

به نام خدا



دانشگاه صنعتی اصفهان

دانشکده مهندسی برق و کامپیوتر

گزارش آزمایش شماره 3

آزمایشگاه پردازش سیگنال های دیجیتال

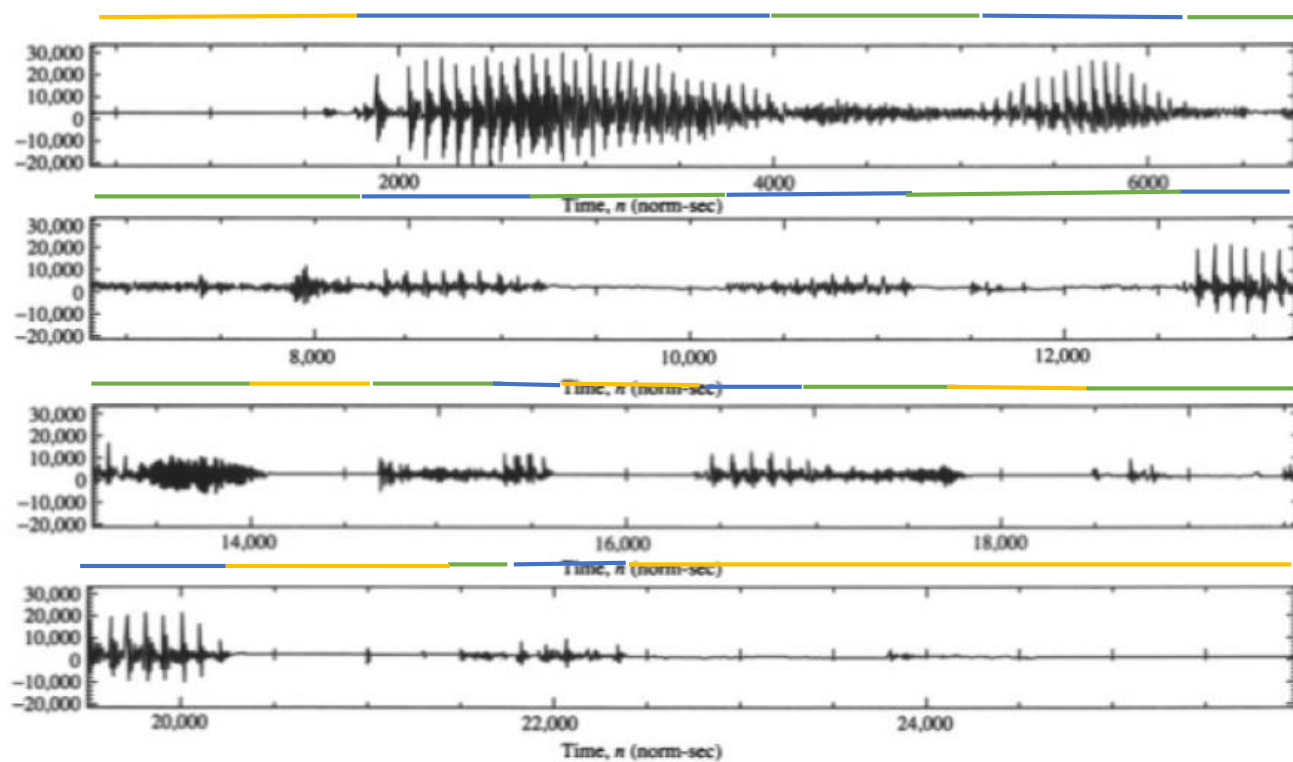
عنوان:

آشنایی با سیگنال صحبت

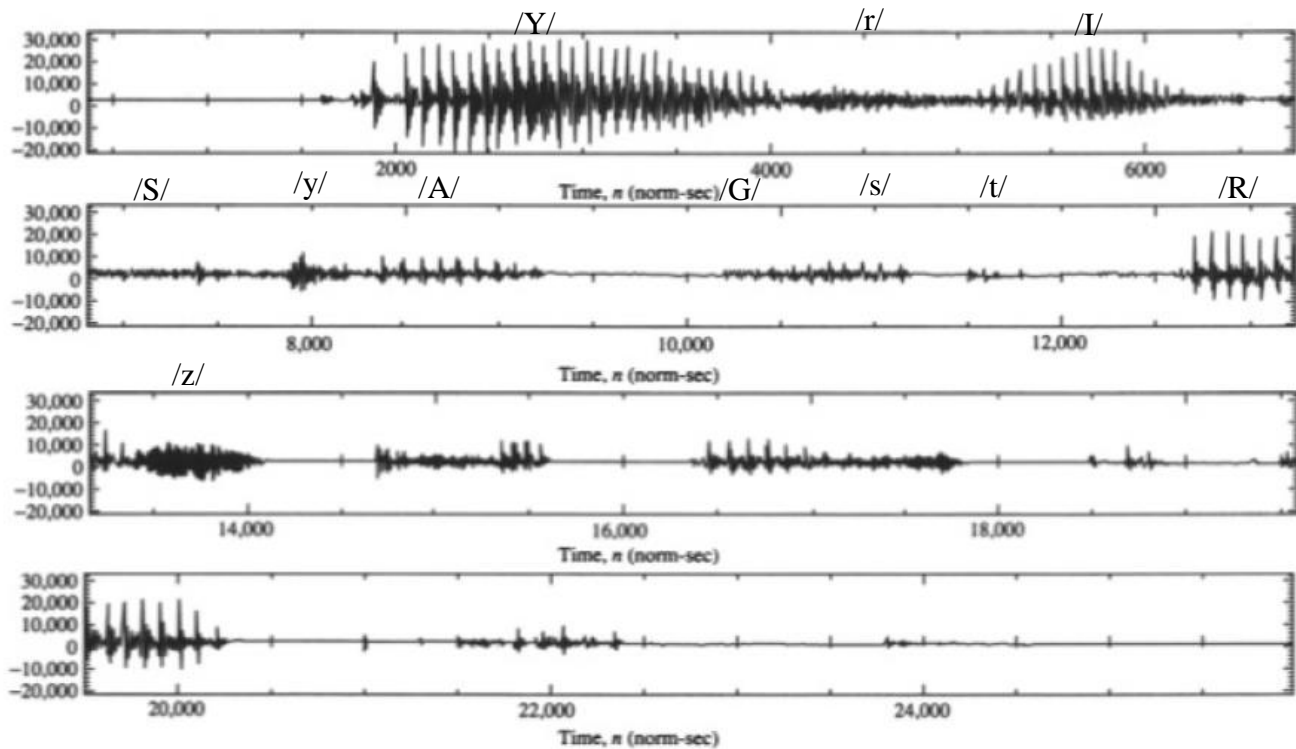
سوال 1)

الف) قسمت های مختلف را به شکل زیر برچسب گذاری می کنیم :

Voiced —————
Silence —————
Unvoiced —————



(ب)



(ج) برای هر کدام از واکه ها خواهیم داشت :

$$/Y/ : 2100-2160 \rightarrow n = 2160-2100 = 60 \rightarrow n = 60 \text{ samples} \rightarrow f_0 = \frac{8000}{60} = 133.33 \text{ Hz}$$

$$/I/ : 5490-5560 \rightarrow n = 5560-5490 = 70 \rightarrow n = 70 \text{ samples} \rightarrow f_0 = \frac{8000}{70} = 114.28 \text{ Hz}$$

$$/A/ : 8500-8590 \rightarrow n = 8590-8500 = 90 \rightarrow n = 90 \text{ samples} \rightarrow f_0 = \frac{8000}{90} = 88.89 \text{ Hz}$$

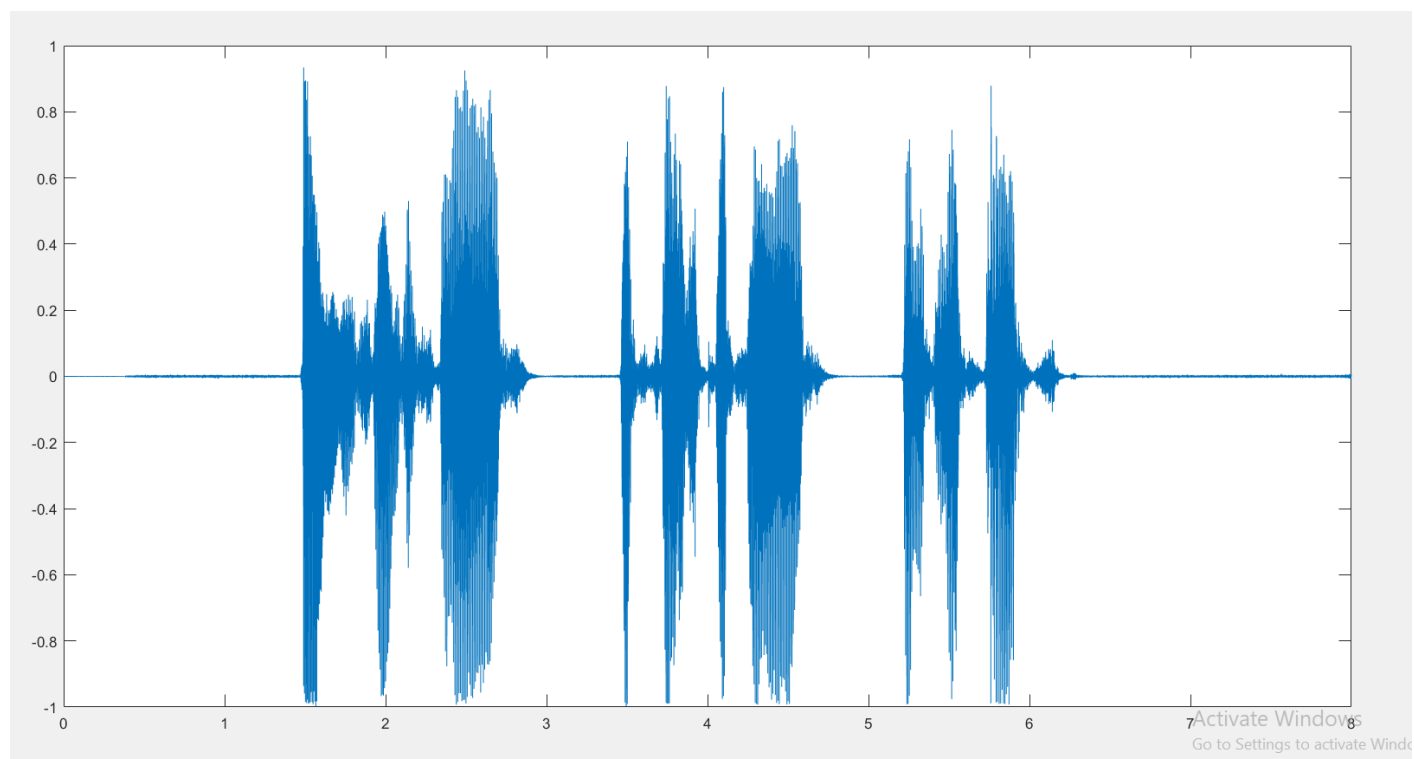
$$/R/ : 12850-12890 \rightarrow n = 12890-12850 = 40 \rightarrow n = 40 \text{ samples} \rightarrow f_0 = \frac{8000}{40} = 200 \text{ Hz}$$

$$f_{0,avg} = \frac{133.33 + 114.28 + 88.89 + 200}{4} = 134.125 \text{ Hz}$$

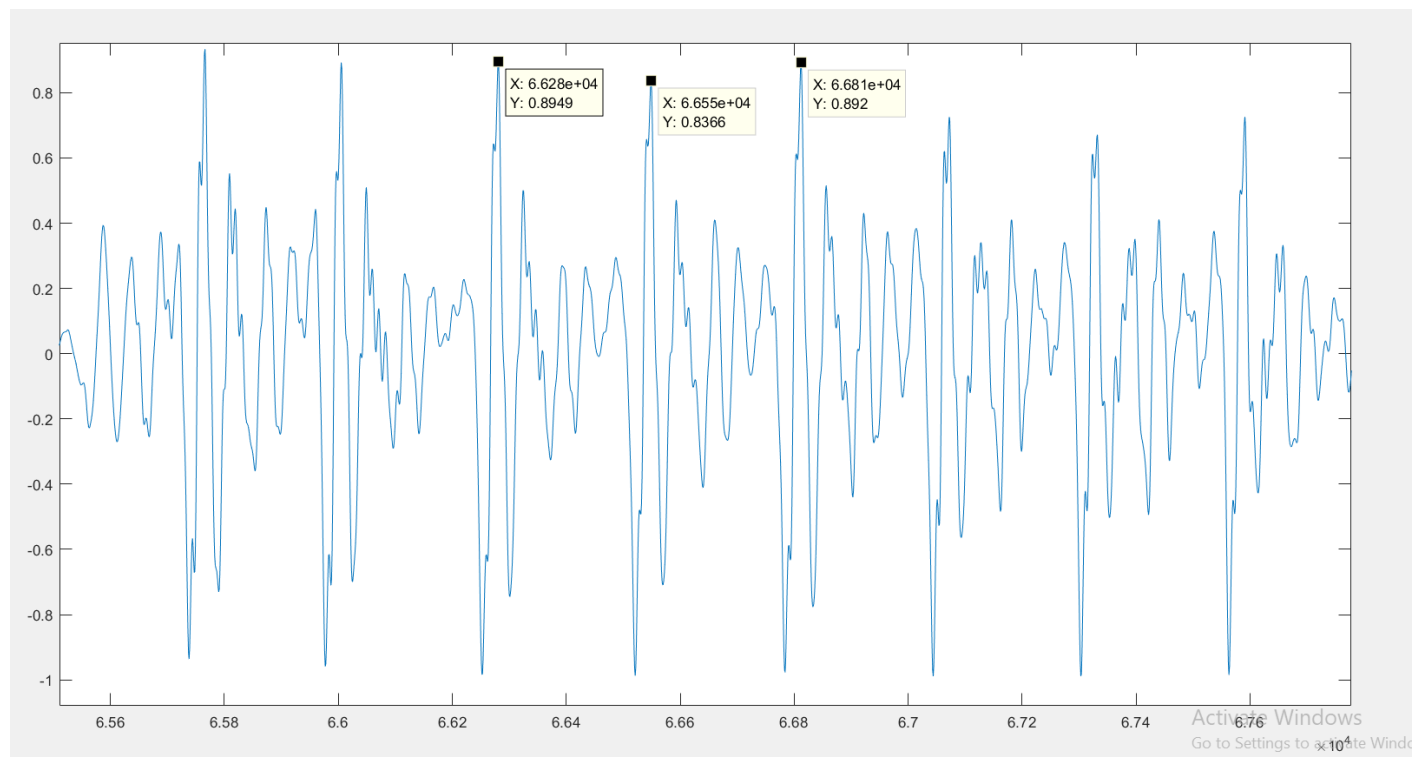
(د) با توجه به عدد بدست آمده می توان نتیجه گرفت گوینده مرد بوده است.

(ه) محاسبات بالا بر اساس همین خواسته مسئله بوده لذا می توان بازه مورد نظر برای صدا را $88.89 - 200$ خواهیم داشت که بسیار به بازه استاندارد مربوط به صدای مرد نزدیک است.

و) شکل سیگنال تلفظ شده توسط خودمان :

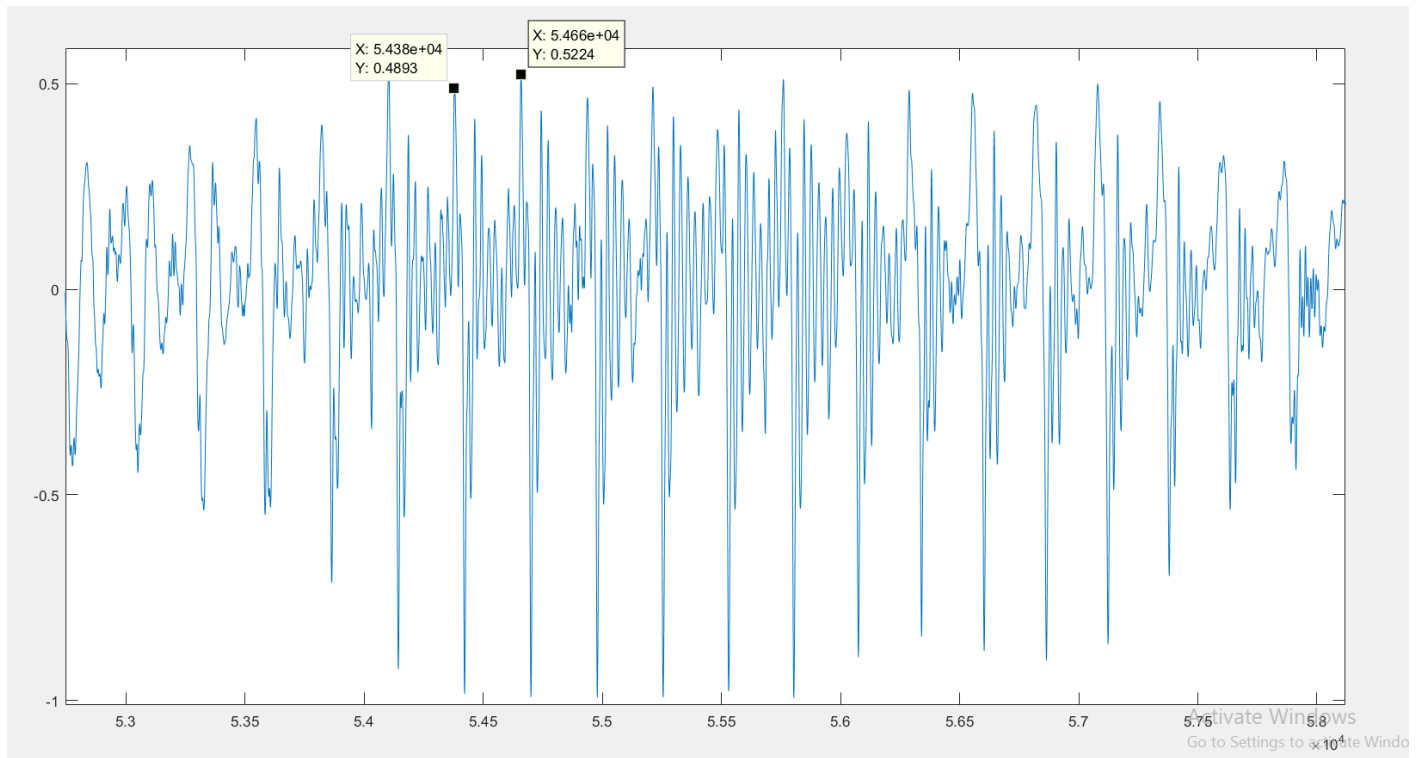


واکه /Y/ :



$$f_0 = \frac{F_s}{\text{samples per period}} = \frac{44100}{(66810 - 66550)} = \frac{44100}{260} = 169.615$$

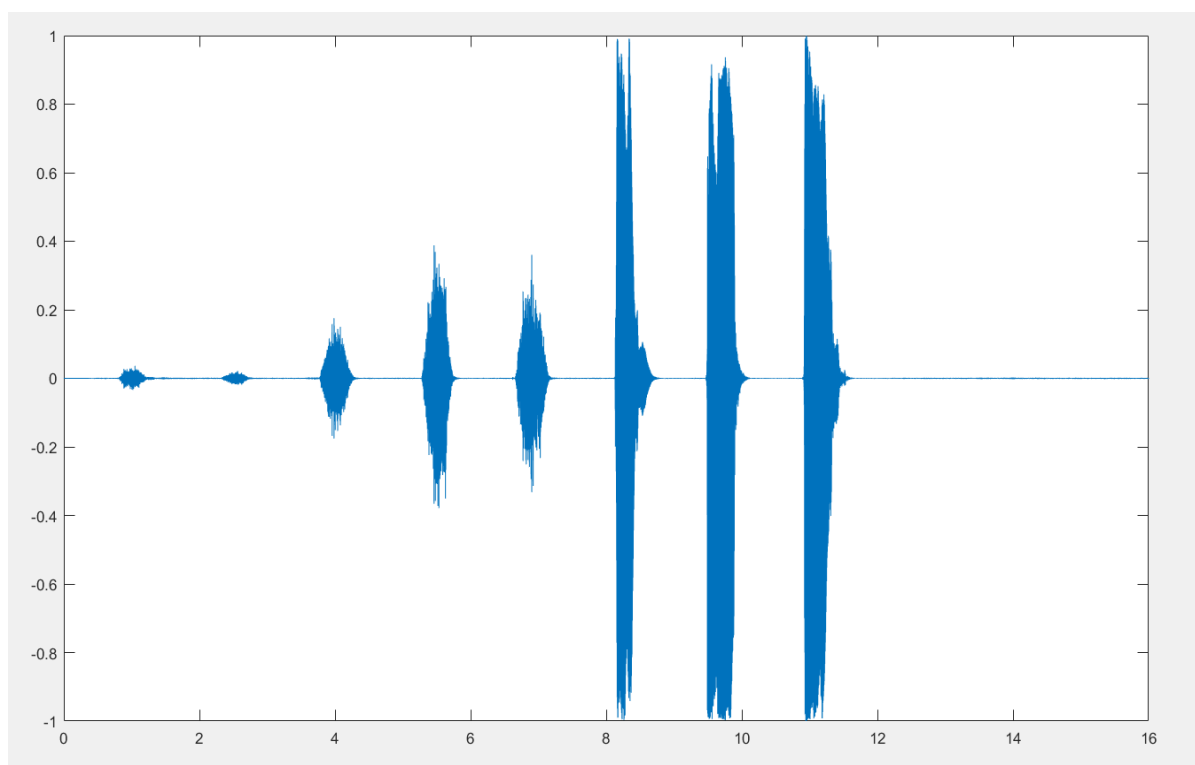
واکه /I/ :



