# Microphone Array Processing
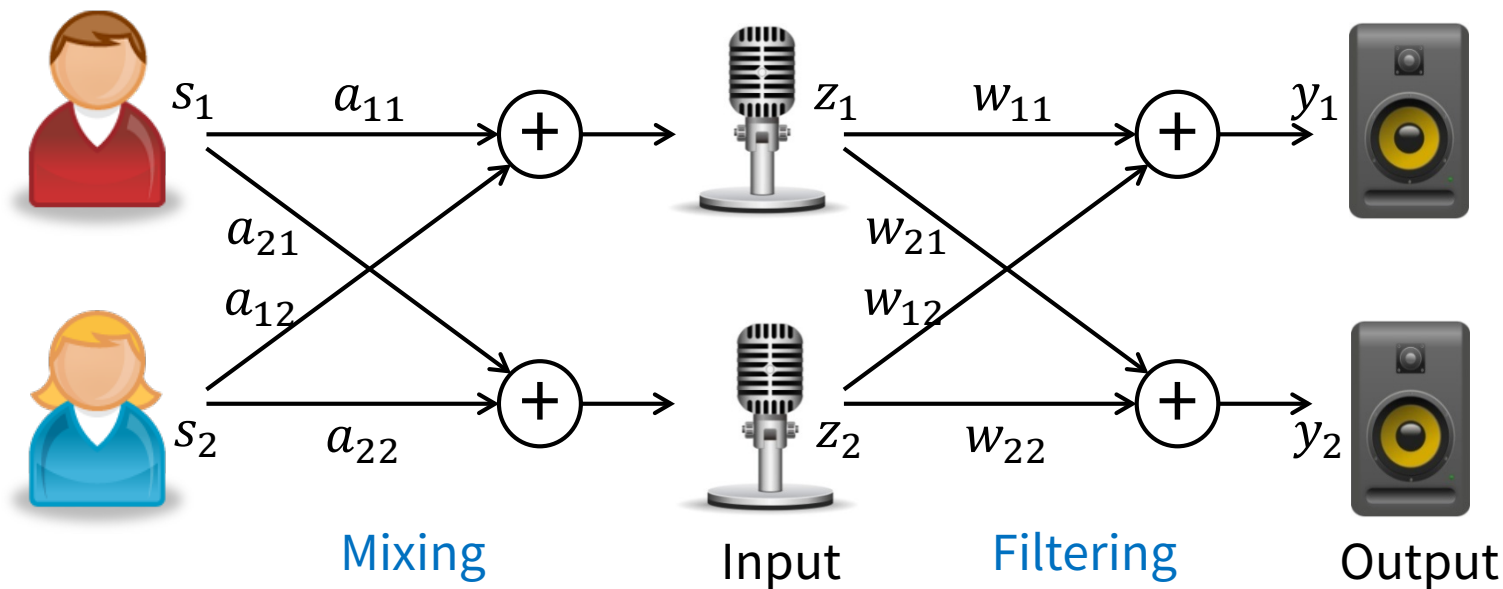
**Graduate School of Informatics**
**Kyoto University**
**Kazuyoshi Yoshii**
**yoshii@kuis.kyoto-u.ac.jp**

- A fundamental technique for various studies
  - Speech recognition & cocktail-party effect
    - It is important to selectively listen to utterances of interest even if we make conversation in a noisy environment
  - Robot audition
    - Robots should use their own ears for listening to sounds
    - Individual sound sources should be localized and separated
  - Analysis of recorded speech communication
    - Speaker identification
    - Voice activity detection for each speaker
    - Noise/reverberation reduction

- We aim at sound source separation and localization
  - Input: $z_1, z_2, \cdots, z_N$  Output: $y_1, y_2, \cdots, y_M$ $(\approx s_1, s_2, \cdots, s_M)$
    - Mixing process: sources $s_1, s_2, \cdots, s_M \rightarrow$ observations $z_1, z_2, \cdots, z_N$
    - Two settings: $A$ is given (non-blind) $\leftrightarrow$ $A$ is not given (blind)



$s_1$   $a_{11}$   $z_1$   $w_{11}$   $y_1$

$a_{21}$   $w_{21}$

$a_{12}$   $w_{12}$

$s_2$   $a_{22}$   $z_2$   $w_{22}$   $y_2$
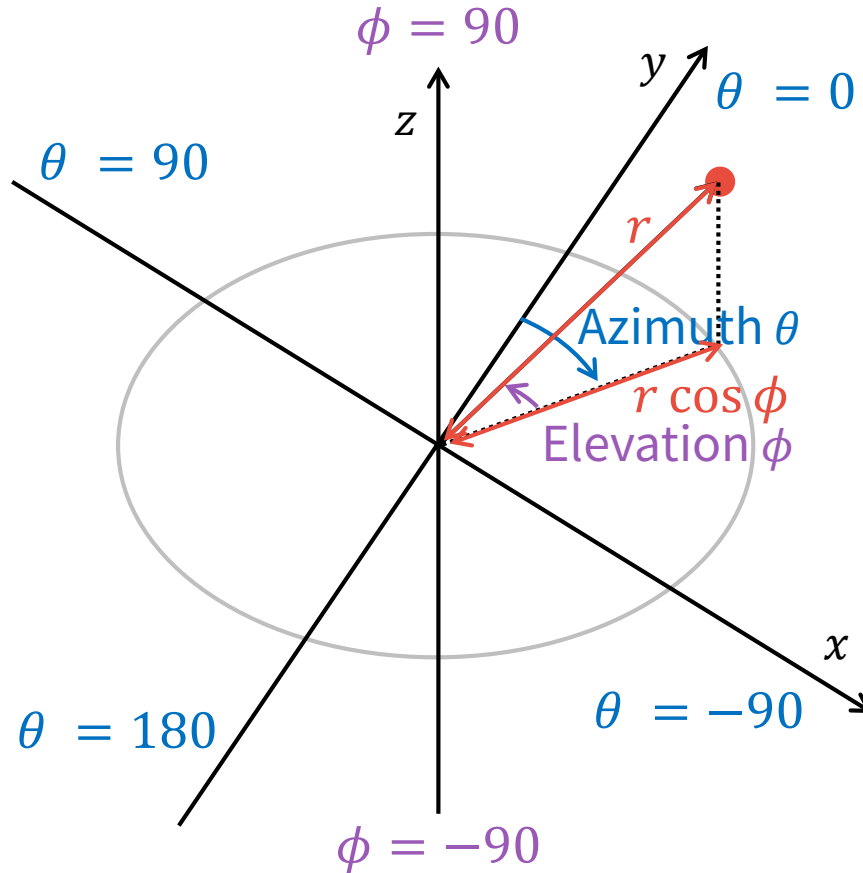
Mixing    Input    Filtering    Output

- Two major approaches to microphone array processing
  - Non-blind setting
    - Beamformer
    - MUSIC (multiple signal classification)
  - Blind setting
    - Independent component/vector analysis (ICA/IVA)
    - Multi-channel nonnegative matrix factorization (NMF)
    - Nonlinear time-frequency masking
  - Advanced topics
    - Bayesian sound source separation and localization
    - Automatic determination of number of sources

- Orthogonal coordinate ↔ Polar coordinate

$\phi = 90$

$\theta = 0$

$y$

$z$

$\theta = 90$

$r$

Azimuth $\theta$

$r \cos \phi$

Elevation $\phi$

$x$

$\theta = -90$

$\theta = 180$

$\phi = -90$

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} r \cos\phi \sin\theta \\ r \cos\phi \cos\theta \\ r \sin\phi \end{bmatrix}$$

$$\theta = \frac{\pi}{2} - \bar{\theta}$$

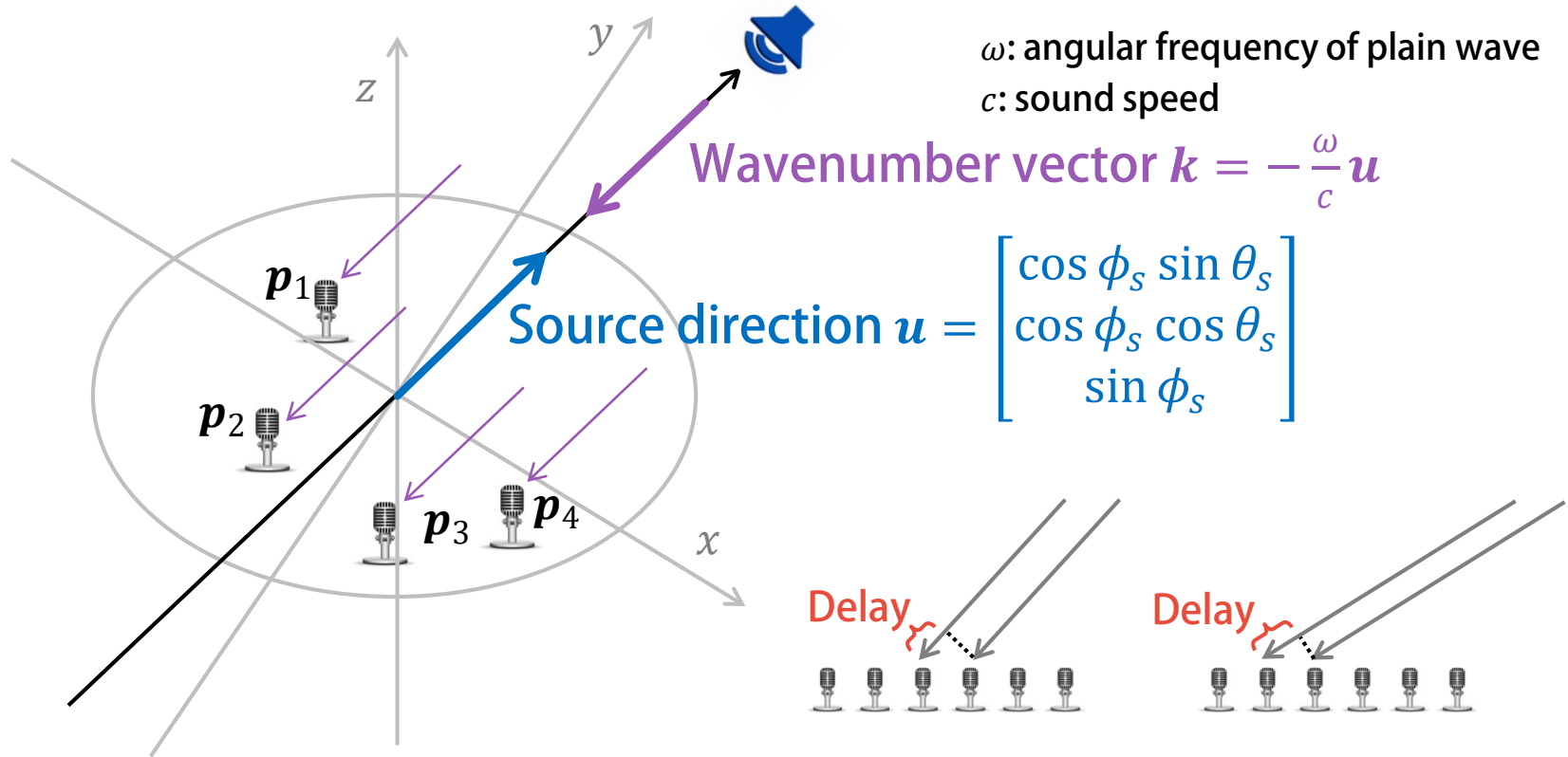Used in this lecture

$$\phi = \frac{\pi}{2} - \bar{\phi}$$

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} r \sin\bar{\phi} \cos\bar{\theta} \\ r \sin\bar{\phi} \sin\bar{\theta} \\ r \cos\bar{\phi} \end{bmatrix}$$

Commonly used in many studies

- A sound wave is observed by using $M$ microphones

  ▪ $\boldsymbol{p}_1, \boldsymbol{p}_2, \cdots, \boldsymbol{p}_M$: the positions of $M$ microhones

$z$

$y$

$x$

$\omega$: angular frequency of plain wave

$c$: sound speed

Wavenumber vector $\boldsymbol{k} = -\dfrac{\omega}{c}\boldsymbol{u}$

$\boldsymbol{p}_1$

Source direction $\boldsymbol{u} = \begin{bmatrix} \cos\phi_s \sin\theta_s \\ \cos\phi_s \cos\theta_s \\ \sin\phi_s \end{bmatrix}$

$\boldsymbol{p}_2$

$\boldsymbol{p}_3$   $\boldsymbol{p}_4$

Delay

Delay

- An observed signal is a delayed version of a source signal

  - Suppose that source signal $s(t)$ is propagated to $M$ microphones

  - Each microphone $m\ (1 \le m \le M)$ has delay time $\tau_m$

$$\mathbf{z}(t) = \begin{bmatrix} z_1(t) \\ z_2(t) \\ \vdots \\ z_M(t) \end{bmatrix} = \begin{bmatrix} s(t-\tau_1) \\ s(t-\tau_2) \\ \vdots \\ s(t-\tau_M) \end{bmatrix} \xrightarrow{\text{Fourier transform}} \mathbf{z}(\omega) = \begin{bmatrix} Z_1(\omega) \\ Z_2(\omega) \\ \vdots \\ Z_M(\omega) \end{bmatrix}$$

