



Computational Intelligence in Electrical Engineering  
Vol. 11, No. 1, 2020  
Research Paper

## Single-channel Speech Enhancement using the Combination of Exponential Deterministic Model and $t$ Location-scale Stochastic Model

Zahra Amini<sup>1</sup>, Neda Faraji<sup>2</sup>

<sup>1</sup> MSc Student, Dept. of Electrical Engineering, Imam Khomeini International University, Qazvin, Iran  
nfaraji@aut.ac.ir

<sup>2</sup> Assistant Professor, Dept. of Electrical Engineering, Imam Khomeini International University, Qazvin, Iran  
zahra.amini56@yahoo.com

### Abstract:

Most speech enhancement algorithms focus on obtaining an estimator relying on stochastic models. In this paper, a minimum mean-square error (MMSE) estimator under a stochastic–deterministic model is proposed where a heavy-tail distribution called  $t$ -Location-Scale (tls) is used for modeling Discrete Fourier Transform coefficients of clean speech signals and exponential and sinusoidal models are employed as deterministic models. In the exponential model, the frequency and damping coefficient are estimated by using the Matrix Pencil method. Also, in previous studies, the number of exponential components in the deterministic model for stochastic–deterministic speech enhancement algorithm has been considered to be one. In this paper, the corresponding exponential model is developed to have an arbitrary number of exponential components. The speech enhancement experiments are performed in three modes, exponential-Gaussian (the first proposed method), exponential-tls (the second proposed method), and sinusoidal-Gaussian. Comparisons are made with the exponential-Gaussian method (with only one exponential component), as well as with the Weiner and tls stochastic estimators. The implementation results in the presence of six noise types from Noisex-92 dataset show that the two proposed methods improve the segSNR values and have quite similar PESQ values comparing with the stochastic based speech enhancement methods.

**Keywords:** Speech Enhancement,  $t$  Location-scale Probability Density Function, Wiener Filter, Minimum Mean Square Error, Exponential Deterministic Model, Sinusoidal Model.



2252-083X/ © 2020 The Authors. Published by University of Isfahan

This is an open access article under the CC BY-NC-ND/4.0/ License (<https://creativecommons.org/licenses/by-nc-nd/4.0/>).



<http://dx.doi.org/10.22108/isee.2019.114459.1171>

## بهسازی گفتار تک کاناله با استفاده از ترکیب مدل قطعی نمایی و مدل تصادفی

### *t* Location-Scale

زهرا امینی<sup>۱</sup>، ندا فرجی<sup>۲</sup>

۱- دانشجوی کارشناسی ارشد، گروه مهندسی برق - دانشگاه بین‌المللی امام خمینی (ره) - قزوین - ایران

zahra.amini56@yahoo.com

۲- استادیار، گروه مهندسی برق - دانشگاه بین‌المللی امام خمینی (ره) - قزوین - ایران

nfaraji@eng.ikiu.ac.ir

**چکیده:** بیشتر روش‌های بهسازی گفتار، تخمین‌گری کاملاً متکی به مدل تصادفی گفتار ارائه می‌دهند. در این مقاله، یک تخمین‌گر کمترین میانگین مربعات خطا تحت یک مدل قطعی - تصادفی پیشنهاد می‌شود که در آن از یک توزیع دنباله - سنگین به نام *t* location-scale (*tls*) برای مدل‌کردن ضرایب تبدیل فوریه گسسته گفتار تمیز و از مدل نمایی و سینوسی به‌عنوان مدل قطعی استفاده شده است. در مدل نمایی به‌کاررفته، تخمین فرکانس و ضریب میرایی به روش ماتریس پینسل انجام می‌شود. همچنین، در پژوهش‌های قبلی تعداد مؤلفه‌های نمایی در ساخت مدل قطعی برای بهسازی گفتار، یک در نظر گرفته شده است که در این مقاله، مدل نمایی به تعداد دلخواه مؤلفه‌های نمایی بسط داده می‌شود. پیاده‌سازی‌ها در سه حالت ترکیبی نمایی - گاوسی (روش پیشنهادی نخست)، نمایی - *tls* (روش پیشنهادی دوم) و سینوسی - گاوسی انجام شده‌اند و با روش موجود نمایی - گاوسی (تنها با یک مؤلفه نمایی) و تخمین‌گرهای تصادفی وینر و مبتنی بر *tls* مقایسه می‌شوند. نتایج پیاده‌سازی در حضور شش نویز از مجموعه داده نویز noise92 نشان می‌دهند که دو روش پیشنهادی در قیاس با روش‌های مبتنی بر مدل تصادفی صرف، به بهبود معیار نسبت سیگنال به نویز قطعه‌ای منجر شده‌اند و در ارزیابی ادراکی کیفیت گفتار عملکرد نسبتاً برابری دارند.

**واژه‌های کلیدی:** بهسازی گفتار، تابع چگالی احتمال *t* Location-Scale، فیلتر وینر، کمترین میانگین مربعات خطا، مدل

قطعی نمایی، مدل سینوسی

### ۱- مقدمه

گرفته می‌شوند، در بسیاری از کاربردهای واقعی به‌طور جدی نقض شده‌اند و برقراری آنها زیر سؤال می‌رود؛ از این‌رو، مبحث بهسازی گفتار، یکی از ضرورت‌های کاربردی و عملی، از زمینه‌های فعال تحقیقاتی در سال‌های اخیر بوده است.

بهسازی گفتار، فرایند بازسازی گفتار تمیز از سیگنال نویزی گفتار است و در مواردی به کار می‌رود که سیگنال گفتار از نویز، انعکاس یا سایر عوامل مخرب تأثیر گرفته است.

در سه دهه اخیر، استراتژی‌های مختلفی برای بهسازی گفتار در حضور نویز جمع‌شونده پیشنهاد شده‌اند [۱]. بیشتر

با رشد روزافزون استفاده از سیستم‌های گفتاری در کاربردهای علمی و روزمره، نیاز به حفظ کیفیت گفتار، امری اجتناب‌ناپذیر مطرح شده است. شرایط ایده‌آل و عاری از نویزی که در کارها و شبیه‌سازی‌های آزمایشگاهی در نظر

<sup>۱</sup> تاریخ ارسال مقاله: ۱۳۹۷/۰۹/۲۰

تاریخ پذیرش مقاله: ۱۳۹۸/۰۸/۱۱

نام نویسنده مسئول: ندا فرجی

نشانی نویسنده مسئول: ایران، قزوین، دانشگاه بین‌المللی امام خمینی (ره)، دانشکده فنی، گروه مهندسی برق

این استراتژی‌ها از الگوریتم‌های بهسازی گفتار در حوزه فرکانس استفاده می‌کنند که باوجود پیچیدگی محاسباتی کمتر، نتایج بهبود کیفیت چشمگیری را به همراه دارند. روش متداول تقریب طیفی [۲] که محاسبات ریاضی ساده دارد، فیلتر وینر و تغییرات آن مانند فیلترینگ تکراری وینر [۳] از روش‌های بهسازی در حوزه فرکانس‌اند.

گروه دیگر از روش‌های بهسازی در حوزه فرکانس روش‌های مبتنی بر مدل آماری‌اند. در این روش‌ها عموماً از روش بیزین برای تخمین سیگنال بهبودیافته استفاده می‌شود. از جمله روش‌های آماری متداول می‌توان به تخمینگر کمترین میانگین مربعات خطا<sup>۱</sup> (MMSE) [۲]، تخمینگر مبتنی بر لگاریتم کمترین میانگین مربعات خطا (Log-MMSE) [۴]، بیشینه درست‌نمایی<sup>۲</sup> (ML) [۵] و تخمینگر بیشینه احتمال پسین<sup>۳</sup> (MAP) [۶] اشاره کرد.

تخمینگر کمترین میانگین مربعات خطا و بیشینه درست نمایی عمدتاً در حوزه تبدیل فوریه گسسته استفاده می‌شوند و با استفاده از این روش‌ها دامنه یا ضرایب مختلط تبدیل فوریه سیگنال گفتار از ضرایب مختلط تبدیل فوریه گسسته گفتار نویزی تخمین زده می‌شود.

در تخمینگرهای کمترین میانگین مربعات خطا، یک تابع هزینه استاندارد ریاضی یا یک معیار پذیرفتنی بهینه می‌شود تا یک تابع بهره غیرخطی برای اصلاح ضرایب تبدیل فوریه گسسته سیگنال نویزی به دست آید. برای پیدا کردن تخمینگرهای کمترین میانگین مربعات خطا و بیشینه درست نمایی به دو تابع چگالی احتمال نیاز است؛ تابع چگالی احتمال پیشین (تابع چگالی احتمال سیگنال تمیز) و تابع چگالی احتمال سیگنال نویز.

با این فرض که ضرایب تبدیل فوریه گسسته سیگنال گفتار تمیز و نویز هر دو گاوسی‌اند، تخمینگر کمترین میانگین مربعات خطا به دست آمده است [۲]. در پژوهش‌های انجام شده درباره تابع چگالی احتمال سیگنال گفتار تمیز در حوزه فرکانس و زمان، نشان داده شد که قسمت‌های حقیقی و موهومی ضرایب تبدیل فوریه گسسته سیگنال تمیز یک توزیع سوپرگاوسی<sup>۴</sup> دارند (یک قله واضح‌تر و دنباله‌های

سنگین در مقایسه با گاوسی) [۶]. تا به امروز عمدتاً فرض بر این بوده است که نویز دارای تابع چگالی احتمال گاوسی و سیگنال گفتار تمیز دارای تابع چگالی احتمال لاپلاس [۷]، گاما [۸]، گامای تعمیم‌یافته<sup>۵</sup> [۹] و [Location-scale  $t$ ] [۱۰] است. یک اصلاح دیگر برای بهبود الگوریتم‌های بهسازی گفتار مبتنی بر روش‌های آماری، در نظر گرفتن احتمال حضور<sup>۶</sup> (غیاب) گفتار است [۱].