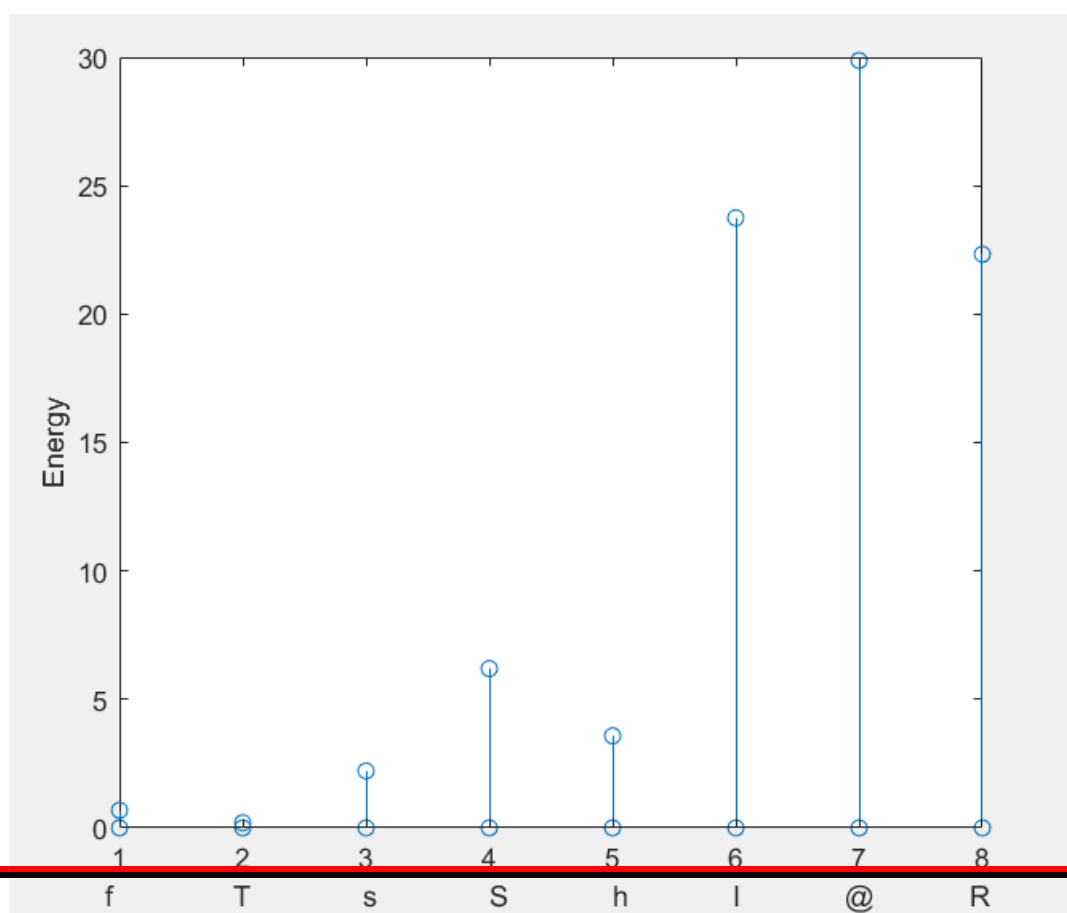
$$f_0 = \frac{F_s}{\text{samples per period}} = \frac{44100}{(54660 - 54380)} = \frac{44100}{280} = 157.5$$

همان طور که مشاهده می شود میانگین pitch ها در دو حالت تقریباً یکسان بدست می آید چون هر دو گوینده مرد می باشند.

سوال 2) شکل حروف ضبط شده در حوزه زمان :



انرژی حروف تلفظ شده را در شکل زیر مشاهده می کنیم :

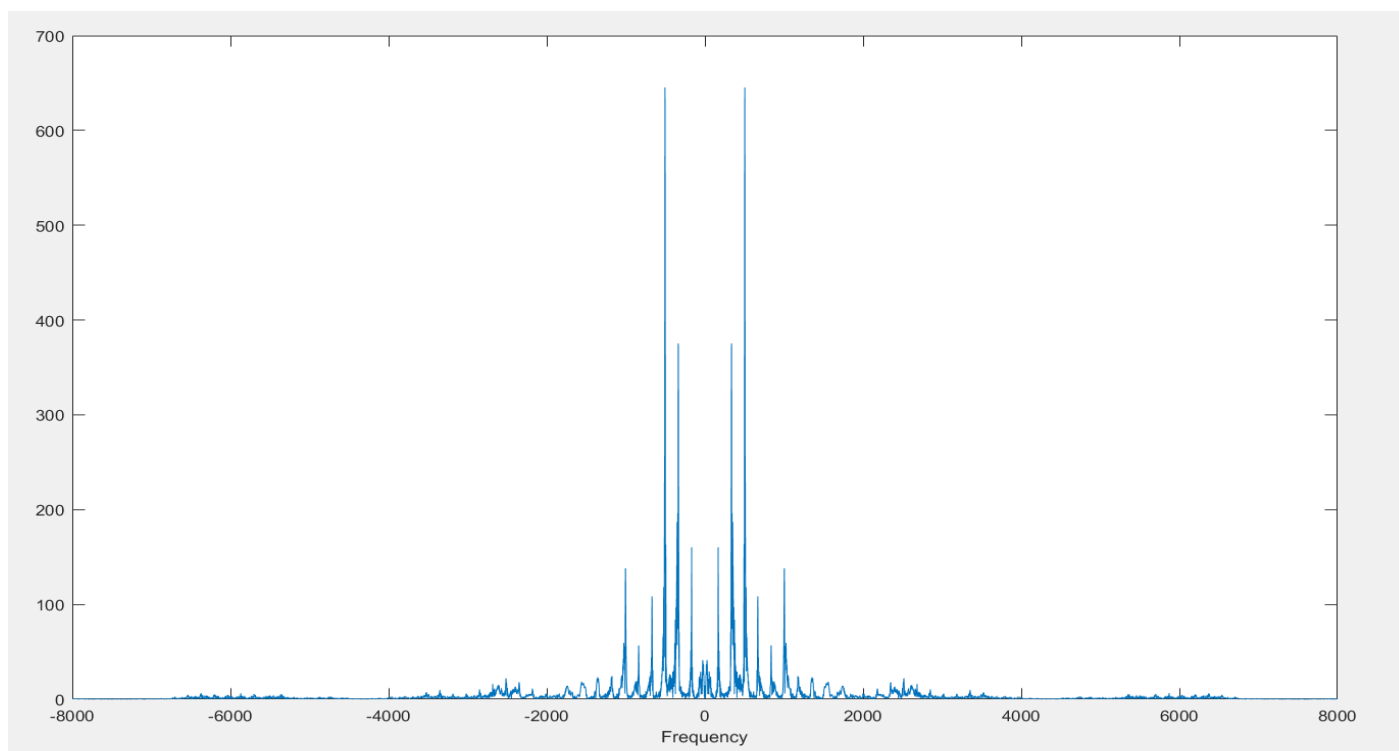
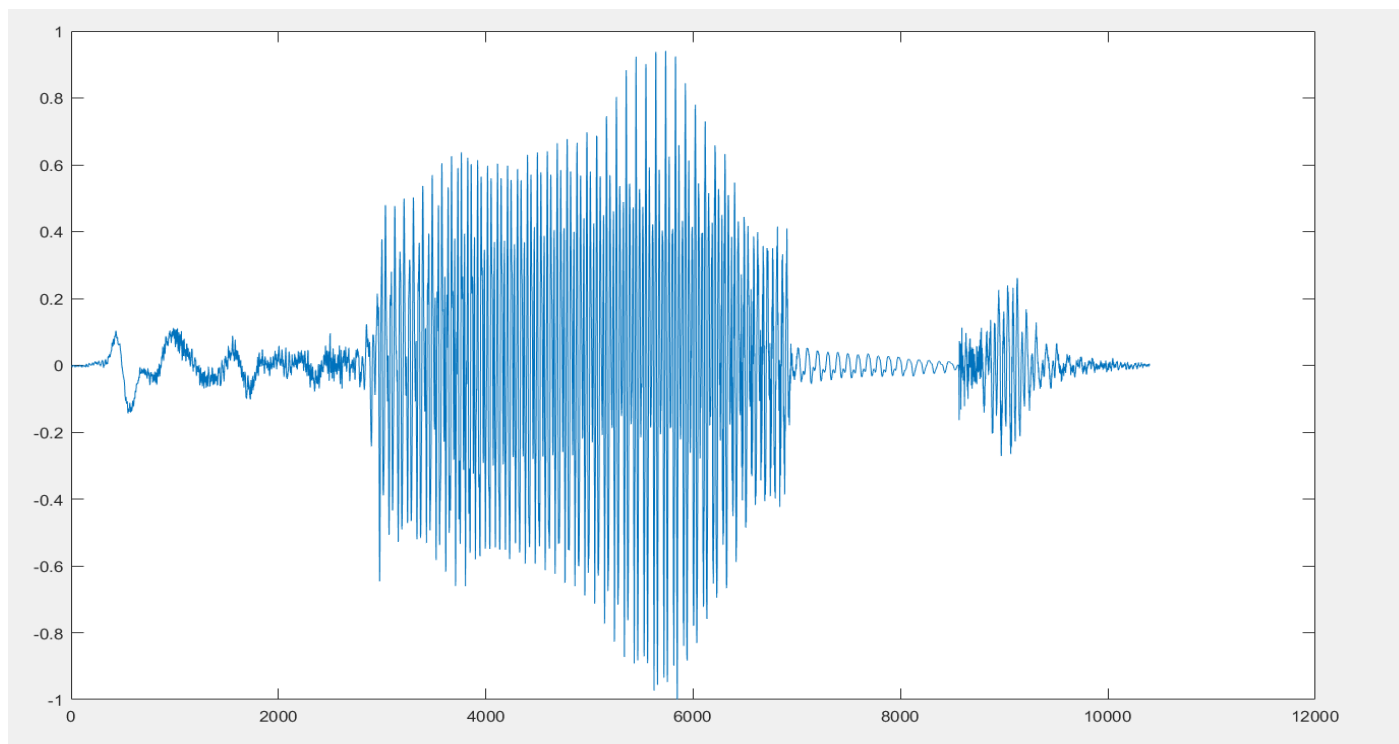


مطابق انتظار انرژی حروف سایشی نسبت به واکه ها همواره کمتر است چون که منشا تولید واکه قطار ضربه اما حروف سایشی معمولاً با هوای داخل دهان (نویز) تولید می شوند که در این بین هم آن هایی که با لب و دندان (قسمت جلویی دهان مثل f,th) نیز انرژی به مراتب کمتری نسبت به آنهایی که با قسمت میانی دهان مثل sh,h تلفظ می شوند ؛ دارند.

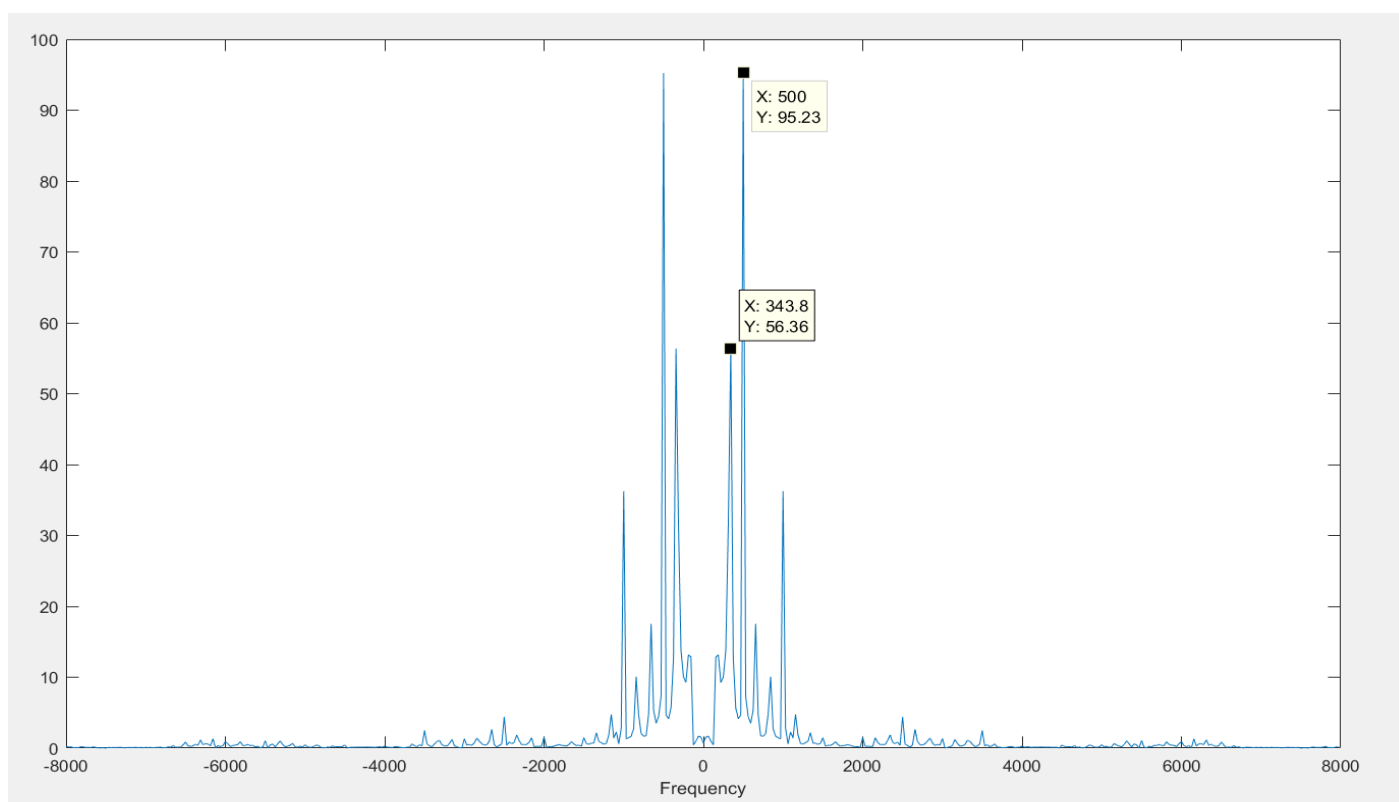
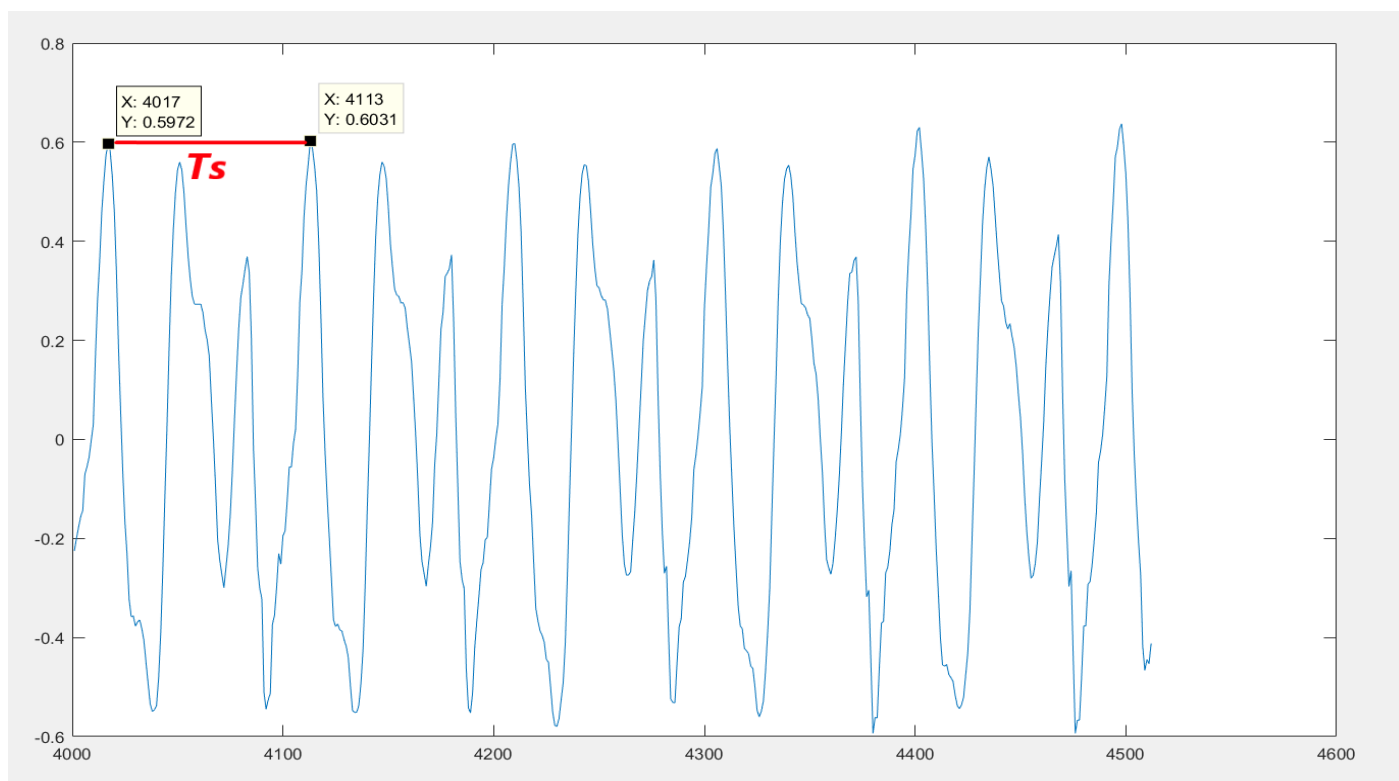
لازم به ذکر است فایل Q2.m در کدهای متلب مربوط به ضبط صدا و پخش صوت ضبط شده به همراه شکل سیگنال در حوزه زمان می باشد و فایل Q2_energy.m به محاسبه انرژی هر کدام از صداها پرداخته است.

سوال 3) در این قسمت برای نمونه سیگنال n0100.wav را به طور کامل با رسم شکل و غیره بررسی می کنیم سپس برای بقیه سیگنال های این فایل نیز مشابه این کارها صورت گرفته و در جدولی بیان خواهد شد.

شکل سیگنال مورد نظر در حوزه زمان و فرکانس :



حال برای آنکه قسمت واکه مربوطه (**oo**) را مشاهده کنیم کافی است بازه انتخابی را بین 4000 تا 4512 امین نمونه قرار دهیم (مطابق صورت سوال می‌بایست بازه ای به طول 512 از این سیگنال را انتخاب نمود).