Observed signals

$$Z_m(\omega) \equiv \int_{-\infty}^{\infty} z_m(t) e^{-j\omega t} dt = \int_{-\infty}^{\infty} s(t-\tau_m) e^{-j\omega t} dt = e^{-j\omega\tau_m} S(\omega)$$
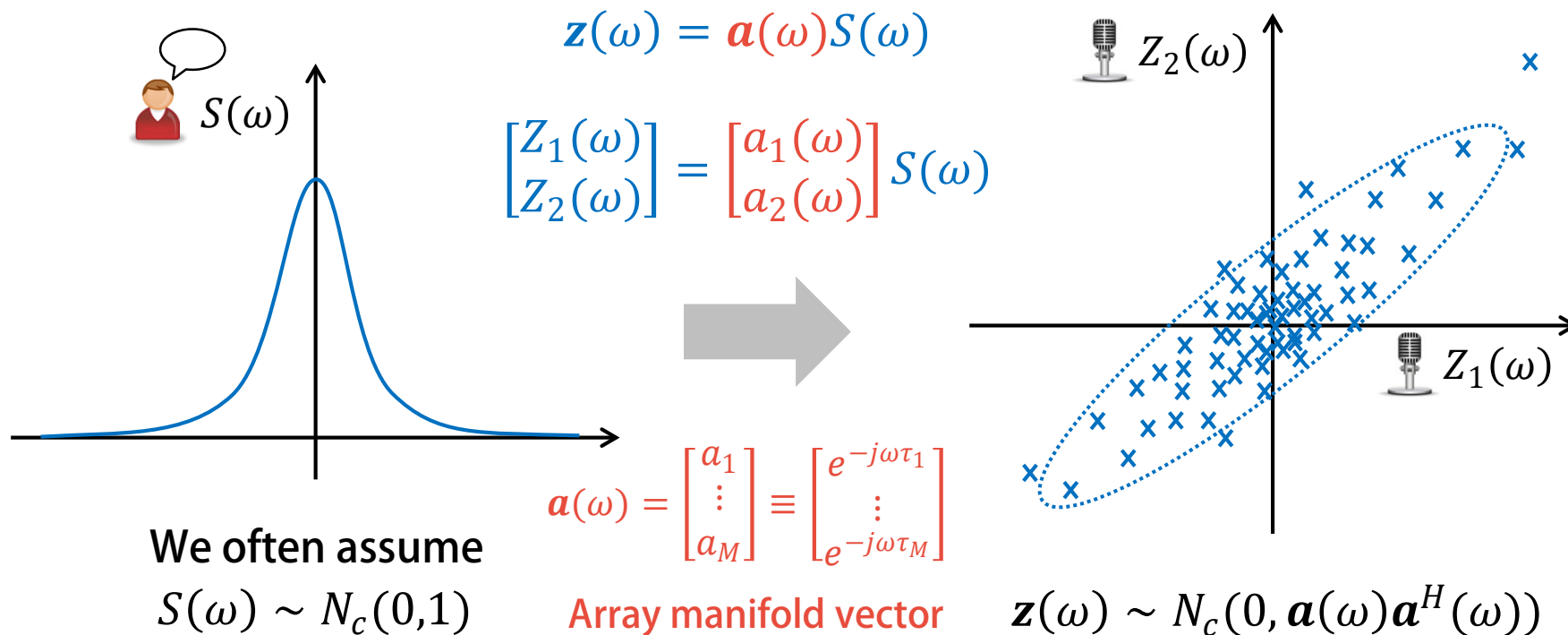
$$S(\omega) \equiv \int_{-\infty}^{\infty} s(t) e^{-j\omega t} dt \qquad \mathbf{a}(\omega) = \begin{bmatrix} a_1 \\ \vdots \\ a_M \end{bmatrix} \equiv \begin{bmatrix} e^{-j\omega\tau_1} \\ \vdots \\ e^{-j\omega\tau_M} \end{bmatrix}$$

Array manifold vector

$$\mathbf{z}(\omega) = \mathbf{a}(\omega) S(\omega)$$

Observation    Source

- **Observed signals are correlated with each other**
  - ▪ The spatial property is determined by an array manifold vector



$$\boldsymbol{z}(\omega) = \boldsymbol{a}(\omega)S(\omega)$$

$$\begin{bmatrix} Z_1(\omega) \\ Z_2(\omega) \end{bmatrix} = \begin{bmatrix} a_1(\omega) \\ a_2(\omega) \end{bmatrix} S(\omega)$$

$$\boldsymbol{a}(\omega) = \begin{bmatrix} a_1 \\ \vdots \\ a_M \end{bmatrix} \equiv \begin{bmatrix} e^{-j\omega\tau_1} \\ \vdots \\ e^{-j\omega\tau_M} \end{bmatrix}$$

We often assume
$$S(\omega) \sim N_c(0,1)$$

Array manifold vector

$$\boldsymbol{z}(\omega) \sim N_c(0, \boldsymbol{a}(\omega)\boldsymbol{a}^H(\omega))$$

- The array manifold vector $\boldsymbol{a}(\omega)$ can be calculated from microphone positions $\boldsymbol{p}_1, \boldsymbol{p}_2, \cdots, \boldsymbol{p}_M$

(azimuth, elevation): $(\theta_s, \phi_s)$     Source direction: $\boldsymbol{u} = \begin{bmatrix} \cos\phi_s \sin\theta_s \\ \cos\phi_s \cos\theta_s \\ \sin\phi_s \end{bmatrix}$

Wave equation: $\dfrac{\partial^2 s}{\partial x^2} + \dfrac{\partial^2 s}{\partial y^2} + \dfrac{\partial^2 s}{\partial z^2} = \dfrac{1}{c^2}\dfrac{\partial^2 s}{\partial t^2}$     $s$: sound pressure
$c$: sound speed

Plain wave with angular frequency $\omega$ that solves the equation:

$$s(\boldsymbol{p}, t) = A\exp(j(\omega t - \boldsymbol{k}^T \boldsymbol{p})) = A\exp(j\omega t)\exp(-j\boldsymbol{k}^T \boldsymbol{p})$$

Source signal     Phase difference

$\boldsymbol{k}$: wavenumber vector

$\boldsymbol{p}$: observation point

$\lambda$: wavelength     $\lambda = \dfrac{2\pi c}{\omega} = \dfrac{c}{f}$

$\boldsymbol{k} \equiv -\dfrac{\omega}{c}\boldsymbol{u} = -\dfrac{2\pi}{\lambda}\boldsymbol{u}$     $|\boldsymbol{k}| \equiv \dfrac{\omega}{c}$

$$a_m(\omega) = \exp(-j\boldsymbol{k}^T \boldsymbol{p}_m) = \exp\left(j\dfrac{2\pi}{\lambda}\boldsymbol{u}^T \boldsymbol{p}_m\right)$$

$$\tau_m = \dfrac{1}{\omega}\boldsymbol{k}^T \boldsymbol{p}_m = -\dfrac{1}{c}\boldsymbol{u}^T \boldsymbol{p}_m$$

- **Microphone position**

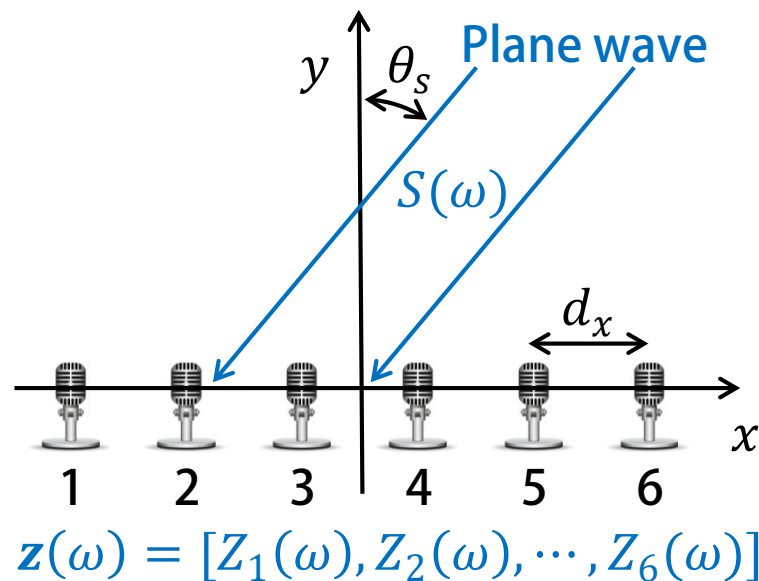$$\boldsymbol{p_m} = \left[ \left( (m-1) - \frac{M-1}{2} \right) d_x, 0, 0 \right]^T$$

- **Time delay**

$$\tau_m = -\left( (m-1) - \frac{M-1}{2} \right) \frac{d_x}{c} \sin \theta_s$$

- **Array manifold vector**

$$a_m(\omega) = \exp\left( j\left( (m-1) - \frac{M-1}{2} \right) \frac{2\pi d_x}{\lambda} \sin \theta_s \right)$$

$$\boldsymbol{a}(\omega) = e^{-\frac{j(M-1)\psi}{2}} \left[ 1, e^{j\psi}, e^{j2\psi}, \cdots, e^{j(M-1)\psi} \right]^T \quad \left( \psi = \frac{2\pi d_x}{\lambda} \sin \theta_s \right)$$

Plane wave

$y$   $\theta_s$

$S(\omega)$

$d_x$

1   2   3   4   5   6   $x$

$$\boldsymbol{z}(\omega) = [Z_1(\omega), Z_2(\omega), \cdots, Z_6(\omega)]$$

$$\boldsymbol{z}(\omega) = \boldsymbol{a}(\omega) S(\omega)$$

- **Microphone position**

$$\boldsymbol{p}_m = [r_c \sin \zeta_m, r_c \cos \zeta_m, 0]^T$$

- **Time delay**

$$\tau_m = -\frac{r_c}{c}(\sin \theta_s \sin \zeta_m + \cos \theta_s \cos \zeta_m)$$

$$= -\frac{r_c}{c}\cos(\theta_s - \zeta_m)$$

- **Array manifold vector**

$$a_m(\omega) = \exp\left(j\frac{2\pi r_c}{\lambda}\cos(\theta_s - \zeta_m)\right)$$



$$\boldsymbol{z}(\omega) = \boldsymbol{a}(\omega)S(\omega)$$

$$\boldsymbol{z}(\omega) = [Z_1(\omega), Z_2(\omega), \cdots, Z_8(\omega)]$$

Such round-shape microphone arrays are often used in practice for localizing and separating sound sources around a robot

- The array manifold vector $a(\omega)$ can be calculated from 3D microphone positions $p_1, p_2, \cdots, p_M$

(azimuth, elevation, distance): $(\theta_s, \phi_s, r)$   Source position: $\begin{bmatrix} r\cos\phi_s\sin\theta_s \\ r\cos\phi_s\cos\theta_s \\ r\sin\phi_s \end{bmatrix}$

Wave equation: $\dfrac{\partial^2(rs)}{\partial r^2} = \dfrac{1}{c^2}\dfrac{\partial^2(rs)}{\partial t^2}$   $s$: sound pressure
$c$: sound speed

**Spherical wave with angular frequency $\omega$ that satisfies the equation:**

$$s(r,t) = \frac{A}{r}\exp(j(\omega t - k_r r)) = A\exp(j\omega t)\frac{1}{r}\exp(-jk_r r)$$

Source signal    Phase/amplitude difference

$k_r$: wavenumber
$\lambda$: wavelength

$$k_r \equiv \frac{2\pi}{\lambda} = \frac{\omega}{c} \qquad \lambda = \frac{2\pi c}{\omega} = \frac{c}{f} \qquad a_m(\omega) = \frac{1}{r_m}\exp(-jk_r r_m) = \frac{1}{r_m}\exp\left(-j\omega\frac{r_m}{c}\right)$$

- ## Geometry-based estimation
  - Use the formula: $a_m(\omega) = \exp(-j\boldsymbol{k}^T\boldsymbol{p}) = \exp\left(j\dfrac{2\pi}{\lambda}\boldsymbol{u}^T\boldsymbol{p}_m\right)$

    Source direction    Microphone position

- ## Recording-based estimation
  - Use only direct sounds for measuring the impulse response
  - Transform the impulse response into the frequency domain



Reflected sound

Direct sound

Windowing

Fourier transform

Amplitude characteristics

[dB]

$\omega$

Phase characteristics

$\pi$

$0$

$-\pi$

$\omega$

- **The observed sound is a mixture of various sounds**
  - Direct sound: $z_s$
  - Reflected sound: $z_r$
  - Spatial colored noise: $v_c$
  - Spatial white noise: $v_w$

Observed sound:
$$z = z_s + z_r + v_c + v_w$$

**Single source**



**Multiple sources**

- Suppose that $N$ sound sources and $M$ microphones



$$\boldsymbol{s}(\omega) = \begin{bmatrix} S_1(\omega) \\ S_2(\omega) \\ \vdots \\ S_M(\omega) \end{bmatrix} \quad \boldsymbol{z}(\omega) = \begin{bmatrix} Z_1(\omega) \\ Z_2(\omega) \\ \vdots \\ Z_N(\omega) \end{bmatrix}$$

- **Sum of direct sounds coming from $N$ sound sources**
  - Suppose that there are $N$ sound sources
  - Each sound source is recorded by each microphone (linear system)

**Single source**

$$\boldsymbol{z}_s(\omega) = \boldsymbol{a}(\omega)S(\omega)$$

**Multiple sources**

$$\boldsymbol{z}_s(\omega) = \sum_{i=1}^{N} \boldsymbol{a}_i(\omega)S_i(\omega) = \boldsymbol{A}(\omega)\boldsymbol{s}(\omega)$$

$$\boldsymbol{z}_s(\omega) = \begin{bmatrix} Z_{s1}(\omega) \\ Z_{s2}(\omega) \\ \vdots \\ Z_{sM}(\omega) \end{bmatrix}$$

**Mixing matrix**
**(array manifold matrix)**

$$\boldsymbol{A}(\omega) = [\boldsymbol{a}_1(\omega), \cdots, \boldsymbol{a}_N(\omega)]$$

$\boldsymbol{a}_n(\omega)$: array manifold vector for each source $n$

$$\boldsymbol{s}(\omega) = \begin{bmatrix} S_1(\omega) \\ S_2(\omega) \\ \vdots \\ S_N(\omega) \end{bmatrix}$$

- Different linear systems are assumed

  - Direct sounds: $z_s = As$

  - Reflected sounds: $z_r = A_r \check{s}$ ($\check{s}$ is <u>highly</u> correlated to $s$)

    - Short direct path $\neq$ Long reflection path $\rightarrow A \neq A_r$

  - Spatial colored noise: $v_c = A_c q$ ($q$ is <u>not</u> correlated to $s$)

    - The elements of $v_c$ are inter-dependent

  - Spatial white noise: $v_w \sim N(0, \sigma^2 I)$

    - The elements of $v_w$ are independent

$$v = v_c + v_w$$

General observation model

$$z = As + v$$

$v_c$ is often assumed to be included in $v_w$

- The spatial correlation matrix $R = E[zz^H]$ represents the spatial characteristics of multi-channel signals $z$

