

Universitatea Politehnica București
Facultatea de Automatică și Calculatoare



Proiect – Procesarea Semnalelor
Comparison of Speech Enhancement Algorithms -
MBSS si PSS

Îndrumător științific: Dr. ing. Marios Choudary

Student: Mărgineanu Nicolae-Vlăduț

Seria și grupa: 341 C2

BUCUREȘTI

2020-2021

Cuprins

Introducere	3
Related Work.....	4
Our Work	4
Rezultate si Experimente	12
Concluzii.....	15
Bibliografie	16

1. Introducere

În acest proiect am realizat o comparație între metodele Multi-Band Spectral Subtraction Method și Power Spectral Subtraction Method.

În lumea reală, zgomotul este în mare parte colorat (**power spectrum of a noise signal**), ceea ce afectează semnalul de vorbire diferit pe tot spectrul. Puține frecvențe sunt mai afectate decât altele, în funcție de caracteristicile spectrale ale zgomotului. Cele două metode sunt folosite pentru îmbunătățirea vorbirii cu zgomot rezidual minimizat.

Semnalul de vorbire degradat în prezența zgomotului de la sol este denumit ca vorbire zgomotoasă, care poate fi modelată ca model de zgomot aditiv, cu presupunerea de bază că atât semnalul de vorbire cât și cel de zgomot sunt necorelate. **Semnalul zgomotos poate fi modelat ca suma semnalului de vorbire curat și zgomotul:**

$$y(n) = x(n) + v(n), n \in (0, N-1);$$

Unde: $y(n)$ – semnalul zgomotos, $x(n)$ – semnalul de vorbire curat și $v(n)$ – zgomotul.

În realizarea algoritmilor, trebuie să ținem cont de următoarele noțiuni teoretice:

1. Vorbirea și alte semnale audio sunt semnale variabile în timp. Dacă luăm spectrul pe întregul semnal, atunci obținem **spectrul mediu**, dar nu putem vedea modificări ale frecvențelor fundamentale.
2. În aplicațiile în timp real, trebuie **să împărțim semnalul în segmente** astfel încât să nu trebuiască să așteptăm finalizarea integrității înainte de a începe procesarea. Dacă luăm spectrul din segmente mici (**windows**) apropiate unul de altul, putem observa apoi **evoluția spectrală a semnalului**. O astfel de reprezentare este cunoscută sub numele de **spectrogramă** a unui semnal.

3. Când spectrograma este calculată utilizând **windowing** și **transformata Fourier discretă** se numește **transformată Fourier de scurtă durată** (STFT). Acum STFT al semnalului zgomotos este reprezentat de: $Y(K) = X(K) + V(K)$.

2. Related Work – algoritmi folosiți în Speech Enhancement

Power Spectral Subtraction PSS - are la bază scăderea zgomotului (power noise spectrum) din semnalul de vorbire. Factorul de “over-subtraction” - α este utilizat pentru a îmbunătăți rezultatul - α depinde de **SNR (signal to noise ratio)** segmentat. Metoda PSS este implementată **cu o singură bandă**.

Multi-Band Spectral Subtraction Method MBSS – este similar cu algoritmul PSS, în care se folosește factorul de “over-subtraction” - α . Una dintre diferențele dintre metode este că această metodă folosește un **coeficient nou, δ band subtraction** – care depinde de frecvență pentru fiecare bandă. Metoda este implementată **pe mai multe benzi**.

3. Our Work – prezentarea algoritmilor PSS și MBSS (implementare)

A. Preprocesarea datelor de intrare

Înainte de aplicarea unei metode de Speech Enhancement, realizăm următoarele operații pe semnalul de intrare:

1. Se dividează vectorul audio în **frame-uri suprapuse** (50% overlapped frames)
2. Se aplică **Hanning window** (20 ms window) pe fiecare frame după care se salvează fiecare frame într-o matrice. Aplicăm Hanning Window pentru a reduce artefactele de la marginile ferestrei.
3. Pentru fiecare frame din matrice, se aplică **funcția FFT** pentru analiză în domeniul frecvență.

4. **Weighted Spectral Average [5]** – aceasta metoda ne ajuta sa netezim variatia zgomotului.

Se realizeaza cu urmatoarea formula:

$$\bar{Y}_j(k) = \sum_{l=-M}^M W_l Y_{j-l}(k)$$

- unde j - reprezinta indexul frame-ului respectiv, M - numarul de frame-uri, W - filtrul weights, iar Y(k) semnalul zgomotos de intrare (obtinut prin FFT) (spectrul de magnitudine)

Filtrul Weights este ales ca un vector de urmatoarele elemente care a fost determinat empiric si setat ca: (valorile sunt preluate din articolul stiintific [5]).

$$W = [0.09, 0.25, 0.32, 0.25, 0.09]$$

Media a fost limitata la 5 frame-uri succesive, dat de l din formula de mai sus, pentru a prevenii alterarea spectrului. Astfel, $-2 \leq l \leq 2$.

