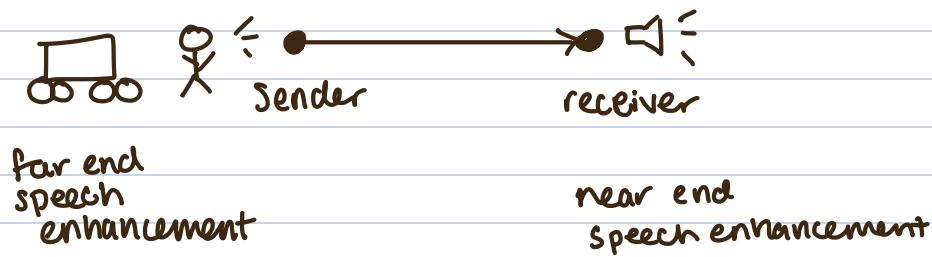


# Microphone Array Processing

Project 2: Design & build a multi-mic speech enhancement / beamforming system for far-end noise reduction



Ex: Single mic far end noise reduction



## Speech Production (light notes due to past coursework)

- unvoiced /s/ and voiced /a/ and mixed /z/ and plosive /p/ (noisy) (periodic) (buildup of pressure)
- speech sound signal = excitation signal
- formants (3-5 formants within Nyquist band)
  - ↳ resonance frequencies give rise to peaks in the overall spectrum
- Spectrograms → "see speech"
- Speech signals
  - = vocal tract excitation + vocal tract filter
- vocal tract changes over time → time-varying
  - only short segments of speech can be assumed to have similar acoustic properties (non-stationary / short term stationary)
- speech is stochastic
  - ↳ 20-30 ms time frames
- max speech BW  $\sim 7\text{-}8 \text{ kHz}$

## Microphone Measurement Model

- Direct path  $x[n] = a(d)s[n - \tau(d)]$ 
  - ↳ orig sound signal delayed by  $\tau(d)$ , time it takes to travel distance d
  - ↳ Scaling factor, depends on dist
  - ↳ Sound goes to mic w/ some attenuations and delay

- reflections are modeled with a room impulse response

$$x[n] = (h * s)[n]$$

room impulse response

↳ n captures reflections

- direct path and early reflections determine intelligibility
- reverb (late refl.) degrades intelligibility

### Single Mic Model:

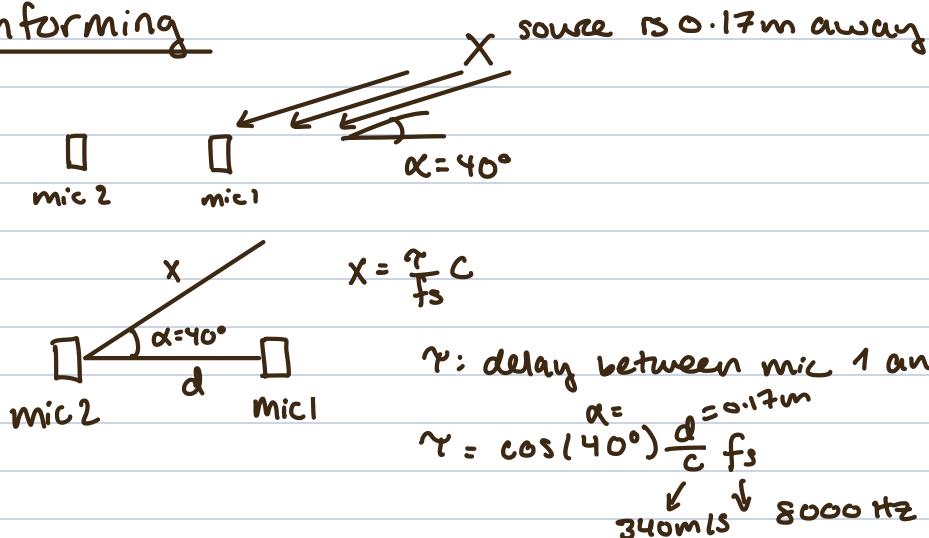
$$x[n] = \underbrace{\sum_{i=1}^d (h_i * s_i)[n]}_{\text{early}} + \underbrace{\sum_{i=1}^d (h_i * s_i)[n]}_{\text{late}} + n[n]$$

- $d$  point sources  $s_i$
- room impulse response  $h_i$  from source position  $i$  to mic
- $n[n]$ : mic self noise & noise (e.g. late reverberation)

Assume

- sources are additive, zero mean & uncorr, short term stationary
- Early & late components correlated via source  $s_i$

### Beamforming



$$\gamma = 3.06 \text{ samples, delay b/w. mics}$$

↳ non-int. → need to use freq. domain phase change  
to best descr.

Narrowband assumption?

- ↳ max delay across array is less than sampling period  $T_s$
- ↳ in audio  $T_s = 1/8000 \rightarrow$  meaning narrowband cannot apply for audio
- ↳ to we, need to process per frequency band
  - ↳ beamformer resp. is freq. dependent
  - ↳ phase shifts become freq. dependent

Spatial aliasing

$$d < \frac{1}{2} \lambda_{\min} < \frac{1}{2} \frac{c}{\frac{1}{2} f_s} = \frac{c}{f_s}$$

\* Narrowband can hold BUT ONLY if signal is processed in narrow frequency bands using the DFT

- delay  $T$  needs to be small enough to satisfy NBG in each f band

## Delay & Sum Beamformer (intro)

K: freq. index in discrete frequency domain

- 2 mics, one mic receives a delayed version
- sum coherently to get better SNR

$$\hat{S}(K, \ell) = \frac{x_1(K, \ell) + x_2(K, \ell) e^{j2\pi K \gamma / N}}{2} = S(K, \ell) + \frac{N_1(K, \ell) + N_2(K, \ell) e^{j2\pi K \gamma / N}}{2}$$

→ you have far field & near field effects to account for

**NEAR**

$$S(K, \ell) = s(K, \ell) a(d) e^{-j2\pi K \gamma(d) / N}$$

- Source close to array center, each mic gets a different  $a(d)$  and phase difference  $\gamma$

**FAR**

$$s(K, \ell) = s(K, \ell) e^{-j2\pi K T(d) / N}$$

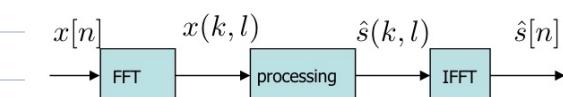
- only  $\gamma$

**FREE FIELD**

- only direct path, no refl.

