

# APPENDIX F

## SUFFICIENT STATISTICS

(from class notes, ENSC 810)

## 2. BASIC ESTIMATION THEORY

2.1.1

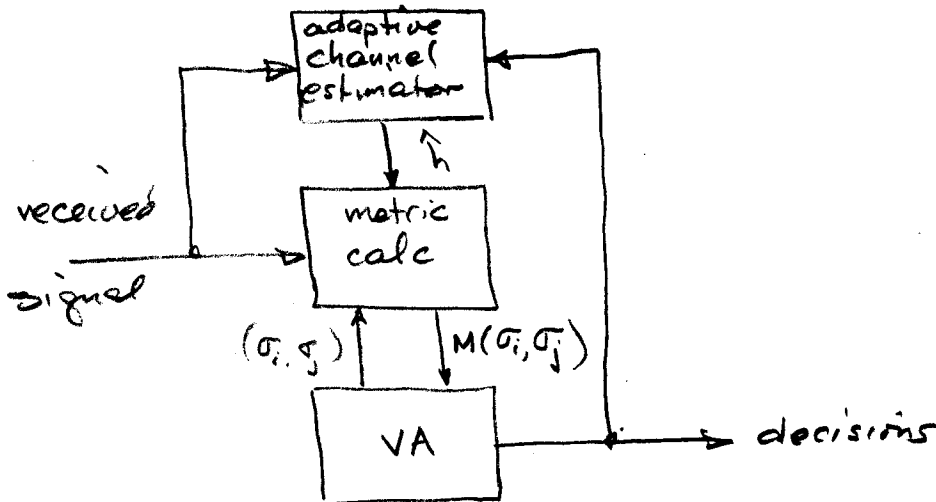
Read Haykin Appendix D

Other references:

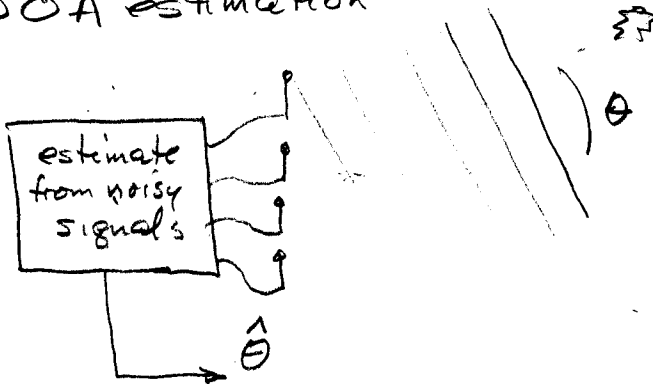
- your undergrad text on statistics
- HL van Trees, Detection Estimation and Modulation Theory, Vol. 1, McGraw Hill, late 60's
- Louis L Scharf, Statistical Signal Processing Addison Wesley 91
- S J Orfanidis, Optimum Signal Processing 2<sup>nd</sup> ed, Macmillan 94

### 2.1 A Few Examples

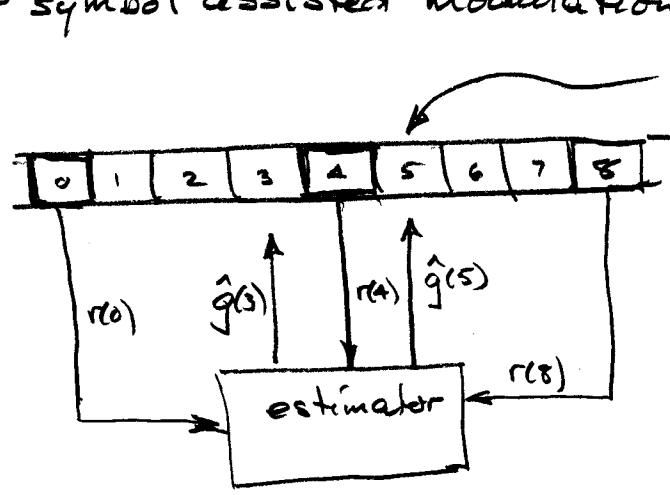
- As we start the first significant topic in the course, ask what good it will do you...
- Adaptive Viterbi Equalizer



