

Digital Audio Signal Processing

Lecture-3 Noise Reduction

Marc Moonen

Dept. E.E./ESAT-STADIUS, KU Leuven

marc.moonen@esat.kuleuven.be

homes.esat.kuleuven.be/~moonen/

Overview

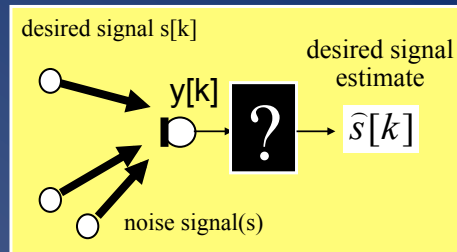
- **Spectral subtraction for single-micr. noise reduction**
 - Single-microphone noise reduction problem
 - Spectral subtraction basics (=spectral filtering)
 - Features: gain functions, implementation, musical noise,...
 - Iterative Wiener filter based on speech signal model
- **Multi-channel Wiener filter for multi-micr. noise red.**
 - Multi-microphone noise reduction problem
 - Multi-channel Wiener filter (=spectral+spatial filtering)
- **Kalman filter based noise reduction**
 - Kalman filters
 - Kalman filters for noise reduction

Single-Microphone Noise Reduction Problem

- **Microphone signal is**

$$y[k] = s[k] + n[k]$$

desired signal contribution
noise contribution



- **Goal:** Estimate $s[k]$ based on $y[k]$
- **Applications:**
Speech enhancement in conferencing, handsfree telephony, hearing aids, ...
Digital audio restoration
- **Will consider speech applications:** $s[k]$ = speech signal

Spectral Subtraction Methods: Basics

$$y[k] = s[k] + n[k]$$

- Signal chopped into 'frames' (e.g. 10..20msec), for each frame a frequency domain representation is

$$Y_i(\omega) = S_i(\omega) + N_i(\omega) \quad (i\text{-th frame})$$

- However, speech signal is an on/off signal, hence some frames have **speech + noise**, i.e.

$$Y_i(\omega) = S_i(\omega) + N_i(\omega) \quad \text{frame}_i \in \{\text{'speech + noise' frames}\}$$

some frames have **noise only**, i.e.

$$Y_i(\omega) = 0 + N_i(\omega) \quad \text{frame}_i \in \{\text{'noise - only' frames}\}$$

- A **speech detection algorithm** is needed to distinguish between these 2 types of frames (based on energy/dynamic range/statistical properties,...)

Spectral Subtraction Methods: Basics

- Definition: $\mu(\omega)$ = average amplitude of noise spectrum

$$\mu(\omega) = E\{|N_i(\omega)|\}$$

- Assumption: noise characteristics change slowly, hence estimate $\mu(\omega)$ by (long-time) averaging over (M) noise-only frames

$$\hat{\mu}(\omega) = \frac{1}{M} \sum_{M \text{ noise-only frames}} |Y_i(\omega)|$$

- Estimate clean speech spectrum $S_i(\omega)$ (for each frame), using corrupted speech spectrum $Y_i(\omega)$ (for each frame, i.e. short-time estimate) + estimated $\mu(\omega)$:

$$\hat{S}_i(\omega) = G_i(\omega) Y_i(\omega)$$

based on 'gain function'

$$G_i(\omega) = f(Y_i(\omega), \hat{\mu}(\omega))$$

Spectral Subtraction: Gain Functions

Magnitude Subtraction	$G_i(\omega) = \left[1 - \frac{\hat{\mu}(\omega)}{ Y_i(\omega) } \right]$
Spectral Subtraction	$G_i(\omega) = \sqrt{1 - \rho \frac{\hat{\mu}^2(\omega)}{ Y_i(\omega) ^2}}$
Wiener Estimation	$G_i(\omega) = 1 - \frac{\hat{\mu}^2(\omega)}{ Y_i(\omega) ^2}$
Maximum Likelihood	$G_i(\omega) = \frac{1}{2} \left[1 + \sqrt{1 - \frac{\hat{\mu}^2(\omega)}{ Y_i(\omega) ^2}} \right]$
Non-linear Estimation	$G_i(\omega) = f(\hat{\mu}(\omega), Y_i(\omega))$
Ephraim-Malah = most frequently used in practice	$G_i(\omega) = f(\text{SNR}_{\text{post}}, \text{SNR}_{\text{prio}})$ see next slide

