

DELFT UNIVERSITY OF TECHNOLOGY

DIGITAL AUDIO AND SPEECH PROCESSING  
IN4182

---

## Course Project

---

*Author:*  
Edoardo Focante (5377544)

August 1, 2021



## Part I

# Single Microphone Speech-Enhancement

## 1 Overview

The purpose of this project is to develop a single-channel speech enhancement system for noise reduction. In order to do this, the system will operate through different blocks. First of all, the signal will be divided into frames, so that in every frame the signal can be considered stationary. The signal frames are then taken with 50% overlap and multiplied with a Hanning window of proper length before further processing. After that, the  $\mathbf{Y}$  matrix is constructed, by taking the FFT of every frame as one of its column.

After that, the system will estimate the PSD of the noise and the signal, and use those values to compute a gain function. After applying the gain to every  $Y_k(l)$ , the de-noised signal is reconstructed by applying an IFFT on each column of  $\mathbf{Y}$  and summing through overlap and add. Due to the properties of the Hanning window, no de-windowing will be necessary.

The final product of the system will hopefully be a less noisy version of the input signal. The improvement will be evaluated via different metrics and through human listening.

## 2 Parameters Description

The data used for experiments is the one provided with the assignment description.

The signals used have been sampled at 16 KHz while the range of stationarity of speech can be taken between 20 and 30 ms. This means that the proper frame length should be taken between 320 and 480 samples. To help with the computation of the FFT, the frame length should also be picked as a power of 2. There is no power of 2 in the desired interval, but the closest one is 512, which has been picked by stretching the restriction a bit.

Due to the fact that the absolute value of the FFT of a real signal is symmetric with respect to the central frequency bin, only one half of the FFT bins are used, in order to reduce computations. The whole  $\mathbf{Y}$  matrix is then reconstructed in view of this property before applying the final IFFT and the overlap and add.

In order to estimate the signal power, a Bartlett periodogram is used. This helped reduce the variance of the classic periodogram through a smoothing over  $M = 32$  Frames.

In order to estimate the noise PSD, instead, a Minimum Statistics approach is used. The noise PSD estimate for a specific frame  $l$  is chosen as the minimum signal PSD over the last  $D$  frames. To reduce the error related to the variance of the signal PSD an exponentially smoothed version of the Periodogram is used. The expression for this smoothing is

$$\hat{P}_{YY,k}(l) = \alpha \hat{P}_{YY,k}(l-1) + (1-\alpha)|Y_k(l)|^2$$

The value of alpha is chosen as suggested in [1] such that

$$\alpha = \frac{T_{SM}f_s/R - 1}{T_{SM}f_s/R + 1} \quad (1)$$

By choosing  $T_{SM} = 0.2s$ ,  $f_s = 16KHz$  and  $R = \frac{L}{2} = 256$ , expression 1 gives a value of  $\alpha$  close to 0.85. Note that according to [1] the optimal  $\alpha$  should be optimized and then updated in order to avoid unwanted effects such as the broadening of peaks in the speech signal. However, for the sake of simplicity,  $\alpha$  will be treated as a constant. In addition to that, by taking the minimum over different frames, the Minimum Statistics estimate is generally biased, and hence has to be corrected accordingly. However, for simplicity, the self-written code will not apply a compensation and will hence work with a slightly biased estimate.

As already mentioned, to achieve our noise PSD estimate via minimum statistics, we will take the minimum signal PSD value over  $D$  frames. With smaller windows the noise estimation will be able to track variations in the noise level faster, but we cannot reduce it too much because we still have to ensure that the window includes a frame where there is no speech, so that the minimum is originated from a frame with noise only. Because of this reason, we would like to have a window between 1 to 2 seconds long, such that the presence of a frame without speech is ensured. Over different tests, smaller values of  $D$  such as 8 and 16 proved to be the best in terms of reduction of the mean squared error (MSE) between the clean and noisy signal. However, these values don't respect our window length constraint, and in fact slightly worsen the intelligibility. The final value of  $D$  is chosen equal to 100, so that the minimum is taken over a time equivalent of  $100/2 \times 32ms = 1.6s$ .

After performing our Minimum Statistics estimation, we compute our Wiener gain as

$$G_k(l) = 1 - \frac{P_{NN,k}(l)}{P_{YY,k}^B(l)} \quad (2)$$

After applying the gain to  $\mathbf{Y}$ , we reconstruct the full matrix through its symmetry properties, and by applying IFFT and overlap and add we get our de-noised signal in the time domain.

