



Books
Publish
About
News
Contact
Author Panel Sign in

Search



What is Open Access?

Open Access is an initiative that aims to make scientific research freely available to all. To date our community has made over 100 million downloads. It's based on principles of collaboration, unobstructed discovery, and, most importantly, scientific progression. As PhD students, we found it difficult to access the research we needed, so we decided to create a new Open Access publisher that levels the playing field for scientists across the world. How? By making research easy to access, and puts the academic needs of the researchers before the business interests of publishers.

Our authors and editors

We are a community of more than 103,000 authors and editors from 3,291 institutions spanning 160 countries, including Nobel Prize winners and some of the world's most-cited researchers. Publishing on IntechOpen allows authors to earn citations and find new collaborators, meaning more people see your work not only from your own field of study, but from other related fields too.

Content Alerts

Brief introduction to this section that describes Open Access especially from an IntechOpen perspective

How it worksManage preferences

Contact

Want to get in touch? Contact our London head office or media team here

Careers

Our team is growing all the time, so we're always on the lookout for smart people who want to help us reshape the world of scientific publishing.

Open access peer-reviewed chapter

Some Commonly Used Speech Feature Extraction Algorithms

By Sabur Ajibola Alim and Nahrul Khair Alang Rashid

Submitted: October 4th 2017 Reviewed: July 20th 2018

Published: December 12th 2018

DOI: 10.5772/intechopen.80419

[Home](#) > [Books](#) > [From Natural to Artificial Intelligence - Algorithms and Applications](#)

Downloaded: 1780



5

Abstract

Speech is a complex naturally acquired human motor ability. It is characterized in adults with the production of about 14 different sounds per second via the harmonized actions of roughly 100 muscles. Speaker recognition is the capability of a software or hardware to receive speech signal, identify the speaker present in the speech signal and recognize the speaker afterwards. Feature extraction is accomplished by changing the speech waveform to a form of parametric representation at a relatively minimized data rate for subsequent processing and analysis. Therefore, acceptable classification is derived from excellent and quality features. Mel Frequency Cepstral Coefficients (MFCC), Linear Prediction Coefficients (LPC), Linear Prediction Cepstral Coefficients (LPCC), Line Spectral Frequencies (LSF), Discrete Wavelet Transform (DWT) and Perceptual Linear Prediction (PLP) are the speech feature extraction techniques that were discussed in these chapter. These methods have been tested in a wide variety of applications, giving them high level of reliability and acceptability. Researchers have made several modifications to the above discussed techniques to make them less susceptible to noise, more robust and consume less time. In conclusion, none of the methods is superior to the other, the area of application would determine which method to select.

Keywords

human speech speech features mel frequency cepstral coefficients (MFCC)
linear prediction coefficients (LPC) linear prediction cepstral coefficients (LPCC)
line spectral frequencies (LSF) discrete wavelet transform (DWT) perceptual linear prediction (PLP)

Chapter and author info

Show +

1. Introduction

Human beings express their feelings, opinions, views and notions orally through speech. The speech production process includes articulation, voice, and fluency [,]. It is a complex naturally acquired human motor abilities, a task categorized in regular adults by the production of about 14 different sounds per

second via the harmonized actions of roughly 100 muscles connected by spinal and cranial nerves. The simplicity with which human beings speak is in contrast to the complexity of the task, and that complexity could assist in explaining why speech can be very sensitive to diseases associated with the nervous system [].

There have been several successful attempts in the development of systems that can analyze, classify and recognize speech signals. Both hardware and software that have been developed for such tasks have been applied in various fields such as health care, government sectors and agriculture. Speaker recognition is the capability of a software or hardware to receive speech signal, identify the speaker present in the speech signal and recognize the speaker afterwards []. Speaker recognition executes a task similar to what the human brain undertakes. This starts from speech which is an input to the speaker recognition system. Generally, speaker recognition process takes place in three main steps which are acoustic processing, feature extraction and classification/recognition [].

The speech signal has to be processed to remove noise before the extraction of the important attributes in the speech [] and identification. The purpose of feature extraction is to illustrate a speech signal by a predetermined number of components of the signal. This is because all the information in the acoustic signal is too cumbersome to deal with, and some of the information is irrelevant in the identification task [,].

Feature extraction is accomplished by changing the speech waveform to a form of parametric representation at a relatively lesser data rate for subsequent processing and analysis. This is usually called the front end signal-processing [,]. It transforms the processed speech signal to a concise but logical representation that is more discriminative and reliable than the actual signal. With front end being the initial element in the sequence, the quality of the subsequent features (pattern matching and speaker modeling) is significantly affected by the quality of the front end [].

Therefore, acceptable classification is derived from excellent and quality features. In present automatic speaker recognition (ASR) systems, the procedure for feature extraction has normally been to discover a representation that is comparatively reliable for several conditions of the same speech signal, even with alterations in the environmental conditions or speaker, while retaining the portion that characterizes the information in the speech signal [,].

Feature extraction approaches usually yield a multidimensional feature vector for every speech signal []. A wide range of options are available to parametrically represent the speech signal for the recognition process, such as perceptual linear prediction (PLP), linear prediction coding (LPC) and mel-frequency cepstrum coefficients (MFCC). MFCC is the best known and very popular [,]. Feature extraction is the most relevant portion of speaker recognition. Features of speech have a vital part in the segregation of a speaker from others []. Feature extraction reduces the magnitude of the speech signal devoid of causing any damage to the power of speech signal [].