استفاده از مدل قطعی در روش‌های بهسازی گفتار نخستین بار در سال ۱۹۸۰ مطرح شد [۱۱]. در این مقاله، سیگنال گفتار تمیز با یک تابع نمایی با دامنه و فاز نامشخص نمایش داده شد. پس از آن در سال ۱۹۹۳، در چارچوب فیلتر وینر، یک الگوریتم بهسازی گفتار ارائه شد که هم‌زمان هر دو جزء گفتار قطعی و تصادفی را در نظر می‌گیرد [۱۲]. در این مقاله، سیگنال گفتار به دو بخش مصوت (بخش پرودیگ سیگنال گفتار) و نامصوت (بخش غیرپرودیگ سیگنال گفتار) تقسیم شد و قسمت مصوت گفتار با مدل قطعی و قسمت نامصوت گفتار نیز تصادفی در نظر گرفته شد. در سال ۲۰۰۷، هندریکس و همکارانش برای ضرایب تبدیل فوریه گفتار، مدل تصادفی یا قطعی در نظر گرفتند و سپس تخمین کمترین میانگین مربعات خطا را مبتنی بر هر دو تصمیم نرم و سخت در انتخاب بین دو مدل تصادفی و قطعی گفتار به دست آوردند [۱۳]؛ البته در مدل نویز به‌علاوه هارمونیک از سیگنال گفتار، ضرایب تبدیل فوریه سیگنال تمیز در هر لحظه از هر دو مدل تصادفی و قطعی پیروی می‌کنند که با یک توزیع با میانگین غیرصفر نمایش داده می‌شوند [۱۴]. مک‌کالوم و همکاران، این ایده را در سال ۲۰۱۲ استفاده کردند. در این مقاله بر خلاف کارهای پیشین که برای نویز یک مدل تصادفی در نظر می‌گرفتند، پژوهشگران یک مدل نویز تصادفی - قطعی در نظر گرفتند [۱۵]. در سال ۲۰۱۳ نیز مک‌کالوم و همکارانش یک الگوریتم بیزین تحت مدل گفتار تصادفی - قطعی با اطلاعات پیشین (استفاده از اطلاعات فریم قبلی برای تخمین فرکانس هر فریم) را پیشنهاد کردند. در این الگوریتم، میانگین توزیع غیرصفر در نظر گرفته می‌شود [۱۶].

پایین تری خواهد داشت [۱۸]. روش‌های مک‌براید<sup>۴</sup>، ماتریس پنسیل<sup>۱</sup> و حداقل مربعات کل<sup>۱۱</sup> در مراجع [۱۹]، [۲۰] و [۲۱]، روش‌های مؤثر در تخمین فرکانس و ضریب میرایی معرفی شده‌اند. در این بخش، این سه روش در تعداد نمونه‌های مختلف از سیگنال، تعداد فرکانس‌های تخمینی و نسبت سیگنال به نویزهای مختلف (از صفر تا ۴۰ دسی بل) با یکدیگر مقایسه می‌شوند. همچنین، برای مقایسه روش‌های مختلف تخمین فرکانس و ضریب میرایی از معیار میانگین قدر مطلق خطا به ترتیب مطابق با روابط (۲) و (۳) استفاده می‌شود.

$$MAE_f = \frac{1}{P} \sum_{p=1}^P |\hat{f}_p - f_p|, \quad (2)$$

$$MAE_d = \frac{1}{P} \sum_{p=1}^P |\hat{d}_p - d_p|, \quad (3)$$

در این دو رابطه،  $P$  تعداد توابع نمایی در معادله (۱) و  $\hat{f}_p, \hat{d}_p, f_p, d_p$  نیز به ترتیب فرکانس تخمینی، ضریب میرایی تخمینی، فرکانس واقعی و ضریب میرایی واقعی  $p$  امین مؤلفه فرکانس مختلط را نشان می‌دهند. سیگنال ورودی مطابق با رابطه (۱) با فرکانس، ضریب میرایی، دامنه، فاز و تعداد نمونه‌های مختلف ساخته شده است و به هرکدام از روش‌های تخمین داده می‌شود. سپس مطابق با روابط (۲) و (۳) مقدار میانگین قدر مطلق خطای تخمین محاسبه می‌شود.

شکل (۱) اثر تعداد فرکانس بر تخمین فرکانس و ضریب میرایی را وقتی تعداد نمونه‌ها برابر ۵۰۰ است، در سه روش مک‌براید (StMcB)، ماتریس پنسیل (MatPen) و حداقل مربعات کل (Total LS) نشان می‌دهد. در شکل (۱.الف) و (۱.ب)، به ترتیب خطای تخمین فرکانس و خطای تخمین ضریب میرایی، وقتی تعداد فرکانس‌ها ۱۰ و در شکل (۱.ب) و (۱.د)، خطای تخمین فرکانس و خطای تخمین ضریب میرایی، وقتی تعداد فرکانس‌ها ۲۰ باشد، نشان داده شده است. طبق این شکل با افزایش تعداد فرکانس، خطای تخمین در روش ماتریس پنسیل کمتر از بقیه روش‌های تخمین است.

به تازگی نیز روش بهسازی مبتنی بر شبکه عصبی عمیق<sup>۷</sup> ارائه شده است که از داده‌های آموزشی برای یافتن نگاشت بین داده‌های گفتار نویزی و تمیز استفاده می‌کند [۱۷].

در این مقاله، روش بهسازی گفتار مبتنی بر مدل قطعی - تصادفی تعمیم داده می‌شود که در [۱۳] بیان شده است. در [۱۳] از توزیع تصادفی گاوسی و لاپلاس به عنوان مدل تصادفی و از مدل نمایی به عنوان مدل قطعی استفاده کردند. در این مقاله از توزیع جدیدی به نام  $t$  location-scale به عنوان مدل تصادفی استفاده می‌شود. همچنین، مدل نمایی استفاده شده در [۱۳]، با افزایش پارامتر  $P$  و تغییر روش تخمین فرکانس تعمیم داده می‌شود.

ساختار مقاله حاضر به شرح زیر است: در بخش دوم، مقوله مدل سازی سیگنال گفتار با استفاده از مدل قطعی - نمایی به همراه روش‌های مختلف تخمین فرکانس و مقایسه آنها بررسی می‌شود. در بخش سوم، چند روش بهسازی گفتار تحت مدل تصادفی و مدل تصادفی - قطعی در حضور نویز جمع شونده معرفی می‌شوند. در بخش چهارم مقاله، روش پیشنهادی مطرح شده است. در بخش پنجم، روش‌های بهسازی گفتار پیشنهاد شده با روش‌های پیشین در این حوزه مقایسه می‌شوند و در بخش پایانی، نتیجه گیری کلی از مقاله ارائه می‌شود.

## ۲- مدل سازی سیگنال گفتار با استفاده از

### مدل قطعی نمایی

مدل قطعی نمایی یک سیگنال، شامل جمع  $P$  تابع تک فرکانس مختلط مطابق رابطه (۱) است:

$$y[n] = \sum_{i=1}^P \frac{A_i}{2} e^{j\phi_i} e^{(d_i + j2\pi f_i)n}, \quad (1)$$

در این رابطه،  $n$  اندیس نمونه‌های زمانی است و  $A_i, d_i, \phi_i, f_i$  نیز به ترتیب دامنه، فاز، ضریب میرایی و فرکانس  $i$  امین تابع فرکانس مختلط را نشان می‌دهند. روش‌های مختلفی برای تخمین پارامترهای فرکانس و ضریب میرایی وجود دارند. روش پرونی<sup>۸</sup> از ساده ترین روش هاست که با افزایش تعداد فرکانس‌های مدل، دقت

برای تخمین فرکانس و ضریب میرایی از این روش استفاده می‌شود.

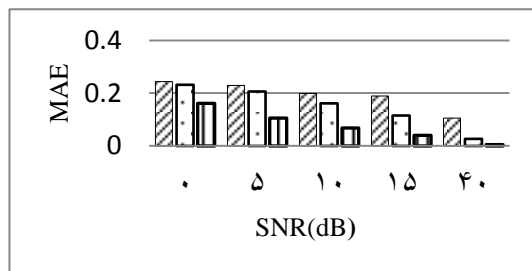
### ۳- بهسازی گفتار در نویز جمع شونده مبتنی بر روش‌های تصادفی و قطعی - تصادفی

اگر  $s(n)$  سیگنال تمیز و  $w(n)$  نویز جمع شونده باشد، سیگنال نویزی  $y(n)$  طبق رابطه (۴) به دست می‌آید:

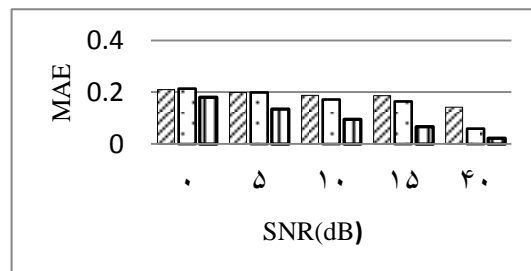
$$y(n) = s(n) + w(n). \quad (4)$$

با فرض ناهمبسته بودن نویز با سیگنال اصلی، در حوزه فوریه رابطه (۵) برقرار است.