مطابق انتظار اگر پنجره 512 نقطه ای را به صورت کامل در واکه مربوطه انتخاب کنیم سیگنال حوزه زمانی پریودیک می باشد و به کمک همین دوره تناوب یا دو پیک متوالی در حوزه فرکانس می توان pitch را محاسبه نمود. برای مثال در مورد بالا مقدار pitch از روی تبدیل فوریه سیگنال برابر با $500 - 343.8 = 156.2$ است و به کمک سیگنال حوزه زمان خواهیم داشت :

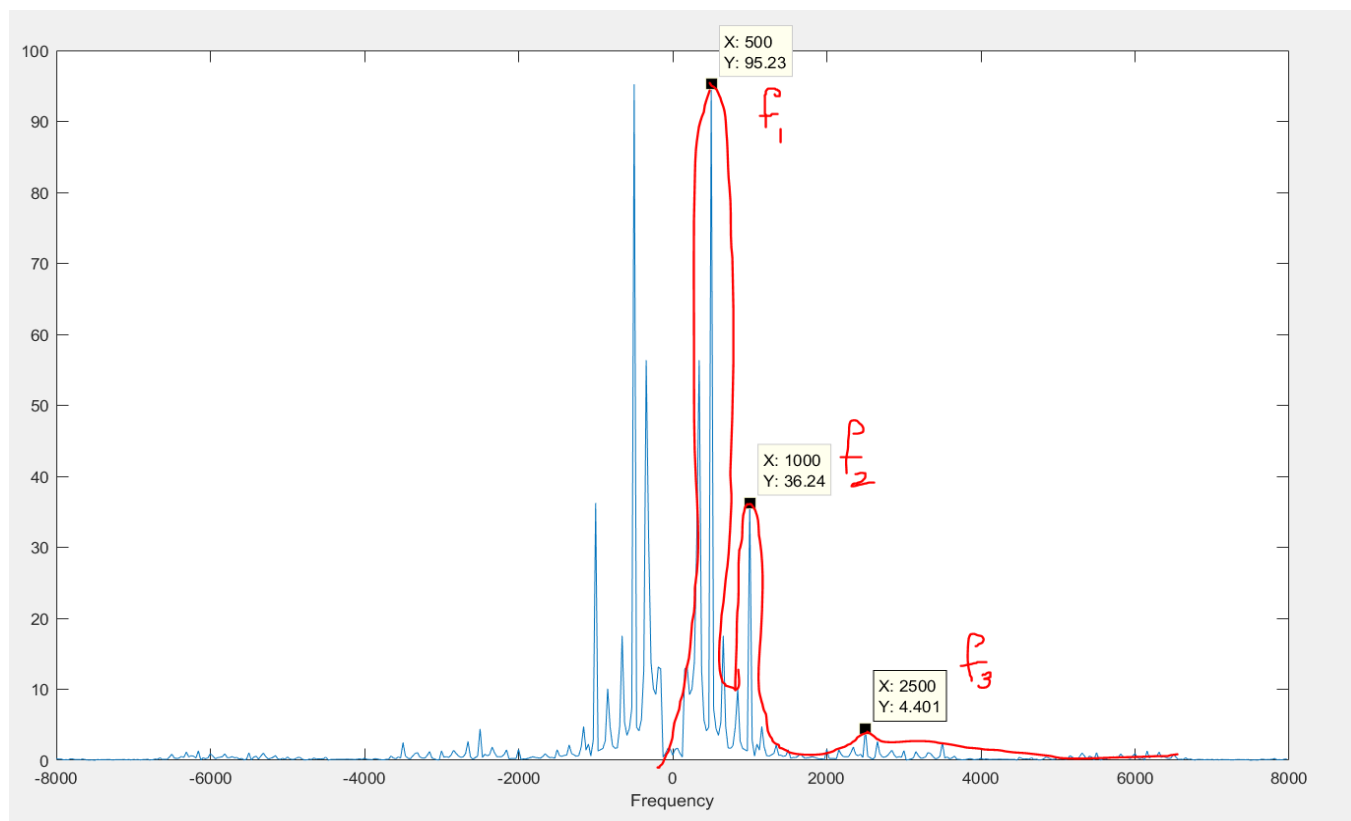
$$4113 - 4017 = 96 \text{ Samples}$$

$$\text{Hz}f_0 = \frac{16000}{96} = 166.67$$

با توجه به اعداد به دست آمده و رنج pitch صدای زن ($\text{pitch} = 120\text{-}500\text{Hz}$) و مرد ($\text{pitch} = 50\text{-}250\text{Hz}$) می توان نتیجه گرفت که گوینده مرد می باشد.

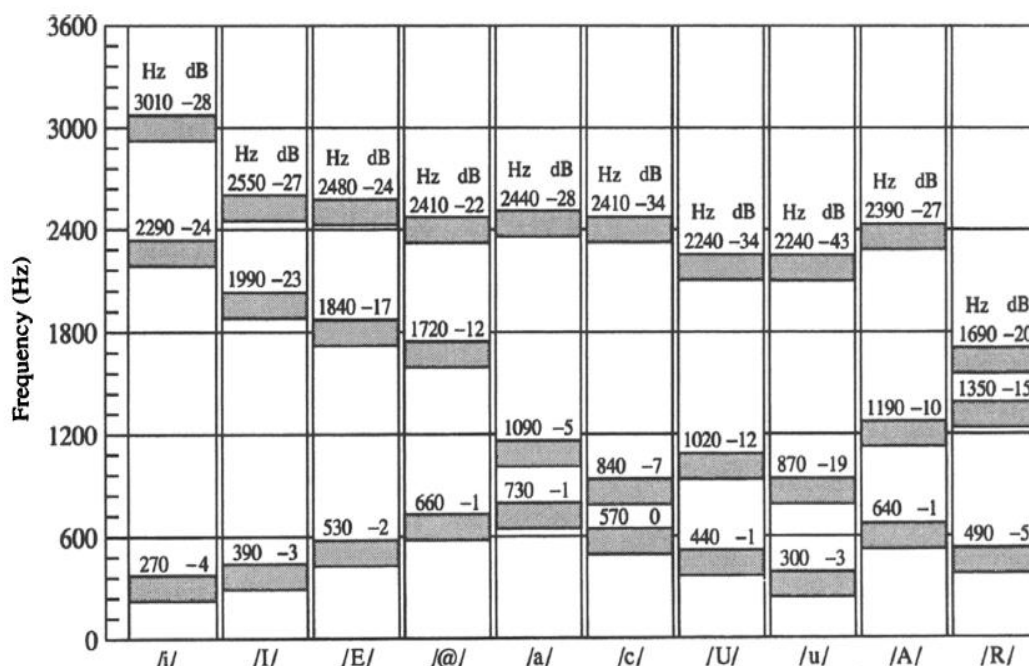
سیگنال ها در حوزه زمان مشخصه ندارند و فقط متناوب هستند.

هم چنین در شکل زیر فرمنت مربوط به این واکه مشخص شده است :



مطابق شکل بالا $f_1 = 500$ و $f_2 = 1000$ و $f_3 = 2500$. این اعداد بسیار به واکه \oo استاندارد نیز نزدیک اند.

جدول مربوط به فرمنت های تمام واکه ها :



روند فوق را برای بقیه واکه های /oo\ نیز تکرار می کنیم. نتایج مربوط به این قسمت در جدول زیر قابل مشاهده است :

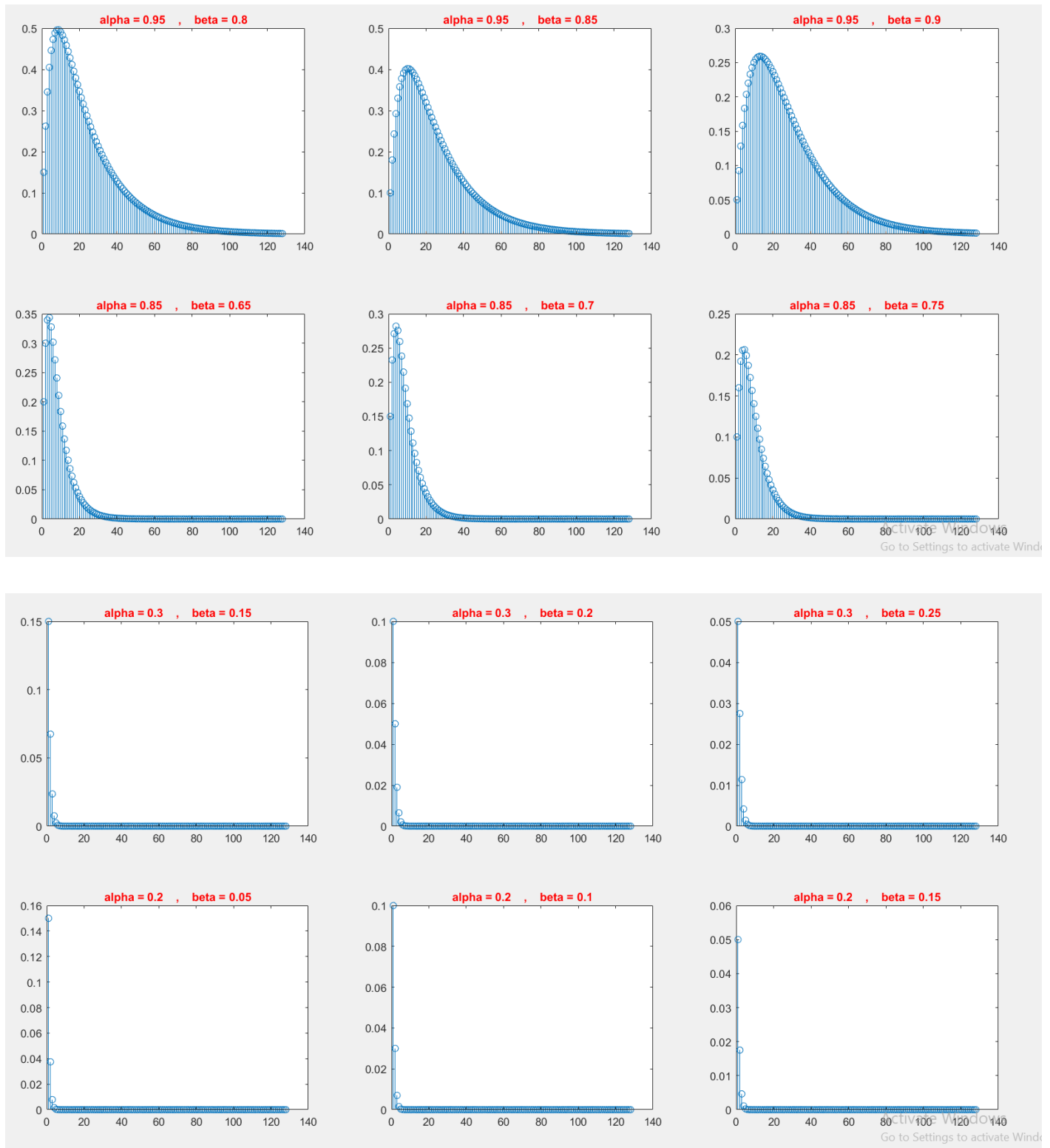
فرمنت	pitch	جنسیت	اسم فایل
437.5-1094	109.35	Male	n02oo
468.8-1188	125	Male	n06oo
437.5-1125	125	Male	n13oo
437.5-1344	125	Male	n14oo
531.3-1063	250	Female	t13oo
468.8-1438	250	Female	t14oo
687.5-1594	250	Female	t15oo
718.8-1688	250	Female	y01oo
687.5-1594	218.7	Female	y05oo
625-1281	187.5	Female	y10oo
468.8-1406	250	Female	y25oo
500 - ×	250	Female	d05oo
593.8-1563	218.8	Female	d11oo
437.5-1531	218.7	Female	d12oo
468.8-1656	250	Female	d13oo
437.5-×	218.7	Female	d20oo
468.8-×	250	Female	d21oo
625-1719	218.8	Female	d25oo
562.5-1313	187.5	Female	d26oo

اسم فایل	فرمت	حدس به کمک جدول
1	843.8-1469	/a/
2	500-1000	/U/
3	500-2781	/i/
4	500-1281	هیچ کدام (با توجه به واکه های موجود)
5	406.3-1031	/U
6	531.3-2156	/E
7	375-781.3	/c/
8	437.5-875	/c/
9	500-1219	/U/
10	437.5-2844	/i/
11	437.5-1536-1938	/R/
12	468.8-2625	/i/
13	343.8-1156	/U/
14	625-1375	/a/
15	437.5-1344	/U/
16	468.8-1063	/U/
17	250-2031	/i/
18	531.3-1813-2469	/I/
19	468.8-1250	/U/
20	468.8-2125-2719	/I/
21	750-1625	/@/
23	656.3-937.5	/a/
24	812.5-1219	/a/
25	687.5-2031	/E/
26	437.5-1563	/U/
27	468.8-937.5-1656	/a/
28	687.5-2063	هیچ کدام (با توجه به واکه های موجود)
29	437.5-1125	/u/
30	281.3-3000	/i/
31	500-2688	/E/
32	406.3-1594	/R/
33	500-2781	/E/
34	718.8-2375	/E/
35	468.8-937.5-2781	/a/
36	437.5-1563	/@/
37	681.3-1875	/@/
38	437.5-x	/U/
39	750-1781	/a/

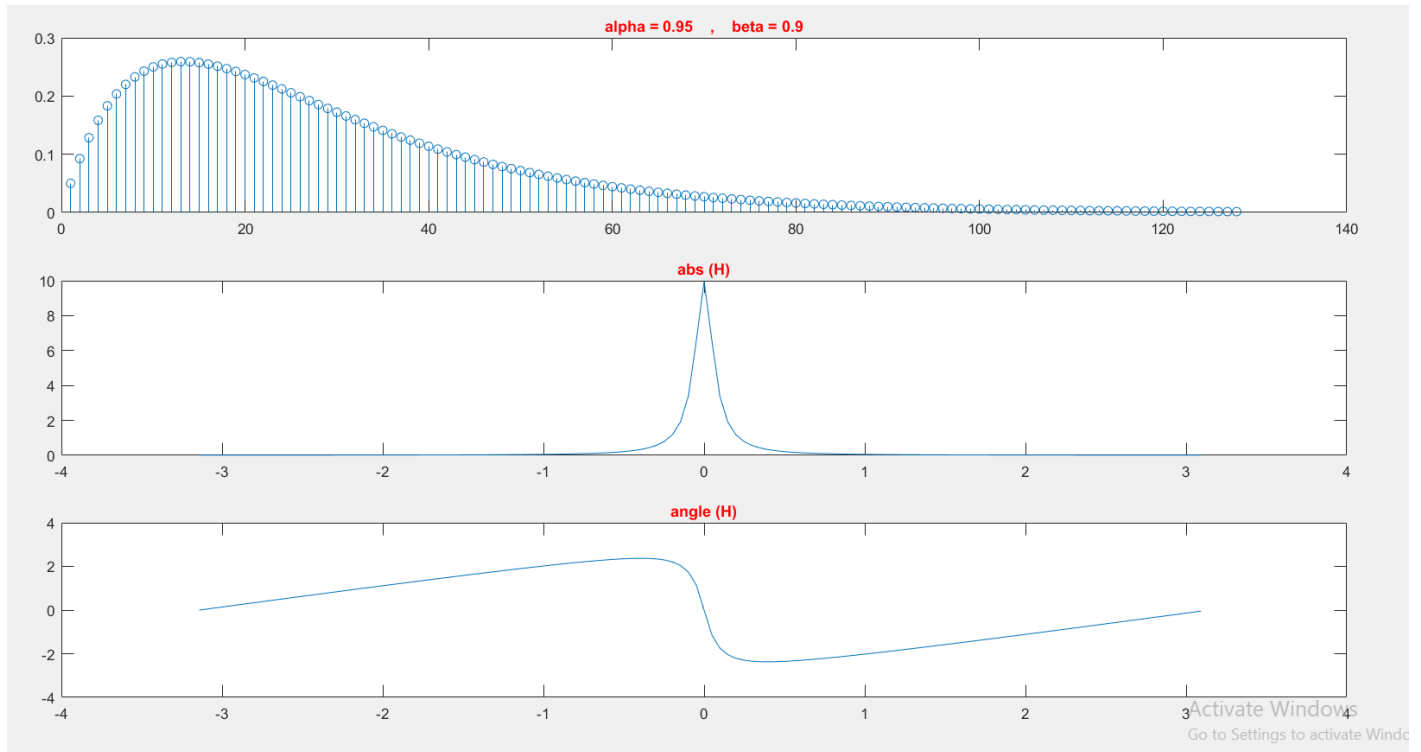
سوال 4

(a) در تمامی شکل های زیر سیگنال حوزه زمان با در نظر گرفتن فرض های مسئله ؛ میرا می شود. مطابق شکل های بالا برای آن که در لحظه $n = 64$ سیگنال زمانی برابر صفر شود می بایست آلفا کوچکتر از 0.95 باشد. اما با توجه به صورت سوال که اشاره به نزدیکی مقدار آلفا به عدد 1 دارد ؛ مقادیر کوچکتر از 0.3 برای بتا مناسب است.

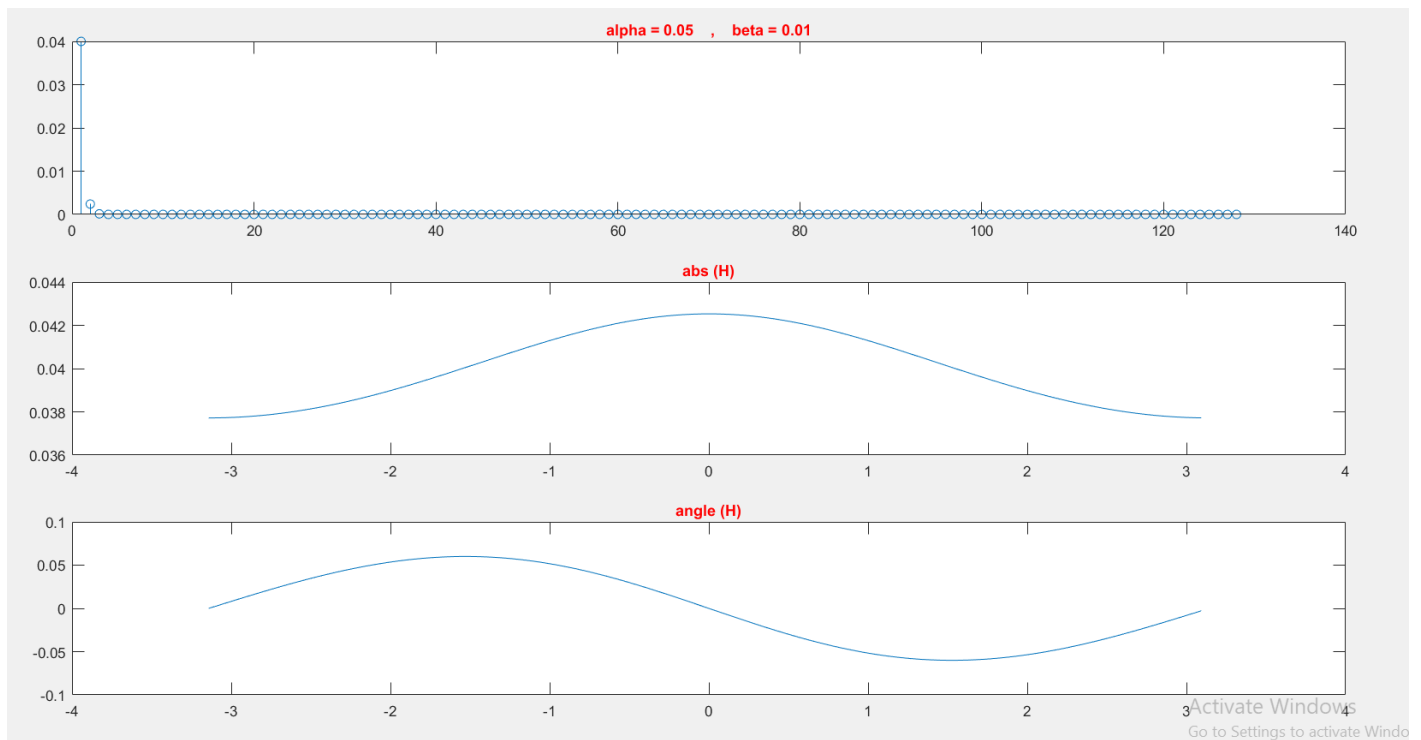
پیاده سازی $g(n) = (\alpha^n - \beta^n) u(n)$ در حوزه زمان :



(b) پیاده سازی $g(n) = (\alpha^n - \beta^n) u(n)$ در حوزه فرکانس :



