PHONEME SELECTIVE SPEECH ENHANCEMENT

USING THE GENERALIZED PARAMETRIC

SPECTRAL SUBTRACTION ESTIMATOR

Amit Das and John H.L. Hansen



Center for Robust Speech Systems (CRSS) Erik Jonsson School of Engineering & Computer Science Department of Electrical Engineering University of Texas at Dallas Richardson, Texas 75083-0688, U.S.A.



ICASSP 2011 May 22-27, 2011 Prague, Czech Republic



Email: das_amit_ece@yahoo.com, John.Hansen@utdallas.edu

Slide 1



- Degree of noise impact is *not uniform* across the speech spectrum or phoneme sequence.
- Regions of low SNR in the spectrum more adversely affected than high SNR.
- SNR variation across the spectrum depends on both: phoneme class and noise characteristics.

Slide 2

Require phoneme class selective enhancement algorithm.



Email: das_amit_ece@yahoo.com, John.Hansen@utdallas.edu



Each phoneme class is distorted differently for each noise type

Sased on Frequency content, articulatory structure, influence of noise for that phoneme, and stationarity of noise

Speech Class type adaptation for enhancement:

McAulay & Malpass (IEEE Trans. ASSP 1980) – softdecision noise suppression

Hansen & Arslan (IEEE Trans. SAP 1995) – HMM based phone class partitioning for AutoLSP

Slide 3

Das & Hansen, (Interspeech-07) - Class Constrained ROVER Based Speech Enhancement



Email: das_amit_ece@yahoo.com, John.Hansen@utdallas.edu







- Gain functions of current algorithms (e.g. MMSE) rely on Apriori/Aposteriori SNR. Gains are limited.
- Can we do better? Could generate a family of gain functions & provide more enhancement coverage.





Modify Error Function to place more weight on spectral valleys (low SNR region) than peaks (high SNR region).

Hence, new error function is Weighted Euclidean distortion (WED), can account for perceptual criterion be $C_{\epsilon}(\vec{X^{\alpha}}, \vec{X^{\alpha}}) = (\vec{X^{\alpha}} - \vec{X^{\alpha}})^T W(\vec{X^{\alpha}} - \vec{X^{\alpha}}),$ where, $W = \text{diag}(X_1^{\beta}, X_2^{\beta}, ..., X_K^{\beta}), K = \text{FFT size.}$ $\alpha > 0, \beta < 0$ ♦ Emphasize errors during valleys: $(X^{\alpha} - \hat{X}^{\alpha})^2$ If $X < 1 \rightarrow C_{\epsilon} = \frac{(X^{\alpha} - \hat{X}^{\alpha})^2}{Y^{\beta}}$ $\alpha > 0, \beta > 0$ ♦ Emphasize errors during peaks: If $X > 1 \rightarrow C_{\epsilon} = (X^{\alpha} - \hat{X}^{\alpha})^2 X^{\beta}$ ↑ CASSP

Email: das_amit_ece@yahoo.com, John.Hansen@utdallas.edu

Slide 6



Generalized Spectral Subtraction (GSS) (Sim, et al [1]) model:

$$\hat{X}_k^{\alpha} = a_k Z_k^{\alpha} - b_k E[D_k^{\alpha}]$$

 X_k = magnitude of clean speech estimate

 $Z_k =$ magnitude of noisy speech

 $D_k =$ magnitude of noise

 $a_k, b_k =$ frequency dependent weighting coefficients

 $\alpha =$ spectrum exponent

Optimize *a_k*, *b_k* to minimize MSE: $(\hat{X}_{k}^{\alpha} - X_{k}^{\alpha})^{2}$ GSS-U (Unconstrained) model: $a_{k} \neq b_{k}$ GSS-C (Constrained) model: $a_{k} = b_{k}$

[1] B. L. Sim, Y. C. Tong, J. S. Chang, and C. T. Tan, "A parametric formulation of the generalized spectral subtraction method," IEEE Trans. Speech Audio Process., vol. 6, pp. 328–337, July 1998.

Czech Republic 2011 International Conference on Acoustics, Speech and Signal Proc

Email: das_amit_ece@yahoo.com, John.Hansen@utdallas.edu

Slide 7



 \diamond Propose optimizing a_k, b_k to minimize WED:

WED optimized coefficients:

$$a_{k} = \frac{E[X_{k}^{\beta+2\alpha}]E[X_{k}^{\beta}] - E^{2}[X_{k}^{\beta+\alpha}]}{E[X_{k}^{\beta+2\alpha}]E[X_{k}^{\beta}] - E^{2}[X_{k}^{\beta+\alpha}] + E^{2}[X_{k}^{\beta}]\left(E[D_{k}^{2\alpha}] - E^{2}[D_{k}^{\alpha}]\right)}$$
$$b_{k} = a_{k} - (1 - a_{k})\frac{E[X_{k}^{\beta+\alpha}]}{E[X_{k}^{\beta}]E[D_{k}^{\alpha}]}$$

Substitute these coefficients in the GSS model to form parametric estimators. Offer better flexibility.