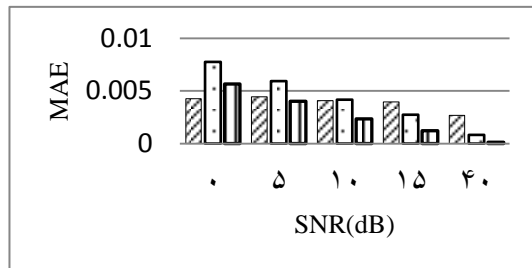
شکل (۲) اثر تعداد نمونه‌ها بر تخمین فرکانس و ضریب میرایی را وقتی تعداد فرکانس‌ها برابر ۲۰ است، در سه روش مک‌براید، ماتریس پنسیل و حداقل مربعات کل نشان می‌دهد. در شکل (۲.الف) و (۲.ج)، به ترتیب خطای تخمین فرکانس و خطای تخمین ضریب میرایی، وقتی تعداد نمونه‌ها ۵۰۰ و در شکل (۲.ب) و (۲.د)، خطای تخمین فرکانس و خطای تخمین ضریب میرایی، وقتی تعداد نمونه‌ها ۱۰۰۰ باشد، نشان داده شده است. مطابق این شکل، روش ماتریس پنسیل با افزایش تعداد نمونه‌ها، به خطای تخمین کمتری در مقایسه با سایر روش‌ها منجر می‌شود. طبق مشاهدات مذکور که حاکی از قدرت روش ماتریس پنسیل است، در بخش‌های بعدی



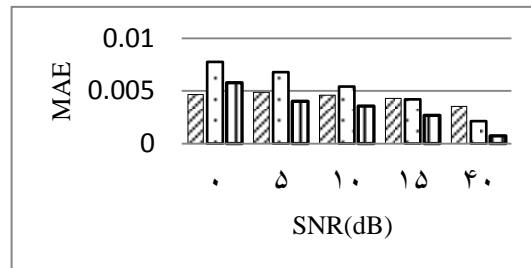
(الف)



(ب)



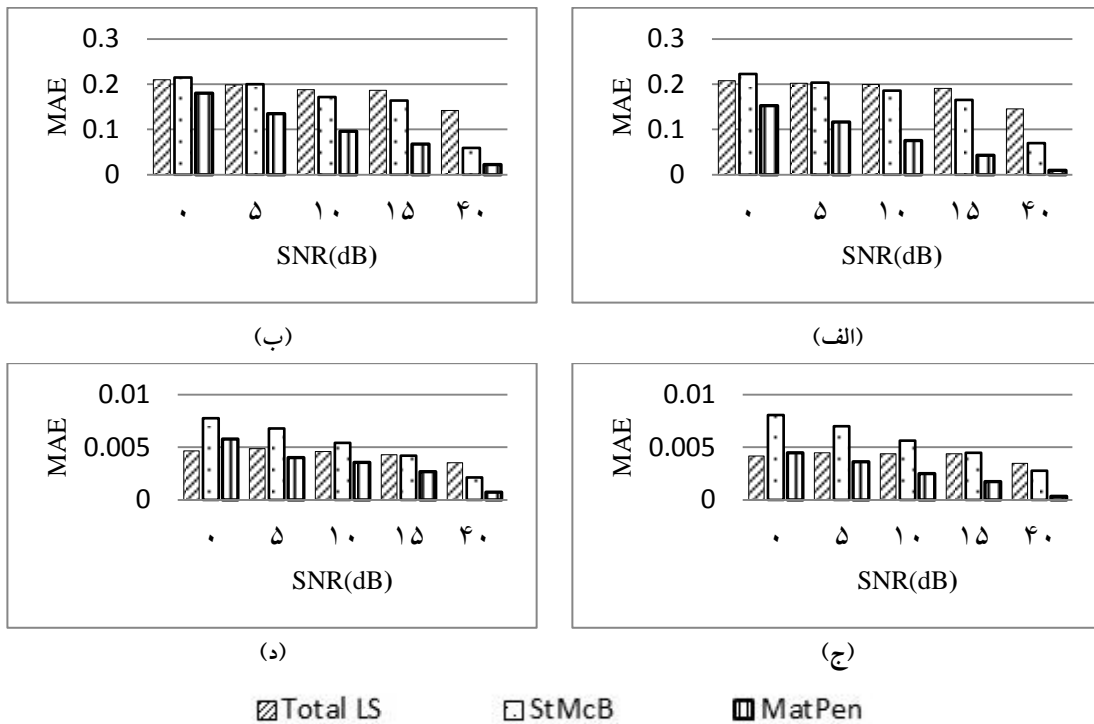
(ج)



(د)

▨ Total LS      □ StMcB      ■ MatPen

شکل (۱): بررسی اثر تعداد فرکانس‌ها بر خطای تخمین فرکانس و ضریب میرایی در روش‌های ماتریس پنسیل (MatPen)، مک‌براید (StMcB) و حداقل مربعات کل (Total LS) در ۵۰۰ نمونه (الف): خطای تخمین فرکانس با ۲۰ فرکانس، (ب): خطای تخمین فرکانس با ۱۰ فرکانس، (ج): خطای تخمین ضریب میرایی با ۲۰ فرکانس، (د): خطای تخمین ضریب میرایی با ۱۰ فرکانس.



شکل (۲): بررسی اثر تعداد نمونه‌ها بر خطای تخمین فرکانس و ضریب میرایی در روش‌های ماتریس پَنسِل (MatPen)، مک‌براید (StMcB) و حداقل مربعات کل (Total LS) در ۲۰ فرکانس (الف): خطای تخمین فرکانس با ۱۰۰۰ نمونه، (ب): خطای تخمین فرکانس با ۵۰۰ نمونه، (ج): خطای تخمین ضریب میرایی با ۱۰۰۰ نمونه، (د): خطای تخمین ضریب میرایی با ۵۰۰ نمونه.

### ۳-۱-۱- بهسازی با فرض مدل تصادفی گاوسی

تحت مدل تصادفی و استفاده از این فرض که ضرایب تبدیل فوریه گسسته سیگنال تمیز، توزیع مختلط گاوسی با میانگین صفر دارند، ضرایب تبدیل فوریه گسسته سیگنال نویزی، توزیع مختلط گاوسی با میانگین صفر خواهند داشت. سیگنال گفتار تمیز با رابطه (۶) تخمین زده می‌شود [۲]:

$$\begin{aligned} \hat{s}(k, i) &= E\{S(k, i) | y(k, i)\} \\ &= \frac{\xi(k, i)}{1 + \xi(k, i)} y(k, i), \end{aligned} \quad (6)$$

$\xi(k, i)$  مقدار سیگنال به نویز پیشین است که در بخش بعد روش تخمین آن توضیح داده می‌شود.

### ۳-۱-۲- بهسازی با فرض مدل تصادفی $t$ location-scale

در این مدل فرض می‌شود ضرایب حقیقی و موهومی تبدیل فوریه زمان کوتاه سیگنال تمیز دارای توزیع  $t$

$$Y(k, i) = S(k, i) + W(k, i), \quad (5)$$

$Y(k, i)$ ،  $S(k, i)$  و  $W(k, i)$  به ترتیب تبدیل فوریه سیگنال نویزی، سیگنال تمیز و نویزند.  $k$  و  $i$  نیز به ترتیب نشان‌دهنده اندیس فرکانس و شماره فریم‌اند.

### ۳-۱-۳- بهسازی گفتار با فرض مدل تصادفی برای گفتار

در این بخش، بهسازی گفتار با استفاده از معیار کمینه میانگین مربعات خطا و با فرض تابع چگالی احتمال گاوسی و location-scale  $t$  برای قسمت حقیقی و موهومی ضرایب فوریه زمان کوتاه سیگنال گفتار تمیز معرفی می‌شود. همچنین، فرض می‌شود قسمت حقیقی و موهومی ضرایب فوریه زمان کوتاه از هم مستقل اند و ضرایب فوریه زمان کوتاه برای سیگنال نویز نیز از توزیع گاوسی مختلط با میانگین صفر تبعیت می‌کند.

در این رابطه،  $\sigma_S^2$  و  $\sigma_W^2$  به ترتیب واریانس گفتار تمیز و نویز را نشان می‌دهند. هرچند واریانس سیگنال گفتار تمیز در دسترس نیست، افرایم و همکارانش روش تصمیم‌گیری مستقیم (Decision-Directed) را برای محاسبه آن پیشنهاد دادند که مطابق با رابطه (۱۱) است [۱]:

$$\xi(k,i) = \alpha \frac{|\hat{s}(k,i-1)|^2}{\sigma_W^2(k,i)} + (1-\alpha) \times \max(\zeta(k,i)-1, 0), \quad (11)$$

$\alpha$  ثابت هموارسازی است و  $||$  عملگر قدر مطلق را نشان می‌دهد.  $\zeta(k,i)$  یعنی مقدار سیگنال به نویز پسین در بین فرکانسی  $k$  ام و فریم  $i$  ام به‌راحتی با رابطه (۱۲) محاسبه می‌شود:

$$\zeta(k,i) = \frac{E[Y(k,i)^2]}{E[W(k,i)^2]} = \frac{\sigma_Y^2(k,i)}{\sigma_W^2(k,i)}. \quad (12)$$

### ۳-۲- بهسازی گفتار به روش قطعی - تصادفی

در این بخش، یک نمونه روش بهسازی گفتار تک‌کاناله با این فرض تشریح می‌شود که گفتار از هر دو مدل قطعی و تصادفی پیروی می‌کند [۱۳].

### ۳-۲-۱- تابع چگالی احتمال ضرایب تبدیل فوریه زمان کوتاه سیگنال نویزی در مدل های تصادفی و قطعی

**مدل تصادفی گاوسی:** تحت مدل تصادفی و استفاده از این فرض که ضرایب تبدیل فوریه گسسته سیگنال تمیز و نویز هر دو توزیع مختلط گاوسی با میانگین صفر دارند، قسمت حقیقی / موهومی ضرایب تبدیل فوریه گسسته سیگنال نویزی دارای توزیع گاوسی با میانگین صفر رابطه (۱۳) است که برای سادگی، رابطه صرفاً برای قسمت حقیقی نوشته شده است.

location-scale و ضرایب تبدیل فوریه زمان کوتاه سیگنال نویز دارای توزیع مختلط گاوسی با میانگین صفرند. تابع چگالی احتمال location-scale  $t$  برای متغیر تصادفی  $S$  با  $\nu$  درجه آزادی و میانگین صفر به صورت زیر تعریف می‌شود [۱۰]:

$$f_S(s) = \frac{\Gamma((\nu+1)/2)}{\sqrt{\nu\pi}\Gamma(\nu/2)} \left(1+s^2/(b^2\nu)\right)^{-\frac{(\nu+1)}{2}}, \quad (7)$$

$\nu > 0$ ،  $b$  نشان دهنده مقیاس،  $f(\cdot)$  نشان دهنده تابع چگالی احتمال و  $\Gamma(\cdot)$  تابع گاما است.