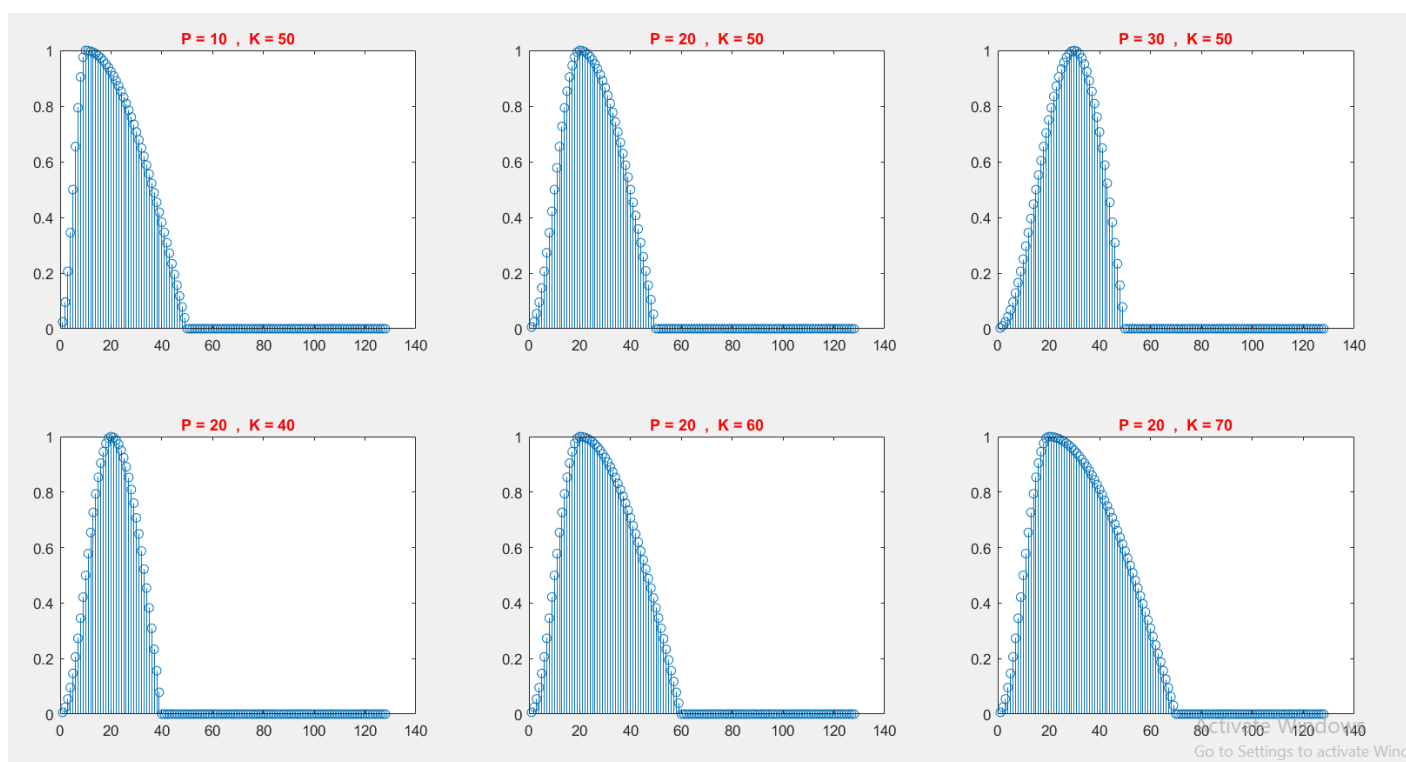
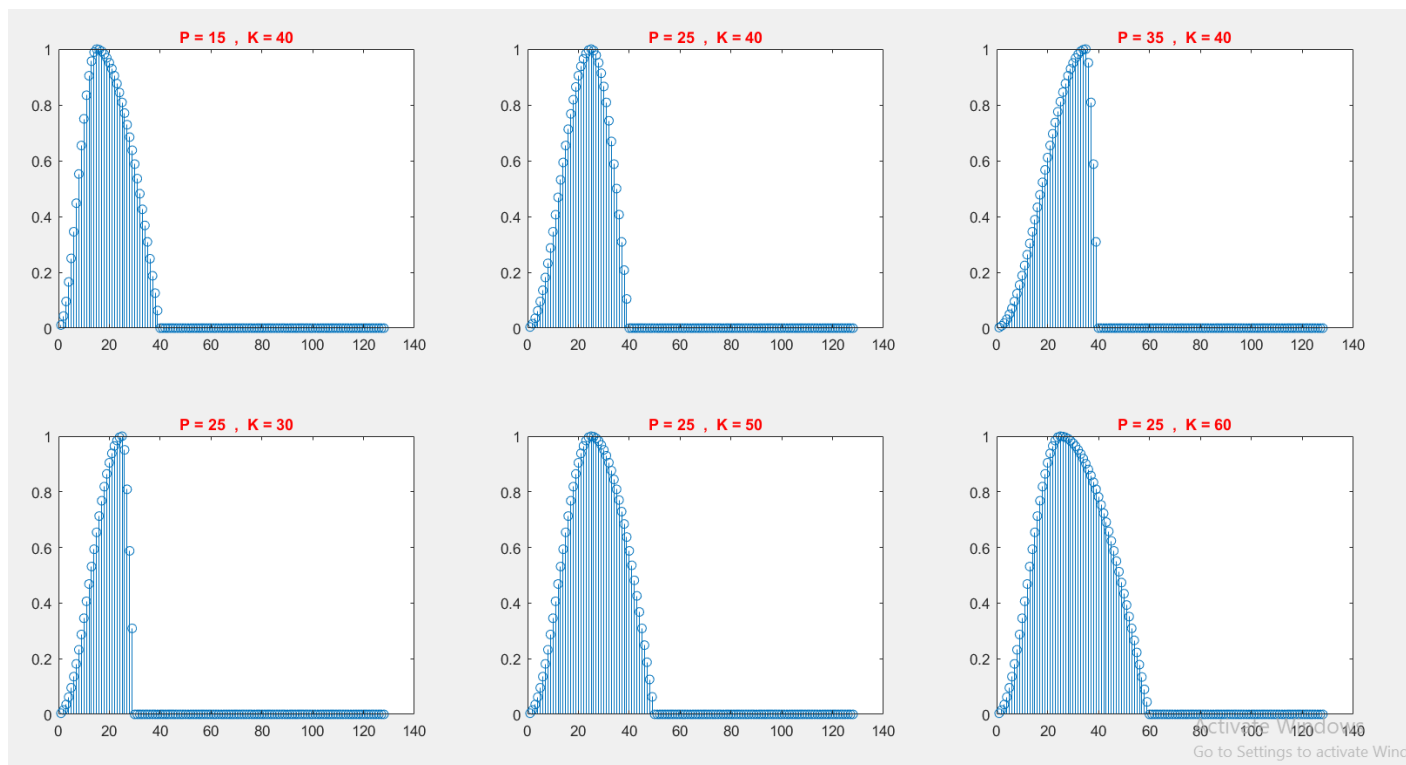
$$\alpha = 0.95 \text{ , } \beta = 0.9$$



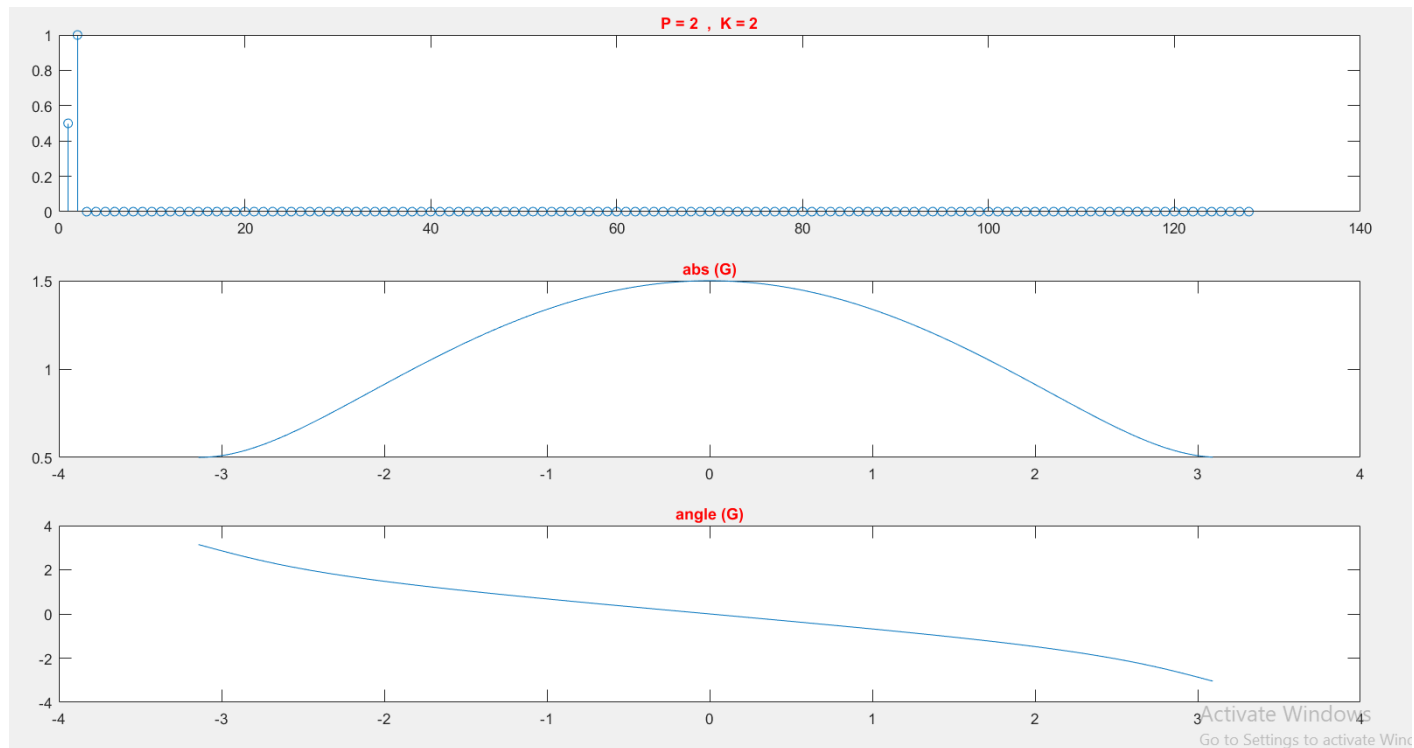
$$\alpha = 0.05 \text{ , } \beta = 0.01$$

(c) از آنجا که $P \leq K$ و طبق ضابطه سوال برای مقادیر بزرگتر از K حاصل تابع برابر صفر می‌باشد لذا کافی است $K < 65$ انتخاب شود تا در لحظه $n = 64$ سیگنال زمانی برابر صفر شود.

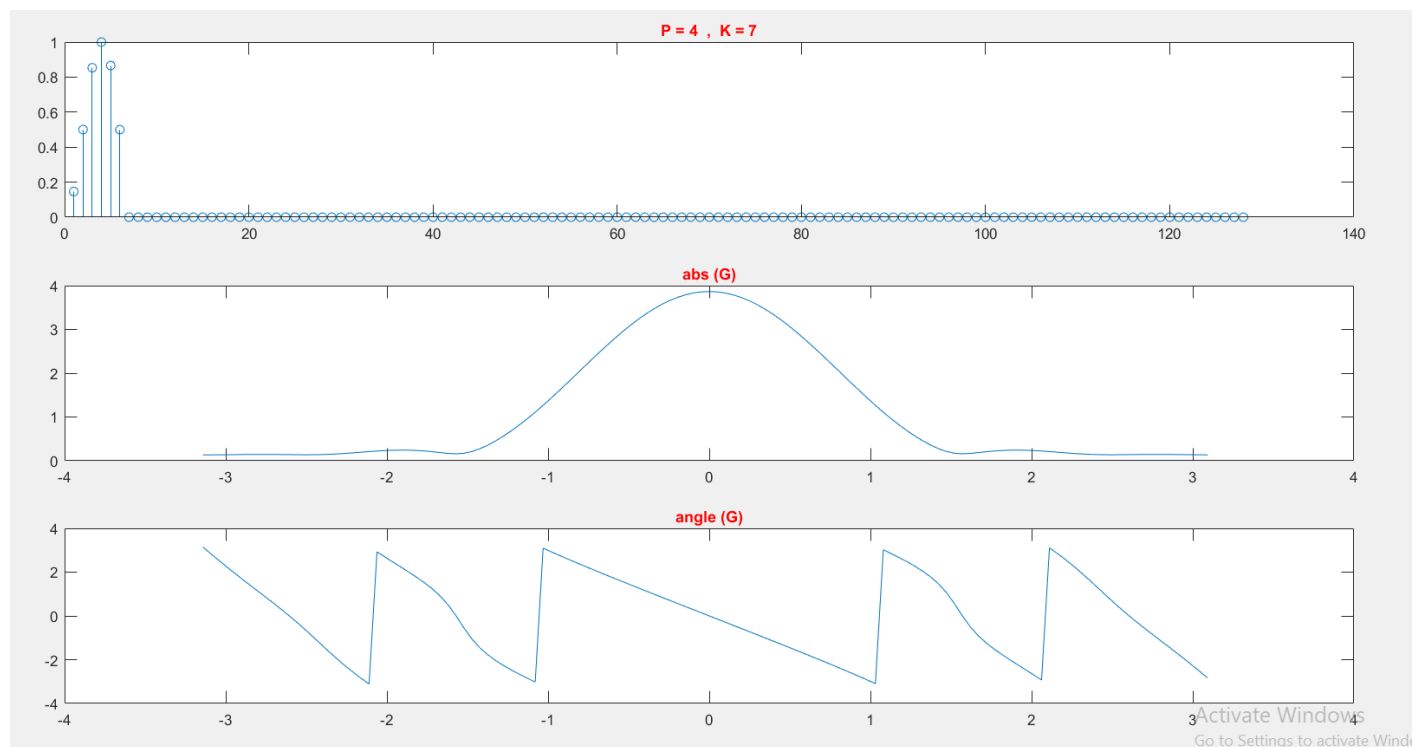
پیاده سازی Rosenberg Pulse در حوزه زمان :



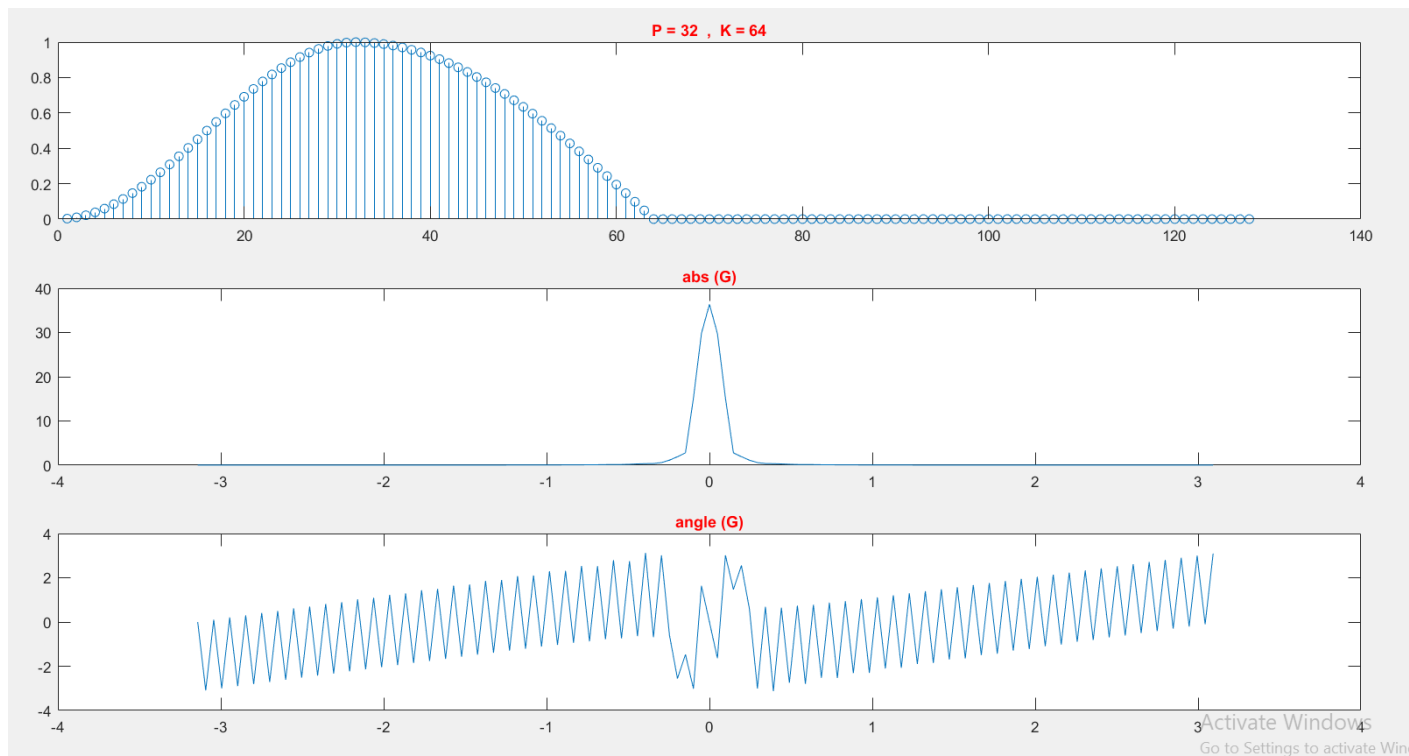
(d) پیاده سازی Rosenberg Pulse در حوزه فرکانس :



$$K = 2, P = 2$$



$$K = 7, P = 4$$



$$K = 64 \text{ و } P = 32$$

مقایسه ویژگی های زمانی و فرکانسی دو فیلتر بالا :

glottal shaping filter فیلتر IIR است در حالی که Rosenberg filter فیلتر FIR می باشد.

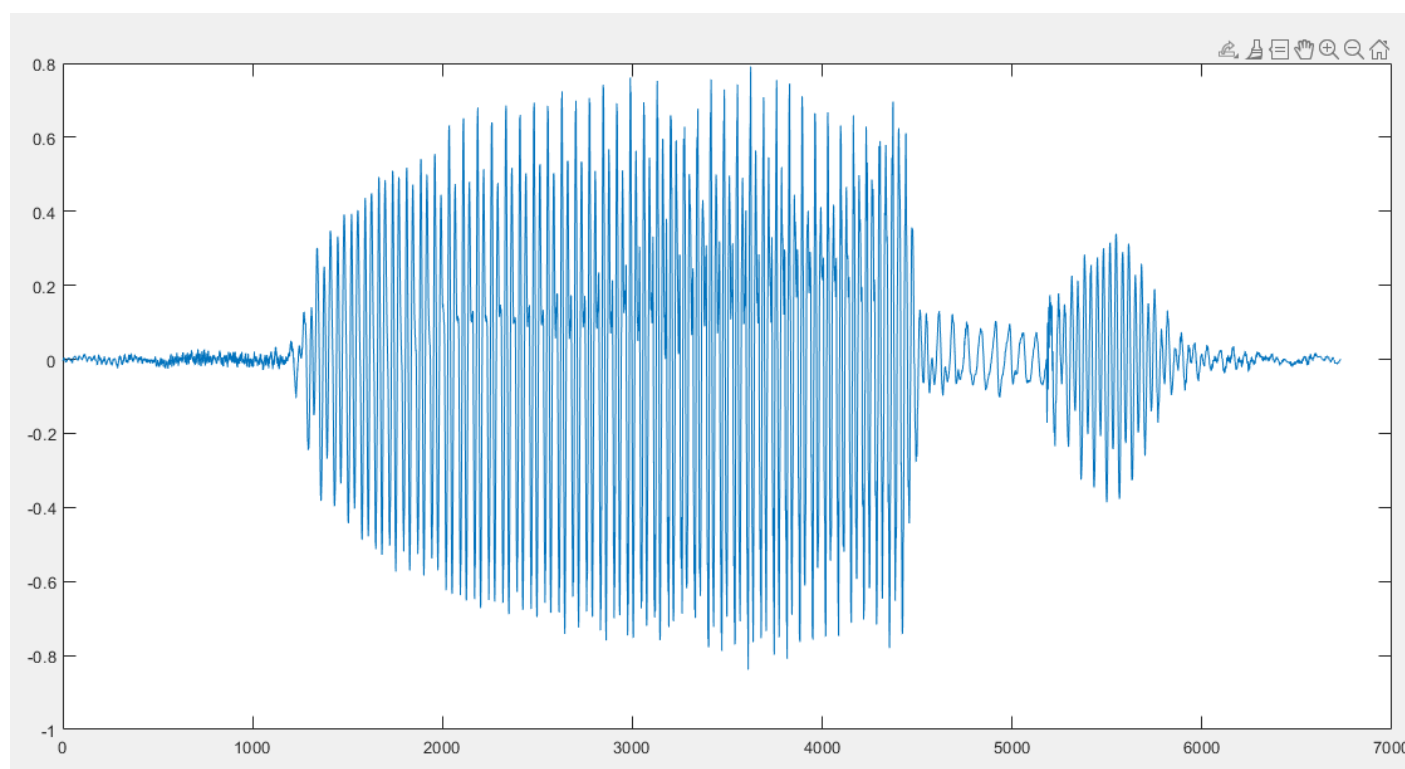
به کمک دو پارامتر α و β می توان glottal shaping filter طراحی کرد که در آنها با افزایش n مقدار سیگنال به سمت صفر میل می کند اما برای آن که در لحظه $n = 64$ خروجی صفر گردد نیاز است با توجه به فرض سوال مقدار α نزدیک 1 انتخاب گردد و به تبع آن مقدار β کوچکتر از 0.3 باشد در حالی که برای Rosenberg filter کافی است $K \leq 64$ باشد. می دانیم هر چه سیگنال در حوزه زمان زودتر میرا شود؛ در حوزه فرکانس آرام تر به سمت صفر میل می کند و از آنجا که ما به دنبال آن هستیم که فیلتر حنجره طول کمی داشته باشد تا بتوانیم فرضیه ای مطرح شده در درس برای صرف نظر از این فیلتر را برآورده سازد.