## STFT

- time frames need to be short due to speech nature
  - STFT makes conv. ≈ mult. & helps satisfy narrowband assumption
- $$x(K, \ell) = \sum_{i=1}^d a_i(K, \ell) s_i(K, \ell) + n(K, \ell)$$
- for M mics  $x(K, \ell) = \sum_{i=1}^d a_i(K, \ell) s_i(K, \ell) + n(K, \ell)$



## Problem Formulation

$$\underline{x}(K, \ell) = \sum_{i=1}^d a_i(K, \ell) s_i(K, \ell) + n(K, \ell)$$

$$\begin{aligned} \mathbf{x}(k, l) &= \underbrace{\mathbf{a}_1(k, l)s_1(k, l)}_{\text{target}} + \underbrace{\sum_{i=2}^d \mathbf{a}_i(k, l)s_i(k, l)}_{\text{interferers+noise}} + \mathbf{n}'(k, l) \\ &= \mathbf{a}(k, l)s(k, l) + \mathbf{n}(k, l) \end{aligned}$$

- Goal: Estimate  $s(k, l)$  given  $\mathbf{x}(k, l)$ : e.g.  $\hat{s}(k, l) = E[s(k, l)|\mathbf{x}(k, l)]$
- 1) Derive beamformers assuming  $\mathbf{a}(k, l)$  is known.
- 2) estimation of  $\mathbf{a}(k, l)$

$\underline{a}(K, \ell)$  : STFT of room impulse resp. per frequency stacked across mics

when normalized wrt ref mic → called RTF → shortens length of resp !  
 (100 ms - 1 s)  
 Room impulse resp. >> frame size (20 ms)  
 → est. target at ref. mic

## DELAY & SUM BEAMFORMER

our proj: multi mic / far end noise reduction

general idea: signal narrowband, time delay is fshift in Fourier domain, sum coherently after multiplying each sensor signal by a compensatory phase shift

$$\underline{x}(k, \ell) = s(k, \ell) \underline{a}(k, \ell) + \underline{n}(k, \ell)$$

↓ models the ATF

(k, ℓ)  
freq., time

$$\hat{s}(k, \ell) = w^H(k, \ell) \underline{x}(k, \ell)$$

$$\underline{a}(k, \ell) = \left[ 1, \frac{a_2 e^{-j 2\pi k \tau_2 / N}}{a_1}, \dots, \frac{a_m e^{-j 2\pi k \tau_m / N}}{a_1} \right]^T \rightarrow \text{free \& near}$$

$$\text{near \& free field: } \underline{w}(k, \ell) = \frac{\underline{a}(k, \ell)}{\underline{a}^H(k, \ell) \underline{a}(k, \ell)}$$

$$\text{far \& free field: } \underline{w}(k, \ell) = \frac{1}{m} \underline{a}(k, \ell) \rightarrow \text{I believe here, } \\ \underline{a}(k, \ell) = \left[ 1, e^{-j 2\pi k \tau_2 / N}, \dots \right]$$

$$\text{general case: } \underline{w}(k, \ell) = \frac{\underline{a}(k, \ell)}{\underline{a}^H(k, \ell) \underline{a}(k, \ell)}$$

because you don't account the  $a(d)$  in far field

### Characteristics

- "preserves the target"
- no explicit knowledge of noise  $\rightarrow$  good esp. if noise changes
- reduces noise var  $\frac{1}{n} = \frac{1}{Z^2} \Rightarrow -P \log_{10}(2) \approx -3P \text{ dB}$

## MVDR BEAMFORMER

Minimum Variance Distortionless Response  
aka Capon beamformer

- also exploits target location info
- additionally exploits position of noise sources

- no change of mag or phase in dir of source
- min. variance of beamformer output in all other directions

$$\text{cost function: } J(\underline{w}(k, \ell)) = \underline{w}^H(k, \ell) \underline{\underline{R}_x}(k, \ell) \underline{w}(k, \ell)$$

↳ minimizes output power s.t. response in desired direction is 1 (this is called the distortionless constraint)

$$\min_{\underline{w}} \underline{w}^H(K, \ell) \underline{\underline{R}_x}(K, \ell) \underline{w}(K, \ell)$$

$$\text{s.t. } \underline{w}(K, \ell)^H \underline{\underline{\alpha}}(K, \ell) = 1$$

\* solve by taking deriv, set = 0,  
plug in constraint and you  
get the following answer

$$\rightarrow \underline{w}(K, \ell) = \frac{(\underline{\underline{R}_x}(K, \ell))^{-1} \underline{\underline{\alpha}}(K, \ell)}{\underline{\underline{\alpha}}^H(K, \ell) (\underline{\underline{R}_x}(K, \ell))^{-1} \underline{\underline{\alpha}}(K, \ell)}$$

- You can also write this result by using  $\underline{\underline{R}_n}(K, \ell)$

$$\underline{w}(K, \ell) = \frac{\underline{\underline{R}_n}^{-1}(K, \ell) \underline{\underline{\alpha}}(K, \ell)}{\underline{\underline{\alpha}}^H(K, \ell) \underline{\underline{R}_n}^{-1}(K, \ell) \underline{\underline{\alpha}}(K, \ell)} \quad \underline{\underline{R}_x}(K, \ell) = \underline{\underline{R}_n}(K, \ell) + \underline{\underline{\alpha}}(K, \ell) \underline{\underline{\alpha}}^H(K, \ell) \sigma_s^2 / (K, \ell)$$

Using MVDR Requires:

①  $\underline{\underline{R}_s}(K, \ell)$  is rank 1

② target & noise are uncorrelated

③ target & noise are additive,

MVDR Spatially Uncorr. Noise: MVDR becomes the delay & sum beamformer

### OPTIMAL LINEAR Multi CHANNEL WIENER

↳ minimize MSE between the estimated signal and the true signal

Signal model:  $\underline{x}(K, \ell) = \underline{s}(K, \ell) \underline{\underline{\alpha}}(K, \ell) + \underline{n}(K, \ell)$

$$\min \mathbb{E} [\|\underline{s}(K, \ell) - \underline{w}^H(K, \ell) \underline{x}(K, \ell)\|_2^2]$$

$$(\rightarrow \underline{w}(K, \ell) = \underline{\underline{R}_x}^{-1}(K, \ell) \sigma_s^2 \underline{\underline{\alpha}}(K, \ell)$$

$$= \frac{\sigma_s^2(K, \ell)}{\underbrace{\sigma_s^2(K, \ell) + (\underline{\underline{\alpha}}^H(K, \ell) \underline{\underline{R}_n}^{-1}(K, \ell) \underline{\underline{\alpha}}(K, \ell))^{-1}}_{\text{single ch w.F.}}} \cdot \frac{\underline{\underline{R}_n}^{-1}(K, \ell) \underline{\underline{\alpha}}(K, \ell)}{\underbrace{\underline{\underline{\alpha}}^H(K, \ell) \underline{\underline{R}_n}^{-1}(K, \ell) \underline{\underline{\alpha}}(K, \ell)}_{\text{MVDR}}}$$