تحت مدل گفتار  $t$  location-scale ضرایب مختلط تبدیل فوریه زمان کوتاه سیگنال گفتار تمیز، با ترکیب تخمین های قسمت های حقیقی و موهومی به فرم رابطه (۸) تخمین زده می‌شوند [۱۰]:

$$\begin{aligned} \bar{s}(k,i) &= E\{S(k,i) | y(k,i)\} \\ &= E\{S_R | y_R\} + jE\{S_I | y_I\}, \end{aligned} \quad (8)$$

در این رابطه، اندیس های  $R$  و  $I$  به ترتیب نشان دهنده قسمت حقیقی و موهومی اند و  $j = \sqrt{-1}$  است. قسمت حقیقی طبق رابطه (۹) تخمین زده می‌شود و جزء موهومی نیز طبق رابطه مشابه و با جایگزینی اندیس  $R$  با  $I$  به دست می‌آید [۱۰].

$$\begin{aligned} E\{S_R | y_R\} &= \frac{\sum_{j=0}^{\infty} \frac{1}{j!} \left(\frac{\zeta_R}{2}\right)^j \Gamma\left(\frac{\nu+1}{2}\right)}{\sum_{j=0}^{\infty} \frac{1}{j!} \left(\frac{\zeta_R}{2}\right)^j \Gamma\left(\frac{\nu+1}{2}\right)} \times \\ &\quad \frac{\psi\left(\frac{\nu+1}{2}, \frac{\nu}{2} - j; \frac{\zeta_R(\nu-2)}{2}\right)}{\psi\left(\frac{\nu+1}{2}, \frac{\nu}{2} + 1 - j; \frac{\zeta_R(\nu-2)}{2}\right)} y_R, \end{aligned} \quad (9)$$

در این رابطه  $\psi(\cdot)$  تابع Tricomi's hypergeometric است و  $\zeta(k,i)$  و  $\xi(k,i)$  به ترتیب نشان دهنده مقدار سیگنال به نویز پیشین و پسین در بین فرکانسی  $k$  ام و فریم  $i$  ام هستند که در ادامه تعریف می‌شوند.

$$\xi(k,i) = \frac{E[S(k,i)^2]}{E[W(k,i)^2]} = \frac{\sigma_S^2(k,i)}{\sigma_W^2(k,i)}, \quad (10)$$

$$\hat{s}(k, i) = E[S(k, i) | y(k, i)] = s'(k, i) \quad (16)$$

در این رابطه،  $s'(k, i) = E[Y(k, i)]$  است.

### ۳-۲-۳- تخمین $E[Y(k, i)]$ با استفاده از مدل قطعی

فرض می شود سیگنال تمیز گفتار با جمع  $P$  تابع نمایی نزولی با فرکانس ثابت بازنمایی می شود:

$$s_t(n) = \sum_{p=1}^P a_p e^{j\phi_p} e^{(-d_p + j\nu_p)n} \quad (17)$$

$s_t(\cdot)$  نمونه ها در حوزه زمان،  $n$  اندیس نمونه زمانی،  $a_p$  دامنه،  $\phi_p$  فاز،  $d_p$  ضریب میرایی و  $\nu_p$  مؤلفه فرکانس  $p$  ام است؛ در نتیجه، ضرایب تبدیل فوریه گسسته در هر بین فرکانسی  $k$  با جمع  $P$  تابع نمایی مختلط توصیف می شوند [۱۳]. بعد از شیفت و پنجره گذاری  $s_t(n)$  رابطه (۱۸) به دست می آید:

$$s_w(n-mM) = \sum_{p=1}^P a_p e^{j\phi_p} e^{(-d_p + j\nu_p)(n-mM)} w(n), \quad (18)$$

$w(n)$  پنجره آنالیز به طول  $K$  و مقدار شیفت هر فریم ( $M \leq K$ ) است. اکنون از سیگنال پنجره شده  $s_w(n)$  تبدیل فوریه گسسته گرفته می شود:

$$s(k, m) = \sum_{n=0}^{K-1} \sum_{p=1}^P a_p e^{j\phi_p} e^{(-d_p + j\nu_p)(n-mM)} \quad (19)$$

$$\times w(n) e^{-j\omega_k n} = s(k, 0) e^{(-d_p + j\nu_p)(mM)},$$

$L$  اندازه تبدیل فوریه گسسته  $\omega_k = 2\pi/L$  و  $m$  اندیس فریم است. رابطه (۱۹) به فرم  $s(k, i) = z^m s(k, 0)$  نوشته می شود؛ با این فرض که  $z = e^{(d_p - j\nu_p)M}$  است. حال اگر نویز در فریم های  $m = i - n_1, \dots, i + n_2$  ایستان باشد و  $M$  به اندازه کافی بزرگ انتخاب شود، نویز در بازه  $m = i - n_1, \dots, i + n_2$  سفید خواهد بود [۱۳].

برای تخمین  $d_p$  و  $\nu_p$  از روش ماتریس پنیسل استفاده می شود که در بخش قبل، بهترین روش تخمین

$$P_{Y_R|ST}(y_R(k, i) | st) = \frac{1}{\sqrt{2\pi \times \sigma_{Y_R}^2(k, i)}} \times \exp\left\{-\frac{|y_R(k, i)|^2}{2 \times \sigma_{Y_R}^2(k, i)}\right\}, \quad (13)$$

در این رابطه،  $ST$  نشان دهنده سیگنال گفتار مبتنی بر مدل تصادفی است. همچنین،  $\sigma_{Y_R}^2(k, i)$  واریانس قسمت حقیقی ضرایب تبدیل فوریه گسسته سیگنال نویزی و برابر با جمع واریانس نویز و واریانس سیگنال تمیز است.

$$\sigma_{Y_R}^2(k, i) = \sigma_{S_R}^2(k, i) + \sigma_{W_R}^2(k, i). \quad (14)$$

**مدل قطعی:** تحت مدل قطعی گفتار فرض می شود  $Y(k, i)$  جمع متغیر قطعی  $S(k, i)$  و متغیر تصادفی  $W(k, i)$  است؛ بنابراین، با فرض توزیع گاوسی مختلط با میانگین صفر برای ضرایب تبدیل فوریه گسسته نویز، ضرایب تبدیل فوریه گسسته سیگنال نویزی، توزیع گاوسی با میانگین غیر صفر خواهد داشت.

$$P_{Y_R|D}(y_R(k, i) | d) = \frac{1}{\sqrt{2\pi \times \sigma_{W_R}^2(k, i)}} \times \exp\left\{-\frac{|y_R(k, i) - E[Y_R(k, i)]|^2}{2 \times \sigma_{W_R}^2(k, i)}\right\}, \quad (15)$$

در این رابطه،  $D$  نشان دهنده قطعی بودن مدل گفتار است. در این صورت  $E[Y_R(k, i)] = S_R(k, i)$  و  $E[Y_R(k, i)] = \text{Re}\{E[Y(k, i)]\}$  است.

### ۳-۲-۲- تخمینگر کمترین میانگین مربعات خطا

**مدل تصادفی گاوسی:** تخمین سیگنال تمیز مطابق بخش ۳-۱-۱ است.

**مدل قطعی:** تحت مدل قطعی گفتار، ضرایب تبدیل فوریه گفتار تمیز، معین اما ناشناخته فرض می شوند. به این معنی که  $P_S(s(k, i)) = \delta(s(k, i) - s'(k, i))$  و  $s'(k, i)$  مقدار ضریب تبدیل فوریه گسسته سیگنال تمیز، معین و تابع ضربه است. پس تخمینگر کمترین میانگین مربعات خطا مطابق با رابطه (۱۶) خواهد بود:



$$P_{ST|Y_R}(st | y_R) = \frac{\Lambda_{STR}}{\Lambda_{DR} + \Lambda_{STR} + 1}. \quad (24)$$

با فرض اینکه

$$\Lambda_{DR} = \frac{P_{Y_R|D}(y_R | d)P_D(d)}{P_{Y_R|A}(y_R | a)P_A(a)}, \quad (25)$$

$$\Lambda_{STR} = \frac{P_{Y_R|ST}(y_R | st)P_{ST}(st)}{P_{Y_R|A}(y_R | a)P_A(a)}, \quad (26)$$

جدول (۱): مقادیر احتمال‌های پیشین طبق مرجع [۱۳]

نوع احتمال پیشین	مقدار احتمال پیشین
$P_D(d)$	۰/۰۲۱
$P_{ST}(st)$	۰/۲۲
$P_A(a)$	۰/۷۵۹

در این معادلات، احتمال‌های  $P_D(d)$ ،  $P_{ST}(st)$  و  $P_A(a)$  به ترتیب نشان دهنده احتمال‌های پیشین در زمانی است که بین فرکانسی گفتار معین، بین فرکانسی گفتار تصادفی و بین فرکانسی گفتار غایب اند (سکوت). برای محاسبه این احتمال‌ها فرض می‌شود برای یک گفتار انگلیسی متوسط دوره گفتار مصوت ۷۸٪ از زمان است، فرکانس اساسی گفتار  $f_0$  نیز بین ۵۰ و ۵۰۰ هرتز است و برای بیشتر صداهای مصوت گفتار انرژی گفتار عمدتاً تا حدود  $f_c = 2000$  Hz خواهد بود. حال احتمال‌ها مطابق رابطه‌های (۲۷)، (۲۸) و (۲۹) محاسبه می‌شوند:

$$P_D(d) = 0.78 \times \frac{f_c}{f_0} \times \frac{K}{2}, \quad (27)$$

$$P_{ST}(st) = 0.22, \quad (28)$$

$$P_A(a) = 1 - P_D(d) - P_{ST}(st), \quad (29)$$

$K$  طول فریم است. برای فرکانس نمونه برداری ۱۶ کیلوهرتز، طول فریم ۴۸۰ نمونه و فرکانس اساسی ۳۰۰ هرتز مقادیر احتمال پیشین طبق جدول (۱) به دست می‌آیند

انتخاب شد. تخمین سیگنال گفتار تمیز در بین فرکانسی  $k$  و فریم  $i$  به صورت رابطه (۲۰) است:

$$\hat{s}(k, i) = E[Y(k, i)] \approx \frac{1}{(n_2 + n_1 + 1)} \times \sum_{m=i-n_1}^{i+n_2} y(k, m) e^{(-d_p + j\nu_p)(i-m)M}. \quad (20)$$

### ۳-۳- تخمین کمترین میانگین مربعات خطا

#### تحت مدل قطعی - تصادفی گفتار

برای پیدا کردن تخمینگر کمترین میانگین مربعات خطا تحت یک مدل ترکیبی قطعی - تصادفی گفتار، یک مدل کاملاً عمومی استفاده می‌شود که در آن تصمیم‌گیری براساس احتمال بین مدل قطعی و تصادفی صورت می‌گیرد و عدم قطعیت حضور گفتار نیز در نظر گرفته شده است. در این مدل ابتدا مجموعه  $\alpha = \{A, D, ST\}$  معرفی می‌شود.  $\alpha = A$ ،  $\alpha = D$  و  $\alpha = ST$  به ترتیب نشان‌دهنده حضور نداشتن گفتار، تولید گفتار با مدل قطعی و تولید گفتار با استفاده از مدل تصادفی است. تمام احتمالات در این بخش در بین فرکانسی  $k$  و فریم  $i$  صادق است که برای سادگی از نوشتن آن صرف نظر شده است.