### 3 Evaluation

To quantify the quality of the results the result we will use two different metrics. First we will compare the MSE before and after denoising. Since a better MSE doesn't necessarily correspond to an improvement in intelligibility, we will also use the short-time objective intelligibility (STOI) measure [2] [3]. This measure was developed with the purpose of measuring intelligibility without relying on human listeners. These type of metrics are called Objective Intelligibility Measures (OIM) and are very much needed to avoid and replace human hearing tests, which are expensive in terms of money and time. Without going too in detail, the STOI measure is based on calculating the correlation coefficient between the norm of the DFT-bins of the clean and noisy signal, averaged over different frames and different frequency octaves in the speech frequency range. Based on the results presented in [2] and [3], it can be said that this metric has a high correlation with speech intelligibility evaluated from human tests. Hence, it can be considered a reliable way to quantify intelligibility. The code for computing the STOI measure has not been self written but has been borrowed from the author of [2] and [3], C. Taal.

In addition to that, remember that our self-implemented version of the Minimum Statistics estimation doesn't take into account the bias and doesn't optimize  $\alpha$ . In order to highlight the limitations of our code and to consolidate our results, we also decided to run a more accurate implementation of the Minimum Statistics, wrote by N. Jakovljević [4], as a reference.

### 4 Results

The final results will include the tests ran using two different speech samples and 4 different types of noise: Gaussian noise, Artificial non-stationary noise, Babble noise and Speech shaped noise. In addition to that, three SNR levels will be tested: -5dB, 0dB and 10DdB. Tables 1, 2 and 3 show the respective results for the different SNRs values. The tests have been carried out using both the speech signals provided, but since changing audio file led to very similar results, only the results achieved using the first speech sample are included in the report, to keep it more compact.

As we can see from all of the tables, the system is able to consistently reduce the MSE between the noisy and clean speech. Under a MSE perspective, the noise reduction works nicely for all the 4 noise types, but struggles a bit more with the non-stationary noise. This can be due to the fact that the noise power varies in the window of  $D$  frames we considered. Due to the fact that a minimum is taken, the Minimum Statistics approach is slow to track increases in noise, which could be making our estimate slightly inaccurate when compared to the cases where other types of noise are present.

By looking at the STOI metric, instead, we cannot see much improvement in the intelligibility of the de-noised signal. This is plausible, as the Wiener gain optimises the MSE, which is not necessarily correlated with the intelligibility.

A personal listening evaluation confirmed what the STOI metric conveys, as the voice in the processed signal is slightly more distinguishable from the noise but is also slightly corrupted by the de-noising process, which leads to barely noticeable improvements of the intelligibility only in the cases with the lowest SNR.

This result is not as surprising as it may seem, since many tests have already pointed out that single-microphone systems struggle in improving intelligibility, as can be read in [5] and [6].

Finally, to determine whether ignoring the bias and the optimization of  $\alpha$  led to a worsened performance we repeated the same tests using the code provided in [4]. The results proven to be similar to the one reported in tables 1, 2 and 3 and have not been included for brevity. Overall, the usage of the more complex algorithm brought only barely significant and not consistent improvements, which were not relevant enough to justify implementing our own unbiased noise estimate.

In conclusion, the implemented system didn't prove to be too effective in increasing intelligibility, but was successful in reducing the mean squared error between the clean and noisy signal under all the considered noise scenarios.

Noise Type	Noisy MSE	Processed MSE	Noisy STOI	Processed STOI
Gaussian Noise	$9.2 \times 10^{-3}$	$2.6 \times 10^{-3}$	0.544	0.556
Artificial Non-stationary Noise	$9.2 \times 10^{-3}$	$5.0 \times 10^{-3}$	0.507	0.497
Babble Noise	$9.2 \times 10^{-3}$	$4.1 \times 10^{-3}$	0.511	0.506
Speech-shaped Noise	$9.2 \times 10^{-3}$	$2.8 \times 10^{-3}$	0.498	0.503

