L,L} \end{bmatrix}$$

where Q_{i_{n-1}, i_n} denotes the conditional pmf of X_n given $(X_0, \dots, X_{n-1}) = (i_0, \dots, i_{n-1})$, i.e.

$$Q_{i_{n-1}, i_n} = P_{X_n | X_{n-1}}(i_n | i_{n-1})$$

\hookrightarrow The initial distribution λ has non-negative entries and they sum to 1.

$$\lambda_i = P_{X_0}(i) \quad i = 1, \dots, L$$

* For MC, we make the notation compact by dropping the subscripts, e.g. $P_{X_0 \dots X_n}(i_0 \dots i_n) \rightarrow P(i_0 \dots i_n)$

For Personal Use Only -bkwk2

Marginals of a Markov chain.

- The marginal of a Markov chain can be expressed as

$$\text{P}_{\Delta}(in) = (\underline{\Delta} \underline{Q}^n)_{in}$$

i.e. $\text{P}_{\Delta}(in) = [\underset{1 \times L}{\Delta}] \left[\underset{L \times L}{Q^n} \right] = [\underset{1 \times L}{\Delta}]$ pick i th entry

- Starting from the joint pmf $P(i_0 \dots i_n)$,

$$P(i_0 \dots i_n) = \frac{P(i_0 \dots i_n)}{P(i_0 \dots i_{n-1})} \cdot \frac{P(i_1 \dots i_n)}{P(i_0 \dots i_{n-1})} \dots \frac{P(i_n)}{P(i_0)} P(i_0) = P(i_0 | i_1 \dots i_{n-1}) P(i_1 | i_0 \dots i_{n-2}) \dots P(i_n | i_0) P(i_0)$$

$$P(i_0 \dots i_n) = P(i_0 | i_1 \dots i_{n-1}) P(i_1 | i_0 \dots i_{n-2}) \dots P(i_n | i_0) = Q_{i_0, i_1} Q_{i_1, i_2} \dots Q_{i_{n-1}, i_n} Q_{i_n, i_0}$$

Marginalising by summing over $i_0 \dots i_{n-1}$,

$$\begin{aligned} P(in) &= \sum_{i_0 \dots i_{n-1} \in S} P(i_0 \dots i_n) = \sum_{i_1 \dots i_{n-1} \in S} Q_{i_0, i_1} \dots Q_{i_{n-2}, i_{n-1}} \sum_{i_0 \in S} Q_{i_0, i_1} \lambda_{i_0} = \sum_{i_1 \dots i_{n-1} \in S} Q_{i_0, i_1} \dots Q_{i_{n-2}, i_{n-1}} (\underline{\Delta} \underline{Q})_{i_1} \\ &= \sum_{i_2 \dots i_{n-1} \in S} Q_{i_1, i_2} \dots Q_{i_{n-1}, i_n} \sum_{i_0 \in S} Q_{i_0, i_1} (\underline{\Delta} \underline{Q})_{i_1} = \sum_{i_2 \dots i_{n-1} \in S} Q_{i_1, i_2} \dots Q_{i_{n-1}, i_n} (\underline{\Delta} \underline{Q})_{i_1} = (\underline{\Delta} \underline{Q})_{i_1}. \end{aligned}$$

dot product b/w $\underline{\Delta}$ and i th col. of \underline{Q} .

Invariant distribution of a Markov chain.

- consider a MC w/ transition probability matrix \underline{Q} w/ state space S . The pmf $\underline{\pi} = (\pi_i : i \in S)$ is invariant for \underline{Q} if $\forall i \in S$,

$$\underline{\pi} \underline{Q} = \sum_{i \in S} \pi_i Q_{ij} = (\underline{\pi} \underline{Q})_j \quad \text{or} \quad \underline{\pi} = \underline{\pi} \underline{Q}.$$

- The MC $(\underline{\pi}, \underline{Q})$ is SSS.

$$\hookrightarrow P(in) = (\underline{\pi} \underline{Q})_{in} = (\underline{\pi} \underline{Q} \underline{Q}^{n-1})_{in} = (\underline{\pi} \underline{Q}^n)_{in} = \dots = \underline{\pi} in$$

$$P(i_0, \dots, i_{n-1}) = Q_{i_0, i_1} \dots Q_{i_{n-2}, i_{n-1}} \pi_{i_{n-1}} \rightarrow \text{we have same joint pmf for any } n \rightarrow \text{SSS}.$$

- The probability distribution $\underline{\pi}$ is known as a stationary distribution.

Ergodic theorem for Markov chain.

- An irreducible MC refers to a chain where all state nodes in S communicate w/ each other.

- This means for any pair of states (i, j) , the MC starting at i will eventually visit j and vice versa, in a finite no. of steps.

- We can identify communicating states by drawing a graph from the transition probability matrix \underline{Q} .

- When the MC $(\underline{\Delta}, \underline{Q})$ w/ stationary distribution $\underline{\pi}$ is irreducible, then for any initial distribution $\underline{\Delta}$,

the sample average converges to $\sum_{i \in S} \pi_i \Delta(i, \cdot)$.

$$\frac{1}{n+1} \sum_{k=0}^n \Delta(k, \cdot) \rightarrow \sum_{i \in S} \pi_i \Delta(i, \cdot)$$

what's
→ DUE

For Personal Use Only -bkwk2

Time series analysis

Time series analysis.

- A time series is a set of observations $y_n, n=0, 1, \dots$ arranged in increasing time.
 - Typically observations are recorded at regular intervals. For start time t_0 and time interval Δ , the n th observation is recorded at $t_0 + n\Delta$.
 - When given a noisy time-series data y_0, y_1, \dots , the workflow for time series analysis is
 - ↳ 1) Select a probability model to represent the data
 - ↳ 2) Estimate model parameters.
 - ↳ 3) Deploy - simulate or forecast.
 - We will consider the probability model $y_n = m_n + s_n + x_n$,
- where m_n is the trend — evolution of mean over time
- s_n is the seasonal component — sinusoids w/ certain period.
- x_n is the residual — zero-mean RV (not necessarily independent over time)

Fitting the trend m_n and seasonal component s_n

- Assuming the trend to be a k th order polynomial,

$$m_n = \alpha_0 + \alpha_1 n + \dots + \alpha_k n^k.$$

We use the data y_0, \dots, y_N to estimate $(\alpha_0, \alpha_1, \dots, \alpha_k)$ via least squares. \rightarrow so we get \hat{m}_n

$$\left(\frac{\partial}{\partial \alpha_0}, \frac{\partial}{\partial \alpha_1}, \dots, \frac{\partial}{\partial \alpha_k} \right) \sum_{n=0}^N (y_n - \alpha_0 - \alpha_1 n - \dots - \alpha_k n^k)^2 = (0, 0, \dots, 0)$$

- Assuming the seasonal component is in the form

$$s_n = A \cos(2\pi f n + \phi) = \alpha_1 \cos(2\pi f n) + \alpha_2 \sin(2\pi f n)$$

We use the data y_0, \dots, y_N to estimate (α_1, α_2) via least squares \rightarrow so we get \hat{s}_n .

$$\left(\frac{\partial}{\partial \alpha_1}, \frac{\partial}{\partial \alpha_2} \right) \sum_{n=0}^N (\hat{y}_n - \alpha_1 \cos(2\pi f n) - \alpha_2 \sin(2\pi f n))^2 = (0, 0)$$

where \hat{y}_n is the detrended data, $\hat{y}_n = y_n - \hat{m}_n$.

* Guessing the trend model will result in a biased estimate \hat{m}_n s.t. $E[\hat{m}_n] \neq 0$.

* Instead of using least squares, the residual is modelled using AR/MA/ARMA models.

For Personal Use Only -bkwk2

Time series models.

Autoregressive (AR) process.

- Let $\{W_n\}_{n \in \mathbb{Z}}$ be a sequence of RVs s.t. $E[W_n] = 0 \forall n$, $E[W_i W_j] = \begin{cases} \sigma^2 & \text{for } i=j \\ 0 & \text{for } i \neq j \end{cases}$.

The AR(p) process $\{X_n\}_{n \in \mathbb{Z}}$ is defined as,

$$X_n = \left(\sum_{i=1}^p a_i X_{n-i} \right) + W_n$$

where a_1, \dots, a_p are const. and p is the order of the process.

- A real time series X_n is modelled as a function of its previous values and random part.

$$X_n = \underbrace{a_1 X_{n-1} + \dots + a_p X_{n-p}}_{\text{predictable part}} + \underbrace{W_n}_{\text{random part}}$$

- Consider the AR(1) process, $X_n = a X_{n-1} + W_n$.

(i) The AR(1) process can be expanded in terms of W_k , $k \leq n$.

$$X_n = a X_{n-1} + W_n = a(a X_{n-2} + W_{n-1}) + W_n = \sum_{k=0}^{\infty} W_{n-k} a^k$$

$$\therefore X_n = \sum_{k=0}^{\infty} W_{n-k} h_k = W_n + h_n, \quad \text{where } h_k \in \mathbb{R}$$

i.e. AR(1) is causal w/ impulse response $\{h_k\}_{k \geq 0}$ [IIR filter]

(ii) The mean $E[X_n]$ is zero, $E[X_n] = 0$.

$$E[X_n] = E\left[\sum_{k=0}^{\infty} W_{n-k} a^k\right] = \sum_{k=0}^{\infty} a^k E[W_{n-k}] = 0$$

$$\begin{aligned} E[X_n] &= E[a X_{n-1} + W_n] \\ E[X_n] &= a E[X_{n-1}] + E[W_n] \\ E[X_n](1-a) &= 0 \rightarrow E[X_n] = 0 \end{aligned}$$

(iii) The variance $\text{Var}[X_n]$ is given by $\text{Var}[X_n] = \frac{\sigma^2}{1-a^2}$

$$\text{Var}[X_n] = E[X_n^2] - E[X_n]^2 = E\left[\left(\sum_{i=0}^{\infty} W_{n-i} a^i\right)^2\right] = E\left[\left(\sum_{i=0}^{\infty} W_{n-i} a^{2i}\right) + \text{cross terms}\right]$$

$$= \sum_{i=0}^{\infty} a^{2i} E[W_{n-i}^2] = \sum_{i=0}^{\infty} a^{2i} \sigma^2 = \frac{\sigma^2}{1-a^2}, \quad \text{provided } |a| < 1. \quad \text{alternatively} \\ \text{Var}[X_n] = \text{Var}[a X_{n-1} + W_n]$$

Note $E[W_i W_j] = 0$ for $i \neq j \rightarrow$ cross-terms are neglected.

$$\text{Var}[X_n] = \text{Var}[a X_{n-1} + W_n] = \text{Var}[a X_{n-1}] + \text{Var}[W_n] = \sigma^2$$

(iv) The autocorrelation $r_{xx}[k]$ is given by $r_{xx}[k] = d k \frac{\sigma^2}{1-a^2}$

$$\begin{aligned} r_{xx}[1] &= E[X_n X_{n-1}] = a E[X_{n-1}^2] + E[X_{n-1} W_n] = a \text{Var}[X_n] = a \frac{\sigma^2}{1-a^2} \\ r_{xx}[0] &= 1 \end{aligned}$$

$$r_{xx}[2] = E[X_n X_{n-2}] = E[X_{n-2} (a X_{n-1} + W_n)] = a E[X_{n-2} X_{n-1}] + E[X_{n-2} W_n] = a r_{xx}[1] = a^2 \frac{\sigma^2}{1-a^2}.$$

By inspection, we see $r_{xx}[k+1] = a r_{xx}[k]$, so $r_{xx}[k] = a^k \frac{\sigma^2}{1-a^2}$, for $k \geq 0$.