برای پیدا کردن تخمینگر بهینه کمترین میانگین مربعات خطا، ابتدا عبارت شرطی  $E[S | y]$  محاسبه می‌شود:

$$\hat{s} = E[S | y] = E[S_R | y_R] + jE[S_I | y_I] = \hat{s}_R + j\hat{s}_I, \quad (21)$$

$$\hat{s}_R = E[S_R | y_R, d]P_{D|Y_R}(d | y_R) + E[S_R | y_R, st]P_{ST|Y_R}(st | y_R). \quad (22)$$

زمانی که  $\alpha = A$  است،  $s_R = 0$  خواهد بود. احتمال‌های شرطی  $P_{D|Y_R}(d | y_R)$  و  $P_{ST|Y_R}(st | y_R)$  طبق قضیه بیز به صورت زیر محاسبه می‌شوند:

$$P_{D|Y_R}(d | y_R) = \frac{\Lambda_{DR}}{\Lambda_{DR} + \Lambda_{STR} + 1}, \quad (23)$$

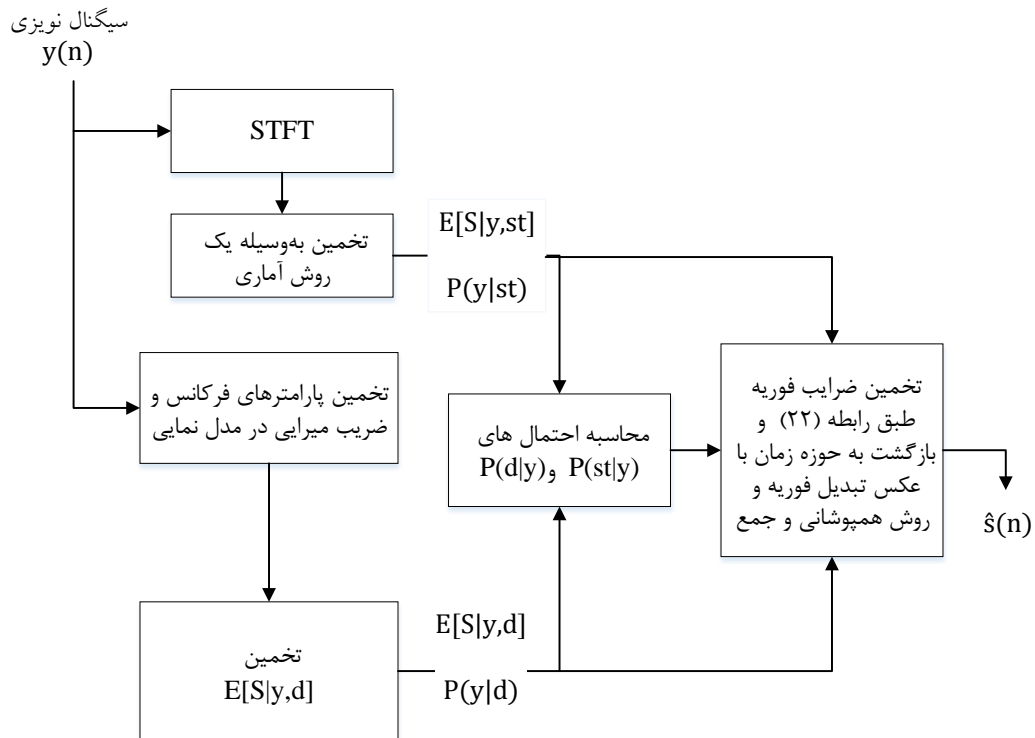
بلوک دیاگرام روش بهسازی به روش ترکیب مدل

[۱۳]. در روابط (۲۵) و (۲۶) احتمال  $P_{Y_R|A}(Y_R|a)$  مطابق

قطعی و تصادفی در شکل (۳) آمده است.

با رابطه (۳۰) است:

$$P_{Y_R|A}(Y_R|a) = \frac{1}{\sqrt{2\pi\sigma_{W_R}^2}} \exp\left\{-\frac{(Y_R)^2}{2\sigma_{W_R}^2}\right\}. \quad (30)$$



شکل (۳): بلوک دیاگرام بهسازی گفتار با روش قطعی - تصادفی

$$P_{Y_R|ST}(Y_R(k,i)|st) = \frac{\exp\left(-\frac{\xi_R}{2}\right) \left(\frac{\xi_R(\nu-2)}{2}\right)^{(\nu+1)/2}}{\sigma_{W_R} \sqrt{\xi_R} \pi^{(\nu-2)} \Gamma\left(\frac{\nu}{2}\right)} \times \sum_{j=0}^{\infty} \frac{1}{j!} \left(\frac{\xi_R}{2}\right)^j \Gamma\left(\frac{\nu+1}{2}\right) \psi\left(\frac{\nu+1}{2}, \frac{\nu}{2}+1-j; \frac{\xi_R(\nu-2)}{2}\right). \quad (31)$$

همچنین، برای توصیف  $E[Y_R(k,i)]$  در رابطه (۱۵)، مدل نمایی تعمیم داده شده در بخش ۳ به کار رفته است؛ با این تفاوت که از روش ماتریس پنسیل برای تخمین فرکانس و ضریب میرایی استفاده شده و  $P$  نیز برابر ۴۰ در نظر گرفته شده است. با این فرضیات، سیگنال گفتار در مدل قطعی مطابق رابطه (۳۲) تخمین زده می‌شود:

#### ۴- روش پیشنهادی بهسازی با ترکیب مدل

##### نمایی و مدل تصادفی location-scale $t$

در این مقاله، از توزیع تصادفی location-scale  $t$  به جای توزیع تصادفی گاوسی در ترکیب با مدل قطعی نمایی استفاده شده است. در این حالت، با توجه به گاوسی بودن نویز، روابط (۱۵) و (۳۰) به قوت خود باقی خواهند ماند و فقط رابطه  $P_{Y|ST}(\cdot)$  با فرض توزیع جدید تغییر می‌کند که مطابق با رابطه (۳۱) استخراج شده است:

زمانی مربعات خطای تخمین تبدیل فوریه زمان کوتاه سیگنال گفتار تمیز را محاسبه می‌کند:

$$MSE = \frac{\sum_{i=1}^{Nframe} \left( \sum_{k=1}^K |E[Y(k,i)] - S(k,i)|^2 \right)}{Nframe}, \quad (33)$$

$K$  تعداد بین فرکانسی،  $Nframe$  تعداد فریم‌ها و  $S(k,i)$  تبدیل فوریه زمان کوتاه سیگنال گفتار تمیز است. سه مدل قطعی بررسی شده در این مقاله، شامل مدل سینوسی (SIN\_model)، مدل نمایی با پارامتر  $P = 1$  و مدل نمایی با پارامتر  $P = 40$  (R\_EXP\_model) و مدل نمایی (EXP\_model) هستند.

در نخستین شبیه‌سازی، سیگنال گفتار یک دقیقه‌ای از دیتابیس TIMIT با فرکانس نمونه برداری ۱۶ کیلوهرتز به کار رفته است [۲۲]. همچنین، طول فریم ۴۸۰ نمونه‌ای با فریم شیفت ۵۰٪، تبدیل فوریه ۲۰۴۸ نقطه‌ای و نویز سفید در سیگنال به نویزهای مختلف استفاده شده است. مقدار  $M$  در رابطه (۳۲) نصف طول فریم،  $n_1 = 1$ ،  $n_2 = 0$  قرار داده شده است و  $v_p$  و  $d_p$  نیز از الگوریتم ماتریس پینیل تخمین زده می‌شوند. همان‌طور که در شکل (۴) مشخص است، مدل پیشنهادی نمایی با پارامتر  $P = 40$ ، به لحاظ کمترین میانگین مربعات خطای تخمین، بهترین نتیجه را دارد.

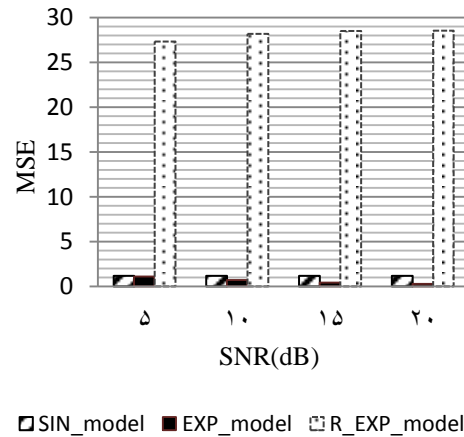
## ۵-۲- مقایسه روش‌های بهسازی گفتار با مدل‌های ترکیبی قطعی - تصادفی و مدل‌های تصادفی

در این بخش، الگوریتم بهسازی گفتار تحت مدل قطعی - تصادفی با سایر روش‌های بهسازی معمول مقایسه می‌شود که متکی به مدل تصادفی اند. برای ارزیابی نیز از معیارهای اندازه‌گیری ارزیابی ادراکی کیفیت گفتار<sup>۱۳</sup> (PESQ) و نسبت سیگنال به نویز قطعه‌ای<sup>۱۴</sup> (segSNR) طبق پارامترهای گفته شده در [۱۰] استفاده می‌شود. داده‌های گفتاری بررسی شده در آزمایشات حدود شش دقیقه سیگنال

$$\hat{s}(k,i) = E[Y(k,i)] \approx \frac{1}{P(n_2 + n_1 + 1)} \times \sum_{m=i-n_1}^{i+n_2} \sum_{i=1}^P y(k,m) e^{(-d_p + jv_p)(i-m)M}. \quad (32)$$

## ۵- شبیه‌سازی و نتایج

در این بخش، ابتدا بهترین مدل قطعی برای ترکیب با مدل تصادفی در بهسازی گفتار بررسی شده است و سپس بهترین مدل قطعی انتخاب شده در ترکیب با تخمینگر کمینه میانگین مربعات خطا مبتنی بر مدل تصادفی location-scale  $t$  برای بهسازی گفتار به کار می‌رود.