↓  
post process MVDR filter to  
further reduce noise variance

↓  
noise var set to remaining noise PSD after beamforming

↓  
provides spatial filtering to  
suppress noise & interference

## Sufficient Statistics

↳ a function of the data that captures all info you need to estimate a parameter of interest

$$\underbrace{T(\underline{x}(k, \ell))}_{\downarrow} = \underbrace{\underline{w}^H_{\text{MVDR}}(k, \ell) \underline{x}(k, \ell)}_{\text{full data}} = \frac{\underline{a}^H(k, \ell) \underline{R}_x^{-1}(k, \ell) \underline{x}(k, \ell)}{\underline{a}^H(k, \ell) \underline{R}_x^{-1}(k, \ell) \underline{a}(k, \ell)}$$

the statistic is sufficient meaning using the statistic  $T(\cdot)$  instead of the full data  $\underline{x}(k, \ell)$  does not lose any information about  $s(k, \ell)$

### LGMV BEAMFORMER

Linearly Constrained Minimum Variance

- generalizes the MVDR by allowing multiple linear constraints
- generalize with d linear constraints

$$\begin{aligned} \min_{\underline{w}(k, \ell)} \quad & \underline{w}^H(k, \ell) \underline{R}_x(k, \ell) \underline{w}(k, \ell) \\ \text{s.t.} \quad & \underline{w}^H(k, \ell) \underline{\Delta}(k, \ell) = \underline{f}^H(k, \ell) \quad \Delta \in \mathbb{C}^{M \times d} \end{aligned}$$

$$\text{when } d < M, \quad \underline{w}(k, \ell) = \underline{R}_x^{-1}(k, \ell) \underline{\Delta}(k, \ell) (\underline{\Delta}^H(k, \ell) \underline{R}_x^{-1}(k, \ell) \underline{\Delta}(k, \ell))^{-1} \underline{f}(k, \ell)$$

How can you use the multiple constraints?

- steer zeros in direction of noise sources
- maintain sig in certain dir.

\* More constr = less deg of freedom to control noise reduction

### ESTIMATION OF THE A.T.F.

- all beamformers require an estimation of the ATF

$$X_m[n] = (h_m * s)[n] \quad s: \text{point source}$$

m: microphone

h\_m: room impulse response

X\_m: noise free source signal given at m-th microphone

$$x_m(n, \ell) = a_m(k, \ell) s(k, \ell)$$

a\_m: acoustic transfer function from source to m-th microphone

↳ temporal-frequency domain representation of the room impulse response

$$\underline{x}(k, \ell) = \underline{a}(k, \ell) \underline{s}(k, \ell) \rightarrow \text{vector form: collecting } M \text{ microphone DFT coefficients in a vector}$$

$$\text{Relative ATF: } \underline{a}'(k, \ell) = \left[ 1, \frac{a_2(k, \ell)}{a_1(k, \ell)}, \dots \right]^T \text{ normalized wrt } a_1.$$

## Gross PSD Matrices

$$\mathbb{E} \left[ \underbrace{\underline{x}(k, \ell)}_{\text{fourier domain}} \underline{x}(k, \ell)^H \right] = \mathbb{E} \left[ \underbrace{\underline{s}(k, \ell) \underline{s}(k, \ell)^H}_{\text{describing how correlated signals are in the frequency domain}} \right] + \mathbb{E} \left[ \underline{n}(k, \ell) \underline{n}(k, \ell)^H \right]$$

describing how correlated signals are in the frequency domain

$$\underline{\underline{R}_x}(k, \ell) = \underline{\underline{R}_s}(k, \ell) + \underline{\underline{R}_n}(k, \ell)$$

Example: one point source

$$\underline{s}(k, \ell) = \underline{a}(k, \ell) \underline{s}_1(k, \ell)$$

$$\underline{a} \in \mathbb{C}^M : \text{RTF}$$

$$\underline{\underline{R}_s} = \mathbb{E} \left[ \underline{s}(k, \ell) \underline{s}(k, \ell)^H \right] = \sigma_{s_1}^2 \underline{a}(k, \ell) \underline{a}(k, \ell)^H$$

$$\mathbb{E} [\underline{s}_1(k, \ell) \underline{s}_1(k, \ell)^H]$$

$$= \mathbb{E} [|\underline{s}_1(k, \ell)|]^2$$

= var of clean sig at ref mic 1

$\underline{\underline{R}_s} = \underline{\underline{I}}$  in this case, should match up to # sources

\* when multiple sources, beamformers a funct. of  $\underline{\underline{R}_s}$  \*

To solve beamforming

① need a good estimate of  $\underline{a}(k, \ell)$  (ATF/RTF)

② if multiple sources, will need a good est. of  $\underline{\underline{R}_s}$

the beamforming gains also need to be adapted for multiple sources

Recall Eigenvalues:

$$A = T \Delta T^{-1}, \text{ if } A \text{ is hermitian } A = T \Delta T^H$$

Estimate  $\underline{a}(k, \ell)$  - No Noise

Using EVD

assume  $\underline{\underline{R}_s}$  is perfectly known,  $\underline{\underline{R}_s} = \sigma_{s_1}^2(k, \ell) \underline{a}(k, \ell) \underline{a}^H(k, \ell)$

$$\underline{\underline{R}_s} = U \Delta U^{-1} = U \Delta U^H$$

The scaled atf is given by  $\underline{u}_i$ :  $\frac{\underline{a}(k, \ell)}{\|\underline{a}(k, \ell)\|} = \underline{u}_i$

## Estimate $\underline{a}(k, \ell)$ - Spatially white noise

using EVD

$$\text{Suppose } \underline{\underline{R}}_x(k, \ell) = \underline{\underline{R}}_s(k, \ell) + \underline{\underline{R}}_n(k, \ell) \\ = \sigma_s^2 \underline{a}(k, \ell) \underline{a}^H(k, \ell) + \sigma_n^2 \underline{\underline{I}}$$

$$\underline{\underline{R}}_x = \underline{U} (\Delta + \sigma_n^2 \underline{\underline{I}}) \underline{U}^H \quad * \underline{\underline{a}} \text{ only effects eigenvalues}$$

$$\text{so still estimate ATF by } \frac{\underline{a}(k, \ell)}{\|\underline{a}(k, \ell)\|} = \underline{u},$$