Table 1: Evaluation of self-implemented noise reduction for SNR=-5 dB

Noise Type	Noisy MSE	Processed MSE	Noisy STOI	Processed STOI
Gaussian Noise	$2.9 \times 10^{-3}$	$9.4 \times 10^{-4}$	0.647	0.656
Artificial Non-stationary Noise	$2.9 \times 10^{-3}$	$1.6 \times 10^{-3}$	0.619	0.616
Babble Noise	$2.9 \times 10^{-3}$	$1.3 \times 10^{-3}$	0.610	0.609
Speech-shaped Noise	$2.9 \times 10^{-3}$	$1.1 \times 10^{-3}$	0.611	0.618

Table 2: Evaluation of self-implemented noise reduction for SNR=0 dB

Noise Type	Noisy MSE	Processed MSE	Noisy STOI	Processed STOI
Gaussian Noise	$2.9 \times 10^{-4}$	$1.4 \times 10^{-4}$	0.844	0.850
Artificial Non-stationary Noise	$2.9 \times 10^{-4}$	$1.9 \times 10^{-4}$	0.824	0.826
Babble Noise	$2.9 \times 10^{-4}$	$1.6 \times 10^{-4}$	0.783	0.785
Speech-shaped Noise	$2.9 \times 10^{-4}$	$1.6 \times 10^{-4}$	0.821	0.827

Table 3: Evaluation of self-implemented noise reduction for SNR=10 dB

# Appendices

## A MATLAB Code

```

1 % DASP Single Channel
2 clear;
3 clear sound;
4 rng(1)
5 [x, fs]=audioread('clean_speech.wav');
6 [noise, fs]=audioread('babble_noise.wav');
7 fplot=0; %plot flag
8
9 SNR=10^(10/10); %linear SNR
10 %noise=randn(length(x),1); % Gaussian Noise
11 noise=noise(1:length(x));
12 comp=0; %More complex algorithm flag (bias compensation and optimized \alpha)
13
14 SNRr=norm(x).^2/norm(noise).^2;
15 noise=noise*sqrt(SNRr/SNR);
16 % NP=norm(noise).^2
17 % SP=norm(x).^2
18 y=x+noise;
19 % sound(y, fs);
20 %% framing
21
22
23 K=512; %512 samples = 32ms with fs=16Khz
24 overlap=0.5;
25 Y=zeros(K, floor(length(y)/K));
26 i=1;
27 for l=1:K*(1-overlap):length(y)-K
28     Y(:, i)=(fft(y(l:l+K-1).*(hann(K))))';
29     i=i+1;
30 end
31 L=size(Y,2);

```

```

32 K=K/2+1;
33 Y=Y(1:K,:);
34
35 %% SIGNAL PSD ESTIMATE
36 Pyy=zeros(K,L);
37 Pyy=abs(Y).^2;
38 Pyyb=Py;
39 M=32; % bartlett average length 0.2ms @16Khz
40
41 %Bartlett averaging
42 %Transient
43 for i=1:M
44     Pyyb(:,i)=1/i*sum(Pyy(:,1:i),2);
45 end
46 for i=M:L %Bartlett averaging
47     Pyyb(:,i)=1/M*sum(Pyy(:,(i-M+1):i),2);
48 end
49
50 %Exponential smoothing
51 Pyyex=Py;
52 alfa=0.85;
53
54 if alfa~1
55     for i=2:L
56         Pyyex(:,i)=alfa*Pyyex(:,i-1)+(1-alfa)*Py(:,i);
57     end
58 end
59
60 if fplot
61     figure;
62     hold on
63     plot(Pyyex(100,:));
64     plot(Pyy(100,:));
65     legend('Py smoothed','Py')
66 end
67
68 %% Minimum Statistics (No bias correction)
69
70
71 D=100; %100*32/2*ms =1600 ms
72 Pnn=zeros(K,L);
73
74
75 for i=1:D
76     Pnn(:,i)=min(Pyyex(:,1:i),[],2);
77 end
78 for i=D:L
79     Pnn(:,i)=min(Pyyex(:,i-D+1:i),[],2);
80 end
81
82
83
84 %Plotting
85 if fplot
86     figure;
87     hold on;
88     plot(Pyyex(100,:));
89     plot(Pnn(100,:));
90     legend('Py','Pnn')
91 end
92