  - For direct sounds: $R_s = E[z_s z_s^H] = AE[ss^H]A^H = A\Gamma A^H$

  - For source signals: $\Gamma = \mathrm{E}[ss^H]$

    - If sound signals are independent, $\Gamma = \mathrm{diag}(\gamma_1, \cdots, \gamma_N)$

    - $\gamma_i = E[S_i(\omega)S_i^*(\omega)]$ is the power of source $i$ at frequency $\omega$

  - For noise: $K = E[vv^H]$

    - If noise $v$ is spatially white, $K = \sigma^2 I$

    - $\sigma^2$ is the power of noise

General observation model: $z = As + v$

- The spectra of each source has a unique spatial property
  - The spatial correlation matrix $\boldsymbol{R}_s$ is determined by the mixing matrix $\boldsymbol{A}$

Observed data: $\boldsymbol{z}_s(\omega) = [Z_{s1}(\omega), \cdots, Z_{sM}(\omega)]$



$Z_{s2}(\omega)$

Covariance matrix
$$\boldsymbol{R}_s = E[\boldsymbol{z}_s \boldsymbol{z}_s^H]$$

$Z_{s1}(\omega)$

$\omega$

$m = 1$

$\omega$

$m = 2$

Source $n = 1$

If $E[\boldsymbol{s}\boldsymbol{s}^H] = \boldsymbol{I}$,
$\boldsymbol{R}_s = \boldsymbol{A}\boldsymbol{A}^H$

- Formulate a probabilistic model of $z = As + v$

  - Deterministic signal model

    $$p(v) = N(v|0, K) \xrightarrow[\begin{array}{c}\text{Linear transform}\\ z = As + v\end{array}]{}$$

    Likelihood: $p(z; \Theta) = N(z|As, K)$

    Find $\Theta = \{A, s, K\}$ that maximizes $p(z; \Theta)$

  - Random signal model

    $A$ is determined by source directions $\{\theta_1, \cdots, \theta_N\}$
    $\Gamma$ is determined by source power $\{\gamma_1, \cdots, \gamma_N\}$

    $$p(v) = N(v|0, K)$$
    $$p(s) = N(s|0, \Gamma) \xrightarrow[\begin{array}{c}\text{Linear transform}\\ z = As + v\\ \Gamma = E[ss^H] (= \text{diag}(\gamma_1, \cdots, \gamma_N))\end{array}]{}$$

    Likelihood: $p(z; \Theta) = N(z|0, A\Gamma A^H + K)$

    Find $\Theta = \{A, \Gamma, K\}$ that maximizes $p(z; \Theta)$

  Bayesian treatment of $\Theta$ is feasible by incorporating a prior $p(\Theta)$

  $$p(\Theta|z) = \frac{p(z|\Theta)p(\Theta)}{p(z)}$$

# Impulse Response

# Impulse Response

- The impulse response is a signal recorded by a microphone when an impulse is emitted from a sound source

  - The source signal is distorted by reflection, noise, and diffraction

  - Impulse response (time domain) = Transfer function (freq. domain)

    - Different rooms have different impulse responses

Impulse $\delta(t)$

$h(t)$

Impulse response $h(t)$

$$h(t) = h(t) * \delta(t)$$

- Room acoustics are often represented as a linear system
  - Source signal + Room acoustics + Additive noise → Observed signal



$h(t)$

Impulse response

$v(t)$

Additive noise

$s(t)$

Source signal

$x(t)$

Observed signal

$$z(t) = h(t) * s(t) + v(t)$$

- **Time-domain convolution ↔ Frequency-domain product**
  - ▪ $s(t)$: source signal → $z(t)$: observed signal
  - ▪ $h(t)$: impulse response that characterizes the linear system

**Continuous time domain**
$$z(t) = h(t) * s(t) = \int_{-\infty}^{\infty} h(t - \tau)s(\tau)d\tau$$

**Discrete time domain**
$$z[t] = h[t] * s[t] = \sum_{\tau=-\infty}^{\infty} h[t - \tau]s[\tau]$$

$n$: time index

**Time domain**
$$z[n] = h[n] * s[n] = \sum_{m=0}^{N-1} h[n - m]s[m]$$

**Freq. domain**
$$Z[k] = H[k] \cdot S[k]$$

$k$: frequency index

- Audio signals recorded in an arbitrary room can be simulated by using the impulse response of the room

| Source signal | | Impulse response | | Observed signal |
|---|---|---|---|---|



Reverb. time: 1.3 s

**Time-domain representation**

$$s(t) \quad * \quad h(t) \quad = \quad z(t)$$

FFT

**Frequency-domain representation**

$$S(\omega) \quad \times \quad H(\omega) \quad = \quad Z(\omega)$$

iFFT

- **The TSP can be easily emitted from a loudspeaker**
  - Frequency characteristics:
    - The impulse contains all frequencies at a moment (huge power)
    - The TSP contains a limited range of frequencies at a moment
  - The impulse is recovered by convoluting two TSPs

| Impulse | | Time stretched pulse | | Inverse TSP |
|---|---|---|---|---|
|  | $=$ |  | $*$ |  |
| $\delta(t)$ | | $\mathrm{TSP}(t)$ | | $\mathrm{iTSP}(t)$ |

- Convolute a TSP response with an inverse TSP
  - The effects of TSP and iTSP are canceled out



$$h(t) * \mathrm{TSP}(t)$$

**Average**

**Impulse response**

$$h(t) * \cancel{\mathrm{TSP}(t)} * \cancel{\mathrm{iTSP}(t)} = h(t)$$

$$h(t) * \mathrm{TSP}(t) \quad * \quad \mathrm{iTSP}(t)$$

- **Prepare devices required for recording TSPs**







Loudspeaker and earplugs:
 The TSPs are emitted multiple times

Microphone array:
 The TSP is recorded by each microphone

Recording device:
 All microphones are synchronized

- **Mark the floor with a certain interval** (5° or 10°)







Angle measurer:
  Laser is emitted while rotating

Markers:
  Stickers are on all directions

Two people:
  Angle measurer control + Marking

# Wiener Filtering

# Wiener Filter

- We aim to learn a linear filter that extracts signals of interest

- **Estimate a linear filter $\boldsymbol{w}_{\mathrm{MF}}$ that extracts $y_k$ from an observed signal $\boldsymbol{u}_k$ such that $y_k$ is close to a given desired response $d_k$**

  - Minimize error $\epsilon_k$ between desired response $d_k$ and filter output $y_k$

  Cost function

  $$J = E\big[|\epsilon_k|^2\big]$$

  $$= E\Big[\big(d_k - \boldsymbol{w}^H\boldsymbol{u}_k\big)\big(d_k - \boldsymbol{w}^H\boldsymbol{u}_k\big)^H\Big]$$

  $$= E[d_k d_k^*] - \boldsymbol{w}^H E[\boldsymbol{u}_k d_k^*] - E\big[d_k^*\boldsymbol{u}_k^H\big]\boldsymbol{w} + \boldsymbol{w}^H E\big[\boldsymbol{u}_k \boldsymbol{u}_k^H\big]\boldsymbol{w}$$

  $$\equiv \boldsymbol{\sigma}_d^2 - \boldsymbol{w}^H \boldsymbol{r}_{ud} - \boldsymbol{r}_{ud}^H \boldsymbol{w} + \boldsymbol{w}^H \boldsymbol{R}_u \boldsymbol{w}$$

  Let the partial derivative be equal to zero

  $$\frac{\partial J}{\partial \boldsymbol{w}^*} = -\boldsymbol{r}_{ud} + \boldsymbol{R}_u \boldsymbol{w} \to 0 \quad \Rightarrow \quad \boldsymbol{R}_u \boldsymbol{w}_{\mathrm{MF}} = \boldsymbol{r}_{ud} \quad \Rightarrow \quad \boldsymbol{w}_{\mathrm{MF}} = \boldsymbol{R}_u^{-1} \boldsymbol{r}_{ud}$$

  Normal equation

- Wiener filter assumes that input $u_k$ is weakly stationary

  - The mean and autocorrelation of $u_k$ are constant for any $k$

  - Auto-correlation: $r_u(n) = E[u_k u_{k-n}^*]$

  - Cross-correlation: $r_{du}(n) = E[d_k u_{k-n}^*]$

Correlation matrix (Toeplitz matrix)

$$\boldsymbol{R}_u = \begin{bmatrix} r_u(0) & & r_u(K-1) \\ r_u(-1) & \cdots & r_u(K-2) \\ \vdots & \ddots & \vdots \\ r_u(1-K) & \cdots & r_u(0) \end{bmatrix} \qquad \boldsymbol{r}_{du} = \begin{bmatrix} r_{du}(0) \\ r_{du}(-1) \\ \vdots \\ r_{du}(K-1) \end{bmatrix}$$

Time-domain Wiener filter

$$\boldsymbol{w}_{\mathrm{MF}} = \boldsymbol{R}_u^{-1} \boldsymbol{r}_{ud} \quad \Rightarrow \quad \boldsymbol{w}_{\mathrm{MF}}^H \boldsymbol{R}_u = \boldsymbol{r}_{du} \quad \Rightarrow \quad \sum_{i=1}^{K} w_i^* r_u(n-i+1) = r_{du}(n)$$
$$(n = 0, \cdots, K-1)$$

- Wiener filter assumes that input $u_k$ is weakly stationary

  - The mean and autocorrelation of $u_k$ are constant for any $k$

  - Auto-correlation: $r_u(n) = E[u_k u_{k-n}^*]$

  - Cross-correlation: $r_{du}(n) = E[d_k u_{k-n}^*]$

$$\sum_{i=1}^{K} w_i^* r_u(n-i+1) = r_{du}(n) \quad \Rightarrow \quad \sum_{i=-\infty}^{\infty} w_i^* r_u(n-i+1) = r_{du}(n)$$

We assume that $u_k$ (input) $= d_k$ (desired response) $+ v_k$ (noise)

$$W(\omega)S_u(\omega) = S_{du}(\omega)$$

Frequency-domain Wiener filer

$$S_u(\omega) = S_d(\omega) + S_v(\omega) \quad \Rightarrow \quad W(\omega) = \frac{S_d(\omega)}{S_d(\omega) + S_v(\omega)}$$

$$S_{du}(\omega) = S_d(\omega)$$

$d_k$ and $v_k$ are independent

- **Fixed filtering**
  - Estimate a Wiener filter from a finite amount of samples
    - Least square method (LS)
- **Adaptive filtering**
  - Estimate a Wiener filter in an online manner
    - Steepest descent method
    - Newton's method
    - Least mean square method (LMS)
    - Affine projection algorithm (APA)
    - Recursive least squares method (RLS)

# Beamforming

- **Extract signals of a particular direction from observations**
  - Assumption: array manifold vectors $a_1, \cdots a_M$ are known

Depend on direction $\theta, \phi$ and frequency $\omega$

Sound source 1

Sound source 2

1   2   $M$

Microphone array

Goal
Design a filter that passes only a signal coming from a particular direction

- **Design filters passing signals of a particular direction**
  - $z_m(t)$: $m^{th}$ **observed signal**   $w_m(t)$: $m^{th}$ **filter**
  - $y(t)$: **output of beamformer**

**Time domain**

$z_1(t)$ — $\boxed{w_1(t)}$

$z_2(t)$ — $\boxed{w_2(t)}$

$\vdots$

$z_M(t)$ — $\boxed{w_M(t)}$

$\to \Sigma \to y(t)$

**Frequency domain**

$Z_1(\omega)$ — $\boxed{W_1^*(\omega)}$

$Z_2(\omega)$ — $\boxed{W_2^*(\omega)}$

$\vdots$

$Z_M(\omega)$ — $\boxed{W_M^*(\omega)}$

$\to \Sigma \to Y(\omega)$

$$y(t) = \sum_{m=1}^{M} w_m(t) * z_m(t)$$

$$Y(\omega) = \sum_{m=1}^{M} W_m^*(\omega) Z_m(\omega)$$

- **Design filters passing components of a particular direction**

  - $Z_m(\omega)$: $m^{th}$ **observed signal**     $W_m(\omega)$: $m^{th}$ **filter**

  - $Y(\omega)$: **output of beamformer**

**Frequency domain**

$$Z_1(\omega) \;\bullet\!\!\!-\!\!\!- \boxed{W_1^*(\omega)}$$

$$Z_2(\omega) \;\bullet\!\!\!-\!\!\!- \boxed{W_2^*(\omega)}$$

$$\vdots$$

$$Z_M(\omega)\bullet\!\!\!-\!\!\!- \boxed{W_M^*(\omega)}$$

$$\Sigma \rightarrow Y(\omega)$$

$$Y(\omega) = \sum_{m=1}^{M} W_m^*(\omega) Z_m(\omega)$$

**Vectorial representation**

$$\boldsymbol{z}(\omega) = \begin{bmatrix} Z_1(\omega) \\ Z_2(\omega) \\ \vdots \\ Z_M(\omega) \end{bmatrix} \quad \boldsymbol{w}(\omega) = \begin{bmatrix} W_1(\omega) \\ W_2(\omega) \\ \vdots \\ W_M(\omega) \end{bmatrix}$$

$$Y(\omega) = \boldsymbol{w}^H(\omega)\boldsymbol{z}(\omega)$$

Estimated source     Observation

- ## Various methods have been proposed for filter estimation

| Method | Filter vector (steering vector) | Beam/filter type and assumptions |
|---|---|---|
| Delay-sum beamformer (DS) | $w = \dfrac{a}{a^H a}$ | Beam (fixed filter) $a$: known |
| Spatial Wiener filter (SWF) | $w = R_Z^{-1} r_{zd}$ | Beam & null (adaptive filter) $d$: known |
| Maximum likelihood (ML) | $w = \dfrac{K^{-1} a}{a^H K^{-1} a}$ | Beam & null (adaptive filter) $a, K$: known |
| Minimum variance (MV) | $w = \dfrac{R^{-1} a}{a^H R^{-1} a}$ | Beam & null (adaptive filter) $a$: known |
| Generalized sidelobe canceller (GSC) | $w = (B^H R B)^{-1} B^H R w_c$ | Beam (fixed) & null (adaptive) $a, K$: known |
| Generalized eigenvalue decomposition (GEVD) | $w = E G E^{-1}$ | Beam & null (adaptive filter) $K$: known |

