

Statistical Signal Processing

Theory

Book : Adaptive Filter, by Simon Haykin

For Random Process - Papoulis
Peebles
Yates

- Probability Theory based Axioms -
sets, Fields and Events

- Sample Description space - Set of all outcomes
- Random Experiment $\rightarrow \Omega$

- Subsets of Ω are Events
 $\Omega \in \{\Xi_i\}$

$\Omega \rightarrow$ Countably Finite or Infinite outcomes
 \hookrightarrow Uncountably infinite

If, $\Omega = \{\Xi_1, \dots, \Xi_N\}$ then the number of subsets is 2^N

- If Ω has a countable number of elements then any subset of Ω may be assigned a probability according to axioms.
- Then the class of all subsets make up a σ -field and each subset is an event.
- σ -field is denoted by F
- We talk of σ -field, F , of events

- When $\Omega = \mathbb{R}$ = the real line, then not every subset of Ω can be assigned a probability. ~~is~~ consistent with the axioms.
- Only those subsets to which a probability can be assigned consistent with the axioms will be called events.
- The smaller collection forms a σ -field, F
- On the real line, \mathbb{R} , the σ -field is sometimes called the BOREL field

$P \rightarrow$ Probability Space
 $\rightarrow (\Omega, F, P)$

Axiomatic Definition

Probability is a set function $P[\cdot]$ that assigns to every event $E \in F$ a number $P[E]$ called the probability of E st.

- i) $P[E] \geq 0$
- ii) $P[\Omega] = 1$
- iii) $P[E \cup F] = P[E] + P[F]$ if $E \cap F = \emptyset$

Random Process

Definition

Let (Ω, F, P) be a probability space.

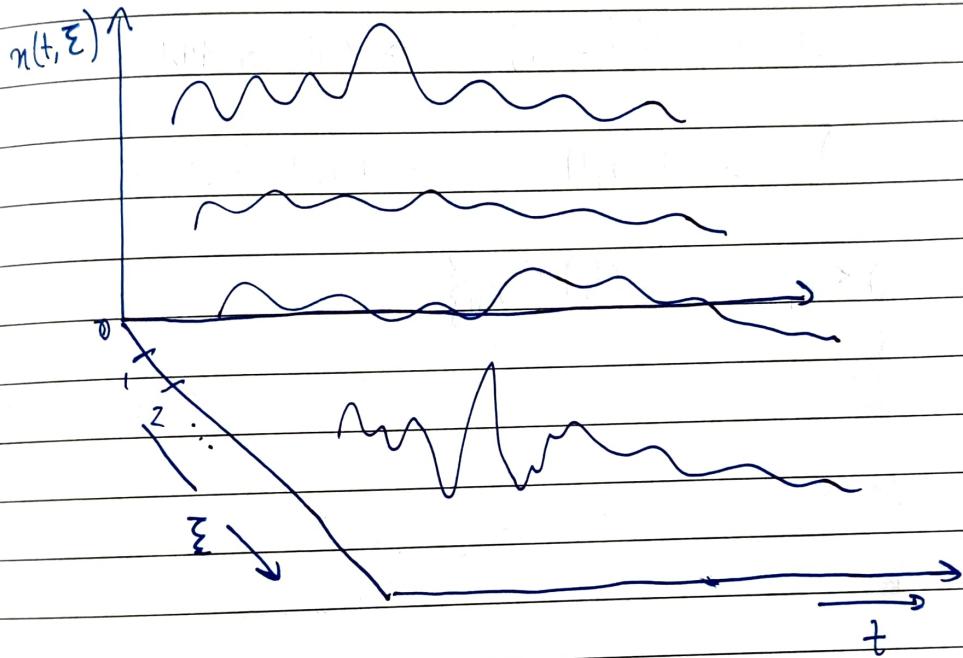
Then define a mapping X from the sample space Ω to a space of continuous time ~~function~~ functions.

The elements in this space will be called sample functions.

This mapping is called a random process if at each fixed time, the mapping is a random variable.

$$X(A, \xi) \in F \rightsquigarrow \{ \xi : X(\xi) \leq x \} \subset F \neq \emptyset$$

for each fixed 't' on the real line $-\infty < t < \infty$



Definition (Random Process).

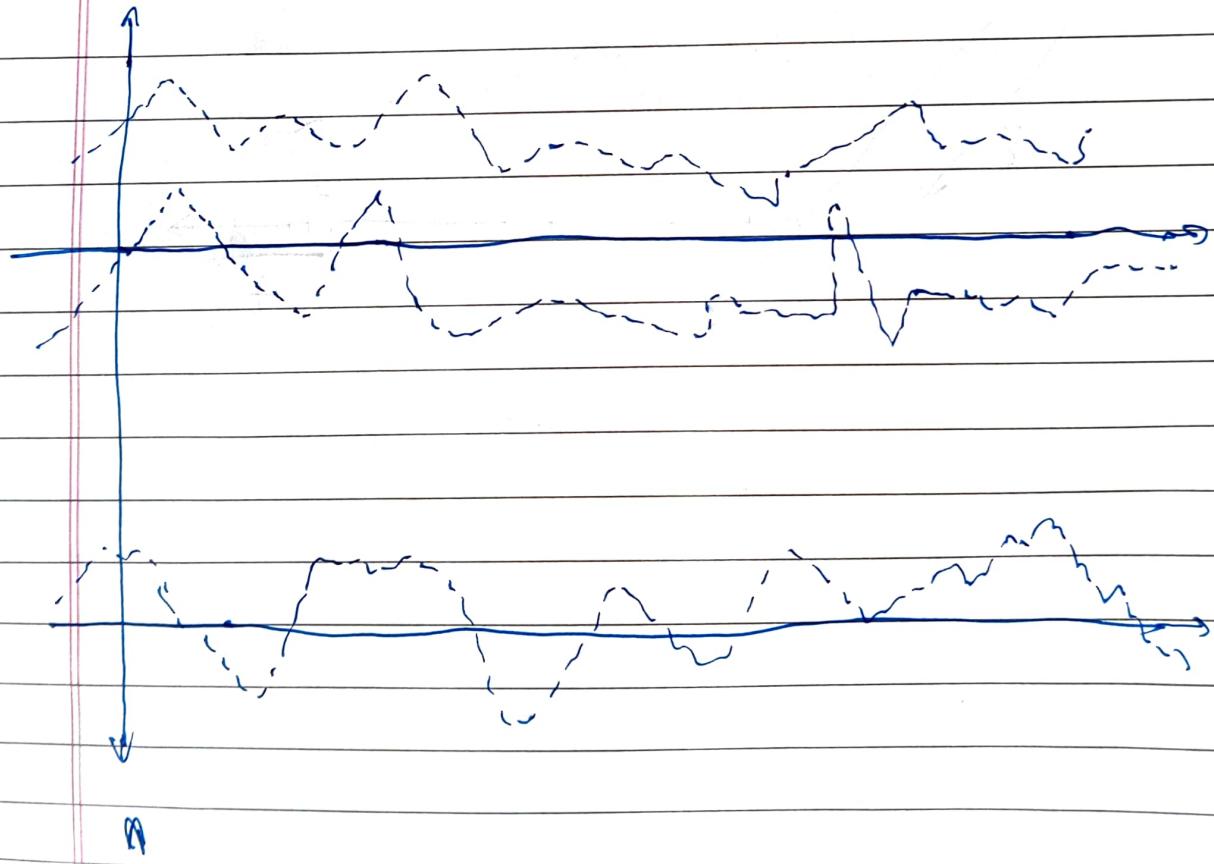
A random process $X(t)$ is statistically specified by its complete set of n^{th} order PDFs for all positive integers 'n',

$$F_x(n_1, n_2, \dots, n_n) ; t_1, t_2, \dots, t_n \quad \forall n_1, n_2, \dots, n_n \text{ and for all } t_1, t_2, \dots, t_n$$

$$F(n, t) = P\{X(t) \leq n\} \rightarrow \text{CDF or PDF}$$

$$f(n, t) = \frac{d}{dn} F(n, t) \rightarrow \text{Probability density function}$$

$$\Rightarrow P\{n(t_1) \leq n_1, \dots, X(t_n) \leq n_n\}$$



Definition

Let $\eta_t(x)$ denote the no. of trials st at time 't' the ordinates of the observed functions do not exceed 'x'

$$F(x, t) \triangleq \frac{\eta_t(x)}{n}$$

$$F(n_1, n_2; t_1, t_2) = P\{X(t_1) \leq n_1, X(t_2) \leq n_2\}$$

Strict Sense Stationary Process

$$f(n_1, n_2, \dots, n_n; t_1, t_2, \dots, t_n) = f(n_1, n_2, \dots, n_n; t_1+c, t_2+c, \dots, t_n+c)$$

n-dimensional

When you sample the random process with a comb at time t and $t+c$, the same random variables are obtained.

Joint Probability Distribution will remain constant over time.

Wide Sense Stationary

Mean and Autocorrelation are time invariant

Mean Ergodicity

Find $E\{x(t)\}$ from a single sample function.

Q When can we substitute time average in place of ensemble average?

$$\hat{x}(t_n) = \frac{1}{n} \sum_i x(t_n, \xi_i)$$

↑
capital x

(Ensemble Average, over
same time stamp in different
ensembles)

- Consider one sample function $n(t)$ of a stationary with the observation interval $-T \leq t \leq T$

- The DC value of $x(t)$ is defined as the time average

$$\mu_x(T) = \frac{1}{2T} \int_{-T}^T n(t) dt$$

$$E[\mu_x(T)] = \frac{1}{2T} \int_T^T E[n(t)] dt$$

↑
small n

$$= \frac{1}{2T} \int_T^T \mu_x dt$$

$$\bar{x} = \mu_x = \text{Ensemble Mean}$$

↓
capital x

- A random process $X(t)$ is ergodic in the mean if

i) The time average $\mu_n(T)$ approaches the ensemble average μ_x in the limit

$$\lim_{T \rightarrow \infty} \mu_n(T) = \mu_x$$

- ii) The variance of $\mu_n(T)$ treated as a random variable approaches zero in the limit as $T \rightarrow \infty$

$$\lim_{T \rightarrow \infty} \text{var} [\mu_n(T)] = 0$$

Ergodicity in Auto-correlation Function

The time averaged auto-correlation function is defined as

$\xrightarrow{\text{Time auto correlation}}$ $R_{\bar{x}}(\tau, T) = \frac{1}{2T} \int_{-T}^T x(t+\tau)x(t) dt$

- i) $\lim_{T \rightarrow \infty} R_{\bar{x}}(\tau, T) = R_x(\tau) \rightarrow \text{Ensemble Auto-correlation}$
- ii) $\lim_{T \rightarrow \infty} \text{var} [R_{\bar{x}}(\tau, T)] = 0$

Gaussian Random Process

Consider a set of r.v.s $X(t_1), \dots, X(t_n)$ obtained by observing a random process $X(t)$.

If this set of r.v.s is jointly Gaussian for any 'n' with the n-fold being completely specified by

$$\mu_X(t_i) = E[X(t_i)]$$

and the of covariance function $C_{XX}(t_n, t_i)$:-

$$C_{XX}(t_n, t_i) = E[\{X(t_n) - \mu_X(t_n)\} \cdot \{X(t_i) - \mu_X(t_i)\}]$$

where $n, i = 1, 2, \dots, n$

Then the random process is a Gaussian process

If $\underline{X} = \begin{bmatrix} X(t_1) \\ X(t_n) \end{bmatrix}$ and \underline{n} denote a realization of \underline{X}

then $f_{\underline{X}(t_1) \dots \underline{n}(t_n)}(n_1, n_2, \dots, n_n)$

$$= \frac{1}{(2\pi)^{n/2} |\Delta|^{\frac{1}{2}}} \exp \left[-\frac{1}{2} (\underline{n} - \underline{\mu}_X)^T \Sigma^{-1} (\underline{n} - \underline{\mu}_X) \right]$$

where $\underline{\mu}_X = \begin{bmatrix} E[X(t_1)] \\ \vdots \\ E[X(t_n)] \end{bmatrix}$, $\sum \{ C_{XX}(t_k, t_i) \}_{k=1}^n$

and Δ is determinant of the covariance matrix Σ

H.V. Show 2-dimensional joint pdf of Gaussian random variable

Read Markovian Process

Multivariate Gaussian PDF :-

$$f_X(n_1, n_2, \dots, n_K) = \frac{1}{\sqrt{(2\pi)^K |\Sigma|}} e^{-\frac{1}{2} (\underline{n} - \underline{\mu})^T \Sigma^{-1} (\underline{n} - \underline{\mu})}$$

where $\sqrt{(\underline{n} - \underline{\mu})^T \Sigma^{-1} (\underline{n} - \underline{\mu})}$ is the Mahalanobis Distance of \underline{n} from $\underline{\mu}$.

Bivariate Gaussian PDF :-

$$f(n, y) = \frac{1}{2\pi \sigma_x \sigma_y \sqrt{1-p^2}} e^{-\frac{1}{2(1-p^2)} \left[\left(\frac{n-\mu_x}{\sigma_x} \right)^2 - 2p \left(\frac{n-\mu_x}{\sigma_x} \right) \left(\frac{y-\mu_y}{\sigma_y} \right) + \left(\frac{y-\mu_y}{\sigma_y} \right)^2 \right]}$$

$$\mu = \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \sigma_x^2 & p\sigma_x \sigma_y \\ p\sigma_x \sigma_y & \sigma_y^2 \end{pmatrix}$$

Markov Process

A random process whose future probabilities are independent of all past values, and depends only on its present state.