Before the features are extracted, there are sequences of preprocessing phases that are first carried out. The preprocessing step is pre-emphasis. This is achieved by passing the signal through a FIR filter [] which is usually a first-order finite impulse response (FIR) filter []. This is succeeded by frame blocking, a method of partitioning the speech signal into frames. It removes the acoustic interface existing in the start and end of the speech signal [].

The framed speech signal is then windowed. Bandpass filter is a suitable window [] that is applied to minimize disjointedness at the start and finish of each frame. The two most famous categories of windows are Hamming and Rectangular windows []. It increases the sharpness of harmonics, eliminates the discontinuous of signal by tapering beginning and ending of the frame zero. It also reduces the spectral distortion formed by the overlap [].

2. Mel frequency cepstral coefficients (MFCC)

Mel frequency cepstral coefficients (MFCC) was originally suggested for identifying monosyllabic words in continuously spoken sentences but not for speaker identification. MFCC computation is a replication of the human hearing system intending to artificially implement the ear's working principle with the assumption that the human ear is a reliable speaker recognizer []. MFCC features are rooted in the recognized discrepancy of the human ear's critical bandwidths with frequency filters spaced linearly at low frequencies and logarithmically at high frequencies have been used to retain the phonetically vital properties of the speech signal. Speech signals commonly contain tones of varying frequencies, each tone with an actual frequency, f (Hz) and the subjective pitch is computed on the Mel scale. The mel-frequency scale has linear frequency spacing below 1000 Hz and logarithmic spacing above 1000 Hz. Pitch of 1 kHz tone and 40 dB above the perceptual audible threshold is defined as 1000 mels, and used as reference point [].

MFCC is based on signal disintegration with the help of a filter bank. The MFCC gives a discrete cosine transform (DCT) of a real logarithm of the short-term energy displayed on the Mel frequency scale []. MFCC is used to identify airline reservation, numbers spoken into a telephone and voice recognition system for security purpose. Some modifications have been proposed to the basic MFCC algorithm for better robustness, such as by lifting the log-mel-amplitudes to an appropriate power (around 2 or 3) before applying the DCT and reducing the impact of the low-energy parts [].

2.1. Algorithm description, strength and weaknesses

MFCC are cepstral coefficients derived on a twisted frequency scale centered on human auditory perception. In the computation of MFCC, the first thing is windowing the speech signal to split the speech signal into frames. Since the high frequency formants process reduced amplitude compared to the low frequency formants, high frequencies are emphasized to obtain similar amplitude for all the formants. After windowing, Fast Fourier Transform (FFT) is applied to find the power spectrum of each frame. Subsequently, the filter bank processing is carried out on the power spectrum, using mel-scale. The DCT is applied to the speech signal after translating the power spectrum to log domain in order to calculate MFCC coefficients []. The formula used to calculate the mels for any frequency is [,]:

$$\text{mel } f = \frac{2595 \times \log_{10} 1 + f / 700}{E1}$$

where $\text{mel}(f)$ is the frequency (mels) and f is the frequency (Hz).

The MFCCs are calculated using this equation [,]:

$$C_n = \sum_{k=1}^K \log S_k \cos(n\pi k - 1.2\pi k) \quad E2$$

where k is the number of mel cepstrum coefficients, S_k is the output of filterbank and C_n is the final mfcc coefficients.

The block diagram of the MFCC processor can be seen in []. It summarizes all the processes and steps taken to obtain the needed coefficients. MFCC can effectively denote the low frequency region better than the high frequency region, henceforth, it can compute formants that are in the low frequency range and describe the vocal tract resonances. It has been generally recognized as a front-end procedure for typical Speaker Identification applications, as it has reduced vulnerability to noise disturbance, with minute session inconsistency and easy to mine []. Also, it is a perfect representation for sounds when the source characteristics are stable and consistent (music and speech) []. Furthermore, it can capture information from sampled signals with frequencies at a maximum of 5 kHz, which encapsulates most energy of sounds that are generated by humans [].

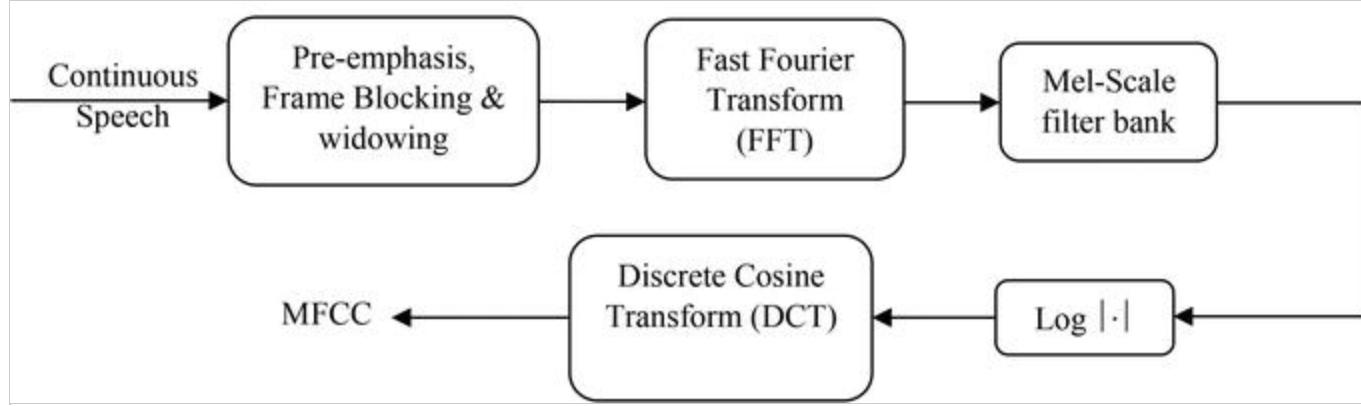


Figure 1.

