# Lecture 3
# MMSE Estimation: & Information Measures

## April 9, 2024

### JOHN M. CIOFFI

Hitachi Professor Emeritus of Engineering

Instructor EE379B – Spring 2024

# Announcements & Agenda

- Announcements
  - Problem Set 1 due Wednesday at 17:00
  - Most relevant reading – 1.5, 2.3, D.1, D.2, 4.1
  - Problem Set 2 due April 17 at 17:00

- Problem Set 2 = PS2 due 4/19 at 17:00
  1. 4.29  biases and error probability
  2. 4.36  MMSE spatial equalizer
  3. 2.10  Entropy and Dimensionality
  4. 2.14  Bandwidth vs Power
  5. 2.20  MMSE and Entropy

- Agenda
  - General MMSE & Gaussian
    - Autocorrelation/Cross-Correlation
  - Linear MMSE & The Orthogonality Principle
    - Biases and SNRs
  - Examples
  - Information Measures – generalizing Gaussian to all distributions
  - Mutual Information and MMSE

**AWGNs, or whited AGNs:**

**MMSE is canonical & asymptotically MAP, and also close @ finite N. Max reliable b is capacity.**
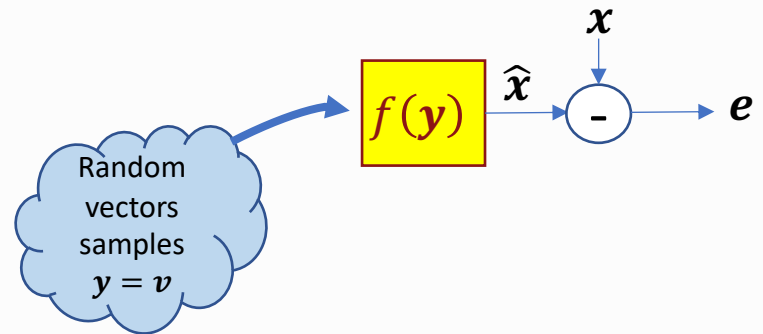
# General MMSE and Gaussian

*Section D.1*

When all random processes are (stationary) Gaussian, Linear MMSE is the overall best MMSE.

# The MMSE Estimation Problem

- Given random $x$ and $y$, design estimates $x$ as $\widehat{x} = f(y)$.
  - Designer knows both $p_{x,y}$ and specific observed $y = v$.
  - $x$ and $y$ 's distributions are continuous. (random code on $x$.)

- The error is:
  - $e = x - f(y)$.

Random vectors samples $y = v$

$f(y)$ $\widehat{x}$ $x$ $-$ $e$

- Its mean-square is
  - $\mathbb{E}_{x,y}[\|e\|^2] = \mathbb{E}[\|x - f(y)\|^2]$.

- Its minimum over $f$ is the MMSE.

$$MMSE = \min_{f} \mathbb{E}\left[\|x - f(y)\|^2\right].$$

- Solution is
  - the conditional mean of $x$ given $y$; $\widehat{x} = \mathbb{E}[x/y]$,
  - From $p_{x/y}$ the *à posteriori* distribution (used for MAP detector).
  - Proof: see Appendix D.1
  - Its linear: $\widehat{(x_1 + x_2)} = \widehat{x}_1 + \widehat{x}_2$.

$$\widehat{x} = \mathbb{E}[x/y] = W \cdot y$$

# Auto- & Cross- correlation

- **Autocorrelation** generalizes mean-square. When stationary,
  - it samples $R_{xx}(\tau) = \mathbb{E}[\boldsymbol{x}(t) \cdot \boldsymbol{x}^*(t-\tau)]$ where time is the dimension and $\tau = kT'$ is correlation interval.
  - Vector process samples may correspond to:
    - frequency-time – time samples $x(kT)$ are the vector elements.
    - space time sets $\tau = 0$ – spatial samples (think antennas),
      - often at each $n \cdot \left(\frac{1}{T}\right)$.

$$R_{xx} = \mathbb{E}[\boldsymbol{x} \cdot \boldsymbol{x}^*] \qquad R_{yy} = \mathbb{E}[\boldsymbol{y} \cdot \boldsymbol{y}^*]$$
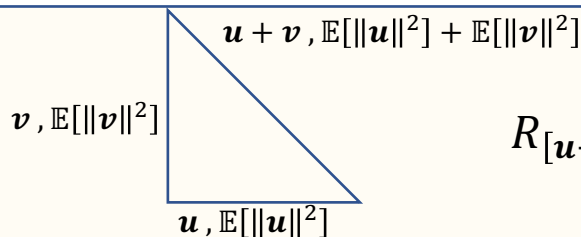
- **Energy** over symbol if FT, over space dimensions:
  - Both cases $\tau = 0$ (non xtalk/energy) terms are on diagonal.

$$\mathcal{E}_x = trace\{R_{xx}\} = \mathbb{E}[\boldsymbol{x}^* \cdot \boldsymbol{x}] = \mathbb{E}[\|\boldsymbol{x}\|^2]$$

- **Cross correlation** "generalizes" inner product:
  - It samples $R_{xx}(\tau) = \mathbb{E}[x(t) \cdot y^*(t-\tau)]$.
  - Vectors can be different lengths $L_x$ and $L_y$.
  - "uncorrelated" (=0) → **orthogonal.**
  - Nondiagonal dimensions.

$$R_{xy} = \mathbb{E}[\boldsymbol{x} \cdot \boldsymbol{y}^*] \qquad R_{yx} = \mathbb{E}[\boldsymbol{y} \cdot \boldsymbol{x}^*]$$

- **Pythagorus** IF uncorrelated $R_{uv} = 0$.
  - Generalizes "variances of uncorrelated random variables add."