سوال 5)

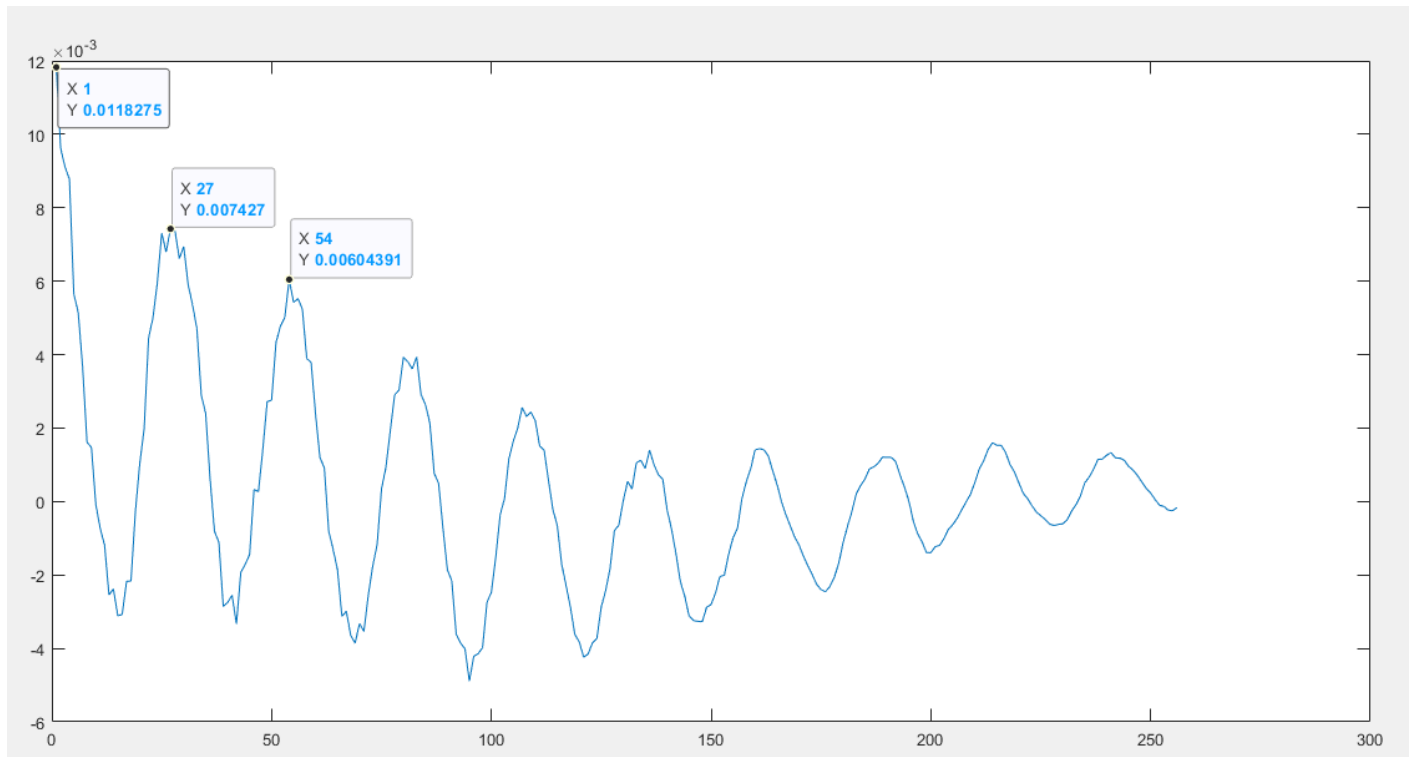
قسمت اول و دوم : محاسبه short – time autocorrelation و DTF

سیگنال صوت “Y0500.wav”

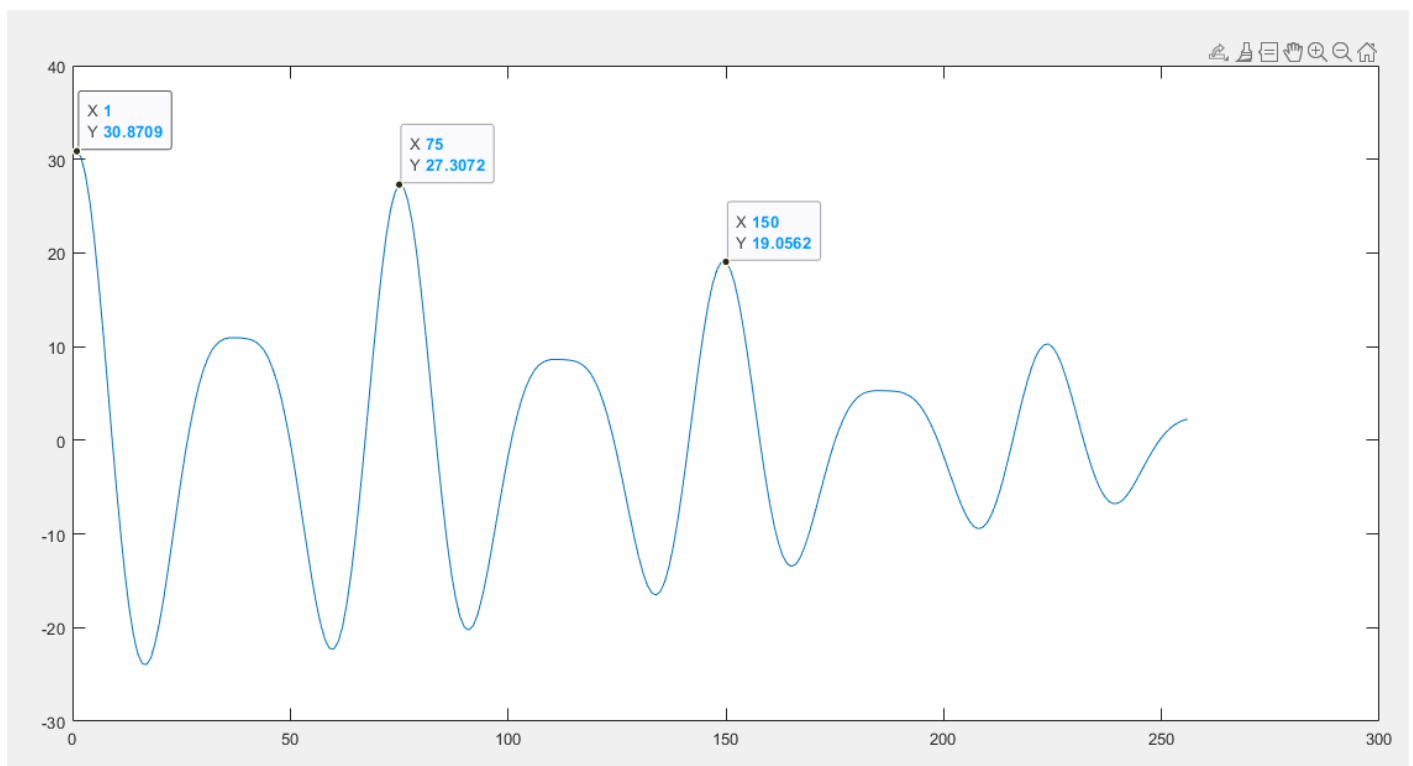
کل سیگنال مورد نظر در حوزه ی زمان



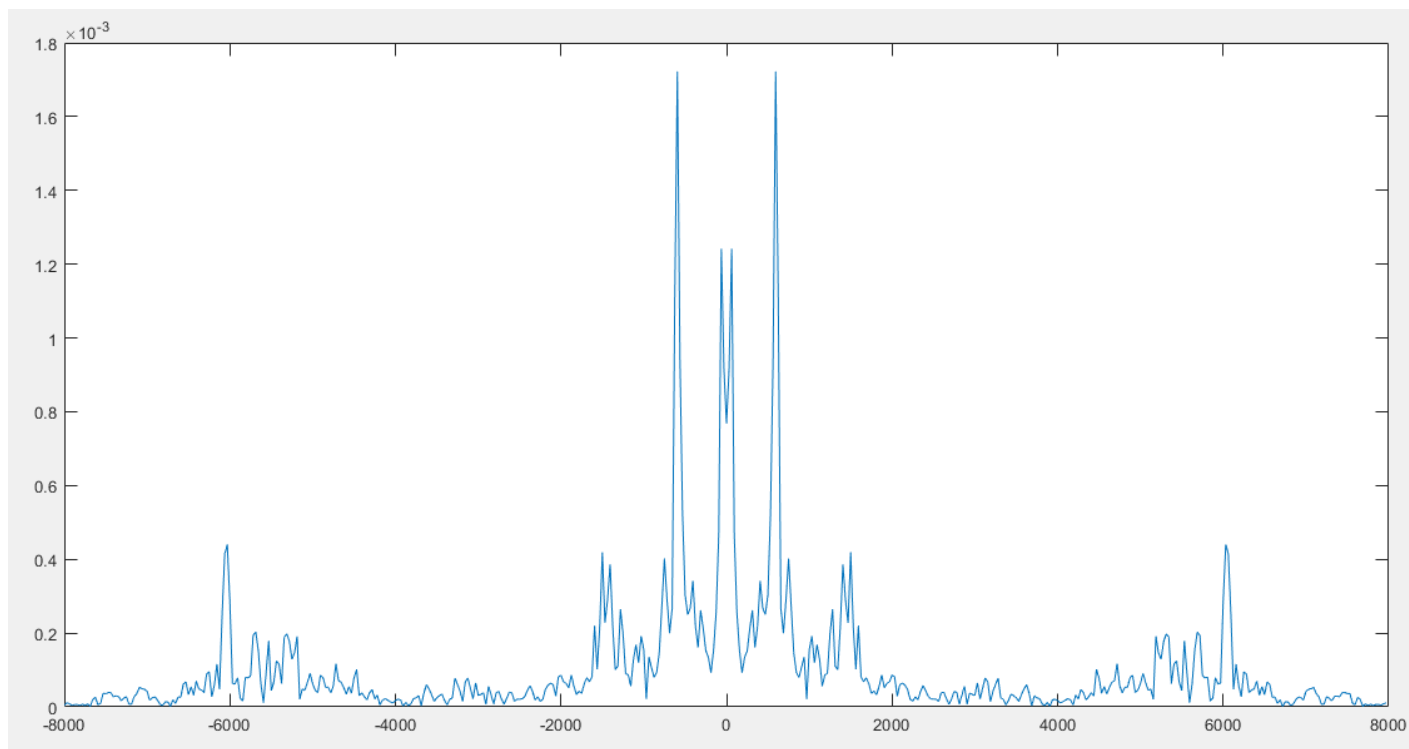
خود همبستگی سیگنال در پنجره ی اول



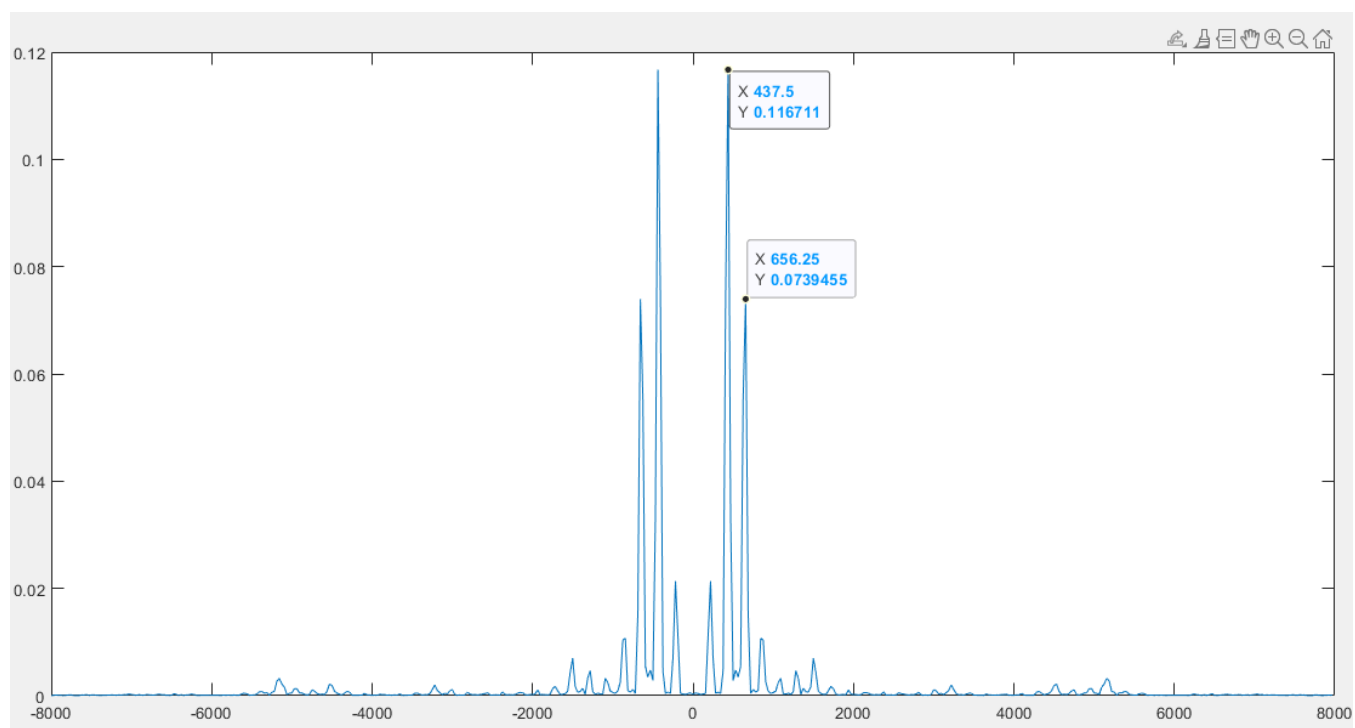
خود همبستگی سیگنال در پنجره پنجم



تبدیل فوریه 512 نقطه ای سیگنال موردنظر در پنجره اول



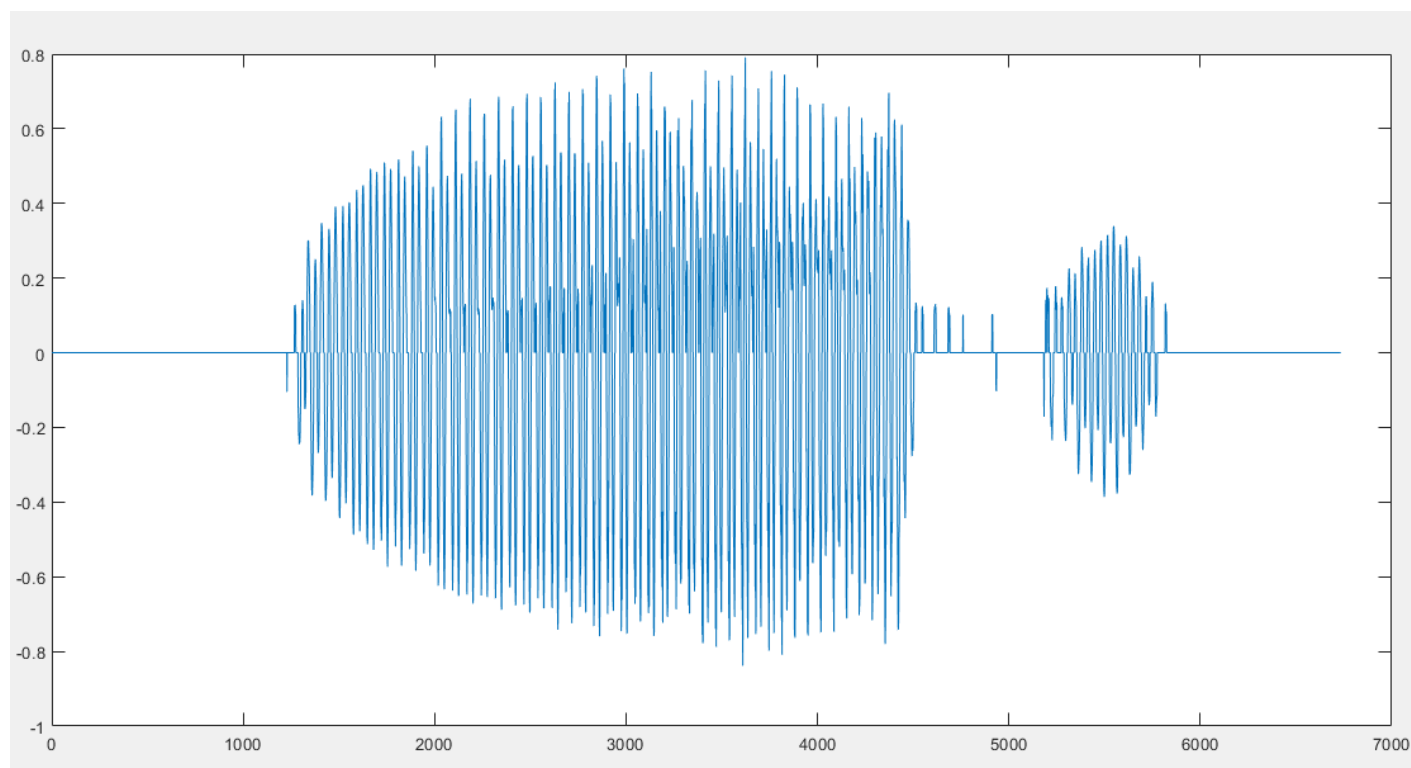
تبدیل فوریه 512 نقطه ای سیگنال موردنظر در پنجره پنجم



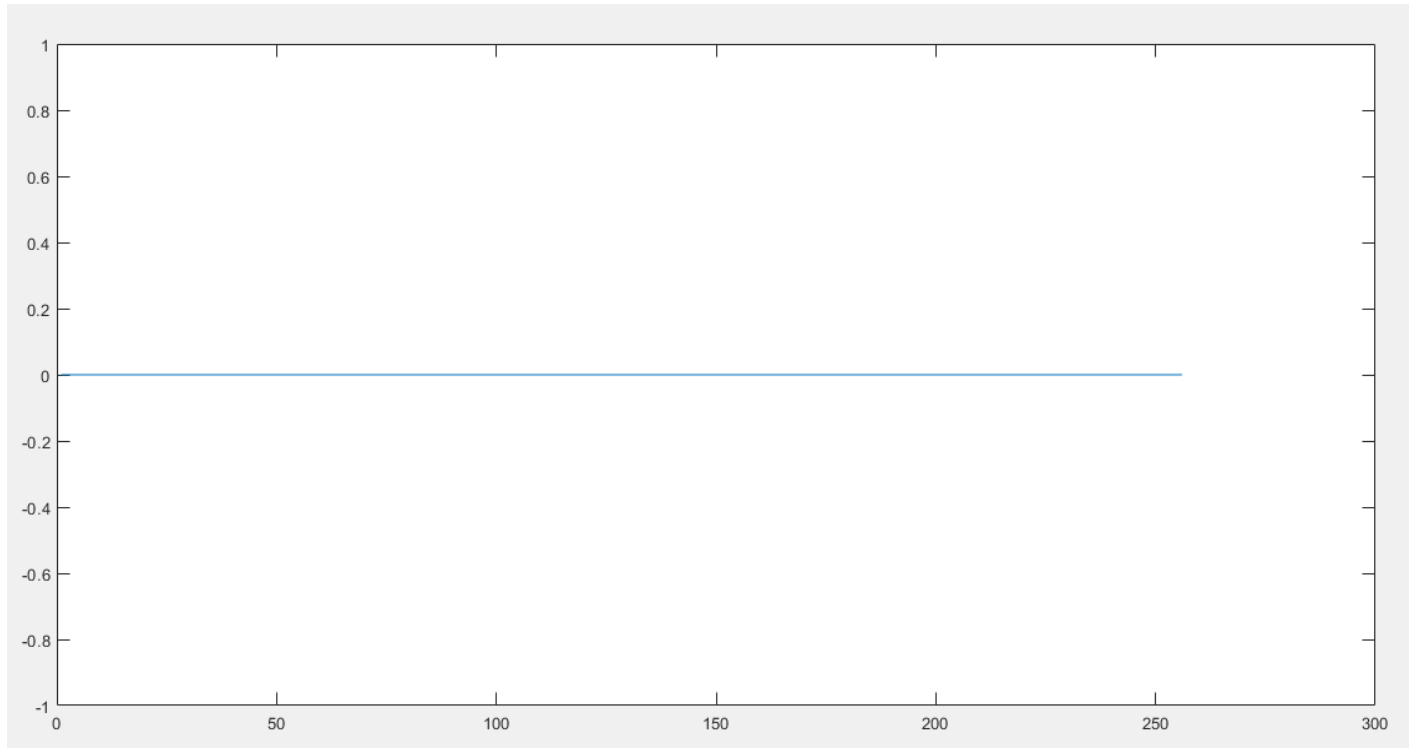
در شکل دیدیم که وقتی ما صدای واکه داریم ؛ تابع خود همبستگی ما انرژی بساری دارد و فاصله ی مبدا تا ماکسیمم مطلق pitch ما رو تعیین می کند.