- ## Take the average of delay-compensated observed signals
  - ▪ The delays are determined by a direction of beamforming

- The filter vector $\boldsymbol{w}$ has a same direction as $\boldsymbol{a}$

Observed signal

Time-domain filter

Source signal of a particular direction

$z_1(t)$ — $\boxed{w_1(t)}$

$z_2(t)$ — $\boxed{w_2(t)}$

$\vdots$

$z_M(t)$ — $\boxed{w_M(t)}$

$\Sigma \rightarrow y(t)$

$$w_m(t) = \frac{1}{M}\delta(t + \tau_m)$$

$$y(t) = \sum_{m=1}^{M} w_m(t) * z_m(t)$$

$$= \frac{1}{M}\sum_{m=1}^{M} z_m(t + \tau_m)$$

Fourier transform

DS beamformer

Steering vector $\quad \boldsymbol{w}^H(\omega) = \frac{1}{M}[e^{j\omega\tau_1}, e^{j\omega\tau_2}, \cdots, e^{j\omega\tau_M}]$

Array manifold vector $\boldsymbol{a}(\omega) = \begin{bmatrix} e^{-j\omega\tau_1} \\ \vdots \\ e^{-j\omega\tau_M} \end{bmatrix}$ $\quad \boldsymbol{w}_{\mathrm{DS}} = \frac{1}{M}\boldsymbol{a} = \frac{\boldsymbol{a}}{\boldsymbol{a}^H\boldsymbol{a}}$

Normalize

- Suppose that a beam with a "wrong" direction is used

  - Source direction $(\theta_s, \phi_s)$
  - Steering-vector direction $(\theta_T, \phi_T)$

  Different!

$$z(\omega) = a(\omega)S(\omega)$$

Observation    Source

$$Y(\omega) = w^H(\omega)z(\omega)$$

Estimated source    Observation

$$Y(\omega) = w(\omega)^H a(\omega)S(\omega) = \Psi(k, \omega)S(\omega)$$

If $(\theta_s, \phi_s) = (\theta_T, \phi_T)$, $Y(\omega) = S(\omega)$

Time delay corresponding to direction $(\theta_T, \phi_T)$

$$\Psi(k, \omega) = \frac{1}{M} \sum_{m=1}^{M} \exp(j\omega\tau_m^{(T)}) \exp(-jk^T p_m)$$

Called a beam pattern
(regarded as a function of $(\theta_s, \phi_s)$)

Wavenumber-frequency response

- Analyze a beam pattern of a straight-shape array

  - Suppose that the steering direction is $\theta_T = 0$

**Array manifold vector**

$$a_m(\omega) = \exp\left(-j\left((m-1) - \frac{M-1}{2}\right)k_x d_x\right)$$

$$\boldsymbol{k} = [k_x, k_y, k_z]$$

$$k_x = -\frac{2\pi}{\lambda}\sin\theta_s$$

**Steering vector**

$$\boldsymbol{w}^H(\omega) = \frac{1}{M}[e^{j\omega\tau_1}, e^{j\omega\tau_2}, \cdots, e^{j\omega\tau_M}] = \frac{1}{M}[1, \cdots, 1]$$

**Visible region**

$$-\frac{2\pi}{\lambda} \le k_x \le \frac{2\pi}{\lambda}$$

**Beam pattern**

$$\Psi(\boldsymbol{k}, \omega) = \boldsymbol{w}^H(\omega)\boldsymbol{a}(\omega) = \frac{1}{M}\sum_{m=1}^{M}\exp\left(-j\left((m-1) - \frac{M-1}{2}\right)k_x d_x\right) = \frac{1}{M}\frac{\sin\left(\frac{Mk_x d_x}{2}\right)}{\sin\left(\frac{k_x d_x}{2}\right)}$$

**Sum of geometric progression**

- Visualize a beam pattern: $20\log_{10}|\Psi(\boldsymbol{k}, \omega)|$



Direction $\theta$

90    0    $-90$

$-\dfrac{2\pi}{d_x}$    0    $\dfrac{2\pi}{d_x}$

Visible region

0

Grating lobe

Mainlobe

Gain
[dB]

Sidelobe

$-\dfrac{2\pi}{\lambda}$    $\dfrac{2\pi}{\lambda}$

Wavenumber $k_x$ [1/m]

If grating lobes are within the visible region, false sources will be detected in sound source localization

- **The microphone interval $d_x$ affects the beam pattern**



$d_x = \lambda$

Visible region

$d_x = \dfrac{\lambda}{2}$

$d_x = \dfrac{\lambda}{4}$

$-\dfrac{2\pi}{\lambda}$     $\dfrac{2\pi}{\lambda}$     Wavenumber $k_x$ [1/m]

- The aperture $Md_x$ affects the beam pattern



$M = 33$

Visible region

$M = 9$

$M = 5$

$-\dfrac{2\pi}{\lambda}$  $\dfrac{2\pi}{\lambda}$  Wavenumber $k_x$ [1/m]

- ## The mic interval must be set for avoiding spatial aliasing

  - Grafting lobes should be without the visible region

$$\frac{2\pi}{d_x} > \frac{4\pi}{\lambda} \quad \Rightarrow \quad d_x \leq \frac{\lambda}{2}$$

Period of beam pattern: $\frac{4\pi}{d_x}$

$-\frac{2\pi}{d_x}$ $\qquad$ $\frac{2\pi}{d_x}$



$-\frac{2\pi}{\lambda}$ $\qquad$ $\frac{2\pi}{\lambda}$ $\qquad$ Wavenumber $k_x$ [1/m]

Visible region: $\frac{4\pi}{\lambda}$

**Sampling theorem in spatial domain**

$$d_x \leq \frac{c}{2f}$$

Mic interval

**Sampling theorem in time domain**

$$T_s \leq \frac{1}{2f}$$

Sampling interval

- **Multichannel Wiener filter in the spatial domain**

## Time-domain Wiener filter

Input
$\boldsymbol{u_k}$ → [ $\boldsymbol{w}$ ] → Output $y_k$

$$y_k = \boldsymbol{w}^H \boldsymbol{u}_k \qquad \boldsymbol{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_K \end{bmatrix} \qquad \boldsymbol{u}_k = \begin{bmatrix} u_k \\ u_{k-1} \\ \vdots \\ u_{k-K+1} \end{bmatrix}$$

$$\boldsymbol{w}_{\mathrm{MF}} = \boldsymbol{R}_u^{-1} \boldsymbol{r}_{ud} \qquad \boldsymbol{R}_u = E[\boldsymbol{u}_k \boldsymbol{u}_k^H] \qquad \boldsymbol{r}_{ud} = E[\boldsymbol{u}_k d_k^*]$$

## Spatial-domain Wiener filter

Input
$\boldsymbol{z}_k$ → [ $\boldsymbol{w}$ ] → Output $y_k$

$$y_k = \boldsymbol{w}^H \boldsymbol{z}_k \qquad \boldsymbol{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_K \end{bmatrix} \qquad z_k = \begin{bmatrix} Z_1(\omega, k) \\ Z_2(\omega, K) \\ \vdots \\ Z_M(\omega, k) \end{bmatrix}$$

$$\boldsymbol{w}_{\mathrm{MF}} = \boldsymbol{R}_z^{-1} \boldsymbol{r}_{zd} \qquad \boldsymbol{R}_Z = E[\boldsymbol{z}\boldsymbol{z}^H] \qquad \boldsymbol{r}_{zd} = E[\boldsymbol{z}d^*]$$

- ## Maximize the likelihood for observed data $z$

Source signal → Observed signals

$$z(\omega) = a(\omega)S(\omega) + v(\omega) \qquad \longrightarrow \qquad z = as + v$$

Observed signals → Source signal

$$y(\omega) = w^H(\omega)z(\omega) \qquad \longrightarrow \qquad y = w^H z$$

Limitation:
we need to estimate $K$
in advance

We assume $v \sim N(v|0, K)$ $\qquad K = \mathrm{E}[vv^H]$: correlation matrix of noise

Linear transformation of $v$

$$z|s \sim N(z|as, K)$$

Log likelihood: $\log p(z|s) = -\log|\pi K| - (z - as)^H K^{-1}(z - as)$

$$\frac{\partial \log p(z|s)}{\partial s^*} = a^H K^{-1}(z - as) \qquad s_{\mathrm{ML}} = \frac{a^H K^{-1} z}{a^H K^{-1} a}(= y) \quad \longrightarrow \quad w_{\mathrm{ML}} = \frac{K^{-1} a}{a^H K^{-1} a}$$

- Combine ML beamformer with voice activity detection
  - Estimate an array manifold vector $a$ for voiced regions
  - Estimate a spatial correlation matrix $K$ for unvoiced (noise) regions

- ## Minimize the output power $|y|^2$

  - The spatial correlation matrix of noise $K$ is not required

  - Constraint: $w^H a = 1$

  - Average output power: $E[|y|^2] = E\left[\left|w^H z\right|^2\right] = w^H E[z z^H] w = w^H R w$
    $\rightarrow$ Minimize

Cost function with a Lagrange multiplier $\lambda$:

$$J = w^H R w + 2\mathrm{Re}\left(\lambda^*(a^H w - 1)\right)$$

$$= w^H R w + \lambda^*(a^H w - 1) + \lambda(w^H a - 1)$$

$$\frac{\partial J}{\partial w^*} = R w + \lambda a \rightarrow 0 \qquad w^* = -\lambda R^{-1} a$$

Noise correlation matrix $K$ is replaced with observed correlation matrix $R$

$$\lambda = -\left(a^H R^{-1} a\right)^{-1} \longrightarrow w_{\mathrm{MV}} = \frac{R^{-1} a}{a^H R^{-1} a} \iff w_{\mathrm{ML}} = \frac{K^{-1} a}{a^H K^{-1} a}$$

- ## We are interested in the power of a signal coming from a steering-vector direction $\theta_T$

Beamformer: $y(\theta_T) = w^H(\theta_T)z$

Average output power:

Spatial spectrum

$$P(\theta_T) = E\big[|y(\theta_T)|^2\big] = \boldsymbol{w}^H(\theta_T)E[\boldsymbol{zz}^H]\boldsymbol{w}(\theta_T) = \boldsymbol{w}^H(\theta_T)\boldsymbol{R}\boldsymbol{w}(\theta_T)$$

Examples:

$$\boldsymbol{w}_{\text{DS}} = \frac{\boldsymbol{a}}{\boldsymbol{a}^H\boldsymbol{a}} \qquad P_{\text{DS}}(\theta) = \frac{\boldsymbol{a}^H(\theta)}{\boldsymbol{a}^H(\theta)\boldsymbol{a}(\theta)}\boldsymbol{R}\frac{\boldsymbol{a}(\theta)}{\boldsymbol{a}^H(\theta)\boldsymbol{a}(\theta)} = \frac{\boldsymbol{a}^H(\theta)\boldsymbol{R}\boldsymbol{a}(\theta)}{|\boldsymbol{a}^H(\theta)\boldsymbol{a}(\theta)|^2}$$

$$\boldsymbol{w}_{\text{MV}} = \frac{\boldsymbol{R}^{-1}\boldsymbol{a}}{\boldsymbol{a}^H\boldsymbol{R}^{-1}\boldsymbol{a}} \qquad P_{\text{MV}}(\theta) = \frac{\boldsymbol{a}^H(\theta)\boldsymbol{R}^{-1}}{\boldsymbol{a}^H(\theta)\boldsymbol{R}^{-1}\boldsymbol{a}(\theta)}\boldsymbol{R}\frac{\boldsymbol{R}^{-1}\boldsymbol{a}(\theta)}{\boldsymbol{a}^H(\theta)\boldsymbol{R}^{-1}\boldsymbol{a}(\theta)} = \frac{1}{\boldsymbol{a}^H(\theta)\boldsymbol{R}^{-1}\boldsymbol{a}(\theta)}$$

- MV gives better spatial resolution than DS
  - MV has a similar property to MUSIC method (explained later)

# Multiple Signal Classification (MUSIC)

- Represent an observed vector $z \in \mathbb{C}^M$ in another space

| Frequency-domain method | Eigenspace method |
|---|---|
| Fourier transform $$y = Fz$$ | Karhunen-Loève transform $$y = E^H z$$ |
| Invserse Fourier transform $$z = F^H y$$ | Karhunen-Loève expansion $$z = Ey$$ |
| $F$ is a discrete transform matrix | $E = [e_1, e_2, \cdots, e_M]$ is a set of eigenvectors of $R = E[zz^H]$ |

PCA

Eigenvalue decomposition
$$Re_i = \lambda_i e_i$$

Spectral decomposition
$$R = E\Lambda E^H$$

$\Lambda = \mathrm{diag}(\lambda_1, \cdots, \lambda_M)$ is a set of the corresponding eigenvalues

The average power of the $i^{th}$ principal component
$$E[|y_i|^2] = E[e_i^H zz^H e_i] = e_i^H Re_i = \lambda_i$$

- Observed signal = Sum of direct signals

  - Suppose that $v = 0$ and $M > N$ (#microphones > #sources)

  Observation model: $z = As + v$     $\Gamma = E[ss^H]$

  $$z = \sum_{i=1}^{N} a_i s_i$$

  $$R = E[zz^H] = E[Ass^H A^H] = A\Gamma A^H$$

  $$\text{rank}(A) = \text{rank}(\Gamma) = N \longrightarrow \text{rank}(R) = N$$

Eigenvalue decomposition: $R = EME^H$

Eigenvalues: $M = \text{diag}(\mu_1, \cdots, \mu_M)$   $\mu_1 > \cdots > \mu_N > 0, \ \mu_{N+1} = \cdots = \mu_M = 0$

Eigenvectors: $E = \{e_1, \cdots, e_M\}$                $e_i^H R e_i = \mu_i$

$$e_i^H R e_i = e_i^H A\Gamma A^H e_i = (A^H e_i)^H \Gamma (A^H e_i) = \mu_i$$

Orthogonal relationships

$$A^H e_i = 0_{N \times 1} \ (N < i \leq M) \longrightarrow a_j^H e_i = 0 \ (1 \leq j \leq N, N < i \leq M)$$

