

Lab 05 probabilistic noise suppression methods

개요: fixed threshold로 잠음구간 (non-speech period)을 추정하면

- ① 입력 gain에 영향받는다.
- ② noise의 상대적인 크기 (SNR)의 변화에 대처하지 못한다.
- ③ 음성구간으로 잘못 추정하면 음성도 차감될 수 있다.

입력 신호에 따라 유연하게 대처하기 위해서 (e.g. adaptive thresholding)
확률 모델을 사용하는 방법은 연습한다.

* 이번 Lab부터는 VAD (voice activity detection)을 주로 수행한다.

- EPD (end point detection): 음성구간의 시작과 끝을 추정
- VAD: 음성의 유무 (active/deactive)를 주로 frame 별로 (주로 10ms 단위) 결정한다. 0/1의 hard decision, $P(\text{voice} | y(t))$ 의 soft decision이 가능하며 확률 모델을 사용하는 경우 hard decision은 soft decision을 thresholding 하여 얻을 수 있다.
- $P(\text{voice} | y(t))$: posterior voicing probability of $y(t)$ being observed
y(t)가 관측되었을 때 후향 확률
- EPD는 VAD 결과를 후처리 (post processing)하여 얻을 수 있다.
thresholding, median filtering, 구간 길이 검증

I. noise spectrum estimation and suppression by probabilistic VAD

- a. 적절한 확률 모델로 voicing probability $P(V | y)$ 구함.
- b. 각 frame 별로 noise 확률 계산, $1 - P(V | y_k)$, $y_k = [y_{k1} \dots y_{kN}]$
frame k
- c. (hard decision)

$$\left. \begin{aligned} 1 - P(V | y_k) &\geq \theta_{\text{noise}} \rightarrow I(k) = 1 \\ &< \theta_{\text{noise}} \rightarrow I(k) = 0 \end{aligned} \right\} k = 1 \dots \underbrace{K}_{\text{\#frames}}$$

$$|\tilde{N}(w)|^2 = E[|N(w)|^2] = \frac{\sum_{k=1}^K I(k) |Y(k, w)|^2}{\sum_{k=1}^K I(k)}$$

- d. (soft decision) MAP estimation (maximum a posteriori estimation)
noise spectrum은 posterior probability의 weighted summation

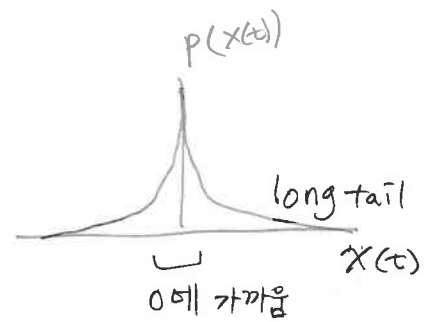
$$|\tilde{N}(w)|^2 = \frac{\sum_{k=1}^K P(\text{noise} | y_k) |Y(k, w)|^2}{\sum_{k=1}^K P(\text{noise} | y_k)}, \quad \begin{aligned} P(\text{noise} | y_k) &= 1 - P(\text{voice} | y_k) \end{aligned}$$

- e. lab04의 time-domain FIR Wiener filtering 이용

II. time-domain VAD

a. 음성신호의 분포는 0이 많은 sparse distribution

Laplacian distribution을 따른다고 알려져 있다. (exponential decay)



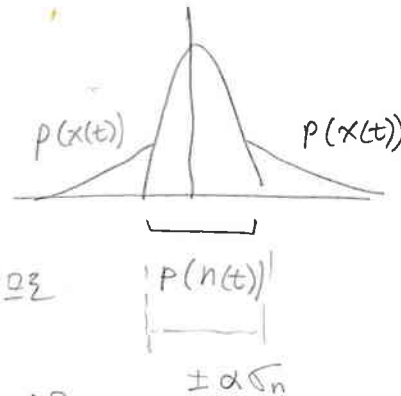
$$p(x(t)) \sim \frac{1}{\sqrt{2}\sigma_x} \exp\left(-\sqrt{2} \frac{|x-u_x|}{\sigma_x}\right)$$

* 가정: noise의 크기가 speech보다 작다.