```

```

93 %% Complete Minimum Statistics
94
95 if comp==1
96 Pnn=noise_est_min_stat(Pyy,D);
97 end
98 %% GAIN CALCULATION
99 G=zeros(K,L);
100
101 G=1-(Pnn./Pyyb);
102 Yout=G.*Y;
103
104 Yout=[Yout; conj(flip(Yout(2:end-1,:)))]; %Reconstruct full Y through simmetry
105
106 %% IFFT
107 S=ifft(Yout);
108
109 %% Overlap and add
110 s=zeros(length(y),1);
111 i=1;
112 K=2*K-2;
113 for l=1:K*(1-overlap):length(y)-K
114     s(l:l+K-1)=s(l:l+K-1)+S(:,i);
115     i=i+1;
116 end
117
118 %% Output and Evaluation
119 clear sound;
120
121 % sound(s,fs);
122 MSEpre=mean((y-x).^2);
123 MSEpost=mean((s-x).^2);
124 STOIpre=stoi(x,y,fs);
125 STOIpost=stoi(x,s,fs);
126
127 function d = stoi(x, y, fs_signal)
128 % d = stoi(x, y, fs_signal) returns the output of the short-time
129 % objective intelligibility (STOI) measure described in [1, 2], where x
130 % and y denote the clean and processed speech, respectively, with sample
131 % rate fs_signal in Hz. The output d is expected to have a monotonic
132 % relation with the subjective speech-intelligibility, where a higher d
133 % denotes better intelligible speech. See [1, 2] for more details.
134 %
135 % References:
136 % [1] C.H.Taal, R.C.Hendriks, R.Heusdens, J.Jensen 'A Short-Time
137 % Objective Intelligibility Measure for Time-Frequency Weighted Noisy
138 % Speech', ICASSP 2010, Texas, Dallas.
139 %
140 % [2] C.H.Taal, R.C.Hendriks, R.Heusdens, J.Jensen 'An Algorithm for
141 % Intelligibility Prediction of Time-Frequency Weighted Noisy Speech',
142 % IEEE Transactions on Audio, Speech, and Language Processing, 2011.
143 %
144 %
145 % Copyright 2009: Delft University of Technology, Signal & Information
146 % Processing Lab. The software is free for non-commercial use. This program
147 % comes WITHOUT ANY WARRANTY.
148 %
149 %
150 %
151 % Updates:
152 % 2011-04-26 Using the more efficient 'taa_corr' instead of 'corr'
153

```

```

154 if length(x)~=length(y)
155     error('x and y should have the same length');
156 end
157
158 % initialization
159 x = x(:); % clean speech column vector
160 y = y(:); % processed speech column vector
161
162 fs = 10000; % sample rate of proposed
163     intelligibility measure
164 N_frame = 256; % window support
165 K = 512; % FFT size
166 J = 15; % Number of 1/3 octave bands
167 mn = 150; % Center frequency of first 1/3
168     octave band in Hz.
169 H = thirdoct(fs, K, J, mn); % Get 1/3 octave band matrix
170 N = 30; % Number of frames for
171     intermediate intelligibility measure (Length analysis window)
172 Beta = -15; % lower SDR-bound
173 dyn_range = 40; % speech dynamic range
174
175 % resample signals if other samplerate is used than fs
176 if fs_signal ~= fs
177     x = resample(x, fs, fs_signal);
178     y = resample(y, fs, fs_signal);
179 end
180
181 % remove silent frames
182 [x y] = removeSilentFrames(x, y, dyn_range, N_frame, N_frame/2);
183
184 % apply 1/3 octave band TF-decomposition
185 x_hat = stdft(x, N_frame, N_frame/2, K); % apply short-time DFT
186     to clean speech
187 y_hat = stdft(y, N_frame, N_frame/2, K); % apply short-time DFT
188     to processed speech
189
190 x_hat = x_hat(:, 1:(K/2+1)).'; % take clean single-sided
191     spectrum
192 y_hat = y_hat(:, 1:(K/2+1)).'; % take processed single-sided
193     spectrum
194
195 X = zeros(J, size(x_hat, 2)); % init memory for clean speech
196     1/3 octave band TF-representation
197 Y = zeros(J, size(y_hat, 2)); % init memory for processed
198     speech 1/3 octave band TF-representation
199
200 for i = 1:size(x_hat, 2)
201     X(:, i) = sqrt(H*abs(x_hat(:, i)).^2); % apply 1/3 octave bands as
202         described in Eq.(1) [1]
203     Y(:, i) = sqrt(H*abs(y_hat(:, i)).^2);
204 end
205
206 % loop al segments of length N and obtain intermediate intelligibility measure
207     for all TF-regions
208 d_interm = zeros(J, length(N:size(X, 2)));
209     % init memory for intermediate intelligibility measure
210 c = 10^(-Beta/20); %
211     constant for clipping procedure
212
213 for m = N:size(X, 2)
214     X_seg = X(:, (m-N+1):m);