Block diagram of MFCC processor.

Cepstral coefficients are said to be accurate in certain pattern recognition problems relating to human voice. They are used extensively in speaker identification and speech recognition []. Other formants can also be above 1 kHz and are not efficiently taken into consideration by the large filter spacing in the high frequency range []. MFCC features are not exactly accurate in the existence of background noise [,] and might not be well suited for generalization [].

3. Linear prediction coefficients (LPC)

Linear prediction coefficients (LPC) imitates the human vocal tract [] and gives robust speech feature. It evaluates the speech signal by approximating the formants, getting rid of its effects from the speech signal and estimate the concentration and frequency of the left behind residue. The result states each sample of the signal as a direct incorporation of previous samples. The coefficients of the difference equation characterize the formants, thus, LPC needs to approximate these coefficients []. LPC is a powerful speech analysis method and it has gained fame as a formant estimation method [].

The frequencies where the resonant crests happen are called the formant frequencies. Thus, with this technique, the positions of the formants in a speech signal are predictable by calculating the linear predictive coefficients above a sliding window and finding the crests in the spectrum of the subsequent linear

prediction filter []. LPC is helpful in the encoding of high quality speech at low bit rate [, ,].

Other features that can be deduced from LPC are linear predication cepstral coefficients (LPCC), log area ratio (LAR), reflection coefficients (RC), line spectral frequencies (LSF) and Arcus Sine Coefficients (ARCSIN) []. LPC is generally used for speech reconstruction. LPC method is generally applied in musical and electrical firms for creating mobile robots, in telephone firms, tonal analysis of violins and other string musical gadgets [].

3.1. Algorithm description, strength and weaknesses

Linear prediction method is applied to obtain the filter coefficients equivalent to the vocal tract by reducing the mean square error in between the input speech and estimated speech []. Linear prediction analysis of speech signal forecasts any given speech sample at a specific period as a linear weighted aggregation of preceding samples. The linear predictive model of speech creation is given as [,]:

$$\hat{s}_n = \sum_{k=1}^p a_k s_{n-k} \quad E_3$$

where \hat{s} is the predicted sample, s is the speech sample, p is the predictor coefficients.

The prediction error is given as [,]:

$$e_n = s_n - \hat{s}_n \quad E_4$$

Subsequently, each frame of the windowed signal is autocorrelated, while the highest autocorrelation value is the order of the linear prediction analysis. This is followed by the LPC analysis, where each frame of the autocorrelations is converted into LPC parameters set which consists of the LPC coefficients []. A summary of the procedure for obtaining the LPC is as seen in []. LPC can be derived by []:

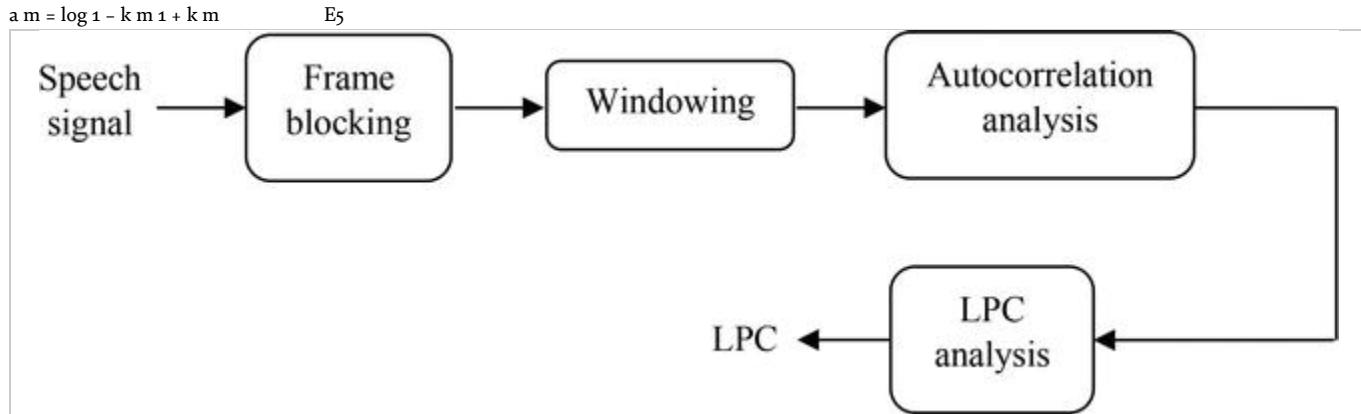


Figure 2.

Block diagram of LPC processor.

where a_m is the linear prediction coefficient, k_m is the reflection coefficient.

Linear predictive analysis efficiently selects the vocal tract information from a given speech []. It is known for the speed of computation and accuracy []. LPC excellently represents the source behaviors that are steady and consistent []. Furthermore, it is also be used in speaker recognition system where the main purpose is to extract the vocal tract properties []. It gives very accurate estimates of speech parameters and is comparatively efficient for computation [,]. Traditional linear prediction suffers from aliased autocorrelation coefficients []. LPC estimates have high sensitivity to quantization noise [] and might not be well suited for generalization [].

