

EE269

Signal Processing for Machine Learning

Cepstrum

Instructor : Mert Pilanci

Stanford University

October 11 2021

Linear systems and additive noise

- ▶ Linear systems, e.g., filters, can easily separate **additive noise** from useful information when we know the frequency range of the noise and information

$$y[n] = x[n] + w[n]$$

- ▶ In vector notation

$$Hy = Hx + Hw$$

Multiplicative or convolutive noise

- ▶ This is harder if the signal and noise are **convoluted**, e.g., in speech processing

$$y[n] = x[n] * w[n]$$

- ▶ $w[n]$ is the flowing air (noise source)
- ▶ $h[n]$ is the vocal tract (filter)

We can develop an operator that can separate convoluted components by **transforming convolution into addition**

Cepstrum

- ▶ Developed to separate convoluted signals

$$y[n] = x[n] * w[n]$$

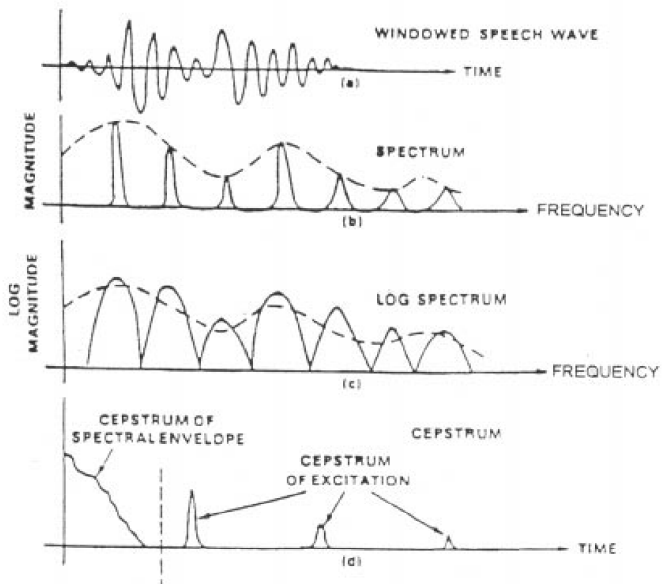
Discrete Fourier Domain:

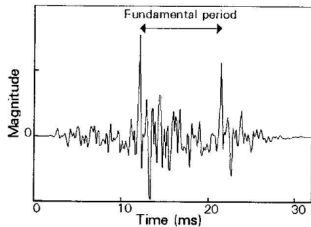
$$Y[k] = X[k]W[k]$$

- ▶ Take logarithms

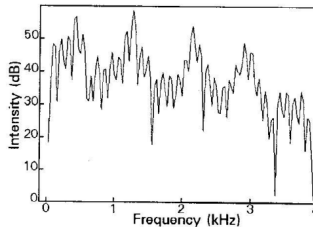
$$\log[Y[k]] = \log X[k] + \log W[k]$$

- ▶ we can apply a linear filter to $\log Y[k]$ to separate
 - ▶ equivalently we can take DFT of $\log Y[k]$ and process in frequency domain
- cepstrum is the DFT (or DCT) of the log spectrum

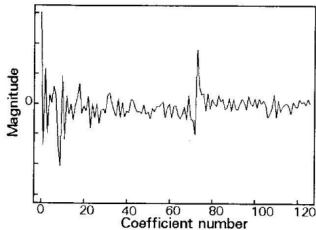




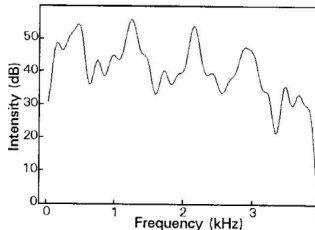
(a) Windowed speech waveform (32 ms at 8 kHz sampling rate).



(b) Log spectrum (from a Fourier transform).