B. Folosirea unei metode in procesarea semnalului de intrare

a) **Multi-Band Spectral Subtraction Method**

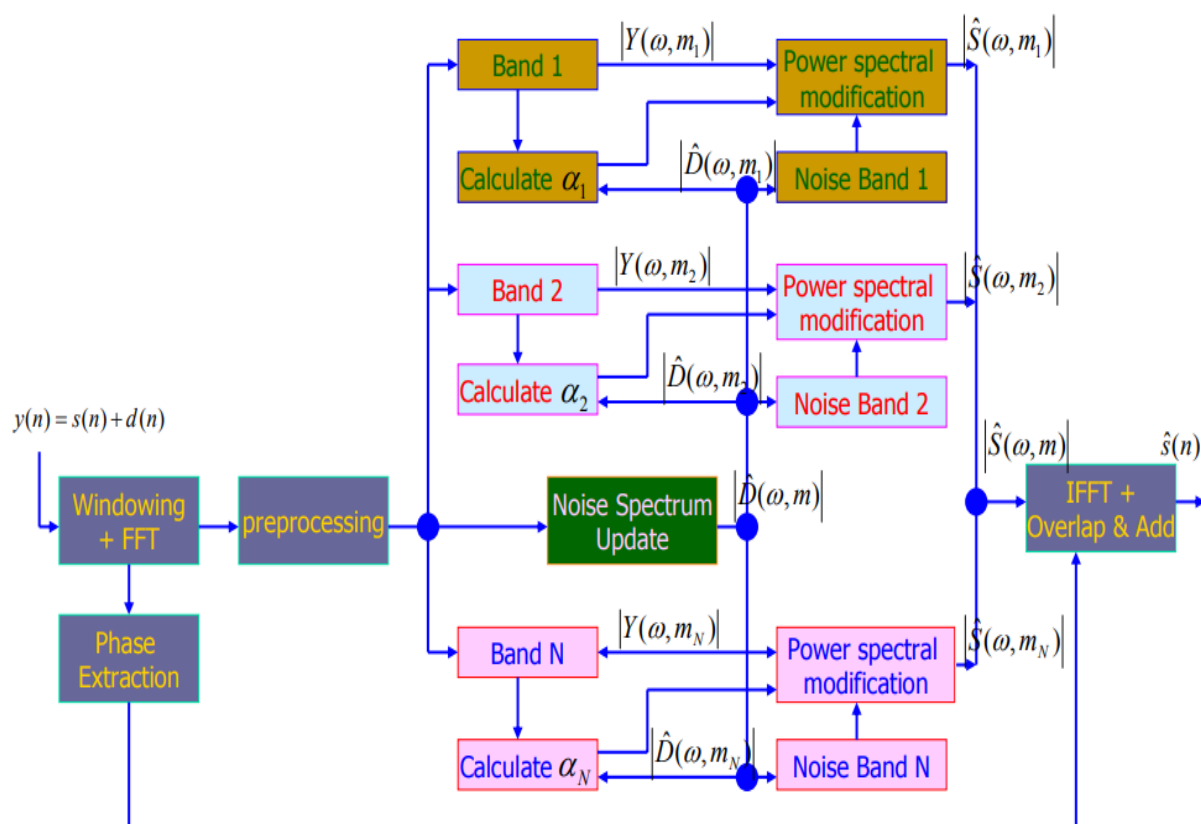
Aceasta metoda MBSS foloseste “**Linear Frequency Spacing**” [1], prin care trecem semnalul de intrare din timp, in frecventa, cu functia FFT, rezultand spectrul semnalului. FFT pentru semnalul de intrare, este dat de formula:

$$X(K) = \sum_{n=0}^{N-1} x(n) e^{-jwnK}$$

$X(K)$ – spectrul de magnitudine

$|X(K)|^2$ – power spectrum

Se trece prin urmatoarele etape in implementare:



1. În prima etapă, semnalul este “**windowed**” și **spectrul de magnitudine** este estimat folosind **functia FFT** (preprocesarea datelor de intrare).
2. În a doua etapa, spectrele de zgomot și vorbire sunt împărțite în **benzi de frecvență** diferite și se calculează **factorul de “over-subtraction”** – α (alfa) pentru fiecare bandă.

3. A treia etapă include etapa de „**band subtraction**” (factorul δ delta) în benzi de frecvență individuale prin reajustarea factorului de „over-subtraction”.
4. În cele din urmă, benzile de frecvență modificate sunt îmbinate (**merged**) și semnalul în domeniul de timp este reconstruit cu ajutorul metodei de „**overlap-add**” și aplicând funcția **IFFT**.