4. Linear prediction cepstral coefficients (LPCC)

Linear prediction cepstral coefficients (LPCC) are cepstral coefficients derived from LPC calculated spectral envelope []. LPCC are the coefficients of the Fourier transform illustration of the logarithmic magnitude spectrum [,] of LPC. Cepstral analysis is commonly applied in the field of speech processing

because of its ability to perfectly symbolize speech waveforms and characteristics with a limited size of features [].

It was observed by Rosenberg and Sambur that adjacent predictor coefficients are highly correlated and therefore, representations with less correlated features would be more efficient, LPCC is a typical example of such. The relationship between LPC and LPCC was originally derived by Atal in 1974. In theory, it is relatively easy to convert LPC to LPCC, in the case of minimum phase signals [].

4.1. Algorithm description, strength and weaknesses

In speech processing, LPCC analogous to LPC, are computed from sample points of a speech waveform, the horizontal axis is the time axis, while the vertical axis is the amplitude axis []. The LPCC processor is as seen in . It pictorially explains the process of obtaining LPCC. LPCC can be calculated using [, ,]:

$$C_m = a_m + \sum k = 1 m - 1 k m c k a m - k \quad E6$$

where a_m is the linear prediction coefficient, C_m is the cepstral coefficient.

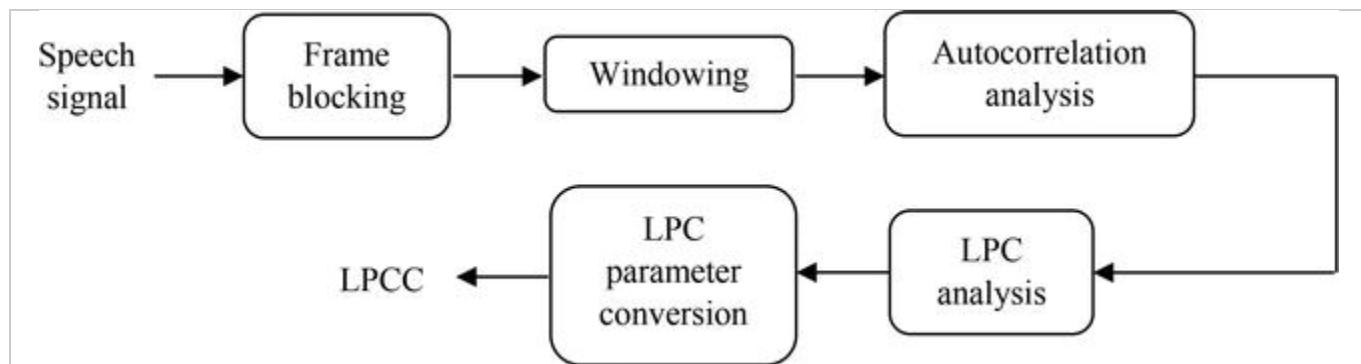


Figure 3.

Block diagram of LPCC processor.

LPCC have low vulnerability to noise []. LPCC features yield lower error rate as compared to LPC features []. Cepstral coefficients of higher order are mathematically limited, resulting in an extremely extensive array of variances when moving from the cepstral coefficients of lower order to cepstral coefficients of higher order []. Similarly, LPCC estimates are notorious for having great sensitivity to quantization noise []. Cepstral analysis on high-pitch speech signal gives small source-filter separability in the quefrency domain []. Cepstral coefficients of lower order are sensitive to the spectral slope, while the cepstral coefficients of higher order are sensitive to noise [].

5. Line spectral frequencies (LSF)

Individual lines of the Line Spectral Pairs (LSP) are known as line spectral frequencies (LSF). LSF defines the two resonance situations taking place in the inter-connected tube model of the human vocal tract. The model takes into consideration the nasal cavity and the mouth shape, which gives the basis for the fundamental physiological importance of the linear prediction illustration. The two resonance situations define the vocal tract as either being completely open or completely closed at the glottis []. The two situations begets two groups of resonant frequencies, with the number of resonances in each group being deduced from the quantity of linked tubes. The resonances of each situation are the odd and even line spectra correspondingly, and are interwoven into a singularly rising group of LSF [].

The LSF representation was proposed by Itakura [,] as a substitute to the linear prediction parametric illustration. In the area of speech coding, it has been realized that this illustration has an improved quantization features than the other linear prediction parametric illustrations (LAR and RC). The LSF illustration has the capacity to reduce the bit-rate by 25–30% for transmitting the linear prediction information without distorting the quality of synthesized speech [, ,]. Apart from quantization, LSF illustration of the predictor are also suitable for interpolation. Theoretically, this can be inspired by the point that the sensitivity matrix linking the LSF-domain squared quantization error to the perceptually relevant log spectrum is diagonal [,].

5.1. Algorithm description, strength and weaknesses

LP is established on the point that a speech signal can be defined by . Recall

$s' n = \sum k=1 p a k s n - k$
where k is the time index and p is the order of the linear prediction, $s' n$ is the predictor signal and $a k$ is the LPC coefficients.