```

```

202     % region with length N of clean TF-units for all j
    Y_seg = Y(:, (m-N+1):m);
203     % region with length N of processed TF-units for all j
    alpha = sqrt(sum(X_seg.^2, 2)./sum(Y_seg.^2, 2)); %
204     % obtain scale factor for normalizing processed TF-region for all j
    aY_seg = Y_seg.*repmat(alpha, [1 N]);
205     % obtain \alpha*Y_j(n) from Eq.(2) [1]
206     for j = 1:J
        Y_prime = min(aY_seg(j, :), X_seg(j, :)+X_seg(j, :)*c); %
207         % apply clipping from Eq.(3)
        d_interm(j, m-N+1) = taa_corr(X_seg(j, :).', Y_prime(:)); %
208         % obtain correlation coefficient from Eq.(4) [1]
    end
209 end
210
211 d = mean(d_interm(:)); %
    % combine all intermediate intelligibility measures as in Eq.(4) [1]
212
213 %%
214 function [A cf] = thirdoct(fs, N_fft, numBands, mn)
215 % [A CF] = THIRDOCT(FS, N_FFT, NUMBANDS, MN) returns 1/3 octave band matrix
216 % inputs:
217 % FS: samplerate
218 % N_FFT: FFT size
219 % NUMBANDS: number of bands
220 % MN: center frequency of first 1/3 octave band
221 % outputs:
222 % A: octave band matrix
223 % CF: center frequencies
224
225 f = linspace(0, fs, N_fft+1);
226 f = f(1:(N_fft/2+1));
227 k = 0:(numBands-1);
228 cf = 2.^(k/3)*mn;
229 fl = sqrt((2.^(k/3)*mn).*2.^(((k-1)/3)*mn));
230 fr = sqrt((2.^(k/3)*mn).*2.^(((k+1)/3)*mn));
231 A = zeros(numBands, length(f));
232
233 for i = 1:(length(cf))
234     [a b] = min((f-fl(i)).^2);
235     fl(ii) = f(b);
236     fl_ii = b;
237
238     [a b] = min((f-fr(i)).^2);
239     fr(ii) = f(b);
240     fr_ii = b;
241     A(i, fl_ii:(fr_ii-1)) = 1;
242 end
243
244 rnk = sum(A, 2);
245 numBands = find((rnk(2:end)>=rnk(1:(end-1))) & (rnk(2:end)~=0)~=0, 1, 'last')+1;
246 A = A(1:numBands, :);
247 cf = cf(1:numBands);
248
249 %%
250 function x_stdft = stdft(x, N, K, N_fft)
251 % X_STDFT = STDFT(X, N, K, N_FFT) returns the short-time
252 % hanning-windowed dft of X with frame-size N, overlap K and DFT size
253 % N_FFT. The columns and rows of X_STDFT denote the frame-index and
254 % dft-bin index, respectively.