Prin urmare, estimarea spectrului de vorbire (clean speech spectrum) în banda a j - a este obținută prin urmatorul **algorithm**:

$$|\widehat{X}_j(K)|^2 = \begin{cases} |Y_j(K)|^2 - \alpha_j \cdot \delta_j \cdot |\widehat{V}_j(K)|^2, \\ \text{if } |\widehat{X}_j(K)|^2 > \beta \cdot |Y_j(K)|^2 \text{ else} \\ \beta \cdot |Y_j(K)|^2 \end{cases} \quad K_j < K < K_{j+1}$$

Unde:

- $Y(K)$ – spectrul de magnitudine al semnalului zgomotos (obținut prin FFT, trecut din timp în frecvență)
- $X^{\wedge}(k)$ – spectrul de magnitudine **estimat** al vorbirii curate
- $V^{\wedge}(n)$ – spectrul de magnitudine **estimat** al zgomotului calculat de-a lungul unei perioade de liniste
- $|X^{\wedge}(K)|^2$ – pătratul din valoarea absolută (abs) din spectrul de magnitudine. Se mai numește și **power spectrum**

Capetele de frecvență de început și sfârșit ale benzii de frecvență j sunt date de K_j și $K_j + 1$.

Banda specifică cu factorul de „**over-subtraction**” este dată de α_j în funcție de SNR-ul segmentat al benzii corespunzătoare. **SNR (signal to noise ratio)** segmentat al benzii j poate fi calculat astfel:

$$SNR_j(db) = 10 \log_{10} \left(\frac{\sum_{K=K_j}^{K_j+1} |Y_j(K)|^2}{\sum_{K=K_j}^{K_j+1} |\bar{V}_j(K)|^2} \right)$$

Factorul α poate fi calculat dupa formula urmatoare:

$$\alpha_j = \begin{cases} \alpha_{max}, & \text{if } SNR_j \leq SNR_{min} \\ \alpha_{max} + (SNR_j - SNR_{min}) \left(\frac{\alpha_{min} - \alpha_{max}}{SNR_{max} - SNR_{min}} \right) & \text{if } SNR_{min} \leq SNR_{max} \\ \alpha_{min}, & \text{if } SNR_j \geq SNR_{max}, \end{cases}$$

In cazul meu, factorul $\alpha_{min} = 1$, $\alpha_{max} = 5$, $SNR_{min} = -5dB$, $SNR_{max} = 20 dB$. [5]

$$\alpha_i = \begin{cases} 4.75, & SNR_i < -5 \\ 4 - \frac{3}{20} (SNR_i), & -5 \leq SNR_i \leq 20 \\ 1, & SNR_i > 20 \end{cases}$$

Factorul δ_j este un factor suplimentar de „**band subtraction**” care oferă un grad suplimentar de control în cadrul fiecărei benzi. Valoarea utilizată în fiecare bandă este calculată empiric deoarece cea mai mare parte a energiei vorbirii este concentrată sub 1 kHz. (fiecare banda are cate un factor suplimentar)

$$\delta_j = \begin{cases} 1 & f_j \leq 1kHz \\ 2.5 & 1kHz < f_j \leq \frac{F_s}{2} - 2 kHz \\ 1.5 & f_j > \frac{F_s}{2} - 2 kHz \end{cases}$$

Primul frame corespunde tăcerii inițiale din speech. Se estimează pragul zgomotului. Aici, pentru estimarea exactă a zgomotului, am calculat în medie zgomotul atunci când a fost detectat cu zgomotul estimat anterior.

Valorile negative ale spectrului estimat sunt eliminate folosind parametrul „**spectral floor**” β , după formula:

$$|\hat{S}_i(k)|^2 = \begin{cases} |\hat{S}_i(k)|^2 & |\hat{S}_i(k)|^2 > \beta |Y_i(k)|^2 \\ \beta |Y_i(k)|^2 & \text{else} \end{cases}$$

Parametrul **spectral floor** a fost setat la $\beta = 0.02$.

Pentru a obține o **distorsiune minimă** a vorbirii în regiunile cu **frecvență joasă** este de preferat să folosim **valori mai mici ale δ_j** și **valori mai mari ale δ_j** în regiunile cu **frecvență înaltă**. Prin reducerea varianței de frecvență al vorbirii, zgomotul rezidual poate fi redus în vorbirea îmbunătățită. Prin urmare, în loc să se utilizeze spectrele de putere ale semnalului, se poate utiliza o **versiune netedă a spectrelor de putere**. Valoarea medie a magnitudinii ajută la îmbunătățirea calității vorbirii procesate.

b) **Power Spectral Subtraction Method**

Descriere succintă a metodei

- Etapele acestei metode sunt la fel ca în metoda MBSS, doar că metoda este implementată **pe o singură bandă**.
- Valoarea factorului α_i „over-subtraction” este calculat pe fiecare ecuație în funcție de SNR, unde i corespunde frame-ului cu indexul respectiv din formula PSS.
- Zgomotul este estimat de lungul perioadei de liniste și apoi „**spectral subtraction**” este aplicat astfel: $|Y(k)|^2 = |S(k)|^2 + |D(k)|^2$.