- Orthogonal-complementary subspaces of $A$

  - Column space: $\mathcal{R}(A) = \text{span}(\boldsymbol{a}_1, \cdots, \boldsymbol{a}_N) \rightarrow$ Signal subspace

  - Left nullspace: $N(A^H) = \text{span}(\boldsymbol{e}_{N+1}, \cdots, \boldsymbol{e}_M) \rightarrow$ Noise subspace

$$\boldsymbol{z} = A\boldsymbol{s}$$

$\mu_i$: the power of signal $\boldsymbol{s}$ in the $i^{th}$ subspace

Eigenvalue decomposition

$$\boldsymbol{R} = E[\boldsymbol{z}\boldsymbol{z}^H] = [\boldsymbol{e}_1, \boldsymbol{e}_2, \cdots, \boldsymbol{e}_M]\text{diag}(\mu_1, \mu_2, \cdots, \mu_M)[\boldsymbol{e}_1, \boldsymbol{e}_2, \cdots, \boldsymbol{e}_M]^H$$

$$\text{span}(\boldsymbol{e}_1, \cdots, \boldsymbol{e}_N) = \text{span}(\boldsymbol{e}_{N+1}, \cdots, \boldsymbol{e}_M)^\perp$$

Orthogonal bases

Result of the previous slide

Identical

$$\boldsymbol{a}_j^H \boldsymbol{e}_i = 0 \ (1 \leq j \leq N, N < i \leq M)$$

$$\text{span}(\boldsymbol{a}_1, \cdots, \boldsymbol{a}_N) = \text{span}(\boldsymbol{e}_{N+1}, \cdots, \boldsymbol{e}_M)^\perp$$

- **Observed signal = Sum of direct signals + White noise**
  - Suppose that $v = v_w$ and $M > N$

Observation model: $z = As + v_w$   $\Gamma = E[ss^H]$   $\sigma^2 I = E[v_w v_w^H]$

$$z = \sum_{i=1}^{N} a_i s_i + v_w$$

$$R = E[zz^H] = A\Gamma A^H + \sigma^2 I$$

$$\text{rank}(A) = \text{rank}(\Gamma) = N \longrightarrow \text{rank}(R) = N$$

**Eigenvalue decomposition: $R = E\Lambda E^H$**

Eigenvalues: $\Lambda = \text{diag}(\lambda_1, \cdots, \lambda_M)$   $\Lambda = M + \sigma^2 I$   No-noise case $+ \sigma^2 I$

Eigenvectors: $E = \{e_1, \cdots, e_M\}$   $e_i^H R e_i = \lambda_i$

$$e_i^H R e_i = e_i^H (A\Gamma A^H + \sigma^2 I) e_i = (A^H e_i)^H \Gamma (A^H e_i) + \sigma^2$$

Orthogonal relationships

$$A^H e_i = 0_{N \times 1} \ (N < i \leq M) \longrightarrow a_j^H e_i = 0 \ (1 \leq j \leq N, N < i \leq M)$$

- Orthogonal-complementary subspaces of $A$

  - Column space: $\mathcal{R}(A) = \mathrm{span}(a_1, \cdots, a_N) \rightarrow$ Signal subspace

  - Left nullspace: $N(A^H) = \mathrm{span}(e_{N+1}, \cdots, e_M) \rightarrow$ Noise subspace

$$z = As + v_w$$

$\lambda_i$: the sum of the power of signal $s$ and noise $v_w$ in the $i^{th}$ subspace

**Eigenvalue decomposition**

$$R = E[zz^H] = [e_1, e_2, \cdots, e_M]\mathrm{diag}(\lambda_1, \lambda_2, \cdots, \lambda_M)[e_1, e_2, \cdots, e_M]^H$$

$$\mathrm{span}(e_1, \cdots, e_N) = \mathrm{span}(e_{N+1}, \cdots, e_M)^\perp$$

Orthogonal bases

**Result of the previous slide**

Identical

$$a_j^H e_i = 0 \ (1 \le j \le N, N < i \le M)$$

$$\mathrm{span}(a_1, \cdots, a_N) = \mathrm{span}(e_{N+1}, \cdots, e_M)^\perp$$

- Observed signal = Sum of direct signals + Colored noise

  - Suppose that $v = v_c$ and $M > N$

    Non-diagonal matrix

    Observation model: $z = As + v_c$    $\Gamma = E[ss^H]$    $K = E[v_c v_c^H]$

$$z = \sum_{i=1}^{N} a_i s_i + v_c$$

$$R = E[zz^H] = A\Gamma A^H + K$$

$$\text{rank}(A) = \text{rank}(\Gamma) = N \longrightarrow \text{rank}(R) = N$$

**Generalized eigenvalue decomp. of $R$**

$$Re_i = \lambda_i Ke_i$$

Eigenvalues: $\Lambda = \{\lambda_1, \cdots, \lambda_M\}$

Eigenvectors: $E = \{e_1, \cdots, e_M\}$

**Eigenvalue decomp. of $\Phi^{-H} R \Phi^{-1}$**

$$(\Phi^{-H} R \Phi^{-1})f_i = \lambda_i f_i$$

Eigenvalues: $\Lambda = \{\lambda_1, \cdots, \lambda_M\}$

Eigenvectors: $F = \{f_1, \cdots, f_M\}$

$$\Phi^H \Phi = K \qquad f_i = \Phi e_i \qquad \text{rank}(\Phi^{-H} R \Phi^{-1}) = N$$

- Orthogonal-complementary subspaces of $A$    $\boldsymbol{\Gamma} = E[\boldsymbol{ss}^H]$

|  | No-noise case $\boldsymbol{z} = \boldsymbol{As}$ | | White noise $\boldsymbol{z} = \boldsymbol{As} + \boldsymbol{v}_w$ | | Colored noise $\boldsymbol{z} = \boldsymbol{As} + \boldsymbol{v}_c$ | |
|---|---|---|---|---|---|---|
|  | Signal power | Noise power | Signal power | Noise power | Signal power | Noise power |
| Signal subspace $(1 \le i \le N)$ | $\mu_i$ | $0$ | $\mu_i$ | $\sigma^2$ | $\check{\mu}_i$ | $1$ |
| Noise subspace $(N < i \le M)$ | $0$ | $0$ | $0$ | $\sigma^2$ | $0$ | $1$ |
| $E[\boldsymbol{zz}^H](= \boldsymbol{R})$ | $\boldsymbol{A\Gamma A}^H$ | | $\boldsymbol{A\Gamma A}^H + \sigma^2 \boldsymbol{I}$ | | $\boldsymbol{A\Gamma A}^H + \sigma^2 \boldsymbol{K}$ | |
| $E[\boldsymbol{vv}^H]$ | $\boldsymbol{0}$ | | $\sigma^2 \boldsymbol{I}$ | | $\boldsymbol{K} = \boldsymbol{\Phi}^H \boldsymbol{\Phi}$ | |
| Eigenvalue decomposition | $\boldsymbol{R} = \boldsymbol{E}\mathsf{M}\boldsymbol{E}^H$ | | $\boldsymbol{R} = \boldsymbol{E}\Lambda\boldsymbol{E}^H$ | | $\boldsymbol{\Phi}^{-H}\boldsymbol{R}\boldsymbol{\Phi}^{-1} = \boldsymbol{F}\check{\Lambda}\boldsymbol{F}^H$ | |

- **Adaptive beamforming based on subspace analysis**
  - Separate signal and nose components into different subspaces
  - Calculate spatial spectrum $P_{\text{MUSIC}}(\theta)$

$$P_{\text{MUS}}(\theta) = \frac{\|\boldsymbol{a}(\theta)\|^2}{\sum_{i=N+1}^{M} |\boldsymbol{a}^H(\theta)\boldsymbol{e}_i|^2} = \frac{\boldsymbol{a}^H(\theta)\boldsymbol{a}(\theta)}{\boldsymbol{a}^H(\theta)\boldsymbol{E}_n\boldsymbol{E}_n^H\boldsymbol{a}(\theta)}$$

$\boldsymbol{E}_n = [\boldsymbol{e}_{N+1}, \cdots, \boldsymbol{e}_M]$: a set of eigenvectors corresponding noise subspaces

$\boldsymbol{a}(\theta)$: array manifold vector ($\theta$: _assumed_ source direction)

If $\theta$ matches a true source direction ($\boldsymbol{a}(\theta) = \boldsymbol{a}_i$),

$$\boldsymbol{a}^H(\theta)\boldsymbol{E}_n = \boldsymbol{0} \ \textit{i.e.,} \ P_{\text{MUS}}(\theta) = \infty$$

Signal and noise subspaces are orthogonal

- Orthogonal-complementary subspaces of $A$    $\mathbf{\Gamma} = E[\boldsymbol{ss}^H]$

| | SEVD-MUSIC $\boldsymbol{z} = \boldsymbol{As} + \boldsymbol{v}_w$ | | GEVD-MUSIC $\boldsymbol{z} = \boldsymbol{As} + \boldsymbol{v}_c$ | | GSVD-MUSIC $\boldsymbol{z} = \boldsymbol{As} + \boldsymbol{v}_c$ | |
|---|---|---|---|---|---|---|
| | Signal power | Noise power | Signal power | Noise power | Signal power | Noise power |
| Signal subspace $(1 \le i \le N)$ | $\mu_i$ | $\sigma^2$ | $\check{\mu}_i$ | $1$ | $\check{\mu}_i$ | $1$ |
| Noise subspace $(N < i \le M)$ | $0$ | $\sigma^2$ | $0$ | $1$ | $0$ | $1$ |
| $E[\boldsymbol{zz}^H](= \boldsymbol{R})$ | $\boldsymbol{A\Gamma A}^H + \sigma^2 \boldsymbol{I}$ | | $\boldsymbol{A\Gamma A}^H + \sigma^2 \boldsymbol{K}$ | | $\boldsymbol{A\Gamma A}^H + \sigma^2 \boldsymbol{K}$ | |
| $E[\boldsymbol{vv}^H]$ | $\sigma^2 \boldsymbol{I}$ | | $\boldsymbol{K} = \boldsymbol{\Phi}^H \boldsymbol{\Phi}$ | | $\boldsymbol{K} = \boldsymbol{U}^H \boldsymbol{V}$ | |
| Eigenvalue decomposition | $\boldsymbol{R} = \boldsymbol{E\Lambda E}^H$ | | $\boldsymbol{\Phi}^{-H} \boldsymbol{R\Phi}^{-1} = \boldsymbol{F\Lambda F}^H$ | | $\boldsymbol{K}^{-1}\boldsymbol{R} = \boldsymbol{U\Lambda V}^{-H}$ | |

- ## Compare MUSIC methods in a <u>simulated</u> environment

  - Assume an observation model: $z = z_s + v_c + v_w$

    - Direct signal: $z_s = a_1 s_1$ (direction $0°$)

    - Colored noise: $v_c = a_1^c s_1^c$ (direction $60°$)



SEVD-MUSIC

GEVD-MUSIC

**Eigenvalues $\Lambda$**



$P_{\mathrm{MUS}}(\theta)$ [dB]

SEVD-MUSIC

Signal

Colored noise

$0°$     $60°$

GEVD-MUSIC

Signal

Colored noise

$0°$     $60°$  **Direction $\theta$**

- **Compare MUSIC methods in a <u>real</u> environment**
  - Assume an observation model: $z = z_s + v_c + v_w$
    - Direct signal: $z_s = a_1 s_1$ (direction $0°$)
    - Colored noise: $v_c = a_1^c s_1^c$ (direction $60°$)



SEVD-MUSIC

GEVD-MUSIC

**Eigenvalues $\Lambda$**



$P_{\mathrm{MUS}}(\theta)$ [dB]

SEVD-MUSIC

Signal     Colored noise

$0°$     $60°$

GEVD-MUSIC

Signal     Colored noise

$0°$     $60°$ **Direction $\theta$**

- ## Take the average of spatial spectra over all frequencies
  - ### Frequency weights $\boldsymbol{\beta}$ are determined according to an application

$P(\theta, \omega_1)$

$P(\theta, \omega_2)$

$\vdots$

$P(\theta, \omega_{N_\omega})$

$P(\theta)$

Source direction

$$P(\theta) = \frac{1}{N_\omega} \sum_{l=1}^{N_w} \beta_l P(\theta, \omega_l)$$

Example: $\beta_l = \left[ \sum_{i=1}^{N} \lambda_i(\omega_i) \right]^\alpha$

$\lambda_i$: sum of eigenvalues of signal subspace

*Comparison of SSL*

*SEVD-MUSIC*
*and*
*GSVD-MUSIC*

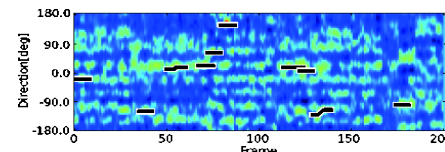Quadrocopter with 16 mics

Work well in a severely noisy environment

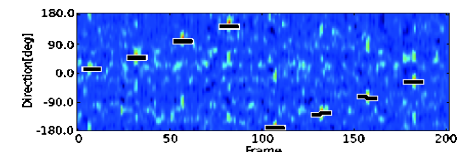Sound captured with Video Camera

SSL with iGSVD-MUSIC

Ground truth

SEVD-MUSIC

iGSVD-MUSIC
with adaptive noise estimation

# Independent Component Analysis

- BSS is a mathematically ill-defined problem
  - We cannot uniquely determine source signals
    if neither prior knowledge nor constraints are taken into account
- Focus on some properties of audio signals
  - Acoustic characteristics
    - Speech: voice timbres, accent, intonation, …
    - Musical instruments: pitches, timbres, rhythms, repetitions, …
  - Spatial characteristics
    - Source direction (angle and elevation)

Linear methods: beamformer, independent component analysis (ICA)
Nonlinear methods: time-frequency masking

- **We aim to sound source separation and localization**
  - Input: $x_1, x_2, \cdots, x_N$  Output: $y_1, y_2, \cdots, y_M$ $(\approx s_1, s_2, \cdots, s_M)$
    - Mixing process: sources $s_1, s_2, \cdots, s_M$ → observations $z_1, z_2, \cdots, z_N$
    - Two settings: $A$ is given (non-blind) ↔ $A$ is not given (blind)