Spectral Subtraction: Gain Functions

• Example 1: Ephraim-Malah Suppression Rule (EMSR)

$$G_i(\omega) = \frac{\sqrt{\pi}}{2} \sqrt{\left(\frac{1}{\text{SNR}_{\text{post}}} \right) \left(\frac{\text{SNR}_{\text{prio}}}{1 + \text{SNR}_{\text{prio}}} \right)} \cdot M \left[\text{SNR}_{\text{post}} \left(\frac{\text{SNR}_{\text{prio}}}{1 + \text{SNR}_{\text{prio}}} \right) \right]$$

with:

$$M[\theta] = e^{-\frac{\theta}{2}} \left[(1-\theta) I_0\left(\frac{\theta}{2}\right) + \theta I_1\left(\frac{\theta}{2}\right) \right]$$

modified Bessel functions

$$\text{SNR}_{\text{post}}(\omega) = \frac{|Y_i(\omega)|^2}{\hat{\mu}(\omega)^2}$$

$$\text{SNR}_{\text{prio}}(\omega) = (1-\alpha) \max(\text{SNR}_{\text{post}} - 1, 0) + \alpha \frac{|G_{i-1}(\omega) Y_{i-1}(\omega)|^2}{\hat{\mu}(\omega)^2}$$

- This corresponds to a **MMSE** (*) estimation of the speech spectral amplitude $|S_i(\omega)|$ based on observation $Y_i(\omega)$ (estimate equal to $\mathbf{E}\{ |S_i(\omega)| \mid Y_i(\omega) \}$) assuming Gaussian a priori distributions for $S_i(\omega)$ and $N_i(\omega)$ [Ephraim & Malah 1984].
- Similar formula for MMSE log-spectral amplitude estimation [Ephraim & Malah 1985].

(*) minimum mean squared error

Spectral Subtraction: Gain Functions

• Example 2: Magnitude Subtraction

– Signal model:

$$\begin{aligned} Y_i(\omega) &= S_i(\omega) + N_i(\omega) \\ &= |Y_i(\omega)| e^{j\theta_{Y_i}(\omega)} \end{aligned}$$

– Estimation of clean speech spectrum:

$$\begin{aligned} \hat{S}_i(\omega) &= [|Y_i(\omega)| - \hat{\mu}(\omega)] e^{j\theta_{Y_i}(\omega)} \\ &= \underbrace{\left[1 - \frac{\hat{\mu}(\omega)}{|Y_i(\omega)|} \right]}_{G_i(\omega)} Y_i(\omega) \end{aligned}$$

– PS: half-wave rectification

$$G_i(\omega) \Leftarrow \max(0, G_i(\omega))$$

Spectral Subtraction: Gain Functions

- **Example 3: Wiener Estimation**

- **Linear MMSE estimation:**

find linear filter $G_i(\omega)$ to minimize MSE

$$= E \left\{ \left| \hat{S}_i(\omega) - \overbrace{G_i(\omega) Y_i(\omega)}^{S_i(\omega)} \right|^2 \right\}$$

- Solution:

$$G_i(\omega) = \frac{E\{S_i(\omega) Y_i^*(\omega)\}}{E\{Y_i(\omega) Y_i^*(\omega)\}} = \frac{P_{sy,i}(\omega)}{P_{yy,i}(\omega)}$$