A stochastic process $x(t)$ is called Markov if for every n and $t_1 < t_2 < \dots < t_n$, we have

$$p(x(t_n) \leq x_n | x(t_{n-1}), \dots, x(t_1)) = p(x(t_n) \leq x_n | x(t_{n-1}))$$

e.g. → Coin Toss, Throwing a die

Non-Markov Process example - Weather

Markov Random Process

a) A continuous valued first order random process satisfies

$$f_x(x_n | x_{n-1}, x_{n-2}, \dots, x_1; t_n, t_{n-1}, \dots, t_1)$$

$$= f_x(x_n | x_{n-1}; t_n, t_{n-1})$$

b) Discrete valued (first order) Markov random process satisfies

$$\begin{aligned} p_x(x_n | x_{n-1}, \dots, x_1; t_n, \dots, t_1) \\ \downarrow \text{probability mass function} \\ = p_x(x_n | x_{n-1}; \cancel{t_{n-2}}, \cancel{t_{n-3}}, t_n, t_{n-1}) \end{aligned}$$

$\forall x_1, x_2, \dots, x_n$ and for all $t_1 < \dots < t_n$ and for all integers n .

$$p_x(x_1, \dots, x_n; t_1, t_2, \dots, t_n) = p_x(x_1) \cdot p_x(x_2|x_1) \cdot p_x(x_3|x_2) \dots p_x(x_n|x_{n-1})$$

for a 1st order Markov Process

Chapman-Kolmogorov Equations

- Consider $t_3 > t_2 > t_1$ and the Markov process random variable at these time instants

$$X(t_3), X(t_2) \text{ and } X(t_1)$$

- We wish to compute the conditional density of $X(t_3)$ given $X(t_1)$

$$f_p(x_3 | x_1) = \int_{-\infty}^{\infty} f_x(x_3 | x_2) \cdot f(x_2 | x_1) dx_2$$

which is known as the Chapman-Kolmogorov equation for the transition density $f_x(x_3 | x_1)$ of a Markov process.

Time Series

- The particular realization of a random process is called a time series.

- The sequence $v(n), v(n-1), \dots, v(n-M)$ represents a time series.

- Mean value function of the process $v(n)$ is defined as

$$\mu(n) = E[v(n)]$$

- Autocorrelation function of the process is defined as

$$r(n, n-k) = E[v(n)v^*(n-k)] \quad \forall k=0, \pm 1, \pm 2, \dots$$

- Autocovariance function of the random process $v(n)$ is

$$c(n, n-k) = E[(v(n) - \mu(n)) \cdot (v(n-k) - \mu(n-k))^*]$$

for $n \neq k = 0, \pm 1, \pm 2$

$$\text{Therefore } c(n, n-k) = r(n, n-k) - \mu(n)\mu^*(n-k)$$

For a Wide Sense Stationary Process :-

$$\mu(n) = \mu \text{ for all } n$$

$$r(n, n-k) = r(k)$$

$$c(n, n-k) = c(k)$$

- Note that when $k=0$, corresponding time lag of zero

$$r(0) = E[|u(n)|^2]$$

$$\text{and } c(0) = \sigma_u^2 \rightarrow \text{Variance of } u(n)$$

Correlation Matrix

Let the $M \times 1$ observation vector $u(n)$ denote the elements of the zero mean Time series

$$u(n), u(n-1), \dots, u(n-M+1)$$

$$\text{So, } \underline{u}(n) = [u(n), u(n-1), \dots, u(n-M+1)]^T$$

matrix

We define the correlation of a stationary discrete-time stochastic process represented by the time series as

$$R = E[\underline{u}(n) \underline{u}(n)^H]$$

' H ' denotes transpose + complex conjugate

$$R = \begin{bmatrix} r(0) & r(1) & \dots & r(M-1) \\ r(-1) & r(0) & \dots & r(M-2) \\ \vdots & \vdots & \ddots & \vdots \\ r(-M+1) & r(-M+2) & \dots & r(0) \end{bmatrix}$$

- The element $r(0)$ on the main diagonal is always real valued

Properties of Correlation Matrix

Assuming a discrete time stochastic process.

1. Correlation matrix of a stationary discrete-time random process is Hermitian : $R = R^H$

Another way of stating the Hermitian property

$$r(-k) = r^*(k)$$

where $r(k)$ is the auto-correlation function of the random process $u(n)$ for a lag of ' k '.

- For a WSS process $u(n)$

$$R = \begin{bmatrix} r(0) & r(1) & \dots & r(m-1) \\ r^*(1) & r(0) & \dots & r(m-2) \\ \vdots & \vdots & \ddots & \vdots \\ r^*(m-1) & r^*(m-2) & \dots & r(0) \end{bmatrix}$$

2. Correlation matrix of a stationary discrete time stochastic process is TOEPLITZ

3. Correlation matrix of a DTSP is always positive semi definite and almost always +ve definite.

Let ' a ' be an arbitrary (non-zero) $m \times 1$ complex-valued vector. Define the scalar random variable $y = \underline{a}^H \underline{u}(n)$

Taking the Hermitian transpose of both sides and noting that y is a scalar, we get

$$y^* = \underline{u^H(n)} \underline{a}$$

The mean square value of the random variable y is

$$\begin{aligned} E[|y|^2] &= E[yy^*] = E[\underline{a}^H \underline{u}(n) \underline{u}^H \underline{a}] \\ &= \underline{a}^H E[\underline{u}(n) \underline{u}^H(n)] \underline{a} \\ &= \underline{a}^H R \underline{a} \end{aligned}$$

Since $E[|y|^2] \geq 0$, it follows that $\underline{a}^H R \underline{a} \geq 0$

$$\underline{u}^B(n) = [u(n-M+1), u(n-M+2), \dots, u(n)]^T \quad [B \leftarrow \text{Backwards}]$$

$$E[\underline{u}^B(n) \underline{u}^B(n)^H] = \begin{bmatrix} r(0) & r^*(1) & \cdots & r^*(M-1) \\ r(1) & r(0) & \cdots & r^*(M-2) \\ \vdots & \vdots & \ddots & \vdots \\ r(M) & r(M-2) & \cdots & r(0) \end{bmatrix} = R^T$$

Correlation Matrix of Sine Wave Plus Noise

- Let α denote the complex amplitude of the complex sinusoid and
- ω denote its angular frequency
- $v(n)$ denote a sample of the noise assumed to have zero mean
- A sample of the time series consisting of the complex sinusoid plus noise can be written as $v(n) = \alpha e^{j\omega n} + v(n)$
 $n=0, 1, \dots, N-1$

Note: The source of the complex sinusoid and the noise are independent of each other.

- By assumption, the noise component $v(n)$ has zero mean, therefore
 $E[v(n)] = \alpha \exp(j\omega n)$
- To find the autocorrelation function of $v(n)$, we need to ~~keep~~ know the auto-correlation function of $v(n)$.
- We assume that the noise has a specialized auto-correlation function
$$E[v(n)v^*(n+k)] = \begin{cases} \sigma_v^2, & k=0 \\ 0, & k \neq 0 \end{cases}$$
- Therefore $v(n)$ is White Noise.

- Since we have assumed that source of complex sinusoid and thermal noise are independent, therefore the auto-correlation function of $v(n)$ equals the sum of the autocorrelation functions of its two individual components.
- Therefore, the auto-correlation function of the process $v(n)$ for lag ' k '

$$r(k) = E[v(n)v^*(n-k)]$$

$$= \begin{cases} |\alpha|^2 + \sigma_v^2 & , k=0 \\ |\alpha|^2 \exp(j\omega k) & , k \neq 0 \end{cases}$$

Where $|\alpha|$ is the magnitude of the complex amplitude α

- Given the series of samples $v(n), v(n-1), \dots, v(n-M+1)$ we may express the correlation matrix of $v(n)$ as

$$R = |\alpha|^2 \begin{bmatrix} 1 + 1/p & e^{j\omega} & \dots & e^{j\omega(M-1)} \\ e^{-j\omega} & 1 + 1/p & \dots & e^{j\omega(M-2)} \\ \vdots & \vdots & \ddots & \vdots \\ e^{-j\omega(M-1)} & e^{-j\omega(M-2)} & \dots & 1 + 1/p \end{bmatrix}_{M \times M}$$

Where $p = \frac{|\alpha|^2}{\sigma_v^2}$ is the signal-to-noise ratio (SNR)

Example

Consider the idealized case of a noiselss sinusoid of angular frequency ω .

Assume that the time series consists of 3 uniformly spaced samples from this sinusoid.

Hence, setting SNR $\rightarrow \rho \rightarrow \infty$ and the number of samples $M=3$

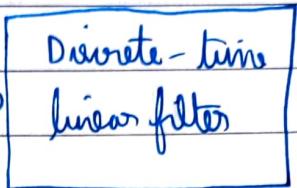
$$R = |\alpha|^2 \begin{bmatrix} 1 & e^{j\omega} & e^{j2\omega} \\ e^{-j\omega} & 1 & e^{j\omega} \\ e^{-j2\omega} & e^{-j\omega} & 1 \end{bmatrix}$$

Note: $|R|=0$

Stochastic Models

"Shock"

Drawn from purely
random process
 $v(n)$



Sample of discrete-time
stochastic process
 $v(n)$

- The representation of a stochastic process by a model dates back to an idea by Yule (1917)
- A time series $v(n)$ consisting of highly correlated observations may be generated by applying a series of statistically independent "shocks" to a linear filter.
- The shocks are random variables drawn from a fixed distribution that is usually assumed to be Gaussian with zero mean and constant variance.
- Input is white Gaussian noise

$$E[v(n)] = 0 \quad \forall n$$

$$\text{and } E[v(n)v^*(\mu)] = \begin{cases} \sigma_v^2, & \mu = n \\ 0, & \text{otherwise} \end{cases}$$

where σ_v^2 is the noise variance

In general, the Time domain input-output relation for the stochastic model shown in the Figure may be described as

$$\begin{bmatrix} \text{Present value} \\ \text{of model output} \end{bmatrix} + \begin{bmatrix} \text{linear combination} \\ \text{of past values} \\ \text{of model output} \end{bmatrix} = \begin{bmatrix} \text{linear combination of} \\ \text{present and past} \\ \text{values of model input} \end{bmatrix}$$

- We may identify 3 types of linear stochastic models
- 1. Auto-regressive (AR) model, in which NO PAST VALUES of the model input are used.
- 2. Moving average (MA) models in which NO PAST VALUES of the model output are used

Auto-Regressive Models (AR)

The time series $v(n), v(n-1), \dots, v(n-M)$ represents the realization of an AR process of order M if it satisfies the difference equation -

$$v(n) + a_1^* v(n-1) + \dots + a_{M-1}^* v(n-M) = v(n)$$

- Reason for calling this auto-regressive -

$$v(n) = w_1^* v(n-1) + w_2^* v(n-2) + \dots + w_M^* v(n-M) + v(n)$$

where $w_k = -a_k$

Thus, $y = \sum_{k=1}^M w_k^* n_k + v$

- LHS of the first equation represents the convolution of the input sequence $\{v(n)\}$ and the sequence of parameters $\{a_n^*\}$

$$\sum_{k=0}^M a_k^* v(n-k) = v(n), \quad (a_0 = 1)$$

Taking Z-transform on both sides, we get

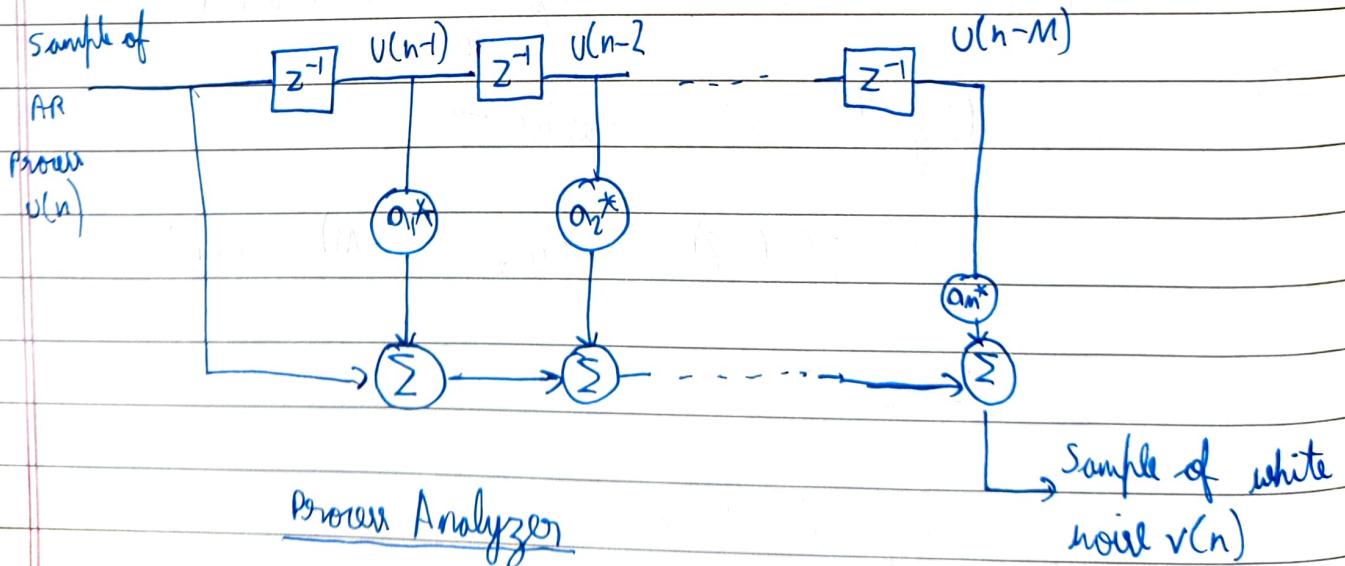
$$H_A(z) U(z) = V(z)$$

where $V(z) = \sum_{n=0}^{\infty} v(n) z^{-n}$

Let $H_A(z) = \sum_{n=0}^M a_n^* z^{-n}$ denote the Z-transform of the ~~sequence~~ $\{a_n^*\}$