b. $y(t) = x(t) + n(t)$

noise의 분포는 Gaussian으로 가정

$$p_n(n(t)) \sim \frac{1}{\sigma_n \sqrt{2\pi}} \exp\left(-\frac{(x-u_n)^2}{2\sigma_n^2}\right)$$



speech의 분포는 sharp peak가 사라지므로 Gaussian으로 modeling

$$p_x(x(t)) \sim \frac{1}{\sqrt{2\pi}\sigma_x} \exp\left(-\frac{(x-u_x)^2}{2\sigma_x^2}\right)$$

dual Gaussian mixture model with mean 0 ($u_x = u_n = 0$)

$$P(y(t)) = p_{\text{voice}} \cdot p_x(y(t)) + (1 - p_{\text{voice}}) \cdot p_n(y(t)) \quad \text{--- ①}$$

voicing probability는? - posterior prob of $y(t)$

$$P(\text{voice} | y(t)) = \frac{P(\text{voice}, y(t))}{P(y(t))} = \frac{P(v) \cdot P(y(t)|v)}{P(y(t))} = \frac{p_{\text{voice}} \cdot p_x(y)}{P(y)} \quad \text{--- ②}$$

$$P(\text{noise} | y(t)) = \frac{P_{\text{noise}} \cdot p_n(y(t))}{P(y(t))} = \frac{(1 - p_{\text{voice}}) p_n(y)}{P(y)} \quad \text{--- ③}$$

$$* P(v|y) + P(n|y) = 1 \rightarrow P(n|y) = 1 - P(v|y)$$

①, ②을 조합하면

$$P(\text{noise} | y(t)) = \frac{(1 - p_v) \cdot p_n(y(t))}{p_v \cdot p_x(y) + (1 - p_v) \cdot p_n(y)} \gtrless 0.5 \begin{matrix} \nearrow \text{noise} \\ \searrow \text{speech} \end{matrix}$$

Soft decision

hard decision

c. sample 별로 하지 않고 frame 별로 decision

$$\text{let } U[k] \equiv \sqrt{\sum_{t \in \text{frame } k} y(t)^2} \quad (\text{frame energy. 평균을 위해 } N_f \text{로 나누어도 된다. 여차피 상수})$$

$$P(\text{frame } k \text{ is noise}) = P(\text{noise} | U[k]), \text{ 당연히 확률도 frame energy로}$$

II. time-domain VAD (continued)

d. dual Gaussian mixture training

EM (expectation maximization) algorithm 사용

<offline training>

① compute $U[k] = \sqrt{\sum_{t \in \text{frame } k} y(t)^2}$ for all $k=1 \dots K$ → $\sqrt{\quad}$ 씌우는 것은 나중에 $U[k]^2$ 으로
확률계산 하기 때문
(편상면이)

② Sort, and split $\left\{ \begin{array}{l} \text{lower } 1 \sim \frac{K}{2} \text{ frames} \rightarrow \text{compute } \sigma_n^2 \\ \text{upper } \frac{K}{2} \sim K \text{ frames} \rightarrow \sigma_x^2 \end{array} \right.$

③ $u_n = u_x = 0, P_{\text{voice}} = 0.5$

④ compute $p_x(u[k]) = \frac{1}{\sqrt{2\pi}\sigma_x} \exp\left(-\frac{u[k]^2}{2\sigma_x^2}\right)$

$$p_n(u[k]) = \frac{1}{\sqrt{2\pi}\sigma_n} \exp\left(-\frac{u[k]^2}{2\sigma_n^2}\right)$$

$$p(\text{voice} | u[k]) = \frac{P_v p_x(u[k])}{P_v p_x(u[k]) + (1-P_v) p_n(u[k])}$$