$\boldsymbol{u} + \boldsymbol{v}, \mathbb{E}[\|\boldsymbol{u}\|^2] + \mathbb{E}[\|\boldsymbol{v}\|^2]$

$\boldsymbol{v}, \mathbb{E}[\|\boldsymbol{v}\|^2]$

$\boldsymbol{u}, \mathbb{E}[\|\boldsymbol{u}\|^2]$

$$R_{[u+v][u+v]} = R_{uu} + R_{vv}$$

# The Joint Gaussian Distribution

- is completely specified by its autocorrelation (and cross correlation):

$$R \triangleq R_{\begin{bmatrix} x \\ y \end{bmatrix}[x^* \quad y^*]} = \begin{bmatrix} R_{xx} & R_{xy} \\ R_{yx} & R_{yy} \end{bmatrix}.$$

- Its marginal distributions for $x$ and $y$
  - are also Gaussian.

$$\text{real: } p(x, y) = (2\pi)^{-\frac{N_x + N_y}{2}} \cdot |R|^{-1/2} \cdot e^{-\frac{1}{2}\left\{[x^* y^*] \cdot R^{-1} \cdot \begin{bmatrix} x \\ y \end{bmatrix}\right\}}$$

- Its conditional distributions are Gaussian.
  - In particular, with non-zero mean $\mathbb{E}[x/y]$.

$$\text{complex: } p(x, y) = (\pi)^{-[N_x + N_y]} \cdot |R|^{-1} \cdot e^{-\left\{[x^* y^*] \cdot R^{-1} \cdot \begin{bmatrix} x \\ y \end{bmatrix}\right\}}$$

- Singularity?
  - $|R_{yy}| > 0$ with nonsingular noise.
  - $|R_{xx}|$ ? $|R|$? – use pseudoinverse and determinant as product of **nonzero** eigenvalues.

$$\boxed{\mathbb{E}[x/y] = \underbrace{R_{xy} \cdot R_{yy}^{-1}} \cdot y}$$

It's linear (for Gaussian).

$$W = R_{xy} \cdot R_{yy}^{+} \text{ if singular.}$$

**MMSE and AWGN's best transmission are fundamentally connected.**

# Linear MMSE & The Orthogonality Principle
*Section D.2*

**Decomposition into pass spaces and null spaces becomes critical in canonical design (think VC).**

# Linear MMSE: <u>any</u> joint distribution of $x$ and $y$

- Given random $\boldsymbol{x}$ and $\boldsymbol{y}$, receiver estimates $\boldsymbol{x}$, $\widehat{\boldsymbol{x}} = W \cdot \boldsymbol{y}$.
  - It knows both $p_{\boldsymbol{x},\boldsymbol{y}}$ and specific observed $\boldsymbol{y} = \boldsymbol{v}$.

- The error is:

$$\boxed{\boldsymbol{e} = \boldsymbol{x} - \sum_{n=1}^{N} \boldsymbol{w}_n \cdot \boldsymbol{y}_n = \boldsymbol{x} - W \cdot \boldsymbol{y}.}$$

- Its mean-square is
  - $\mathbb{E}_{\boldsymbol{x},\boldsymbol{y}}[\|\boldsymbol{e}\|^2] = \mathbb{E}[\|\boldsymbol{x} - W \cdot \boldsymbol{y}\|^2]$ .

- Its minimum occurs when $\boxed{\mathbb{E}[\boldsymbol{e} \cdot \boldsymbol{y}_n^*] = 0}$ for all $n$:
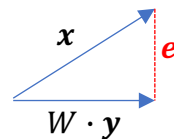  - Proof is in Appendix D.2.
  - That is, the error and the estimator's input are uncorrelated.
- Minimum $\widehat{\boldsymbol{x}} = \underbrace{R_{\boldsymbol{xy}} \cdot R_{\boldsymbol{yy}}^{-1}}_{W} \cdot \boldsymbol{y}$ , linear in $\boldsymbol{y}$ , **so true MMSE if Gaussian**.

  - The true MMSE estimator may not be linear if non-Gaussian.
- Also again: $(\widehat{\boldsymbol{x}_1 + \boldsymbol{x}_2}) = \widehat{\boldsymbol{x}}_1 + \widehat{\boldsymbol{x}}_2$ or $\widehat{A \cdot \boldsymbol{x}} = A \cdot \widehat{\boldsymbol{x}}$ .

**Orthogonality Principle**



MMSE Matrix: $R_{\boldsymbol{ee}} = R_{\boldsymbol{xx}} - R_{\boldsymbol{xy}} \cdot R_{\boldsymbol{yy}}^{-1} \cdot R_{\boldsymbol{yx}} = R_{\boldsymbol{x/y}}^{\perp}$
$= R_{\boldsymbol{xx}} - W \cdot R_{\boldsymbol{yx}}$

# Vector and Matrix Norms

- The trace of an autocorrelation matrix is its norm (and also equal to mean-squared length of random vector).

- MMSE $= \mathbb{E}[\|\boldsymbol{e}\|^2] = trace\{R_{\boldsymbol{ee}}\}$ .

- The **trace** of a square autocorrelation matrix is also equal to the **sum** of its eigenvalues $\mathcal{E}_{e',n}$ :
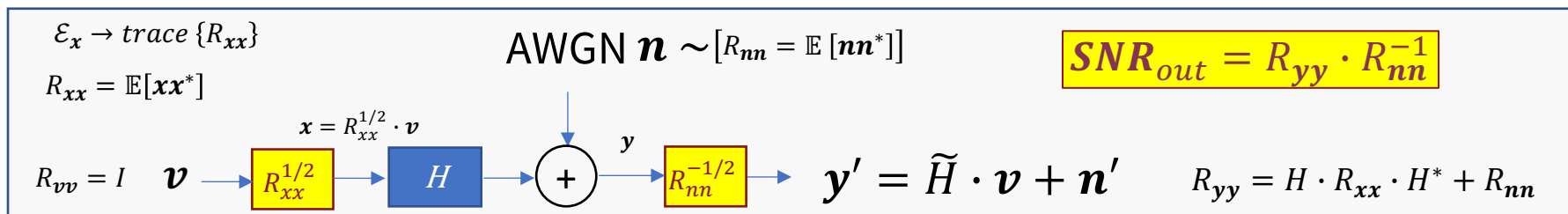
$$\boldsymbol{e}' = Q \cdot \boldsymbol{e} \qquad diagonal: R_{\boldsymbol{e'e'}} = Q \cdot R_{\boldsymbol{ee}} \cdot Q^* .$$

$$\|\boldsymbol{e}\|^2 = \|\boldsymbol{e}'\|^2 \text{ because } QQ^* = Q^*Q = I .$$