Let $U(z) = \sum_{n=0}^{\infty} u(n) z^{-n}$ denote the Z-transform of the input sequence $\{u(n)\}$

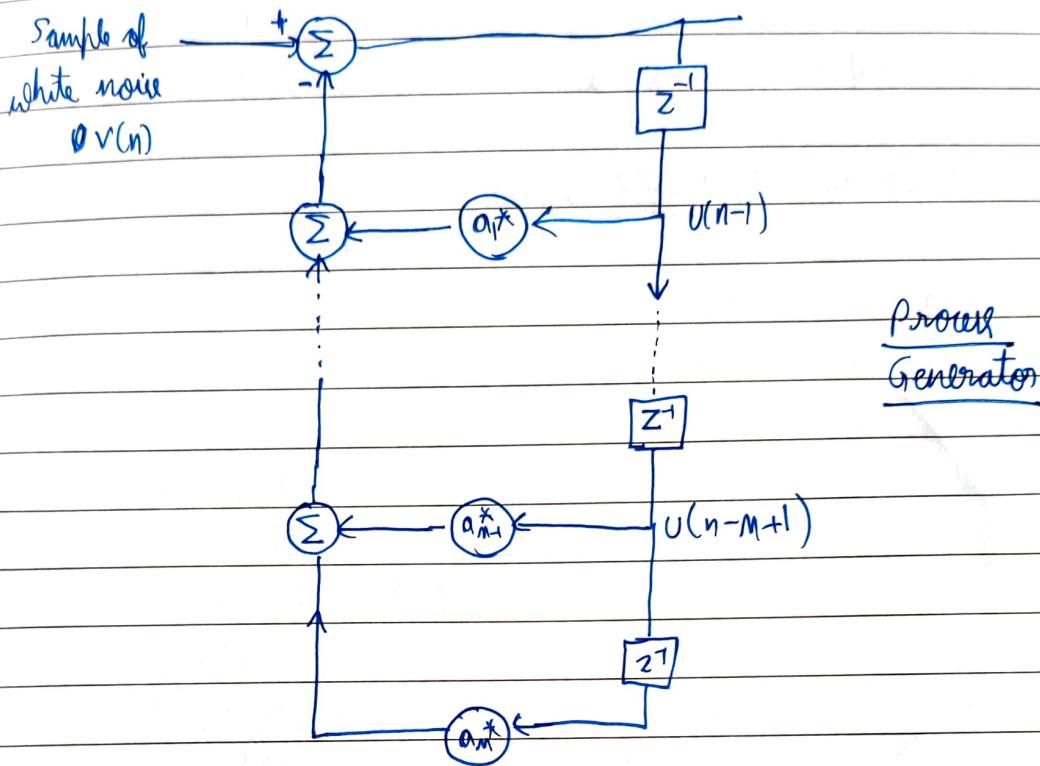
Case I



The above filter has a transfer function $H_A(z) = \frac{V(z)}{U(z)}$

Impulse response of the AR process analyzer has finite duration.

Case II



- The filter represents a process generator where IF:

$$\Rightarrow H_g(z) = \frac{U(z)}{V(z)} = \frac{1}{H_a(z)} = \frac{1}{\sum_{n=0}^M a_n z^{-n}}$$

- The ~~infinite~~ impulse response of an AR process generator has infinite duration.
- The AR process generator is an All-Pole filter since $H_g(z)$ can be completely defined by specifying the location of its poles.

$$H_g(z) = \frac{1}{(1-p_1 z^{-1})(1-p_2 z^{-1}) \dots (1-p_m z^{-1})}$$

Where p_1, p_2, \dots, p_m are the poles of $H_g(z)$

- The poles are defined by the roots of the characteristic equation

$$\cancel{1 + a_1 z^{-1} + a_2 z^{-2} + \dots + a_n z^{-N}} = 0$$

- For the all-pole AR process generator to be stable, the roots of the above equation must all lie inside the unit circle in the z -plane.

Study WOLD DECOMPOSITION yourself

Moving Average Model

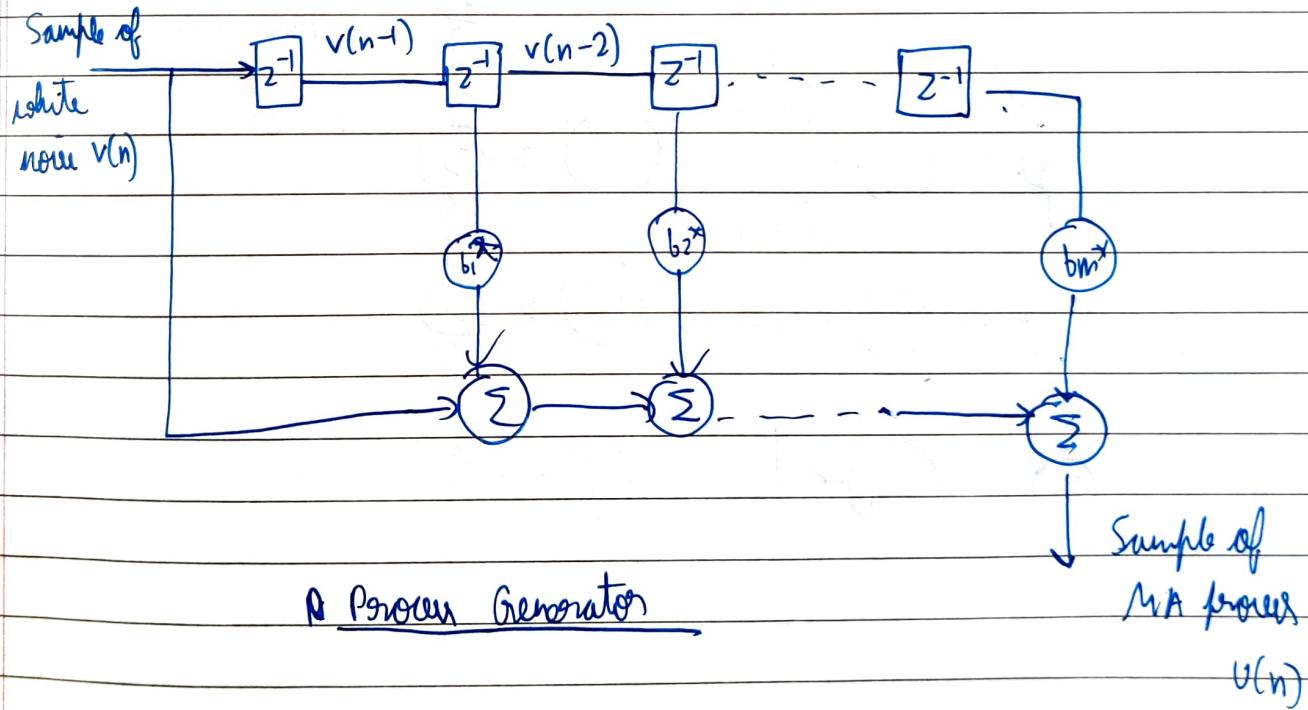
In the MA model, the time linear filter consists of an all-zero filter driven by white noise.

The resulting random $v(n)$ produced at the output can be ~~wrote~~ written as

$$\text{u} - \\ v(n) = v(n) + b_1 * v(n-1) + \dots + b_K * v(n-K)$$

Where b_1, b_2, \dots, b_K are called MA ~~for~~ parameters and $v(n)$ is white noise of zero mean and variance σ_v^2 .

- The order of the MA process is K



- Filters used in the generation and analysis of an MA process are the opposite of those used in the case of an AR process

Auto-Regressive Moving Average (ARMA) Model

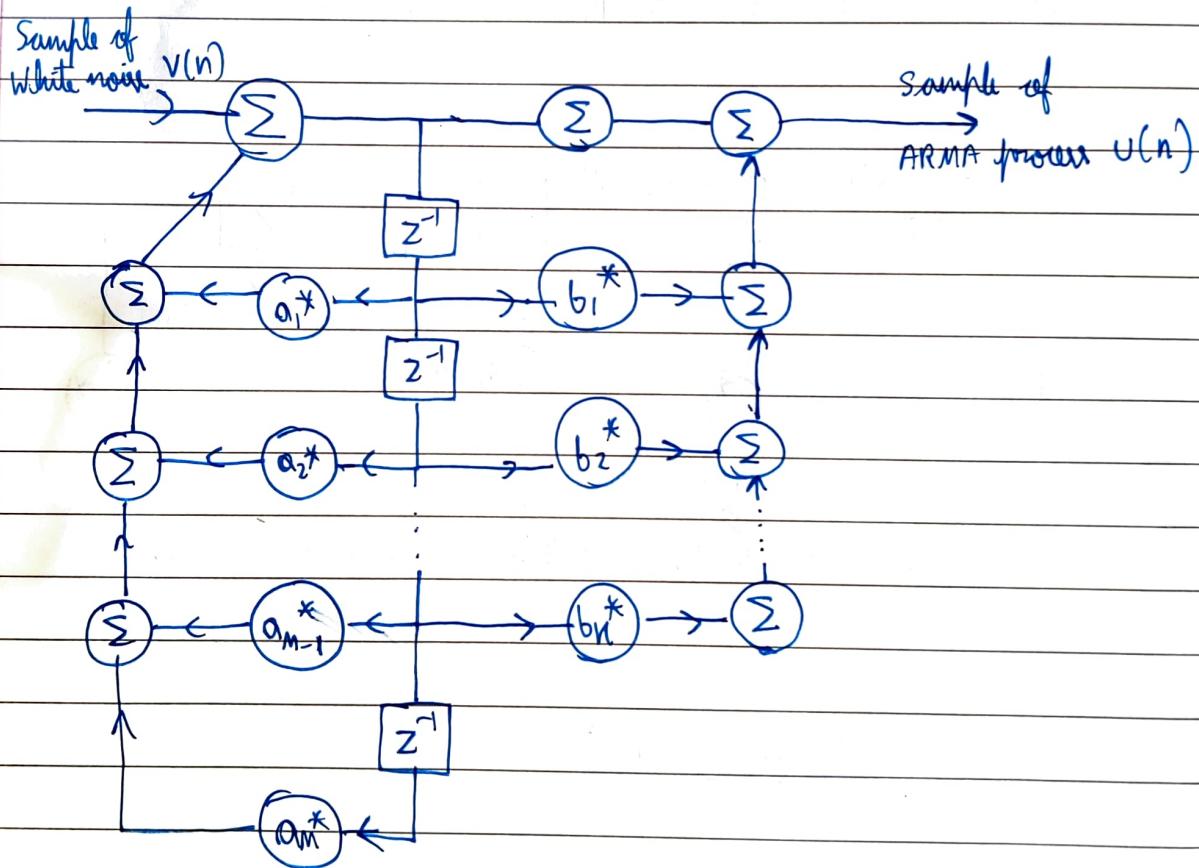
$$v(n) + a_1^* v(n-1) + \dots + a_m^* v(n-m) \\ = v(n) + b_1^* v(n-1) + \dots + b_n^* v(n-k)$$

When a_1, a_2, \dots, a_m and b_1, b_2, \dots, b_n are called ARMA parameters

Order of ARMA process = (m, k)

Assuming $m > k$, ARMA model of order (m, k)

Process Generator



- From a computational point of view, the AR models have an advantage over MA and ARMA models.
- Computation of the AR coefficients involves a system of linear equations known as YULE-WALKER equations.

On the other hand, computation of MA and ARMA coefficients require solving
of

WOLD Decomposition (1938)

The wide application of AR model is fixed, by a fundamental theorem
of Time-series analysis.

Theorem : Any stationary discrete-time stochastic process may be
decomposed into a sum of a general linear process and a
regular predictable process with these two process being
uncorrelated with each other.

$$x(n) = \cancel{v(n)} + s(n)$$

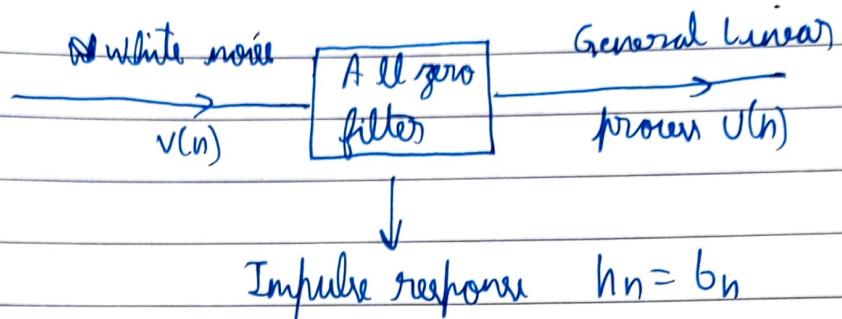
Where $v(n)$ is a general linear process represented by the MA model.

$$v(n) = \sum_{k=0}^{\infty} b_k v(n-k)$$

with $b_0 = 1$ and $\sum_{k=0}^{\infty} |b_k|^2 < \infty$ where $v(n)$ is white noise

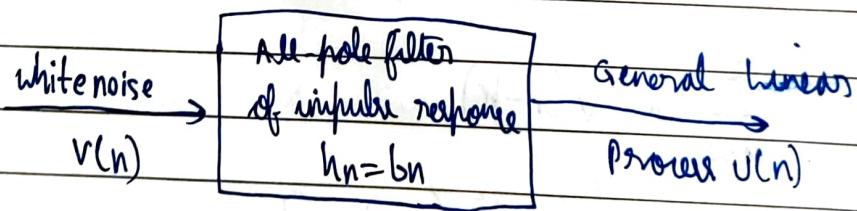
uncorrelated with $s(n)$ i.e., $E(v(n)s^*(k)) = 0 \forall n, k$

$s(n)$ is a predictable process, that means, it can be predicted from its own past with zero prediction variance.