\* For spatially white noise,  $\underline{\underline{R}}_x$  and  $\underline{\underline{R}}_s$  share same eigenvectors  $\underline{u}$

Ultimately, you want to estimate  $\underline{\underline{R}}_s$  b/c  $\underline{\underline{R}}_s = \underline{U} \Delta \underline{U}^H$  and

$$\underline{\underline{U}}_1 = \frac{\underline{a}(k, \ell)}{\|\underline{a}(k, \ell)\|} \quad \text{so by getting } \hat{\underline{\underline{R}}}_s, \text{ you can get } \underline{a}(k, \ell)$$

New problem formulation  $\Rightarrow$  estimate  $\underline{\underline{R}}_s$

Estimate  $\underline{\underline{R}}_s$  under spatially white noise

$$\text{rank}(\underline{\underline{R}}_s) = r < M$$

$$\underline{\underline{R}}_x = (\underline{U}_1 \quad \underline{U}_2) \begin{pmatrix} \Delta + \sigma_n^2 \underline{\underline{I}}_r & \underline{0} \\ \underline{0} & \sigma_n^2 \underline{\underline{I}}_{M-r} \end{pmatrix} \begin{pmatrix} \underline{U}_1^H \\ \underline{U}_2^H \end{pmatrix}$$

$\underline{U}_1$  to  $\underline{U}_r$  span speech + noise subspace  
 $\underline{U}_{r+1}$  to  $\underline{U}_M$  span only noise subspace

$$\text{so } \hat{\underline{\underline{R}}}_s = \underline{U}_1 (\Delta_1 + \sigma_n^2 \underline{\underline{I}}_r) \underline{U}_1^H$$

simplifies to

$$\hat{\underline{\underline{R}}}_s = \underline{U}_1 \Delta_1 \underline{U}_1^H$$

theoretically exact but i/r noise introduces error

Estimate  $\underline{\underline{R}}_s$  under spatially colored data

Using EVD

$$\textcircled{1} \text{ Pre-whiten data: } \underline{\underline{R}}_n^{1/2} : \tilde{\underline{x}} = \underline{\underline{R}}_n^{1/2} \underline{x}$$

$$\textcircled{2} \text{ } \underline{\underline{R}}_x = \tilde{\underline{U}} (\tilde{\Delta} + \underline{\underline{I}}_m) \tilde{\underline{U}}^H, \text{ truncate } M-r \text{ smallest eig. to get } \tilde{\underline{U}}, \tilde{\Delta}, \tilde{\underline{U}}^H = \hat{\underline{\underline{R}}}_s$$

$$\textcircled{3} \text{ Deconviten result: } \hat{\underline{\underline{R}}}_s = \underline{\underline{R}}_n^{1/2} \hat{\underline{\underline{R}}}_s \underline{\underline{R}}_n^{1/2}$$

\* Note: using  $R_n^{1/2}$  to decenter may result in a loss in accuracy

Therefore we do another method: generalized EVD

### GENERALIZED EVD FOR A.T.F. ESTIMATION

- hermitian matrices  $\underline{A}$  and  $\underline{B}$
- diagonalize both with  $\underline{U}$  → columns of  $\underline{U}$  are the generalized eigenvectors of  $\underline{A}$  and  $\underline{B}$
- the generalized eigenvalues and eigenvectors of  $\underline{A}, \underline{B}$  are the ordinary eigenvalues and eigenvectors of  $\underline{B}^{-1}\underline{A}$

Derivation

$$\begin{aligned} \underline{U}^H \underline{A} \underline{U} &= \Lambda_A \\ \underline{U}^H \underline{B} \underline{U} &= \Lambda_B \end{aligned} \quad \left. \begin{array}{l} \underline{A} \underline{U} = \underline{B} \underline{U} \Lambda \\ \Lambda = \Lambda_B^{-1} \Lambda_A \end{array} \right\} \quad \rightarrow \text{the generalized eigenvalues in } \Lambda \text{ are } \lambda_i = \frac{a_i}{b_i}$$

$$\underline{R}_x = \underline{U}^{-H} (\Lambda + I_M) \underline{U}^{-1} \quad \rightarrow \text{then assume rank}(R_s) = r < M$$

→ then partition,

$$\underline{R}_x = (Q_1, Q_2) \begin{pmatrix} \Lambda_r + I_r & 0 \\ 0 & I_{M-r} \end{pmatrix} \begin{pmatrix} Q_1^H \\ Q_2^H \end{pmatrix} \quad \begin{array}{l} Q_1 \in \mathbb{C}^{M \times r} \\ Q_2 \in \mathbb{C}^{M \times (M-r)} \end{array}$$

→ Estimate  $R_s$  from GEVD of  $R_x$  :  $\hat{R}_s = Q_1 (\Lambda_r + I_r) Q_1^H$

→ Simplify to (when one source) :  $\hat{R}_s = Q_1 \Lambda_r Q_1^H$

If  $\text{rank } R_s = 1$ , then ATF can be obtained by selecting  $q_1$ , the principle generalized eigenvector between  $R_s(n, e)$  and  $R_n(n, e)$

### GEVD application to Beamforming

One way to write  $R_x$  is,  $R_x = Q_1 (\Lambda_r + I_r) Q_1^H + Q_2 Q_2^H$

Recall beamformer takes linear combos of microphone signals ( $\hat{S} = W^H X$ )

$$\text{Thus, } \hat{R}_s = W^H R_x W = W^H Q_1 (\Lambda_r + I_r) Q_1^H W + W^H Q_2 Q_2^H W$$

\* A good beamformer can be expressed as linear combos of  $U$ , cols  $w = U_1 b$ ,  $b \in \mathbb{C}^r$

## PERFORMANCE METRICS

- ① output SNR
- ② MSE
- ③ noise reduction
- ④ speech distortion

### Output SNR

$$\text{SNR}_{\text{out}}(\underline{w}) = \frac{\underline{w}^H \underline{R}_S \underline{w}}{\underline{w}^H \underline{R}_n \underline{w}}$$

Notably,

$$\nabla_{\underline{w}^H} \text{SNR}_{\text{out}}(\underline{w}) = \mathbf{0} \text{ meaning } \underline{w} \text{ is a stationary point}$$

(as per operator theory)

So if we consider  $\underline{w}$  to be a generalized eigenvector, then,

$\underline{w} = \underline{w}_1$  maximizes output SNR

$$\frac{\underline{w}^H \underline{R}_S \underline{w}}{\underline{w}^H \underline{R}_n \underline{w}} = \lambda \text{ where } \lambda \text{ is a generalized eigenvalue}$$