<- cross-correlation in i-th frame

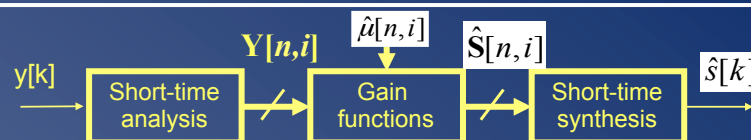
<- auto-correlation in i-th frame

Assume speech $s[k]$ and noise $n[k]$ are uncorrelated, then...

$$G_i(\omega) = \frac{P_{ss,i}(\omega)}{P_{yy,i}(\omega)} = \frac{P_{yy,i}(\omega) - P_{nn,i}(\omega)}{P_{yy,i}(\omega)} = \frac{|Y_i(\omega)|^2 - \hat{\mu}(\omega)^2}{|Y_i(\omega)|^2} = 1 - \frac{\hat{\mu}(\omega)^2}{|Y_i(\omega)|^2}$$

- PS: half-wave rectification

Spectral Subtraction: Implementation



→ Short-time Fourier Transform (=uniform DFT-modulated analysis filter bank)

$$Y[n,i] = \sum_{k=0}^{K-1} h[k] y[iD - k] e^{-j2\pi kn/N} \quad \text{= estimate for } Y(\omega_n) \text{ at time } i \text{ (i-th frame)}$$

N=number of frequency bins (channels) $n=0..N-1$

D=downsampling factor

K=frame length $h[k]$ = length-K analysis window (=prototype filter)

→ frames with 50%...66% overlap (i.e. 2-, 3-fold oversampling, $N=2D..3D$)

→ subband processing: $\hat{S}[n,i] = G[n,i] \cdot Y[n,i]$

→ synthesis bank: matched to analysis bank (see DSP-CIS)

Spectral Subtraction: Musical Noise

- Audio demo: car noise



- Artifact: **musical noise**

What?

Short-time estimates of $|Y_i(\omega)|$ fluctuate randomly in noise-only frames, resulting in random gains $G_i(\omega)$

→ *statistical analysis shows that broadband noise is transformed into signal composed of short-lived tones with randomly distributed frequencies (=musical noise)*

Spectral Subtraction: Musical Noise

Solutions?

- Magnitude averaging: replace $Y_i(\omega)$ in calculation of $G_i(\omega)$ by a local average over frames

$$\hat{S}_i(\omega) = \overbrace{G(\omega)}^{\text{average}} \overbrace{Y_i(\omega)}^{\text{instantaneous}}$$

- EMSR (p7)
- augment $G_i(\omega)$ with soft-decision VAD:

$$G_i(\omega) \rightarrow P(H_1 | Y_i(\omega)) \cdot G_i(\omega)$$

...

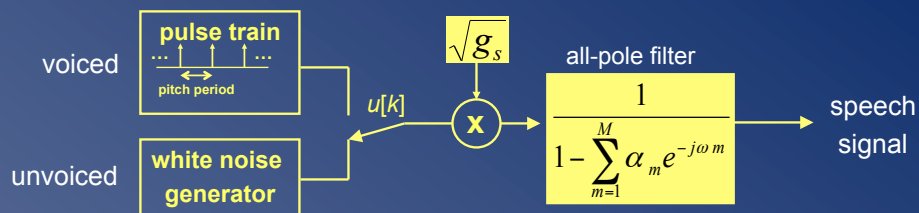
probability that speech is present, given observation

Spectral Subtraction: Iterative Wiener Filter

Example of signal model-based spectral subtraction...