⑤ update $P_v \Leftarrow \frac{\sum_{k=1}^K P(v | u[k])}{K}$

$$\sigma_x^2 = \frac{\sum_{k=1}^K P(v | u[k]) \cdot u[k]^2}{\sum_{k=1}^K P(v | u[k])}$$

$$\sigma_n^2 = \frac{\sum_k P(\text{noise} | u[k]) \cdot u[k]^2}{\sum_k P(\text{noise} | u[k])} = \frac{\sum_k (1 - P(v | u[k])) \cdot u[k]^2}{\sum_k (1 - P(v | u[k]))}$$

e. noise spectrum estimation

$$\begin{aligned} |\tilde{N}(w)|^2 &= E_N[N(k, w)] = \frac{\sum_k P(\text{noise} | u[k]) \cdot u[k]^2}{\sum_k P(\text{noise} | u[k])} = \frac{\sum_k (1 - P(v | u[k])) \cdot u[k]^2}{\sum_k (1 - P(v | u[k]))} \quad (\text{soft}) \\ &= \frac{\sum_k (1 - I_x[k]) \cdot u[k]^2}{\sum_k (1 - I_x[k])} \quad (\text{hard}) \end{aligned}$$

where $I_x[k] = \begin{cases} 1 & \text{if } P(v | u[k]) \geq 0.5 \\ 0 & \text{if } \quad \quad \quad < 0.5 \end{cases}$

$I_x[k]$ is voice의 posterior 확률이 더 클때 1이 되는 indicator function

III. frequency domain VAD using Rayleigh mixture model

만약 noise 특성이 (세로, 들어) 저주파에 집중되어 있고 amplitude가 매우 크다면 time-domain VAD에 매우 불리함



noise 크기가 매우 크지만 주파수 영역에서는 narrow band 이와 같은 colored noise에 대해서 신뢰성 있는 VAD 결과를 얻기 위해 frequency domain에서 확률 모델을 사용하여 본다.

$Y(k, \omega)$: frame k 의 주파수 ω 성분

$$|Y(k, \omega)|^2 = Y(k, \omega) Y^*(k, \omega) \quad - \text{ } ^* \text{는 complex conjugate (결제 복소수)}$$

$$|Y(k, \omega)| \sim \text{Rayleigh}(\sigma), \text{ pdf: } f_X(x|\sigma) = \frac{x}{\sigma^2} \exp\left(-\frac{x^2}{2\sigma^2}\right)$$

Rayleigh Distribution

두개의 Random Variable X, Y 가 평균이 0이고 같은 variance σ^2 의 Gaussian distribution을 따른다면 ($X \sim N(0, \sigma^2), Y \sim N(0, \sigma^2)$)

random variable $R = \sqrt{X^2 + Y^2}$ 은 Rayleigh distributed (wikipedia)

*단, $X \perp Y \leftrightarrow E[XY] = 0$, 두 변수 X 와 Y 는 independent (i.i.d.) (uncorrelated)

따라서 $Y(k, \omega) = A + jB$ 이고 A, B 가 independent Gaussian with same variance σ^2

$$\Rightarrow |Y(k, \omega)| \sim f_X(x|\sigma) = \frac{x}{\sigma^2} \exp\left(-\frac{x^2}{2\sigma^2}\right)$$

$$\textcircled{1} f_X(0|\sigma) = \frac{0}{\sigma^2} \cdot \exp(-0) = 0.1 = 0 //$$

$$\textcircled{2} x < 0 \text{는 정의되지 않음 } f_X(x|\sigma) = 0 \text{ for } x < 0$$

$$\textcircled{3} y = x \text{ 와 } y = \exp(-x^2) \text{ 이 곱해진 형태}$$



$$\textcircled{4} x = \sigma \text{는 mode (local peak)} \rightarrow \leftarrow \text{중요함!!}$$

$$\textcircled{5} \text{parameter는 } \sigma^2 \text{ 하나밖에 없으며 (variance는 아님)}$$

다음과 같이 추정한다.

$$\hat{\sigma}^2 = \frac{1}{2N} \sum_{i=1}^N x_i^2 \quad \text{maximum likelihood, unbiased estimation}$$

$$\textcircled{6} E[X] = \sigma \sqrt{\frac{\pi}{2}}$$

$$\text{Var}[X] = E[(X - E_X)^2] = \frac{4 - \pi}{2} \sigma^2, \text{ 즉 } \sigma^2 \text{은 variance가 아님}$$

III-2. frequency domain VAD part 2

multivariate extension of Rayleigh distribution

* 매우 복잡하다. 따라서 다음과 같이 변수들이 (dimensions들이) independent 하다고 가정.