(v) The PSD $S_x(f)$ is given by $S_x(f) = \frac{\sigma^2}{1+a^2 - 2a \cos(2\pi f)}$

$$\begin{aligned} S_x(f) &= \sum_{k=-\infty}^{\infty} d k \frac{\sigma^2}{1-a^2} e^{j 2\pi f k} = \frac{\sigma^2}{1-a^2} \sum_{k=-\infty}^{\infty} a^k e^{j 2\pi f k} = \frac{\sigma^2}{1-a^2} \left[\sum_{k=0}^{\infty} a^k e^{j 2\pi f k} + \sum_{k=-\infty}^{-1} a^k e^{j 2\pi f k} - 1 \right] \\ &= \frac{\sigma^2}{1-a^2} \left[\frac{1}{1-a e^{j 2\pi f}} + \frac{1}{1-a e^{-j 2\pi f}} - 1 \right] = \frac{\sigma^2}{1-a^2} \left[\frac{1-a e^{-j 2\pi f}}{1-a e^{j 2\pi f}} + \frac{1-a e^{j 2\pi f}}{1-a e^{-j 2\pi f}} - \frac{1-a^2 + a e^{j 2\pi f} + a e^{-j 2\pi f}}{(1-a e^{j 2\pi f})(1-a e^{-j 2\pi f})} \right] \\ &= \frac{\sigma^2}{(1-a e^{-j 2\pi f})(1-a e^{j 2\pi f})} = \frac{\sigma^2}{1+a^2 - 2a \cos(2\pi f)} = \frac{\sigma^2}{1+a^2 - 2a \cos(2\pi f)} \end{aligned}$$

Alternatively, the PSD $S_x(f)$ can be found using $S_x(f) = |H(f)|^2 S_w(f)$, where $H(f) = DTFT[h_k] = DTFT[a^k]$

$$H(f) = \sum_{k=0}^{\infty} h_k e^{-j 2\pi f k} = \sum_{k=0}^{\infty} a^k e^{-j 2\pi f k} = \frac{1}{1-a e^{-j 2\pi f}}, \quad S_w(f) = \sum_{k=0}^{\infty} h_k |k| e^{j 2\pi f k} = \sigma^2 (r_{ww}[k]) = \begin{cases} \sigma^2 f k & k \neq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\therefore S_x(f) = H(f) H^*(f) S_w(f) = \frac{\sigma^2}{(1-a e^{-j 2\pi f})(1-a e^{j 2\pi f})}$$

* (ii), (iii), (iv) \rightarrow AR(1) process is WSS,

- Higher order AR models can capture finer details and short-term fluctuations.

For Personal Use Only -bkwk2

Moving average process (MA)

- Let $\{w_n\}_{n=-\infty}^{\infty}$ be a sequence of RNS s.t. $E[w_n] = 0 \quad \forall n$, $E[w_i w_j] = \begin{cases} \sigma^2 & \text{for } i=j \\ 0 & \text{for } i \neq j \end{cases}$

The MA(q) process $\{x_n\}_{n=-\infty}^{\infty}$ is defined as

$$x_n = \left(\sum_{i=1}^q b_i w_{n-i} \right) + w_n$$

where b_1, \dots, b_q are const. and q is the order of the process

- A real time series x_n is modelled as a function of its previous values and a random part.

$$x_n = b_1 w_{n-1} + \dots + b_q w_{n-q} + w_n$$

- consider the MA(q) process, $x_n = \left(\sum_{i=1}^q b_i w_{n-i} \right) + w_n$

(i) The MA(q) process can be written in terms of w_k , $k \leq n$.

$$x_n = \sum_{k=0}^{\infty} h_k w_k = w_k * h_k, \text{ where } h_k = \begin{cases} 1 & \text{if } k=0 \\ b_k & \text{if } k \in [1, q] \\ 0 & \text{o/w} \end{cases}$$

i.e. MA(q) is causal w/ impulse response $\{h_k\}_{k=0}^q$ [FIR filter]

(ii) The mean $E[x_n]$ is zero, $E[x_n] = 0$.

$$E[x_n] = E\left[\sum_{k=0}^q h_k w_{n-k}\right] = \sum_{k=0}^q h_k E[w_{n-k}] = 0.$$

(iii) The variance $\text{Var}[x_n]$ is given by $\text{Var}[x_n] = \sigma^2 (1 + b_1^2 + \dots + b_q^2)$

$$\begin{aligned} \text{Var}[x_n] &= E[x_n^2] - E[x_n]^2 = E\left[\left(\sum_{i=1}^q b_i w_{n-i}\right)^2 + w_n^2\right] = E\left[\sum_{i=1}^q b_i^2 w_{n-i}^2 + w_n^2 + \text{cross terms}\right] \\ &= \left(\sum_{i=1}^q b_i^2 E[w_{n-i}^2]\right) + E[w_n^2] = \left(\sum_{i=1}^q b_i^2 \sigma^2\right) + \sigma^2 = \sigma^2 (1 + b_1^2 + \dots + b_q^2) \end{aligned}$$

(iv) The autocorrelation $r_{xx}(k)$ is given by $r_{xx}(k) = \sigma^2 \sum_{i=0}^{\infty} h_i h_{i+k}$

$$r_{xx}(k) = E[x_n x_{n+k}] = E\left[\sum_{i=0}^{\infty} h_i w_i\right] \left[\sum_{j=0}^{\infty} h_{j+k} w_{j+k}\right] = E\left[\sum_{i=0}^{\infty} \sum_{j=0}^{\infty} h_i h_{j+k} w_i w_{j+k}\right] = \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} h_i h_{j+k} E[w_i w_{j+k}]$$

$$E[w_i w_{j+k}] = \begin{cases} \sigma^2 & \text{if } i=j+k \\ 0 & \text{o/w} \end{cases} \text{ so } r_{xx}(k) = \sigma^2 \sum_{i=0}^{\infty} h_i h_{i+k}$$

Note $r_{xx}(k) = 0$ for $k > q$ (since either h_i or h_{i+k} is zero).

(v) The PSD $S_{xx}(f)$ is given by $S_{xx}(f) = (1 + \sum_{k=1}^q b_k e^{-j2\pi fk})(1 + \sum_{k=1}^q b_k e^{j2\pi fk}) \sigma^2$

$$S_{xx}(f) = |H(f)|^2 S_{ww}(f), \text{ where } H(f) = DTFT\{h_k\} = DTFT\{b_1, b_2, \dots, b_q\}$$

$$H(f) = \sum_{k=0}^{\infty} h_k e^{-j2\pi fk} = \sum_{k=0}^q h_k e^{-j2\pi fk} = 1 + \sum_{k=1}^q b_k e^{-j2\pi fk}, S_{ww}(f) = \sum_{k=0}^{\infty} r_{ww}(k) e^{-j2\pi fk} = \sigma^2 (r_{ww}(k)) = \begin{cases} \sigma^2 & \text{if } k=0 \\ 0 & \text{o/w} \end{cases}$$

$$S_{xx}(f) = H(f) H^*(f) S_{ww}(f) = (1 + \sum_{k=1}^q b_k e^{-j2\pi fk})(1 + \sum_{k=1}^q b_k e^{j2\pi fk}) \sigma^2$$

* (i),(ii),(iii),(iv) \rightarrow MA(q) process is WSS

Autoregressive-moving average process (ARMA).

- Let $\{w_n\}_{n=-\infty}^{\infty}$ be a sequence of RNS s.t. $E[w_n] = 0 \quad \forall n$, $E[w_i w_j] = \begin{cases} \sigma^2 & \text{for } i=j \\ 0 & \text{for } i \neq j \end{cases}$

The ARMA(p,q) process $\{x_n\}_{n=-\infty}^{\infty}$ is defined as.

$$x_n = \sum_{i=1}^p a_i x_{n-i} + w_n + \sum_{i=1}^q b_i w_{n-i}$$

where a_1, \dots, a_p , b_1, \dots, b_q are const. and p/q are the order of the IIR/FIR respectively.

- The ARMA(p,q) process can be written in terms of w_k , $k \leq n$.

$$x_n = \sum_{k=0}^{\infty} w_k * h_k = w_k * h_k, \text{ where } h_k = \begin{cases} 1 & \text{if } k=0 \\ a_k + b_k & \text{if } k \in [1, q] \\ 0 & \text{o/w} \end{cases}$$

i.e. ARMA(p,q) is causal w/ impulse response $\{h_k\}_{k=0}^q$ [IIR filter]

- The ARMA(p,q) process is WSS.

For Personal Use Only -bkwk2

AR, MA, ARMA processes as causal filters.

- AR, MA, ARMA process can act as causal filters w/ impulse response $\{h_n\}$, subject to the i/p $\{w_n\}_{n=0}^{\infty}$, where $E[w_n] = 0$ and $E[w_n w_j] = \begin{cases} \sigma^2 & \text{for } j=0 \\ 0 & \text{for } j \neq 0 \end{cases}$.
- consider the random process $\{W_n\}_{n=0}^{\infty}$.
 - (i) Mean $E[W_n]$: $E[W_n] = 0$, i.e. const. W_n .
 - (ii) Power $E[W_n^2]$: $E[W_n^2] = \sigma^2 < \infty$ i.e. finite power
 - (iii) Autocorrelation $r_{WW}[k]$: $r_{WW}[k] = E[W_n W_{n+k}] = 0 = h_n(k)$, i.e. only depends on k
 $\rightarrow \{W_n\}_{n=0}^{\infty}$ is WSS.
- provided for the causal filter, its impulse response $\{h_n\}$ satisfies $\sum h_n < \infty$, then given a WSS i/p $\{x_n\}$, its o/p $\{y_n\}$ is also WSS \rightarrow AR, MA, ARMA processes must be WSS.

Random phase cosine (harmonic process)

- For the random phase cosine, the amplitude a and frequency f_0 are known, but the phase ϕ is unknown and random (could correspond to an intrinsic delay). It has the form

$$X_n = a \cos(2\pi f_0 n + \phi)$$

where a, f_0 are const. and $\phi \sim \text{unif}[0, 2\pi]$, i.e. its pdf is $P_\phi(\phi) = \begin{cases} 1/2\pi & \text{if } \phi \in [0, 2\pi) \\ 0 & \text{otherwise} \end{cases}$

+ To generate $\{x_n\}$, we need to first sample ϕ to get $\phi = \hat{\phi}$, then set $X_n = \cos(2\pi f_0 n + \hat{\phi})$, for $n=0, \dots$

- consider the harmonic process, $X_n = a \cos(2\pi f_0 n + \phi)$

(i) The mean $E[X_n]$ is zero, $E[X_n] = 0$

$$E[X_n] = E[a \cos(2\pi f_0 n + \phi)] = E[a \cos(2\pi f_0 n) \cos(\phi) - a \sin(2\pi f_0 n) \sin(\phi)] = a \cos(2\pi f_0 n) E[\cos(\phi)] - a \sin(2\pi f_0 n) E[\sin(\phi)] = 0$$

(since $E[\cos(\phi)] = E[\sin(\phi)] = 0$ for uniform distribution ϕ).

(ii) The variance $\text{Var}[X_n]$ is given by $\text{Var}[X_n] = \frac{1}{2} a^2$.

$$\text{Var}[X_n] = E[X_n^2] - E[X_n]^2 = E[a^2 \cos^2(2\pi f_0 n + \phi)] = E\left[\frac{a^2}{2}(1 + \cos(4\pi f_0 n + 2\phi))\right] = \frac{a^2}{2} + \frac{a^2}{2} E[\cos(4\pi f_0 n + 2\phi)] = \frac{a^2}{2}$$

(iii) The autocorrelation $r_{XX}[k]$ is given by $r_{XX}[k] = \frac{1}{2} a^2 \cos(2\pi f_0 k)$

$$r_{XX}[k] = E[X_n X_{n+k}] = E[(a \cos(2\pi f_0 n + \phi)) (a \cos(2\pi f_0 (n+k) + \phi))] = E\left[\frac{a^2}{2} (\cos(2\pi f_0 (2n+k) + 2\phi) + \cos(2\pi f_0 k))\right] = \cancel{\frac{a^2}{2} E[\cos(2\pi f_0 (2n+k) + 2\phi)]} + \frac{a^2}{2} \cos(2\pi f_0 k) = \frac{a^2}{2} \cos(2\pi f_0 k)$$

(iv) The PSD $S_X(f)$ is given by $S_X(f) = \frac{a^2 \pi}{2} (\delta(2\pi(f-f_0)) + \delta(2\pi(f+f_0)))$

$$S_X(f) = \sum_{n=-\infty}^{\infty} r_{XX}[k] e^{-j2\pi f n} = \sum_{k=-\infty}^{\infty} \frac{a^2}{2} \cos(2\pi f_0 k) e^{-j2\pi f n} = \frac{a^2}{2} \sum_{k=-\infty}^{\infty} (e^{-j2\pi(f-f_0)k} + e^{j2\pi(f+f_0)k}).$$

$$= \frac{a^2}{2} [2\pi \delta(2\pi(f-f_0)) + 2\pi \delta(2\pi(f+f_0))] = \frac{a^2 \pi}{2} (\delta(2\pi(f-f_0)) + \delta(2\pi(f+f_0))) \quad \begin{matrix} \text{over } f \in [-f_0, f_0], \\ \text{and periodic w/} \\ \text{period of } 1. \end{matrix}$$

* (i), (ii), (iii) \rightarrow harmonic process is WSS.

