

Speech Information Processing

Final Project

Speech Dereverberation

Based on the paper of: "Speech dereverberation based on variance-normalized delayed linear prediction," by Tomohiro Nakatani, Takuya Yoshioka, Keisuke Kinoshita, Masato Miyoshi, and Biing-Hwang Juang, in "IEEE Transactions on Audio, Speech, and Language Processing, vol. 18, no. 7, pp. 1717-1731, Sep. 2010.

Made By: Rohit Kumar

MTech SP(16100)

1 Fundamentals

If we consider one source, than signal at the m-th microphone at time n is :

$$x_m(n) = \sum_{m=1}^M r_m(n)s(n) + e_m(n)$$

where s(n) is Clean Speech Signal

$r_m(n)$ is Room Impulse response between source and m-th microphone

Now s(k,n) is STFT of clean Speech

where $n \in 1, \dots, N$

where $k \in 1, \dots, K$, then reverberant speech signal at m-th microphone is

$$x_m(k, n) = \sum_{l=0}^{L_h-1} h_m(k, l)s(k, n-l) + e_m(k, n)$$

where the first part: $d_m(k, n)$ is the desired speech which is composed of direct speech component and early reverberation part.

and second term is $r_m(k, n)$ is the late reverberation which we need to remove

$$x_m(k, n) = \sum_{l=0}^{\tau-1} h_m(k, l)s(k, n-l) + \sum_{l=\tau}^{L_g-1} h_m(k, l)s(k, n-l)$$

Considering the response of all M microphone

$$d_m(k, n) = \sum_{i=1}^M \sum_{l=0}^{L_g-1} x_i(k, n-\tau-l)g_{m,i}(k, n-l) + x_m(k, n)$$

$$D(k) = X(k) - X_{\tau}(k)G(k),$$

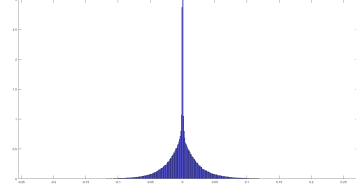


Figure 1: Histogram of reverb speech

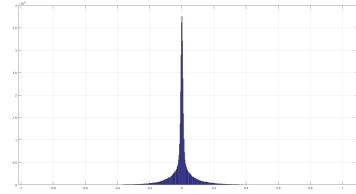


Figure 2: Histogram of original speech

where $D(k) = [d_1(k), \dots, d_M(k)] \in (N * M)$
 $d_m(k) = [d_m(k, 1), \dots, d_m(k, N)]' \in (N * 1)$
 $X(k) = [x_1(k), \dots, x_M(k)] \in (N * M)$
 $x_m(k) = [x_m(k, 1), \dots, x_m(k, N)]' \in (N * 1)$
 $X_{\tau}(k) = [X_{\tau, 1}(k), \dots, X_{\tau, M}(k)] \in (N * ML_g)$ convolution matrix
 $G(k) = [g_1(k), \dots, g_M(k)] \in (ML_g * M)$

2 Concept of Sparsity

Idea Behind for Using Sparse Concept in Dereverberation

1. Aim: to obtain sparse residual rather than a residual with a min variance.

2. L1 norm outperform the traditional L2 norm for Linear Prediction Algo.

3.Originally L2 norm \Rightarrow for an excitation signal with iid and Gaussian assumptions.

4.However speech is quasi periodic.

5.Therefor L2 norms suffers from an overemphasis on peaks, which is avoided by using l1 norm which give less emphasis on the outliers of the spiky excitation associated with speech.

2.1 Maths behind the idea

The prediction coefficient matrix G is solved by

$$G = \operatorname{argmin}_G ||X - X||_{pq}^q + \alpha ||G||_{rs}^s$$

Different type of p,q,r,s will result is different solution.

In my experiment the form of the equation is (when cvx toolbox is used)

$$G = \operatorname{argmin}_G ||X - X||_1 + \alpha ||G||_1$$

while using the IRLS, value of

$$\alpha = 0$$

$$G = \operatorname{argmin}_G ||X - X||_{22}^2$$

3 Algorithm:For dereverberation of speech

for each k

input:Reverberant speech $x_{nk}^m \forall n, m$

set parameters: $\tau, L_g, \epsilon, \delta$

initialize the variances $\lambda_{nk} \Leftarrow |x_{nk}^1|^2$

repeat

$$A_k \Leftarrow \sum_{n=1}^N \frac{(x_n - \tau_k)(x_n^H - \tau_k^H)}{\lambda_{nk}}$$

$$b_k \Leftarrow \sum_{n=1}^N \frac{(x_n - \tau_k)(x_{nk})^*}{\lambda_{nk}}$$

$$g_k \Leftarrow (pinv A)_k b_k$$

$$d_{nk} \Leftarrow (g^H)_k x_n - \tau_k \lambda_{nk} \Leftarrow \max |d_{nk}|^2, \epsilon$$

until g_k converges as $||g^{(i+1)}_k - g^{(i)}_k|| / ||g^{(i)}_k|| < \delta$, pr a maximum number of iterations