- Spectral flooring (eliminarea valorilor negative), cu ajutorul parametrului „spectral floor” β , se face la fel ca în MBSS, astfel:

$$|\hat{S}_i(k)|^2 = \begin{cases} |\hat{S}_i(k)|^2 & |\hat{S}_i(k)|^2 > \beta |Y_i(k)|^2 \\ \beta |Y_i(k)|^2 & \text{else} \end{cases}$$

Algoritmul:

1. Principiul de bază al acestei metode este de a scădea spectrul de magnitudine al zgomotului $d(n)$ din vorbirea zgomotoasă $y(n)$

$$y(n) = s(n) + d(n)$$

- unde $s(n)$ este vorbirea curată “clean speech”. Se presupune că zgomotul $d(n)$ este necorelat și aditiv la semnalul de vorbire. Estimarea zgomotului $d(n)$ este măsurată în timpul tăcerii sau al activității non-verbale în semnal.
2. Spectrul de putere al semnalului zgomotos poate fi scris ca:

$$|Y(k)|^2 = |S(k)|^2 + |D(k)|^2$$

$D(k)$ – spectrul de magnitudine al zgomotului

$S(k)$ – spectrul de magnitudine al vorbirii curate

$Y(k)$ – spectrul de magnitudine al semnalului de intrare

3. Deoarece spectrul de zgomot $D(k)$ nu poate fi obținut direct, “time-average” al spectrului de putere $|D^{\wedge}(k)|$ este calculat în timpul unei perioade de tăcere, o estimare a spectrului de vorbire modificat poate fi dată ca:

$$|\hat{S}(k)|^2 = |Y(k)|^2 - \alpha |\hat{D}(k)|^2$$

- α – factorul de „over-substraction” (functie cu SNR segmentat)
- $S^{\wedge}(k)$ – spectrul de magnitudine **estimat** al vorbirii curate
- $D^{\wedge}(k)$ – spectrul de magnitudine **estimat** al zgomotului calculat de-a lungul unei perioade de liniste

Metoda PSS este implementată cu o singură bandă.

SNR-ul si **Factorul α** sunt calculate ca in metoda MBSS, prezentate mai sus.

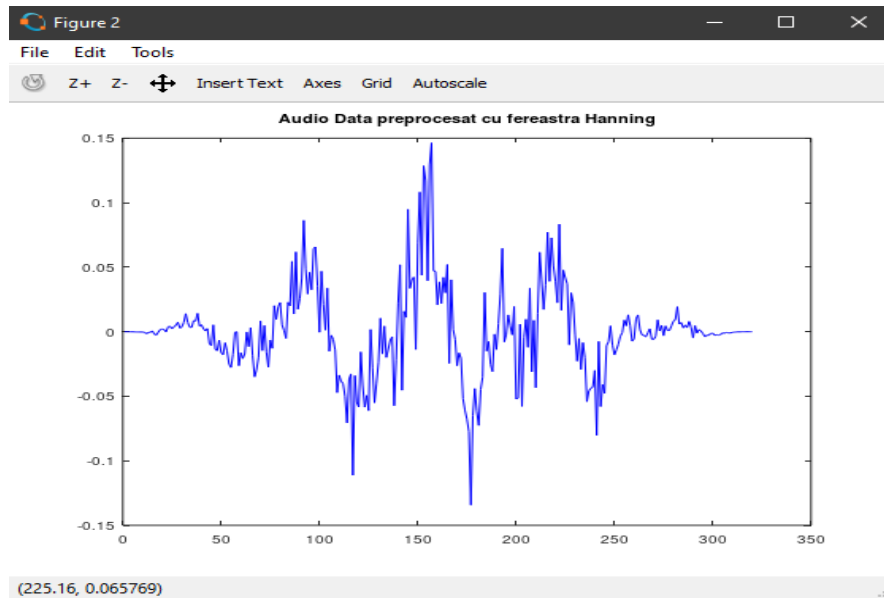
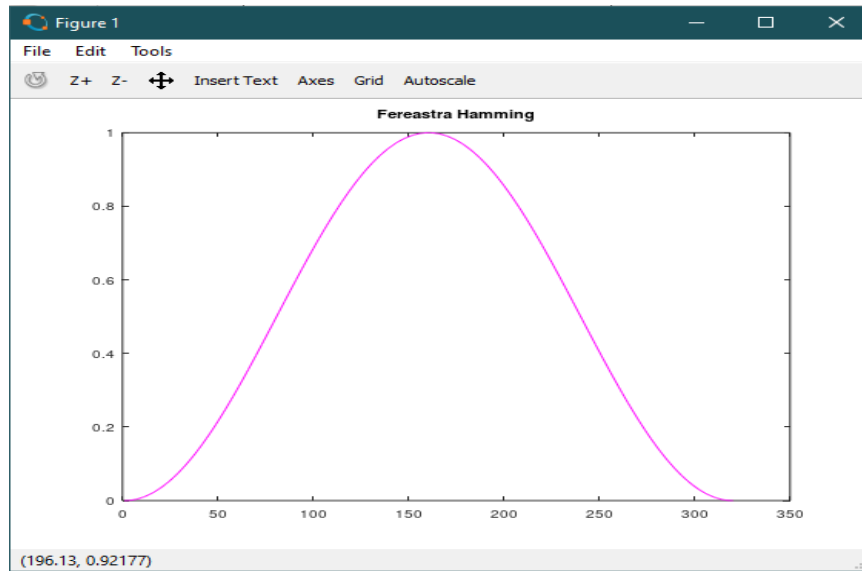
4. Spectrul modificat poate conține unele **valori negative** din cauza erorilor din spectrul de zgomot estimat. Aceste valori sunt rectificate folosind parametrul „spectral floor” β .