- consider the inversion formula for PSD, $r_{XX}[k] = \int_{-1/2}^{1/2} S_X(f) e^{j2\pi f k} df = \int_{-1/2}^{1/2} S_X(f) \cos(2\pi f k) df + \int_{-1/2}^{1/2} S_X(f) \sin(2\pi f k) df$

$$\therefore r_{XX}[k] = 2 \int_{-1/2}^{1/2} S_X(f) \cos(2\pi f k) df \quad \begin{matrix} \text{divide into } K \\ \text{intervals, with } V_{2k} \end{matrix} \quad 2 \sum_{i=1}^K S_X(f_i) \cos(2\pi f_i k) \frac{1}{V_{2k}} = \frac{1}{K} \sum_{i=1}^K S_X(f_i) \cos(2\pi f_i k)$$

i.e. we have a sum of random phase cosines w/ freq. f_i and amplitude a s.t. $\frac{a^2}{2} = \frac{S_X(f)}{K} \rightarrow a = \sqrt{\frac{2S_X(f)}{K}}$

$$\therefore X_n = \sum_{i=1}^K \left[\frac{2S_X(f_i)}{K} \right] \cos(2\pi f_i n + \phi_i) \quad \text{where } \phi_i \sim \text{unif}[0, 2\pi]$$

For Personal Use Only -bkwk2

white noise process,

- A WSS process $\{x_n\}$ is termed white noise if

$$c_{xx}[k] = \sigma^2 \delta[k]$$

where $\delta[k]$ is the discrete impulse function, $\delta[k] = \begin{cases} 1 & \text{if } k=0 \\ 0 & \text{o/w.} \end{cases}$

- The PSD $S_x(f)$ for white noise is given by $S_x(f) = \sigma^2 + M^2 2\pi \delta(f)$

$$S_x(f) = \sum_{k=-\infty}^{\infty} c_{xx}[k] e^{j2\pi fk} = \sum_{k=-\infty}^{\infty} (c_{xx}[k] + M^2) e^{-j2\pi fk} = \sum_{k=-\infty}^{\infty} \sigma^2 \delta[k] e^{-j2\pi fk} + \sum_{k=-\infty}^{\infty} M^2 e^{-j2\pi fk} = \sigma^2 + M^2 2\pi \delta(f).$$

for zero mean white noise, $M=0$, so $S_x(f) = \sigma^2$, i.e. flat across all freq. ✓ white noise for ARMA is zero mean white noise

- If the values x_n are drawn independently from a Gaussian distribution $N(0, \sigma^2)$, we have

white Gaussian noise (WGN), the n th order pdf for WGN process is

$$f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n) = \prod_{i=1}^n N(x_i | 0, \sigma^2) \quad \begin{matrix} \leftarrow \text{simple product of} \\ \text{x}_i \text{ are independent} \end{matrix}$$

where $N(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu)^2}{2\sigma^2})$ is the univariate normal pdf.

- The WGN process is SSS, since

$$f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n) = \prod_{i=1}^n N(x_i | 0, \sigma^2) = f_{X_1 \perp X_2 \perp \dots \perp X_n}(x_1, x_2, \dots, x_n)$$

optimal filtering and signal detection

Optimal filtering

PSD and transfer function in terms of S_x .

- The PSD and TF can be expressed in terms of the normalised freq. ω , (in radians per sample).

where $S_x = \frac{1}{T} \int_{-T/2}^{T/2} S_x(t) dt$, and T is the sampling interval of the discrete-time process.

$$S_x(e^{j\omega}) = \sum_{k=-\infty}^{\infty} S_x(kT) e^{-jk\omega}$$

$$S_{xy}(e^{j\omega}) = \sum_{k=-\infty}^{\infty} S_{xy}(kT) e^{-jk\omega}$$

$$P_{x,y}(e^{j\omega}) = r_{xy}(0) = \frac{1}{2\pi} \int_{-\pi}^{\pi} S_{xy}(e^{j\omega}) d\omega$$

$$r_{xx}(kT) = \frac{1}{2\pi} \int_{-\pi}^{\pi} S_x(e^{j\omega}) e^{jk\omega} d\omega$$

$$r_{yy}(kT) = \frac{1}{2\pi} \int_{-\pi}^{\pi} S_y(e^{j\omega}) e^{jk\omega} d\omega$$

$$P_{y,y}(e^{j\omega}) = \frac{1}{2\pi} \int_{-\pi}^{\pi} S_y(e^{j\omega}) d\omega$$

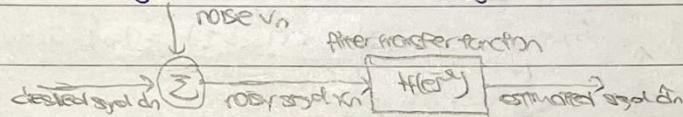
+ In this section, we explicitly write the PSD/TF as a function of $e^{j\omega}$ to stress that it is always periodic w/ a period of 2π . \rightarrow looks like a Z-transform, w/ $z = e^{j\omega}$

Optimal filtering.

- Optimal filtering is where we design filters that are optimally adapted to the statistical characteristics of a random process.

- The Wiener filter can optimally estimate a desired signal d_n given the noisy observations x_n , w/ some assumptions about the statistics of the random signal and noise processes.

- The general Wiener filtering problem can be expressed as a block diagram.



i.e., we have $x_n = d_n + v_n$ and $\hat{d}_n = x_n + h_n$.

- We define the error signal e_n as

$$e_n = d_n - \hat{d}_n = d_n - x_n + h_n$$

and the mean-squared error (MSE) J as

$$J = E[e_n^2]$$

where the expectation is wrt the random signal d_n and random noise v_n .

- The Wiener filter minimises the MSE J wrt the filter coefficients $\{h_k\}$.

- The Wiener filter assumes the following.

(i) $\{x_n\}$ and $\{d_n\}$ are jointly WSS, i.e. $r_{xx}(m) \rightarrow r_{xx}(\omega)$; $r_{dd}(m) \rightarrow r_{dd}(\omega)$, $r_{xd}(m) \rightarrow r_{xd}(\omega)$

(ii) $\{d_n\}$ and $\{v_n\}$ have zero mean, i.e. $E[d_n] = 0$, $E[v_n] = 0$

+ The Wiener filter is only the optimal linear estimator for stationary signals

For Personal Use Only -bkwk2

The general Wiener filter.

- In the general case, we can filter the observed signal x_n w/ an infinite-dimensional filter, having a non-causal impulse response $\{h_p\}$. Our estimate of the desired signal \hat{d}_n is therefore

$$\hat{d}_n = \sum_{p=-\infty}^{\infty} h_p x_{n-p}.$$

- We want to minimise the MSE J wrt the impulse response values h_p .

A sufficient condition for a min. is that simultaneously $\forall q \in (-\infty, \infty)$,

$$\frac{\partial J}{\partial h_q} = \frac{\partial E[\epsilon_n^2]}{\partial h_q} = E\left[\frac{\partial \epsilon_n^2}{\partial h_q}\right] = E[2\epsilon_n \delta_{n,q}] = 0.$$

Noting that $\frac{\partial \epsilon_n}{\partial h_q} = \frac{\partial}{\partial h_q} [d_n - \sum_{p=-\infty}^{\infty} h_p x_{n-p}] = x_{n-q}$, we get the condition

$$E[2\epsilon_n \delta_{n,q}] = E[2\epsilon_n x_{n-q}] = 0 \rightarrow E[\epsilon_n x_{n-q}] = 0, \forall q \in (-\infty, \infty)$$

Substituting ϵ_n in $E[\epsilon_n x_{n-q}]$ gives

$$E[\epsilon_n x_{n-q}] = E[(d_n - \sum_{p=-\infty}^{\infty} h_p x_{n-p}) x_{n-q}] = E[d_n x_{n-q}] - \sum_{p=-\infty}^{\infty} h_p E[x_{n-p} x_{n-q}] = r_{dd}[q] - \sum_{p=-\infty}^{\infty} h_p r_{xx}[q-p]$$

$$r_{dd}[q] - \sum_{p=-\infty}^{\infty} h_p r_{xx}[q-p] = 0 \rightarrow \sum_{p=-\infty}^{\infty} h_p r_{xx}[q-p] = r_{dd}[q], \forall q \in (-\infty, \infty)$$

which are collectively known as the Wiener-Hopf eqns.

- The simplest way to solve for the infinite no. of unknowns h_p in the Wiener-Hopf eqns. is to transform everything into the frequency domain.

Rewriting the Wiener-Hopf eqns. as a convolution,

$$h_q * r_{xx}[q] = r_{dd}[q] \quad \forall q \in (-\infty, \infty).$$

Taking the DTFT on both sides,

$$H(e^{j\omega}) S_x(e^{j\omega}) = S_d(e^{j\omega}) \rightarrow H(e^{j\omega}) = \frac{S_d(e^{j\omega})}{S_x(e^{j\omega})}$$

- The MSE of the optimal filter, J_{min} , can be found using (i) time or (ii) frequency domain

$$(i). J = E[\epsilon_n^2] = E[\epsilon_n(d_n - \sum_{p=-\infty}^{\infty} h_p x_{n-p})] = E[d_n \epsilon_n] - \sum_{p=-\infty}^{\infty} h_p E[\epsilon_n x_{n-p}]$$

For the optimal filter, the orthogonality principle states $E[\epsilon_n x_{n-p}] = 0 \forall p \in (-\infty, \infty)$, so

$$J_{min} = E[d_n \epsilon_n] = E[(d_n - \sum_{p=-\infty}^{\infty} h_p x_{n-p}) d_n] = E[d_n^2] - \sum_{p=-\infty}^{\infty} h_p E[d_n x_{n-p}] = r_{dd}[0] - \sum_{p=-\infty}^{\infty} h_p r_{dd}[p]$$

$$(ii) \epsilon_n x_{n-k} = (d_n - \sum_{p=-\infty}^{\infty} h_p x_{n-p})(d_{n-k} - \sum_{j=-\infty}^{\infty} h_j x_{n-j}) = d_n d_{n-k} - d_n \sum_{j=-\infty}^{\infty} h_j x_{n+j} - d_{n-k} \sum_{i=-\infty}^{\infty} h_i x_{n-i} + \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} h_i h_j x_{n-i} x_{n+j}$$

$$\begin{aligned} r_{dd}[k] &= E[\epsilon_n x_{n-k}] = r_{dd}[k] - \sum_{j=-\infty}^{\infty} h_j r_{xx}[k-j] - \sum_{i=-\infty}^{\infty} h_i r_{xx}[-i+k] + \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} h_i h_j r_{xx}[k+i-j] \\ &= r_{dd}[k] - h_k + r_{xx}[k] - \tilde{h}_k r_{dd}[k] + h_k + \tilde{h}_k * r_{xx}[k] \end{aligned}$$

$$S_d(e^{j\omega}) = S_d(e^{j\omega}) - H(e^{j\omega}) S_{xx}(e^{j\omega}) - H^*(e^{j\omega}) S_{dd}(e^{j\omega}) + H(e^{j\omega}) H^*(e^{j\omega}) S_x(e^{j\omega})$$

$$J = E[\epsilon_n^2] = r_{dd}[0] = \frac{1}{2\pi} \int_{-\pi}^{\pi} S_d(e^{j\omega}) d\omega = \frac{1}{2\pi} \int_{-\pi}^{\pi} [S_d(e^{j\omega}) - H(e^{j\omega}) S_{xx}(e^{j\omega}) - H^*(e^{j\omega}) S_{dd}(e^{j\omega}) + H(e^{j\omega})^2 S_x(e^{j\omega})] d\omega$$