 \diamond MSE optimized coefficients = Special case of WED optimized coefficients when $\beta=0$.



Email: das_amit_ece@yahoo.com, John.Hansen@utdallas.edu

 $X^{\beta}(\hat{X}^{\alpha}_{k} - X^{\alpha}_{k})^{2}$







ROVER Using MIXMAX

Three Broad Phoneme Class (BPC) types:

- Sonorants (vowels, nasals, semivowels).
- Obstruents (fricatives, affricates, stops).
- Silence.
- Parametric beta estimators (GSS-BU, GSS-BC) may be tuned to adapt to each BPC.
 - Sonorant estimator Best enhances sonorants only.
- Outputs from each estimator converted to MFCC. Decide weights of each estimator at each frame. Soft combine weights and generate composite utterance. Mechanism is similar to ROVER (Das & Hansen, Interspeech-07).
- Noisy speech can be modeled by the MIXMAX (mixture maximum) model and can be used to classify sonorants/obstruents.



Email: das_amit_ece@yahoo.com, John.Hansen@utdallas.edu

Slide 11





♦ GMMs for sonorants (X), obstruents (Y): MFCCs, N mixtures, K components: $f(\mathbf{x}) = \sum_{i} c_x(i) \prod_k \mathcal{N}(\mu_x(i,k), \sigma_x(i,k))$ $t(\mathbf{y}) = \sum_{j} c_y(j) \prod_k \mathcal{N}(\mu_y(j,k), \sigma_y(j,k))$ ♦ GMM for silence (D): MFCC,1 mixture, K components:

 $g(\mathbf{d}) = \prod_k \mathcal{N}(\mu_d(k), \sigma_d(k))$

 MIXMAX model for noisy speech (Z): $\mathbf{Z} \approx \max(\mathbf{X}, \mathbf{Y}, \mathbf{D})$

[2] A. Nadas, D. Nahamoo, M.A. Picheny, ``Speech recognition using noise-adaptive prototype," IEEE Trans. Speech & Audio Proc., 37(10):1495-1505, Oct. 1989.

Email: das_amit_ece@yahoo.com, John.Hansen@utdallas.edu

Slide 12





- Results are based on 32 tokens from TIMIT test corpus
- Metrics: Segmental-SNR, Itakura-Saito (IS) Distortion
- ♦ GMMs trained from 300 tokens, # mixtures = 16 (sonorants, obstruents),

1 (silence); 39-dim MFCC based GMMs.

- Noise Types: Flat communications channel noise (FLN, mostly stationary), large crowd noise (LCR, mostly non-stationary).
- Acronyms used in figures (next slides):

Baseline	Parametric	ROVER			
MMSE (Ephraim-Malah) [3]	WC (WED Chi) [5]	RWC (ROVER WED Chi) [5]			
JMAP (Wolfe-Godsil Joint- MAP) [4]	JC (JMAP Chi) [5]	RJC (ROVER JMAP Chi) [5]			
GU (Sim et al., GSS Unconstrained)	GBU (GSS-Beta unconstrained)	RGBU (ROVER GBU)			
GC (Sim et al., GSS Constrained)	GBC (GSS-Beta Constrained)	RGBC (ROVER GBC)			

[3] Y. Ephraim, D. Malah, "Speech enhancement using a minimum mean-square error short-time amplitude estimator," IEEE Trans. ASSP, 32(6):443-445, Dec 1984.

[4] P. J. Wolfe, S. J. Godsill, "Efficient alternatives to the Ephraim and Malah suppression rule for audio signal enhancement," *EURASIP J. App. Sig. Process.*, vol. 10.

[5] A. Das. J.H.L. Hansen, "Broad phoneme class based speech enhancement using the Mixture Maximum Model," ICASSP-10, 4762-4765, March 2010.

Email: das_amit_ece@yahoo.com, John.Hansen@utdallas.edu

SegSNR Increase, FLN (OdB SNR) DCRSS								
Baseline(GU) vs Parametric (GBU) vs ROVER (RGBU)								
ſ	Enhancement $\langle \alpha, \beta \rangle$	Rise in SegSNR (dB)						
		Son	Obs	Sil	Ovl	Green:		
Parametric	$GU\langle 1.00, 0.00 \rangle$	4.39	3.36	0.20	3.85	Good Yellow:		
	$\mathrm{GBU}_S\langle 1.50, 1.00 \rangle$	4.85	5.25	0.88	4.79	Reasonable		
	$\mathrm{GBU}_O\langle 1.00, -1.75\rangle$	3.94	7.77	4.23	5.41	Poor		
	$\operatorname{GBU}_N\langle 2.00, -1.75\rangle$	-2.70	7.12	7.45	1.45	Goal:		
	RGBU	5.18	7.89	4.97	6.58	More greens and/or fewer		