- **Basic:**
Wiener filtering based spectral subtraction (p.9), with (improved) spectra estimation based on parametric models
- **Procedure:**
 1. Estimate parameters of a speech model from noisy signal $y[k]$
 2. Using estimated speech parameters, perform noise reduction (e.g. Wiener estimation, p. 9)
 3. Re-estimate parameters of speech model from the speech signal estimate
 4. Iterate 2 & 3

Spectral Subtraction: Iterative Wiener Filter



frequency domain:

$$S(\omega) = \frac{\sqrt{g_s}}{1 - \sum_{m=1}^M \alpha_m e^{-j\omega m}} U(\omega)$$

time domain:

$$s[k] = \sum_{m=1}^M \alpha_m s[k-m] + \sqrt{g_s} u[k]$$

$$\mathbf{a} = [\alpha_1 \quad \dots \quad \alpha_M]^T = \text{linear prediction parameters}$$

Spectral Subtraction: Iterative Wiener Filter

For each frame (vector) $y[m]$ (i=iteration nr.)

1. Estimate $g_{s,i}$ and $\alpha_i = [\alpha_{1,i} \ \dots \ \alpha_{M,i}]^T$
2. Construct Wiener Filter (p.9)

$$G(\omega) = \dots = \frac{P_{ss}(\omega)}{P_{ss}(\omega) + P_{nn}(\omega)}$$

with:

- $P_{nn}(\omega)$ estimated during noise-only periods

$$P_{ss}(\omega) \approx \frac{g_{s,i}^2}{\left| 1 - \sum_{m=1}^M \alpha_{m,i} e^{-j\omega m} \right|^2}$$

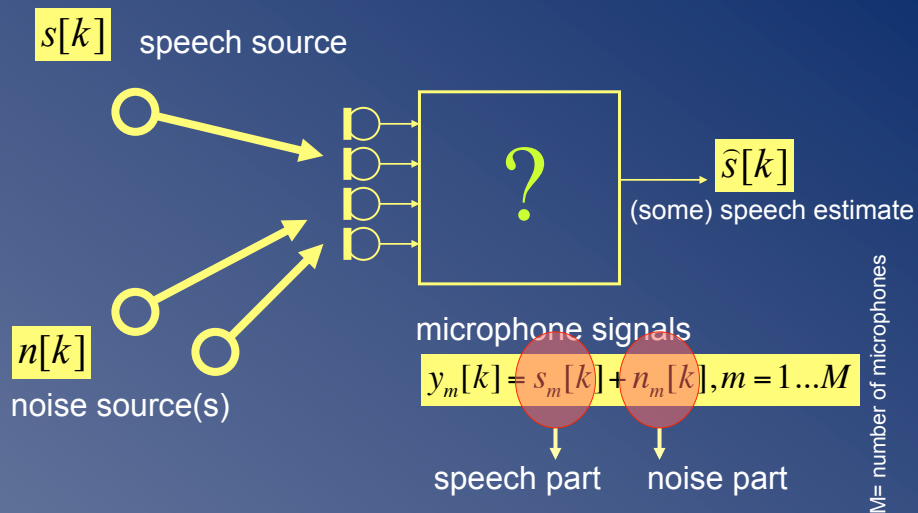
3. Filter speech frame $y[m] \rightarrow \hat{s}_i[m]$