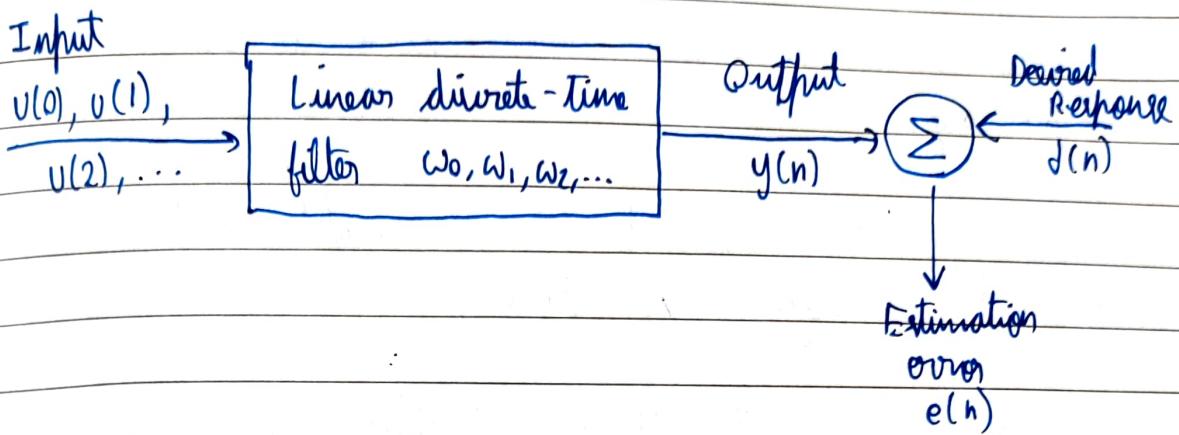
- The general linear process $u(n)$ may be generated by feeding ~~an~~ an all-zero filter with white noise $v(n)$.
- The zeroes of the transfer function of the filter equal the roots of the equation

$$B(z) = \sum_{n=0}^{\infty} b_n * z^{-n} = 0$$
- A solution of particular interest is an all-zero filter that is 'minimum phase', i.e., all the zeroes of the polynomial $B(z)$ lie inside the unit circle.
- In such a case, we may replace the all-zero filter with an equivalent all-pole filter that has the same impulse response $h_n = b_n$



Therefore except for a predictable component, a stationary, discrete-time stochastic process may ~~not~~ also be represented as an AR process of the appropriate order subject to the restriction on $B(z)$.

Wiener Filter



Linear Optimum Filtering

- Filter input consists of a time series $u(0), u(1), u(2), \dots$ and the linear discrete-time filter is itself characterized by the impulse response $w(0), w(1), \dots$.
- At time instant 'n', the filter produces an output $y(n)$. This output provides an estimate of the desired response $d(n)$.
- Since the filter input and the desired response are single realizations of respective stochastic processes, the estimation is accompanied by an error with its own statistical characterization.
- Requirement is to make the estimation error as small as possible in some statistical sense.

$$e(n) = d(n) - y(n)$$

≈

Two restrictions are placed on the filter :

1. The filter is assumed to be linear to make mathematical analysis easy.
2. The filter operates in discrete-time to enable the filter to be implemented in software/hardware.

The final details of the filter specification depends upon -

1. Whether the impulse response has finite or infinite duration.
↳ FIR or IIR filter?
2. The type of statistical criterion used for optimization.
↳ Mathematical tractability

Criterion to choose for Statistical Optimization

We consider optimizing the filter design by minimizing a cost function or index of performance selected from -

1. Mean-square value of the estimation error.
2. Expectation of the absolute value of the estimation error.
3. Expectation of third or higher power of the estimation error.

This ~~choice~~ choice results in a second order dependence for the cost function on the unknown coefficients in the impulse response of the filter.

Also, the cost function has a single minimum that UNIQUELY defines the optimum statistical design of the filter.

We summarize the problem as -

Design a linear discrete-time filter whose output $y(n)$ provides an estimate of the desired response $d(n)$, given a set of input samples $v(0), v(1), \dots$ such that the mean-square value of the estimation error $e(n) = d(n) - y(n)$ is minimized.

Principle of Orthogonality

The input is denoted by the time series $v(0), v(1), v(2), \dots$ and the impulse response of the filter is denoted by w_0, w_1, w_2, \dots both of which are assumed to be complex valued and of infinite duration.

The filter output at discrete-time 'n' is defined by the linear convolution sum -

$$y(n) = \sum_{k=0}^{\infty} w_k^* v(n-k) , \quad n=0, 1, 2, \dots$$

- In the complex terminology $w_k^* v(n-k)$ represents the scalar version of an inner product
- The purpose of the filter is to produce an estimate of the desired response $d(n)$.
- We assume that the filter input and the desired response ~~are~~ are single realizations of jointly WSS process, both with zero mean.

Estimation of $d(n)$ is naturally accompanied by an error defined as $e(n) = d(n) - y(n)$

- We define the cost function as the mean squared error.

$$\begin{aligned} J &= E [e(n) e^*(n)] \\ &= E [|e(n)|^2] \end{aligned}$$

- We have to minimize the above cost function.

- For example input data, the filter coefficients are in general complex functions.

- Let the k^{th} filter coefficient $w_k = a_k + j b_k$, $k=0, 1, 2, \dots$

- we define a gradient operator the k^{th} element of which is written in terms of first order partial derivative

$$\nabla_k = \frac{\partial}{\partial a_k} + j \frac{\partial}{\partial b_k}, \quad k=0,1,2$$

- Applying this gradient operator ∇ to the cost function J , we obtain a multi-dimensional complex gradient vector ∇J , the k^{th} element of which is

$$\nabla_k J = \frac{\partial J}{\partial a_k} + j \frac{\partial J}{\partial b_k}, \quad k=0,1,2.$$

- For the cost function J to attain its minimum value, all the elements of the gradient vector ∇J must be simultaneously set to zero

$$\nabla_k J = 0, \quad k=0,1,2, \dots$$

- Under this set of conditions, the filter is said to be optimum in the mean square sense. MMSE - Minimum Mean Squared Error Optimization

$$\text{Therefore, } \nabla_k J = E \left[\frac{\partial c(n)}{\partial a_k} e^*(n) + \frac{\partial}{\partial a_k} e^*(n) e(n) \right. \\ \left. + \frac{\partial}{\partial b_k} e(n) j e^*(n) + \frac{\partial}{\partial b_k} e^*(n) j e(n) \right]$$

H.W. Derive these 4 terms J and report the answers

Answer in next page -

$$\frac{\partial}{\partial a_n} e(n) = -v(n-k)$$

$$\frac{\partial}{\partial b_n} e(n) = j v(n-k)$$

$$\frac{\partial}{\partial a_n} e^*(n) = -v^*(n-k)$$

$$\frac{\partial}{\partial b_n} e^*(n) = -j v^*(n-k)$$

Substituting the above partial derivatives into $\nabla_n J$ and cancelling common terms.

$$\nabla_n J = -2E[v(n-k)e^*(n)]$$

- Let e_0 denote the spherical value of the estimation error that results when the filter operates in the optimum condition.

Therefore $E[v(n-k)e_0^*(n)] = 0, \quad n=0,1,2,\dots$

- In other words, the necessary and sufficient condition for the cost function J to attain its minimum value is for the corresponding value of the estimation error $e_0(n)$ to be orthogonal to each input sample that enters into the estimation of the desired response at time 'n'.

- This constitutes the Principle of Orthogonality.

Corollary to 1 Principle of Orthogonality

- Examining the correlation between filter output $y(n)$ and the estimation error $e(n)$.

$$\begin{aligned} E[y(n)e^*(n)] &= E\left[\sum_{n=0}^{\infty} w_n^* u(n-k) e^*(n)\right] \\ &= \sum_{n=0}^{\infty} w_n^* E[u(n-k)e^*(n)] \\ &= 0 \end{aligned}$$

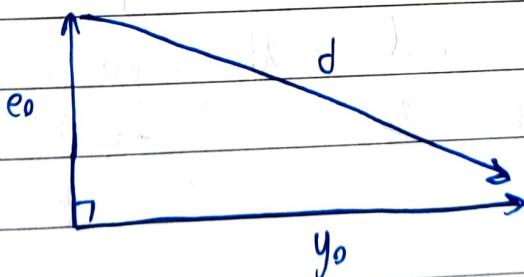
- Let $y_0(n)$ denote the output produced by the filter optimized in the MSE sense, with $e_0(n)$ denoting the corresponding estimation error.

- Hence, using the principle of orthogonality $E[y_0(n)e^*(n)] = 0$

- Let $\hat{d}(n|U_n)$ denote the estimate of the desired response in the MSE sense given the input data that span the space U_n upto the time n ,
then

$$\hat{d}(n|U_n) = y_0(n)$$

Geometric Interpretation



Statistician's Pythagorean Theorem

Wiener - Hopf Equations

- The principle of orthogonality specifies the necessary and sufficient condition for the optimum operation of the filter.

We may write : $E[u(n-k) \left(d^*(n) - \sum_{i=0}^{\infty} w_i u^*(n-i) \right)] = 0$
for $k=0, 1, 2, \dots$

where w_i is the i^{th} coefficient in the impulse response of the optimum filter.

Re-arranging the terms, we get $\sum_{i=0}^{\infty} w_i E[u(n-k) u^*(n-i)]$
 $= E[u(n-k) d^*(n)]$ for $k=0, 1, 2, \dots$

The two expectations in the above equation can be interpreted as -

1. $E[u(n-k) u^*(n-i)]$ is equal to the auto-correlation function of the filter input by a lag of $i-k$

$$r(i-k) = E[u(n-k) u^*(n-i)]$$

2. $E[u(n-k) d^*(n)]$ is equal to the cross-correlation between the filter input $u(n-k)$ and the desired output response $d(n)$ for a lag of $-k$.

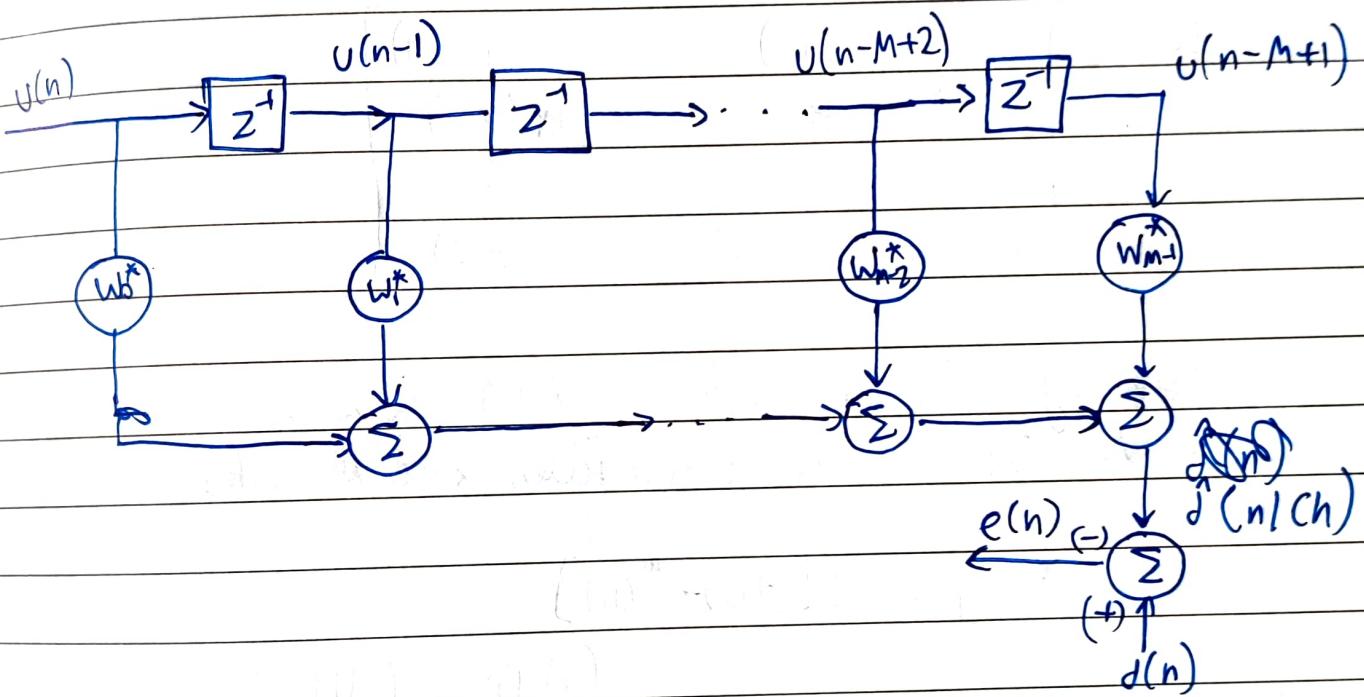
We may thus express $p(-k) = E[u(n-k) d^*(n)]$

Therefore, we get an infinitely large system of equations as the necessary and sufficient condition for the optimality of

$$\sum_{i=0}^{\infty} w_i r(i-k) = p(-k) \quad , \quad k=0,1,2,\dots$$

These equations are called Wiener-Hopf ~~equations~~ equations

Transversal Filter



Accordingly, the Wiener-Hopf equations reduce to a system of M simultaneous linear equations.

$$\sum_{i=0}^{M-1} w_i r(i-k) = p(-k) \quad , \quad k=0,1,2,\dots, M-1$$

where w_i are the optimum weights of the filter.