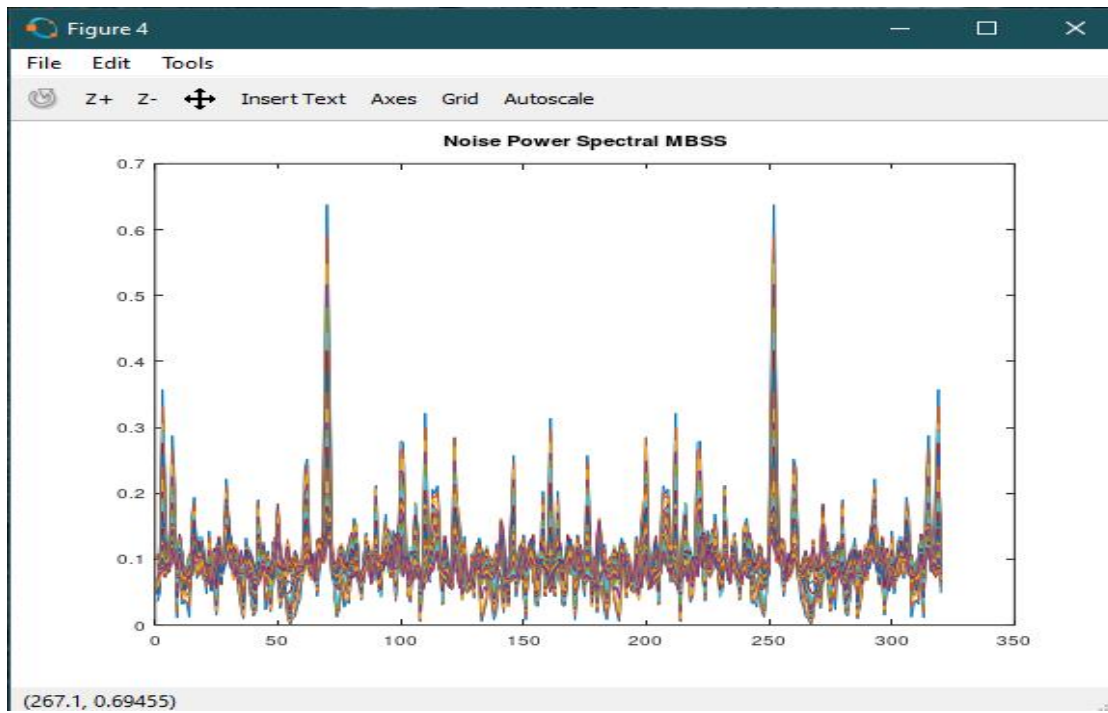
C. Reconstructia semnalului dupa aplicarea metodei PSS sau MBSS