- The **determinant** of an autocorrelation matrix is the **product** of its eigenvalues $\mathcal{E}_{e',n}$:
- MMSE $= \mathbb{E}[\|\boldsymbol{e}\|^2]$ and $\ln|R_{\boldsymbol{ee}}| = \sum_n \ln \mathcal{E}_{e',n}$ .
- The minimization of each component of $\boldsymbol{e}$ is *variables separable* (has its own row of $W$),
  - so then the sum is minimized,
  - but this means each of the $\boldsymbol{e}'$ also ($W \rightarrow Q \cdot W$) minimized,
  - so then $|R_{\boldsymbol{ee}}| = |R_{\boldsymbol{e'e'}}|$ is also minimized
  - → Minimizing sum (trace) here **is same as minimizing product (determinant).**

# Matrix SNR?

$\mathcal{E}_x \rightarrow trace\{R_{xx}\}$

$R_{xx} = \mathbb{E}[\boldsymbol{xx}^*]$

AWGN $\boldsymbol{n} \sim [R_{nn} = \mathbb{E}[\boldsymbol{nn}^*]]$

$\boxed{SNR_{out} = R_{yy} \cdot R_{nn}^{-1}}$

$\boldsymbol{x} = R_{xx}^{1/2} \cdot \boldsymbol{v}$

$R_{vv} = I$ $\quad \boldsymbol{v} \rightarrow \boxed{R_{xx}^{1/2}} \rightarrow \boxed{H} \rightarrow (+) \xrightarrow{y} \boxed{R_{nn}^{-1/2}} \rightarrow \boldsymbol{y}' = \widetilde{H} \cdot \boldsymbol{v} + \boldsymbol{n}'$

$R_{yy} = H \cdot R_{xx} \cdot H^* + R_{nn}$

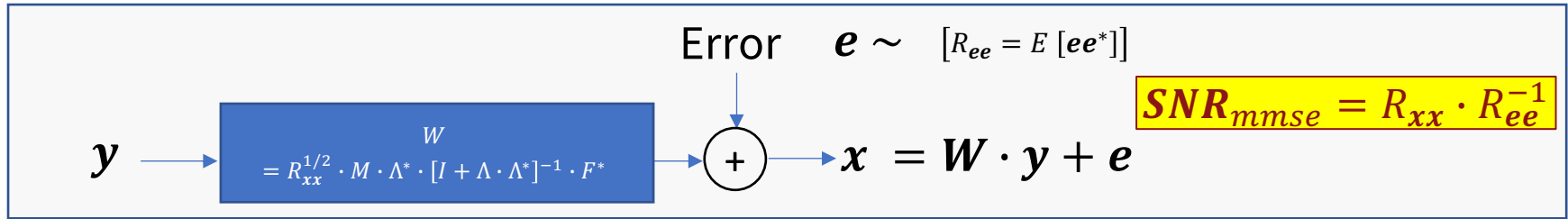$$\widetilde{H} \triangleq R_{nn}^{-1/2} \cdot H \cdot R_{xx}^{1/2} = F \cdot \Lambda \cdot M^*$$

- It's like the parallel channels (take determinants) $SNR_{out} = |R_{yy}| \cdot |R_{nn}^{-1}| = \frac{|R_{yy}|}{|R_{nn}|}$.
  - Vector code from $\boldsymbol{v}$ to $\boldsymbol{y}'$.
  - Bit rate is then $b = \log_2(SNR_{out})$.

$$SNR_{out} = \frac{|R_{yy}|}{|R_{nn}|} = \frac{|H \cdot R_{xx} \cdot H^* + R_{nn}|}{|R_{nn}|} = |\widetilde{H} \cdot \widetilde{H}^* + I| = |\Lambda^2 + I| = \prod_{n=1}^{N}(SNR_n + 1)$$

- This set depends on $R_{xx}$ choice, likely with fixed $trace\{R_{xx}\}$.
  - Water-fill $R_{xx} = M \cdot diag\{\mathcal{E}_{water-fill}\} \cdot M^*$ maximizes the matrix SNR or effectively its determinant.

Error $e \sim \quad [R_{ee} = E[ee^*]]$

$$SNR_{mmse} = R_{xx} \cdot R_{ee}^{-1}$$

$y \longrightarrow$ $W = R_{xx}^{1/2} \cdot M \cdot \Lambda^* \cdot [I + \Lambda \cdot \Lambda^*]^{-1} \cdot F^*$ $\longrightarrow (+) \longrightarrow x = W \cdot y + e$

- How about "backward channel's," $x = W \cdot y + e$ , (MMSE) $SNR$?

$$SNR_{mmse} = \frac{|R_{xx}|}{|R_{ee}|} = |W \cdot R_{yy} \cdot W^* + R_{ee}|/|R_{ee}| = |\Lambda^2 + I| = \prod_{n=1}^{N}(1 + SNR_n)$$
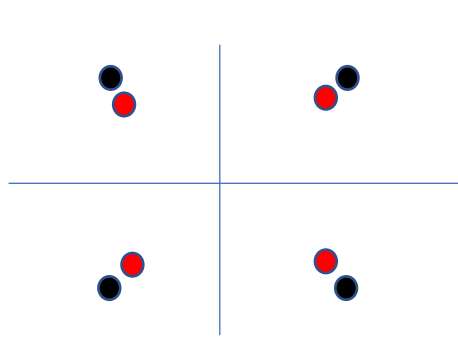
$$= SNR_{out}$$

- Bit rate is again $b = \log_2(SNR_{mmse})$, $\mathbb{I}(x; y)$ is symmetric.

- $M^* \cdot W$ will estimate $v$ (linearity of MMSE estimates).

- Optimizing determinants is same as optimizing MSE/traces.