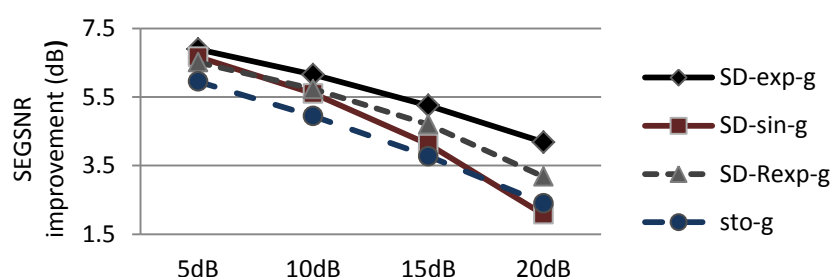
شکل (۴): خطای تخمین  $E[Y(k,i)]$  در مدل‌های قطعی مختلف.

## ۵-۱- مقایسه مدل‌های قطعی مختلف در بهسازی گفتار قطعی - تصادفی

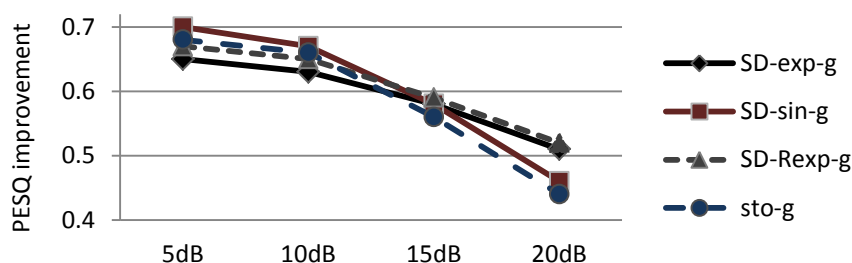
هدف در این بخش، پیدا کردن بهترین مدل قطعی برای استفاده در روش بهسازی گفتار قطعی - تصادفی است. مدل قطعی  $E[Y(k,i)]$  را تخمین می‌زند که با فرض سیگنال گفتار قطعی،  $E[W(k,i)] = 0$  و طبق رابطه (۱)،  $E[Y(k,i)] = E[S(k,i)]$  خواهد بود. برای مقایسه بهترین مدل، از رابطه (۳۳) استفاده می‌شود که میانگین

از مجموعه دادگان TIMIT هستند که به صورت تصادفی از این پایگاه داده انتخاب شده‌اند و مشتمل بر ۵۶ فایل صوتی از گویندگان مرد و ۵۸ فایل صوتی از گویندگان خانم است [۲۲]. فرکانس نمونه برداری ۱۶ کیلوهرتز، طول فریم ۴۸۰ نمونه با فریم شیفت ۵۰٪، تبدیل فوریه ۲۰۴۸ نقطه و نویزهای HF-، m109، F16، volvo، pink، white، channel از مجموعه دادگان نویز Noisex-92 برای آزمایش‌ها انتخاب شده‌اند [۲۳]. مقدار احتمال‌های اولیه

از مجموعه دادگان TIMIT هستند که به صورت تصادفی از این پایگاه داده انتخاب شده‌اند و مشتمل بر ۵۶ فایل صوتی از گویندگان مرد و ۵۸ فایل صوتی از گویندگان خانم است [۲۲]. فرکانس نمونه برداری ۱۶ کیلوهرتز، طول فریم ۴۸۰ نمونه با فریم شیفت ۵۰٪، تبدیل فوریه ۲۰۴۸ نقطه و نویزهای HF-، m109، F16، volvo، pink، white، channel از مجموعه دادگان نویز Noisex-92 برای آزمایش‌ها انتخاب شده‌اند [۲۳]. مقدار احتمال‌های اولیه



(الف)



(ب)

شکل (۵): مقایسه عملکرد مدل تصادفی گاوسی و مدل ترکیبی گاوسی / قطعی برای بهسازی در حضور نویز سفید (الف) بهبود معیار segSNR نسبت به حالت نویزی، (ب) بهبود معیار PESQ نسبت به حالت نویزی.

روش‌های بهسازی گفتار مختلف در این مقاله با علائم اختصاری زیر معرفی می‌شوند. بهسازی گفتار تحت مدل تصادفی گاوسی با sto-g [۲]، بهسازی گفتار تحت مدل تصادفی location-scale  $t$  با sto-tls [۱۰]، بهسازی گفتار تحت مدل گاوسی - سینوسی با SD-sin-g، بهسازی گفتار تحت مدل گاوسی - نمایی هندریکس با پارامتر  $P = 1$  با SD-Rexp-g [۱۳]، بهسازی گفتار تحت مدل گاوسی - نمایی تعمیم‌یافته با  $P = 40$  با SD-exp-g (روش

همچنین، در تخمین  $\hat{\xi}(k, i)$ ، ضریب هموارکنندگی  $\hat{\xi}(k, i) = \min(\hat{\xi}(k, i), -15 \text{ dB})$  و  $\alpha = 0.98$  در نظر گرفته می‌شود. گفتنی است پارامترهای طول فریم، فریم شیفت و  $M$ ، طبق مرجع [۱۳] و پارامترهای  $P$ ،  $n_1$  و  $n_2$  نیز با سعی و خطا برای دستیابی به عملکرد مناسب در بهسازی انتخاب شده‌اند.

می‌شود. در آزمایش بعد، مدل قطعی که به کار می‌رود، مدل نمایی تعمیم یافته پیشنهادی است که ترکیب آن با دو مدل تصادفی گاوسی و  $t$  location-scale بررسی می‌شود.

در شکل (۶) و (۷)، عملکرد الگوریتم SD-exp-tls (روش پیشنهادی دوم) به‌ازای  $v=3.1$ ، با SD-exp-g (روش پیشنهادی نخست)، sto-g [۲] و sto-tls [۱۰] مقایسه می‌شود. معیار مقایسه الگوریتم‌ها، متوسط بهبود معیار segSNR و PESQ در شش نویز مختلف و در رنج SNRهای صفر تا ۲۰ دسی بل است. مطابق با شکل (۶)، استفاده از مدل نمایی تعمیم‌یافته پیشنهادی در ترکیب با مدل تصادفی گاوسی، یعنی روش SD-exp-g، به بهبود چشمگیری در معیار سیگنال به نویز قطعه‌ای نسبت به مدل sto-g منجر شده است که اهمیت ترکیب یک مدل قطعی مناسب با مدل تصادفی را در افزایش بازدهی بهسازی گفتار نشان می‌دهد. روش sto-tls صرفاً مبتنی بر یک مدل تصادفی است و از توزیع منطبق‌تری با داده‌های گفتار تمیز بهره می‌جوید؛ بنابراین، دقت درخور قیاسی با روش SD-exp-g دارد و البته ترکیب مدل قطعی نمایی با آن، طبق روش پیشنهادی SD-exp-tls در برخی نویزها به بهبود در معیار سیگنال به نویز قطعه‌ای منجر شده است. در شکل (۷) بهبودی معیار PESQ ارزیابی شده است که در مجموع، تفاوت چشمگیری بین انواع روش‌های بهسازی مشهود نیست و از این لحاظ روش‌ها می‌توانند با یکدیگر مقایسه شوند.

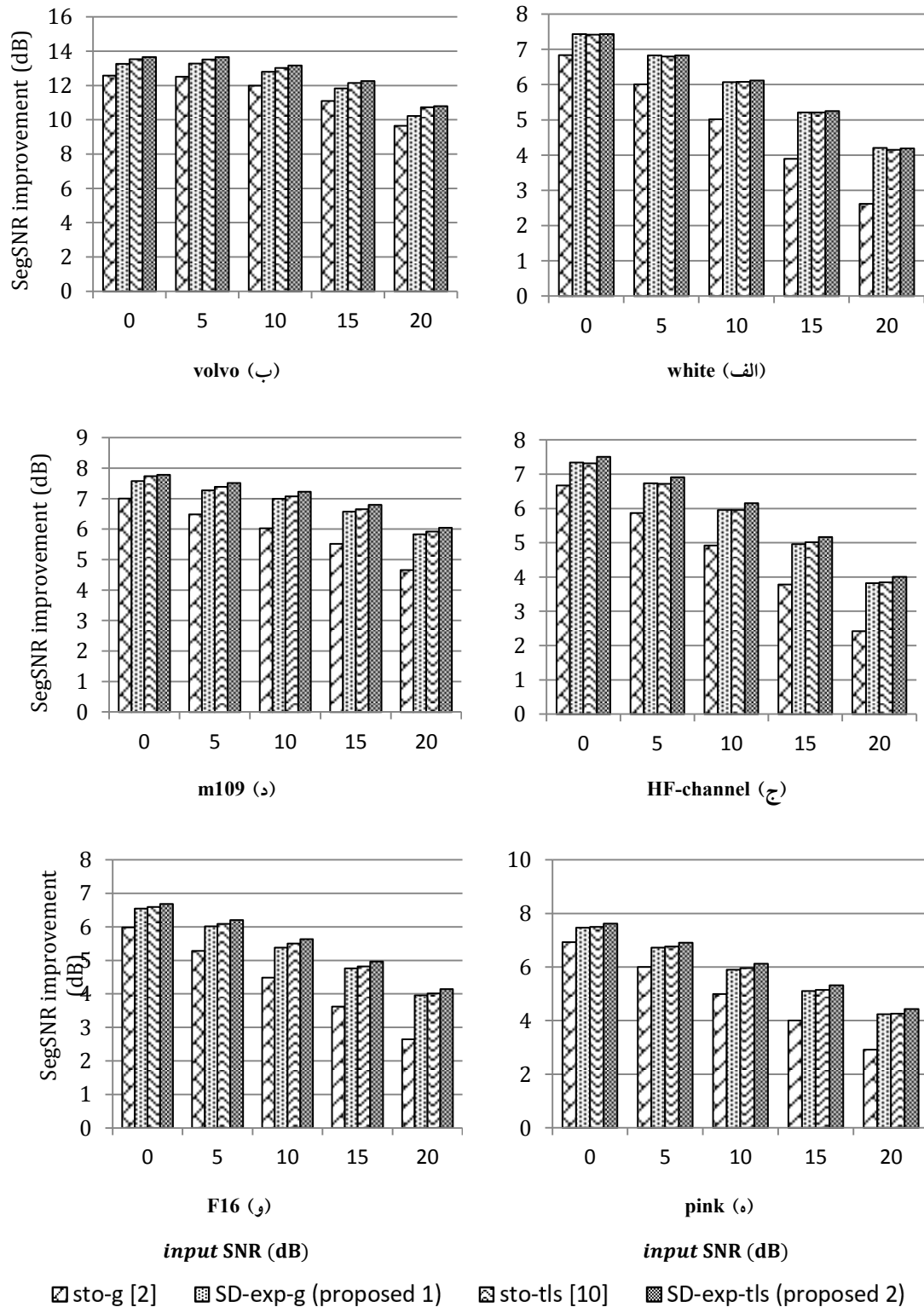
پیشنهادی نخست) و درنهایت، بهسازی گفتار با مدل  $t$  location-scale - نمایی با پارامتر  $P=40$  با SD-exp-tls (روش پیشنهادی دوم) نمایش داده می‌شوند.