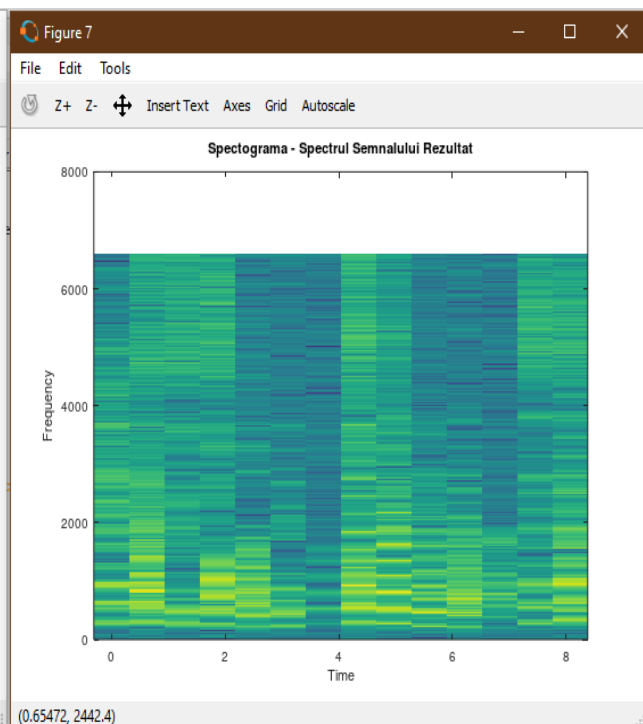
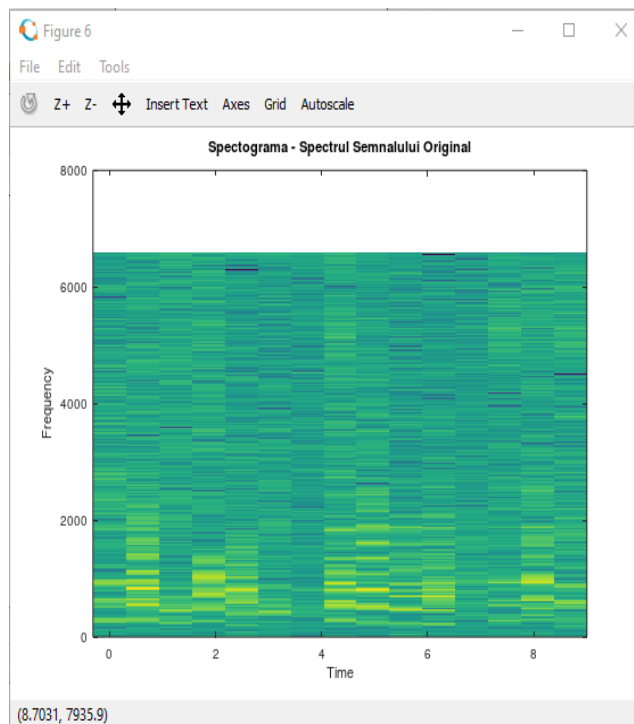
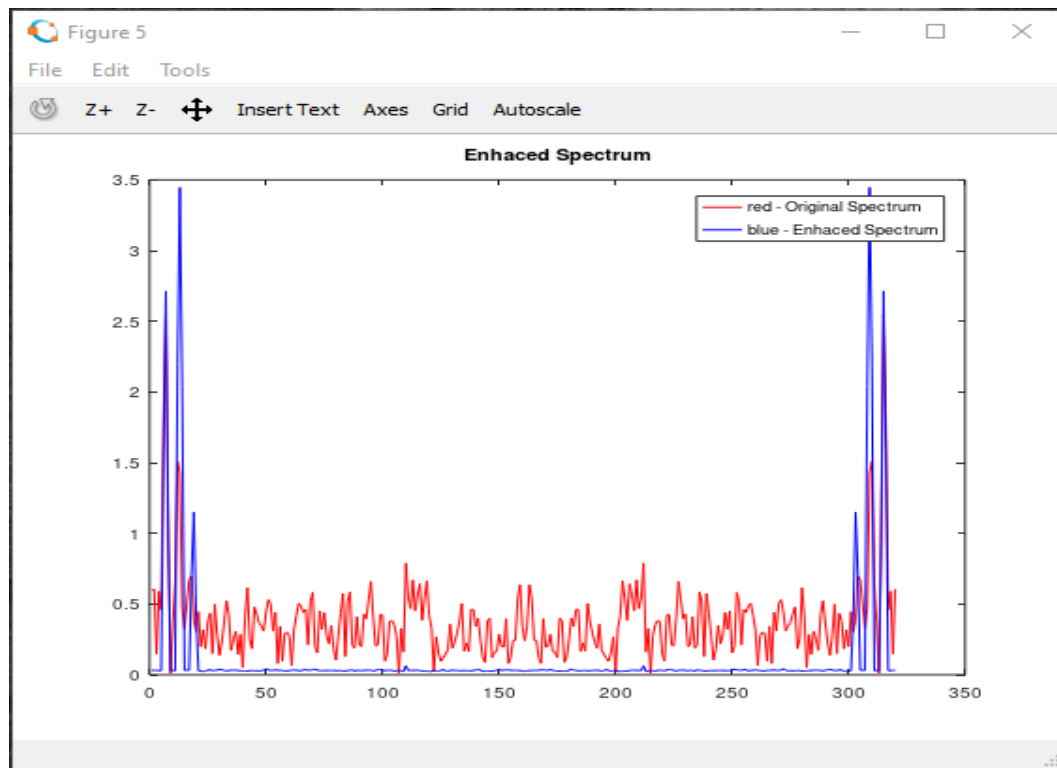
Semnalul de vorbire obtinut din una din cele doua metode, **este reconstruit pentru a obtine semnalul de vorbire direct** si cu **aceeasi frecventa de esantionare (sampling frequency)**. **Spectrul de vorbire** din fiecare cadru este transformat în **domeniu timp** folosind **functia IFFT** și apoi semnalul de vorbire complet este reconstruit folosind **Overlap Add**.