Input                                          Output

# Beamforming vs. Blind Source Separation

|  | Beamforming | Blind source separation |
|---|---|---|
| Transfer functions | Required | Not necessary |
| Performance | Low | High |
| Reverberation | Can be suppressed to some extent | Included in separated signals |
| Issues |  | Permutation problem Scaling problem |

Beamformer

Blind source separation

- **Formulate a mixing process in the frequency domain**
  - $N$ sound sources are observed by $M$ microphones

  $M = N$ is assumed

  $$z = As = \sum_{i=1}^{N} a_i s_i \qquad y = Wz = WAs \qquad \Longrightarrow \qquad \text{if } W = A^{-1}, \, y \approx s$$

  Mixing system: $z = As$  Separating process: $y = Wz$



$s_1$ $a_{11}$ $a_{21}$ $a_{12}$ $s_2$ $a_{22}$ $z_1$ $w_{11}$ $w_{21}$ $w_{12}$ $z_2$ $w_{22}$ $y_1$ $y_2$

Source signals      Observed signals      Separated signals

- Linearly transform an observed space into a latent space

$$\boldsymbol{z} \longrightarrow \boxed{\boldsymbol{y} = \boldsymbol{W}\boldsymbol{z}} \longrightarrow \boldsymbol{y}$$

Observed vector        Output vector

$$\boldsymbol{W} = [\boldsymbol{w}_1, \boldsymbol{w}_2, \cdots, \boldsymbol{w}_M]^T$$

First eigenvector $\boldsymbol{e}_1$ of $\boldsymbol{R}_z$       First principal component

Estimate $\boldsymbol{w}_1$ such that the variance of $y_1 = \boldsymbol{w}_1^H \boldsymbol{z}$ is maximized

$$E\big[|y_1|^2\big] = \boldsymbol{w}_1^H E[\boldsymbol{z}\boldsymbol{z}^H]\boldsymbol{w}_1 = \boldsymbol{w}_1^H \boldsymbol{R}_Z \boldsymbol{w}_1 \qquad \|\boldsymbol{w}_1\| = 1$$

Cost function: $J = \boldsymbol{w}_1^H \boldsymbol{R}_Z \boldsymbol{w}_1 + \lambda_1 (1 - \boldsymbol{w}_1^H \boldsymbol{w}_1)$

$$\frac{\partial J}{\partial \boldsymbol{w}_1^*} = \boldsymbol{R}_z \boldsymbol{w}_1 - \lambda_1 \boldsymbol{w}_1 \to 0 \qquad E\big[|y_1|^2\big] = \boldsymbol{w}_1^H \boldsymbol{R}_Z \boldsymbol{w}_1 = \lambda \boldsymbol{w}_1^H \boldsymbol{w}_1 = \lambda_1$$

$\lambda_1$ is the maximum eigenvalue & $\boldsymbol{w}_1$ is the corresponding eigenvector

- ## The dimensions of a latent space should be orthogonal

Second eigenvector $e_2$ of $R_z$     Second principal component

Estimate $w_2$ such that the variance of $y_2 = w_2^H z$ is maximized

$$E[|y_2|^2] = w_2^H E[zz^H] w_2 = w_2^H R_Z w_2 \qquad \|w_2\| = 1 \ \& \ w_1^H w_2 = 0$$

Third eigenvector $e_3$ of $R_z$     Third principal component

Estimate $w_3$ such that the variance of $y_3 = w_3^H z$ is maximized

$$E[|y_3|^2] = w_3^H E[zz^H] w_3 = w_3^H R_Z w_3 \qquad \|w_3\| = 1 \ \& \ w_1 \perp w_2 \perp w_3$$

Eigenvalue decomposition $\longrightarrow$ Eigenvectors: $E = [e_1, \cdots, e_M]$

$$R_z = E[zz^H]$$ Eigenvalues: $\Lambda = [\lambda_1, \cdots, \lambda_M]$

PCA: $y = E^H z$     PCA with dimensionality reduction: $y = E_{1:p}^H z$

- **Perform linear transform** $\boldsymbol{y} = \boldsymbol{W}\boldsymbol{z}$ **such that** $E[\boldsymbol{y}\boldsymbol{y}^H] = 0$

  - Input space: $E[\boldsymbol{z}\boldsymbol{z}^H] = \boldsymbol{R}_z \rightarrow$ Output space: $E[\boldsymbol{y}\boldsymbol{y}^H] = \boldsymbol{I}$

  $$E[\boldsymbol{y}\boldsymbol{y}^H] = E[\boldsymbol{W}\boldsymbol{z}\boldsymbol{z}^H\boldsymbol{W}^H] = \boldsymbol{W}E[\boldsymbol{z}\boldsymbol{z}^H]\boldsymbol{W}^H = \boldsymbol{W}\boldsymbol{R}_z\boldsymbol{W}^H$$

  If $\boldsymbol{W} = \boldsymbol{\Lambda}^{-\frac{1}{2}}\boldsymbol{E}^H$, $E[\boldsymbol{y}\boldsymbol{y}^H] = \boldsymbol{\Lambda}^{-\frac{1}{2}}\boldsymbol{E}^H \boldsymbol{R}_z \boldsymbol{E}\boldsymbol{\Lambda}^{-\frac{1}{2}} = \boldsymbol{\Lambda}^{-\frac{1}{2}}\boldsymbol{\Lambda}\boldsymbol{\Lambda}^{-\frac{1}{2}} = \boldsymbol{I}$

  Scaling | Transform    Eigenvalue decomposition: $\boldsymbol{R}_z = \boldsymbol{E}\boldsymbol{\Lambda}\boldsymbol{E}^H$



Observed space

Latent space discovered by PCA

$\boldsymbol{y} = \boldsymbol{W}\boldsymbol{z}$

$\boldsymbol{y} = \boldsymbol{W}_{\text{PCA}}\boldsymbol{z}$

$\boldsymbol{z}$

$\boldsymbol{e}_2$

$\boldsymbol{e}_1$

Orthogonal

Variances along axes are normalized

- **PCA achieves second-order decorrelation**
  - ▪ The dimensions of a latent space are diagonal

**Sufficient condition**

- **ICA achieves higher-order decorrelation**
  - ▪ The dimensions of a latent space are independent
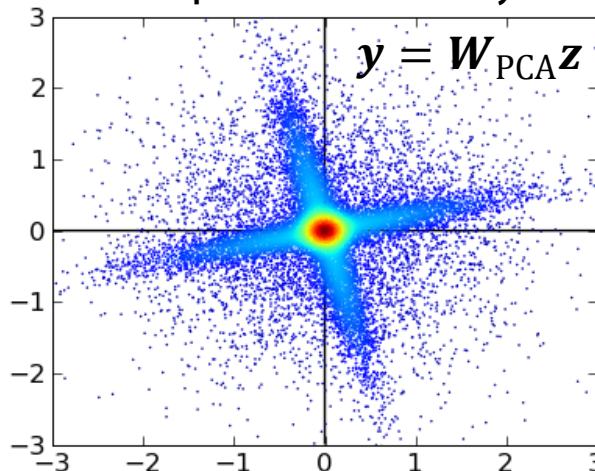
$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \boldsymbol{w}_1 z_1 + \boldsymbol{w}_2 z_2$$
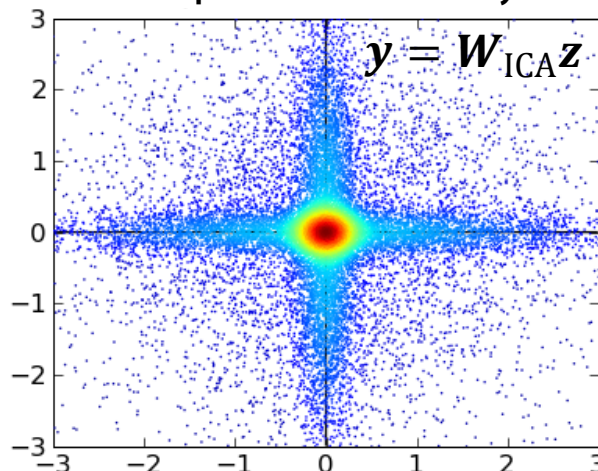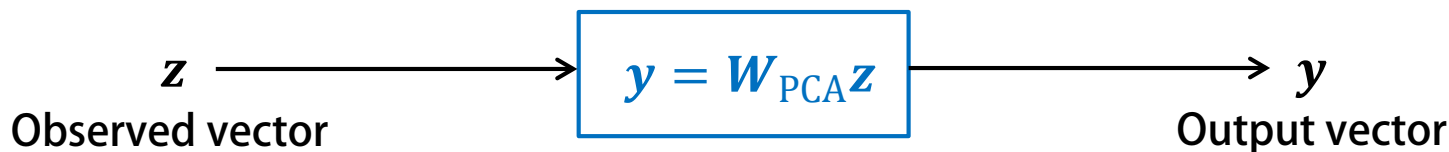


Observed space — $\boldsymbol{z}$

ICA axes

PCA axes

Latent space discovered by PCA — $\boldsymbol{y} = \boldsymbol{W}_{\mathrm{PCA}} \boldsymbol{z}$

Latent space discovered by ICA — $\boldsymbol{y} = \boldsymbol{W}_{\mathrm{ICA}} \boldsymbol{z}$

- PCA can be used as preprocessing of ICA

  ▪ ICA filter $W_{\mathrm{ICA}}$ becomes unitary after performing PCA

$$z \longrightarrow \boxed{y = W_{\mathrm{PCA}} z} \longrightarrow y$$

Observed vector · · · · · · · · · · · · · · · · · · · · Output vector

The requirement of PCA: $E[yy^H] = W_{\mathrm{PCA}} E[zz^H] W_{\mathrm{PCA}}^H = I$

If we multiply any unitary matrix $U^H$ $(U^H U = I, W_{\mathrm{PCA}} \leftarrow U^H W_{\mathrm{PCA}})$

$$y = U^H W_{\mathrm{PCA}} z \longrightarrow E[yy^H] = U^H W_{\mathrm{PCA}} E[zz^H] W_{\mathrm{PCA}}^H U = U^H U = I$$



Latent space discovered by PCA

$U^H$

$y' = W_{\mathrm{PCA}} z$

$y = U^H y'$

Latent space discovered by ICA

- **Make the dimensions of a latent spaces independent**
  - Minimize the KL divergence between $p(\boldsymbol{y})$ and $\prod_{i=1}^{N} p(y_i)$
    - If the dimensions of $\boldsymbol{y}$ are independent, $p(\boldsymbol{y}) = \prod_{i=1}^{N} p(y_i)$
    - We aim to make $p(\boldsymbol{y})$ as close to $\prod_{i=1}^{N} p(y_i)$ as possible

$$D_{KL}\left(p(\boldsymbol{y}) \middle\| \prod_{i=1}^{N} p(y_i)\right) = \int p(\boldsymbol{y}) \log \frac{p(\boldsymbol{y})}{\prod_{i=1}^{N} p(y_i)} d\boldsymbol{y}$$

$$= -\int p(\boldsymbol{y}) \log p(\boldsymbol{y}) \, d\boldsymbol{y} + \sum_{i=1}^{N} \int p(y_i) \log p(y_i) \, dy_i$$

$$= -H(\boldsymbol{y}) + \sum_{i=1}^{N} H(y_i)$$

$$\boldsymbol{y} = \boldsymbol{W}\boldsymbol{z} \longrightarrow H(\boldsymbol{y}) = H(\boldsymbol{z}) + \log|\det(\boldsymbol{W})|$$

$$D_{KL} = -H(\boldsymbol{z}) - \log|\det(\boldsymbol{W})| - \sum_{i=1}^{N} E[\log p(y_i)]$$

- **Minimize the cost function by using a gradient method**

**Cost function**

$$D_{KL} = -H(\boldsymbol{z}) - \log|\det(\boldsymbol{W})| - \sum_{i=1}^{N} E[\log p(y_i)]$$

$$\frac{\partial}{\partial w_{ij}} \sum_{i=1}^{N} E[\log p(y_i)] = E\left[\frac{\partial \log p(y_i)}{\partial y_i}\frac{\partial y_i}{\partial w_{ij}}\right] = E[-\varphi(y_i)z_j]$$

**Score function**

$$\varphi(y_i) = -\frac{\partial \log p(y_i)}{\partial y_i}$$
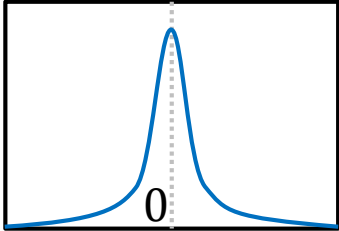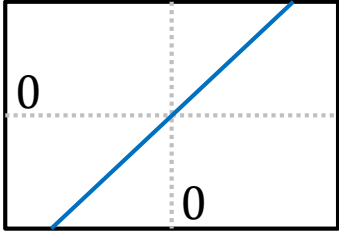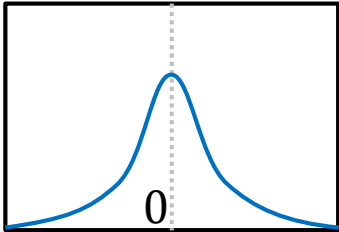$$\boldsymbol{\varphi}(\boldsymbol{y}) = [\varphi(y_1), \cdots, \varphi(y_N)]^T$$
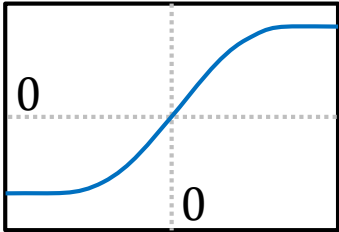
**Gradient**

$$\frac{\partial D_{KL}}{\partial \boldsymbol{W}} = -\boldsymbol{W}^{-H} + E[\boldsymbol{\varphi}(\boldsymbol{y})\boldsymbol{z}^H] = \left(\boldsymbol{I} - E[\boldsymbol{\varphi}(\boldsymbol{y})\boldsymbol{y}^H]\right)\boldsymbol{W}^{-H}$$