Matrix Formulation

Let R denote the $M \times M$ correlation matrix of the top infinite input $v(n), v(n-1), \dots, v(n-M+1)$

$$\text{in the transversal filter } R = E[\underline{v}(n) \underline{v}^H(n)]$$

$$\text{where } \underline{v}(n) = [v(n), v(n-1), \dots, v(n-M+1)]^T$$

$$R = \begin{bmatrix} r(0) & r(1) & \cdots & r(M-1) \\ r^*(1) & r(0) & \cdots & r(M-2) \\ \vdots & \vdots & \ddots & \vdots \\ r^*(M-1) & r^*(M-2) & \cdots & r(0) \end{bmatrix}$$

Correspondingly, let p denote the $M \times 1$ cross-correlation vector between the top infinite of the filter and the desired response of $d(n)$

$$p = E[\underline{v}(n) d^*(n)]$$

In the expanded form,

$$p = [p(0), p(-1), \dots, p(1-M)]$$

Therefore, the Wiener-Hopf equation can be written in the matrix form as

$$R \underline{w}_0 = \underline{p}$$

$$\text{where } \underline{w}_0 = [w_{00}, w_{01}, \dots, w_{0M-1}]^T$$

$$\therefore \underline{w}_0 = R^{-1} \underline{p}$$

Error Performance Surface

We can also derive the Wiener-Hopf equations by examining the dependence of the cost function J on the tap weights of the transversal filter.

We write the estimation error as

$$e(n) = d(n) - \sum_{n-k=0}^{m-1} w_k^* u(n-k)$$

The cost function of the transversal filter is

$$J = E[e(n)e^*(n)]$$

$$= E[|d(n)|^2] - \sum_{n=0}^{m-1} w_k^* E[u(n-k) d^*(n)]$$

$$- \sum_{n=0}^{m-1} w_n E[u^*(n-k) d(n)] + \sum_{n=0}^{m-1} \sum_{i=0}^{m-1} w_n^* w_i E[u(n-k) u(n-i)^*]$$

For the first expectation we have $\sigma_d^2 = E[|d(n)|^2]$

where σ_d^2 is the variance of the desired response $d(n)$, assumed to be of zero mean.

For the 2nd and 3rd terms in RHS we have

$$p(-k) = E[u(n-k) d^*(n)]$$

$$p^*(-k) = E[u^*(n-k) d(n)]$$

For the 4th expectation, $r(i-k) = E[u(n-k) u^*(n-i)]$

Thus we may write the cost function as

P.T.O.

$$J = \sigma_d^2 - \sum_{n=0}^{M-1} w_n^* p(-n) - \sum_{n=0}^{M-1} w_n p(-n) + \sum_{n=0}^{M-1} \sum_{i=0}^{M-1} w_n^* w_i r(i-n)$$

tap

- When the tap inputs of the transversal filter and the desired response are jointly stationary, the cost function or mean-squared error is precisely a second-order function of the tap weights for the filter.
- We may visualize the dependency of J on the tap weights w_0, w_1, \dots, w_{M-1} as a bowl-shaped $(M+1)$ dimensional surface with M degrees of freedom represented by tap weights.
- This is known as error performance surface, which has a ^(EPS) UNIQUE minimum.
- At the bottom or minimum point of the EPS, the cost J attains its minimum value J_{\min} .
- At this point, the gradient vector ∇J is identically zero.

$$\nabla_k J = 0, \quad k=0, 1, \dots, M$$

$$- \text{We write } w_k = a_k + j b_k$$

$$\begin{aligned} \nabla_k J &= \frac{\partial J}{\partial a_k} + j \frac{\partial J}{\partial b_k} \\ &= -2p(-k) + 2 \sum_{i=0}^{M-1} w_i r(i-k) \end{aligned}$$

$$- \text{Putting } \nabla_k J = 0, \text{ we obtain } \sum_{i=0}^{M-1} w_i r(i-k) = p(-k), \quad k=0, 1, \dots, M$$

Minimum Mean-Square Error

Let $\hat{f}(n|U_n)$ denote the estimate of the desired signal response $f(n)$, produced at the output of the transversal filter that is optimum in the mean square sense, given the tap inputs $u(n), u(n-1), \dots, u(n-M+1)$, that span the space U_n .

$$\begin{aligned}\hat{f}(n|U_n) &= \sum_{k=0}^{M-1} w_0^* u(n-k) \\ &= \underline{w}_0^H \underline{u}(n)\end{aligned}$$

Where $\underline{w}_0 = [w_0, w_1, \dots, w_{M-1}]^T$
 and $\underline{u}(n) = [u(n), u(n-1), \dots, u(n-M+1)]^T$

Hence, the variance of $\hat{f}(n|U_n)$

$$\begin{aligned}\sigma_{\hat{f}}^2 &= E[\underline{w}_0^H \underline{u}(n) \underline{u}^H(n) \underline{w}_0] \\ &= \underline{w}_0^H E[\underline{u}(n) \underline{u}^H(n)] \underline{w}_0 \\ &= \underline{w}_0^H R \underline{w}_0\end{aligned}$$

$$\sigma_{\hat{f}}^2 = p^H \underline{w}_0 = p^H R^{-1} p$$

$$\begin{aligned}J_{\min} &= \sigma_d^2 - \sigma_{\hat{f}}^2 \\ &= \sigma_d^2 - \underline{w}_0^H R \underline{w}_0 \\ &= \sigma_d^2 - p^H \underline{w}_0 \\ &= \sigma_d^2 - p^H R^{-1} p \quad \rightarrow \text{Error for optimal weight}\end{aligned}$$

Canonical Form of the Error Performance Surface

$$J = \sigma_d^2 - \sum_{n=0}^{M-1} w_n^* p(-n) - \sum_{k=0}^{M-1} w_k p^*(-k) + \sum_{n=0}^{M-1} \sum_{i=0}^{n-1} w_i^* w_i r(i-k)$$

$$J(\underline{w}) = \sigma_d^2 - \underline{w}^H p - p^H \underline{w} + \underline{w}^H R \underline{w}$$

Express $J(\underline{w})$ as a "perfect square" in \underline{w} by

$$J(\underline{w}) = \sigma_d^2 - p^H R^{-1} p + (\underline{w} - R^{-1} p)^H \cdot R (\underline{w} - R^{-1} p)$$

$$\text{Therefore, } \min J(\underline{w}) = \sigma_d^2 - p^H R^{-1} p \quad \text{for } \underline{w}_0 = R^{-1} p$$

$$\text{We can also write } J(\underline{w}) = J_{\min} + (\underline{w} - \underline{w}_0)^H R (\underline{w} - \underline{w}_0)$$

$$\text{Therefore } J(\underline{w}_0) = J_{\min}$$

- Using eigen-decomposition, express $R = Q \Lambda Q^H$ where Λ is a diagonal matrix consisting of eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_M$ of R and the matrix Q has for its columns the eigenvectors q_1, q_2, \dots, q_M associated with these eigenvalues.

$$\text{Therefore, } J = J_{\min} + (\underline{w} - \underline{w}_0)^H Q \Lambda Q^H (\underline{w} - \underline{w}_0)$$

- Define a transformed version of the difference between the optimum solution \underline{w}_0 and the tap weight vector \underline{w} as

$$\underline{v} = Q^H (\underline{w}_0 - \underline{w})$$

$$\text{Hence we can now write } J = J_{\min} + \underline{v}^H \Lambda \underline{v}$$

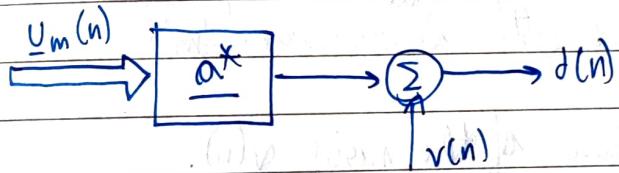
This new formulation of the mean-square error contains no cross-product becomes

$$\begin{aligned} J &= J_{\min} + \sum_{n=1}^M \lambda_n v_n v_n^* \\ &= J_{\min} + \sum_{n=1}^M \lambda_n |v_n|^2 \end{aligned}$$

where v_n is the n^{th} component of v

This representation of the error performance surface is useful because all the components of the transformed coefficient vector v constitute the Principal Axes of the surface.

Multiple Linear Regression Model



The maximum mean-square error of the Wiener filter

$$J_{\min} = \sigma_d^2 - f^H R^{-1} f$$

applies to a transversal filter configuration of a prescribed length (i.e., no. of tap weights)

Suppose, it is possible to adjust the length of the transversal filter as a design parameter, can we reduce J_{\min} to an irreducible level?

The answer to this fundamental question depends on the model describing the relationship between $d(n)$ and the input vector $u(n)$.

- We make the following 3 assumptions -

- i) The model is linear
- ii) The observable data are noisy (measurable)
- iii) The ~~noise~~ noise is additive and white.

$$d(n) = \underline{\alpha}^H \underline{u}_m(n) + v(n)$$

where $\underline{\alpha}$ denotes an unknown parameter vector of the model, $\underline{u}_m(n)$ denotes the input vector (regressor) and $v(n)$ is the additive white noise that is statistically independent and independent of $\underline{u}_m(n)$.

Parameter $\underline{\alpha}$ is also referred to as the regression vector. Vectors $\underline{u}_m(n)$ and $\underline{\alpha}$ are both of dimension 'm'.

Let σ_v^2 denote the variance of the noise $v(n)$.

Then the variance of the observable data $d(n)$ supplying the desired response is given by

$$\begin{aligned}\sigma_d^2 &= E[d(n) d^*(n)] \\ &= \sigma_v^2 + \underline{\alpha}^H R_m \underline{\alpha}\end{aligned}$$

where $R_m = E[\underline{u}_m(n) \underline{u}_m^H(n)]$

If the $m \times m$ correlation matrix of the I/P vector.

Now consider a Wiener filter that operates on an input vector $\underline{u}(n)$ and desired response $d(n)$ to produce a minimum mean square error $J_{\min}(M)$ which is adjustable by varying filter length 'M'.

The input vector $\underline{v}(n)$ and the corresponding tap-weight vector \underline{w}_0 are both $M \times 1$ vectors.

Substituting $\sigma_d^2 = \sigma_v^2 + \underline{\alpha}^H R_m \underline{\alpha}$

into $J_{\min} = \sigma_d^2 - \underline{w}_0^H R \underline{w}_0$ we get

$$J_{\min}(M) = \sigma_v^2 + (\underline{\alpha}^H R_m \underline{\alpha} - \underline{w}_0^H R \underline{w}_0)$$

Therefore we have 3 regimes of model selection

Underfitted model ($M < m$)

In this regime, the Wiener filter exhibits better performance with increasing M for a prescribed value of ' m '.

In particular, the minimum MSE decreases quadratically with filter length M and

$$J_{\min}(0) = \sigma_v^2 + \underline{\alpha}^H R_m \underline{\alpha}$$

is the worst possible value for J_{\min} from which the MMSE falls as M is increased.

Critically fitted model ($M=m$)

At the critical point, $M=m$, Wiener filter is perfectly matched to the regression model in the sense

$$\underline{w}_0 = \underline{\alpha} \text{ and } R = R_m$$

Correspondingly, the MMSE of the Wiener filter attains the irreducible value.

$$J_{\min}(M) = \sigma_v^2$$

3. Overfitted Model : ($M > m$)

When the length of the Wiener filter is greater than the model order 'm', the tail end of the tap-weight vector of the Wiener filter is zero.

$$\text{Therefore, } \underline{w}_0 = \begin{bmatrix} \underline{a} \\ \underline{0} \end{bmatrix}$$

Correspondingly, the tap-weight input vector of the Wiener filter takes the form

$$\underline{v}(n) = \begin{bmatrix} \underline{v}_m(n) \\ \underline{v}_{m-m}(n) \end{bmatrix}$$

Where $\underline{v}_{m-m}(n)$ is an $(M-m) \times 1$ vector made up of part data samples immediately preceding the $m \times 1$ vector $\underline{v}_m(n)$.

- The net result is that, in theory, the same MMSE performance as the critically fitted case is achieved but with a longer filter.

ExampleWhen a regression

$$\underline{\alpha} = \begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \alpha_2 \end{bmatrix}$$

Statistical characterization of the assumed to be real valued is as follows.