4. Rezultate si Experimente

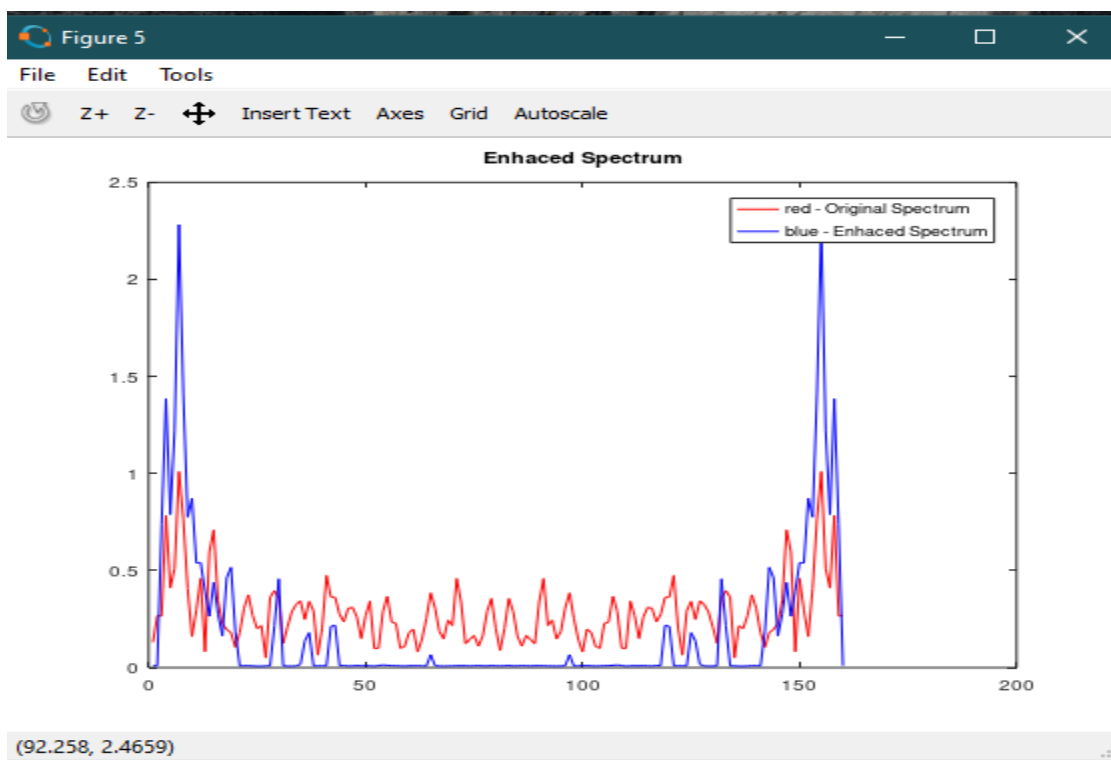
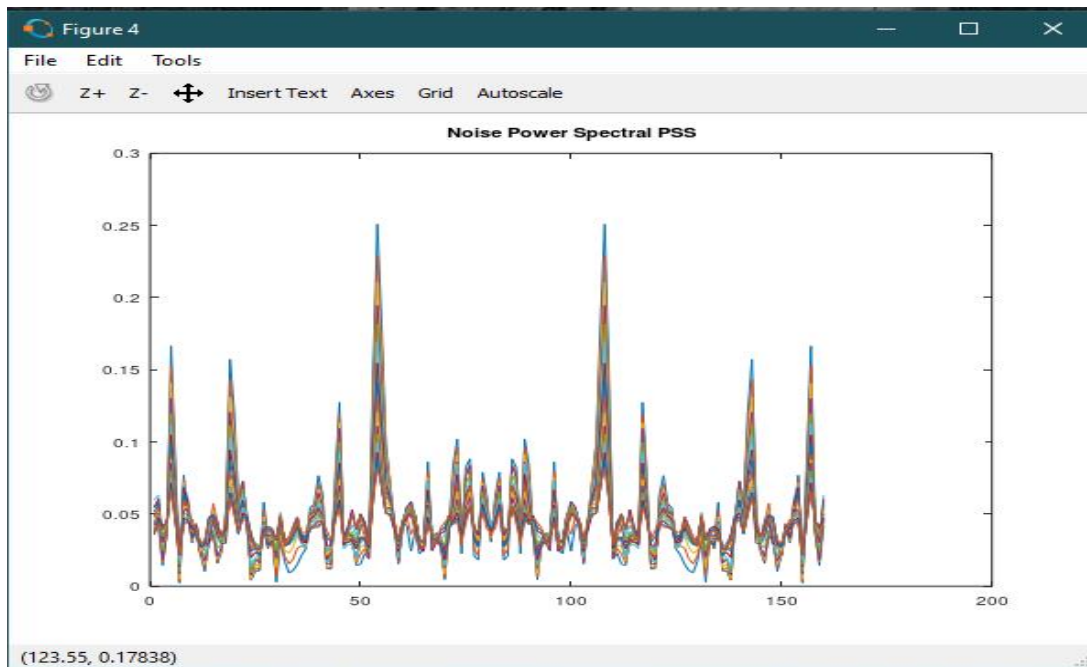
1. MBSS method