**Forward and backward have same SNR and "bit rate" (continuous $x$ distribution, or Random AEP good-code sense.)**
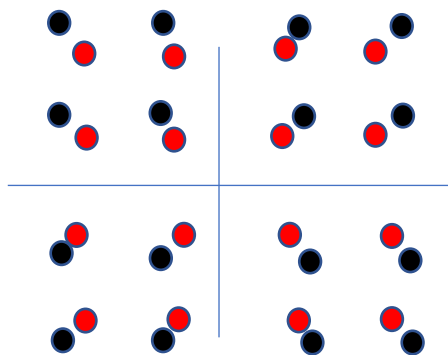
# MMSE is always a Biased Estimate

- Biased-Estimate Definition: $\mathbb{E}\left[\hat{x}/x\right] \neq x.$

- MMSE estimates always have bias (if noise is nonzero), See Appendix D.2:
  - $\mathbb{E}\left[\hat{x}/x\right] = (I - R_{ee} \cdot R_{xx}^{-1}) \cdot x = (I - SNR^{-1}) \cdot x$



decision regions same

decision regions change

MMSE trades a little signal reduction for simultaneous noise reduction when minimizing the error, now on every dimension.
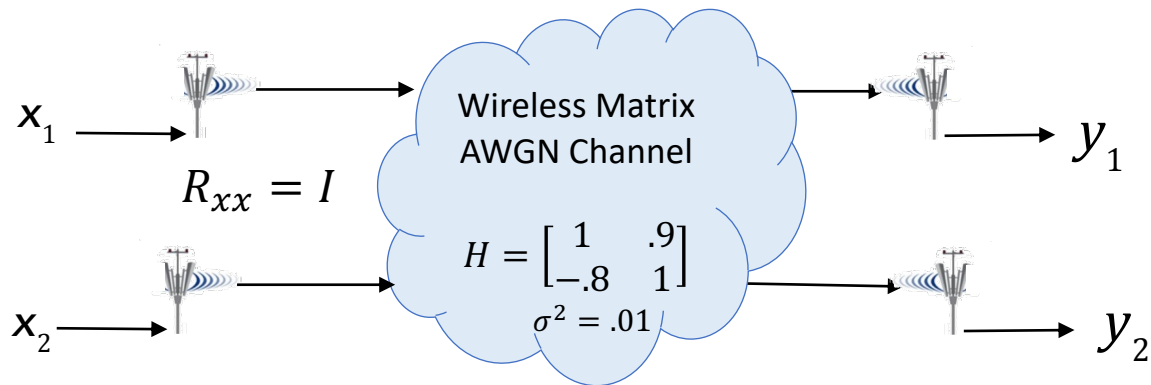
*See PS2.1 (Prob 4.29)*

- For scalar case above, removal is scale up (by $\frac{SNR_{mmse}}{SNR_{mmse}-1}$ ).

- MIMO case, same per dimension, scale up (by $\frac{SNR_{mmse,n}}{SNR_{mmse,n}-1}$ ) **IF** MMSE $R_{ee}$ is diagonal (vector coding).
  - IF not diagonal? (we'll learn what to do in later lectures.)

# Linear Matrix MMSE Examples

*See PS2.2 (Prob 4.36)*

# 2 x 2 Antenna System



$x_1$

$R_{xx} = I$

Wireless Matrix
AWGN Channel

$H = \begin{bmatrix} 1 & .9 \\ -.8 & 1 \end{bmatrix}$

$\sigma^2 = .01$

$y_1$

$x_2$

$y_2$

```
>> H=[1 .9
-.8 1];
>> Rxx=eye(2);
>> Rnn=.01*eye(2)
>> Ryy=H*Rxx*H'+Rnn;
>> Ryx=H;
>> W=(Ryx')*inv(Ryy) =
   0.5780  -0.5199
   0.4627   0.5780
>> W*H =
   0.9939   0.0003
   0.0003   0.9945
>> Ree=Rxx-W*Ryx =
   0.0061  -0.0003
  -0.0003   0.0055
>> snr=det(Rxx)/det(Ree);
>> log2(snr) =  14.8693
```

- Strong Crosstalk case from Chapter 1

  Basically, the same L1:31 result, even without the "$M$" discrete modulator
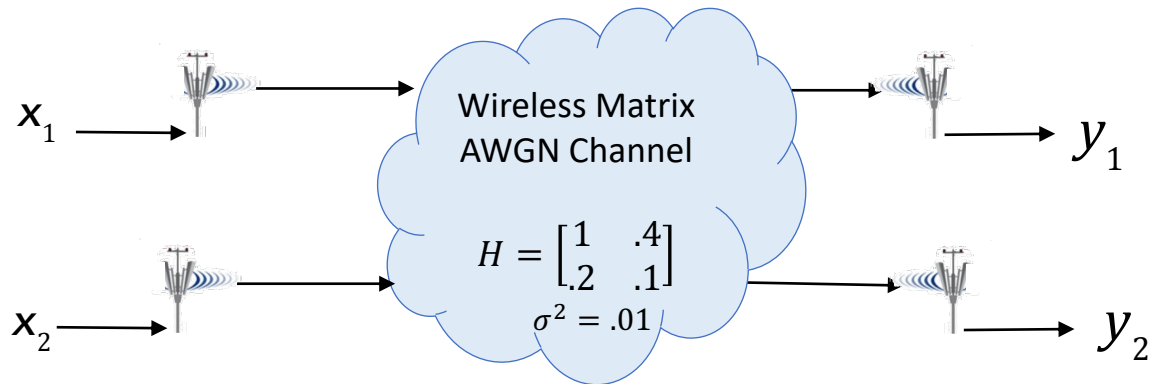  but why with no transmit $M$ matrix?

```
>> Mstar =

   0.4197   0.9076
   0.9076  -0.4197
```

$R_{xx} = I$ is close to water-fill (equal energy this channel);
$R_{xx} = M \cdot I \cdot M^*$ ; so "lucky" that its already close to best.