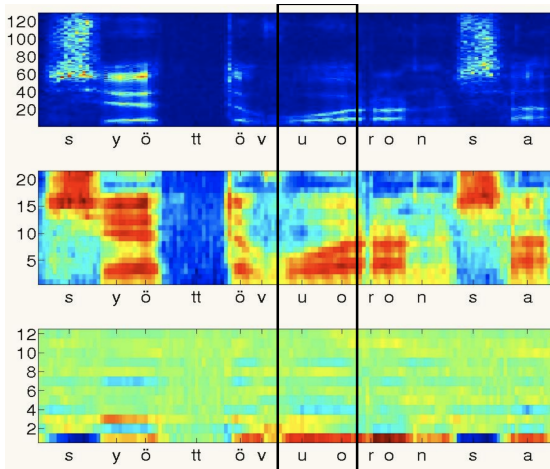


(c) Cepstrum computed from the log spectrum shown in (b).



(d) Log spectrum reconstructed from the first 40 cepstral coefficients in (c).

Figure 10.3 Analysing a section of speech waveform to obtain the cepstrum and then to reconstruct a cepstrally smoothed spectrum.



1. Frames:
short 10ms
windows
2. FFT:
power spectrum
spectrogram

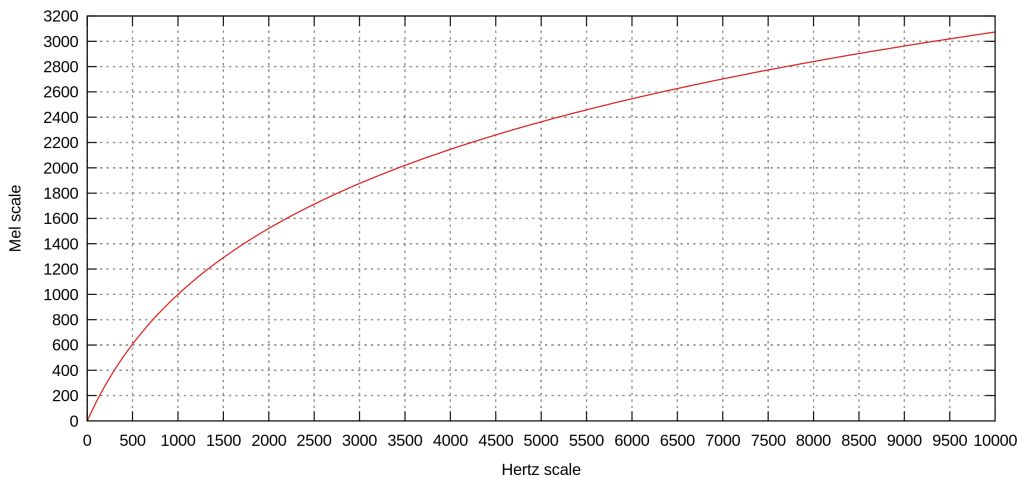
3. Filtering:
mel filter
motivated by
human ear
“essential data”

4. Features:
DCT transform
mel cepstrum
MFCC
-less features
-less correlation

Application: Mel-frequency spectrum

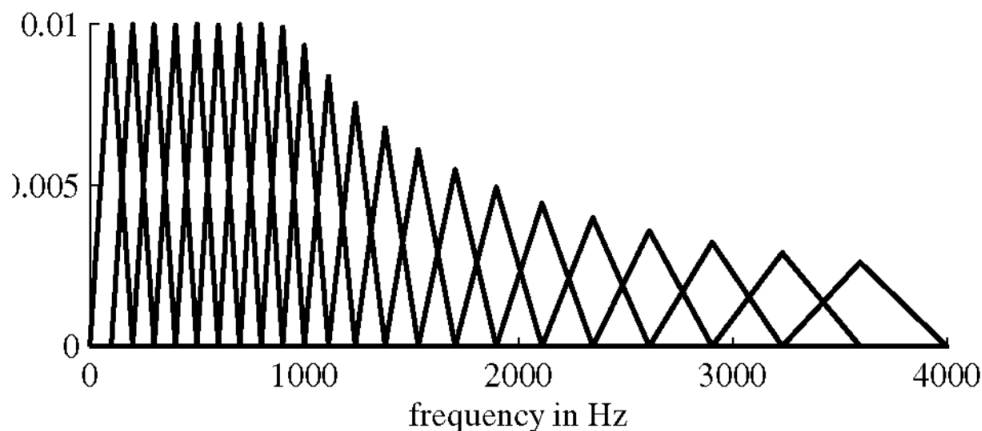
- ▶ perceptual scale of pitches
- ▶ 1 mels = 1000 Hz
- ▶ a formula to convert f hertz into m mels