The $a k$ coefficients are determined in order to reduce the prediction error by method of autocorrelation or covariance. can be modified in the frequency domain with the z -transform. As such, a small part of the speech signal is anticipated to be given as an output to the all-pole filter $H z$. The new equation is

$$H z = 1 A z = 1 1 - \sum i=1 p a i z^{-1} \quad E7$$

where $H z$ is the all-pole filter and $A z$ is the LPC analysis filter.

In order to compute the LSF coefficients, an inverse polynomial filter is split into two polynomials $P z$ and $Q z$ [, , ,]:

$$P z = A z + z - p + 1 A z - 1 \quad E8$$

$$Q z = A z - z - p + 1 A z - 1 \quad E9$$

where $P z$ is the vocal tract with the glottis closed, $Q z$ is the LPC analysis filter of order p .

In order to convert LSF back to LPC, the equation below is used [, , ,]:

$$A z = 0.5 P z + Q z \quad E10$$

The block diagram of the LSF processor is as seen in . The most prominent application of LSF is in the area of speech compression, with extension into the speaker recognition and speech recognition. This technique has also found restricted use in other fields. LSF have been investigated for use in musical instrument recognition and coding. LSF have also been applied to animal noise identification, recognizing individual instruments and financial market analysis. The advantages of LSF include their ability to localize spectral sensitivities, the fact that they characterize bandwidths and resonance locations and lays emphasis on the important aspect of spectral peak location. In most instances, the LSF representation provides a near-minimal data set for subsequent classification [].

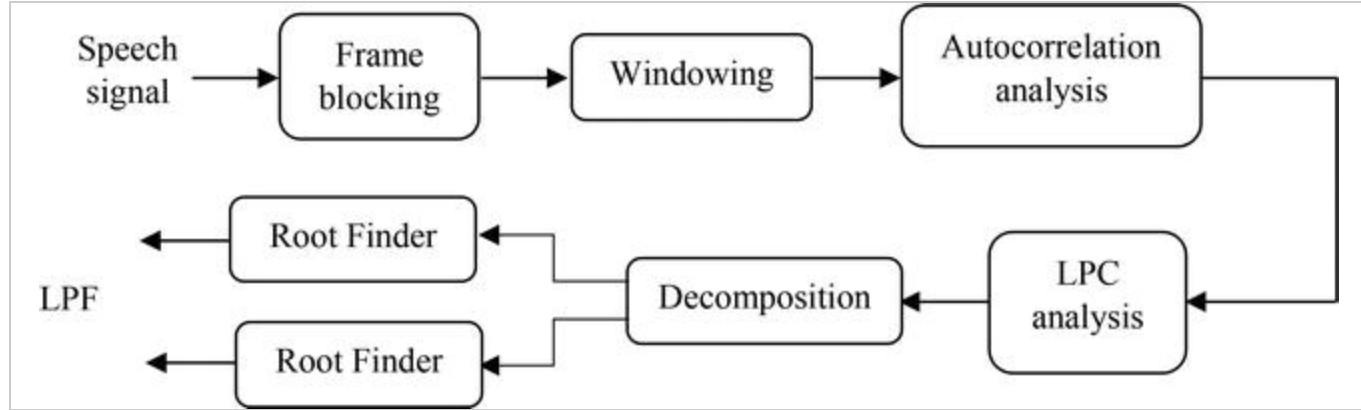


Figure 4.

Block diagram of LSF processor.

Since LSF represents spectral shape information at a lower data rate than raw input samples, it is reasonable that a careful use of processing and analysis methods in the LSP domain could lead to a complexity reduction against alternative techniques operating on the raw input data itself. LSF play an important role in the transmission of vocal tract information from speech coder to decoder with their widespread use being a result of their excellent quantization properties. The generation of LSP parameters can be accomplished using several methods, ranging in complexity. The major problem revolves around finding the roots of the P and Q polynomials defined in and . This can be obtained through standard root solving methods, or more obscure methods and it is often performed in the cosine domain [].

6. Discrete wavelet transform (dwt)

Wavelet Transform (WT) theory is centered around signal analysis using varying scales in the time and frequency domains []. With the support of theoretical physicist Alex Grossmann, Jean Morlet introduced wavelet transform which permits high-frequency events identification with an enhanced temporal resolution [,]. A wavelet is a waveform of effectively limited duration that has an average value of zero. Many wavelets also display orthogonality, an ideal feature of compact signal representation []. WT is a signal processing technique that can be used to represent real-life non-stationary signals with high efficiency [,]. It has the ability to mine information from the transient signals concurrently in both time and frequency domains [, ,].

Continuous wavelet transform (CWT) is used to split a continuous-time function into wavelets. However, there is redundancy of information and huge computational efforts is required to calculate all likely scales and translations of CWT, thereby restricting its use []. Discrete wavelet transform (DWT) is an extension of the WT that enhances the flexibility to the decomposition process []. It was introduced as a highly flexible and efficient method for sub band breakdown of signals [,]. In earlier applications, linear discretization was used for discretizing CWT. Daubechies and others have developed an orthogonal DWT specially designed for analyzing a finite set of observations over the set of scales (dyadic discretization) [].

6.1. Algorithm description, strength and weaknesses

Wavelet transform decomposes a signal into a group of basic functions called wavelets. Wavelets are obtained from a single prototype wavelet called mother wavelet by dilations and shifting. The main characteristic of the WT is that it uses a variable window to scan the frequency spectrum, increasing the temporal resolution of the analysis [, ,].