**ML detector is only per-dimension independent if $R_{xx}$ and $R_{ee}$ are diagonal.**

# 2 x 2 Antenna System



$x_1$

$x_2$

Wireless Matrix AWGN Channel

$H = \begin{bmatrix} 1 & .4 \\ .2 & .1 \end{bmatrix}$

$\sigma^2 = .01$

$y_1$

$y_2$

- This channel water-filled nonzero energy only on 1 dimension in L1.

```
>>H = 1.0000   0.4000
       0.2000   0.1000
>> Rxx=eye(2);
>> Ryy=H*Rxx*H'+Rnn;
>> Ryx=H;
>> W=(Ryx')*inv(Ryy) =
   0.9524  -0.4762
   0.0000   1.6667
>> W*H = 0.8571   0.3333
         0.3333   0.1667
>> Ree=Rxx-W*Ryx =
   0.1429  -0.3333
  -0.3333   0.8333
>> snr=det(Rxx)/det(Ree) = 126.0000
>> b=log2(snr) =  6.9773  (only for VC)
```

not diagonal; ML detect is NOT parallel

```
>> SNR=inv(diag(diag(Ree))) = 7.0000      0
                                  0   1.2000

>> log2(diag(SNR)) =
   2.8074
   0.2630
>> sum(log2(diag(SNR))) =  3.0704
```

Previously in L1, $b = \log_2(1 + 2 \cdot g_2) = 6.93$ bits/subsymbol
But this time, two dimensions are used, and the ML detectors are interdependent

**Thus, data rate loss can occur
with independent detectors and MMSE.**

*(This loss can be recovered with Chapter 5's MMSE GDFE,
in addition to using vector coding, so more than 1 solution)*

*Section 4.1.2*

*See PS2.2 (Prob 4.36)*

**Stanford University**

```
>> H=toeplitz([1 zeros(1,7)]',[1 .9 zeros(1,7)]);
>> Rxx=eye(9);
>> Rnn=.181*eye(8);
>> Ryy=H*Rxx*H'+Rnn;
>> Ryx=H;
>> W=(Ryx')*inv(Ryy);
>> P=W*H;
>> size(P) % =  9   9

>> Ree=Rxx-W*Ryx;
>> snr=det(Rxx)/det(Ree) = 2.4089e+07
>> SNR=inv(diag(diag(Ree)));
>> bn = 0.5*log2(diag(SNR))' =
 0.8769  1.0096  1.0691  1.0907  1.0902  1.0673  1.0054  0.8681  0.6085
>> sum(bn/9) =   0.9651
>> 10*log10(2^(2*ans)-1) = 4.4885 dB
```

Repeat for 8→32:

>> sum(bn/33)  =  1.0753

>> 10*log10(2^(2*ans)-1) =   5.3654 dB

Best infinite length is 5.7 dB.
(with dimension-by-dimension linear)
- See Chapter 3, 379A MMSE-LE example

Best with full ML is 8.8 dB, but requires
input WF energy distribution .

*Might want to make your own homework*

**Stanford University**

# Entropy and Estimation:
## generalizing energy to all distributions

*See PS2.3 (Prob 2.10)*

**Rest of L3 repeats EE379A early coding,
but here emphasizes the MMSE-canonical connection.**

# Information Measures and Energy/Mean-Square

| Gaussian Distribution | Any Distribution |
|---|---|
| Mean-square energy $\mathcal{E}_x = \mathbb{E}[|\boldsymbol{x}|^2]$ | Entropy $\mathcal{H}_{\boldsymbol{x}}$    *Section 2.3.1* |
| Mean-square error $\sigma_e^2 = \mathbb{E}[|\boldsymbol{e}|^2]$ | Conditional Entropy $\mathcal{H}_{\boldsymbol{x}/\boldsymbol{y}}$ *Section 2.3.2* |
| Signal-to-Noise $\mathcal{E}_x / \sigma_e^2$ | Mutual Information $\mathbb{I}(\boldsymbol{x}; \boldsymbol{y}) = \mathcal{H}_{\boldsymbol{x}}\text{-}\mathcal{H}_{\boldsymbol{x}/\boldsymbol{y}}$ *Section 2.3.2* |

The information-carried by random variable/process generalizes the energy concepts from MMSE/Gaussian analysis to a general distribution.

These information measures correspond to bits/symbol quantities, and for the Gaussian case are basically the $\log_2$ of the corresponding energy measure.

- **Entropy:** $$\mathcal{H}_{\widetilde{x}} = \mathbb{E}\left[\log_2\left(\frac{1}{p_{\widetilde{x}}}\right)\right] = \sum_{i=0}^{|C|-1} p_{\widetilde{x}}(i) \cdot \log_2\left(\frac{1}{p_{\widetilde{x}}(i)}\right)$$   Discrete $p_{\widetilde{x}}(i)$

- Measures a distribution's many values, its ***information***, by probability (think subsymbols).

- Generalizes bits/subsymbol, especially when the constellation size $|C| \geq M^{1/\overline{N}} = 2^{\tilde{b}}$.

$$\text{example: } p_{\widetilde{x}}(i) = \frac{1}{M} \text{ (uniform) } \rightarrow |C| = 2^{\tilde{b}}$$

$$\text{Uniform} \rightarrow \mathcal{H}_{\widetilde{x}} = \log_2\left(M^{1/\overline{N}}\right) = \tilde{b} \quad (|C| = 2^{\tilde{b}+\tilde{\rho}}) \; \tilde{\rho} = 0 \text{ ; uncoded)}$$

- Uniform distribution has maximum entropy

$$\mathcal{H}_{\widetilde{x}} \leq \log_2|C| \qquad \text{Binary example: } p_{\widetilde{x}}(0) = \frac{1}{128} \text{ and } p_{\widetilde{x}}(1) = \frac{127}{128}$$

$$\mathcal{H}_{\widetilde{x}} = \frac{\log_2(128)}{128} + \frac{127}{128} \cdot \log_2\left(\frac{128}{127}\right) = .06 < 1 \qquad \textit{See PS2.3 (Prob 2.10)}$$

# Information left after given another random vector

- Conditional entropy

$$\mathcal{H}_{\widetilde{x}/\widetilde{y}} = \mathbb{E}\left[\log_2\left(\frac{1}{p_{\widetilde{x}/\widetilde{y}}}\right)\right] = \sum_{j=0}^{|Y|-1}\sum_{i=0}^{|C|-1} p_{\widetilde{x}\widetilde{y}}(i,j)\cdot\log_2\left(\frac{1}{p_{\widetilde{x}/\widetilde{y}}(i,j)}\right)$$

$$\mathcal{H}_{\widetilde{x}/\widetilde{y}} = \mathcal{H}_{\widetilde{x}\widetilde{y}} - \mathcal{H}_{\widetilde{y}}$$