For the optimal filter, solving the Wiener-Hopf eqns. since $H(e^{j\omega}) = \frac{S_d(e^{j\omega})}{S_x(e^{j\omega})}$.

$$H(e^{j\omega}) S_{dd}(e^{j\omega}) = \frac{S_d(e^{j\omega}) S_{dd}(e^{j\omega})}{S_x(e^{j\omega})} = \frac{|S_d(e^{j\omega})|^2}{S_x(e^{j\omega})}, \therefore |H(e^{j\omega})|^2 S_{dd}(e^{j\omega}) = \frac{|S_d(e^{j\omega})|^2}{|S_x(e^{j\omega})|^2} = \frac{(S_d(e^{j\omega}))^2}{(S_x(e^{j\omega}))^2}$$

$$\rightarrow J_{min} = \frac{1}{2\pi} \int_{-\pi}^{\pi} [S_d(e^{j\omega}) - H(e^{j\omega}) S_{xx}(e^{j\omega})] d\omega$$

* If the random process is ergodic, we estimate $r_{dd}[k]$, $r_{xx}[k]$ from the average

* In general, the optimal filter yielded is noncausal \rightarrow not implementable in practice

For Personal Use Only -bkwk2

THE FIR Wiener filter

- For the FIR Wiener filter, it has a causal p -th order finite impulse response $\{h_p\}_{p=0}^P$. Our estimate of the desired signal, d_n , is therefore

$$\hat{d}_n = \sum_{p=0}^P h_p x_{n-p}$$

- As before, we want to minimize the MSE J wrt the impulse response values h_q .

Orthogonality principle: $E[\epsilon_n x_m] = 0, \forall m \in [0, P]$

Wiener-Hopf eqns: $\sum_{p=0}^P h_p r_{xd}[q-p] = r_{xd}[q], \forall q \in [0, P]$

- The Wiener-Hopf eqns are now a finite set of simultaneous eqns that can be solved in the time domain. The eqns may be written in matrix form as

$$R_x h = r_{xd}$$

where $h = \begin{bmatrix} h_0 \\ h_1 \\ \vdots \\ h_P \end{bmatrix}, R_{xd} = \begin{bmatrix} r_{xd}[0] & & & \\ r_{xd}[1] & \ddots & & \\ \vdots & & \ddots & \\ r_{xd}[P] & & & r_{xd}[P] \end{bmatrix}, R_x = \begin{bmatrix} r_{xx}[0] & r_{xx}[1] & \cdots & r_{xx}[P] \\ r_{xx}[1] & r_{xx}[0] & \cdots & r_{xx}[P-1] \\ \vdots & \vdots & \ddots & \vdots \\ r_{xx}[P] & r_{xx}[P-1] & \cdots & r_{xx}[0] \end{bmatrix}$

(Note the correlation matrix R_x is symmetric and has zeros diagonals).

The coefficient vector can be found by matrix inversion,

$$h = R_x^{-1} r_{xd}$$

- As before, the MSE of the optimal filter J_{min} is given by

$$J_{min} = G_d[0] - \sum_{p=0}^P h_p G_d[p] = G_d[0] - \underline{r_{xd}}^\top \underline{h} = G_d[0] - \underline{r_{xd}}^\top \underline{R_x^{-1} r_{xd}}$$

Uncorrelated signal and noise processes.

- Typically, the environmental noise v_n is independent of the desired signal $d_n \rightarrow$ uncorrelated, i.e.

$$r_{dv}[k] = E[d_n v_{nk}] \quad k \in (-\infty, \infty)$$

- In this situation, the correlation functions $r_{xd}[k]$ and $r_{dv}[k]$ can be simplified.

(i) $r_{xd}[k] = E[x_n d_{n+k}] = E[(d_n + v_n)d_{n+k}] = E[d_n d_{n+k}] + E[v_n d_{n+k}] = r_{dd}[k]$.

$$r_{xd}[k] = r_{dd}[k] \leftrightarrow S_{xd}(e^{j\omega}) = S_d(e^{j\omega})$$

(ii) $r_{dv}[k] = E[x_n v_{n+k}] = E[(d_n + v_n)(d_{n+k} + v_{n+k})] = E[d_n d_{n+k}] + E[d_n v_{n+k}] + E[v_n d_{n+k}] + E[v_n v_{n+k}] = r_{dd}[k] + r_{nv}[k]$

$$r_{dv}[k] = r_{dd}[k] + r_{nv}[k] \leftrightarrow S_{dv}(e^{j\omega}) = S_d(e^{j\omega}) + S_v(e^{j\omega})$$

- For the general Wiener filter, the optimal filter TF $H(e^{j\omega})$ becomes.

$$H(e^{j\omega}) = \frac{S_d(e^{j\omega})}{S_d(e^{j\omega}) + S_v(e^{j\omega})} = \frac{S_d(e^{j\omega})}{1 + \frac{S_v(e^{j\omega})}{S_d(e^{j\omega})}}$$

where $\rho(\omega) = \frac{S_v(e^{j\omega})}{S_d(e^{j\omega})}$ is the (freq.-dependent) SNR power ratio.

- The MSE of the optimal filter, J_{min} is now reduced to,

$$J_{min} = \frac{1}{2\pi} \int_{-\pi}^{\pi} [P_d(e^{j\omega}) + H(e^{j\omega}) S_v(e^{j\omega})] d\omega = \frac{1}{2\pi} \int_{-\pi}^{\pi} S_d(e^{j\omega}) \left(1 + \frac{1}{1 + \rho(\omega)} \right) d\omega = \frac{1}{2\pi} \int_{-\pi}^{\pi} S_d(e^{j\omega}) \frac{1}{1 + \rho(\omega)} d\omega$$

- In practice, the correlation functions are estimated as (assuming S_{xd} is ergodic).

$$r_{xd}[k] \approx \frac{1}{N} \sum_{n=0}^{N-1} x_n d_{n+k}, \quad r_{dd}[k] = r_{xx}[k] - r_{vv}[k] = \begin{cases} r_{xx}[k] - \sigma_v^2 & \text{if } k=0 \\ r_{xx}[k] & \text{o.w.} \end{cases}$$

* We should make sure $r_{xx}[k]$ and $r_{vv}[k]$ are valid autocorrelation sequences, since they are only empirical estimates.

For Personal Use Only -bkwk2

Validity of autocorrelation sequence

- For the autocorrelation $R_{xx}[k]$ to form a valid sequence for constructing the correlation matrix \underline{R}_x , the necessary condition is the resulting \underline{R}_x is positive-semidefinite, i.e.

$$\underline{a}^T \underline{R}_x \underline{a} \geq 0$$

for any length $P+1$ vector \underline{a} .

(since \underline{R}_x is symmetric, this means its eigenvalues λ_i must satisfy $\lambda_i \geq 0 \ \forall i$)

- consider the vector $\underline{x}_n = [x_n \ x_{n-1} \ \dots \ x_{n-P}]^T$ for real-valued x_n , $\underline{a}^T \underline{x}_n = \underline{x}_n^T \underline{a}$

$$0 \leq (\underline{a}^T \underline{x}_n)^2 = (\underline{a}^T \underline{x}_n)(\underline{x}_n^T \underline{a}) = \underline{a}^T (\underline{x}_n \underline{x}_n^T) \underline{a}$$

The (i,j) th element of $\underline{x}_n \underline{x}_n^T$ is $x_{n-i+1} x_{n+j}$ and $E[x_{n-i+1} x_{n+j}] = R_{xx}[i-j] \rightarrow E[\underline{x}_n \underline{x}_n^T] = \underline{R}_x$.

$$\therefore 0 \leq E[\underline{a}^T (\underline{x}_n \underline{x}_n^T) \underline{a}] = \underline{a}^T E[(\underline{x}_n \underline{x}_n^T)] \underline{a} = \underline{a}^T \underline{R}_x \underline{a} \rightarrow \underline{R}_x \text{ must be positive-semidefinite.}$$

Extending the Wiener filter

- To extend the Wiener filter beyond the regular noise reduction case, we simply replace the desired signal d w/ what we want to predict/estimate, then rederive the filter.
- In general, the FIR formula

$$\underline{b} = \underline{R}_x^{-1} \underline{r}_{xd}$$

applies for any form of the (unobserved) desired signal d and the (observed) "noisy" signal x , provided we can calculate the necessary correlation functions.

- Some standard examples include:

↳ Prediction of a signal, $x_n = u_n + v_n$, $d_n = u_{n+p}$

↳ Smoothing of a signal, $x_n = u_n + v_n$, $d_n = u_{n-p}$

↳ Deconvolution, $x_n = h_n * u_n + v_n = \sum_{q=0}^Q h_q u_{n-q} + v_n$, $d_n = u_n$.

→ In all cases, the Wiener filter estimate takes the same form as before

$$\hat{d}_n = \sum_{p=0}^P h_p x_{n-p}$$

and the error is defined as before

$$e_n = d_n - \hat{d}_n$$

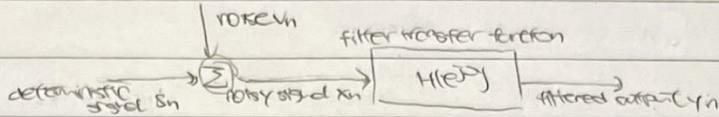
Note that only d is different, and we need to rederive an expression for $R_{dd}[k]/S_{dd}(e^{j\omega k})$.

For Personal Use Only -bkwk2

Signal detection

Signal detection

- The matched filter is an optimal FIR filter that can detect a deterministic signal $s_n, n=0, \dots, N-1$, from the noisy observations x_n .
- The general signal detection problem can be expressed as a block diagram



$$\text{i.e., we have } x_n = s_n + v_n \text{ and } y_n = x_n * h_n$$

- In vector form, we have

$$\underline{x} = \underline{s} + \underline{v}$$

$$\text{where } \underline{s} = [s_0, s_1, \dots, s_{N-1}]^T \text{ and } \underline{x} = [x_0, x_1, \dots, x_{N-1}]^T.$$

- Using a $(N+1)$ th order FIR filter, the o/p of the filter at time $N-1$ is given by

$$y_{N-1} = \sum_{m=0}^{N-1} h_m x_{N-1-m} = \underline{h}^T \underline{x} = \underline{h}^T (\underline{s} + \underline{v}) = \underline{h}^T \underline{s} + \underline{h}^T \underline{v} = \sum_{m=0}^{N-1} h_m s_{N-1-m} + \sum_{m=0}^{N-1} h_m v_{N-1-m} = y_{N-1}^s + y_{N-1}^v$$

where y_{N-1}^s is the o/p from the signal only part ; y_{N-1}^v is the o/p from the noise only part .

- We define the o/p SNR as

$$\text{SNR} = \frac{E[y_{N-1}^s]}{E[y_{N-1}^v]} = \frac{E[\underline{h}^T \underline{s}]}{E[\underline{h}^T \underline{v}]} = \frac{\underline{h}^T \underline{s}^T}{\underline{h}^T \underline{v}^T}$$

- The matched filter maximizes the o/p SNR, w.r.t the filter coefficients $\{h_m\}$ to give the best possible chance of detecting the signal s_n .

Signal output energy.

- The signal o/p energy is given by $y_{N-1}^s = \underline{h}^T \underline{s}^T = (\underline{h}^T \underline{s})(\underline{s}^T \underline{h}) = \underline{h}^T (\underline{s} \underline{s}^T) \underline{h}$.