4 Algorithm:Dereverberation using the Sparsity Constraint

for each k input:Reverberant speech $x_{nk}^m \forall n, m$

set parameters: τ, L_g, ϵ

initialize the desired signal matrix $D = x$ and the correlation structure $\theta = I$

repeat

$$w^{(i)}_n \Leftarrow ||d^{(i-1)}_n + \epsilon||_1$$

$$G^{(i)} \Leftarrow (X^H \tau W(i) X \tau) X^H \tau W(i) X$$

$$D \Leftarrow X - X \tau G$$

until g_k converges or a maximum number of iterations

5 Experimental Setup

1.Dataset: The Dataset has been taken from Reverb Challenge Site, where in the project tackle the problem all type of reverberation that can occur namely very far field ,medium field and near the microphone(considered both male and female.).

2.Sampling Frequency=16Khz

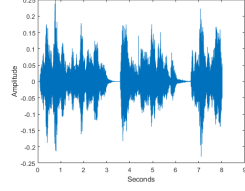


Figure 3: Reverberated Speech Signal

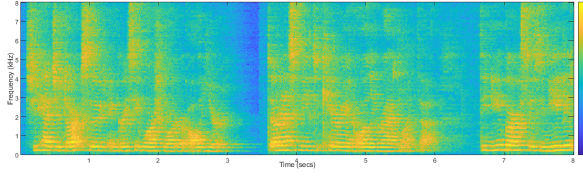


Figure 4: spectrogram of Reverberated sound

3.Window Size=512 sample

4.Window Shift=256

5.Window Used=Hann window

6. $\tau=2$ samples

7. $L_g=30$

8. $\epsilon=1e-4$

9.max iterations=2

10.CVX toolbox is used for obtaining the L1 minimisation of the prediction coefficient.

11. $\alpha = 0.1$

6 Observation and Results

Figure 3 and 4 of the reverberated speech, figure 5 and 6 of the de reverberated speech, and the figure 7 and 8 is the spectrogram of reverberated speech and histogram of sparse residual de reverb speech

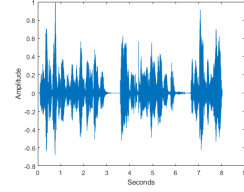


Figure 5: Dereverbrated SPEech Signal

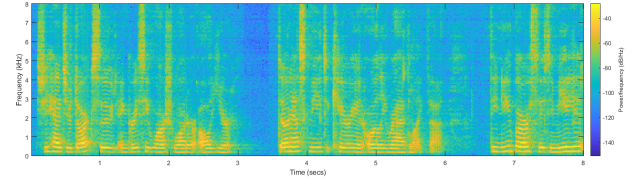


Figure 6: Dereverbrated SPEech spectrogram

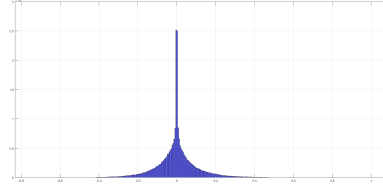


Figure 7: Sparse Signal Histogram

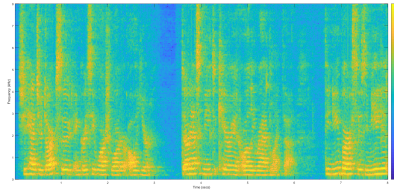


Figure 8: Sparse Signal Spectrogram

7 Performance Analysis

group sparse Linear Prediction by Daniele Giacobello and Tobias Lindstrøm Jensen.

7.1 Performance of Dereverb sound using WPE

Performance Parameters	Reverb Speech	Dereverb Speech
1.Cepstral Distance	2.72	1.23
2.Spectral Distance	1.062	0.96
3.Frequency Weighted SNR	8.29	32.02
4.LPC based Log Likelihood Ratio	0.416	0.0182

7.2 Performance of Dereverb sound using sparsity constraint in wpe

Performance Parameters	Reverb Speech	Dereverb Speech
1.Cepstral Distance	2.72	2.017
2.Spectral Distance	1.062	0.727
3.Frequency Weighted SNR	8.29	13.05
4.LPC based Log Likelihood Ratio	0.416	0.2099

8 References

- 1.Tomohiro Nakatani, Takuya Yoshioka, Keisuke Kinoshita, Masato Miyoshi, and Biing-Hwang Juang, "Speech dereverberation based on variance-normalized delayed linear prediction," IEEE Transactions on Audio, Speech, and Language Processing, vol. 18, no. 7, pp. 1717-1731, Sep. 2010.
- 2.P. C. Loizou, Speech Enhancement: Theory and Practice. CRC Press, 2013.
- 3.P. Naylor and N. D. Gaubitch, Speech Dereverberation. Springer Science Business Media, 2010.
- 4.Speech Dereverberation Based on Convex Optimisation Algorithm for