This lets us say,

$$\text{SNR}_{\text{out}}(\underline{w}) \leq \max_{\underline{w}} \frac{\underline{w}^H \underline{R}_S \underline{w}}{\underline{w}^H \underline{R}_n \underline{w}} = \lambda_1$$

\* Note: result independent of  $r = \text{rank}(\underline{R}_S)$

### MSE

$$\begin{aligned} \mathbb{E} [|\underline{w}^H \underline{x} - \underline{s}_1|^2] &= \mathbb{E} [|\underline{w}^H \underline{s}_1 + \underline{w}^H \underline{n} - \underline{s}_1|^2] \\ &= \underbrace{\mathbb{E} [|\underline{w}^H \underline{s}_1 - \underline{s}_1|^2]}_{\text{signal distortion}} + \underbrace{\mathbb{E} [|\underline{w}^H \underline{n}|^2]}_{\text{residual noise var}} \end{aligned}$$

$$\begin{aligned} \min \quad & \mathbb{E} [|\underline{w}^H \underline{s}_1 - \underline{s}_1|^2] \\ \text{s.t.} \quad & \mathbb{E} [|\underline{w}^H \underline{n}|^2] \leq C \quad 0 \leq C \leq \sigma_n^2 \rightarrow \text{noise at ref. mic before beamforming} \end{aligned}$$

$$\begin{aligned} \text{solution: } \underline{w}^* &= \underline{U} \underline{b}^* \\ &= \underline{U} (\underline{\Lambda}_1 + \mu \underline{I}_M)^{-1} \underline{U}^H \underline{R}_S \underline{e}_1 \end{aligned}$$

Lagrange mult.  $\mu \geq 0$  chosen such that  $\underline{b}^H \underline{b} = 0$

→ when we have  $\text{rank}(\underline{R}_S) = r < M$ , then  $\underline{R}_S = \underline{Q}_1 \underline{\Lambda}_1 \underline{Q}_1^H$   
so

$$\underline{w}^* = \underline{U}_1 (\underline{\Lambda}_1 + \mu \underline{I}_r)^{-1} \underline{\Lambda}_1^{-1} \underline{Q}_1^H \underline{e}_1$$

∴ MMSE opt. beamformer can be expressed as lin combo of  $\underline{U}_1$  columns

$$\rightarrow \text{another math simplification,}$$

$$W^* = (\underline{R}_S + M \underline{R}_n)^{-1} \underline{R}_S \underline{q}_1$$

\*  $\mu$  is a "trade off param" controlling signal dist. & noise reduction

\* when  $M = 1 \rightarrow$  multi-ch W.F.  $W_{MWF} = \underline{R_x^{-1}} \underline{R_s} \underline{e_1} = \underline{\tau_s}^2 \underline{R_x^{-1}} \underline{a}$

\* when  $\mu = 0$  and rank  $r$ ,  $W^* = U Q_1^{-H} e_1 = U Q_1 e_1 \rightarrow$  the MVDR beamformer  
 |  $\downarrow$  response distortionless

When  $r = 1$ ,  $w^* = \lambda^{-1} \sigma_s^2 R_n^{-1} \underline{a}$

$$\text{and doing more simp., } W^* = \frac{\underline{R}^n \underline{\alpha}}{\underline{\alpha}^* \underline{R}^n \underline{\alpha}}$$

\* General Case again but with  $r = 1$ ,

$$W^* = \left( \underline{\underline{R}}_S + \mu \underline{\underline{R}}_n \right)^{-1} \underline{\underline{R}}_S \underline{\underline{q}}_1$$

$$= \frac{\sigma_{S_1}^2}{\sigma_{S_1}^2 + \mu (\underline{\underline{a}}^H \underline{\underline{R}}_n^{-1} \underline{\underline{a}})^{-1}} \cdot \frac{\underline{\underline{R}}_n^{-1} \underline{\underline{a}}}{\underline{\underline{a}}^H \underline{\underline{R}}_n^{-1} \underline{\underline{a}}}$$

MVRD  
singulär W.F.

\* SDW MWF can be  
MVDR and then S.C.W.F.

## ESTIMATE A.T.F. VIA CRAMER RAD BOUND

→ presented by PhD student

\* Using C.R. L.B. to do ATF estimation

\* seems like a side quest, a potential option of further implementation

↳ a challenge

↳ will skip details for now

## ESTIMATE A.T.F. WITH INTER FREQUENCY CORRELATION

- improve acoustic parameter est. by exploiting hidden corr. across freq.
    - ↳ More state of the art research → same, challenge

## Projects: Multi-Mic Noise Reduction (class notes)

↳ at some point you need to estimate A.T.F.  $\rightarrow$  do via EVD or GEVD

下

- ! only accurate when  $R_x$  and  $R_n$  known! estimation errors severely effect results

\*Grad project: combine DL / GEVD to estimate ATF

## PERCEPTION & BINAURAL SP

\*Binaural typically refers to hearing aids  
and is out of scope of project

### Measuring intelligibility

- matrix test
- speech recognition threshold
- audiogram

SII : speech intel. index (weighted SNR to intelligibility)

STI, CSII, STOI, HASPI

### Hints for Project 2

Signal model:  $x_m[t] \xrightarrow{\text{we called } n \text{ in the notes}} (s_* * h_{1,m})[t] + \sum_{p=2}^P (s_p * h_{p,m})[t]$

Constructing the noisy signal:

- For sources  $s_p$  and microphone  $m$ :  $x_m[t] = (s_1 * h_{1,m})[n] + \sum_{p=2}^P (s_p * h_{p,m})[t]$
- Processing using STFT (i.e., using short time frames of 20 ms): window and FFT the samples  $x_m[\text{overlap}(l-1) + 1 : \text{overlap}(l-1) + \text{frsize}]$

Estimating Correlation matrices:  $\mathbf{R}_n(k, l) = E[\mathbf{n}(k, l)\mathbf{n}^H(k, l)]$  and  $\mathbf{R}_x(k, l) = E[\mathbf{x}(k, l)\mathbf{x}^H(k, l)]$

- Assuming ergodicity (sources are spatially invariant) you can estimate  $\mathbf{R}_n(k, l)$  e.g. as

$$\hat{\mathbf{R}}_n(k, l) = \frac{1}{N} \sum_{p=l-M_1}^{l+M_2} \mathbf{n}(k, p)\mathbf{n}^H(k, p)$$

or as

$$\begin{aligned}\hat{\mathbf{R}}_n(k, l) &= \begin{cases} \hat{\mathbf{R}}_n(k, l-1)\alpha + \mathbf{n}(k, l)\mathbf{n}^H(k, l)(1-\alpha) & \text{target not present} \\ \hat{\mathbf{R}}_n(k, l-1) & \text{target is present} \end{cases} \\ \hat{\mathbf{R}}_x(k, l) &= \hat{\mathbf{R}}_x(k, l-1)\alpha + \mathbf{x}(k, l)\mathbf{x}^H(k, l)(1-\alpha)\end{aligned}$$

- How to know whether the target is present or not? Either cheat by using directly the mix of interferers, or build a detector.