$$m = 2595 \log_{10} \left(1 + \frac{f}{700} \right)$$



Application: Mel-frequency spectrum

- ▶ weighted DFT magnitude
- ▶ mel-frequency spectrum $MF[r]$ is defined as

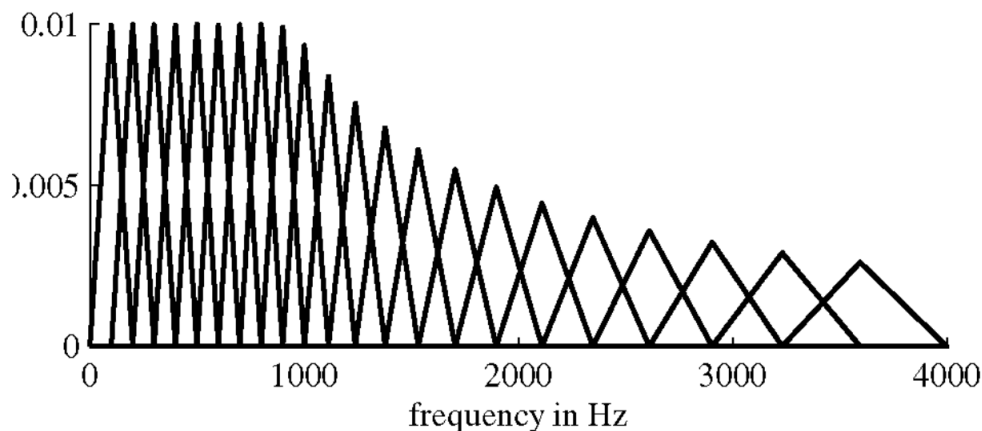


$$MF[r] = \sum_k |V_r[k]X[k]|^2$$

- ▶ $V_r[k]$ is the triangular weighting function for the r th filter.
- ▶ bandwidths are constant for center frequencies ≤ 1 kHz and then increase exponentially
- ▶ identical to convolutions with 22 filters

Application: Mel-frequency spectrum

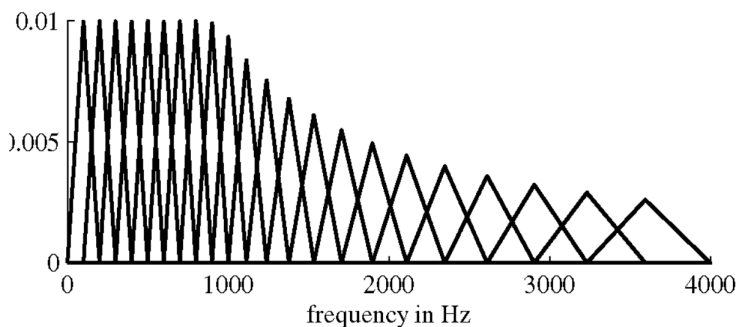
- ▶ weighted DFT magnitude
- ▶ mel-frequency spectrum $MF[r]$ is defined as



$$MF[r] = \sum_k |V_r[k]X[k]|^2$$

- ▶ $V_r[k]$ is the triangular weighting function for the r th filter.
- ▶ bandwidths are constant for center frequencies ≤ 1 kHz and then increase exponentially
- ▶ identical to convolutions with 22 filters