```



```

255
256 frames      = 1:K:(length(x)-N);
257 x_stdft     = zeros(length(frames), N_fft);
258
259 w            = hanning(N);
260 x            = x(:);
261
262 for i = 1:length(frames)
263     ii        = frames(i):(frames(i)+N-1);
264     x_stdft(i, :) = fft(x(ii).*w, N_fft);
265 end
266
267 %%
268 function [x_sil y_sil] = removeSilentFrames(x, y, range, N, K)
269 % [X_SIL Y_SIL] = REMOVESILENTFRAMES(X, Y, RANGE, N, K) X and Y
270 % are segmented with frame-length N and overlap K, where the maximum energy
271 % of all frames of X is determined, say X_MAX. X_SIL and Y_SIL are the
272 % reconstructed signals, excluding the frames, where the energy of a frame
273 % of X is smaller than X_MAX-RANGE
274
275 x      = x(:);
276 y      = y(:);
277
278 frames = 1:K:(length(x)-N);
279 w      = hanning(N);
280 msk     = zeros(size(frames));
281
282 for j = 1:length(frames)
283     jj      = frames(j):(frames(j)+N-1);
284     msk(j)   = 20*log10(norm(x(jj).*w)/sqrt(N));
285 end
286
287 msk      = (msk-max(msk)+range)>0;
288 count    = 1;
289
290 x_sil    = zeros(size(x));
291 y_sil    = zeros(size(y));
292
293 for j = 1:length(frames)
294     if msk(j)
295         jj_i      = frames(j):(frames(j)+N-1);
296         jj_o      = frames(count):(frames(count)+N-1);
297         x_sil(jj_o) = x_sil(jj_o) + x(jj_i).*w;
298         y_sil(jj_o) = y_sil(jj_o) + y(jj_i).*w;
299         count      = count+1;
300     end
301 end
302
303 x_sil = x_sil(1:jj_o(end));
304 y_sil = y_sil(1:jj_o(end));
305
306 %%
307 function rho = taa_corr(x, y)
308 % RHO = TAA_CORR(X, Y) Returns correlation coefficient between column
309 % vectors x and y. Gives same results as 'corr' from statistics toolbox.
310 xn      = x-mean(x);
311 xn      = xn/sqrt(sum(xn.^2));
312 yn      = y-mean(y);
313 yn      = yn/sqrt(sum(yn.^2));
314 rho     = sum(xn.*yn);
315

```

```

316
317 function [sigma_N_2, P, alpha] = noise_est_min_stat(abs_Y_2, param)
318 % Estimate noise spectrum using minimum statistics as proposed in [1]
319 %
320 % Inputs:
321 %   abs_Y_2 – noisy signal PSD (one column per frame)
322 %   param – parameters of the algorithm
323 %       + D – number of frames in analysing window for minimum search
324 %       + V – number of frames in analysing subwindow
325 %       + tshift – frame shift in seconds
326 %       + tdecay – decay time from signal peak to the noise level
327 %       + alpha_max – maximum value of smoothing factor
328 %       + alpha_min – minimum value of smoothing factor
329 %       + alpha_c_min – minimum value of alpha_c
330 %       + min_y_energy – minimum energy
331 %       + inv_q_eq_max – maximum value of  $Q_{eq}^{-1}$ 
332 %       + beta_max / maximum value of beta
333 %       + noise_slope_max_table – for comparison of  $Q^{-1}$ 
334 %       + a_v – factor used in  $B_c$  calculation
335 %
336 % Outputs:
337 %   sigma_N_2 – estimated noise PSD
338 %   P – smoothed periodogram
339 %   alpha – optimal smoothing factor
340
341 % [1] Martin, R.
342 %   "Noise power spectral density estimation based on optimal smoothing
343 %   and minimum statistics"
344 %   IEEE Transactions on Speech and Audio Processing, 2001, 9, 504–512
345
346 % author: jakovnik@gmail.com
347 % date: 2016/08/23
348
349 if or(nargin < 1, nargin > 2)
350     error('invalid number of input arguments');
351 end
352 %% Set default values
353 if nargin == 1
354     param = [];
355 end
356 D = 96;
357 if isfield(param, 'D')
358     D = param.D;
359 end
360 V = 12;
361 if isfield(param, 'V')
362     V = param.V;
363 end
364 tshift = 16e-3;
365 if isfield(param, 'tshift')
366     tshift = param.tshift;
367 end
368 tdecay = 64e-3;
369 if isfield(param, 'tdecay')
370     tdecay = param.tdecay;
371 end
372 alpha_max = 0.96;
373 if isfield(param, 'alpha_max')
374     alpha_max = param.alpha_max;
375 end
376 alpha_min = 0.30;