### ۵-۲-۱- مقایسه روش‌های بهسازی معرفی شده تحت

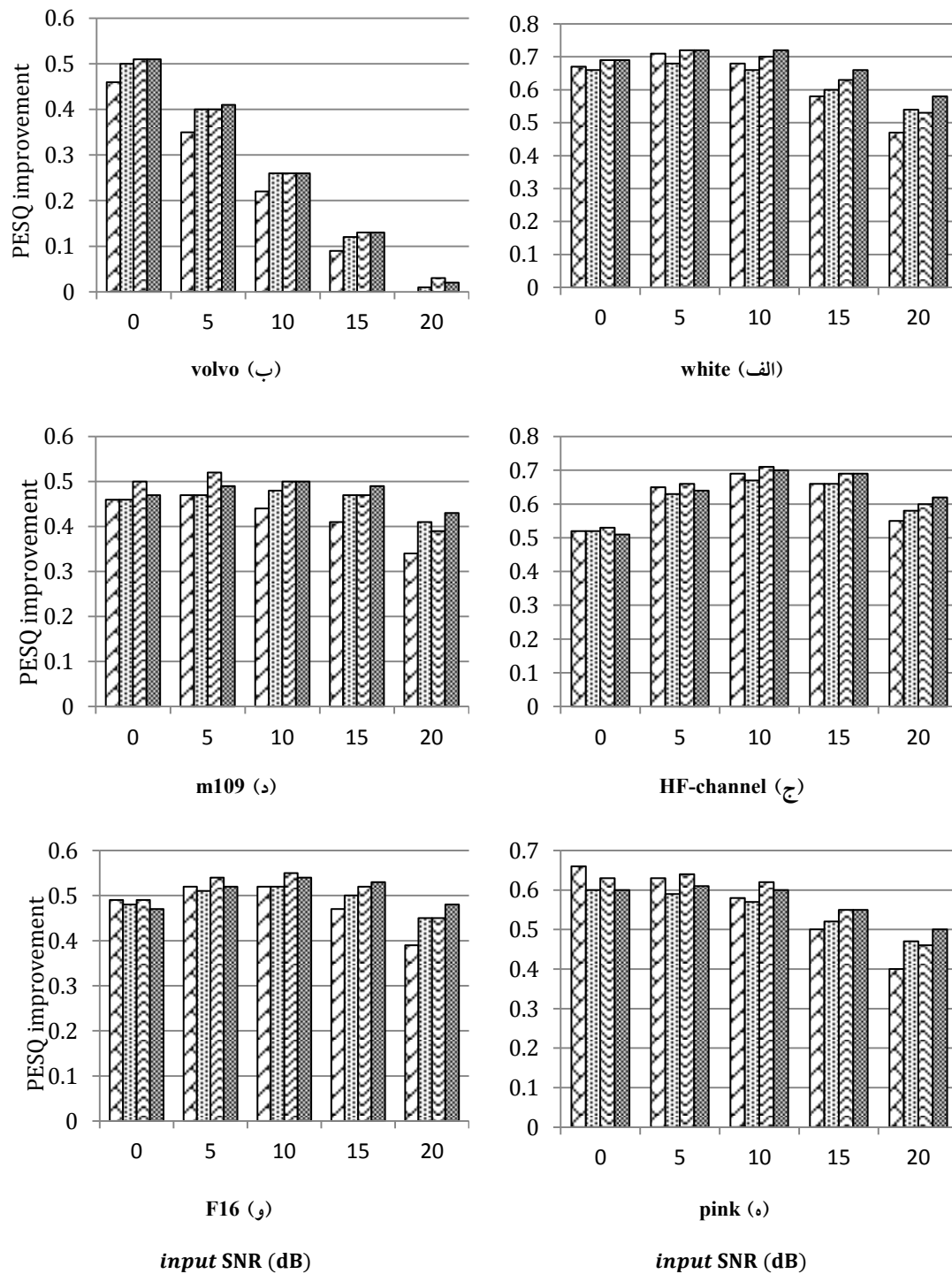
#### نویز سفید

در شکل (۵) عملکرد الگوریتم‌های SD-exp-g، SD-sin-g، SD-Rexp-g و sto-g با هم مقایسه می‌شوند که همگی بر مبنای مدل تصادفی گاوسی‌اند. در این شکل، بهبود معیارهای segSNR و PESQ در حالتی بررسی می‌شود که سیگنال گفتار شش دقیقه‌ای با نویز سفید در سیگنال به نویز ورودی در رنج ۵ تا ۲۰ دسی بل آلوده شده است.

طبق شکل (۵الف)، ترکیب مدل نمایی تعمیم یافته پیشنهادی با مدل تصادفی گاوسی، روش SD-exp-g، نسبت به روش بهسازی SD-Rexp-g به بهبود ۰/۳ دسی بلی در معیار segSNR منجر شده است. طبق شکل (۵ب)، الگوریتم SD-sin-g، بهبود ۰/۲ در معیار PESQ نسبت به الگوریتم SD-Rexp-g دارد. با توجه به اینکه معیار segSNR و PESQ را نمی‌توان هم‌زمان با هم بهتر کرد و همیشه یک چالش در بهبود هم‌زمان این دو معیار وجود دارد و البته بهبود معیار PESQ در الگوریتم SD-sin-g کمتر از ۰/۳ است، در مجموع، روش SD-exp-g روش بهتری نسبت به دو روش دیگر ارزیابی



شکل (۶): بهبود معیار segSNR نسبت به حالت نویزی در روش مدل تصادفی گاوسی (sto-g)، مدل ترکیبی گاوسی - نمایی (SD-exp-g)، مدل تصادفی t location-scale (sto-tls) و مدل ترکیبی t location-scale - نمایی (SD-exp-tls) برای بهسازی در حضور شش نویز ایستان از مجموعه داده‌های نویز Noisex-92.



sto-g [2]
  SD-exp-g (proposed 1)
  sto-tls [10]
  SD-exp-tls (proposed 2)

شکل (۷): بهبود معیار PESQ نسبت به حالت نویزی در روش مدل تصادفی گاوسی (sto-g)، مدل ترکیبی گاوسی - نمایی (SD-exp-g)، مدل تصادفی  $t$  location-scale (sto-tls) و مدل ترکیبی  $t$  location-scale - نمایی (SD-exp-tls) برای بهسازی در حضور شش نویز ایستان از مجموعه داده‌های نویز Noisex-92.

نرم افزار MATLAB روی رتبه روش‌ها در نتایج سیگنال به نویز قطعه‌ای اجرا شده، مقدار  $p$  برابر با  $10^{-17} \times 1/7$  نشان‌دهنده تفاوت معنادار چهار روش به لحاظ آماری است. همچنین، برای تأیید مؤثر بودن ترکیب مدل نمایی با مدل تصادفی تست فریدمن روی دوبه‌دوی روش‌ها انجام شده و مقدار  $p$  برابر با  $10^{-8} \times 4/3$  در مقایسه دو روش SD-exp-g و sto-g و نیز  $10^{-8} \times 4/3$  در مقایسه دو روش SD-exp-tls و sto-tls حاصل شده است. مقدار  $p$  کمتر از  $0/01$  نشان‌دهنده تفاوت معنادار روش‌های بهسازی ارزیابی شده و مؤثر بودن ترکیب مدل نمایی تعمیم‌یافته پیشنهادی در هر دو مدل تصادفی است. همچنین، انجام تست فریدمن روی دو مدل SD-exp-tls و SD-exp-g و حصول مقدار  $p$  برابر با  $10^{-7} \times 8/9$  کارایی بالاتر مدل تصادفی  $t$  location-scale را در مقابل مدل گاوسی در ترکیب با مدل نمایی پیشنهادی نشان می‌دهد. انجام تست فریدمن در معیار PESQ و مقدار  $p$  به دست آمده، تفاوت آماری معناداری را بین دوبه‌دوی روش‌ها به روال بالا نشان‌دهنده؛ به این ترتیب، برابری نسبی میانگین رتبه‌های عملکرد دو روش sto-tls و SD-exp-tls (روش پیشنهادی دوم) و نیز sto-g و SD-exp-g (روش پیشنهادی نخست) طبق جدول (۳)، کارایی نداشتن مدل ترکیبی قطعی - نمایی را در قیاس با مدل تصادفی صرف در بهبود معیار PESQ نشان می‌دهد.