- Consider the matrix $\underline{M} = \underline{s} \underline{s}^T$. We want to find its eigenvectors \underline{u}_i and eigenvalues λ_i .

$$\underline{u} = \underline{s} : \quad \underline{M} \underline{u} = (\underline{s} \underline{s}^T) \underline{u} = \underline{s} (\underline{s}^T \underline{s}) \underline{u} = (\underline{s}^T \underline{s}) \underline{u}$$

i.e. $\underline{u}_0 = \frac{\underline{s}}{\|\underline{s}\|}$ is an eigenvector and $\lambda_0 = \frac{\underline{s}^T \underline{s}}{\|\underline{s}\|^2}$ is the corresponding eigenvalue .

$$\underline{u} = \{\underline{u}_i : \underline{s}^T \underline{u}_i = 0\} : \quad \underline{M} \underline{u}_i = (\underline{s} \underline{s}^T) \underline{u}_i = \underline{s} (\underline{s}^T \underline{u}_i) = 0$$

i.e. $\underline{u}_i = \underline{u}'$ is an eigenvector w/ $\lambda = 0$. (There are $N-1$ orthogonal vectors orthogonal to $\underline{s}, \underline{u}_1, \underline{u}_2, \dots, \underline{u}_{N-1}$)

- We can represent any filter coefficient vector \underline{h} as a linear combination of these eigenvectors.

$$\underline{h} = \alpha \underline{u}_0 + \beta \underline{u}_1 + \gamma \underline{u}_2 + \dots + \underline{u}_{N-1}$$

$$\rightarrow \underline{M} \underline{h} = \underline{M}(\alpha \underline{u}_0 + \beta \underline{u}_1 + \gamma \underline{u}_2 + \dots + \underline{u}_{N-1}) = \alpha \underline{M} \underline{u}_0 + \beta \underline{M} \underline{u}_1 + \gamma \underline{M} \underline{u}_2 + \dots + \underline{M} \underline{u}_{N-1} = \alpha \underline{u}_0^T \underline{s} \underline{s}^T \underline{u}_0$$

- The signal o/p energy is therefore .

$$\underline{h}^T \underline{s} \underline{s}^T \underline{h} = \underline{h}^T \underline{M} \underline{h} = \underline{h}^T (\alpha \underline{u}_0^T \underline{s} \underline{s}^T \underline{u}_0) = (\alpha \underline{u}_0^T \underline{s}) \underline{h}^T \underline{u}_0 = (\alpha \underline{u}_0^T \underline{s}) (\alpha u_0 + \beta u_1 + \gamma u_2 + \dots + u_{N-1})^T u_0$$

$$\therefore \underline{u}_0^T \underline{s}^2 = \underline{h}^T \underline{u}_0^2 = \alpha^2 \underline{u}_0^T \underline{s}$$

since $\underline{u}_0^T \underline{u}_0 = 1$ and $\underline{u}_0^T \underline{u}' = 0$ (orthogonality).

For Personal Use Only -bkwk2

Noise output energy

- The expected noise o/p energy is given by $E[Y_{\text{out}}^2] = E[h^T Y^2] = E[h^T h] = h^T E[Y^2] h$

- considering the case where the noise is white, for time indices $i=0, \dots, N-1$ and $j=0, \dots, N-1$,

$$E[v_i v_j] = \begin{cases} \sigma_v^2 & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}$$

$\therefore E[Y^2] = \sigma_v^2 I$, where I is the $N \times N$ identity matrix.

- The expected noise o/p energy is therefore

$$E[Y_{\text{out}}^2] = h^T E[Y^2] h = \sigma_v^2 h^T h$$

Expanding h in terms of the eigenvectors of $S = \frac{Y^T Y}{N}$, ($h = \alpha y_0 + \beta y_1 + \gamma y_2 + \dots + \nu y_N$), we have

$$\boxed{E[Y_{\text{out}}^2] = \sigma_v^2 h^T h = \sigma_v^2 (\alpha^2 + \beta^2 + \gamma^2 + \dots)}$$

SNR maximisation

- The SNR may now be expressed as

$$\text{SNR} = \frac{|h^T S Y|^2}{E[h^T Y^2]} = \frac{\alpha^2 \frac{S^T S}{N}}{\sigma_v^2 (\alpha^2 + \beta^2 + \gamma^2 + \dots)}$$

- By inspection, scaling h by some factor ρ will not change the SNR, since both the numerator and denominator will both scale equally by ρ^2 . \rightarrow we can fix $|h| = |\alpha y_0 + \beta y_1 + \gamma y_2 + \dots + \nu y_N| = 1$, then max.

- Setting $|h|=1$, we have $(\alpha^2 + \beta^2 + \gamma^2 + \dots) = 1$, so the SNR becomes

$$\text{SNR} = \frac{\alpha^2 \frac{S^T S}{N}}{\sigma_v^2}$$

- The largest possible value of α given that $|h|=1$ corresponds to $\alpha=1$ (i.e. $\beta=\gamma=\dots=0$).

This means the optimal filter coefficient vector h is

$$\boxed{h = 1 \cdot y_0 = \frac{y_0}{\|y\|}}$$

the received signal normalized.

and the corresponding SNR is

$$\boxed{\text{SNR} = \frac{S^T S}{\sigma_v^2}} \quad \text{depends on energy of signal s and noise v, or expected.}$$

Practical implementation of the matched filter.

- For a signal s with length N , to optimise a filter h , we considered a batch of data X of length N .

- In practice, we would run this optimised filter over a much longer length of data $\{x_n\}$, which contains s at some unknown position, then find the time at which max. energy occurs.

- The time at which max. energy occurs is the pt at which s can be detected, based on some threshold.

- Matched filtering increases the best SNR attainable significantly. Before filtering,

$$\text{SNR}_{\text{raw}} = \frac{\text{max signal value}^2}{\text{average noise energy}} = \frac{\text{max } s_n}{\sigma_v^2}$$

Estimation theory and inference methods.**Estimation and inference**

- In estimation theory, we start w/ a vector of signal measurements $\mathbf{x} = [x_0 \ x_1 \ \dots \ x_{N-1}]^T$ and some unknown quantities (parameters) that we wish to infer, $\boldsymbol{\theta} = [\theta_0 \ \theta_1 \ \dots \ \theta_{P-1}]^T$ (usually $P < N$)
- We typically suppose the probability distribution of the data \mathbf{x} can be expressed as a likelihood function $p(\mathbf{x}|\boldsymbol{\theta})$, which represents how "likely" diff. realizations of the observed data would be if we knew $\boldsymbol{\theta}$.
- In many problems, we may treat $\boldsymbol{\theta}$ as a random vector, and a prior pdf can be formulated for $\boldsymbol{\theta}$, $p(\boldsymbol{\theta})$. The prior represents prior belief about likely parameter configurations prior to any data observations.
- The prior distribution can be used to regularize the inference problem by constraining the parameter search to reasonable parts of the domain of $\boldsymbol{\theta}$ — Bayesian inference.
- In some cases, we may not wish to consider $\boldsymbol{\theta}$ as a random quantity. \rightarrow we can only rely on the likelihood function — classical inference.
- Estimation is the task of formulating an estimate $\hat{\boldsymbol{\theta}}$ which is in some sense close to the true $\boldsymbol{\theta}$.

Inference has estimation as a special case, and studies the whole prob. distribution of the unknown.

Simple estimators and their properties

- It is useful to consider some fundamental properties of estimators — their mean and variance.
 - An estimator $\hat{\boldsymbol{\theta}}$ is said to be unbiased if its expectation satisfies
- $$E[\hat{\boldsymbol{\theta}}] = \boldsymbol{\theta}$$
- An estimator $\hat{\boldsymbol{\theta}}$ is said to be consistent if it is unbiased and its variance satisfies
- $$\lim_{N \rightarrow \infty} \text{Var}[\hat{\boldsymbol{\theta}}] = 0$$
- The defn of unbiased estimators and their variance can be used to measure the performance of any proposed estimation scheme.
 - Estimators can be designed for a given problem specifically to lead to no bias and min. variance.
 - Consider a set of N independent samples of a RV X which has mean μ and standard deviation or $\{x_1, x_2, x_3, \dots, x_N\}$.

An estimator for the mean $\hat{\mu}$, is obtained by taking the arithmetic average of the samples.

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$$

- The expectation of the estimator $\hat{\mu}$, $E[\hat{\mu}]$ is given by

$$E[\hat{\mu}] = E\left[\frac{1}{N} \sum_{i=1}^N x_i\right] = \frac{1}{N} \sum_{i=1}^N E[x_i] = \mu \rightarrow \text{unbiased.}$$

The variance of the estimator $\hat{\mu}$, $\text{Var}[\hat{\mu}]$ is given by

$$\text{Var}[\hat{\mu}] = E[(\hat{\mu} - E[\hat{\mu}])^2] = E[\hat{\mu}^2] - E[\hat{\mu}]^2 = (\mu^2 + \frac{\sigma^2}{N}) - \mu^2 = \frac{\sigma^2}{N} \rightarrow \text{tends to 0 as } N \rightarrow \infty.$$

where $E[\hat{\mu}^2] = E\left[\left(\frac{1}{N} \sum_{i=1}^N x_i\right)^2\right] = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N E[x_i x_j]$, and $E[x_i x_j] = [E[x_i] E[x_j]] = \mu^2 \quad i \neq j$.

$$\therefore E[\hat{\mu}^2] = \frac{1}{N^2} [(N^2 - N)\mu^2 + N(\mu^2 + \sigma^2)] = \frac{1}{N^2} (N\mu^2 + N\sigma^2) = \mu^2 + \frac{\sigma^2}{N}$$

For Personal Use Only -bkwk2

The general linear model / linear regression model.

- In the linear model, it is assumed that the data x are generated as a linear function of the parameters θ w/ an additive random modelling error term e_n .

$$x_n = g_n^T \theta + e_n$$

where g_n is a p -dimensional column vector.

- This can be written for the whole vector x as

$$x = g\theta + e, \quad \text{where } g = \begin{bmatrix} g_1^T \\ g_2^T \\ \vdots \\ g_N^T \end{bmatrix}$$

- The choice of the design matrix g will lead to a wide range of possible models, for example,

① constant level in noise

$$x_n = \theta + e_n \quad \leftrightarrow \quad x = g\theta + e, \quad \text{where } g = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

② sinusoidal model (1 frequency)

$$x_n = a \cos(nw) + b \sin(nw) + e_n \quad \leftrightarrow \quad x = g\theta + e, \quad \text{where } g = [c(w) \ s(w)], \theta = [a \ b]$$

$$(\text{Note } c(w) = [\cos(0) \ \cos(w) \dots \cos((N-1)w)]^T, \ s(w) = [\sin(0) \ \sin(w) \dots \sin((N-1)w)]^T)$$

③ sinusoidal model (J frequencies)

$$x_n = \sum_{j=1}^J a_j \cos(nw_j) + b_j \sin(nw_j) + e_n \quad \leftrightarrow \quad x = g\theta + e, \quad \text{where } g = [c(w_1) \ s(w_1) \dots c(w_J) \ s(w_J)], \theta = [a_1 \ b_1 \ \vdots \ a_J \ b_J]$$

$$(\text{Note } c(w_j) = [\cos(0) \ \cos(w_j) \dots \cos((N-1)w_j)]^T, \ s(w_j) = [\sin(0) \ \sin(w_j) \dots \sin((N-1)w_j)]^T)$$

④ autoregressive model.

$$x_n = \sum_{k=1}^p a_k x_{n-k} + e_n \quad \leftrightarrow \quad x = g\theta + e, \quad \text{where } g = \begin{bmatrix} x_1 & x_2 & \dots & x_{p+1} & x_p \\ x_2 & x_3 & \dots & x_{p+2} & x_{p+1} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{p+1} & x_{p+2} & \dots & x_{N-p} & x_{N-p+1} \end{bmatrix}$$

(Note g is not fixed before making the measurements, but is made up of observed data pts.)

→ cannot directly generate a vector of data using the formula $x = g\theta + e \rightarrow$ generated sequentially)

Einstein-Wiener-Khinchin theorem.