- DOA estimation

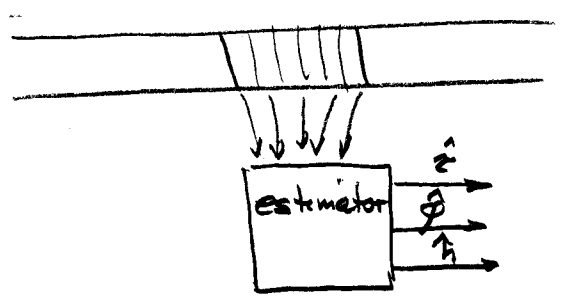


• pilot symbol assisted modulation



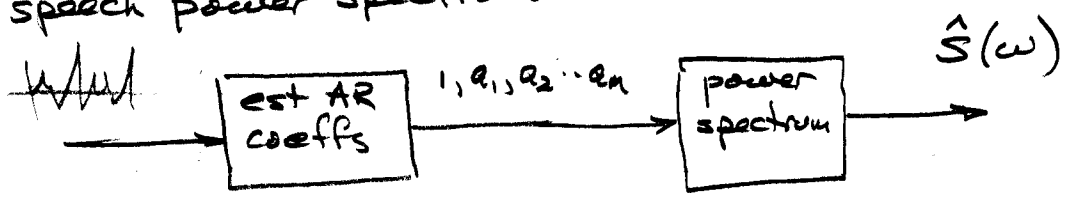
$r(k) = g(k)c(k) + n(k)$   
 fading and noise  
 some signals  $c(0), c(4), c(8)$   
 are known, so estimate  
 $\hat{g}(3), \hat{g}(5)$  etc  
 from  $r(0), r(4), r(8)$  etc

• burst channel estimation

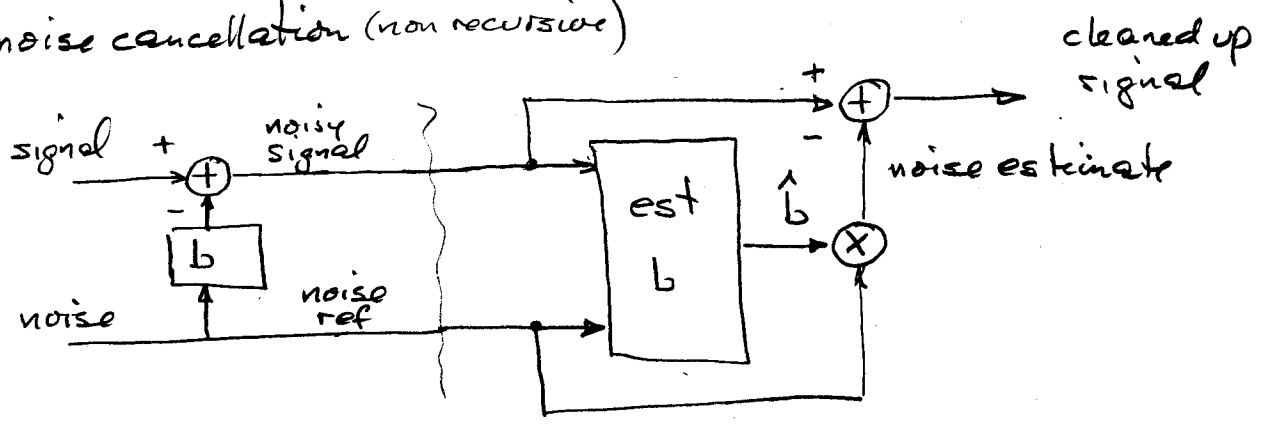


known training sequence  
 estimate delay, carrier phase,  
 impulse response for use in  
 demodulator