Repeat
until
some error
criterion is
satisfied

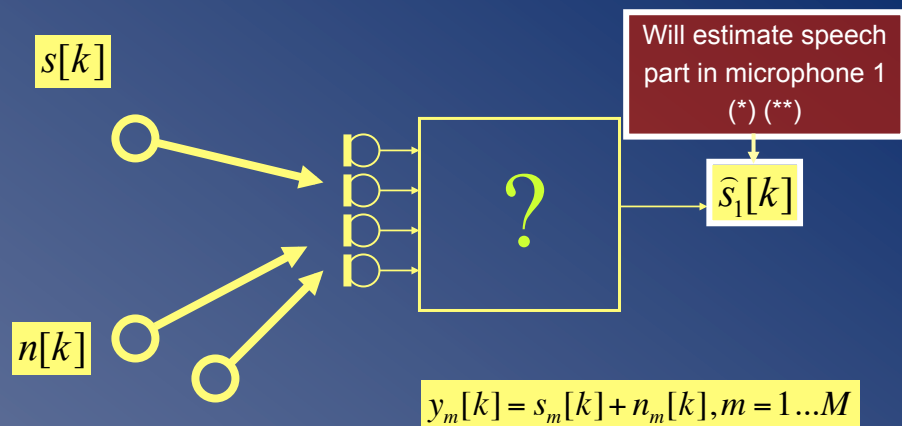
Overview

- **Spectral subtraction for single-micr. noise reduction**
 - Single-microphone noise reduction problem
 - Spectral subtraction basics (=spectral filtering)
 - Features: gain functions, implementation, musical noise,...
 - Iterative Wiener filter based on speech signal model
- **Multi-channel Wiener filter for multi-micr. noise red.**
 - Multi-microphone noise reduction problem
 - Multi-channel Wiener filter (=spectral+spatial filtering)
- **Kalman filter based noise reduction**
 - Kalman filters
 - Kalman filters for noise reduction

Multi-Microphone Noise Reduction Problem



Multi-Microphone Noise Reduction Problem



(*) Estimating $s[k]$ is more difficult, would include dereverberation...

(**) This is similar to single-microphone model (p.3), where additional microphones ($m=2..M$) help to get a better estimate

Multi-Microphone Noise Reduction Problem

- Data model:

$$\begin{aligned} \mathbf{Y}(\omega) &= \mathbf{S}(\omega) + \mathbf{N}(\omega) \\ &= \mathbf{d}(\omega).S(\omega) + \mathbf{N}(\omega) \end{aligned}$$

$$\begin{bmatrix} Y_1(\omega) \\ Y_2(\omega) \\ \vdots \\ Y_M(\omega) \end{bmatrix} = \begin{bmatrix} H_1(\omega) \\ H_2(\omega) \\ \vdots \\ H_M(\omega) \end{bmatrix} . S(\omega) + \begin{bmatrix} N_1(\omega) \\ N_2(\omega) \\ \vdots \\ N_M(\omega) \end{bmatrix}$$

See Lecture-2 on multi-path propagation, with q left out for conciseness.

$H_m(\omega)$ is complete transfer function from speech source position to m-th microphone

Multi-Channel Wiener Filter (MWF)

- Data model:

$$\mathbf{Y}(\omega) = \mathbf{d}(\omega).S(\omega) + \mathbf{N}(\omega)$$

- Will use linear filters to obtain speech estimate (as in Lecture-2)

$$\hat{S}_1(\omega) = \sum_{m=1}^M F_m^*(\omega).Y_m(\omega) = \mathbf{F}^H(\omega).\mathbf{Y}(\omega)$$

- Wiener filter (=linear MMSE approach)

$$\min_{\mathbf{F}(\omega)} E\left\{ \left| S_1(\omega) - \mathbf{F}^H(\omega).\mathbf{Y}(\omega) \right|^2 \right\}$$

Note that (unlike in DSP-CIS) 'desired response' signal $S_1(\omega)$ is **unknown** here (!), hence solution will be 'unusual' ...

Multi-Channel Wiener Filter (MWF)

- Wiener filter solution is (see DSP-CIS)

$$\begin{aligned}
 \mathbf{F}(\omega) &= \underbrace{E\{\mathbf{Y}(\omega) \cdot \mathbf{Y}^H(\omega)\}}_{\text{autocorrelation}}^{-1} \underbrace{E\{\mathbf{Y}(\omega) \cdot S_1^*(\omega)\}}_{\text{crosscorrelation}} \\
 &= \dots \quad (\text{with } E\{S(\omega) \cdot N_1^*(\omega)\} = 0) \\
 &= \underbrace{E\{\mathbf{Y}(\omega) \cdot \mathbf{Y}^H(\omega)\}}_{\text{compute during speech+noise periods}}^{-1} \cdot \left(\underbrace{E\{\mathbf{Y}(\omega) \cdot Y_1^*(\omega)\}}_{\text{compute during noise-only periods}} - \underbrace{E\{\mathbf{N}(\omega) \cdot N_1^*(\omega)\}}_{\text{compute during noise-only periods}} \right)
 \end{aligned}$$

- All quantities can be computed !
- Special case of this is single-channel Wiener filter formula on p.9
- In practice, use alternative to 'subtraction' operation (see literature)

Multi-Channel Wiener Filter (MWF)

- Note that...