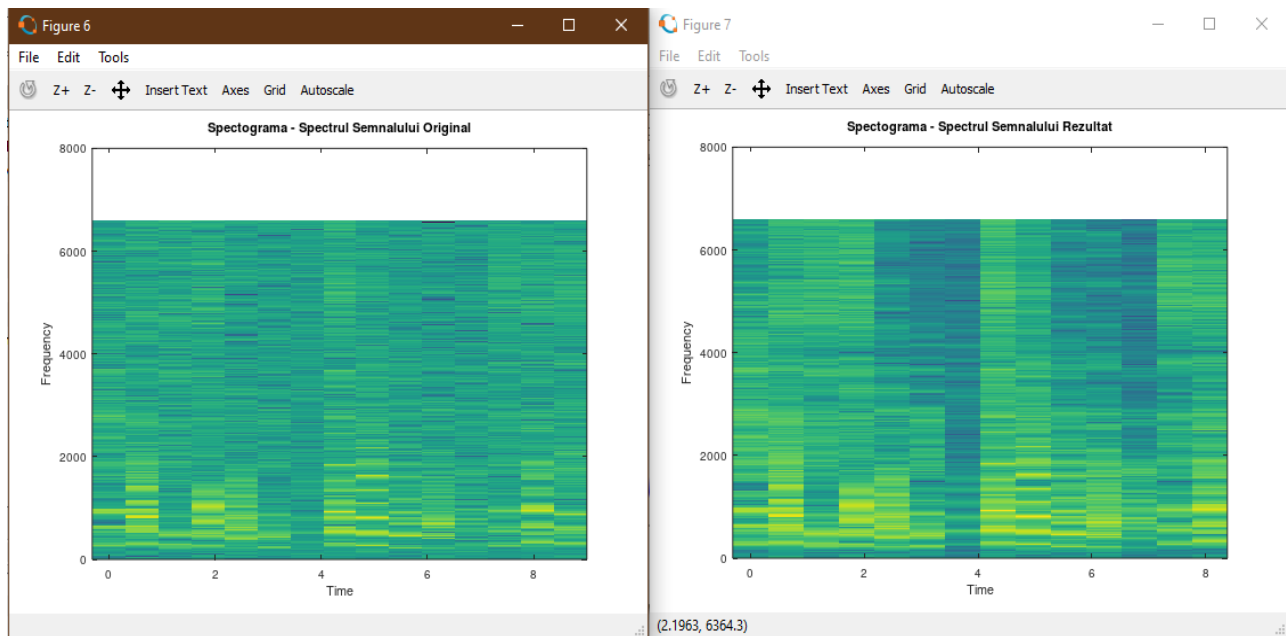






2. PSS method





5. Concluzii

1. Pentru metoda **MBSS**, **zgomotul este updatat pentru fiecare frame**. Se verifica pentru fiecare frame in parte daca este predominant speech sau zgomot si se actualizeaza vectorul noise.
2. Pentru metoda PSS, **vectorul noise se initializeaza la inceput**.
3. Metoda MBSS este mai buna **pentru speech-urile de mai lunga durata**, deoarece se **actualizeaza mreu noise vector**.
4. Factorul **subtraction-band** (delta) ne ajuta sa avem o estimare mai buna a zgomotului prezent in MBSS.

Avantajul metodei MBSS – avantajul pre-procesarii în această metodă este de a reduce variatia estimării spectrale și, ulterior, de a reduce zgomotul rezidual.

Acest algoritm reajustează factorul de „**over-subtraction**”(alfa) din fiecare bandă.

6. Bibliografie

- [1] P. Sunitha, Dr. K. Satya Prasad „Multi Band Spectral Subtraction for Speech Enhancement with Different Frequency Spacing Methods and their Effect on Objective Quality Measures”
<http://j.mecs-press.net/ijigsp/ijigsp-v11-n5/IJIGSP-V11-N5-6.pdf>
- [2] D. Deepa, A. Shanmugam „SPECTRAL SUBTRACTION METHOD OF SPEECH ENHANCEMENT USING ADAPTIVE ESTIMATION OF NOISE WITH PDE METHOD AS A PREPROCESSING TECHNIQUE”
https://www.researchgate.net/publication/271770420_SPECTRAL_SUBTRACTION_METHOD_OF_SPEECH_ENHANCEMENT_USING_ADAPTIVE_ESTIMATION_OF_NOISE_WITH_PDE_METHOD_AS_A_PREPROCESSING_TECHNIQUE
- [3] Siddala Vihari, A. Sreenivasa Murthy, Priyanka Soni and D. C. Naik “Comparison of Speech Enhancement Algorithms”
https://www.researchgate.net/publication/306362918_Comparison_of_Speech_Enhancement_Algorithms
- [4] Yasser Ghanbari, Mohammad Reza Karami-Mollaei, Behnam Amelifard “IMPROVED MULTI-BAND SPECTRAL SUBTRACTION METHOD FOR SPEECH ENHANCEMENT”
https://www.researchgate.net/profile/Yasser_Ghanbari/publication/228609880_Improved_multi-band_spectral_subtraction_method_for_speech_enhancement/links/00463517956744dd9200000.pdf
- [5] Sunil D. Kamath and Philipos C. Loizou “A MULTI-BAND SPECTRAL SUBTRACTION METHOD FOR ENHANCING SPEECH CORRUPTED BY COLORED NOISE”
<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.7.4102&rep=rep1&type=pdf>