WT decomposes signals over translated and dilated mother wavelets. Mother wavelet is a time function with finite energy and fast decay. The different versions of the single wavelet are orthogonal to each other. The continuous wavelet transform (CWT) is given by [, ,]:

$$W(x, a, b) = \int_{-\infty}^{\infty} x(t) \psi(t - b/a) dt \quad E11$$

where $\psi(t)$ is the mother wavelet, a and b are continuous parameters.

The WT coefficient is an expansion and a particular shift represents how well the original signal corresponds to the translated and dilated mother wavelet. Thus, the coefficient group of $CWT(a, b)$ associated with a particular signal is the wavelet representation of the original signal in relation to the mother wavelet []. Since CWT contains high redundancy, analyzing the signal using a small number of scales with varying number of translations at each scale, i.e. discretizing scale and translation parameters as $a = 2^j$ and $b = 2^j k$ gives DWT. DWT theory requires two sets of related functions called scaling function and wavelet function given by []:

$$\phi(t) = \sum n = 0 N - 1 h(n) \phi_2(t - n) \quad E12$$

$$\psi(t) = \sum n = 0 N - 1 g(n) \phi_2(t - n) \quad E13$$

where $\phi(t)$ is the scaling function, $\psi(t)$ is the wavelet function, $h[n]$ is the an impulse response of a low-pass filter, and $g[n]$ is an impulse response of a high-pass filter.

There are several ways to discretize a CWT. The DWT of the continuous signal can also be given by []:

$$DWT(m, p) = \int_{-\infty}^{\infty} x(t) \psi(m, p) dt \quad E14$$

where $\psi(m, p)$ is the wavelet function bases, m is the dilation parameter and p is the translation parameter.

Thus, $\psi(m, p)$ is defined as:

$$\psi(m, p) = \int_{-\infty}^{\infty} \phi(t) \psi(t - p) dt \quad E15$$

The DWT of a discrete signal is derived from CWT and defined as:

$$DWT(m, k) = \sum n = 0 m - 1 x(n) \psi(m, k) \quad E16$$

where $\psi(m, k)$ is the mother wavelet and $x(n)$ is the discretized signal. The mother wavelet may be dilated and translated discretely by selecting the scaling parameter $a = a_0 m$ and translation parameter $b = n b_0 a_0 m$ (with constants taken as $a > 1$, $b > 1$, while m and n are assigned a set of positive integers).

The scaling and wavelet functions can be implemented effectively using a pair of filters, $h[n]$ and $g[n]$, called quadrature mirror filters that confirm with the property $g(n) = -1 \cdot h(n)$. The input signal is filtered by a low-pass filter and high-pass filter to obtain the approximate components and the detail components respectively. This is summarized in []. The approximate signal at each stage is further decomposed using the same low-pass and high-pass

filters to get the approximate and detail components for the next stage. This type of decomposition is called dyadic decomposition [1].

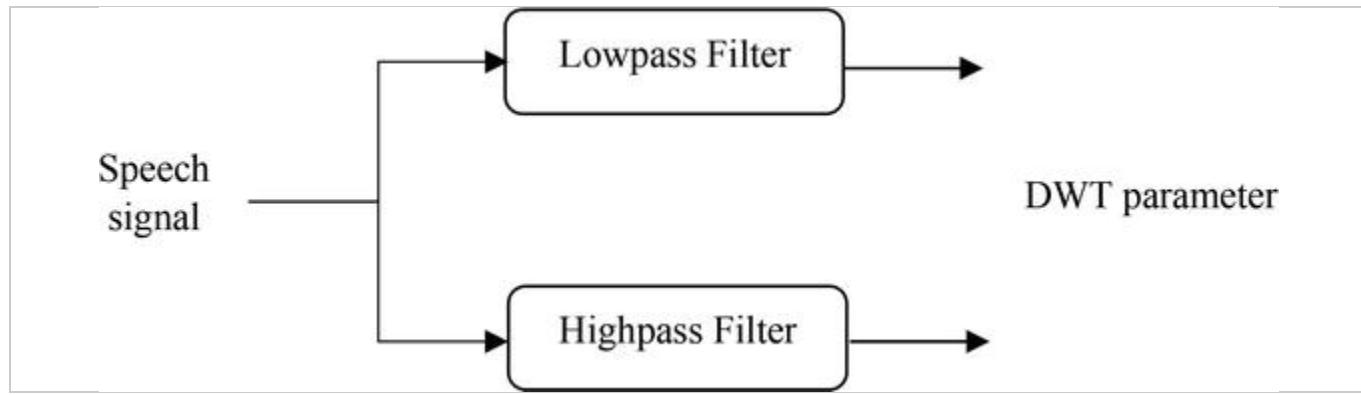


Figure 5.

Block diagram of DWT.

The DWT parameters contain the information of different frequency scales. This enhances the speech information obtained in the corresponding frequency band [2]. The ability of the DWT to partition the variance of the elements of the input on a scale by scale basis is an added advantage. This partitioning leads to the opinion of the scale-dependent wavelet variance, which in many ways is equivalent to the more familiar frequency-dependent Fourier power spectrum [3]. Classic discrete decomposition schemes, which are dyadic do not fulfill all the requirements for direct use in parameterization. DWT does provide adequate number of frequency bands for effective speech analysis [4]. Since the input signals are of finite length, the wavelet coefficients will have unwantedly large variations at the boundaries because of the discontinuities at the boundaries [5].