- The Einstein-Wiener-Khinchin thm. states that the power spectrum is equal the expected value of the time-averaged DTFT-squared of the signal values, i.e.

$$S_x(e^{j\omega}) = \lim_{N \rightarrow \infty} E\left[\frac{1}{2N+1} |X^N(e^{j\omega})|^2\right]$$

- consider a time-windowed version of a signal x_n , w/ duration of $2N+1$ samples denoted by x_n^N ,

$$x_n^N = w_n^N x_n, \quad \text{where } w_n^N = \begin{cases} 1 & \text{if } n \in [-N, N] \\ 0 & \text{otherwise} \end{cases}$$

Its DTFT $X^N(e^{j\omega})$ is just the DFT of a frame of $2N+1$ data pts + linear phase shift

$$\text{DTFT}[x_n^N] = X^N(e^{j\omega}) = \sum_{n=-N}^N w_n^N x_n e^{-jn\omega}$$

$$\therefore |X^N(e^{j\omega})|^2 = X^N(e^{j\omega}) X^N(e^{j\omega})^* = \text{DTFT}[x_n^N + \tilde{x}_n^N] = \text{DTFT}\left[\sum_{n=-N}^N x_n w_n^N \tilde{x}_{n+N} w_{n+N}^*\right]$$

ensures we have finite energy
Dividing both sides by the window duration $2N+1$ and taking expectation,

$$E[|X^N(e^{j\omega})|^2] = E\left[\frac{1}{2N+1} \text{DTFT}\left[\sum_{n=-N}^N x_n w_n^N \tilde{x}_{n+N} w_{n+N}^*\right]\right] = \frac{1}{2N+1} \text{DTFT}\left[\sum_{n=-N}^N E[x_n w_n^N \tilde{x}_{n+N} w_{n+N}^*]\right] = \text{DTFT}[R_{xx}(t)]$$

where $R_{xx}(t) = \sum_{n=-N}^N w_n^N w_{n+N}^*$ is the deterministic autocorrelation function of the window function w_n .

For Personal Use Only -bkwk2

Noting that $\text{DTFT}[x[n]] = S_x(e^{j\omega})$ and defining $\text{DTFT}[e[n]] = T(e^{j\omega})$, we have

$$E[\sum_{n=1}^N |X(e^{j\omega})|^2] = \text{DTFT}[x[n]^2] = S_x(e^{j\omega}) * T(e^{j\omega})$$

As N increases, $e[n]$ become flatter and $T(e^{j\omega})$ tends to a delta function, so the limit as $N \rightarrow \infty$ is

$$\lim_{N \rightarrow \infty} E[\sum_{n=1}^N |X(e^{j\omega})|^2] = \lim_{N \rightarrow \infty} \text{DTFT}[x[n]^2] = \text{DTFT}[x[n]^2] = S_x(e^{j\omega})$$

i.e. power spectrum is prop. to expected value of the DTFT squared of the data.

Useful linear algebra/vector calculus properties

- Denote $\nabla J = \frac{\partial J}{\partial \underline{x}}$ as the vector gradient for a column vector \underline{x} , we have

① $J = b^T \underline{x} \rightarrow \nabla J = b$ for a column vector b ,

$$J = b^T \underline{x} = \sum_i b_i x_i \rightarrow \frac{\partial J}{\partial x_i} = b_i \rightarrow \frac{\partial J}{\partial \underline{x}} = b.$$

② $J = \underline{x}^T B \underline{x} \rightarrow \nabla J = (B + B^T) \underline{x}$ for a matrix B .

$$J = \underline{x}^T B \underline{x} = \sum_i x_i \sum_j B_{ij} x_j \rightarrow \frac{\partial J}{\partial x_i} = 2B_i^T \underline{x} + \sum_j (x_j B_{ji} + B_{ij} x_j) = \sum_j (B_{ij} + B_{ji}) x_j \\ \therefore \frac{\partial J}{\partial \underline{x}} = (B + B^T) \underline{x}$$

* For symmetric B , $B^T = B$, so $\frac{\partial J}{\partial \underline{x}} = (B + B) \underline{x} = 2B \underline{x}$

- For $\mathbf{0}$ matrix in the form $B = G^T G$, it is positive semi-definite, i.e. $\underline{x}^T B \underline{x} \geq 0$ for any \underline{x} .

$$0 \leq (\underline{G}\underline{x})^2 = (\underline{G}\underline{x})^T (\underline{G}\underline{x}) = \underline{x}^T \underline{G}^T \underline{G} \underline{x} = \underline{x}^T B \underline{x}$$

If G is full rank, then $\underline{G}\underline{x} \neq 0$ for any non-zero \underline{x} , so B is positive definite, i.e. $\underline{x}^T B \underline{x} > 0$ for any $\underline{x} \neq 0$

- The global min. of $J = \underline{x}^T B \underline{x} - 2b^T \underline{x}$ wrt. \underline{x} for a symmetric positive definite B is $\underline{x} = B^{-1} b$

(i) differentiating J using vector calculus.

$$\frac{\partial J}{\partial \underline{x}} = 2B \underline{x} - 2b = 0 \rightarrow \underline{x} = B^{-1} b$$

(ii) completing the square

$$J = \underline{x}^T B \underline{x} + 2b^T \underline{x} = (\underline{x} - B^{-1} b)^T B (\underline{x} - B^{-1} b) - b^T B^{-1} b \rightarrow \underline{x} = B^{-1} b$$

Linear estimation of the general linear model

Least squares (LS) estimator $\underline{\theta}^{\text{OLS}}$

- For the general linear model $\underline{y} = G\underline{\theta} + \underline{e}$, linear estimator finds the best fit by min. the error J ,

$$J = \sum_{n=1}^{N-1} e_n^2 = \underline{e}^T \underline{e}.$$

Expanding using $\underline{e} = \underline{y} - G\underline{\theta}$,

$$J = \underline{e}^T \underline{e} = (\underline{y} - G\underline{\theta})^T (\underline{y} - G\underline{\theta}) = \underline{y}^T \underline{y} + \underline{G}^T \underline{G} \underline{\theta} - 2 \underline{\theta}^T \underline{G}^T \underline{y}$$

The global min. of J wrt. $\underline{\theta}$ can be found to be $\underline{\theta}^{\text{OLS}} = (\underline{G}^T \underline{G})^{-1} \underline{G}^T \underline{y}$

(i) $\frac{\partial J}{\partial \underline{\theta}} = 0 + 2 \underline{G}^T \underline{G} \underline{\theta} - 2 \underline{G}^T \underline{y} = 0 \rightarrow \underline{\theta}^{\text{OLS}} = (\underline{G}^T \underline{G})^{-1} \underline{G}^T \underline{y}$

(ii) $J = \underline{y}^T \underline{y} + \underline{G}^T \underline{G} \underline{\theta} - 2 \underline{\theta}^T \underline{G}^T \underline{y} = (\underline{\theta} - (\underline{G}^T \underline{G})^{-1} \underline{G}^T \underline{y})^T \underline{G}^T \underline{G} (\underline{\theta} - (\underline{G}^T \underline{G})^{-1} \underline{G}^T \underline{y}) - ((\underline{G}^T \underline{G})^{-1} \underline{G}^T \underline{y})^T \underline{G}^T \underline{G} \underline{\theta} + \underline{y}^T \underline{y}$

$$\therefore \underline{\theta}^{\text{OLS}} = (\underline{G}^T \underline{G})^{-1} \underline{G}^T \underline{y}$$

For Personal Use Only -bkwk2

Properties of the linear estimator

- The OLS estimator of the general linear model is a linear estimator since it is in the form $\hat{\theta} = \underline{C} \underline{x}$.
- To investigate its properties, we consider its mean and variance.

① Mean, $E[\hat{\theta}^{OLS}]$

$$E[\underline{x}] = E[\underline{G}\underline{\theta} + \underline{\epsilon}] = E[\underline{G}\underline{\theta}] + E[\underline{\epsilon}]^T = \underline{G}\underline{\theta}$$

$$\therefore E[\hat{\theta}^{OLS}] = E[(\underline{G}^T \underline{G})^{-1} \underline{G}^T \underline{x}] = (\underline{G}^T \underline{G})^{-1} \underline{G}^T E[\underline{x}] = (\underline{G}^T \underline{G})^{-1} \underline{G}^T \underline{G}\underline{\theta} = \underline{\theta} \rightarrow \text{unbiased}.$$

② Covariance, $\text{cov}[\hat{\theta}^{OLS}] = E[(\hat{\theta}^{OLS} - E[\hat{\theta}^{OLS}])(\hat{\theta}^{OLS} - E[\hat{\theta}^{OLS}])^T]$

Define the OLS matrix term as $\underline{C} = (\underline{G}^T \underline{G})^{-1} \underline{G}^T$, and consider a general unbiased linear estimator

$$\hat{\theta} = \underline{D}\underline{x} \quad \text{where } \underline{D} = \underline{C} + \underline{\Delta}.$$

For $\hat{\theta} = \underline{D}\underline{x}$ to be unbiased, we req.

$$E[\hat{\theta}] = E[\underline{D}\underline{x}] = \underline{\theta} \rightarrow E[\underline{D}\underline{x}] = E[(\underline{C} + \underline{\Delta})\underline{x}] = (\underline{C} + \underline{\Delta})E[\underline{x}] = (\underline{C} + \underline{\Delta})\underline{G}\underline{\theta} = \underline{\theta} + \underline{\Delta}\underline{G}\underline{\theta} = \underline{\theta}$$

$\therefore \underline{\Delta}\underline{G}\underline{\theta} = \underline{0}$, since it must work for any $\underline{\theta}$.

Consider the outer product $\underline{x}\underline{x}^T$, and take its expectation,

$$\underline{x}\underline{x}^T = (\underline{G}\underline{\theta} + \underline{\epsilon})(\underline{G}\underline{\theta} + \underline{\epsilon})^T = \underline{G}\underline{\theta}\underline{\theta}^T\underline{G}^T + \underline{\epsilon}\underline{\epsilon}^T + \underline{\epsilon}\underline{\theta}^T\underline{G}^T + \underline{\theta}\underline{\epsilon}^T \rightarrow E[\underline{x}\underline{x}^T] = \underline{G}\underline{\theta}\underline{\theta}^T\underline{G}^T + \underline{\sigma}^2 \underline{I}$$

$\underline{\theta}$ is zero mean white noise
w variance σ^2 .

For the linear estimator $\hat{\theta} = \underline{D}\underline{x}$, the expectation of the outer product $\hat{\theta}\hat{\theta}^T$, $E[\hat{\theta}\hat{\theta}^T]$, is

$$E[\hat{\theta}\hat{\theta}^T] = E[\underline{D}\underline{x}\underline{x}^T\underline{D}^T] = D E[\underline{x}\underline{x}^T] D^T = \underline{D}(\underline{G}\underline{\theta}\underline{\theta}^T\underline{G}^T + \underline{\sigma}^2 \underline{I}) \underline{D}^T = \underline{\theta}\underline{\theta}^T + \underline{\sigma}^2 \underline{D}\underline{D}^T.$$

$$\text{since } \underline{D}\underline{G}\underline{\theta} = (\underline{C} + \underline{\Delta})\underline{G}\underline{\theta} = \underline{G}\underline{\theta}\underline{\theta}^T + \underline{\Delta}\underline{G}\underline{\theta}^T = \underline{\theta} \rightarrow \underline{D}\underline{G}\underline{\theta}\underline{G}^T\underline{D}^T = \underline{\theta}\underline{\theta}^T. \quad \underline{D}\underline{G}\underline{\theta} = \underline{D}\underline{\Delta}\underline{\theta}^T \rightarrow \underline{D}\underline{G}\underline{\theta}\underline{G}^T\underline{D}^T = \underline{D}\underline{\Delta}\underline{\theta}^T \underline{D}^T = 0$$

$$\therefore \text{cov}[\hat{\theta}] = E[\hat{\theta}\hat{\theta}^T] - \underline{\theta}\underline{\theta}^T = \underline{\theta}\underline{\theta}^T + \underline{\sigma}^2 \underline{D}\underline{D}^T - \underline{\theta}\underline{\theta}^T = \underline{\sigma}^2 (\underline{C} + \underline{\Delta}) (\underline{C} + \underline{\Delta})^T = \underline{\sigma}^2 (\underline{C}\underline{C}^T + \underline{\Delta}\underline{C}^T + \underline{C}\underline{\Delta}^T + \underline{\Delta}\underline{\Delta}^T) \\ = \underline{\sigma}^2 (\underline{C}\underline{C}^T + \underline{\Delta}\underline{\Delta}^T) = \underline{\sigma}^2 ((\underline{G}^T \underline{G})^{-1} \underline{G}^T + \underline{\Delta}\underline{\Delta}^T) = \text{cov}[\hat{\theta}^{OLS}] + \underline{\sigma}^2 \underline{\Delta}\underline{\Delta}^T$$

$$\text{since w/ } \text{cov}[\hat{\theta}^{OLS}] = \text{cov}[\hat{\theta}]|_{\underline{\Delta}=0} = \underline{\sigma}^2 (\underline{G}^T \underline{G})^{-1}$$

\rightarrow Denoting $\hat{\theta}_{ir}$ as the i th diagonal element of $\text{cov}[\hat{\theta}]$, and by construction, the diagonal elements

$$\text{of } \underline{\Delta}\underline{\Delta}^T \text{ are also } \geq 0 \rightarrow \text{var}[\hat{\theta}_{ii}] \geq \text{var}[\hat{\theta}_{ij}^{OLS}] \text{ for each } i=0, \dots, p-1, \text{ w/ equality when } \underline{\Delta}=0.$$

- This means the OLS estimator $\hat{\theta}^{OLS}$ is the MVN estimator of θ . Such an estimator is a best linear unbiased estimator (BLUE).