$$\mathbf{x} = [x_1, x_2, \dots, x_D]$$

$$f_{\mathbf{x}}(\mathbf{x} | \Sigma) \cong f_{x_1}(x_1 | \sigma_1^2) \cdot f_{x_2}(x_2 | \sigma_2^2) \cdots f_{x_D}(x_D | \sigma_D^2)$$

covariance matrix

$$= \prod_{d=1}^D f_{x_d}(x_d | \sigma_d^2) \quad \text{각 individual pdf들의 곱으로 표현}$$

$$= \prod_{d=1}^D \frac{x_d}{\sigma_d^2} \cdot \exp\left(-\frac{x_d^2}{2\sigma_d^2}\right) = \left(\prod_{d=1}^D \frac{x_d}{\sigma_d^2}\right) \cdot \exp\left(-\sum_{d=1}^D \frac{x_d^2}{2\sigma_d^2}\right)$$

$$\log f_{\mathbf{x}}(\mathbf{x} | \Sigma) = \sum_d \log x_d - \sum_d 2 \log \sigma_d - \sum_d \frac{x_d^2}{2\sigma_d^2}$$

$$Y(k, \omega) = X(k, \omega) + N(k, \omega) \quad X \text{ and } N \text{ are independent (i.i.d.) distinctively distributed}$$

$|Y(k, \omega)| = |X(k, \omega) + N(k, \omega)| \Rightarrow$ dual Rayleigh mixture model 적용.

* 엄밀히 말하면 Y 는 두 complex Random Variables X 와 N 의 합의 절대값 mixture model은 X 와 N 이 변할아 관측되는 것을 가정하므로 오류가 있다. 하지만 speech는 sparse 하므로 적용해 보도록 한다

* $\sigma_x^2 > \sigma_n^2$ 으로 가정한다



$$f_Y(y) = P_X f_X(y | \sigma_x^2) + P_N f_N(y | \sigma_n^2)$$

$$= P_X \frac{y}{\sigma_x^2} \exp\left(-\frac{y^2}{2\sigma_x^2}\right) + P_N \frac{y}{\sigma_n^2} \exp\left(-\frac{y^2}{2\sigma_n^2}\right)$$

$$= y \left(\frac{P_X}{\sigma_x^2} \exp\left(-\frac{y^2}{2\sigma_x^2}\right) + \frac{P_N}{\sigma_n^2} \exp\left(-\frac{y^2}{2\sigma_n^2}\right) \right)$$

* posterior probability

$$P(\text{voice} | y) = \frac{P_X f_X(y | \sigma_x^2)}{f_Y(y)} = \frac{P_X f_X(y | \sigma_x^2)}{P_X f_X(y) + P_N f_N(y)}$$

$$P(\text{noise} | y) = \frac{P_N f_N(y)}{P_X f_X(y) + P_N f_N(y)} = 1 - P(\text{voice} | y)$$

* EM algorithm으로 학습한다.

* $\sigma_x^2 < \sigma_n^2 \rightarrow \text{swap}$ (학습하다가 일어날 수 있다).

* 학습이 끝나면 $E[|\tilde{N}(\omega)|^2] = \sigma_n^2$ (noise Rayleigh의 mode)

III-3. frequency domain VAD part 3

$|\tilde{N}(w)|^2$ 은 주파수 w 에 따라 Rayleigh mixture model을 하지 않고 multivariate Rayleigh로 추정

$$\text{let } \mathbf{y}(k) = [|Y(k, w_1)| \quad |Y(k, w_2)| \quad \dots \quad |Y(k, w_D)|]^T$$

주파수 성분들의 절대값들로 구성된 column vector.

speech와 noise의 pdf는 각

주파수 값들이 independent하다고 가정하고 각 dimension의 scalar pdf의 곱으로 표현

$$f_x(\mathbf{y}) = f_{x_1}(y_1) \cdot f_{x_2}(y_2) \cdot \dots \cdot f_{x_D}(y_D) = \prod_d f_{x_d}(y_d)$$

$$= \prod_d \frac{y_d}{\sigma_{x_d}^2} \cdot \exp\left(-\frac{y_d^2}{2\sigma_{x_d}^2}\right) = \left(\prod_d \frac{y_d}{\sigma_{x_d}^2}\right) \cdot \exp\left(-\sum_d \frac{y_d^2}{2\sigma_{x_d}^2}\right)$$

$$f_n(\mathbf{y}) = \prod_d \frac{y_d}{\sigma_{nd}^2} \cdot \exp\left(-\frac{y_d^2}{2\sigma_{nd}^2}\right) = \left(\prod_d \frac{y_d}{\sigma_{nd}^2}\right) \cdot \exp\left(-\sum_d \frac{y_d^2}{2\sigma_{nd}^2}\right)$$