- Measures $\widetilde{x}$'s residual randomness/info when $\widetilde{y}$ is known/given

| $\widetilde{x}$ ; $\widetilde{y}$ | 0 | 1 | $p_{\widetilde{x}}$ |
|---|---|---|---|
| 0 | 3/8 | 1/8 | 1/2 |
| 1 | 1/8 | 3/8 | 1/2 |
| $p_{\widetilde{y}}$ | 1/2 | 1/2 | |

$$\mathcal{H}_{\widetilde{x}\widetilde{y}} = \frac{6}{8}\cdot\log_2\frac{8}{3} + \frac{2}{8}\cdot\log_2 8 = 1.811$$

$$\mathcal{H}_{\widetilde{x}} = 1 = \mathcal{H}_{\widetilde{y}}$$

$$\mathcal{H}_{\widetilde{x}/\widetilde{y}} = 1.811 - 1 = .811 \; bits/subsymbol$$

- If $x$ and $y$ are independent, then $\mathcal{H}_{\widetilde{x}/\widetilde{y}} = \mathcal{H}_{\widetilde{x}}$

# Continuous Distribution – DIFFERENTIAL Entropy

- Differential Entropy

$$\mathscr{H}_{\widetilde{x}} = \mathbb{E}\left[\log_2\left(\frac{1}{p_{\widetilde{x}}}\right)\right] = -\int_{-\infty}^{\infty} p_{\widetilde{x}}(u) \cdot \log_2\left(\frac{1}{p_{\widetilde{x}}(u)}\right) \cdot du$$

- Differential Entropy $\mathscr{H}_{\widetilde{x}}$ is not same as an integral-to-sum via a discrete approximation of $p_{\widetilde{x}}(u)$.
  - They differ by a constant that depends on the approximation-interval size.

- Differential Entropy $\mathscr{H}_{\widetilde{x}}$ does still however measure information content when subsymbols in codewords are chosen (usually at random) from $p_{\widetilde{x}}(u)$.

- Maximum $\mathscr{H}_{\widetilde{x}}$ occurs when $p_{\widetilde{x}}(u)$ is **Gaussian** (any mean), with constant average energy.

$$\int_{-\infty}^{\infty} p_{\widetilde{x}}(u) \cdot \|u\|^2 \cdot du = \mathcal{E}_{\widetilde{x}}$$

Complex  $\quad \mathscr{H}_{\widetilde{x}} = \log_2(\pi e \mathcal{E}_{\widetilde{x}})$ bits/clpx-subsymbol

Real  $\quad \mathscr{H}_x = \frac{1}{2}\log_2(2\pi e \bar{\mathcal{E}}_x)$ bits/dimension

- More generally, $\text{trace}\{R_{\widetilde{x}\widetilde{x}}\} = \mathcal{E}_{\widetilde{x}}$.

$$\mathscr{H}_{\widetilde{x}} = \log_2|\pi e R_{\widetilde{x}\widetilde{x}}| \text{ bits/cplx-subsymbol}$$

# Gaussian MMSE & conditional entropy

- **Complex scalar Gaussian** $x$

$$p_x(u) = \frac{1}{\pi \sigma_x^2} e^{-\frac{|x|^2}{\sigma_x^2}}$$

$$\mathscr{H}_x = \log_2\{\pi \cdot e \cdot \sigma_x^2\}$$

- **Conditional** $x/y$ ?

$$\mathscr{H}_{x/y} = \log_2\{\pi \cdot e \cdot \sigma_{x/y}^2\}$$

$$\sigma_{x/y}^2 = \sigma_x^2 - \frac{r_{xy}^2}{\sigma_y^2} = \text{MMSE}$$

- **Vector** $\boldsymbol{x}$?

$$\mathscr{H}_{\boldsymbol{x}} = \log_2\{(\pi e)^{\bar{N}} \cdot |R_{\boldsymbol{xx}}|\}$$

$$\mathscr{H}_{\boldsymbol{x/y}} = \log_2\{(\pi e)^{\bar{N}} \cdot |R_{\boldsymbol{x/y}}^{\perp}|\}$$

$$R_{\boldsymbol{x/y}}^{\perp} = R_{\boldsymbol{xx}}^2 - R_{\boldsymbol{x/y}} \cdot R_{\boldsymbol{yy}}^{-1} \cdot R_{\boldsymbol{x/y}} = \text{MMSE}$$

(Appendix D on MMSE)

$\bar{N}$ is the number of complex dimensions = $N/2$

- If $\widetilde{x}$ and $\widetilde{y}$ are jointly Gaussian, then $p_{\widetilde{x}/\widetilde{y}}$ is also Gaussian and has mean as MMSE estimate $\mathbb{E}[\widetilde{x}/\widetilde{y}]$ and autocorrelation $R_{ee} = R_{\widetilde{x}\widetilde{x}} - R_{\widetilde{x}\widetilde{y}} \cdot R_{\widetilde{y}\widetilde{y}}^{-1} \cdot R_{\widetilde{y}\widetilde{x}}$ .

- $\mathcal{H}_{\widetilde{x}/\widetilde{y}} = \log_2 |\pi e R_{ee}|$ - that is, the entropy is essentially just the log of the MMSE (Gaussian).

  - Entropy generalizes MMSE to any probability distribution.
  - Measures the information content of the "miss" in estimating $\widetilde{x}$ from $\widetilde{y}$ for any $p_{\widetilde{x}\widetilde{y}}$.

# Mutual Information and SNR
## *Subsection 2.3.2*

*See PS2.5 (Prob 2.20)*

For Gaussian, $\mathcal{I}$ and (geo) SNR are in 1-to-1 relationship
MMSE and best rate are essentially same thing.