قسمت سوم : وقتی که تابع غیر خطی cutter اضافه می کنیم

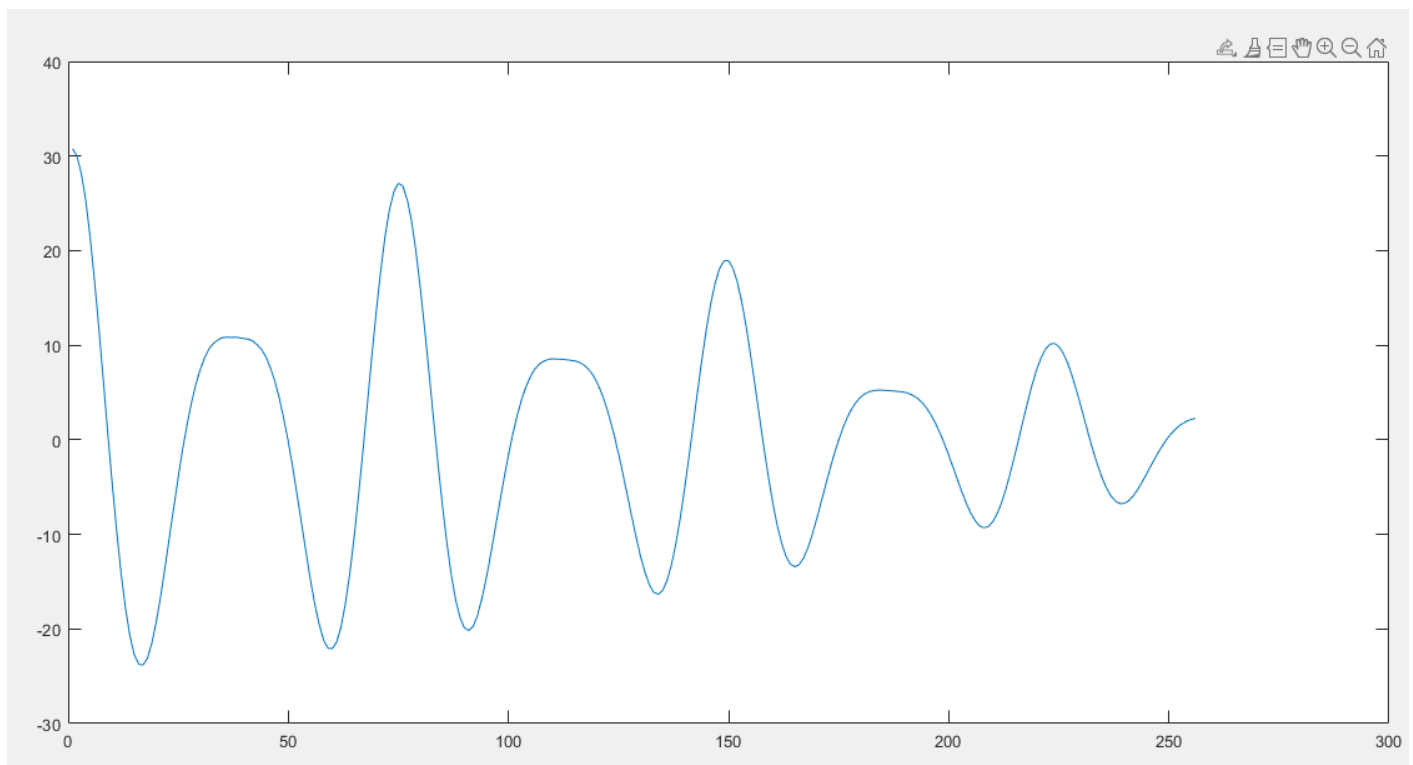
کل سیگنال مورد نظر در حوزه ی زمان



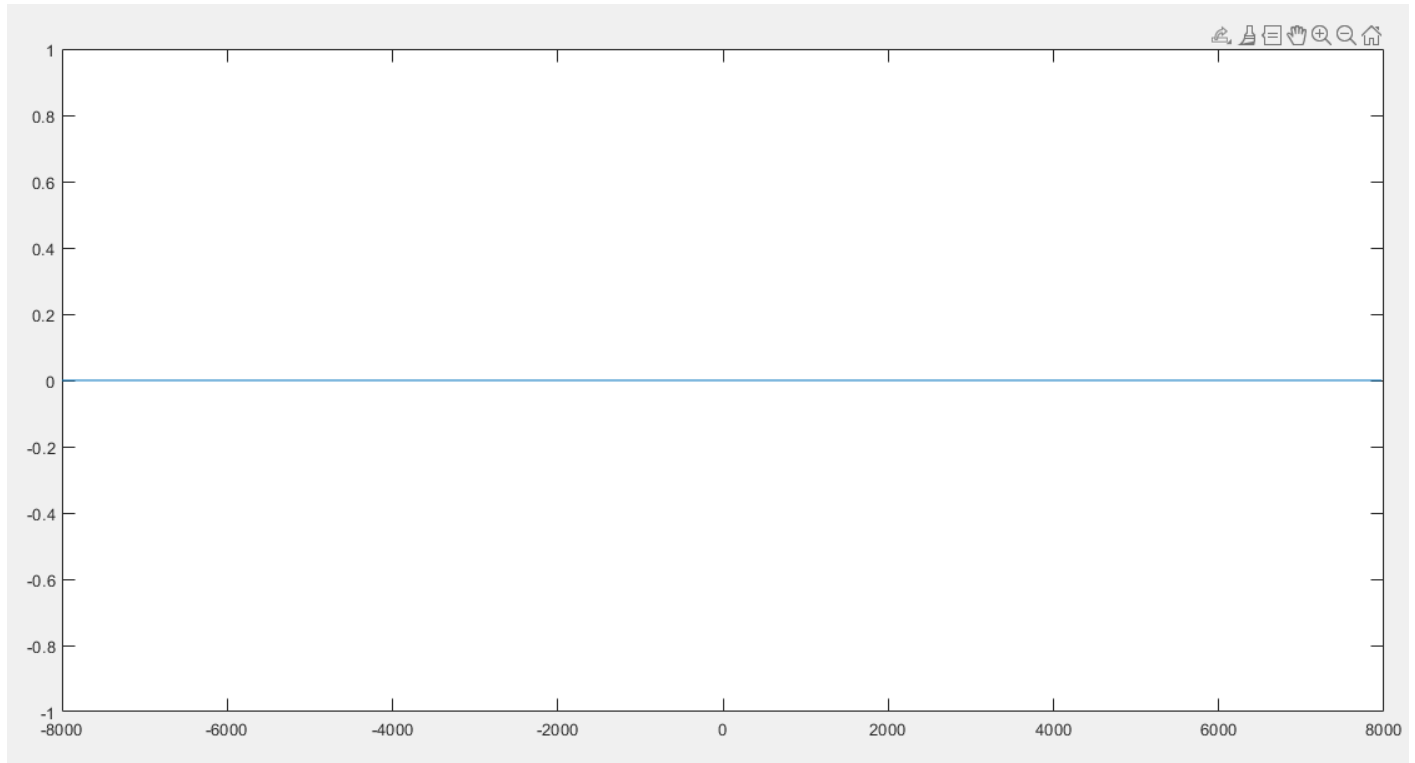
خودهمبستگی سیگنال در پنجره ی اول



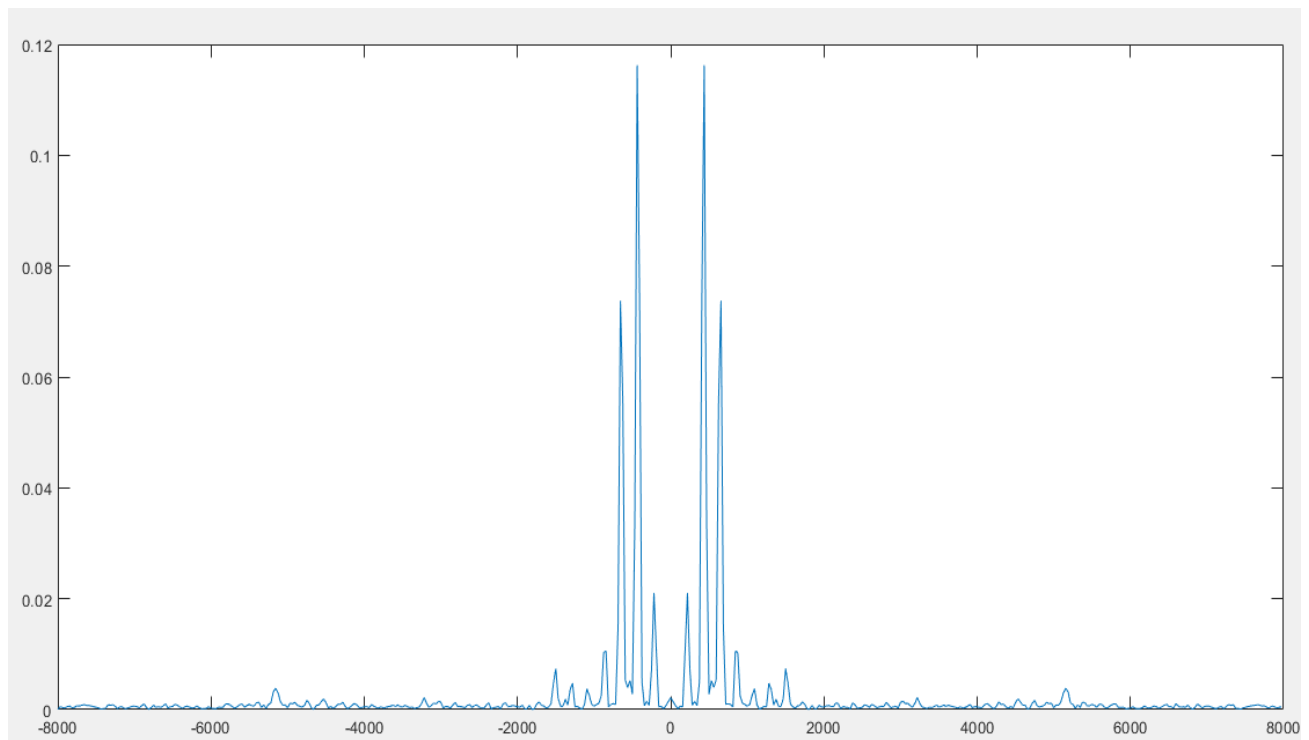
خودهمبستگی سیگنال در پنجره پنجم



تبدیل فوریه 512 نقطه ای سیگنال موردنظر در پنجره اول



تبدیل فوریه 512 نقطه ای سیگنال موردنظر در پنجره پنجم

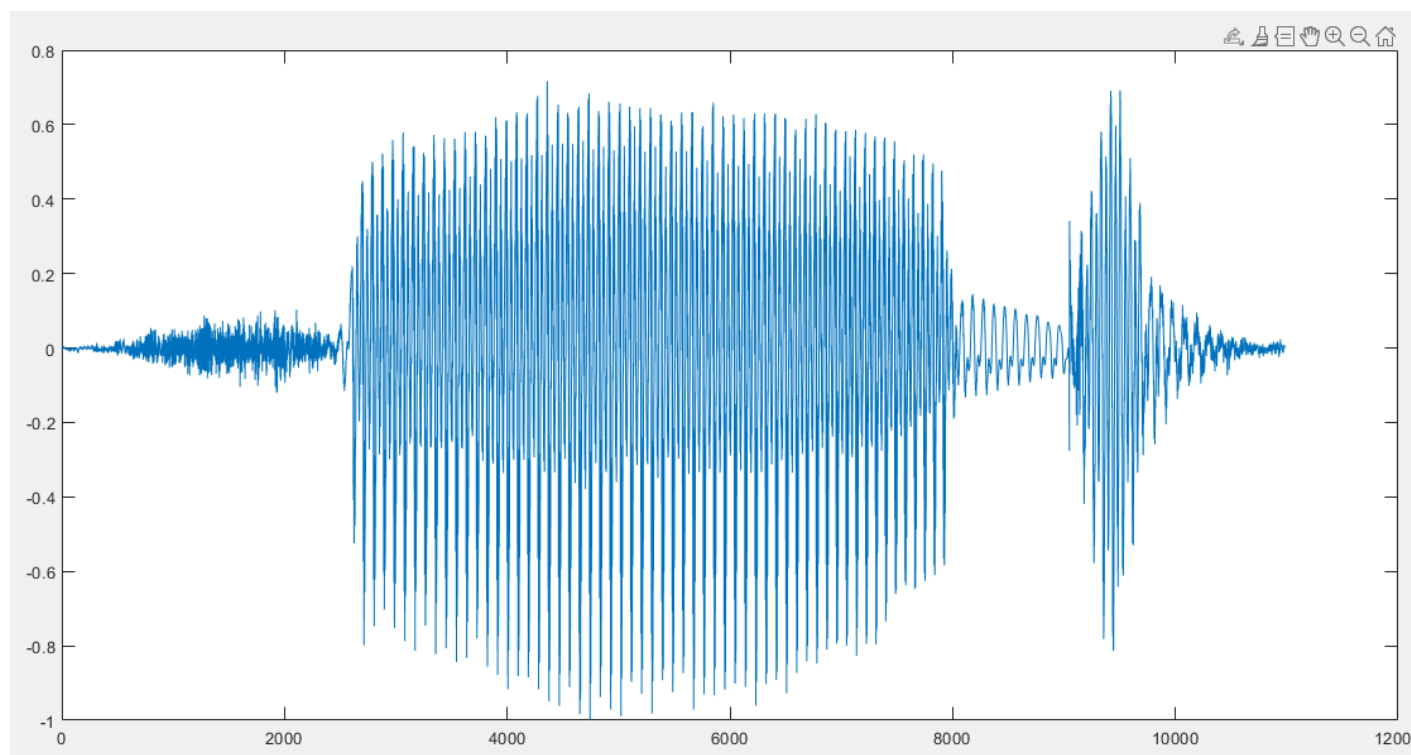


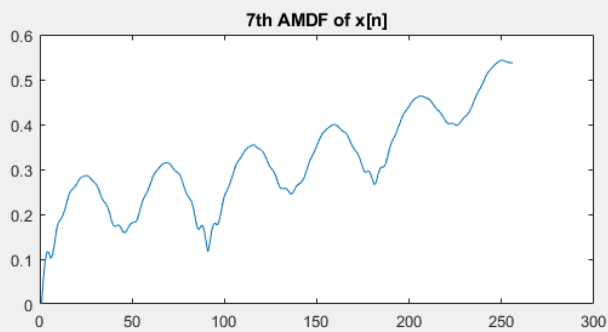
تابع غیر خطی cutter تابع غیر خطی خوبی است که البته می توان ان را با تقریب خوبی همانند خطی دانست که سبب گردیده صدای هایی که انرژی فوق العاده کمی نسبت به واکه ها دارند حذف گردند .

سوال 6)

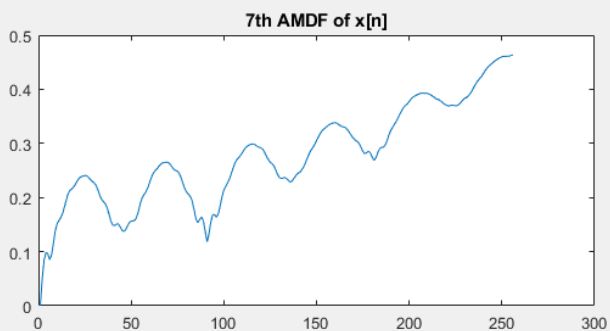
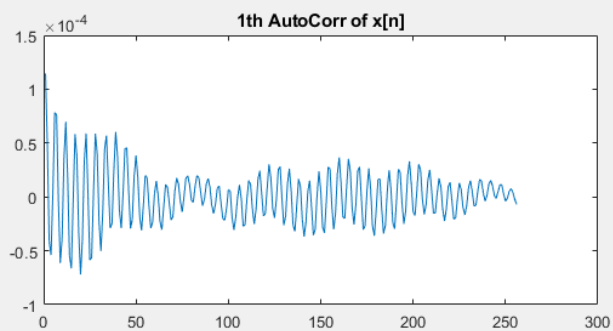
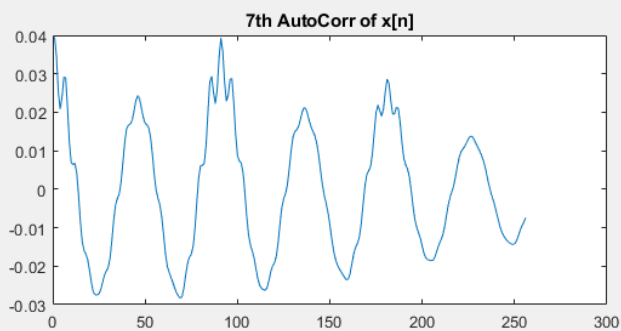
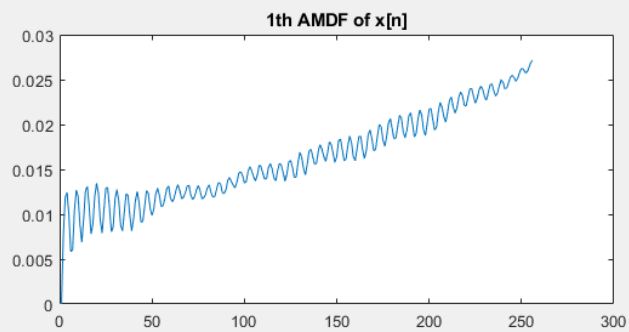
فایل صوتی مورد نظر "m01iy.wav" می باشد

سیگنال مورد نظر در حوزه زمان

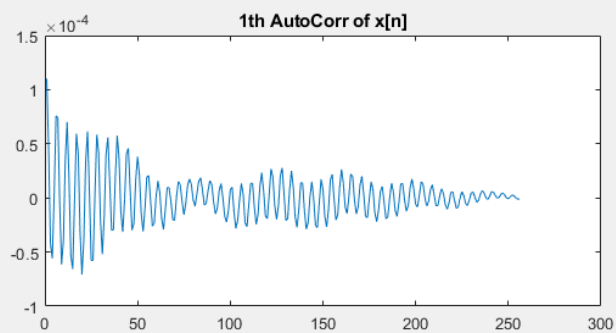
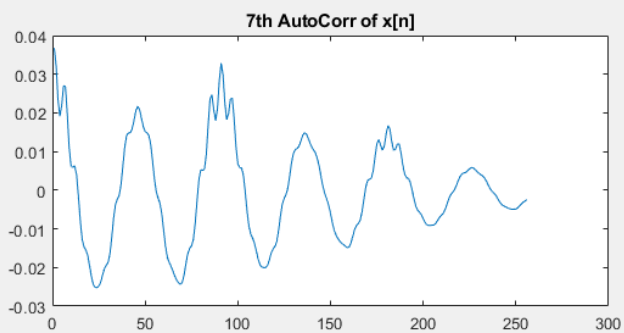
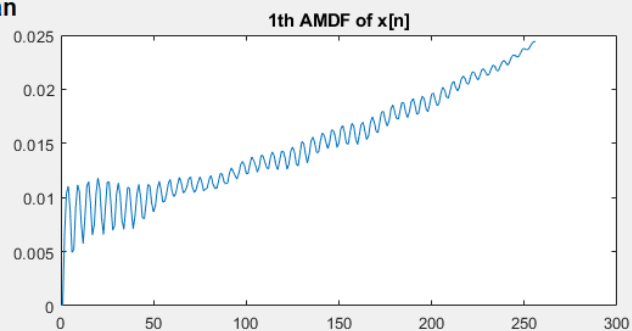




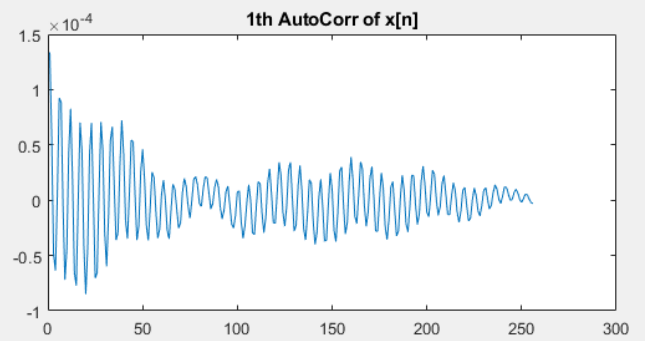
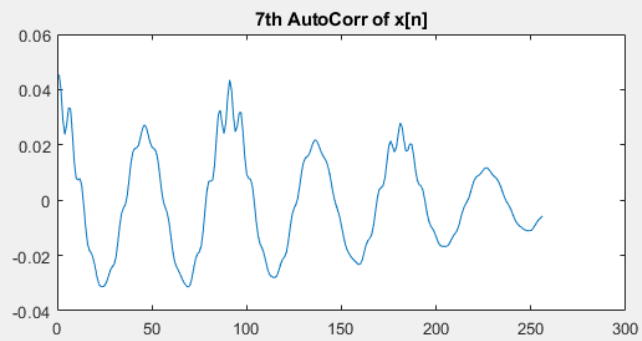
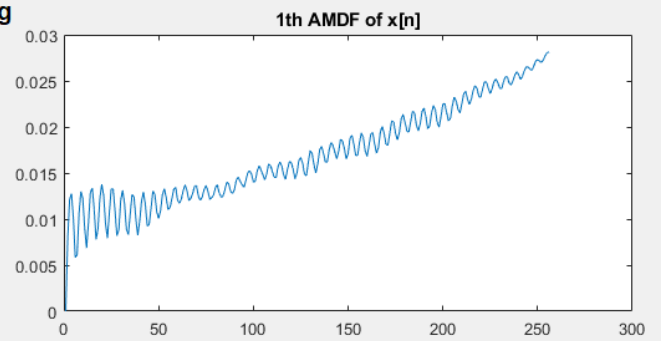
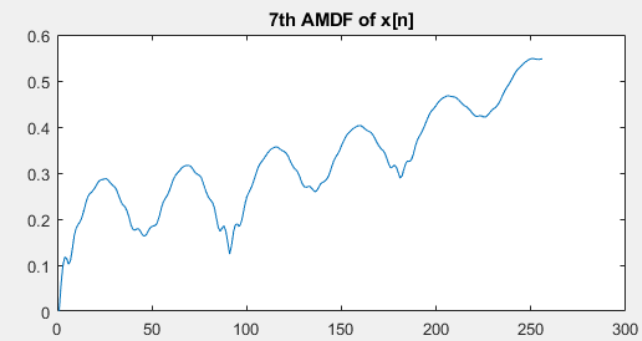
Bartlett