GU: Sim et al. Baseline GSS-U estimator

 GBU_S : GSS-BU estimator with best configurable parameters for sonorants GBU_O : GSS-BU estimator with best configurable parameters for obstruents GBU_N : GSS-BU estimator with best configurable parameters for silence RGBU: ROVER GSS-BU estimator

Email: das_amit_ece@yahoo.com, John.Hansen@utdallas.edu

Slide 14



Email: das amit ece@vahoo.com, John.Hansen@utdallas.edu

Slide 15



Email: das amit ece@yahoo.com, John.Hansen@utdallas.edu

Slide 16



	BPC	SegSNR			IS Distortion					
		-5 dB	0 dB	$5 \mathrm{dB}$	10 dB	-5 dB	$0 \mathrm{dB}$	$5 \mathrm{dB}$	10 dB	
Noise FLN Noise	Son 🐥	RGBU	RGBU	RGBU	\mathbf{GC}	RGBU	RGBU	RGBU	RGBU	1st
	Son 🌲	RGBC	RGBC	RGBC	GU	RJC	RGBC	RGBC	RGBC	2nd
	Obs 🐥	RJC	RJC	RJC	RJC	RGBU	RGBU	RGBU	RGBU	
	Obs 🌲	RWC	RWC	RWC	RWC	RWC	RWC	RWC	RWC	
	Sil 🐥	RJC	RJC	RJC	RJC	RGBU	RGBU	RGBU	RGBU	
	Sil 🌲	RWC	RWC	RWC	RWC	RJC	RJC	RJC	RWC	
	Ovl 🐥	RJC	RJC	RJC	RGBU	RGBU	RGBU	RGBU	RGBU	
	Ovl 🏟	RWC	RWC	RGBU	RGBC	RWC	RWC	RWC	RGBC	j
	BPC	SegSNR			IS Distortion					
		-5 dB	0 dB	5 dB	10 dB	-5 dB	0 dB	5 dB	10 dB	
	Son 🌲	RJC	RGBC	GC	GC	RGBC	RGBC	RGBC	RGBC	
	Son 🌲	RGBC	RGBU	\mathbf{GU}	\mathbf{GU}	RGBU	RGBU	RGBU	RGBU	
	Obs 🐥	RJC	RJC	RJC	RJC	RGBU	RGBU	RGBU	RGBU	
CR	Obs 🌲	RWC	RWC	RWC	RGBU	RWC	RWC	RWC	RGBC	
ICASSP	Sil 🐥	RJC	RJC	RJC	RJC	RGBU	RGBU	RGBU	RGBU	
	Sil 🌲	RWC	RWC	RWC	RWC	RJC	RWC	RWC	RJC	
	Ovl 🐥	RJC	RJC	RJC	RGBU	RGBU	RGBU	RGBU	RGBU	
	Ovl 🌲	RWC	RWC	RGBU	RGBC	RWC	RGBC	RGBC	RGBC	
May 22–27, 2011 Prague Czech Republic	2011 International Conference on Acoustics, Speech and Signal Proces	ling							UT	D

Email: das_amit_ece@yahoo.com, John.Hansen@utdallas.edu

Slide 17

conclusions

- \diamond Parametric GSS- β U, GSS- β C estimators.
- Parametric estimators can be pre-tuned per phoneme class but may not perform well across all classes.
- ROVER based paradigm to pick phoneme class segments from parametric pre-tuned estimators and form a single composite utterance.
- ROVER based estimators outperform baseline and parametric estimators across most combinations of FLN/LCR noise types and global SNRs.



Email: das_amit_ece@yahoo.com, John.Hansen@utdallas.edu

Slide 18



May 22-27, 2011



- J. Deller, J.H.L. Hansen, J. Proakis, Discrete Time Processing of Speech Signals, Prentice-Hall Publishers, NY, 2000.
- [2] J.H.L. Hansen, L. Arslan, "Markov model based phoneme class partitioning for improved constrained iterative speech enhancement", IEEE Trans. Speech & Aud. Proc., 3(1):98-104, Jan 1995.
- [3] B. L. Sim, Y. C. Tong, J. S. Chang and C. T. Tan, "A parametric formulation of the generalized spectral subtraction method," IEEE Trans. Speech & Audio Proc., 6(4):328-337, July 1998.
- [4] A. Das, J.H.L. Hansen, "Generalized parametric spectral subtraction using weighted Euclidean distortion," Interspeech-2008, pp. 399-402, Sept. 2008.
- [5] A. Nádas, D. Nahamoo, M.A. Picheny, "Speech recognition using noise-adaptive prototype," IEEE Trans. Speech & Aud. Proc., 37(10):1495-1505, Oct. 1989.
- [6] D. Burshtein, S. Gannot, "Speech enhancement using a mixturemaximum model," IEEE Trans. Speech & Aud. Proc., 10(6):341-351, Sept. 2002.
- [7] A. Das, J.H.L. Hansen, "Broad phoneme class based speech enhancement using mixture maximum model," Proc. IEEE ICASSP, 4762-4765, Mar 2010.

Email: das_amit_ece@yahoo.com, John.Hansen@utdallas.edu

Slide 19