7. Perceptual linear prediction (PLP)

Perceptual linear prediction (PLP) technique combines the critical bands, intensity-to-loudness compression and equal loudness pre-emphasis in the extraction of relevant information from speech. It is rooted in the nonlinear bark scale and was initially intended for use in speech recognition tasks by eliminating the speaker dependent features [6]. PLP gives a representation conforming to a smoothed short-term spectrum that has been equalized and compressed similar to the human hearing making it similar to the MFCC. In the PLP approach, several prominent features of hearing are replicated and the consequent auditory like spectrum of speech is approximated by an autoregressive all-pole model [7]. PLP gives minimized resolution at high frequencies that signifies auditory filter bank based approach, yet gives the orthogonal outputs that are similar to the cepstral analysis. It uses linear predictions for spectral smoothing, hence, the name is perceptual linear prediction [8]. PLP is a combination of both spectral analysis and linear prediction analysis.

7.1. Algorithm description, strength and weaknesses

In order to compute the PLP features the speech is windowed (Hamming window), the Fast Fourier Transform (FFT) and the square of the magnitude are computed. This gives the power spectral estimates. A trapezoidal filter is then applied at 1-bark interval to integrate the overlapping critical band filter responses in the power spectrum. This effectively compresses the higher frequencies into a narrow band. The symmetric frequency domain convolution on the bark warped frequency scale then permits low frequencies to mask the high frequencies, concurrently smoothing the spectrum. The spectrum is subsequently pre-emphasized to approximate the uneven sensitivity of human hearing at a variety of frequencies. The spectral amplitude is compressed, this reduces the amplitude variation of the spectral resonances. An Inverse Discrete Fourier Transform (IDCT) is performed to get the autocorrelation coefficients. Spectral smoothing is performed, solving the autoregressive equations. The autoregressive coefficients are converted to cepstral variables [9]. The equation for computing the bark scale frequency is:

$$\text{bark } f = 26.81 f_{1960} + f - 0.53 \quad E17$$

where $\text{bark}(f)$ is the frequency (bark) and f is the frequency (Hz).

The identification achieved by PLP is better than that of LPC [10], because it is an improvement over the conventional LPC because it effectively suppresses the speaker-dependent information [11]. Also, it has enhanced speaker independent recognition performance and is robust to noise, variations in the channel and microphones [12]. PLP reconstructs the autoregressive noise component accurately [13]. PLP based front end is sensitive to any change in the formant frequency.

shows the PLP processor, showing all the steps to be taken to obtain the PLP coefficients. PLP has low sensitivity to spectral tilt, consistent with the findings that it is relatively insensitive to phonetic judgments of the spectral tilt. Also, PLP analysis is dependent on the result of the overall spectral balance (formant amplitudes). The formant amplitudes are easily affected by factors such as the recording equipment, communication channel and additive noise []. Furthermore, the time-frequency resolution and efficient sampling of the short-term representation are addressed in an ad-hoc way [].

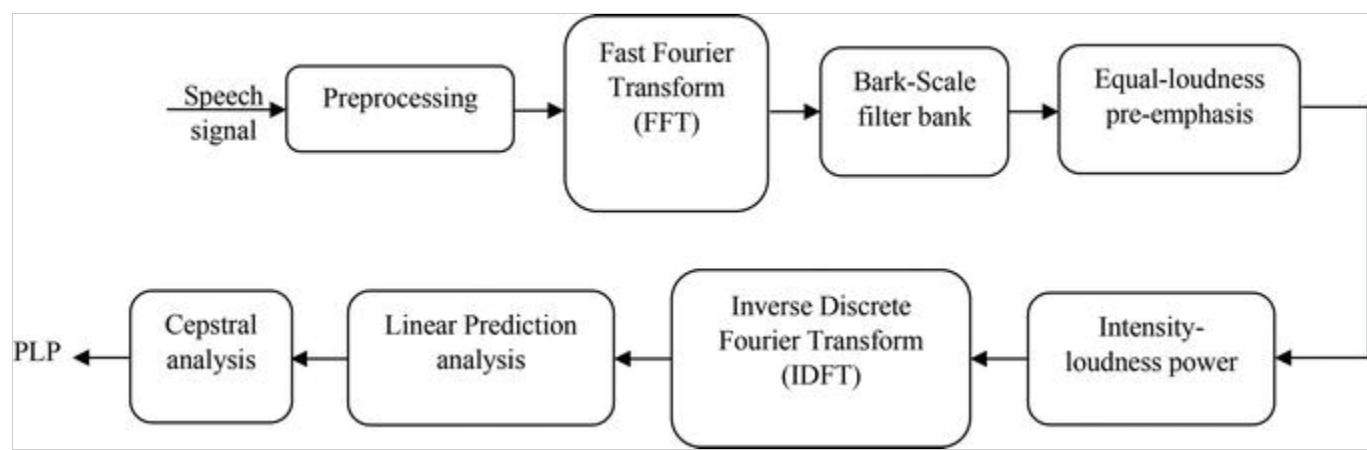


Figure 6.

Block diagram of PLP processor.

shows a comparison between the six feature extraction techniques that have been explicitly described above. Even though the selection of a feature extraction algorithm for use in research is individual dependent, however, this table has been able to characterize these techniques based on the main considerations in the selection of any feature extraction algorithm. The considerations include speed of computation, noise resistance and sensitivity to additional noise. The table also serves as a guide when considering the selection between any two or more of the discussed algorithms.