- a) The correlation matrix of input vector is

$$R_u = \begin{bmatrix} 1.1 & 0.5 & 0.1 & -0.05 \\ 0.5 & & & \\ 0.1 & & & \\ -0.05 & & & \end{bmatrix}$$

- b) The cross-correlation vector between the input vector $\underline{u}(n)$ and observable ~~st~~ $d(n)$ is

$$\underline{P} = [0.5272, -0.4458,$$

- c) The variance of the observable state $\sigma_d^2 = 0.9486$

- d) Variance of the additive white noise $\sigma_v^2 = 0.1066$

Method of Steepest Descent

- Consider a cost function $J(\underline{w})$ that is continuously differentiable function of some unknown weight vector \underline{w}
- We want to find an optimal solution \underline{w}_0 that satisfies the condition

$$J(\underline{w}_0) \leq J(\underline{w}) \text{ for all } \underline{w}$$

Local Iterative Descent

Starting with an initial guess denoted by $\underline{w}(0)$, generate a sequence of weights $\underline{w}(1), \underline{w}(2), \dots$ such that the cost function $J(\underline{w})$ is reduced at each iteration of the algorithm

$$J(\underline{w}(n+1)) < J(\underline{w}(n))$$

Where $\underline{w}(n+1)$ is the updated value and $\underline{w}(n)$ is the old value of the weight vector.

We "hope" that this algorithm will eventually converge on the optimal value \underline{w}_0 although there is a ~~distinct~~ distinct possibility that the algorithm will diverge unless certain precautions are taken.

- In a simple form of iterative descent known as the "method of steepest descent" successive adjustments made to the weight vector \underline{w} are in the direction of steepest descent - opposite to gradient vector of $J(\underline{w})$

$$\text{Let } \mathbf{g} = \nabla J(\underline{w}) = \frac{\partial}{\partial \underline{w}} J(\underline{w})$$

~~Therefore~~, the steepest descent algorithm is formally described as

$$w(n+1) = w(n) - \frac{1}{2} \mu g(n)$$

where 'n' denotes the iteration, μ is a positive constant called the 'step-size' parameter.

- To go from iteration 'n' to 'n+1', the algorithm

Taylor Series

- Taylor series of a function is an infinite sum of terms that are expressed in terms of the function's derivatives at a single point.
- If zero is the point where the derivatives are considered a Taylor series is also called Maclaurin series.

History, Definition

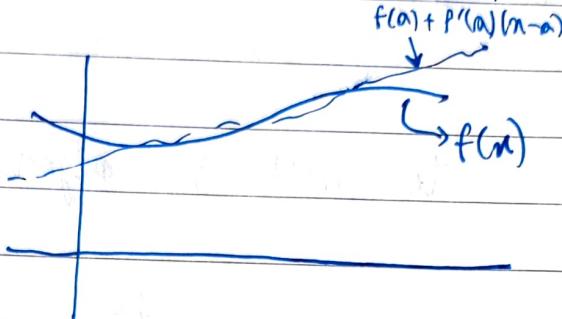
- The Taylor series of a real or complex valued function $f(n)$ that is infinitely differentiable . . .

$$\begin{aligned} f(n) &\cong f(a) + \frac{f'(a)}{1!} (n-a) + \frac{f''(a)}{2!} (n-a)^2 + \dots \\ &\Rightarrow \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (n-a)^n \end{aligned}$$

Visualization

Best linear approximation for $f(\underline{x})$ at $\underline{x} = \underline{a}$ is

$$f(\underline{x}) \approx f(\underline{a}) + f'(\underline{a})(\underline{x} - \underline{a})$$



Multiple Variables

$$f(\underline{x}, \underline{y}) \approx f(\underline{a}, \underline{b}) + (\underline{x} - \underline{a}) f_x(\underline{a}, \underline{b}) + (\underline{y} - \underline{b}) f_y(\underline{a}, \underline{b})$$

$$+ \frac{1}{2} \left\{ (\underline{x} - \underline{a})^2 f_{xx}(\underline{a}, \underline{b}) + 2(\underline{x} - \underline{a})(\underline{y} - \underline{b}) f_{xy}(\underline{a}, \underline{b}) + (\underline{y} - \underline{b})^2 f_{yy}(\underline{a}, \underline{b}) \right\}$$

+

For more than 2 variables, Taylor series of scalar valued fun is

$$f(\underline{x}) \approx f(\underline{a}) + (\underline{x} - \underline{a})^T Df(\underline{a}) + \frac{1}{2} (\underline{x} - \underline{a})^T D^2 f(\underline{a}) (\underline{x} - \underline{a})$$

+

where $Df(\underline{a})$ is the gradient evaluated at $\underline{x} = \underline{a}$ and

$D^2 f(\underline{a})$ is the Hessian matrix

Steepest Descent technique approach of Wiener filter

$$\begin{aligned} \text{The estimation error } e(n) &= d(n) - \hat{d}(n|U_n) \\ &= d(n) - \underline{\omega}^H(n) \underline{u}(n) \end{aligned}$$

$$\text{Where } \underline{\omega}(n) = [w_0(n), w_1(n), \dots, w_{m-1}(n)]^T$$

$$\text{and } \underline{u}(n) = [u(n), u(n-1), \dots, u(n-M+1)]^T$$

If the tap input $\underline{u}(n)$ and the desired response $d(n)$ are jointly stationary, then the MSE or cost function $J(\underline{\omega}(n))$ or simply $J(n)$ at time 'n' is a quadratic function of the tap weight vector.

$$J(n) = \sigma_d^2 - \underline{\omega}^H(n) p - p^H \underline{\omega}(n) + \underline{\omega}^H(n) R \underline{\omega}$$

(from Canonical form of Error Performance Surface)

$$\underline{\nabla} J(\underline{\omega}(n)) = \frac{d}{d\underline{\omega}} J(\underline{\omega}(n))$$

$$\underline{\nabla} J(n) = \left[\begin{array}{c} \frac{\partial J(n)}{\partial a_0(n)} + j \frac{\partial J(n)}{\partial b_0(n)} \\ \frac{\partial J(n)}{\partial a_1(n)} + j \frac{\partial J(n)}{\partial b_1(n)} \\ \vdots \\ \frac{\partial J(n)}{\partial a_{m-1}(n)} + j \frac{\partial J(n)}{\partial b_{m-1}(n)} \end{array} \right] = -2p + 2R\underline{\omega}(n)$$

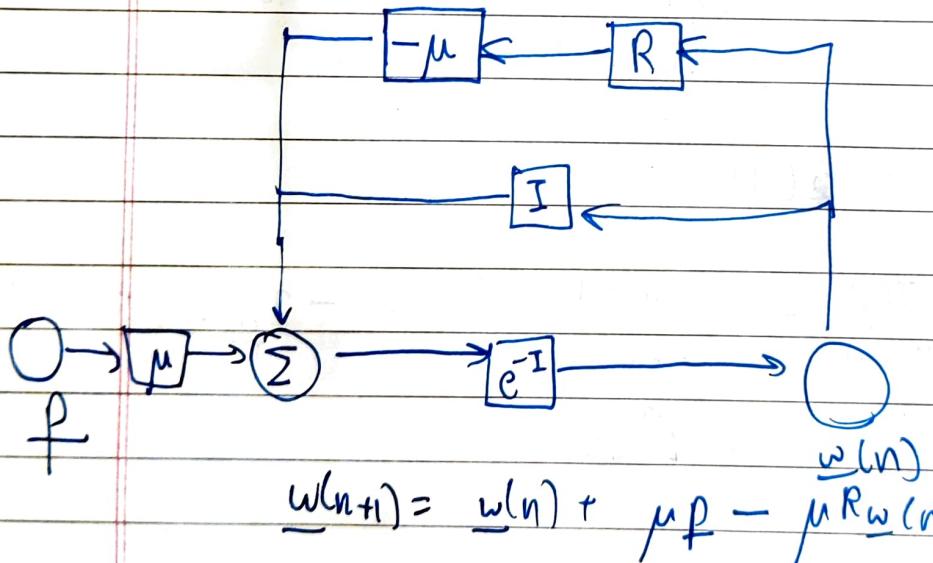
~~scribbled text~~

- We assume that the correlation matrix R and the cross-correlation vector p are known, so that we may compute the gradient vector $\nabla J(n)$ for a given value of the tap weight vector $w(n)$
- Thus, we may ~~not~~ compute the updated value of the tap weight vector $w(n+1)$ by using the recursive relation

$$\underline{w}(n+1) = \underline{w}(n) + \mu [p - R \underline{w}(n)] \quad \text{for } n=0,1,2,\dots$$

which describes the mathematical form of the steepest descent algorithm for Wiener filtering

Signal Flow Graph Representation



Stability of Steepest Descent

To determine the ~~stable~~ condition for the stability of the gradient descent algorithm, we examine Natural Model of it.

In particular, we use the representation of the correlation matrix, R

We know $p = R \underline{w}_0$ and

$$\underline{w}(n+1) = \underline{w}(n) + \mu [p - R \underline{w}(n)]$$

$$\text{which can be written as } \underline{w}(n+1) = \underline{w}(n) + \mu [R \underline{w}_0 - R \underline{w}(n)]$$

$$= \underline{w}(n) + \mu R c(n)$$

$$\begin{aligned} \text{and } \underline{c}(n+1) &= \underline{w}_0 - \underline{w}(n+1) \\ &= \underline{w}_0 - \underline{w}(n) - \mu R c(n) \\ &= \underline{c}(n) - \mu R \underline{c}(n) \\ &= [I - \mu R] \underline{c}(n) \end{aligned}$$

Expressing $R = Q \Lambda Q^H$

Q is called the UNITARY matrix of the orthogonal set of eigenvalues associated with eigen values of R .

Λ is a diagonal matrix containing $\lambda_1, \lambda_2, \dots, \lambda_m$ eigenvalues of the correlation matrix R .

These eigenvalues are positive and real.

We can write $\underline{c}(n+1) = (I - \mu Q \Lambda Q^H) \underline{c}(n)$ since Q is unitary,

$$Q^H = Q^{-1}$$

$$\text{Therefore, } Q^H \underline{c}(n+1) = (I - \mu \Lambda) Q^H \underline{c}(n)$$

Define a new set of coordinates

$$\underline{v}(n) = Q^H \underline{c}(n) = Q^H [\underline{w}_0 - \underline{w}(n)]$$

$$\text{So, } \underline{v}(n+1) = (I - \mu \Lambda) \underline{v}(n)$$

The initial value of $\underline{v}(n)$ is

$$\underline{v}(0) = Q^H [\underline{w}_0 - \underline{w}(0)]$$

Assuming that the initial tap weight vector is zero, i.e., $\underline{w}(0) = 0$
We get

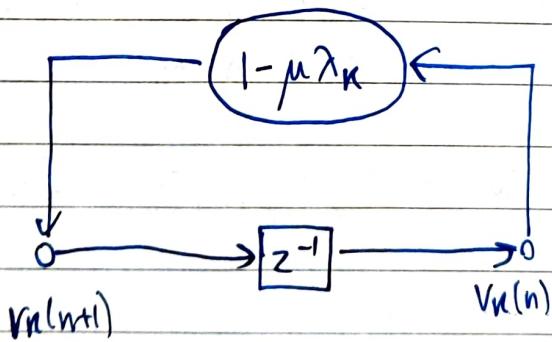
$$\Rightarrow \underline{v}(0) = Q^H [\underline{w}_0]$$

- For the k^{th} natural mode of the steepest descent algorithm we have

$$v_n(n+1) = (1 - \mu \lambda_k) v_n(n)$$

where $k=0, 1, 2, \dots, M$

and λ_k is the k^{th} eigenvalue of R



Signal flow graph representation of the k^{th} natural mode of the steepest descent algorithm.

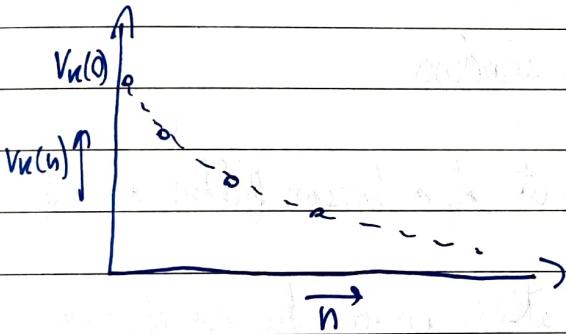
- The above equation is a first order homogeneous difference equation

- Assuming that $v_n(n)$ has the initial value $v_n(0)$, we readily obtain the solution

$$v_n(n) = (1 - \mu \lambda_n)^n v_n(0)$$

where $\mu = 0, 1, 2, \dots, M$

- For stability of convergence of the steepest descent algorithm, the magnitude of the geometric ratio must be less than UNITY for all values of 'n'.
- If $-1 < 1 - \mu \lambda_n < 1$ for all n then as the number of iteration 'n' approaches infinity, the natural mode of the steepest descent algorithm die out irrespective of the ~~and~~ initial conditions.