MWF combines spatial filtering (as in Lecture-2) with single-channel spectral filtering (as in single-channel noise reduction)

if

$$\begin{bmatrix} Y_1(\omega) \\ Y_2(\omega) \\ \vdots \\ Y_M(\omega) \end{bmatrix} = \underbrace{\mathbf{d}(\omega)}_{\text{steering vector}} \cdot S(\omega) + \underbrace{\mathbf{N}(\omega)}_{\text{noise}}$$

$$E\{\mathbf{N}(\omega) \cdot \mathbf{N}^H(\omega)\} = \mathbf{\Phi}_{NN}(\omega)$$

then...

Multi-Channel Wiener Filter (MWF)

...then it can be shown that

$$\mathbf{F}(\omega) = \underbrace{\alpha(\omega)}_{\text{scalar}} \cdot \Phi_{NN}^{-1}(\omega) \cdot \mathbf{d}(\omega)$$

- ① $\Phi_{NN}^{-1}(\omega) \cdot \mathbf{d}(\omega)$ represents a spatial filtering (*)

Compare to superdirective & delay-and-sum beamforming (Lecture-2)

- Delay-and-sum beamf. maximizes array gain in white noise field
- Superdirective beamf. maximizes array gain in diffuse noise field
- MWF maximizes array gain in unknown (!) noise field.

MWF is operated without invoking any prior knowledge (steering vector/noise field) ! (the secret is in the voice activity detection... (explain))

(*) Note that spatial filtering can improve SNR, spectral filtering never improves SNR (at one frequency)

Multi-Channel Wiener Filter (MWF)

...then it can be shown that

$$\mathbf{F}(\omega) = \underbrace{\alpha(\omega)}_{\text{scalar}} \cdot \Phi_{NN}^{-1}(\omega) \cdot \mathbf{d}(\omega)$$

- ① $\Phi_{NN}^{-1}(\omega) \cdot \mathbf{d}(\omega)$ represents a spatial filtering (*)

- ② $\alpha(\omega)$ represents an additional 'spectral post-filter'

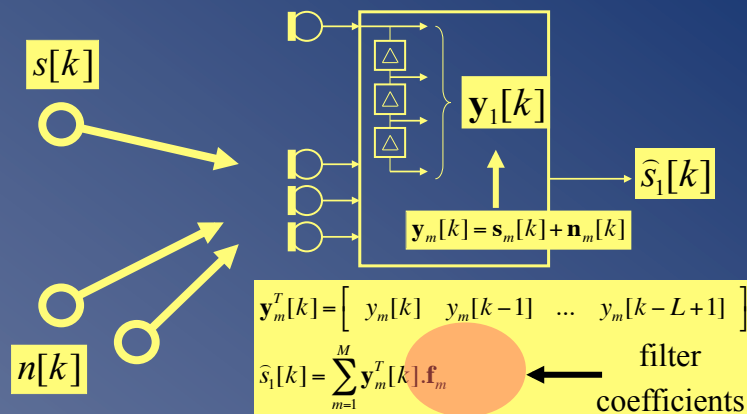
i.e. single-channel Wiener estimate (p.9), applied to output signal of spatial filter

$$\alpha(\omega) = \dots = \frac{|S(\omega)|^2 \cdot H_1^*(\omega)}{(\mathbf{d}^H(\omega) \cdot \Phi_{NN}^{-1}(\omega) \cdot \mathbf{d}(\omega)) |S(\omega)|^2 + 1}$$

(prove it!)