Application: Mel-frequency spectrum



$$\text{MF}[r] = \sum_k |V_r[k]X[k]|^2$$

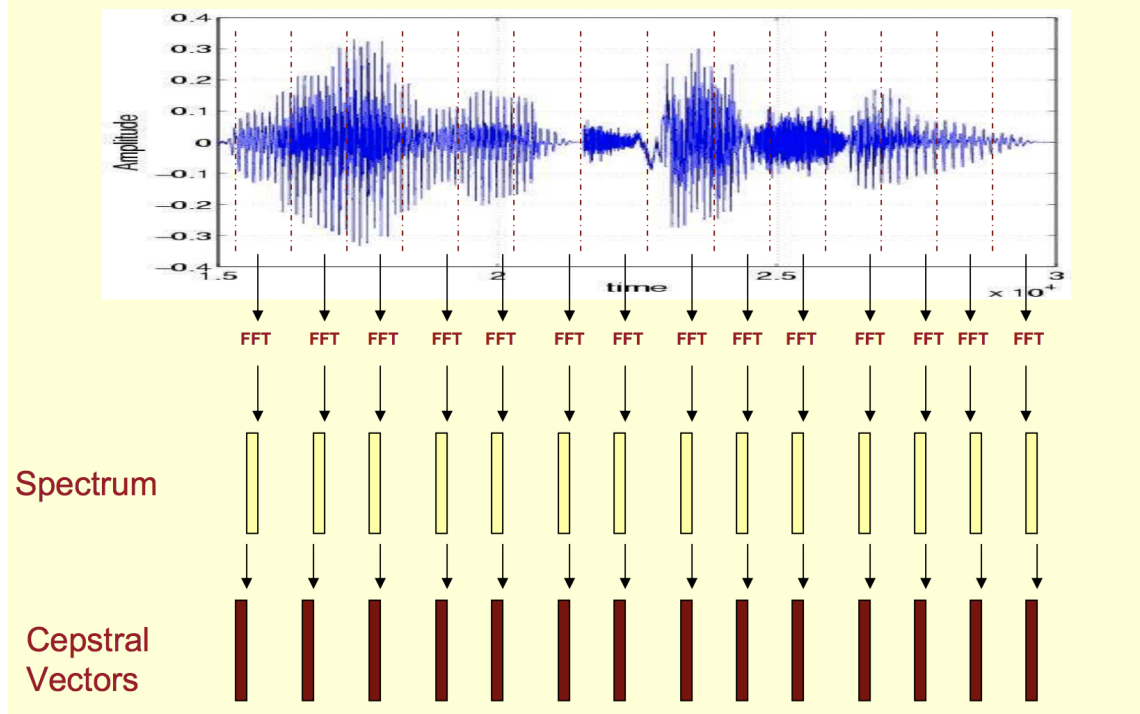
- Mel Frequency Cepstral Coefficient (MFCC)

$$\text{MFCC}[m] = \sum_{r=1}^R \log(\text{MF}[r]) \cos \left[\frac{2\pi}{R} \left(r + \frac{1}{2} \right) m \right] \quad (1)$$

- i.e., inner-product with cosines $\text{MFCC}[m] = \langle \log \text{MF}[r], c_m[r] \rangle$

Application: Speaker Identification

Speech signal represented as a sequence of CEPSTRAL vectors



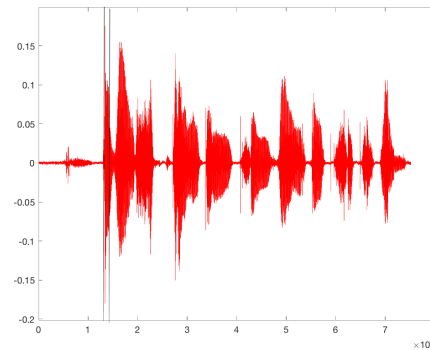
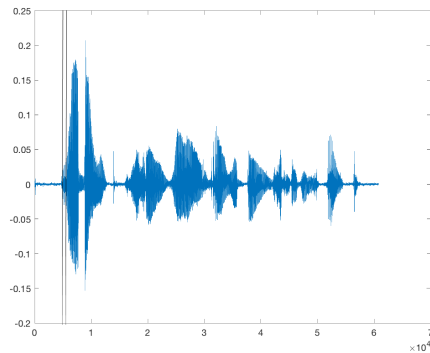
- train a k-Nearest Neighbor classifier to classify frames