- Referring to $v(n) = Q^H (\underline{w}_0 - \underline{w}(n))$ it is equivalent to saying that $\underline{w}(n)$ approaches the optimal solution \underline{w}_0 as $n \rightarrow \infty$

- Since eigenvalues of R are real and positive, it follows that a necessary and sufficient condition for the convergence or stability of the steepest descent algorithm is

$$0 < \mu < \frac{2}{\lambda_{\max}}$$

~~CT - 10%~~
~~Assignments - 10% (maybe more)~~

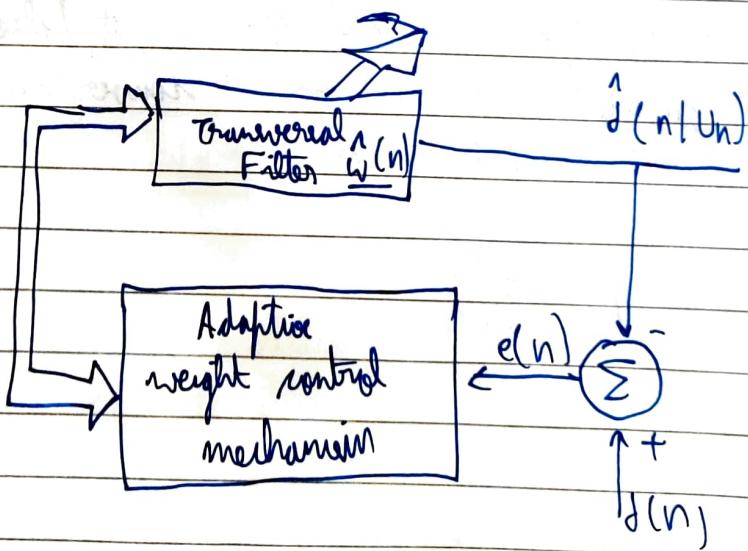
LMS Adaptive Filter

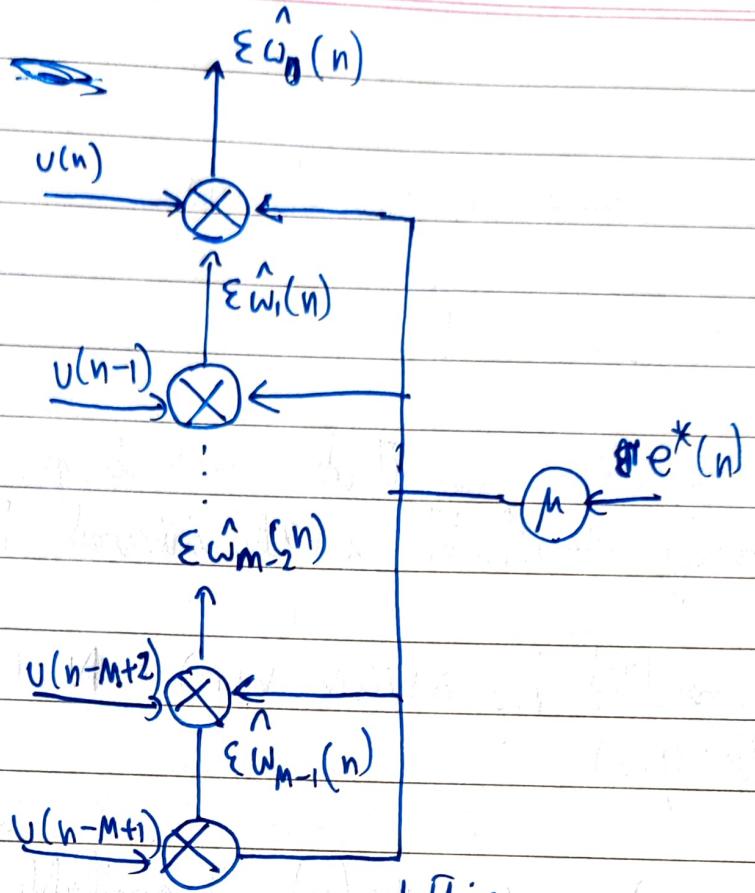
- Do not require measurement of R and f and also does not require any matrix inversion.

Overview

LMS algorithm is a linear adaptive filtering algo which has two processes -

- i) A filtering process which involves
 - a) computing the output of a linear filter in response to an input signal
 - b) generating an estimation error by comparing this output with a desired response.
- ii) An adaptive process, which involves automatic adjustment of the parameters of the filter in accordance with the estimation error



DetailedDetermined structure of the filter weight control mechanism

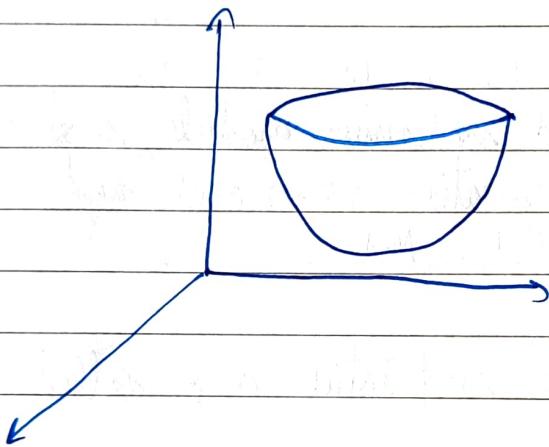
- The tap weights $u(n), u(n-1), \dots, u(n-M+1)$ form the component compensators of the $M \times 1$ tap ~~the~~ input vector $u(n)$.
- Their inputs span a multi ~~the~~ dimensional array shaped denoted by U_n

NOTE The value computed for the vector using the LMS algorithm represents an estimate whose expected value may come close to the Wiener-Solution w_0 (for a wss measurement) as the no. of iterations 'n' approaches infinity.

- The estimation error is defined by $e(n) = d(n) - \hat{e}(n|U_n)$
- The estimation errors $e(n)$ are the tap input vector $\mathbf{U}(n)$ are applied to the control mechanism and the feedback loop around the tap weight is then closed.
- In the adaptive weight-control mechanism module, a scalar version of the inner product of the estimation error $e(n)$ and the tap input $U(n-k)$ computed for $n = 0, 1, \dots, M-1$
- The scaling factor used in ^{this} computation is called the step size
- The result of the scalar inner-product ~~is~~ the correlation $S_{w_k}(n)$ applied to the tap-weight $w_k(n)$ at the iteration $(n+1)$

$$\begin{array}{c} n(n) \\ \downarrow \\ n(n-M+1) \end{array}$$

- Comparing the control mechanism for the LMS algorithm with the method of steepest descent, we observe that the LMS algorithm uses the product $v(n-k) e^*(n)$ as an estimate of descent in the gradient vector $\nabla J(n)$ that characterizes the method of steepest descent.
- Now there is NO expectation operator in all the paths. Therefore, the recursive computation of each tap weight in the LMS algorithm suffers from a gradient noise.



For such an ~~enviroment~~ environment, we know that the method of steepest descent computes a tap vector $w(n)$ that moves down the ensemble average error performance surface along a deterministic trajectory that terminate on the Wiener sol w_0 .

The LMS algo behaves differently due to presence of gradient noise. Rather than terminating on Wiener sol, the tap weight vector $w(n)$ computed by LMS algorithm execute a random motion around the minimum point of the Error Performance Surface. It can be shown that this motion is a form of Brownian motion for small values of μ .

LMS Algorithm

- If it were possible to make exact measurements of the gradient vector $\nabla J(n)$ at each iteration 'n' and if ' μ ' is suitably chosen, then the tap-weight vector computed using the steepest descent algorithm would indeed converge to the optimum Wiener solution.
- However in practice, exact measurements of the gradient vector are not possible as this requires prior knowledge of R and p .
- The simplest choice of estimation is to use instantaneous estimates of R and p that are based on sample values of the tap-input vector and the desired response.

$$\hat{R}(n) = \underline{v}(n) \underline{v}^H(n)$$

$$\hat{p}(n) = \underline{v}(n) p^*(n)$$

- Therefore the instantaneous estimate of the gradient vector is

$$\hat{\nabla} J(n) = \mu \cdot -\frac{2}{\hat{R}(n)} \hat{p}(n) + \frac{2}{\hat{R}(n)} \hat{R} \hat{w}(n)$$

$$= -2 \mu(n) \hat{p}(n) + 2 \underline{v}(n) \underline{v}^H(n) \hat{w}(n)$$

- Substituting this estimate of the gradient vector we get a new recursive solution for updating the tap weight vector

$$\hat{w}(n+1) = \hat{w}(n) + \mu \underline{v}(n) [\hat{p}(n) - \underline{v}^H(n) \hat{w}(n)]$$

- A hat operator over the symbol to distinguish it from the value obtained from steepest descent algorithm.

~~Skipped~~

- We may write the above result in the form of 3 basic relations

1. Filter output $y(n) = \hat{w}^H(n) \underline{u}(n)$

2. Estimation error or error $e(n) = d(n) - y(n)$

$$e^*(n) = d^*(n) - y^*(n)$$

$\leftarrow d^*(n) - \underline{u}^H(n) \hat{w}(n)$

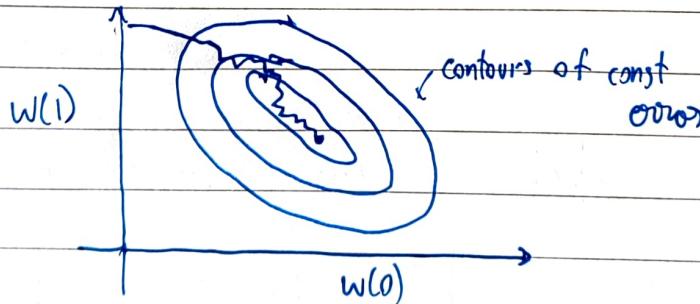
Note: $\langle \underline{a}, \underline{b} \rangle = \langle \underline{b}, \underline{a} \rangle$

so $\underline{y}^H \underline{x} = (\underline{x}^H \underline{y})^* = \underline{y}^H \underline{x}$

3. Tap-weight adjustment

$$\hat{w}(n+1) = \hat{w}(n) + \mu \underline{u}(n) e^*(n)$$

- The term $\mu \underline{u}(n) e^*(n)$ represents the adjustment that is applied to the current estimate of the tap weight vector $\hat{w}(n)$.



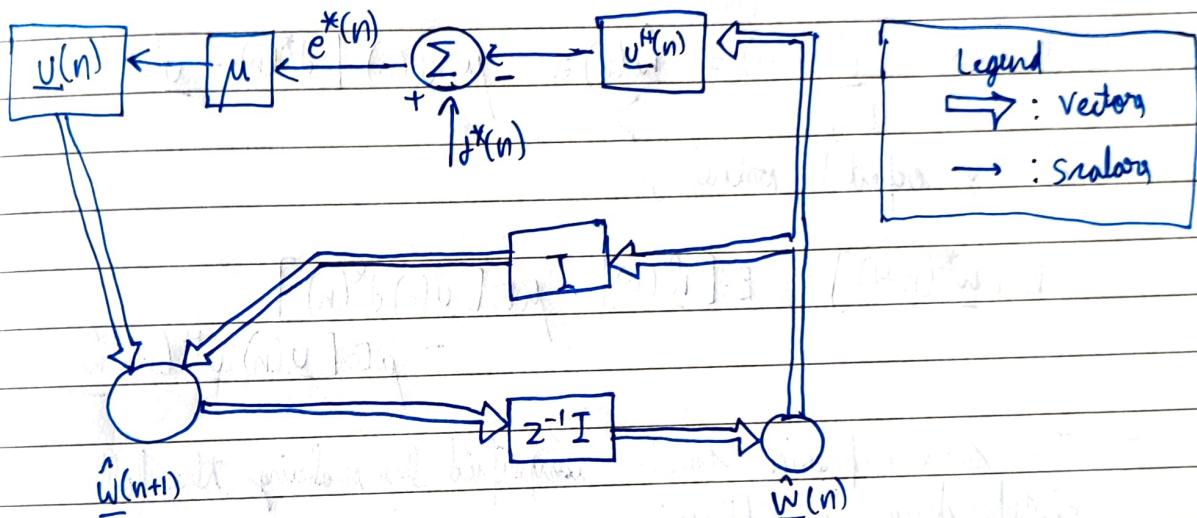
- The LMS algorithm is a member of the family of stochastic gradient algorithms

- In particular, when the LMS Algorithm operates on stochastic inputs, the allowed set of directions along which we step

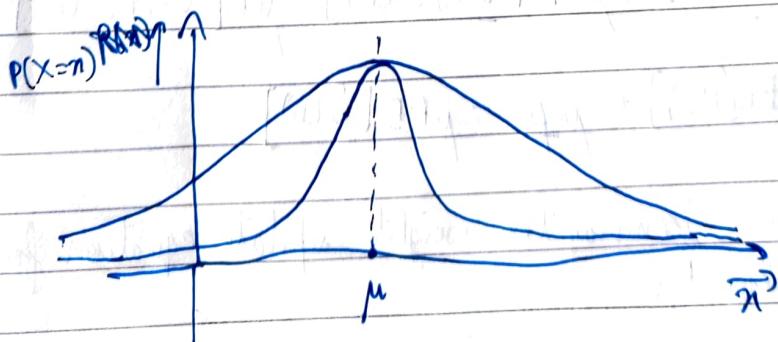
from one iteration cycle to the next is quite random and cannot therefore be thought of consisting of true gradient direction.

- Although the correction (update) that is applied to $\hat{w}(n)$ is generally not aligned with the direction of steepest descent, the correction is on the average in the direction of the steepest descent.

Signal Flow Graph Representation



- The computational complexity of the LMS Algorithm is $O(M)$.
- The instantaneous estimates of R and p have relatively large variances but since LMS algorithm is recursive in nature, the algorithm itself effectively averages these estimates in some sense during the course of adaptation.



Convergence of LMS Algorithm

We begin by assuming that $\underline{v}(n)$ and $d(n)$ are jointly stationary processes.