**Natural gradient**

$$\frac{\partial D_{KL}}{\partial \boldsymbol{W}}\boldsymbol{W}^H\boldsymbol{W} = \left(\boldsymbol{I} - E[\boldsymbol{\varphi}(\boldsymbol{y})\boldsymbol{y}^H]\right)\boldsymbol{W}$$

**Updating formula**

$$\boldsymbol{W}_{t+1} = \boldsymbol{W}_t + \eta\left(\boldsymbol{I} - E[\boldsymbol{\varphi}(\boldsymbol{y})\boldsymbol{y}^H]\right)\boldsymbol{W}_t$$

- A distribution of source signal $s \approx y$ is required

| | $p(y)$ | | $\varphi(y)$ | |
|---|---|---|---|---|
| Gaussian |  | $\dfrac{1}{\sqrt{2\pi\sigma^2}}\exp\left(-\dfrac{|y|^2}{2\sigma^2}\right)$ |  | $\dfrac{y}{\sigma^2}$ |
| Hyperbolic cosine |  | $\dfrac{1}{\pi\cosh\dfrac{y}{\sigma^2}}$ |  | $\tanh\left(\dfrac{y}{\sigma^2}\right)$ |
| Laplacian |  | $\dfrac{1}{2\sigma}\exp\left(-\dfrac{|y|}{\sigma}\right)$ |  | $\dfrac{1}{\sigma}\mathrm{sgn}(y)=\dfrac{1}{\sigma}\dfrac{y}{|y|}$ |

- ICA assumes sound sources are NOT Gaussian distributed

  ▪ The Gaussian distribution cannot be used as $p(y)$ in ICA

  > Score function: $\boldsymbol{\varphi}(\boldsymbol{y}) = [\varphi(y_1), \cdots, \varphi(y_N)]^T$
  >
  > Updating formula: $\boldsymbol{W}_{t+1} = \boldsymbol{W}_t + \eta(\boldsymbol{I} - E[\boldsymbol{\varphi}(\boldsymbol{y})\boldsymbol{y}^H])\boldsymbol{W}_t$

**Gaussian case** $\qquad \boldsymbol{\varphi}(\boldsymbol{y}) = \dfrac{\boldsymbol{y}}{\sigma^2} \qquad E[\boldsymbol{\varphi}(\boldsymbol{y})\boldsymbol{y}^H] = \dfrac{1}{\sigma^2}E[\boldsymbol{y}\boldsymbol{y}^H] = \dfrac{1}{\sigma^2}\boldsymbol{R}_y$

→ The updating formula is depend on only second-order statistics

→ ICA reduces to PCA

**Laplacian case** $\qquad \varphi(y_i) = \dfrac{1}{\sigma}\dfrac{y_i}{|y_i|}$

→ Widely used for modeling speech and music signals

- Estimate $\boldsymbol{W}$ such that $p(\boldsymbol{Z}|\boldsymbol{W})$ is maximized

ICA formulation: $\boldsymbol{y} = \boldsymbol{W}\boldsymbol{z}$

Independence of ICA outputs: $p(\boldsymbol{y}) = \prod_{i=1}^{N} p_i(y_i)$

$p(\boldsymbol{z}) = |\det(\boldsymbol{W})| p(\boldsymbol{y})$

$\boldsymbol{Z} = [\boldsymbol{z}_1, \cdots, \boldsymbol{z}_K]$    $\boldsymbol{z}_k$: observation at time $k$

$\boldsymbol{Y} = [\boldsymbol{y}_1, \cdots, \boldsymbol{y}_K]$    $\boldsymbol{y}_k$: ICA output at time $k$

$p(\boldsymbol{Z}|\boldsymbol{W}) = \prod_{k=1}^{K} |\det(\boldsymbol{W})| \prod_{i=1}^{N} p_i(y_{i,k})$

$\dfrac{\partial p(\boldsymbol{Z}|\boldsymbol{W})}{\partial \boldsymbol{W}} = \left(\boldsymbol{I} - E[\boldsymbol{\varphi}(\boldsymbol{y})\boldsymbol{y}^H]\right)\boldsymbol{W}^{-H} \to \boldsymbol{0}$

The same updating formulate is derived

- ICA variant with a constraint $\boldsymbol{W}_{\mathrm{ICA}}^{H}\boldsymbol{W}_{\mathrm{ICA}} = \boldsymbol{I}$
  - PCA is used as preprocessing
  - Fewer iterations are required for convergence

Cost function

**Restricted to be unitary**

$$D_{KL} = -H(\boldsymbol{z}) - \log|\det(\boldsymbol{W})| - \sum_{i=1}^{N} E[\log p(y_i)] = -H(\boldsymbol{z}) - 1 - \sum_{i=1}^{N} E[\log p(y_i)]$$

Minimize

**Optimization problem**

$$\min_{W} \sum_{i=1}^{N} E[G(y_i)] \quad \text{subject to} \quad \boldsymbol{W}^{H}\boldsymbol{W} = \boldsymbol{I}$$

$G(y_i) \equiv -\log p(y_i)$

We have to design $G(y_i)$ such that $G(y_i) = -\log p(y_i) \approx -\log p(s_i)$

- ## Choice of function $G(y_i)$

  **Example: generalized Laplacian:** $p(y_i) \propto \exp\left(-\frac{\sqrt{|y_i|^2+\alpha}}{\sigma}\right)$ [Sawada 2004]

  $$G(y_i) = \sqrt{|y_i|^2 + \alpha} \qquad g(y_i) = \frac{\partial G(y_i)}{\partial y_i} = \frac{y_i^*}{2\sqrt{|y_i|^2 + \alpha}}$$

  $$g'(y_i) = \frac{\partial g(y_i)}{\partial y_i^*} = \frac{1}{2\sqrt{|y_i|^2 + \alpha}}\left(1 - \frac{1}{2}\frac{|y_i|^2}{|y_i|^2 + \alpha}\right)$$
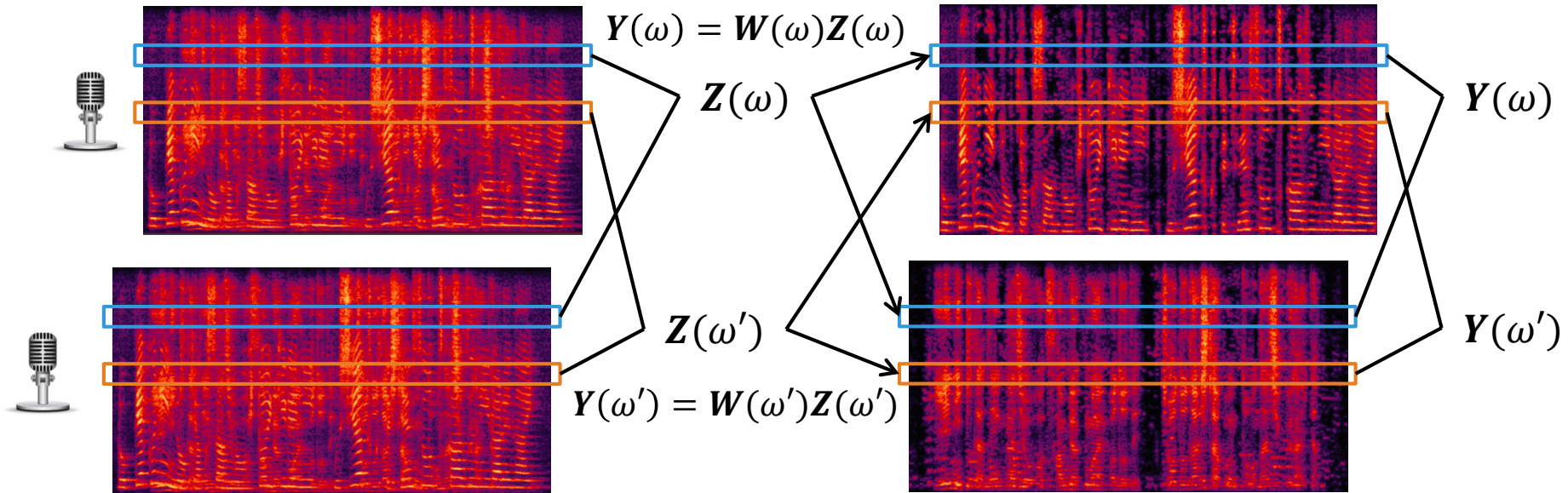
- ## Updating formula of $W$

  $$y = Wz \qquad W \equiv [w_1, w_2, \cdots, w_M]^T$$

  Update a filter: $W \leftarrow E[g(y_i)z] - E[g'(y_i)]w_m$

  Unitarize a filter: $W \leftarrow W(W^H W)^{-\frac{1}{2}}$

  Iterate until convergence

- **Permutation ambiguity**

  ▪ The dimension order of $Y$ cannot be determined uniquely

- **Amplitude ambiguity**

  ▪ The dimension amplitude of $Y$ cannot be determined uniquely



$$Y(\omega) = W(\omega)Z(\omega)$$

$$Z(\omega)$$

$$Y(\omega)$$

$$Z(\omega')$$

$$Y(\omega')$$

$$Y(\omega') = W(\omega')Z(\omega')$$

- **Solve permutation ambiguity**
  - Focus on $y$
    - Temporal power envelopes
  - Focus on $W$
    - Directional patters of $W$
    - Relative delay times from sources to microphones
    - Column vectors of $W^{-1}$
- **Solve amplitude ambiguity**
  - Recover observed signals
    - Use the invserse of $W$ for filtering each $y_i$
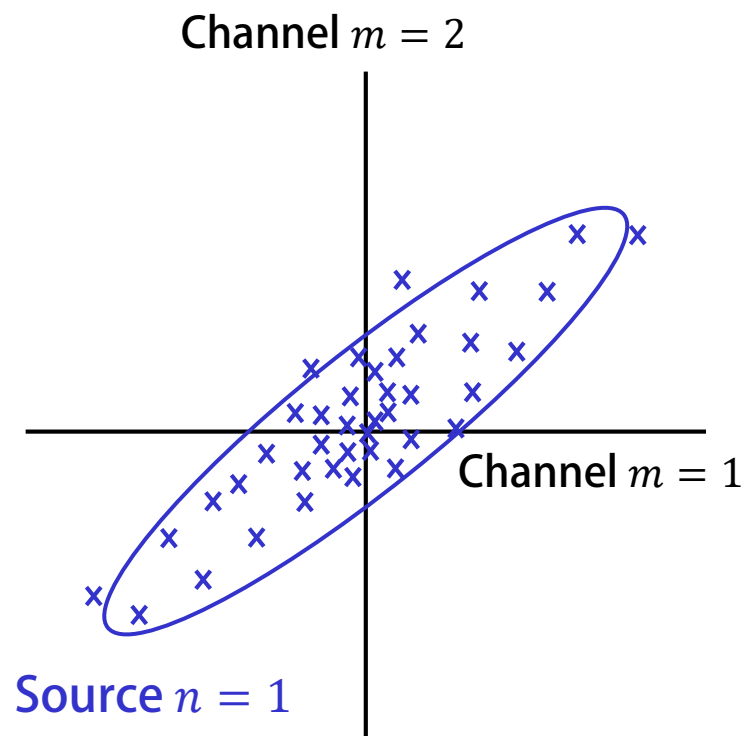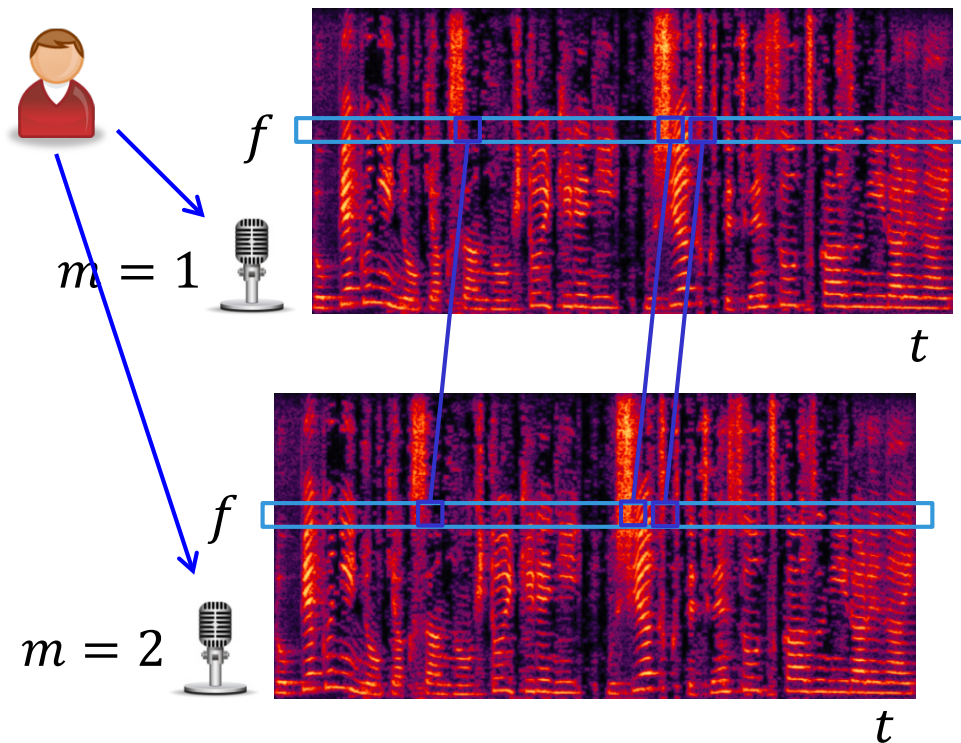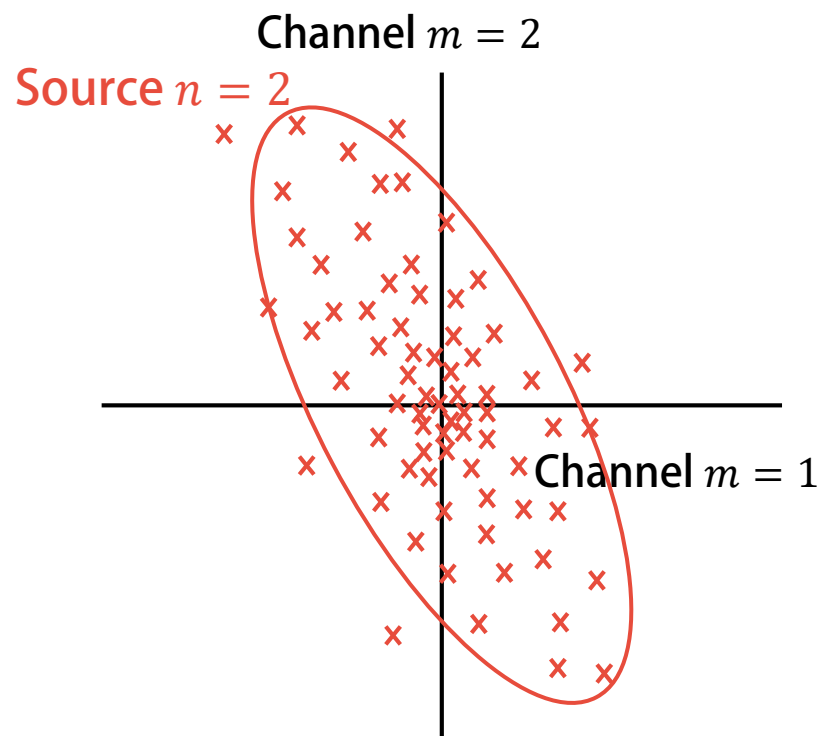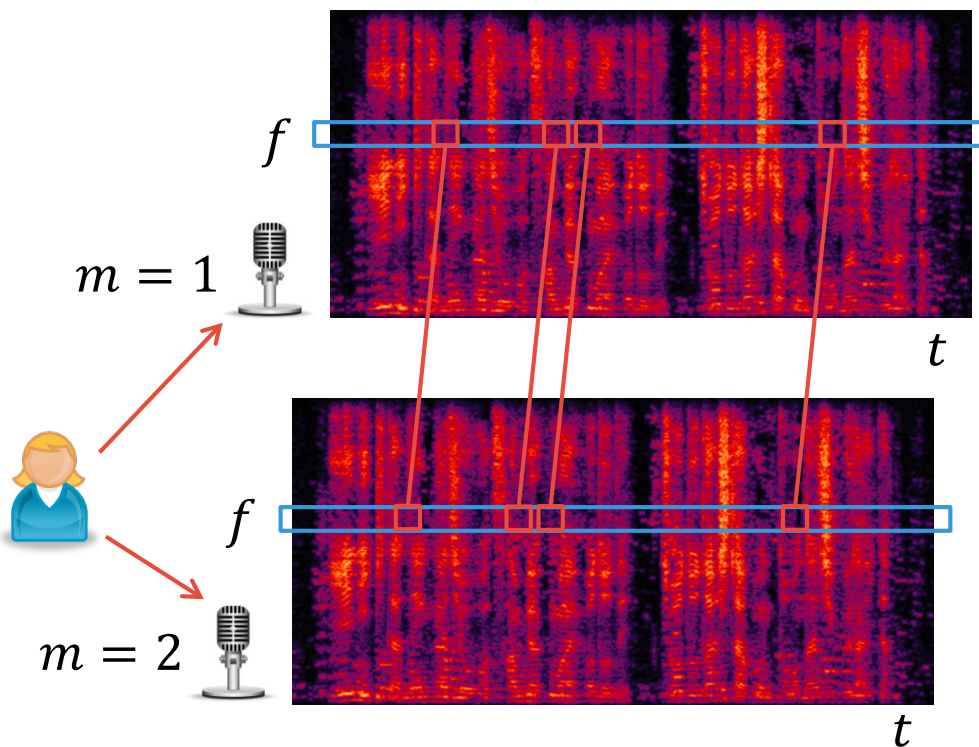      $$z_i = W^{-1}[0, \cdots, 0, y_i, 0, \cdots, 0]^T$$