# Project TODOS / Outline

## Part 1 : generate signals

- use audio files given + impulse responses

$$x_m[t] = (s_* * h_{,,m})[t] + \sum_{p=2}^P (s_p * h_{p,m})[t]$$

$m$ : can be arbitrary but let's say  $m=4$  for now  $\rightarrow$  enough to get better SNR if beamforming correct

going to assume only one intended speaker per setup. means only one intended source to reconstruct

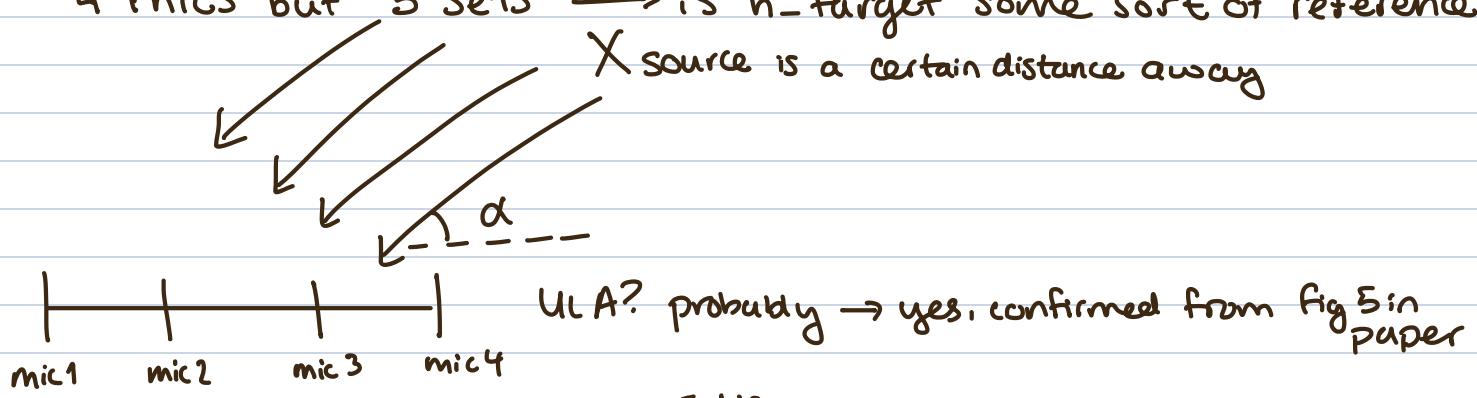
two clean speech signals given, microphone impulse resp given as mat file

Matlab file:

h-inter1	$4 \times 400$
h-inter2	$4 \times 400$
h-inter3	$4 \times 400$
h-inter4	$4 \times 400$
h-target	$4 \times 400$

}  $\star$  let's see if we can plot these convolve with sound signal to get signal at mic locations  
 → assumption: sampled at the same sampling frequency as the speech signals  
 $\hookrightarrow$  16K sampl.  $\rightarrow$  16 samples per ms

"4 mics" but "5 sets"  $\rightarrow$  is h-target some sort of reference?



$$x_{\text{mic}1} = s_* * h_{,,1} + \sum_{p=2}^5 (s_p * h_{p,\text{mic}1})[t]$$

$\star$  In report also probably need to write something to defend the narrowband assumption

$\star$   $P = 5$  would imply 5 sources but two clean speech signals gotten?

$\downarrow$   
 you can make any combination of the 5 sources

## About the Speech files given:

- 16-bit sampled at 16kHz
- Signal BW is at maximum  $\frac{16\text{kHz}}{2} = 8\text{kHz} \rightarrow$  wideband speech

S1 : clean speech 1 .wav	36.1 sec	one female speaker
S2 : clean speech 2 .wav	33.61 sec	4 female, 4 male speakers
S3 : bubble noise .wav	37 sec	
S4 : artif nonstat noise .wav	36.1 sec	
S5 : Speech shaped noise.wav	90 sec	

\* S2 can have a direction so becauf. should work for that

"constant interferers"  $\rightarrow$  S4, S5  
 "nonst. interferers"  $\rightarrow$  S2, S3

- \* Start slow with only one clean, one noise
- \* make all same length
- \* can also start with even shorter segments to make data more workable at the start like 5s / 10s

\* Need to also check the SNR of the created signals

$\hookrightarrow$  negative SNRs may not be recoverable ... consider setting up a test with S1 pos. SNR then S1 negative SNR maybe

Setup 1  
 S1: clean 1  
 S2: clean 2  
 S3: nonstat n  
 S4: speech sh. #1 (splitup)  
 S5: speech sh #2

"S1 is only clean,  
 this is ideal"  
 "no bubble"

Setup 2  
 S1: clean 1 + speech sh. noise 1  
 S2: clean 2  
 S3: nonstat n  
 S4: speech sh. #1 (splitup)  
 S5: speech sh #2

"S1 is now noisy"  
 "no bubble"

Setup 3  
 S1: same as 2  
 S2: clean 2  
 S3: bubble  
 S4: nonstat  
 S5: speech sh #2

"S1 noisy"  
 "bubble"

} not enough "base cases"

IDEAL: clean 1, all other sources off

Scen. 1 clean 1 + clean 2 (noise)

• high, medium low SNR

Scen. 2 clean 1 + bubble (noise)

• H, M, L SNR

Scen. 3 clean 1 + art. non. stat (noise)

• H, M, L SNR

Scen. 4 clean 1 + speech shaped (noise)

• H, M, L SNR

Scen. 5 (clean 1 + non-st.) + (clean 2 + speech-like noise)

Scen. 6 (clean 1) + (bubble) + (clean 2)

& more later maybe

high SNR: 25 dB

med JNR: 15

low SNR: 5

} do this to start, maybe just start w/ Med. to figure out what works  
 $\hookrightarrow$  Eval. on STOI



\* Pay attn to SNR, it needs to be normalized  
 (b/c bubble is too strong)

## PART 1 STEPS:

⑥ Given  $S_i$

⑦ Generate noisy signals ( $S_2$  thru  $S_5$ )

↳ you need to make  $S_i$  evidently the intended speaker. This can be done by increasing the "volume" of  $S_i$ , and decr. the volume of the other clean speech signal which is not the intended signal

Eg.