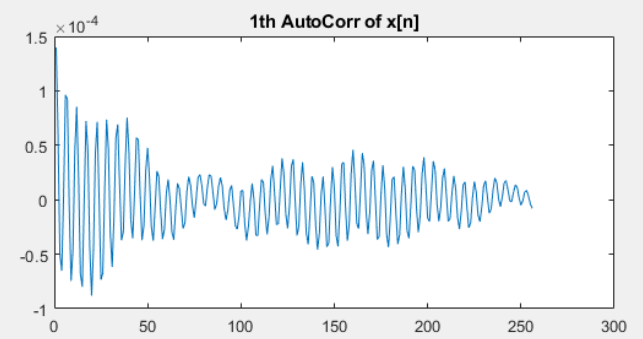
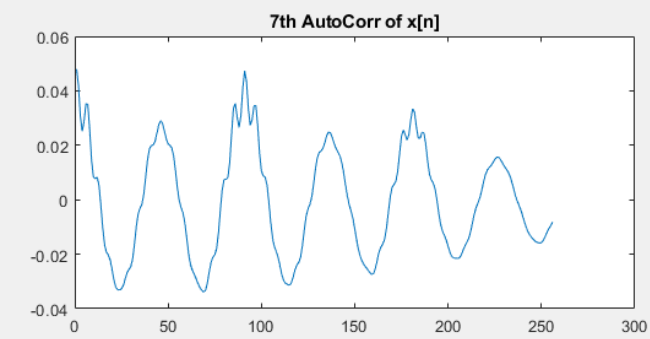
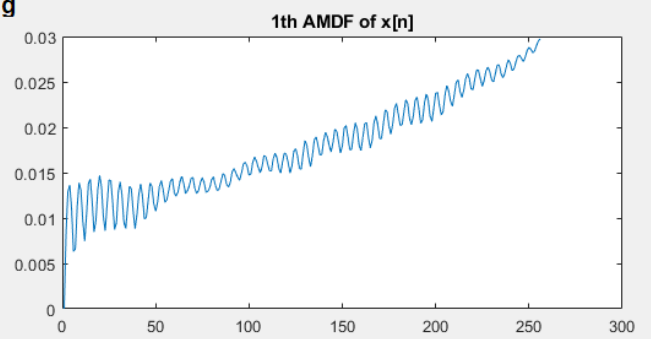
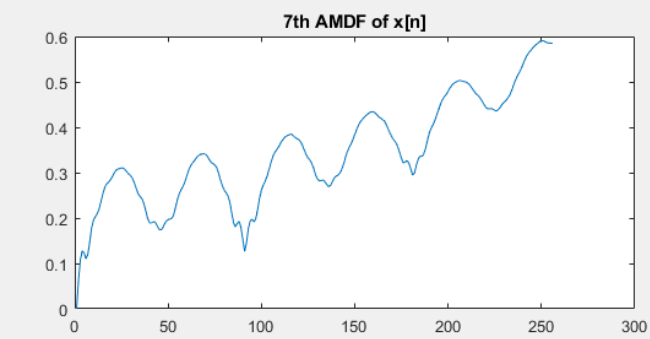
Blackman

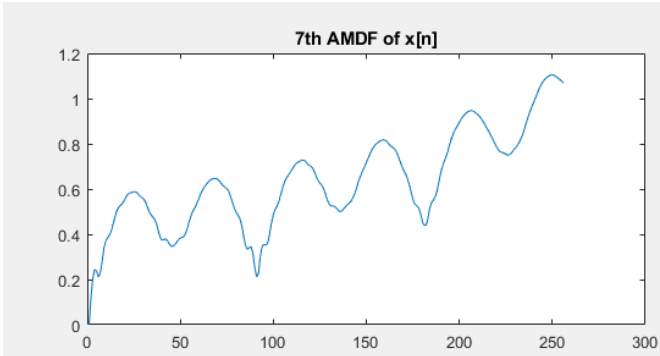


Hanning

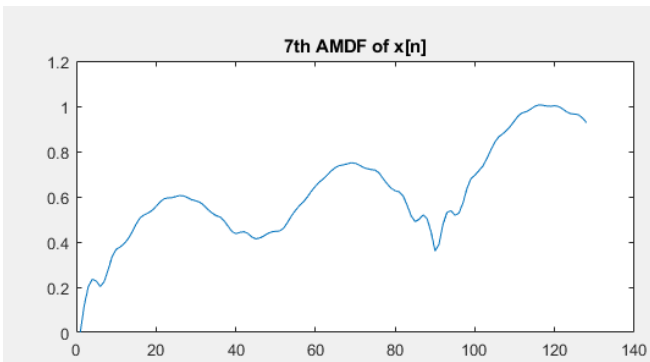
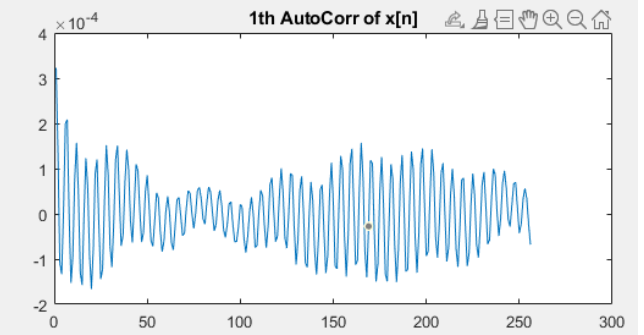
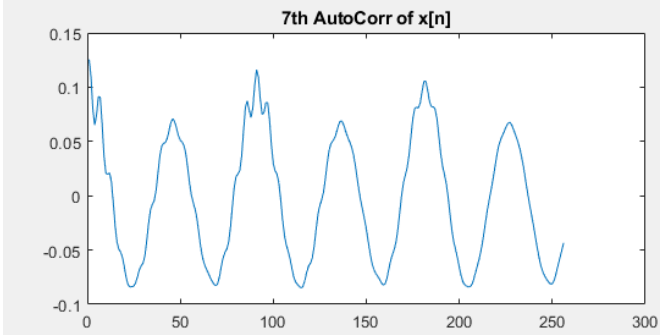
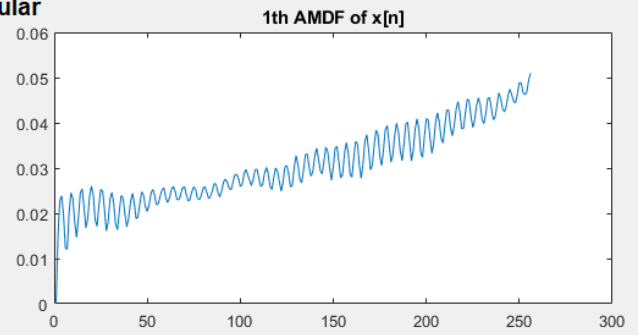


Hamming

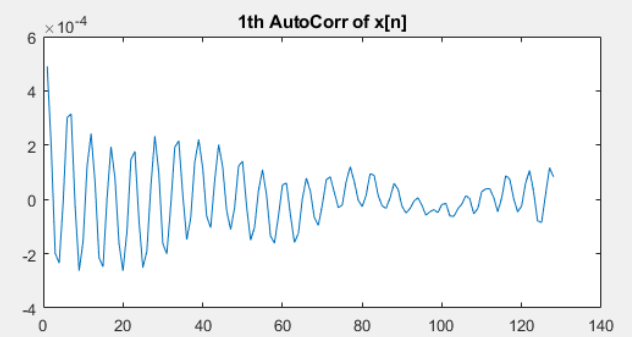
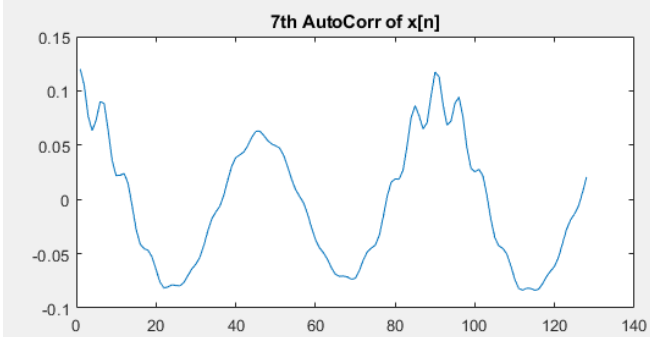
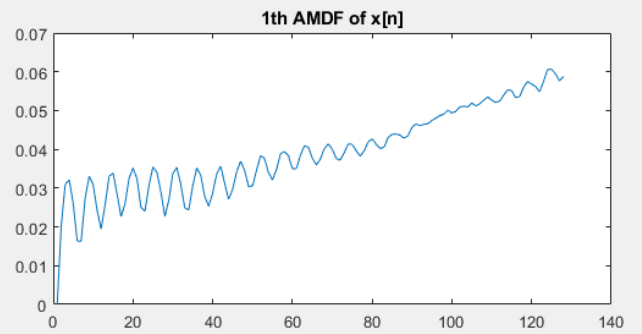


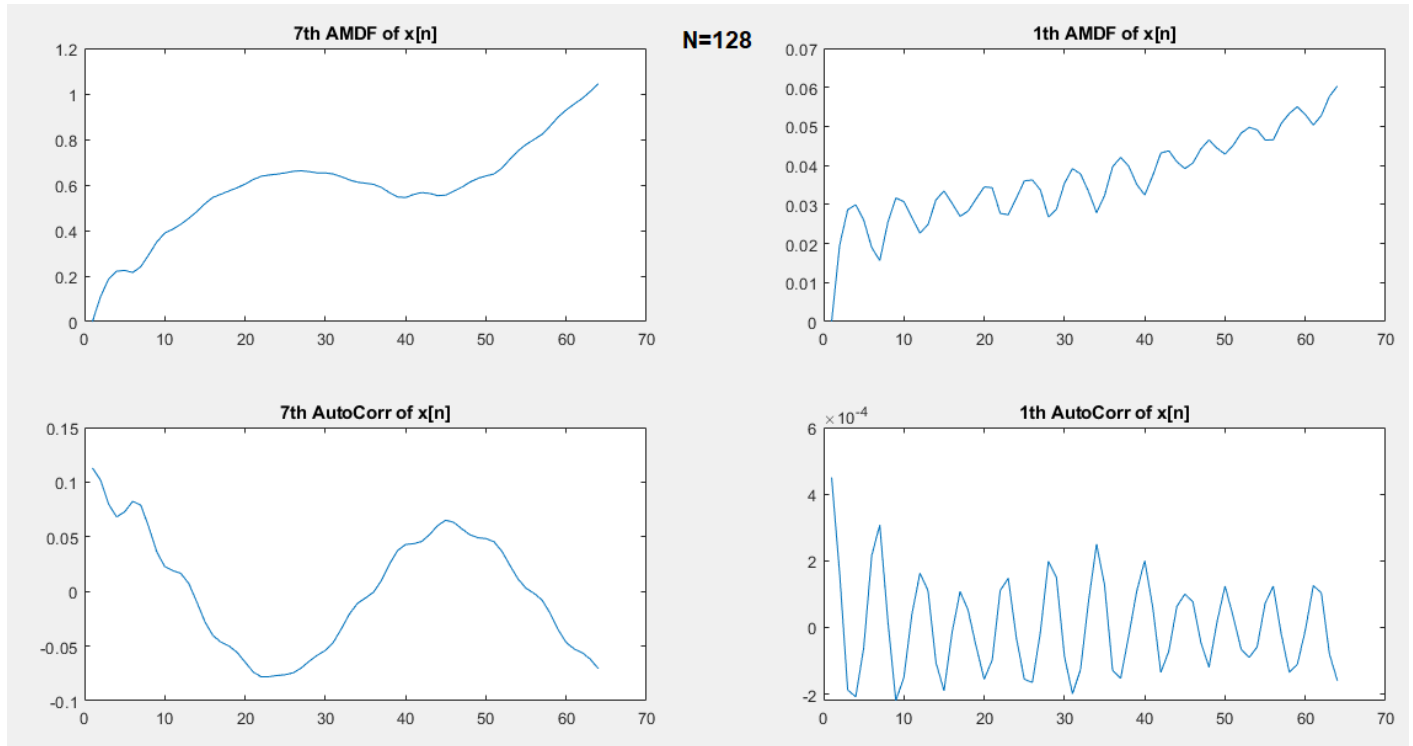


Rectangular



N=256





فرمول های پنجره زمانی

Rectangular

$$w[n] = \begin{cases} 1, & 0 \leq n \leq M, \\ 0, & \text{otherwise} \end{cases} \quad (7.47a)$$

Bartlett (triangular)

$$w[n] = \begin{cases} 2n/M, & 0 \leq n \leq M/2, \\ 2 - 2n/M, & M/2 < n \leq M, \\ 0, & \text{otherwise} \end{cases} \quad (7.47b)$$

Hanning

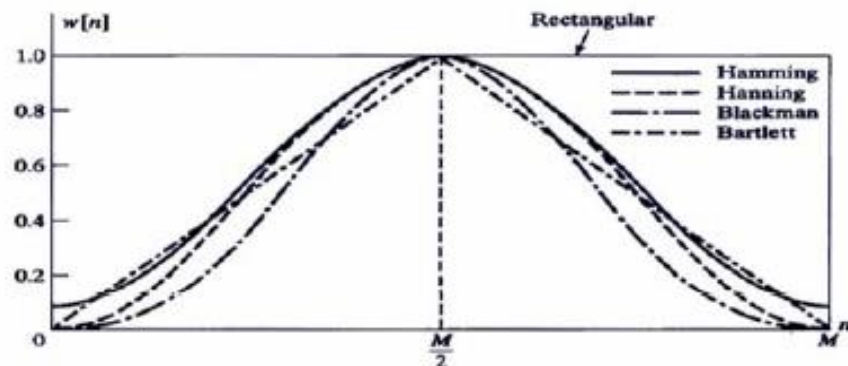
$$w[n] = \begin{cases} 0.5 - 0.5 \cos(2\pi n/M), & 0 \leq n \leq M, \\ 0, & \text{otherwise} \end{cases} \quad (7.47c)$$

Hamming

$$w[n] = \begin{cases} 0.54 - 0.46 \cos(2\pi n/M), & 0 \leq n \leq M, \\ 0, & \text{otherwise} \end{cases} \quad (7.47d)$$

Blackman

$$w[n] = \begin{cases} 0.42 - 0.5 \cos(2\pi n/M) + 0.08 \cos(4\pi n/M), & 0 \leq n \leq M, \\ 0, & \text{otherwise} \end{cases} \quad (7.47e)$$



در شکل های بالا از پنجره های هایی که در درس DSP تحت عناوین rectangular و hamming و hanning و blackman و bartlett استفاده گردیده است.

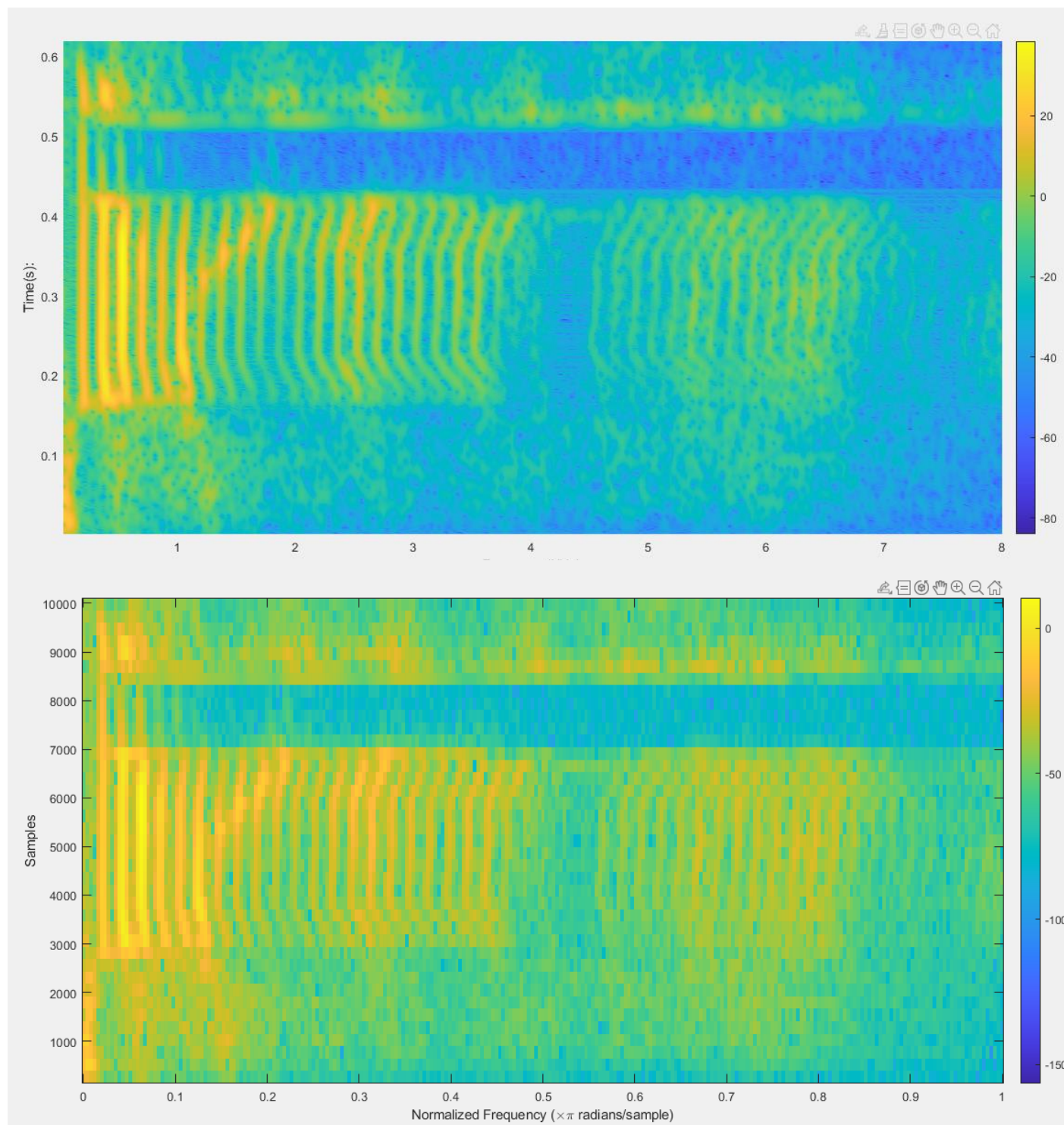
طول پنجره ها هم 512 و 256 و 128 است.