```

```

377 if isfield(param, 'alpha_min')
378     alpha_min = param.alpha_min;
379 end
380 alpha_c_min = 0.7;
381 if isfield(param, 'alpha_c_min')
382     alpha_c_min = param.alpha_c_min;
383 end
384 min_y_energy = 1e-9;
385 if isfield(param, 'min_y_energy')
386     min_y_energy = param.min_y_energy;
387 end
388 inv_q_eq_max = 1/2;
389 if isfield(param, 'inv_q_eq_max')
390     inv_q_eq_max = param.inv_q_eq_max;
391 end
392 beta_max = 0.8;
393 if isfield(param, 'beta_max')
394     beta_max = param.beta_max;
395 end
396 noise_slope_max_table = [0.03 8;
397                           0.05 4;
398                           0.06 2];
399 if isfield(param, 'noise_slope_max_table')
400     noise_slope_max_table = param.noise_slope_max_table;
401 end
402 a_v = 2.12;
403 if isfield(param, 'a_v')
404     a_v = param.a_v;
405 end
406 MAX_VALUE = 1e9;
407
408 U = round(D/V);
409 M_D = mean_of_the_min_table3(D);
410 M_V = mean_of_the_min_table3(V);
411 [L, n_frames] = size(abs_Y_2);
412 snr_exp = -tshift/tdecay;
413
414 % initialize
415 sigma_N_2 = zeros(L, n_frames);
416 P = zeros(L, n_frames);
417 alpha = zeros(L, n_frames);
418 % calculation
419 buff_idx = 1;
420 subwc = V;
421 P_prev = abs_Y_2(:, 1);
422 alpha_c_prev = 1; %set initial value (it is suppose noise at begin)
423 sigma_N_2_prev = abs_Y_2(:, 1);
424 sigma_N_2(:, 1) = abs_Y_2(:, 1);
425 mean_P = abs_Y_2(:, 1);
426 mean_P_2 = mean_P.^2;
427 actmin_buff = MAX_VALUE*ones(L, U);
428 actmin = MAX_VALUE*ones(L, 1);
429 actmin_sub = MAX_VALUE*ones(L, 1);
430 lmin_flag = zeros(L, 1);
431 for lambda = 1:n_frames
432     % compute the smoothing parameter alpha
433     y_energy = sum(abs_Y_2(:, lambda));
434     p_energy = sum(P_prev);
435     if y_energy == 0
436         y_energy = min_y_energy;
437     end

```

```

438 alpha_c = 1/(1 + (p_energy/y_energy - 1)^2); %eq. 9
439 alpha_c = 0.7*alpha_c_prev + 0.3*max(alpha_c, alpha_c_min); %eq. 10
440 alpha_c_prev = alpha_c;
441 alpha_opt = alpha_max * alpha_c ./...
442 (1 + (P_prev./sigma_N_2_prev - 1).^2); %eq. 11
443 snr = p_energy/sum(sigma_N_2_prev);
444 alpha_opt = max(alpha_opt, min(alpha_min, snr^snr_exp)); %eq. 12
445 alpha(:, lambda) = alpha_opt;
446 % compute smoothed power P(\lambda, k)
447 P(:, lambda) = alpha_opt .* P_prev + (1 - alpha_opt) ...
448 .*abs_Y_2(:, lambda); %eq. 4
449 % compute bias correction
450 beta = min(alpha_opt.^2, beta_max);
451 mean_P = beta.*mean_P + (1-beta).*P(:, lambda); %eq. 20
452 mean_P_2 = beta.*mean_P_2 + (1-beta).*P(:, lambda).^2; %eq. 21
453 var_P = mean_P_2 - mean_P.^2; %eq. 22
454 inv_Q_eq = var_P./(2*sigma_N_2_prev.^2); %eq. 23
455 %inv_Q_eq = max(min(inv_Q_eq, INV_Q_EQ_MAX), INV_Q_EQ_MIN/lambda);
456 inv_Q_eq = max(min(inv_Q_eq, inv_q_eq_max), 1/20/lambda);
457 tilda_Qeq_D = (1./inv_Q_eq - 2*M_D)/(1-M_D); %eq. 16
458 tilda_Qeq_V = (1./inv_Q_eq - 2*M_V)/(1-M_V); %eq. 16
459 B_min = 1 + 2*(D-1)./tilda_Qeq_D; %eq. 17
460 B_min_sub = 1 + 2*(V-1)./tilda_Qeq_V; %eq. 17
461 % compute inv Q_eq
462 inv_Q_eq = mean(inv_Q_eq);
463 % minimum search
464 B_c = 1 + a_v*sqrt(inv_Q_eq);
465 P_x_B_min_x_B_c = P(:, lambda).*B_min*B_c;
466 k_mode = P_x_B_min_x_B_c < actmin;
467 actmin(k_mode) = P_x_B_min_x_B_c(k_mode);
468 actmin_sub(k_mode) = P(k_mode, lambda).*B_min_sub(k_mode)*B_c;
469 if subwc == V
470     lmin_flag(k_mode) = 0;
471     actmin_buff(:, buff_idx) = actmin;
472     if buff_idx == U
473         buff_idx = 1;
474     else
475         buff_idx = buff_idx + 1;
476     end
477 P_min_u = min(actmin_buff, [], 2);
478 p = find(inv_Q_eq < noise_slope_max_table(:, 1));
479 if isempty(p)
480     noise_slope_max = 1.2;
481 else
482     noise_slope_max = noise_slope_max_table(p(1), 2);
483 end
484 p = and(and(lmin_flag, actmin_sub < noise_slope_max.*P_min_u), ...
485         actmin_sub > P_min_u);
486 if any(p)
487     P_min_u(p) = actmin_sub(p);
488     actmin(p) = actmin_sub(p);
489 end
490 lmin_flag = zeros(L, 1);
491 subwc = 1;
492 actmin(:) = MAX_VALUE;
493 actmin_sub(:) = MAX_VALUE;
494 sigma_N_2(:, lambda) = P_min_u;
495 else
496     if subwc > 1
497         lmin_flag(k_mode) = 1;
498         sigma_N_2(:, lambda) = min(actmin_sub, P_min_u);