# Mutual Information ~ SNR

- Mutual Information is:

$$\mathcal{I}(\widetilde{\boldsymbol{x}}; \widetilde{\boldsymbol{y}}) = \mathbb{E}\left[\log_2\left(\frac{p_{\widetilde{x}\widetilde{y}}}{p_{\widetilde{x}} \cdot p_{\widetilde{y}}}\right)\right] = \mathscr{H}_{\widetilde{x}} - \mathscr{H}_{\widetilde{x}/\widetilde{y}} = \mathscr{H}_{\widetilde{y}} - \mathscr{H}_{\widetilde{y}/\widetilde{x}}.$$

- For discrete example $\mathcal{I}$ = 1-.811 = .189 bits/subsymbol.

$$= \mathcal{H}_{\widetilde{x}} - \mathcal{H}_{\widetilde{x}/\widetilde{y}} = \mathcal{H}_{\widetilde{y}} - \mathcal{H}_{\widetilde{y}/\widetilde{x}}.$$

- $\mathcal{I}(\widetilde{\boldsymbol{x}}; \widetilde{\boldsymbol{y}})$ is symmetric in $\widetilde{\boldsymbol{x}}$ and $\widetilde{\boldsymbol{y}}$ (MMSE forward and backward channel).

- $\mathcal{I}(\widetilde{\boldsymbol{x}}; \widetilde{\boldsymbol{y}})$ measures common ("mutual") information between $\widetilde{\boldsymbol{x}}$ and $\widetilde{\boldsymbol{y}}$ , $\mathbb{E}\left[\log_2\left(\frac{p_{\widetilde{x}/\widetilde{y}}}{p_{\widetilde{y}}}\right)\right]$.

  - On average, $\mathcal{I}(\widetilde{\boldsymbol{x}}; \widetilde{\boldsymbol{y}})$ measures how much bigger is unconditional info versus conditional info, in bits.

- $\mathcal{I}(\widetilde{\boldsymbol{x}}; \widetilde{\boldsymbol{y}}) = \log_2\frac{|R_{\widetilde{x}\widetilde{x}}|}{|R_{ee}|} = \log_2\frac{|R_{\widetilde{y}\widetilde{y}}|}{|R_{nn}|} = \log_2\left(\left(1 + SNR_{geo}\right)^{\overline{N}}\right)$ for the matrix AWGN.

  - OR as earlier for vector coding $\mathcal{I}(\widetilde{\boldsymbol{x}}; \widetilde{\boldsymbol{y}}) = \sum_{n=1}^{\overline{N}} \log_2 SNR_{mmse,n}$ for the matrix AWGN.

**Theorem 2.1.1 (Law of Large Numbers (LLN))** *The LLN observes that a stationary random variable $z$'s sample average over its observations $\{z_n\}_{n=1,\ldots,N}$ converges to its mean with large $N$ such that*

$$\lim_{N \to \infty} \quad Pr\left\{\left|\left(\frac{1}{N}\sum_{n=1}^{N} z_n\right) - \mathbb{E}[z]\right| > \epsilon\right\} \to 0 \quad weak\ form \qquad (2.13)$$

$$\lim_{N \to \infty} \quad Pr\left\{\frac{1}{N}\sum_{n=1}^{N} z_n = \mathbb{E}[z]\right\} = 1 \quad strong\ form\ . \qquad (2.14)$$

- Distribution of $z$ must be the same (stationary) for all random selections.

- The random $z$ can be function of random variable ( $z = f(x)$ )and the sample mean converges to $\mathbb{E}[f(x)]$.
  - E.g., $z_n = \|\boldsymbol{x}_n\|^2$ where the vector $\boldsymbol{x}_n$ might also have (a growing) $N$ components (energy sample or length of the vector).
  - LLN then states that all the energy (really points in selection from any distribution with $\mathbb{E}[\|\boldsymbol{x}\|^2] \leq \mathcal{E}_x$) of a hypersphere are are at its surface with probability 1. Points on the interior have probability zero. It is also a sum of independent terms, and thus Gaussian (central limit theorem).
  - The marginal distributions for the vector $\boldsymbol{x}_n$'s element selections, and thus for $\boldsymbol{x}_n$ also, would be Gaussian if this $N \to \infty$-sequence has max entropy (uniform).

- The function of most interest in coding is $f(x)$ = $-\log_2[p_x(x)]$- that is the function itself is probability distribution's log.
  - The **sample average** of this function **converges** to the **entropy**.
  - This suggests choosing codewords (this means each subsymbol in the codeword) at random from stationary distribution,
    - and then repeat at higher level for several codes chosen at random.
  - These are discrete codes, even when $\boldsymbol{x}$ is continuous, but their average follows the entropy (and mutual information).
  - Generalizes **sphere-packing** (which was for the AWGN only).

# Random coding generalizes 379A sphere packing

- Pick subsymbols $\boldsymbol{x}_n$ randomly (independently) from (stationary) distribution $p_{\widetilde{\boldsymbol{x}}}$ for each of $M = 2^b$ c'words.
  - This is one **random code**.

- Repeat the exercise for another code, and …. many more.

- Compute the average performance of all these random selected codes:
  - As $\overline{N} \to \infty$ , this average performance is outstanding (as we'll see), as long as $\tilde{b} < \mathcal{I}(\widetilde{\boldsymbol{x}}; \widetilde{\boldsymbol{y}})$.
  - So at least one good one must exist.

- Entropy per subsymbol is

- LLN with function $\log_2\left[p_{\tilde{z}}^{-1}\right]$ is the sample-average entropy estimate.

$$\widetilde{\mathcal{H}}_{\boldsymbol{x}} = \frac{-1}{\overline{N}} \cdot E\left[\log_2(p_{\boldsymbol{x}})\right]$$

$$= \frac{-1}{\overline{N}} \sum_{n=1}^{\overline{N}} E\left[\log_2(p_{\tilde{\boldsymbol{x}}_n})\right] \quad ,$$

$$\hat{\widetilde{\mathcal{H}}}_{\tilde{\boldsymbol{x}}} = \frac{-1}{\overline{N}} \cdot \sum_{n=1}^{\overline{N}} \log_2\left[p(\tilde{\boldsymbol{x}}_n)\right] = \frac{-1}{\overline{N}} \cdot \log_2\left[p(\boldsymbol{x})\right] \quad \cdot \quad$$ , which converges to (constant) $\widetilde{\mathcal{H}}_x$

- The constant means the ave code has uniform distribution of codewords (asymptotically), $2^{\overline{N} \cdot \widetilde{\mathcal{H}}_x}$ of them.