* Computation of multivariate pdf using log

$$f_x(\mathbf{y}) = \exp\left(\sum_d \log y_d - 2 \sum_d \log \sigma_{x_d} - \sum_d \frac{y_d^2}{2\sigma_{x_d}^2}\right)$$

exponential function을

매우 많이 곱하면 수가 많이 작아져서 underflow가 날 수 있다. 따라서 곱하기를 줄이고 더하기로 대체한다.

$$f_n(\mathbf{y}) = \exp\left(\sum_d \log y_d - 2 \sum_d \log \sigma_{nd} - \sum_d \frac{y_d^2}{2\sigma_{nd}^2}\right)$$

* dual mixture multivariate Rayleigh

$$P(\text{voice} | \mathbf{y}) = \frac{P_x \cdot f_x(\mathbf{y})}{P_x \cdot f_x(\mathbf{y}) + P_n \cdot f_n(\mathbf{y})} \triangleq \gamma_x(\mathbf{y})$$

효율적인 계산을 위하여 분모/분자를 $P_x \cdot f_x(\mathbf{y})$ 로 나눔

$$= \frac{1}{1 + \frac{P_n}{P_x} \cdot \frac{f_n(\mathbf{y})}{f_x(\mathbf{y})}}$$

$$\begin{aligned} \log \frac{f_x(\mathbf{y})}{f_n(\mathbf{y})} &= \sum_d \left(\underbrace{\log y_d - 2 \log \sigma_{x_d} - \frac{y_d^2}{2\sigma_{x_d}^2}}_{\text{voice part}} - \underbrace{\log y_d + 2 \log \sigma_{nd} + \frac{y_d^2}{2\sigma_{nd}^2}}_{\text{noise part}} \right) \\ &= \sum_d \left(2 \log \frac{\sigma_{nd}}{\sigma_{x_d}} + \frac{y_d^2}{2} \left(\frac{1}{\sigma_{nd}^2} - \frac{1}{\sigma_{x_d}^2} \right) \right) \triangleq Z_x(\mathbf{y}) \end{aligned}$$

\Rightarrow log 1번, exp 1번으로 $\gamma_x(\mathbf{y})$ 계산 가능

$$\gamma_x(\mathbf{y}) = \frac{1}{1 + \frac{P_n}{P_x} \cdot \exp Z_x(\mathbf{y})}$$

$$\gamma_n(\mathbf{y}) = P(\text{noise} | \mathbf{y}) = 1 - P(\text{voice} | \mathbf{y}) = 1 - \gamma_x(\mathbf{y})$$

III-4. frequency domain VAD using Rayleigh mixture model part 4

< EM algorithm >

① Initialization

초기 10~30개의 frame으로 energy 큰 것 절반 $\rightarrow f_x(y)$ ^{initial} parameters
작은 것 절반 $\rightarrow f_n(y)$ initial parameters

② Expectation

compute $\gamma_x(y(k))$ for all $k = 1 \dots \#frames$

③ maximization

update σ_{xd} and σ_{nd} , probabilistically (maximum likelihood estimate)

$$\sigma_{xd}^2 = \frac{\sum_{k=1}^K \gamma_x(y(k)) \cdot y_d^2(k)}{2 \sum_{k=1}^K \gamma_x(y(k))}, \quad \sigma_{nd}^2 = \frac{\sum_{k=1}^K (1 - \gamma_x(y(k))) \cdot y_d^2(k)}{2 \sum_{k=1}^K (1 - \gamma_x(y(k)))}$$

for all $d = 1 \dots D$

④ Repeat 3~4

< Noise Spectrum Estimation >

Nois Rayleigh의 mode 위치인 σ_{nd} 로 예측

$$|\tilde{N}(w_d)|^2 = \sigma_{nd}^2$$

* implementation tips.