```

```

499     P_min_u = sigma_N_2(:,lambda);
500     else
501         sigma_N_2(:,lambda) = min(actmin_sub,P_min_u);
502     end
503     subwc = subwc + 1;
504 end
505 P_prev = P(:,lambda);
506 sigma_N_2_prev = sigma_N_2(:,lambda);
507 end
508 end
509
510 function [M_D, H_D] = mean_of_the_min_table3(D)
511 DMH = [ 1 0.000 0.000; 2 0.260 0.150; 5 0.480 0.480; 8 0.580 0.780; ...
512        10 0.610 0.980; 15 0.668 1.550; 20 0.705 2.000; 30 0.762 2.300; ...
513        40 0.800 2.520; 60 0.841 2.900; 80 0.865 3.250; 120 0.890 4.000; ...
514        140 0.900 4.100; 160 0.910 4.100];
515
516 p = find(DMH(:,1) == D);
517 if ~isempty(p)
518     M_D = DMH(p,2);
519     H_D = DMH(p,3);
520 else
521     p = find(DMH(:,1) < D);
522     if isempty(p)
523         error('invalid value for: D');
524     end
525     p = p(end);
526     if p > size(DMH,1)-1;
527         error('invalid value for: D');
528     end
529     delta_p_rat = (D - DMH(p,1))/(DMH(p+1,1)-DMH(p,1));
530     M_D = DMH(p,2) + (DMH(p+1,2)-DMH(p,2))*delta_p_rat;
531     H_D = DMH(p,3) + (DMH(p+1,3)-DMH(p,3))*delta_p_rat;
532 end
533 end

```

## References

- [1] R. Martin. Noise power spectral density estimation based on optimal smoothing and minimum statistics. *IEEE Transactions on Speech and Audio Processing*, 9(5):504–512, 2001.
- [2] Cees H. Taal, Richard C. Hendriks, Richard Heusdens, and Jesper Jensen. A short-time objective intelligibility measure for time-frequency weighted noisy speech. In *2010 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 4214–4217, 2010.
- [3] Cees H. Taal, Richard C. Hendriks, Richard Heusdens, and Jesper Jensen. An algorithm for intelligibility prediction of time-frequency weighted noisy speech. *IEEE Transactions on Audio, Speech, and Language Processing*, 19(7):2125–2136, 2011.
- [4] Niksa Jakovljevic, Dragisa Miskovic, and Željko Trpovski. Evaluation of noise estimation algorithms based on minimum statistics and signal to noise ratio. pages 1–4, 11 2016.
- [5] Yi Hu and Philippos C. Loizou. A comparative intelligibility study of speech enhancement algorithms. In *2007 IEEE International Conference on Acoustics, Speech and Signal Processing - ICASSP '07*, volume 4, pages IV–561–IV–564, 2007.
- [6] C. Ludvigsen, C. Elberling, and G. Keidser. Evaluation of a noise reduction method—comparison between observed scores and scores predicted from sti. *Scandinavian audiology. Supplementum*, 38:50–5, 1993.