- By considering the expression for $\text{cov}[\hat{\theta}]$, we can see that $\hat{\theta}^{OLS}$ is the unique BLUE for the general linear model.

- The only assumption used is that the noise process $\{\epsilon_i\}$ is zero mean and white (we did not assume any probability distribution for $\{\epsilon_i\}$).

- However, if we do know the distribution of $\{\epsilon_i\}$, it is possible that the linear estimator can be beaten by a nonlinear estimation method.

- If in addition, we have prior probability information about θ , $p(\theta)$, then a Bayesian estimator can give even better mean squared performance at the cost of some additional bias in the estimate.

For Personal Use Only -bkwk2

Likelihood-based inference

Error sequence ε and the linear Gaussian model

- We assume the error sequence ε to be drawn from an iid noise distribution w/ pdf $p_\varepsilon(\cdot)$, so

$$P(\varepsilon) = p_\varepsilon(e_0)p_\varepsilon(e_1)\dots p_\varepsilon(e_{n-1})$$

and we also assume the noise process is white (though not necessarily zero mean).

- For the linear Gaussian model (Gauss-Markov model), p_ε is the zero-mean Gaussian distribution.

$$p_\varepsilon(x) = N(x | \mu=0, \sigma^2=\sigma_e^2) = \frac{1}{\sqrt{2\pi\sigma_e^2}} e^{-\frac{x^2}{2\sigma_e^2}}$$

The error sequence ε thus has a joint pdf $p(\varepsilon)$, given by

$$p(\varepsilon) = \prod_{n=0}^{N-1} N(e_n | \mu=0, \sigma^2=\sigma_e^2) = \prod_{n=0}^{N-1} \frac{1}{\sqrt{2\pi\sigma_e^2}} e^{-\frac{e_n^2}{2\sigma_e^2}} = \frac{1}{(2\pi\sigma_e^2)^{N/2}} e^{-\frac{\sum_{n=0}^{N-1} e_n^2}{2\sigma_e^2}} = \frac{1}{(2\pi\sigma_e^2)^{N/2}} e^{-\frac{\varepsilon^T \varepsilon}{2\sigma_e^2}} = N(\varepsilon | \mu=0, Q_e = \sigma_e^2 I)$$

- Note the multivariate Gaussian density function for a length N random column vector X , $f_X(X)$ is

$$f_X(X) = \frac{1}{(2\pi)^{N/2} |C_X|^{1/2}} \exp\left(-\frac{1}{2}(X-\mu)^T C_X^{-1}(X-\mu)\right)$$

where $\mu = E[X]$ is the mean vector and $C_X = E[(X-\mu)(X-\mu)^T]$ is the covariance matrix.

(Elementwise, we have $[C_X]_{ij} = C_{X_i X_j}$, the covariance b/w elements X_i and X_j .)

Maximum likelihood (ML) estimator of θ

- The ML estimator treats the parameters θ as unknown const. dat which we incorporate no prior information.

The observed data X is considered random and we can often obtain the pdf for X when θ is known.

- This pdf is termed the likelihood $L(X; \theta)$, defined as

$$L(X; \theta) = p(X | \theta).$$

- The ML estimate for θ is then the value of θ which max the likelihood for given observation X .

$$\theta_{ML} = \arg \max_{\theta} L(X; \theta).$$

i.e. the ML soln corresponds to the parameter vector which would have generated X w/ highest probability.

- we typically max. the log-likelihood function $I(X; \theta) = \log(L(X; \theta))$ rather than $L(X; \theta)$ itself for convenience. since $\log(\cdot)$ is a monotonically increasing function, the two solns are identical.

- For the linear Gaussian model, $X = G\theta + \varepsilon$, we do a vector change of variable $\varepsilon \rightarrow X$ to get the likelihood $p(X | \theta)$. (we are conditioning on $\theta \rightarrow$ treat $G\theta$ as a const. term for change of variable)

$$\frac{\partial X_i}{\partial \theta_j} = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{otherwise} \end{cases} \rightarrow \vec{I} = \vec{G} \rightarrow p(X | \theta) = p_\varepsilon(X - G\theta) | \vec{I}^\top \vec{I}^{-1} = p_\varepsilon(X - G\theta)$$

Expanding this out, we get

$$L(X; \theta) = p_\varepsilon(X - G\theta) = \frac{1}{(2\pi\sigma_e^2)^{N/2}} \exp\left(-\frac{1}{2\sigma_e^2} (X - G\theta)^T (X - G\theta)\right)$$

$$\therefore I(X; \theta) = \log L(X; \theta) = -\frac{N}{2} \log 2\pi\sigma_e^2 - \frac{1}{2\sigma_e^2} (X - G\theta)^T (X - G\theta) = -\frac{N}{2} \log 2\pi\sigma_e^2 - \frac{1}{2\sigma_e^2} \varepsilon^T \varepsilon = -\frac{N}{2} \log 2\pi\sigma_e^2 - \frac{1}{2\sigma_e^2} \varepsilon^T \varepsilon$$

\rightarrow Maximising $I(X; \theta)$ wrt θ is equivalent to minimising $\varepsilon^T \varepsilon$ wrt $\theta \rightarrow$ some o/ optimisation, so

$$\theta^{ML} = \theta^{OLS} = (G^T G)^{-1} G^T X$$

- The OLS soln is a special case of the ML soln where the error process ε is zero-mean, independent and Gaussian w/ fixed variance.

For Personal Use Only -bkwk2

Estimating the noise variance

- The noise variance σ_e^2 of the linear Gaussian model can be estimated by ML.

- consider the log-likelihood function of $\theta = \theta^{ML}$, now also considered as a function of σ_e^2 .

$$I(\underline{x}; \theta = \theta^{ML}, \sigma_e^2) = \log L(\underline{x}; \theta = \theta^{ML}, \sigma_e^2) = -\frac{N}{2} \log(2\pi\sigma_e^2) - \frac{1}{2\sigma_e^2} (\underline{x} - \underline{\theta}^{ML})^T (\underline{x} - \underline{\theta}^{ML}) = -\frac{N}{2} \log(2\pi\sigma_e^2) - \frac{J^{ML}}{2\sigma_e^2}$$

where $J^{ML} = (\underline{x} - \underline{\theta}^{ML})^T (\underline{x} - \underline{\theta}^{ML})$ is the m.s. squared error term corresponding to ML optimisation.

Differentiating wrt σ_e^2 and setting to zero,

$$\frac{\partial I(\underline{x}; \theta = \theta^{ML}, \sigma_e^2)}{\partial \sigma_e^2} = -\frac{N\sigma_e^2}{\sigma_e^4} + \frac{J^{ML}}{2\sigma_e^4} = 0 \quad \rightarrow \quad \boxed{\sigma_e^{ML} = \sqrt{\frac{J^{ML}}{N}}}$$

i.e. the noise variance σ_e^{ML} is estimated as the mean squared error at the ML parameter θ^{ML} .

Bayesian inference

Bayesian methods

- The ML method treats parameters θ as unknown const. If we instead treat the parameters θ as RVs, we can assign a prior pdf $p(\theta)$. key diff b/w likelihood-based / Bayesian inference

- The prior pdf $p(\theta)$ should ideally express some prior knowledge abt. the rel. probability of diff. parameter values θ before the data \underline{x} are observed. (If nothing abt the parameters θ is known a priori, $p(\theta)$ should not express initial preference for a set of parameters over others).

- The precise form assigned to $p(\theta)$ req. careful consideration since mis-leading results can be obtained from erroneous priors.

- We can use Bayes thm. to find the posterior probability $p(\theta | \underline{x})$, which can be used to estimate random parameters θ from a random vector \underline{x} of observations.

$$p(\theta | \underline{x}) = \frac{p(\underline{x} | \theta) p(\theta)}{p(\underline{x})}$$

(Note $p(\underline{x} | \theta)$ is the likelihood ; $p(\theta)$ is the prior ; $p(\underline{x})$ is the marginal likelihood)

- The generation of the posterior from the prior when data \underline{x} is observed can be thought of as a refinement to any previous knowledge abt. the parameters θ . ($p(\theta) \xrightarrow{\text{data } \underline{x}} p(\theta | \underline{x})$).

- The marginal likelihood $p(\underline{x})$ is useful in model selection problems, but is const. for any given observation $\underline{x} \rightarrow$ ignored if we only want to optimise the posterior probabilities wrt θ .

$$\boxed{p(\theta | \underline{x}) \propto p(\underline{x} | \theta) p(\theta)}$$

- $p(\underline{x})$ could be calculated using $p(\underline{x}) = \int p(\underline{x}, \theta) d\theta = \int p(\underline{x} | \theta) p(\theta) d\theta$, and effectively serves as the normalising const. for the posterior density $p(\theta | \underline{x})$.

- All the distributions above are implicitly conditioned upon all prior modelling assumptions, so some texts adopt the notation $p(\theta | \underline{x}, M)$, where M denotes all of the modelling and distributional assumptions that are being made (not just the assumed prior $p(\theta)$).

- LS / ML methods only provide a pt-estimate of θ , whereas Bayesian methods produce a posterior pdf w/ values defined for all $\theta \rightarrow$ fully interpretable probability distribution

For Personal Use Only -bkwk2

Maximum a posteriori (MAP) estimator $\underline{\theta}^{\text{MAP}}$

- The MAP estimate for $\underline{\theta}$ is the value of $\underline{\theta}$ which max. the posterior for given observations \underline{x} and prior $p(\underline{\theta})$

$$\underline{\theta}^{\text{MAP}} = \arg \max p(\underline{\theta}|\underline{x})$$

i.e., the MAP soln corresponds to the parameter vector $\underline{\theta}$ inferred from \underline{x} w/ highest probability

- For the linear Gaussian model, suppose the prior $p(\underline{\theta})$ is the multivariate Gaussian, i.e.

$$p(\underline{\theta}) = N(\underline{\theta} | \underline{m}_\theta, \underline{\Sigma}_\theta) = \frac{1}{(2\pi)^{n/2} |\underline{\Sigma}_\theta|^{1/2}} \exp(-\frac{1}{2} (\underline{\theta} - \underline{m}_\theta)^T \underline{\Sigma}_\theta^{-1} (\underline{\theta} - \underline{m}_\theta))$$

where n is the no. of parameters in $\underline{\theta}$.