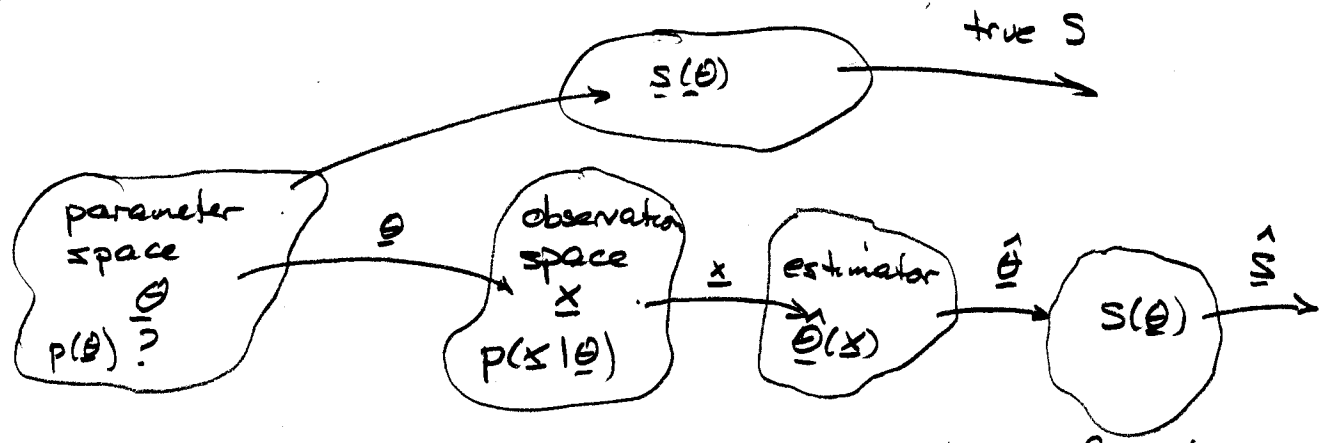
• speech power spectrum estimation



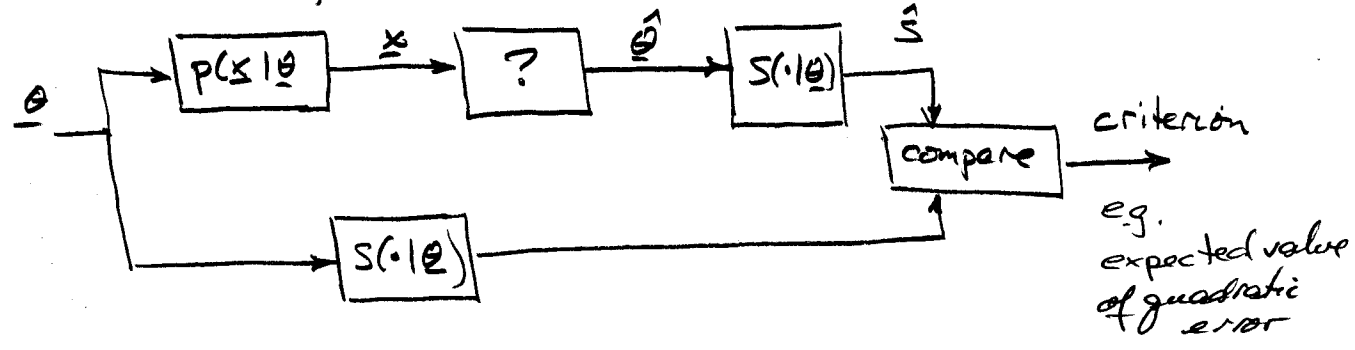
• noise cancellation (non recursive)



# 2.2 A Model of What We're Doing



We want to estimate a random or deterministic function  $S(\cdot|\theta)$  of an unknown parameter  $\theta$ . All we have is an observation  $x$  with pdf dependent on  $\theta$ . Use  $x$  to obtain  $\hat{\theta}$ , hence  $\hat{S}$ , using a criterion of goodness



Examples (Schof):

- $S(\cdot|\theta) = \theta$  Here we estimate  $\theta$  itself
- $S(\cdot|\theta) = S(e^{j\omega}|\theta)$  Estimate power spectral density

Haykin 3.5

- Two important cases
  - $\theta$  is random, pdf  $p(\theta)$  known
  - $\theta$  is non random, but unknown

• Note this is parameter estimation, not interval estimation;  
 "The poll is considered accurate to within 3% 19 times out of 20."