همان طور که شکل ها مشاهده می شود ؛ نوع پنجره اختلاف زیادی در تشخیص درست pitch ایجاد نمی کند. اما در مورد طول پنجره این اختلاف در تشخیص درست با طول 128 دیده شد. البته طبیعی و منطقی به نظر می رسد که هر چه طول پنجره کم گردد از دقت ایده آل فاصله می گیریم

سوال (7)

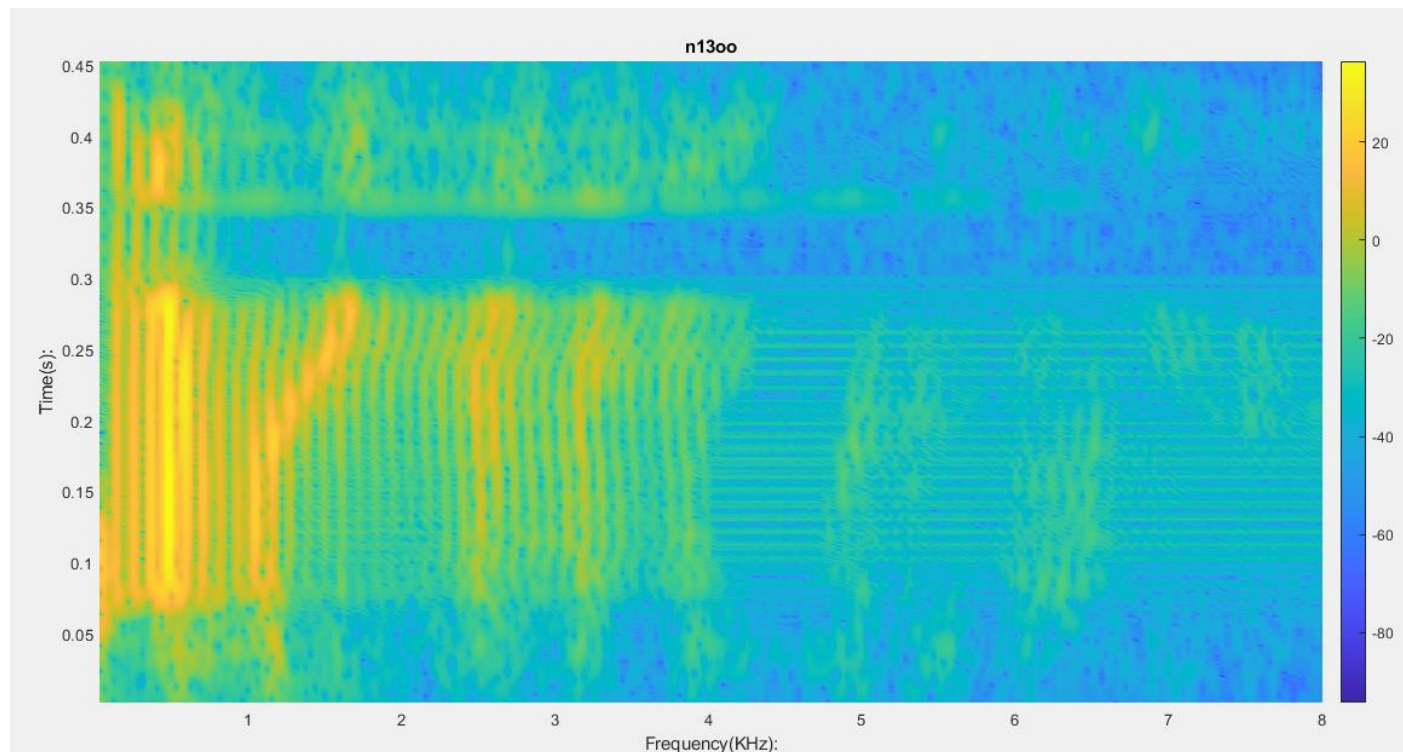
فایل صوتی مورد نظر “y05oo.wav” می باشد

قسمت اول

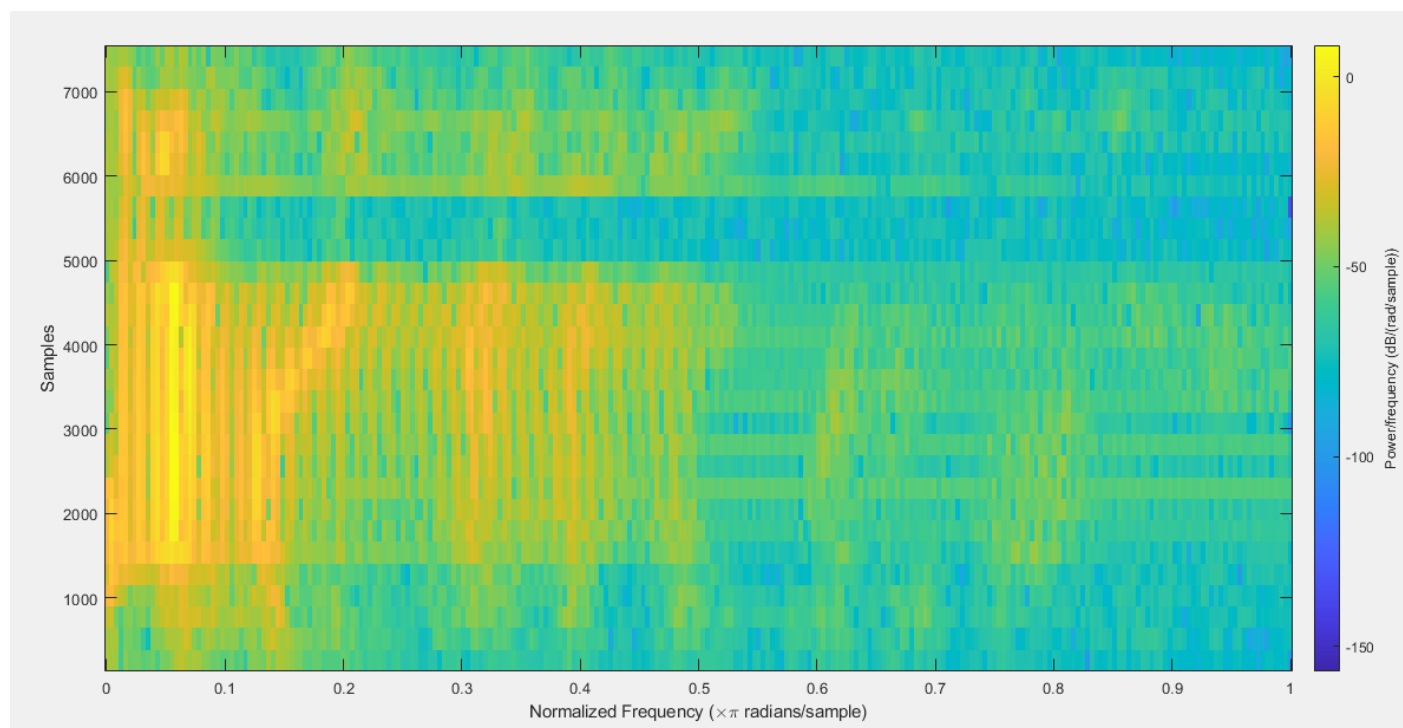


قسمت ب)

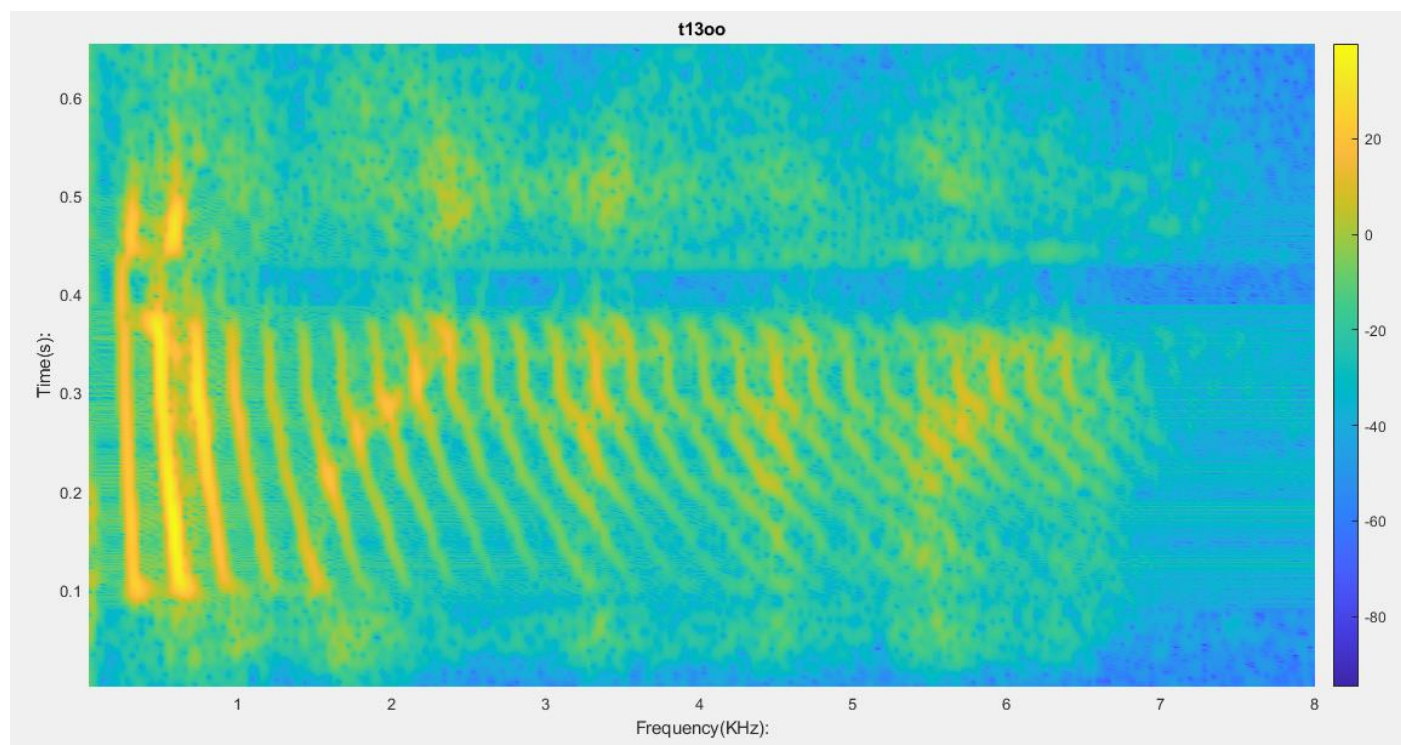
شکل رسم شده توسط خودمان :



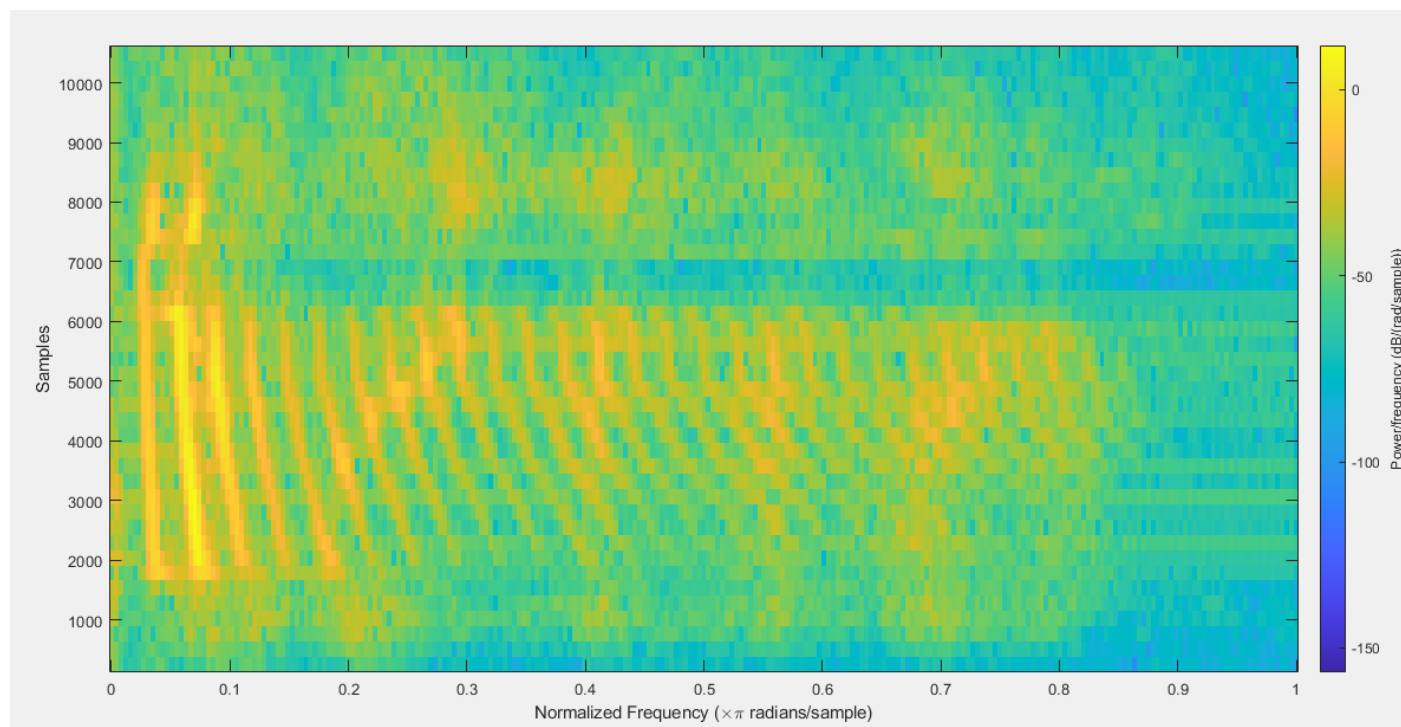
تابع Spectrogram متلب :



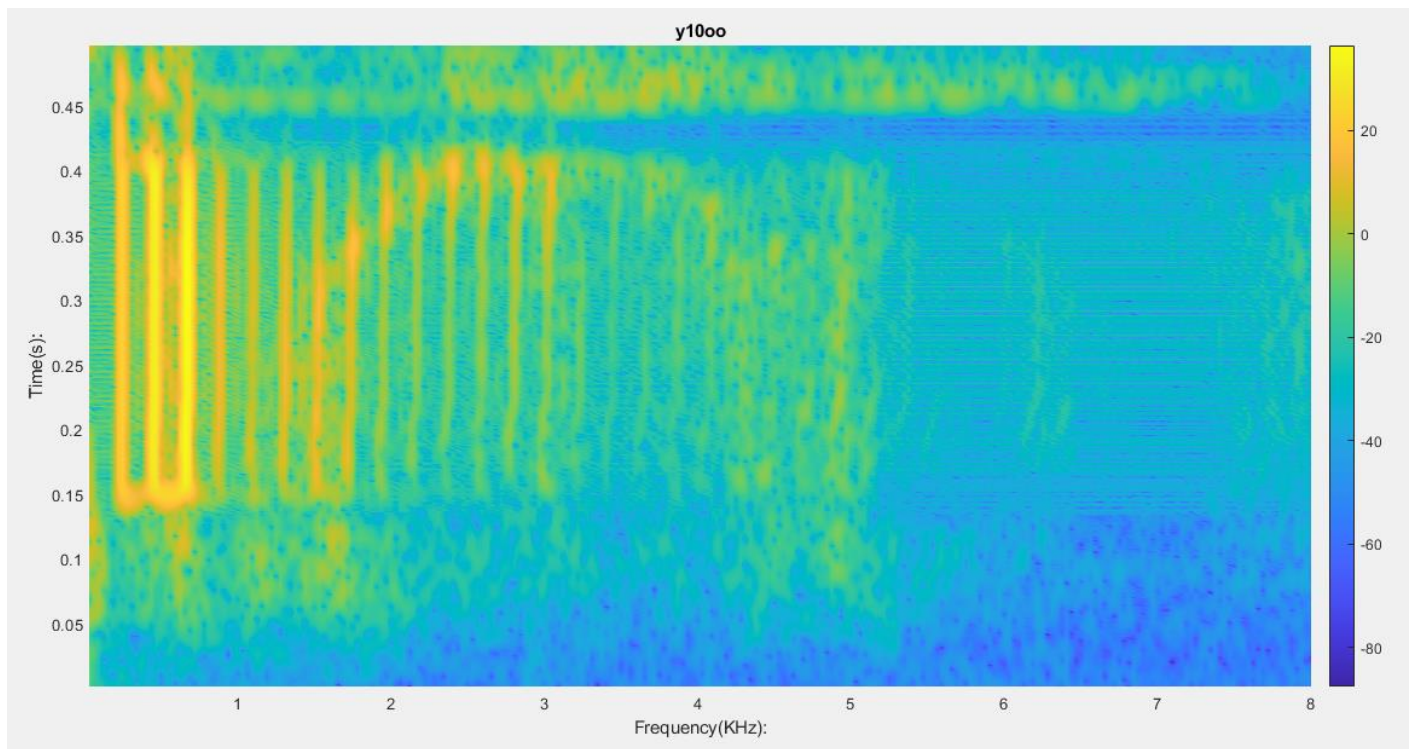
شکل رسم شده توسط خودمان :



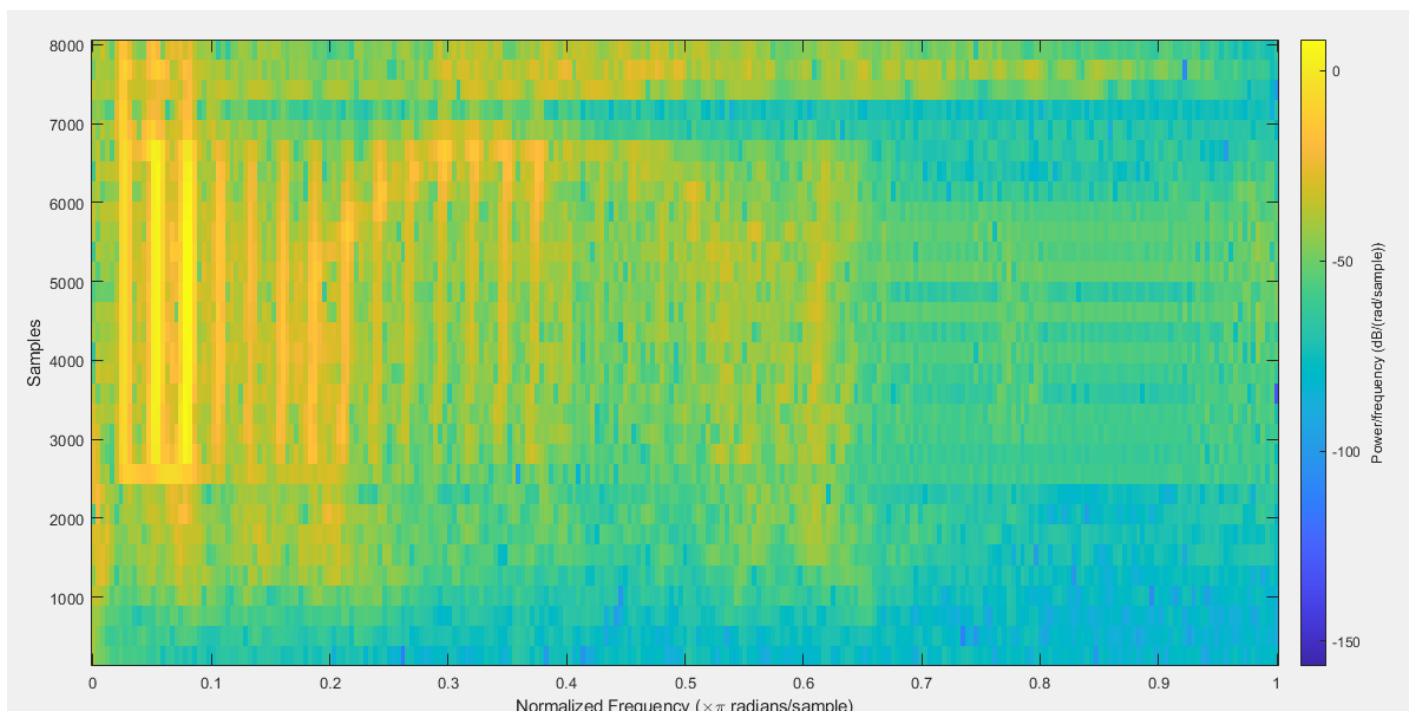
تابع Spectrogram متلب :



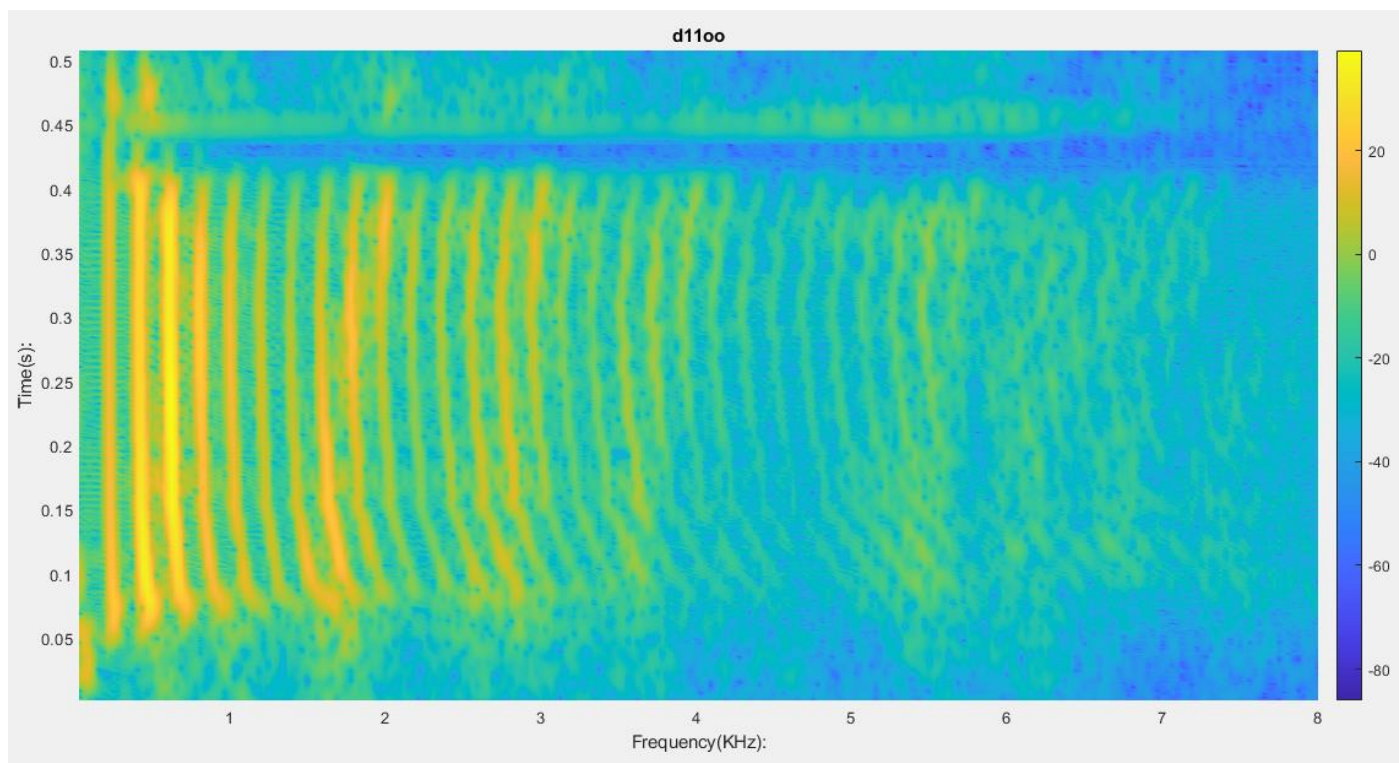
شکل رسم شده توسط خودمان :



تابع Spectrogram متلب :



شکل رسم شده توسط خودمان



تابع Spectrogram متلب