(Note if we don't know much a priori, we usually select $\underline{m}_\theta = 0$ and $\underline{\Sigma}_\theta = \underline{\Omega}_\theta^2 \underline{I}$)

The likelihood $p(\underline{x}|\underline{\theta})$ takes the same form as before for the ML estimator, so the posterior $p(\underline{\theta}|\underline{x})$ is

$$p(\underline{\theta}|\underline{x}) \propto p(\underline{\theta}) p(\underline{x}|\underline{\theta}) = \frac{1}{(2\pi)^{n/2} |\underline{\Sigma}_\theta|^{1/2}} \exp(-\frac{1}{2} (\underline{\theta} - \underline{m}_\theta)^T \underline{\Sigma}_\theta^{-1} (\underline{\theta} - \underline{m}_\theta)) \cdot \frac{1}{(2\pi)^{n/2} |\underline{\Sigma}_x|^{1/2}} \exp(-\frac{1}{2} (\underline{x} - \underline{G}\underline{\theta})^T \underline{\Sigma}_x^{-1} (\underline{x} - \underline{G}\underline{\theta}))$$

Considering $-2\Omega_e^2 \log(p(\underline{\theta}|\underline{x}))$, we have

$$-2\Omega_e^2 \log(p(\underline{\theta}|\underline{x})) = \Omega_e^2 (\underline{\theta} - \underline{m}_\theta)^T \underline{\Sigma}_\theta^{-1} (\underline{\theta} - \underline{m}_\theta) + (\underline{x} - \underline{G}\underline{\theta})^T (\underline{x} - \underline{G}\underline{\theta}) + \text{const.}$$

$$= \Omega_e^2 (\underline{\theta}^T \underline{\Sigma}_\theta^{-1} \underline{\theta} + \underline{m}_\theta^T \underline{\Sigma}_\theta^{-1} \underline{m}_\theta - 2 \underline{m}_\theta^T \underline{\Sigma}_\theta^{-1} \underline{\theta}) + \underline{x}^T \underline{x} + \underline{x}^T \underline{\Sigma}_x^{-1} \underline{x} - 2 \underline{x}^T \underline{\Sigma}_x^{-1} \underline{x} + \text{const.}$$

(i) Differentiating wrt $\underline{\theta}$ and setting to zero,

$$\frac{\partial (-2\Omega_e^2 \log(p(\underline{\theta}|\underline{x})))}{\partial \underline{\theta}} = \Omega_e^2 (2\underline{\Sigma}_\theta^{-1} \underline{\theta} + \underline{m}_\theta^T \underline{\Sigma}_\theta^{-1} \underline{m}_\theta) + 0 + 2\underline{\Sigma}_x^{-1} \underline{x} - 2\underline{\Sigma}_x^{-1} \underline{x} = 0$$

$$\therefore (\Omega_e^2 \underline{\Sigma}_\theta^{-1} + \underline{\Sigma}_x^{-1}) \underline{\theta} = \Omega_e^2 \underline{\Sigma}_\theta^{-1} \underline{m}_\theta + \underline{\Sigma}_x^{-1} \underline{x} \rightarrow \underline{\theta}^{\text{MAP}} = (\underline{\Sigma}_\theta^{-1} \underline{\Sigma}_\theta + \Omega_e^2 \underline{\Sigma}_\theta^{-1})^{-1} (\underline{\Sigma}_\theta^{-1} \underline{x} + \Omega_e^2 \underline{\Sigma}_\theta^{-1} \underline{m}_\theta)$$

(ii) Grouping terms in the form $\underline{\theta}^T \underline{\Sigma}_\theta \underline{\theta} \rightarrow \underline{\theta}^T \underline{\theta}$. Then completing the square

$$-2\Omega_e^2 \log(p(\underline{\theta}|\underline{x})) = \underline{\theta}^T (\underline{\Sigma}_\theta^{-1} \underline{\theta} + \Omega_e^2 \underline{\Sigma}_\theta^{-1}) \underline{\theta} - 2(\underline{x}^T \underline{\Sigma}_x^{-1} + \Omega_e^2 \underline{\Sigma}_\theta^{-1} \underline{m}_\theta) \underline{\theta} + \text{const.}$$

$$= \underline{\theta}^T (\underline{\Sigma}_\theta^{-1} \underline{\theta} + \Omega_e^2 \underline{\Sigma}_\theta^{-1}) \underline{\theta} - 2(\underline{x}^T \underline{\Sigma}_x^{-1} + \Omega_e^2 \underline{\Sigma}_\theta^{-1} \underline{m}_\theta) \underline{\theta} + \text{const.}$$

$$= \underline{\theta}^T \underline{\Sigma}_\theta \underline{\theta} - 2 \underline{\theta}^T \underline{\Sigma}_x^{-1} \underline{x} + \text{const.} = (\underline{\theta} - \underline{\beta})^T \underline{\Sigma}_\theta (\underline{\theta} - \underline{\beta}) + \text{const.}$$

$$\therefore \underline{\theta}^{\text{MAP}} = \underline{\beta} + \underline{b}, \quad \text{where } \underline{\beta} = (\underline{\Sigma}_\theta^{-1} \underline{\theta} + \Omega_e^2 \underline{\Sigma}_\theta^{-1}), \quad \underline{b} = (\underline{\Sigma}_\theta^{-1} \underline{x} + \Omega_e^2 \underline{\Sigma}_\theta^{-1} \underline{m}_\theta)$$

(Note we assumed $\underline{\Sigma}_\theta$ is invertible, which is the case if $\underline{\Sigma}_\theta$ is full rank.)

- By completing the square, for $-2\Omega_e^2 \log(p(\underline{\theta}|\underline{x}))$, and considering the exponent of $p(\underline{\theta}|\underline{x})$,

$$-2\Omega_e^2 \log(p(\underline{\theta}|\underline{x})) = (\underline{\theta} - \underline{\theta}^{\text{MAP}})^T \underline{\Sigma}_\theta (\underline{\theta} - \underline{\theta}^{\text{MAP}}) + \text{const.}$$

\therefore Exponent of $p(\underline{\theta}|\underline{x}) = -\frac{1}{2\Omega_e^2} (\underline{\theta} - \underline{\theta}^{\text{MAP}})^T \underline{\Sigma}_\theta (\underline{\theta} - \underline{\theta}^{\text{MAP}}) \rightarrow$ multivariate Gaussian $N(\underline{m}_\theta^{\text{post}}, \underline{\Sigma}_\theta^{\text{post}})$,

$$\underline{m}_\theta^{\text{post}} = \underline{\theta}^{\text{MAP}}, \quad \underline{\Sigma}_\theta^{\text{post}} = \left(\frac{\underline{\beta}}{\Omega_e^2} \right)^{-1} = \Omega_e^2 \underline{\beta}^{-1}$$

since the remaining terms do not depend on $\underline{\theta}$, so the multivariate density function must be proper

we conclude the posterior $p(\underline{\theta}|\underline{x})$ is still a multivariate Gaussian.

$$p(\underline{\theta}|\underline{x}) = N(\underline{\theta} | \underline{\theta}^{\text{MAP}}, \Omega_e^2 \underline{\beta}^{-1})$$

- comparing $\underline{\theta}^{\text{MAP}} = (\underline{\Sigma}_\theta^{-1} \underline{\theta} + \Omega_e^2 \underline{\Sigma}_\theta^{-1})^{-1} (\underline{\Sigma}_\theta^{-1} \underline{x} + \Omega_e^2 \underline{\Sigma}_\theta^{-1} \underline{m}_\theta) \sim \underline{\theta}^{\text{ML}} = (\underline{\Sigma}_\theta^{-1} \underline{\Sigma}_\theta)^{-1} \underline{\Sigma}_\theta^{-1} \underline{x}$, we can see in the limit $\Omega_e^2 = 0$, $\underline{\theta}^{\text{MAP}} = \underline{\theta}^{\text{ML}}$

($\underline{\Sigma}_\theta^{-1} \rightarrow \underline{\Sigma}_\theta \rightarrow \infty$, i.e. the prior tends to be flat/uniform (equivalent $\underline{\theta}$) \rightarrow impose no prior information)

- The ML soln is a special case of the MAP soln, where a uniform prior is assigned to $\underline{\theta}$.

- The MAP estimate also tends to the ML estimate when the likelihood is sharply peaked around the max, compared w/ the prior (i.e. $\Omega_e^2 \ll \Omega_\theta^2$). This is the case for large N .

For Personal Use Only -bkwk2

Minimum mean-squared error (MMSE) estimator. $\hat{\theta}^{\text{MMSE}}$

- consider a cost function $C(\hat{\theta}, \theta)$ which expresses the cost of estimating the parameter as $\hat{\theta}$ when the true value is θ . ($C(\hat{\theta}, \theta)$ should be nonnegative and usually satisfies $C(\hat{\theta}, \theta) = 0$).
- the form of the cost function $C(\hat{\theta}, \theta)$ depends on the req. of the problem. A cost of 0 indicates that the estimate is perfect for our req. (not necessarily $\hat{\theta} = \theta$), while the costs indicate poorer estimate.
- The expected cost over all of the unknown parameters, conditioned on the observed data x is

$$E[C(\hat{\theta}, \theta) | X=x] = \int_0 C(\hat{\theta}, \theta) p(\theta | x) d\theta$$

- The MMSE estimate for θ is the value of $\hat{\theta}$ which min. the average squared error for given data x .

$$\hat{\theta}^{\text{MMSE}} = \underset{\hat{\theta}}{\operatorname{arg\,min}} E[(\hat{\theta} - \theta)^2 | X=x]$$

(Note the cost function $C(\hat{\theta}, \theta)$ is $C(\hat{\theta}, \theta) = (\hat{\theta} - \theta)^2$)

Any small bias introduced can be tolerated provided the MSE is low.

- In general, the MSE can be expressed as

$$J = E[(\hat{\theta} - \theta)^2 | X=x] = \int_0 (\hat{\theta} - \theta)^2 p(\theta | x) d\theta$$

Differentiating wrt $\hat{\theta}$ and setting to zero,

$$\frac{dJ}{d\hat{\theta}} = \int_0 (\hat{\theta} - \theta)^2 p'(\theta | x) d\theta = \int_0 (\hat{\theta} - \theta)^2 p(\theta | x) d\theta = \int_0 2(\hat{\theta} - \theta) p(\theta | x) d\theta = 0$$

$$\therefore \int_0 \hat{\theta} p(\theta | x) d\theta = \int_0 \theta p(\theta | x) d\theta \rightarrow \hat{\theta}^{\text{MMSE}} = E[\theta | X=x] = \int_0 \theta p(\theta | x) d\theta$$

since $\int_0 \theta^2 p(\theta | x) d\theta = \hat{\theta} \int_0 \theta p(\theta | x) d\theta = \hat{\theta}^2$

- For the linear Gaussian model, the posterior distribution is $p(\theta | x) = N(\theta | \underline{x}^{\text{MAP}}, \sigma_x^2 \underline{I}^{-1})$, so

$$\hat{\theta}^{\text{MMSE}} = E[\theta | X=x] = \underline{\theta}^{\text{MAP}} = (\underline{x}^T \underline{Q} + \sigma_x^2 \underline{C}^{-1})^{-1} (\underline{x}^T \underline{x} + \sigma_x^2 \underline{C}^{-1} \underline{m}_0)$$

Summary of estimators

① Least squares ($\hat{\theta}^{\text{LS}}$)

- ✓: Req. no knowledge of probability distributions ✗: cannot incorporate prior knowledge of $p(\theta)$
- ✓: Usually simplest to implement
- ✓: Guarantee of performance as BLUE estimator (not necessarily better than non-Gaussian estimators)

② Maximum likelihood ($\hat{\theta}^{\text{ML}}$)

- ✓: Guarantee of performance when N is large
- ✗: Req. knowledge of noise probability distribution
- ✗: Cannot incorporate prior knowledge of $p(\theta)$.
- ✗: Harder to implement than LS for non-Gaussian case.

③ Bayesian ($\hat{\theta}^{\text{MAP}}, \hat{\theta}^{\text{MMSE}}$)

- ✓: Incorporates prior knowledge of $p(\theta)$.
- ✗: Req. knowledge of noise probability distribution
- ✓: Guarantee of performance for any N (given correct $p(\theta)$)
- ✗: Req. knowledge of prior distribution $p(\theta)$.
- ✗: Can be harder to implement than LS/ML