## 2.3 Properties of Estimators

There are certain key properties we look for in estimators.  
Here are some common ones

### • Unbiased estimators

An estimator  $\hat{\theta}(x)$  is unbiased if  $E[\hat{\theta}] = E[\theta]$ .

example  $\theta$  is mean  $\mu$ ,  $x = (x_1, \dots, x_N)$   $N$  measurements

Estimate  $\mu$ . Easy to show that

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i \quad (\text{sample mean})$$

is unbiased — just take expectation

$$E[\hat{\mu}] = \frac{1}{N} N \mu = \mu$$

example sample variance.

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2 \quad \text{if we know true mean } \mu$$

$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \hat{\mu})^2 \quad \text{to remove bias if we use sample mean}$$

Since

$$\begin{aligned} E\left[\sum_{i=1}^N (x_i - \hat{\mu})^2\right] &= E\left[\sum_i (x_i - \frac{1}{N} \sum_j x_j)^2\right] \\ &= E\left[\sum_i x_i^2 + \frac{1}{N} \sum_i \sum_j x_i x_j - \frac{2}{N} \sum_i \sum_j x_i x_j\right] \\ &= E\left[\sum_i x_i^2 + \frac{1}{N} \sum_i \sum_j x_i x_j\right] \\ &= N \overline{x^2} - \frac{1}{N} (N \overline{x^2} + N(N-1)\mu^2) = (N-1) \overline{x^2} - (N-1)\mu^2 \\ &= (N-1) \sigma^2 \quad \text{so divide by } N-1 \end{aligned}$$

Sometimes we accept a little bias to gain less scatter.

## • Consistent estimators

- unbiasedness is nice, but all it says is that the mean of the estimator is right - doesn't say the estimate is close.
- estimator is consistent if prob  $\hat{\theta}$  takes a value farther away from  $\theta$  than some arbitrary  $c$  approaches 0 when  $N \rightarrow \infty$  ( $N$  is # observations).
- e.g. the sample mean is consistent, since sample variance  $\rightarrow 0$ , and use Chebyshev upper bound
- we need two properties
  - a)  $\hat{\theta}$  is unbiased
  - b)  $\text{var}(\hat{\theta}) \rightarrow 0$  when  $N \rightarrow \infty$  (or cov matrix  $E[(\theta - \hat{\theta})(\theta - \hat{\theta})^T]$ )

## • Efficient estimators

- if there is more than one unbiased consistent estimator available, the one with the smallest variance is said to be the most efficient of them
- a lower bound on the variance of any unbiased consistent estimator is the Cramér Rao bound (more about it later). If it meets the C-R bound it is an efficient estimate

- Why the restriction to unbiased estimators?  
It's because there is often a tradeoff between bias and variance; e.g.

$$\hat{\theta}(x) = \underline{\theta} \quad \text{zero variance!}$$

$$E[\hat{\theta}] - \underline{\theta} = -\underline{\theta} \quad \text{bias is } -\underline{\theta}$$

but it's not much use.

• Sufficient estimator

- This is linked to sufficient statistics - and to data reduction.
- example: When you estimate  $\mu$  from a set of measurements, do you really need all of  $x_1, \dots, x_N$  in detail, or is the sample mean  $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$  sufficient?
- An estimator  $\hat{\theta}$  is sufficient if it uses all the info relevant to estimation of  $\theta$  in the observations.