Application: Speaker Identification

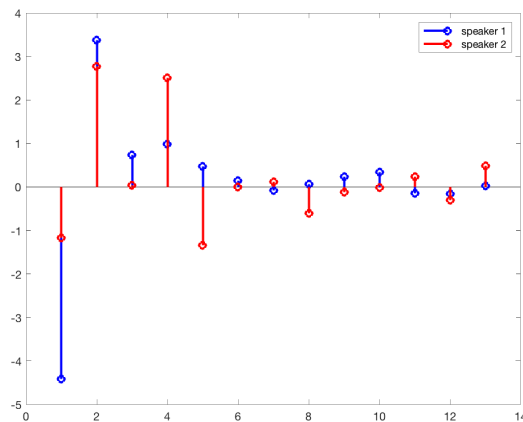
- ▶ AN4 dataset (CMU): 5 male and 5 female subjects speaking words and numbers
- ▶ collect the training samples into frames of 30 ms with an overlap of 75%
- ▶ calculate MFCC
- ▶ train a k-Nearest Neighbor classifier on the frames
- ▶ for a given test signal, predictions are made every frame
- ▶ most frequently occurring label is declared as the speaker

Application: Speaker Identification

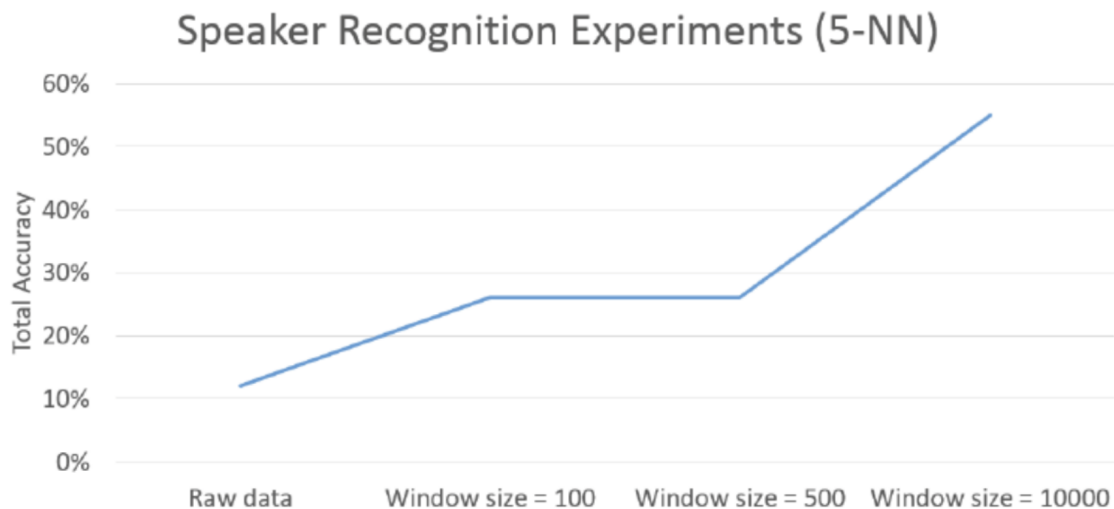
speaker 1 (blue) and speaker 2 (red) time domain signals



frame based MFCC features



Application: Speaker Identification



Application: Speaker Identification

Validation Accuracy													
True Class	fejs	1806	29	27	18	6	6	2	2		5	95.0%	5.0%
	fmjd	32	2137	35	55	25	4		3	1		93.2%	6.8%
	fsrb	50	35	2018	22	19	15	1	4	5	5	92.8%	7.2%
	ftmj	35	71	28	1796	20	6	3	7	4	5	90.9%	9.1%
	fwxs	26	55	17	25	1908	4	2	16	1	8	92.5%	7.5%
	mcen	11	8	2	7	7	1461	19	9	10	13	94.4%	5.6%
	mrcb	23	5	5	8	6	42	1285	5	18	7	91.5%	8.5%
	msjm	12	15	5	16	28	26	3	1262	1	21	90.9%	9.1%
	msjr	15		8		3	16	30	1	1256	3	94.3%	5.7%
	msmn	14	9	7	7	18	21	1	17	2	1404	93.6%	6.4%
		89.2%	90.4%	93.8%	91.9%	93.5%	91.3%	95.5%	95.2%	96.8%	95.4%		
		10.8%	9.6%	6.2%	8.1%	6.5%	8.7%	4.5%	4.8%	3.2%	4.6%		
		fejs	fmjd	fsrb	ftmj	fwxs	mcen	mrcb	msjm	msjr	msmn		
Predicted Class													

- average accuracy is 92.93%