## ۶- نتیجه‌گیری

در این مقاله، یک روش نوین بهسازی گفتار در حالت تک‌کاناله با استفاده از ترکیب مدل قطعی نمایی و مدل تصادفی ارائه شد. روش پیشنهاد شده، تعمیمی بر روش معرفی شده هندریکس و همکارانش در سال ۲۰۰۷ است که از توزیع تصادفی گاوسی و لاپلاس به‌عنوان مدل تصادفی و از مدل نمایی به‌عنوان مدل قطعی استفاده کردند. در این مقاله از توزیع جدید  $t$  location-scale به‌عنوان مدل تصادفی استفاده شد و همچنین، مدل نمایی استفاده شده در مقاله هندریکس با افزایش پارامتر  $P$  و تغییر روش تخمین فرکانس از اسپریت<sup>۱۶</sup> به ماتریس پینیل بهبود داده شد. نتایج پیاده‌سازی در شش نویز مختلف نشان داد روش جدید ارائه شده، یعنی ترکیب مدل نمایی

در جدول (۲) میزان کارایی متوسط چهار روش بهسازی نمایش داده شده که در شش نویز مختلف و پنج مقدار سیگنال به نویز به دست آمده است. مطابق با این جدول، روش ترکیبی گاوسی - نمایی (روش پیشنهادی نخست) در مجموع، موجب بهبود چشمگیری در حدود  $0/9$  دسی‌بل در معیار سیگنال به نویز قطعه‌ای شده است؛ البته بهبود حدود  $0/1$  دسی‌بلی روش ترکیب مدل نمایی با مدل تصادفی tls (روش پیشنهادی دوم) نسبت به مدل تصادفی صرف نیز مشهود است.

جدول (۲): مقایسه عملکرد میانگین الگوریتم‌های sto-g، sto-exp-tls و SD-exp-tls در حضور شش نویز مختلف از دادگان Noisex-92 و پنج مقدار سیگنال به نویز

الگوریتم	معیار متوسط بهبود معیار segSNR (dB)	معیار متوسط بهبود معیار PESQ
sto-g	۶/۰۸	۰/۴۸۶
SD-exp-g (proposed 1)	۷/۰۱	۰/۴۹۹
sto-tls	۷/۰۹	۰/۵۲
SD-exp-tls (proposed 2)	۷/۲۱	۰/۵۲۱

جدول (۳): رتبه میانگین الگوریتم‌های sto-tls، sto-g، SD-exp-tls و SD-exp-g در حضور شش نویز مختلف از دادگان Noisex-92 و پنج مقدار سیگنال به نویز

الگوریتم	معیار segSNR	معیار PESQ
sto-g	۳/۸۶	۲/۸۳
SD-exp-g (proposed 1)	۲/۶۶	۲/۶۶
sto-tls	۲/۱	۱/۴۶
SD-exp-tls (proposed 2)	۱/۰۳	۱/۵۶

تفاوت در معیار PESQ به‌طور متوسط حدود  $0/01$  تا  $0/03$  بوده است که از این لحاظ، روش‌ها تفاوت معناداری ندارند. برای بررسی علمی‌تر، تست فریدمن<sup>۱۵</sup> روی رتبه چهار روش بهسازی در ۳۰ شرایط مختلف آزمایش، یعنی شش سیگنال نویز مختلف و پنج مقدار سیگنال به نویز اجرا شده است. رتبه متوسط روش‌ها در هر دو معیار segSNR و PESQ در جدول (۳) نشان داده شده است. در تست فریدمن که با استفاده از تابع  $p = \text{friedman}(\cdot)$  در



- Signal Process., Vol. 28, No. 2, pp. 137–145, Apr. 1980.
- [12] J. Hardwick, C. Yoo, and J. Lim, “Speech enhancement using the dual excitation speech model”, in Proc. IEEE Int. Conf. Acoust., Speech, Signal Process., Vol. 2, pp. 367–370, 27-30 Apr 1993.
- [13] R. Hendriks, R. Heusdens, J. Jensen, “An MMSE estimator for speech enhancement under a combined stochastic-deterministic speech model”, IEEE Trans. Audio Speech Lang. Process., Vol. 15, No. 2, pp. 406–415, Jan 2007.
- [14] J. Laroche, Y. Stylianou, and E. Moulines, “HNS: Speech modification based on a harmonic+noise model”, in Proc. IEEE Int. Conf. Acoust., Speech, Signal Process., Vol. 2, pp. 550–553, 27-30 Apr 1993.
- [15] M. C. McCallum and B. J. Guillemin, “Accounting for deterministic noise components in a MMSE STSA speech enhancement framework,” in Proc. 12th Int. Symp. Commun. Inf. Technol., pp. 174–179, 2-5 Oct 2012.
- [16] M. McCallum, B. Guillemin, “Stochastic-deterministic MMSE STFT speech enhancement with general a priori information”, IEEE Trans. Audio, Speech, Lang. Process., Vol. 21, No. 7, pp. 1445–1457, July 2013.
- [17] Y. Du, J. Du, L.R. Dai, et al., “A regression approach to speech enhancement based on deep neural networks”, IEEE/ACM Trans. Audio Speech, Lang. Process., Vol. 23, No. 1, pp. 7–19, Jan 2015.
- [18] J.G. Proakis, D.G. Manolakis, Digital Signal Processing: Principles, Algorithms and Applications, Prentice Hall, 3rd edition, 1995.
- [19] K. Duda, T. P. Zielinski, “Efficiency of the frequency and damping estimation of a real value sinusoid,” IEEE Instrumentation & Measurement Magazine, Vol. 16, No. 2, pp. 48–58, Apr 2013.
- [20] T.K. Sarkar, O. Pereira, “Using the Matrix Pencil Method to Estimate the Parameters of a Sum of Complex Exponentials”, IEEE Antennas and Propagation Magazine, Vol. 37, No. 1, pp. 48-55, Feb 1995.
- [21] T.K. Moon, W.C. Stirling, Mathematical Methods and Algorithms for Signal Processing, Pearson, PAP/CDR edition, 1999.
- [22] W.M. Fisher, G.R. Doddington, K.M. Goudie-Marshall, “The DARPA speech recognition research database: specifications and status”, in Proceedings of DARPA workshop on speech recognition, pp. 93–99, 1986.
- [23] A. Varga, and H.J.M. Steeneken, “Assessment for automatic speech recognition II: NOISEX-92: a database and an experiment to study the effect of additive noise on speech recognition systems”, Speech Communication, Vol. 12, No. 3, pp. 247-251, 1993.
- تعمیم‌یافته و مدل تصادفی  $t$  location-scale به بهبود معیار segSNR می‌تواند منجر شود و کارایی درخوردی قیاسی را در معیار PESQ در مقایسه با روش‌های بهسازی دیگر نتیجه دهد.

## مراجع

- [1] Y. Ephraim, D. Malah, “Speech enhancement using a minimum mean-square error short-time spectral amplitude estimator”, IEEE Trans. on Acoust., Speech, Signal Process, Vol. 32, No. 6, pp. 1109–1121, Dec 1984.
- [2] S. Boll, “Suppression of acoustic noise in speech using spectral subtraction”, IEEE Trans. Acoust., Speech, Signal Process, Vol. 27, No. 2, pp. 113–120, Apr 1979.
- [3] K. Funaki, “Speech enhancement based on iterative Wiener filter using complex speech analysis”, 16th European Signal Processing Conference, pp. 1–5, 25-29 Aug 2008.
- [4] Y. Ephraim, D. Malah, “Speech enhancement using a minimum mean-square error log-spectral amplitude estimator”, IEEE Trans. Acoustic., Speech, Signal Process., Vol. 33, No. 2, pp. 443–445, May 1985.
- [5] R.J. Macaulay, M.L. Malpass, “Speech enhancement using a soft decision noise suppression filter”, IEEE Trans. Acoustic., Speech, Signal Process, Vol. 28, No. 2, pp. 137–145, Apr 1980.
- [6] T. Lotter, P. Vary, “Speech enhancement by MAP spectral amplitude estimation using a super-Gaussian speech model”, EURASIP Journal on Advances in Signal Processing, pp. 1110–1126, Dec 2005.
- [7] B. Chen, P.C. Loizou, “A Laplacian-based MMSE estimator for speech enhancement”, Speech Communication., Vol. 49, No. 2, pp. 134–143, Feb 2007.
- [8] R. Martin, “Speech enhancement based on minimum mean-square error estimation and super Gaussian priors”, IEEE Trans. Speech, Audio Process., Vol. 13, No. 5, pp. 845–856, Aug 2005.
- [9] J.S. Erkelens, R.C. Hendriks, R. Heusdens, et al. , “Minimum mean-square error estimation of discrete Fourier coefficients with generalized gamma priors”, IEEE Trans. Audio, Speech, Lang. Process., Vol. 15, No. 6, pp. 1741–1752, July 2007.
- [10] N. Faraji, A. Kohansal, “MMSE and maximum a posteriori estimators for speech enhancement in additive noise assuming a  $t$ -location-scale clean speech prior”, IET Signal Processing, Vol. 12, No. 4, pp. 532-543, June 2018.
- [11] R. McAulay and M. Malpass, “Speech enhancement using a soft-decision noise suppression filter”, IEEE Trans. Acoust., Speech,

<sup>1</sup> Minimum mean squared error

- <sup>2</sup> Maximum likelihood
- <sup>3</sup> Maximum a posteriori
- <sup>4</sup> Super Gaussian
- <sup>5</sup> Generalized gamma
- <sup>6</sup> Presence probability
- <sup>7</sup> Deep Neural Network
- <sup>8</sup> Prony method
- <sup>9</sup> Steiglitz-McBride method
- <sup>10</sup> Matrix pencil
- <sup>11</sup> Total least squares
- <sup>12</sup> Frequency bin
- <sup>13</sup> Perceptual Evaluation of Speech Quality
- <sup>14</sup> Segmental Signal to Noise Ratio
- <sup>15</sup> Friedman test
- <sup>16</sup> ESPRIT