Multi-Channel Wiener Filter: Implementation

- Implementation with short-time Fourier transform: see p.10
- Implementation with time-domain linear filtering:



Multi-Channel Wiener Filter: Implementation

- Implementation with time-domain linear filtering:

$$\min_{\mathbf{f}} E \left\{ \left| s_1[k] - \mathbf{y}^T[k] \mathbf{f} \right|^2 \right\}$$

$$\mathbf{f} = [\mathbf{f}_1^T \quad \mathbf{f}_2^T \quad \dots \quad \mathbf{f}_M^T]^T$$

$$\mathbf{y}[k] = [\mathbf{y}_1^T[k] \quad \mathbf{y}_2^T[k] \quad \dots \quad \mathbf{y}_M^T[k]]^T$$

Solution is...

$$\mathbf{f} = \left[E \{ \mathbf{y}[k] \mathbf{y}[k]^T \} \right]^{-1} \cdot E \{ \mathbf{y}[k] s_1[k] \}$$

$$= \left[E \{ \mathbf{y}[k] \mathbf{y}[k]^T \} \right]^{-1} \cdot \left[E \{ \mathbf{y}[k] y_1[k] \} - E \{ \mathbf{n}[k] n_1[k] \} \right]$$

compute during speech+noise periods

compute during noise-only periods

Overview

- **Spectral subtraction for single-micr. noise reduction**
 - Single microphone noise reduction problem
 - Spectral subtraction basics (=spectral filtering)
 - Features: gain functions, implementation, musical noise,...
 - Iterative Wiener filter based on speech signal model
- **Multi-channel Wiener filter for multi-micr. noise red.**
 - Multi-microphone noise reduction problem
 - Multi-channel Wiener filter (=spectral+spatial filtering)
- **Kalman filter based noise reduction**
 - Kalman filters : See Lecture-6
 - Kalman filters for noise reduction

Kalman filter for Speech Enhancement

- Assume AR model of speech and noise

$$s[k] = \sum_{\bar{n}=1}^{N_s} \alpha_{\bar{n}} s[k - \bar{n}] + \sqrt{g_s} u[k]$$

$$n[k] = \sum_{\bar{n}=1}^{N_n} \beta_{\bar{n}} n[k - \bar{n}] + \sqrt{g_n} w[k]$$

$u[k], w[k]$ = zero mean, unit variance, white noise

- Equivalent state-space model is...

$$\begin{cases} \mathbf{x}[k+1] &= \mathbf{A}\mathbf{x}[k] + \mathbf{v}[k] \\ y[k] &= \mathbf{c}^T \mathbf{x}[k] \end{cases}$$

y =microphone signal

Kalman filter for Speech Enhancement

with:

$$\mathbf{x}^T[k] = \begin{bmatrix} s[k - N_s + 1] & \cdots & s[k] & n[k - N_n + 1] & \cdots & n[k] \end{bmatrix}$$

$$\mathbf{v}[k] = \mathbf{G} \begin{bmatrix} u[k] & w[k] \end{bmatrix}^T \longrightarrow \mathbf{Q} = \mathbf{G} \mathbf{G}^T$$

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_s & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_n \end{bmatrix}; \mathbf{C}^T = \begin{bmatrix} \overbrace{0 \cdots 0}^M & 1 & \overbrace{0 \cdots 0}^N & 1 \end{bmatrix}; \mathbf{G} = \begin{bmatrix} \mathbf{g}_s & \mathbf{0} \\ \mathbf{0} & \mathbf{g}_n \end{bmatrix}$$