**Asymptotic Equal Partition (AEP)**

- The set is $A_{\overline{N}}^{\epsilon}(\boldsymbol{x}) \triangleq \left\{ \boldsymbol{x} = [\tilde{\boldsymbol{x}}_1, \tilde{\boldsymbol{x}}_2, ..., \tilde{\boldsymbol{x}}_{\overline{N}}] \,\middle|\, 2^{-\overline{N} \cdot \mathcal{H}_{\tilde{x}} - \epsilon} \le p(\tilde{\boldsymbol{x}}_1, \tilde{\boldsymbol{x}}_2, ..., \tilde{\boldsymbol{x}}_{\overline{N}}) \le 2^{-\overline{N} \cdot \mathcal{H}_{\tilde{x}} + \epsilon} \right\}$

**Lemma 2.3.6 [AEP Lemma]** *For a typical set with $\overline{N} \to \infty$, the following are true:*

- *$Pr\{A_{\overline{N}}^{\epsilon}(\boldsymbol{x})\} \to 1$*

- *for any codeword $\boldsymbol{x} \in A_{\overline{N}}^{\epsilon}$, $Pr\{\boldsymbol{x}\} \to 2^{-\overline{N} \cdot \mathcal{H}_{\tilde{x}}}$*

Decoder works well
if only one codeword
in conditional set for each
$\boldsymbol{y}$ value, so good code spreads
them uniformly.



$A_{\overline{N}}^{\epsilon}$

Prob $\to 1$

typical set

$2^{\overline{N} \cdot \overline{\mathcal{H}}_x}$

number of
sequences

all sequences



$A_{\overline{N}}^{\epsilon}(\boldsymbol{x})$

$2^{\overline{N} \cdot H_{\tilde{x}}}$

unconditional typical set

redundancy $\tilde{\rho} = H_{\tilde{x}/y}$ bits

$A_{\overline{N}}^{\epsilon}(\boldsymbol{x}/\boldsymbol{y})$

conditional
Typical set

$2^{\overline{N} \cdot H_{\tilde{x}/\tilde{y}}}$

all sequences

$\mathcal{H}_{\tilde{x}} = \mathcal{I}(\tilde{x}; \tilde{y}) + \mathcal{H}_{\tilde{x}/\tilde{y}}$
bits/subsymbol
$\le \tilde{b} + \tilde{\rho}$
So pick $\tilde{\rho}$ large enough.

There are $2^{N \cdot H_{\tilde{x}}} \cdot 2^{-N \cdot H_{\tilde{x}/\tilde{y}}} = 2^{N \cdot \mathcal{I}(\tilde{x}; \tilde{y})}$ little sets
In the big set if "equally partitioned"

# General Capacity Theorem

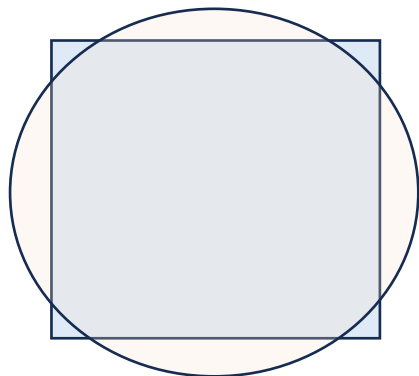$$\frac{|A_N^\epsilon(\boldsymbol{x})|}{|A_N^\epsilon(\boldsymbol{x}/\boldsymbol{y})|} \to 2^{\mathcal{I}(\boldsymbol{x};\boldsymbol{y})} \qquad \text{since } \mathcal{I}(\boldsymbol{x};\boldsymbol{y}) = \mathcal{H}_{\boldsymbol{x}} - \mathcal{H}_{\boldsymbol{x}/\boldsymbol{y}}$$

- Good codes will have only 1 codeword per conditional entropy subset.

- MAP detector decision region is then ~ $A_N^\epsilon(\boldsymbol{x}/\boldsymbol{y})$ - on average; but we can find it for one good code.

- If $A_N^\epsilon(\boldsymbol{x})$ were any larger, all codes (good or bad) will have at least one $A_N^\epsilon(\boldsymbol{x}/\boldsymbol{y})$ that contains 2+ codewords, which mean the MAP has to "flip a coin" – not good (high error prob).

- SHANNON's CAPACITY THEOREM
  - Number of codewords is limited by mutual info $b \leq \mathcal{I}(\boldsymbol{x};\boldsymbol{y})$.

  - Which is per-subsymbol equivalent with random code $\tilde{b} \leq \mathcal{I}(\widetilde{\boldsymbol{x}};\widetilde{\boldsymbol{y}})$.

  - If maximized over input distributions $\tilde{b} < \tilde{\mathcal{C}} \leq \max_{p_{\widetilde{x}}} \mathcal{I}(\widetilde{\boldsymbol{x}};\widetilde{\boldsymbol{y}}) \frac{\text{bits}}{\text{subsymbol}}$.

# The uniform part is most important.

- The Gaussian distribution corresponds to marginal of uniform distribution over a hypersphere.
  - This uniform distributions marginals are asymptotically Gaussian.
  - This is a special case where uniform and Gaussian are basically the same.
  - Because all the Gaussian infinite-length vectors (codewords) have same energy (zero variance of the energy).

- All the points (really volume) are (is) at the surface.

- The Gaussian marginal dist'n is important only for shaping gain (< 1.53 dB).

- The (AEP) uniform spacing of points (no matter where the majority of them sit, surface or otherwise) remains for the fundamental gain.

The uniform spacing separates codewords in the union of the hypersquare (orthotope) and hypersphere.

Thus, good codes can be based on sequences from uniformly spaced PAM/QAM subsymbols.

And the rest is MMSE Estimation,

With a chain-rule twist in Some situations

Vector Coding is always all MAP, All ML, all MMSE special case.

Stanford University

# End Lecture 3