We will determine when the coefficients $\hat{\underline{w}}(n)$ converge in the mean to $\underline{w}_0 = R^{-1}f$

$$\lim_{n \rightarrow \infty} E[\hat{\underline{w}}(n)] = \underline{w}_0 = R^{-1}f$$

- We have $\hat{\underline{w}}(n+1) = \hat{\underline{w}}(n) + \mu \underline{v}(n) [d^*(n) - \underline{v}^H(n) \hat{\underline{w}}(n)]$

Taking expected values,

$$E[\underline{w}^*(n+1)] = E[\hat{\underline{w}}(n)] + \mu E[\underline{v}(n) d^*(n)] - \mu E[\underline{v}(n) \underline{v}^H(n) \hat{\underline{w}}(n)]$$

- The above expression can be simplified by making the following independence assumption:

- * The data $\underline{v}(n)$ and the weight vector $\hat{\underline{w}}(n)$ are statistically independent

The above equation can now be simplified as -

$$\begin{aligned} E[\hat{\underline{w}}(n+1)] &= E[\hat{\underline{w}}(n)] + \mu E[\underline{v}(n) d^*(n)] - \mu E[\underline{v}(n) \underline{v}^H(n)] E[\hat{\underline{w}}(n)] \\ &= E[\hat{\underline{w}}(n)] + \mu f - \mu R E[\hat{\underline{w}}(n)] \end{aligned}$$

which is in the same form as the weight update equation for the steepest descent algorithm

- Therefore, the analysis for the steepest descent algorithm is applicable to $E[\hat{w}(n+1)]$

- We already know that $\underline{v}(n) = (I - \mu A)^n v(0)$

$$\text{where } \underline{v}(n) = Q^T [w_0 - \underline{w}(n)]$$

- If $E[\underline{w}(n) \underline{v}(n)]$ converges to 0, then $\hat{w}(n)$ will converge to mean.

- Thus we have the following property :

For jointly WSS process, the LMS algo converges in the mean if

$$0 < \mu < \frac{2}{\lambda_{\max}}$$

and the independence assumption is satisfied.

Note that although this bound on μ ensures that $E[\hat{w}(n)]$ converges, it places no constraint on how large the variance of $\hat{w}(n)$ may become.

- The upper bound on μ is expressed in terms of the highest eigenvalue of R but if this matrix is unknown, the λ_{\max} has to be estimated as

$$\lambda_{\max} \leq \sum_{n=0}^{M-1} \lambda_n = \text{tr}[R]$$

- If $\underline{v}(n)$ is WSS, then R is Toeplitz and the trace becomes

$$\text{tr}[R] = M r_v(0) = M E[|v(n)|^2]$$

Therefore the bounds on μ can be replaced as

$$0 < \mu < \frac{2}{M E\{|v(n)|^2\}}$$

Where $E[|v(n)|^2]$ represents the power in $v(n)$

- $\hat{E}[|v(n)|^2] = \frac{1}{M} \sum_{n=0}^{N-1} |v(n)|^2$

- Therefore the following bound step size parameters, μ , for mean-square convergence

$$0 < \mu < \frac{2}{\underline{v}^H(n) \underline{v}(n)}$$

Normalized LMS Algorithm

- A convenient way to incorporate this bound into the LMS adaptive filter is to use a time varying step size

$$\mu(n) = \frac{\beta}{\underline{v}^H(n) \underline{v}(n)} = \frac{\beta}{\|\underline{v}(n)\|^2}$$

where $0 < \beta < 2$

- The weight update relationship is given as

$$\hat{w}(n+1) = \hat{w}(n) + \frac{\beta v(n)}{\|\underline{v}(n)\|^2} \cdot c^*(n)$$

- Note: The effect of the normalization by $\|v(n)\|^2$ is to alter the magnitude but not the direction of the estimated gradient vector.
- we know that in the LMS algorithm the ~~step-size~~ proportion which is applied to $\hat{w}(n)$ is proportional to the ~~step-size~~ input vector $v(n)$.
- Therefore, when $v(n)$ is large, the LMS algorithm experiences a problem called "gradient noise amplification".
- With the normalization of the step-size by $\|v(n)\|^2$ in the LMS algorithm, this gradient ~~noise~~ noise amplification problem is diminished.
- When $\|v(n)\|$ becomes too small,

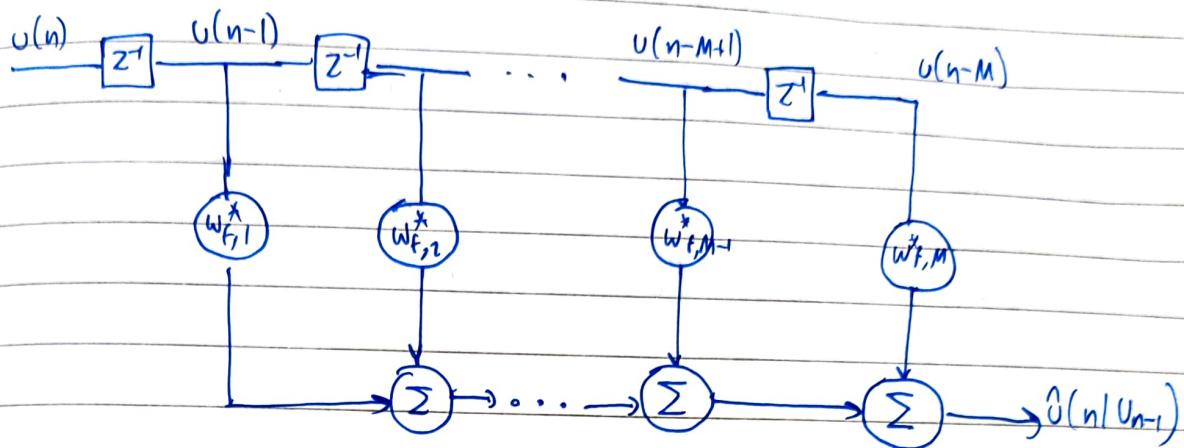
$$\hat{w}(n+1) = \hat{w}(n) + \beta \frac{v(n)}{\epsilon + \|v(n)\|^2} \cdot e^*(n)$$

where ϵ is a small tre no number

Linear Prediction

- Consider a time series $v(n), v(n-1)$ representing $(M+1)$ samples of a stationary discrete-time stochastic process.
- The operation of prediction involves using $v(n-1), v(n-2), \dots, v(n-M)$ to make an estimate of $v(n)$.
- Let U_{n-1} denote the M -dimensional space spanned by the samples $v(n-1), v(n-2), \dots, v(n-M)$. Then $\hat{v}(n, U_{n-1})$ denotes the predicted value of $v(n)$ given these samples.
- In linear prediction, we express this predicted value as a linear combination of $v(n-1), v(n-2), \dots, v(n-M)$.
- The operation corresponds to one-step prediction of the future measured with respect to time $(n-1)$.
- Accordingly, we refer to this form of prediction as one-step linear prediction, simply, forward linear prediction.
- In another form of prediction, we use the samples $v(n), v(n-1), \dots, v(n-M+1)$ to make a prediction of past sample $v(n-M)$.
- We refer to the second form of prediction as backward linear prediction.
- We will use the Wiener filter theory to optimize the design of a forward or backward predictor in the mean-square error sense, for the case of WSS discrete-time stochastic process.

Forward Linear Prediction



- The figure shows a forward prediction that consists of a linear transversal filter with M tap weights $w_{f,1}, w_{f,2}, \dots, w_{f,M}$ and tap inputs $u(n-1), u(n-2), \dots, u(n-M)$ respectively.
- Assume that the tap inputs are drawn from a WSS process with zero mean.
- We further assume that the tap-weights are optimized in the mean-square error sense in accordance with Wiener filter theory.

$$\text{The predicted value is } \hat{u}(n|U_{n-1}) = \sum_{k=1}^M w_{f,k}^* u(n-k)$$

Note that for the situation described above the desired response $d(n)$ equals $u(n)$

$$d(n) = u(n)$$

The forward prediction error equals

$$e_M(n) = d(n) - \hat{u}(n|U_{n-1})$$

signifies

- The subscript M for the forward prediction error signifies the order of the predictor defined as the number of unit-delay elements needed to store the given set of samples used to make the prediction.

- Let $P_m = E[|f_{\mu}(n)|^2]$ for all n denote the minimum mean square prediction error.
- With the tap inputs assumed to have 0 mean, the forward prediction error $f_{\mu}(n)$ will likewise have 0 mean.
- Under this condition, P_m equals the variance of the forward prediction error.
- Let w_f denote the $M \times 1$ optimum tap-weight vector of the forward prediction

$$w_f = [w_{f,1}, w_{f,2}, \dots, w_{f,M}]^T$$

- Let the $M \times 1$ tap-input vector

$$v(n-1) = [v(n-1), v(n-2), \dots, v(n-M)]^T$$

Hence the correlation matrix of the tap inputs equal

$$R = E[v(n-1) v^{*H}(n-1)]$$

$$= \begin{bmatrix} r(0) & r(1) & \cdots & r(M-1) \\ r^*(1) & r(0) & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ r^*(M-1) & \ddots & \ddots & r(0) \end{bmatrix}$$

where $r(k)$ denotes the auto-correlation of the input process for lag ' k ' where $k = 0, 1, \dots, (M-1)$

The error correlation vector b/w those tap inputs $v(n-1), v(n-2), \dots, v(n-M)$ and the desired response $v(n)$ is

$$\underline{r} = E[\underline{v}(n-1) \underline{v}^*(n)]$$

$$= \begin{bmatrix} r^*(1) \\ r^*(2) \\ \vdots \\ r^*(M) \end{bmatrix} = \begin{bmatrix} r(-1) \\ r(-2) \\ \vdots \\ r(-M) \end{bmatrix}$$

- The variance of $v(n)$ equals $r(0)$ since $v(n)$ has 0 mean.
- We may adopt the Wiener-Hopf equation to solve the forward linear prediction problem for stationary inputs.

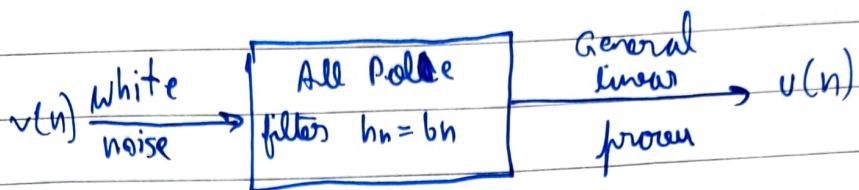
$$\underline{R} \underline{w}_f = \underline{r}$$

$$\begin{aligned} - \text{Using } J_{\min} &= \sigma_j^2 - \underline{w}_f^H \underline{R} \underline{w}_f = \sigma_j^2 - \underline{f}^H \underline{w}_f \\ &= \sigma_j^2 - \underline{p}^H \underline{R} \underline{p} \end{aligned}$$

^{we}
~~we~~ get the forward prediction error power

$$P_{e,M} = r(0) - \underline{r}^H \underline{w}_f$$

Relation between Linear Prediction and Auto-Regressive modeling



In order to uniquely define the AR model of order M, we need to specify two sets of model parameters

1. The AR coefficients a_1, a_2, \dots, a_M
2. The variance σ_v^2 of the white noise $v(n)$ used as excitation.

The time series $v(n), v(n-1), \dots, v(n-M)$ represents the realization of an auto-regressive process of order M if it satisfies

$$v(n) + a_1^* v(n-1) + \dots + a_M^* v(n-M) = v(n)$$

$$\text{or } \sum_{k=0}^M a_k^* v(n-k) = v(n)$$

$$\text{where } a_0 = 1$$

Multiply the above equation by $v^*(n-l)$ and taking expectation

$$E \left[\sum_{k=0}^M a_k^* v(n-k) v^*(n-l) \right] = E [v(n) v^*(n-l)]$$

We may simplify the above equation to

$$\sum_{k=0}^M a_k^* r(l-k) = 0, l > 0, \text{ where } a_0 = 1$$

Thus we see that the auto-correlation function of the AR process satisfies the difference eqⁿ

$$r(l) = w_1^* r(l-1) + w_2^* r(l-2) + \dots + w_M^* r(l-M)$$

where $l > 0$

$$\text{and } w_k = -\alpha_k, \quad k=1, 2, \dots, M$$

Yule-Walker Equations

Writing the above equation for $l=1, 2, \dots, M$ we get a set of M simultaneous equations with the values $r(0), r(1), \dots, r(M)$ of the auto-correlation function of the AR process as known quantities and the AR parameters $\alpha_1, \alpha_2, \dots, \alpha_M$ as the unknowns.

This set of equations can be expanded -

$$\begin{bmatrix} r(0) & r(1) & \dots & r(M-1) \\ r^*(1) & r(0) & \dots & r(M-2) \\ \vdots & \vdots & \ddots & \vdots \\ r^*(M-1) & r^*(M-2) & \dots & r(0) \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_M \end{bmatrix} = \begin{bmatrix} r^*(1) \\ r^*(2) \\ \vdots \\ r^*(M) \end{bmatrix}$$

$$\text{where } w_k = -\alpha_k$$

In a compact form,

$$R_w = \underline{r}$$

Assuming that the auto-correlation matrix is non-singular,

$$\underline{w} = R^{-1} \underline{r}$$

We may uniquely determine the matrix R and the vector \underline{r} given the auto-correlation sequence $r(0), \dots, r(M)$.

- Hence, using the above equation, we can compute the AR coefficients,

$$a_n = -w_n, \quad n=1, 2, \dots, M$$

For $l=0$, we can re-write

$$E \left[\sum_{n=0}^M a_n^* v(n) v^*(n-l) \right] = E [v(n) v^*(n-l)]$$

$$\text{or, } E [v(n) v^*(n)] = E [v(n) v^*(n)] = \sigma_v^2$$

Using the fact that

$$v(n) = u(n) + a_1^* u(n-1) + \dots + a_M^* u(n-M)$$

where σ_v^2 is the variance of the zero mean white noise $u(n)$.

Setting $l=0$ and performing complex conjugation on both sides

$$\sigma_v^2 = \sum_{n=0}^M a_n r(n) \quad \text{not with } a_0 = 1$$