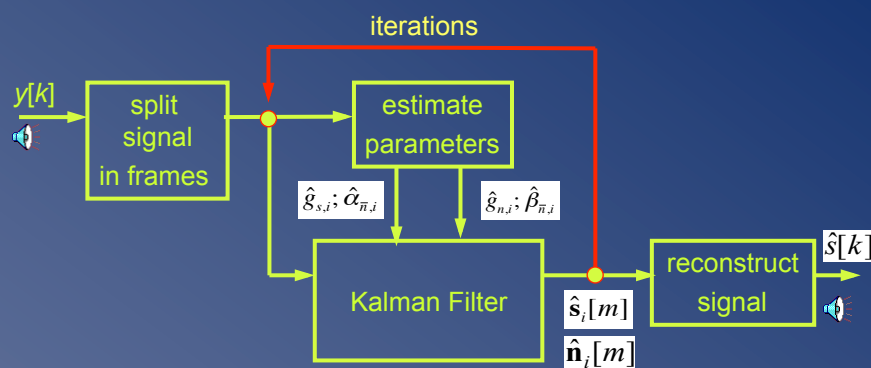
$$\mathbf{A}_s = \begin{bmatrix} 0 & 1 & \cdots & 0 \\ \vdots & 0 & \ddots & 0 \\ 0 & \cdots & 0 & 1 \\ \alpha_{N_s} & \alpha_{N_s-1} & \cdots & \alpha_1 \end{bmatrix}; \mathbf{A}_n = \begin{bmatrix} 0 & 1 & \cdots & 0 \\ \vdots & 0 & \ddots & 0 \\ 0 & \cdots & 0 & 1 \\ \beta_{N_n} & \beta_{N_n-1} & \cdots & \beta_1 \end{bmatrix}$$

$$\mathbf{g}_s^T = \begin{bmatrix} 0 & \cdots & 0 & \sqrt{g_s} \end{bmatrix}; \mathbf{g}_n^T = \begin{bmatrix} 0 & \cdots & 0 & \sqrt{g_n} \end{bmatrix};$$

$s[k]$ and $n[k]$ are included in state vector, hence can be estimated by Kalman Filter

Kalman filter for Speech Enhancement

Iterative algorithm (details omitted)



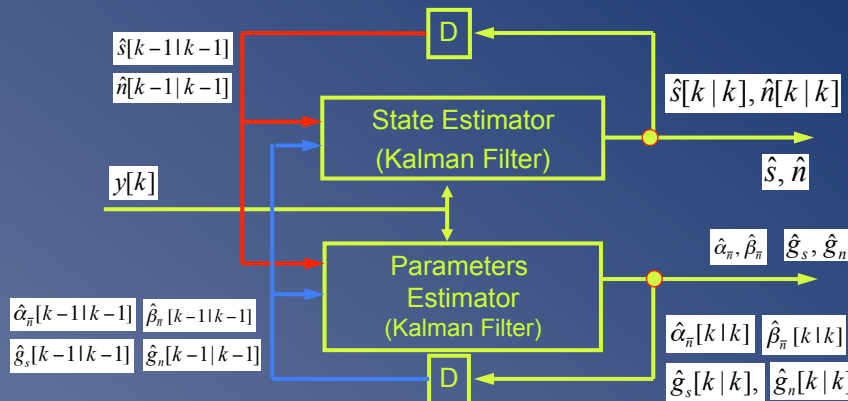
Disadvantages iterative approach:

- complexity
- delay

Kalman filter for Speech Enhancement

Sequential algorithm (details omitted)

iteration index \rightarrow time index (no iterations)



CONCLUSIONS

- **Single-channel noise reduction**
 - Basic system is spectral subtraction
 - Only spectral filtering, hence can only exploit differences in spectra between noise and speech signal:
 - noise reduction at expense of speech distortion
 - achievable noise reduction may be limited
- **Multi-channel noise reduction**
 - Basic system is MWF,
 - Provides spectral + spatial filtering (links with beamforming!)
- **Iterative Wiener filter & Kalman filtering**
 - Signal model based approach