# Nonlinear Time-Frequency Masking

- **The spectra of each source has a unique spatial property**
  - The spectra are assumed to be Gaussian distributed

Observed data: $\boldsymbol{x}_{tf} = [x_{tf1}, x_{tf2}, \cdots, x_{tfM}]$



$m = 1$

$f$

$t$

$m = 2$

$f$

$t$

Channel $m = 2$

Channel $m = 1$

Source $n = 1$

- ## The spectra of each source has a unique spatial property
  - ▪ The spectra are assumed to be Gaussian distributed

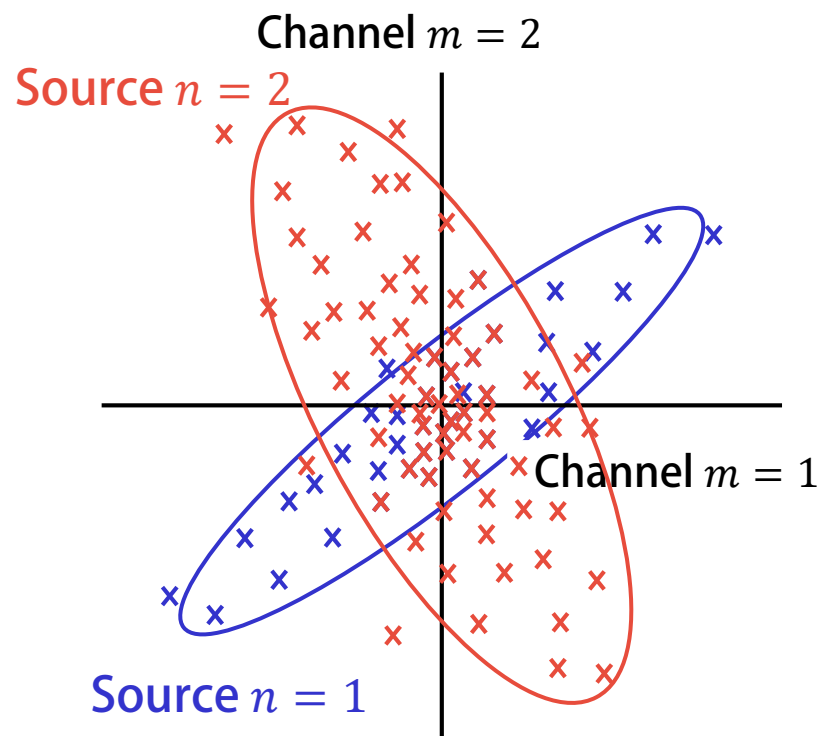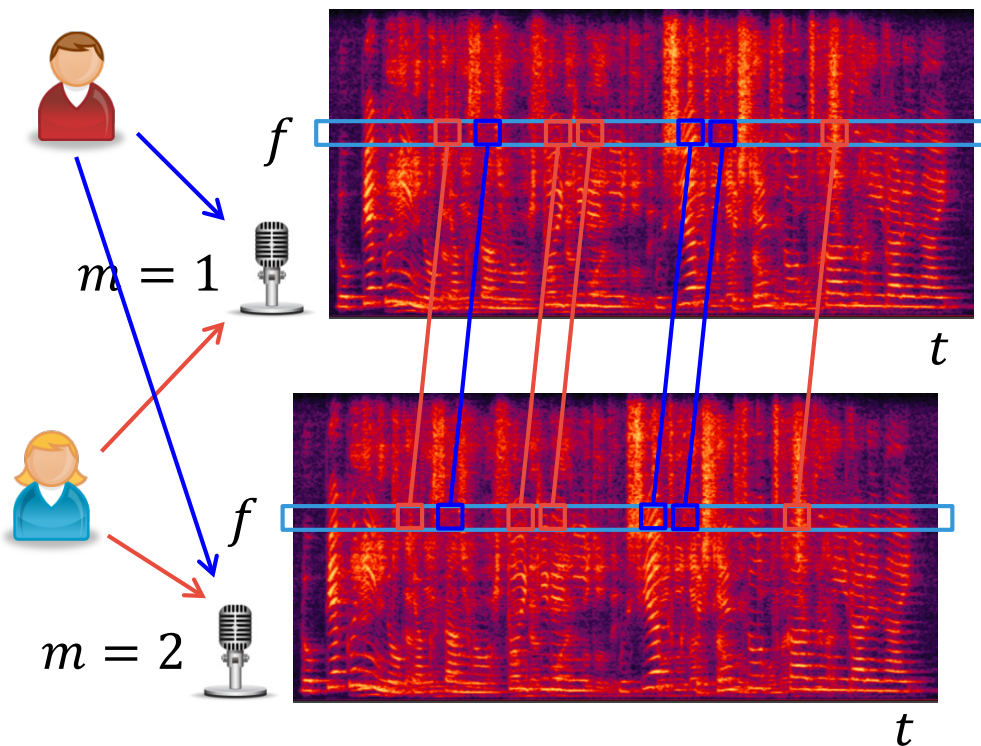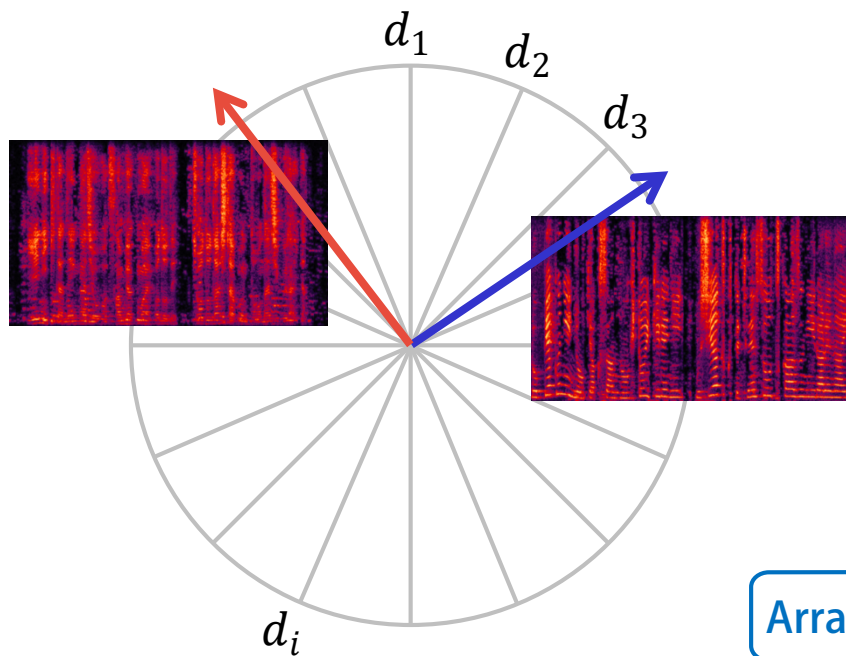Observed data: $x_{tf} = [x_{tf1}, x_{tf2}, \cdots, x_{tfM}]$

- ## The observed scatter plot is a mixture of spatial properties
  - Assume that source spectra are sparse (disjoint with each other)

Observed data: $\boldsymbol{x}_{tf} = [x_{tf1}, x_{tf2}, \cdots, x_{tfM}]$



$f$

$m = 1$

$f$

$m = 2$

$t$

$t$

Channel $m = 2$

Source $n = 2$

Channel $m = 1$

Source $n = 1$

- **Classify each frequency bin into one of sound sources**

  - $z_{tf} = k$ indicates (time $t$, frequency $f$) is classified into source $k$

  - $\boldsymbol{H}_{fd}$: spatial correlation matrix for frequency $f$ and direction $d$



**Observation model** [Duong 2010]

$$\boldsymbol{x}_{tf} \sim N_c \left( \boldsymbol{x}_{tf} \middle| \boldsymbol{0}, \left( \lambda_{tf} \boldsymbol{H}_{fd_{z_{tf}}} \right)^{-1} \right)$$

> Source direction of time $t$ and frequency $f$

**Bayesian formulation** [Otsuka 2014]

$$\boldsymbol{H}_{fd} \sim W_c \left( \left( \boldsymbol{a}_{fd} \boldsymbol{a}_{fd}^H + \epsilon \boldsymbol{I} \right)^{-1}, \nu_0 \right)$$

> Array manifold vector for frequency $f$ and direction $d$

- Automatically estimate the number of sound sources
  - Assume that infinitely many sound sources exist in theory

Channel $m = 2$

Channel $m = 1$

**Observation model** [Duong 2010]

$$\boldsymbol{x}_{tf} \sim N_c\left(\boldsymbol{x}_{tf}\middle|\boldsymbol{0}, \left(\lambda_{tf}\boldsymbol{H}_{fd_{z_{tf}}}\right)^{-1}\right)$$
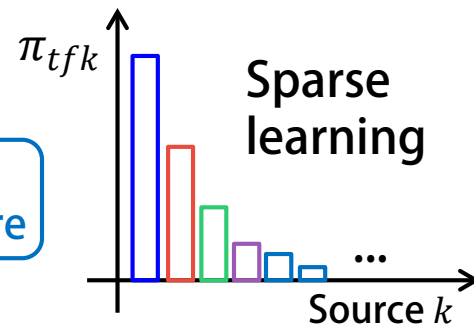
Source direction of time $t$ and frequency $f$

**Hierarchical Dirichlet process prior** $(k \to \infty)$
[Otsuka 2014]

$$\boldsymbol{\pi}_{tf} \sim \text{HDP}(\alpha, \gamma, \boldsymbol{\beta})$$

Concentration parameters

Base measure

$$z_{tf} \sim \text{Categorical}(\boldsymbol{\pi}_{tf})$$

$\pi_{tfk}$

Sparse learning

...

Source $k$

# Advantages

- **Simultaneous localization and separation**
  - Improved performance of each task
    - Integration based on a probabilistic model
  - Automatic estimation of the number of sound sources
    - Nonparametric Bayesian formulation
  - Solving permutation problems
    - All frequency bins are simultaneously analyzed
- **Various extensions feasible**
  - Simultaneous dereverberation, localization, and separation [Otsuka 2014]
  - Analyzing moving sound sources [Otsuka 2014]
  - Real-time online inference (future work)

- Questions
  - Explain delay-sum (DS) and minimum-variance (MV) beamforming methods using equations and why MV is better than DS.
  - Describe the relationships (differences) between PCA and ICA and how to estimate the parameters.
  - Report how microphone array processing is used in practice.
- How to submit
  - Submit a PDF file to "Assignment (Yoshii)" on PandA.
  - Deadline: 2018/01/30 23:59