$S_1$ : clean speech 1 (high SNR)

$S_3$ : babble speech

$S_5$ : noise signal 2

$S_2$ : clean speech 2 (low SNR)

$S_4$ : noise signal 1

↳ accomplish by cs 2 + noise-loud

⑧ Generate  $\underline{x}[t]$  array

$$\begin{bmatrix} x_1[t] \\ x_2[t] \\ x_3[t] \\ x_4[t] \end{bmatrix}$$

I think angles you can self choose

↳ probably just a good idea to sketch out setup so that you're always consistent

## Part 2 - Estimate the ATF

Methods we have learned:

→ EVD

→ GEVD

→ Cramer Rao Bound (likely out of scope but can consider as a challenge)

→ Inter Frequency Correlation (also)

\* Due to the various setups in part 1 → just do GEVD

\* This replaces the h files

↳ first you test beamformer with h then do GEVD which simulates real world

↳ safe to assume here

Estimate  $R_s$  under spatially colored data

Using EVD

① Pre-whiten data:  $\underline{\underline{R}}_n^{1/2} : \underline{\underline{\tilde{x}}} = \underline{\underline{R}}_n^{1/2} \underline{\underline{x}}$

②  $\underline{\underline{R}}_x = \underline{\underline{U}} (\underline{\Lambda} + \underline{\underline{I}}_m) \underline{\underline{U}}^H$ , truncate  $M-r$  smallest eig. to get  $\underline{\underline{U}}, \underline{\Lambda}, \underline{\underline{U}}^H = \underline{\underline{\tilde{R}}}_s$

③ Deconviten result:  $\underline{\underline{\tilde{R}}}_s = \underline{\underline{R}}_n^{1/2} \underline{\underline{\tilde{R}}}_s \underline{\underline{R}}_n^{1/2}$

why?

\* Note: using  $\underline{\underline{R}}_n^{1/2}$  to deconviten may result in a loss in accuracy

Therefore we do another method: generalized EVD

Gen Sig.  
Use h files  
Beamformers  
Evaluate

Gen Sig.  
GEVD for ATF  
Beamformers  
Evaluate

\* Here you need the detector to get  $R_n$  and  $R_s$

\* Should be easy for clean 1

## GENERALIZED EVD FOR A.T.F. ESTIMATION



- hermitian matrices  $\underline{A}$  and  $\underline{B}$
- diagonalize both with  $\underline{U} \rightarrow$  columns of  $\underline{U}$  are the generalized eigenvectors of  $\underline{A}$  and  $\underline{B}$
- the generalized eigenvalues and eigenvectors of  $\underline{A}, \underline{B}$  are the ordinary eigenvalues and eigenvectors of  $\underline{B}^{-1}\underline{A}$

### Derivation

$$\begin{aligned} \underline{U}^H \underline{A} \underline{U} &= \Lambda_A \\ \underline{U}^H \underline{B} \underline{U} &= \Lambda_B \end{aligned} \quad \left. \begin{array}{l} \underline{A} \underline{U} = \underline{B} \underline{U} \Lambda \\ \Lambda = \Lambda_B^{-1} \Lambda_A \end{array} \right\} \rightarrow \text{the generalized eigenvalues in } \Lambda \text{ are } \lambda_i = \frac{a_i}{b_i}$$

$$\underline{R}_x = \underline{U}^{-H} (\Lambda + I_M) \underline{U}^{-1} \quad \rightarrow \text{then assume rank}(R_s) = r < M$$

$\rightarrow$  then partition,

$$\underline{R}_x = (Q_1, Q_2) \begin{pmatrix} \Lambda_r + I_r & 0 \\ 0 & I_{M-r} \end{pmatrix} \begin{pmatrix} Q_1^H \\ Q_2^H \end{pmatrix} \quad \begin{array}{l} Q_1 \in \mathbb{C}^{M \times r} \\ Q_2 \in \mathbb{C}^{M \times (M-r)} \end{array}$$

$\rightarrow$  Estimate  $R_s$  from GEVD of  $R_x$  :  $\hat{R}_s = Q_1 (\Lambda_r + I_r) Q_1^H$

$\rightarrow$  Simplify to (when one source) :  $\hat{R}_s = Q_1 \Lambda_r Q_1^H$

If  $\text{rank } R_s = 1$ , then ATF can be obtained by selecting  $q_1$ , the principle generalized eigenvector between  $R_s(n, e)$  and  $R_s(k, e)$

## Part 3 - Creation of the Beamformers

Objective: multi-mic for far end noise reduction.

### Beamformer Types We Have Learned In Class Which Are Applicable:

- delay & sum
- MVDR
- MCWF
- LCMV



These 3 only take spatial sep. into account  
 ↳ need to do MMSE noise reduction after  
 ↳ do w/ single channel WF  
 ↳ but then MVDR + S.C.W.F. is  
 ↳ equal to the lin. constr. minimum power beamf. already the MWF?? no?

equally weights spatial and noise filtering

↳ variation is the SDW-MWF  $\rightarrow$  with param  $M$

$M=1$  : SDW-MWF becomes MWF

$M \rightarrow 0$  : SDW-MWF becomes MVDR

\* Easiest, just do  
 Single Ch WF as  
 Post-processing  
 ↳ for noise reduction

if satisfies constr. (see notes) it becomes the max. SNR beamformer

↳ can direct almost a perfect null towards interference source if spatial cov. matrix is rank 1

• MVDR and LCMV are considered distortionless beamformers

→ Technically as per the instructions we only have to implement a single beamforming system.

→ However, we should probably implement the 4 learned in class and intelligently highlight their strengths / weaknesses

## DELAY & SUM BEAMFORMER

pros / comments  
cons

### Characteristics

- "preserves the target"
- no explicit knowledge of noise → good esp. if noise changes
- reduces noise var  $\frac{1}{n} = \frac{1}{2^P} \Rightarrow -P \log_2(2) \approx -3P \text{ dB}$

## MVDR BEAMFORMER

Minimum Variance Distortionless Response  
aka Capon beamformer

- also exploits target location info
- additionally exploits position of noise sources

- no change of mag or phase in dir of source
- min. variance of beamformer output in all other directions

cost function:  $J(\underline{w}(k, \ell)) = \underline{w}^H(k, \ell) \underline{\underline{R}_x}(k, \ell) \underline{w}(k, \ell)$

↳ minimizes output power s.t. response in desired direction is 1 (this is called the distortionless constraint)

Using MVDR Requires:

①  $\underline{\underline{R}_s}(k, \ell)$  is rank 1

② target & noise are uncorrelated

③ target & noise are additive

MVDR Spatially Uncorr. Noise: MVDR becomes the delay & sum beamformer

## OPTIMAL LINEAR MULTI CHANNEL WIENER

↳ minimize MSE between the estimated signal and the true signal

signal model:  $\underline{x}(k, \ell) = s(k, \ell) \underline{a}(k, \ell) + \underline{n}(k, \ell)$

$$\min \mathbb{E} \left[ \| s(k, \ell) - \underline{w}^H(k, \ell) \underline{x}(k, \ell) \|_2^2 \right]$$

$$(\rightarrow \underline{w}(k, \ell) = \underline{R}_x^{-1}(\ell) \underline{\sigma}_{s,k}^2 \underline{a}(k, \ell)$$

$$= \frac{\underline{\sigma}_s^2(k, \ell)}{\underline{\sigma}_s^2(k, \ell) + (\underline{a}^H(k, \ell) \underline{R}_n^{-1}(k, \ell) \underline{a}(k, \ell))^{-1}} \cdot \frac{\underline{R}_n^{-1}(k, \ell) \underline{a}(k, \ell)}{\underline{a}^H(k, \ell) \underline{R}_n^{-1}(k, \ell) \underline{a}(k, \ell)}$$

single ch w.F.

MVDR

↓  
post process MVDR filter to  
further reduce noise variance

↓  
provides spatial filtering to  
suppress noise & interference

↓  
noise var set to remaining noise PSD after beamforming

## LGMV BEAMFORMER

### Linearly Constrained Minimum Variance

- generalizes the MVDR by allowing multiple linear constraints
- generalize with  $d$  linear constraints

$$\min_{\underline{w}(k, \ell)} \underline{w}^H(k, \ell) \underline{R}_x(k, \ell) \underline{w}(k, \ell)$$

$$\text{s.t. } \underline{w}^H(k, \ell) \underline{\Delta}(k, \ell) = \underline{f}^H(k, \ell) \quad \underline{\Delta} \in \mathbb{C}^{M \times d}$$

$$\text{when } d < M, \quad \underline{w}(k, \ell) = \underline{R}_x^{-1}(k, \ell) \underline{\Delta}(k, \ell) (\underline{\Delta}^H(k, \ell) \underline{R}_x^{-1}(k, \ell) \underline{\Delta}(k, \ell))^{-1} \underline{f}(k, \ell)$$

How can you use the multiple constraints?

- steer zeros in direction of noise sources → requires knowing where noise sources coming from
- maintain sig in certain dir.

\* More constr = less deg of freedom to control noise reduction

This specifically reviews pros and cons of the beamformers and discusses what functions to call in matlab

A lot of this is directly copied from the website so it needs to be reformatted before including snippets

- General
  - Conventional beamforming techniques include delay-and-sum beamforming, phase-shift beamforming, subband beamforming, and filter-and-sum beamforming. These beamformers are similar because the weights and parameters that define the beampattern are fixed and do not depend on the array input data.

- Delay & Sum explanation
  - plane waves arriving at a linear array have a time delay that is a linear function of distance along the array. Delay-and-sum beamforming compensates for these delays by applying a reverse delay to each sensor. If the time delay is accurately computed, the signals from each sensor add constructively. When the signal is narrowband, time delay becomes a phase shift in the frequency domain and is implemented by multiplying each sensor signal by a frequency-dependent compensatory phase shift. This algorithm is implemented in the `phased.PhaseShiftBeamformer`.

- Delay & Sum cons
  - Finding the compensating delay at each sensor requires accurate knowledge of the sensor locations and signal direction.

- Delay & Sum pros
  - Another advantage is its robustness against pointing errors and signal direction errors.

- Delay & Sum cons
  - A disadvantage is its broad main lobe which decreases resolution of closely spaced sources or targets. A second disadvantage is that it has large sidelobes that allow interference sources to leak into the main beam.

- General
  - Also have Optimal and Adaptive Beamforming — data dependent beamformers

- MVDR general
  - Optimal beamformers apply weights that are determined by optimizing some quantity. The MVDR beamformer determines the beamforming weights,  $w$ , by maximizing the signal-to-noise+interference ratio of the array output

- MVDR general
  - Because of the constraint, beamformer preserves the desired signal while minimizing contributions to the array output due to noise and interference. The MVDR beamformer is implemented in `phased.MVDRBeamformer`.

- MVDR pros
  - The beamformer incorporates the noise and interference into an optimal solution. The beamformer has higher spatial resolution than a conventional beamformer. The beamformer puts nulls in the direction of any interference sources. Sidelobes are smaller and smoother.

- MVDR cons
  - The MVDR beamformer is sensitive to errors in either the array parameters or arrival direction. The MVDR beamformer is susceptible to self-nulling. In addition, trying to use MVDR as an adaptive beamformer requires a matrix inversion every time the noise and interference statistics change. When there are many array elements, the inversion can be computationally expensive.

- MVDR general — becoming MPDR be careful
  - it often turns out that the noise is not separable from the signal and it is impossible to determine  $R_n$ . In that case, you can estimate a sample covariance matrix from the data. and minimizes  $w'R_w w$  instead. Minimizing this quantity leads to the minimum power distortionless response (MPDR) beamformer. If the data vector,  $x$ , contains the signal and the estimated data covariance matrix is perfect and the steering vector of the desired signal is known exactly, the MPDR

beamformer is equivalent to the MVDR beamformer. However, MPDR degrades more severely when Rx is estimated from insufficient data or the signal arrival vector is not known precisely.

- LCMV general
  - The LCMV beamformer is a generalization of MVDR beamforming and is implemented in `phased.LCMVBeamformer` and `phased.TimeDelayLCMVBeamformer`. There are several different approaches to specifying constraints such as amplitude and derivative constraints. You can, for example, specify weights that suppress interfering signals arriving from a particular direction while passing signals from a different direction without distortion.
- LCMV pros and cons
  - The advantages and disadvantages of the MVDR beamformer also apply to the LCMV beamformer.
- Improved LCMV / MVDR by recursive implementation (save computation)
  - While MVDR and LCMV are adaptive in principle, re-computation of the weights requires the inversion of a potentially large covariance matrix when the array has many elements. The Frost and generalized sidelobe cancelers are reformulations of LCMV that convert the constrained optimization into minimizing an unconstrained form and then compute the weights recursively. This approach removes any need to invert a covariance matrix.  
See `phased.FrostBeamformer` and `phased.GSCBeamformer`.

## PART 4 - EVALUATION

Evaluation Metrics (base SP & perception metrics)

- output SNR
- MNSE
- SII : speech intel. index (weighted SNR to intelligibility)
- STI, CSII, STOI, HASPI