① $\sigma_n < \sigma_x$ 가 보장되도록 유지

② $\sum \gamma_x > 0$, $\sum (1 - \gamma_x) > 0$ 인지 확인하기 너무 작으면 update 하지 않음 (이미 diverged)

③ $|Y(k, w)|^2 = 0$ 이면 Wiener filter gain 계산할 때 "divide by zero" exception 발생.

$|Y(k, w)|^2 \leftarrow \max(\epsilon_y, |Y(k, w)|^2)$ 등의 방법 사용.

④ $y(t)$ 에 매우 작은 크기 ($\sigma = 10^{-10} \sim 10^{-6}$)의 Gaussian random noise를 더해주는 방법도 있다.

< Online EM algorithm >

실시간으로 들어오는 입력에는 offline training 불가능, 수정 필요

① Initialization: 초기 10~30 frame으로 초기화, Wiener filtering 하지 않음

② Expectation at current frame

compute $\gamma_x(y(k))$ for current frame k .

③ Adaptive maximization

$$\sigma_{xd}^2 = (1 - \alpha \gamma_x(k)) \sigma_{xd}^{2(\text{old})} + \frac{1}{2} \cdot \alpha \gamma_x(k) \cdot y_d^2(k), \quad 0 \leq \alpha \leq 1$$

$$\sigma_{nd}^2 = (1 - \alpha (1 - \gamma_x(k))) \sigma_{nd}^{2(\text{old})} + \frac{1}{2} \alpha (1 - \gamma_x(k)) \cdot y_d^2(k)$$

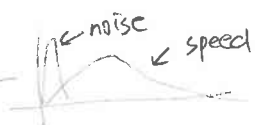
④ repeat ②-③

IV. frequency domain VAD using LogNormal distribution

Rayleigh 분포는 parameter가 하나밖에 없기 때문에 표현력이 떨어짐.

특히 mixture modeling을 하면 이론과 달리 잘 안 맞을 수 있다.

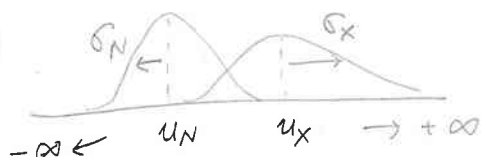
또한, noise와 speech의 scale 차이가 크면 한쪽으로 수렴될 수 있다.



noise의 분포가 sharp하지만
scale이 speech에 비해 너무 작다

< LogNormal >

Let $y(k, \omega) = \log |Y(k, \omega)|^2$, 그리고 Gaussian mixture model 사용



⇒ scale mismatch et
표현력 ($\sigma \rightarrow u, \sigma$) 증가

PDF

$$f_x(y) = \frac{1}{(2\pi)^{D/2} \cdot \sqrt{|\det \Sigma|}} \cdot \exp\left(-\frac{1}{2} (y - u_x)^T \Sigma^{-1} (y - u_x)\right)$$

u_x : mean vector, Σ : covariance matrix

너무 복잡하기 때문에 independent 가정

$$\approx \prod_d \frac{1}{\sqrt{2\pi} \sigma_{xd}} \cdot \exp\left(-\frac{(y_d - u_{xd})^2}{2\sigma_{xd}^2}\right)$$

$f_n(y)$ 도 같은 방법으로 modeling 하라

$$\gamma_x(y) = \frac{P_x f_x(y)}{P_x f_x(y) + P_n f_n(y)} = \frac{1}{1 + \frac{P_n}{P_x} \frac{f_n(y)}{f_x(y)}}$$

EM update

$$u_x = \frac{\sum_k \gamma_x(y) \cdot y}{\sum_k \gamma_x(y)}, \quad \sigma_{xd}^2 = \frac{\sum_k \gamma_x(y) \cdot (y_d - u_{xd})^2}{\sum_k \gamma_x(y)}$$

$$u_n = \frac{\sum_k (1 - \gamma_x(y)) \cdot y}{\sum_k (1 - \gamma_x(y))}, \quad \sigma_{nd}^2 = \frac{\sum_k (1 - \gamma_x(y)) \cdot (y_d - u_{nd})^2}{\sum_k (1 - \gamma_x(y))}$$

$$|\hat{N}(w_d)|^2 \approx \exp(u_{nd}) \quad \because \text{Gaussian의 mode는 mean vector}$$