An informal example is an easier way to understand than mathematics (that comes later).

- Measurements  $x$  are useful for estimating  $\theta$  only to the extent that  $\theta$  influences the pdf of  $x$ . If  $P_{x|\theta}(x|\theta)$  is not a function of  $\theta$ , then  $x$  gives no information about  $\theta$ .

- Consider  $N$  Bernoulli trials,  $x_n = \begin{cases} 1 & \text{prob } \theta \\ 0 & \text{prob } 1-\theta \end{cases}$

In  $x$ , get  $k$  ones.

- Question: Does the sample mean  $\hat{\theta} = k/N$  capture all the info about  $\theta$  in  $\underline{x}$ ?

Equivalent question: Given  $k$ , does the conditional pdf of  $\underline{x}$  still depend on  $\theta$ ?

- Example. Suppose  $N=4$  and you know there are  $k$  ones. Then the possible patterns are

1100    1010    1001  
0110    0101    0011

Nothing about  $\theta$  would make one of them more probable than another, so they are equiprobable. Consequently, knowing which  $\underline{x}$  pattern it is tells nothing about  $\theta$ . Therefore,  $k$  alone is sufficient, not the detailed pattern.

Formally,

$\hat{\theta}(\underline{x})$  is a sufficient estimator of  $\theta$  if  $p(\underline{x}|\theta)$  does not depend on the true value  $\theta$ .

In effect, the value of  $\hat{\theta}$  summarizes all the information in  $\underline{x}$  relevant to estimation of  $\theta$



- Example:  $N$  Bernoulli trials, observe  $x_n = \begin{cases} 1, & \text{prob } \theta \\ 0, & \text{prob } 1-\theta \end{cases}$

$$P_{x|\theta}(\underline{x}|\theta) = \theta^k (1-\theta)^{N-k} \quad k = \# \text{ of ones}$$

$$P_k(k) = \binom{N}{k} \theta^k (1-\theta)^{N-k} \quad \text{prob of } k \text{ ones}$$

Then  $k$  is a sufficient statistic for  $\theta$ , and

$\hat{\theta} = k/N$  is a sufficient estimator.

why?

$$P_{x|\theta}(\underline{x}|\hat{\theta}=k) = \frac{P_{x,\theta}(\underline{x}, k)}{P_k(k)} = \frac{\theta^k (1-\theta)^{N-k}}{\binom{N}{k} \theta^k (1-\theta)^{N-k}} = \frac{1}{\binom{N}{k}}$$

a uniform pdf indep of  $\theta$

- Example: binomial trials with same mean  $\theta$

$N_1$  has  $x_1$  ones;  $N_2$  has  $x_2$  ones.

To estimate mean, we don't need  $x_1, x_2$  individually — the sum will do:

$$\hat{\theta} = \frac{x_1 + x_2}{N_1 + N_2}$$

$$\begin{aligned} P_{\underline{x}|\theta}(\underline{x}|\hat{\theta} = \frac{x_1+x_2}{N_1+N_2}) &= \frac{P_{\underline{x}, x_1+x_2}(\underline{x}, x_1+x_2)}{P_{\hat{\theta}}((x_1+x_2)/(N_1+N_2))} \\ &= \frac{P_{\underline{x}}(x_1, x_2)}{P_{x_1+x_2}(x_1+x_2)} = \frac{\binom{N_1}{x_1} \theta^{x_1} (1-\theta)^{N_1-x_1} \binom{N_2}{x_2} \theta^{x_2} (1-\theta)^{N_2-x_2}}{\binom{N_1+N_2}{x_1+x_2} \theta^{x_1+x_2} (1-\theta)^{N_1+N_2-x_1-x_2}} \\ &= \frac{\binom{N_1}{x_1} \binom{N_2}{x_2}}{\binom{N_1+N_2}{x_1+x_2}} \quad \theta \text{ has no further influence} \end{aligned}$$