	Type of Filter	Shape of filter	What is modeled	Speed of computation	Type of coefficient	Noise resistance	Sensitivity to quantization/additional noise
Mel frequency cepstral coefficient (MFCC)	Mel	Triangular	Human Auditory System	High	Cepstral	Medium	Medium
Linear prediction coefficient (LPC)	Linear Prediction	Linear	Human Vocal Tract	High	Autocorrelation Coefficient	High	High
Linear prediction cepstral coefficient (LPCC)	Linear Prediction	Linear	Human Vocal Tract	Medium	Cepstral	High	High
Line spectral frequencies (LSF)	Linear Prediction	Linear	Human Vocal Tract	Medium	Spectral	High	High
Discrete wavelet transform (DWT)	Lowpass & highpass	—	—	High	Wavelets	Medium	Medium

	Type of Filter	Shape of filter	What is modeled	Speed of computation	Type of coefficient	Noise resistance	Sensitivity to quantization/addition noise
Perceptual linear prediction (PLP)	Bark	Trapezoidal	Human Auditory System	Medium	Cepstral & Autocorrelation	Medium	Medium

Table 1.

Comparison between the feature extraction techniques.

8. Conclusion

MFCC, LPC, LPCC, LSF, PLP and DWT are some of the feature extraction techniques used for extracting relevant information from speech signals for the purpose of speech recognition and identification. These techniques have stood the test of time and have been widely used in speech recognition systems for several purposes. Speech signal is a slow time varying signal, quasi-stationary, when observed over an adequately short period of time between 5 and 100 msec, its behavior is relatively stationary. As a result of this, short time spectral analysis which includes MFCC, LPCC and PLP are commonly used for the extraction of important information from speech signals. Noise is a serious challenge encountered in the process of feature extraction, as well as speaker recognition as a whole. Subsequently, researchers have made several modifications to the above discussed techniques to make them less susceptible to noise, more robust and consume less time. These methods have also been used in the recognition of sounds. The extracted information will be the input to the classifier for identification purposes. The above discussed feature extraction approaches can be implemented using MATLAB.

Download for free
 SHARE THIS CHAPTER

 DOWNLOAD FOR FREE

SECTIONS

 Chapter and author info



2. Mel frequency cepstral coefficients (MFCC)

Share

3. Linear prediction coefficients (LPC)

 cepstral coefficients (LPCC)

5. Line spectral frequencies (LSF)

 transform (dwt)

 prediction (PLP)

8. Conclusion

>



More

[PRINT CHAPTER](#)

© 2018 The Author(s). Licensee IntechOpen. This chapter is distributed under the terms of the [Creative Commons Attribution 3.0 License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

How to cite and referenceLink to this chapter [Copy to clipboard](#)

<https://www.intechopen.com/books/from-natural-to-artificial-intelligence-algorithms-and-applications/some-commonly-used-speech-feature-extraction-algorithms>

Cite this chapter [Copy to clipboard](#)

Sabur Ajibola Alim and Nahrul Khair Alang Rashid (December 12th 2018). Some Commonly Used Speech Feature Extraction Algorithms, From Natural to Artificial Intelligence - Algorithms and Applications, Ricardo Lopez-Ruiz, IntechOpen, DOI: 10.5772/intechopen.80419. Available from: <https://www.intechopen.com/books/from-natural-to-artificial-intelligence-algorithms-and-applications/some-commonly-used-speech-feature-extraction-algorithms>

Over 21,000 IntechOpen readers like this topic

Help us write another book on this subject and reach those readers

[SUGGEST A BOOK TOPIC](#)

BOOKS OPEN FOR SUBMISSIONS

Chapter statistics

1780 total chapter downloads

4 Crossref citations



More statistics for editors and authors

Login to your personal dashboard for more detailed statistics on your publications.

ACCESS PERSONAL REPORTING

Related Content

This Book



InTechOpen

From Natural to Artificial Intelligence
Algorithms and Applications

Edited by Ricardo Lopez-Ruiz



From Natural to Artificial Intelligence

Edited by

Convolutional Neural Networks for Raw Speech Recognition

By Vishal Passricha and Rajesh Kumar Aggarwal

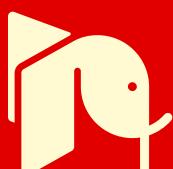
Related Book



InTechOpen

Numerical Simulation
From Brain Imaging to Turbulent Flows

Edited by Ricardo Lopez-Ruiz



Numerical Simulation

Edited by

BOLD fMRI Simulation

By Zikuan Chen and Vince Calhoun

We are IntechOpen, the world's leading publisher of Open Access books. Built by scientists, for scientists. Our readership spans scientists, professors, researchers, librarians, and students, as well as business professionals. We share our knowledge and peer-reviewed research papers with libraries, scientific and engineering societies, and also work with corporate R&D departments and government entities.

5,100

+108,170

Open Access BooksCitations in Web of Science

126,000

+560,000

IntechOpen Authors and Academic EditorsUnique Visitors per Month

MORE ABOUT US

Your recently viewed content

Chapter

Some Commonly
Used Speech Feature
Extraction Algor...

By Sabur Ajibola Alim
and Nahrul Khair
Alang Rashid

Call for Authors
Submit your work to IntechOpen

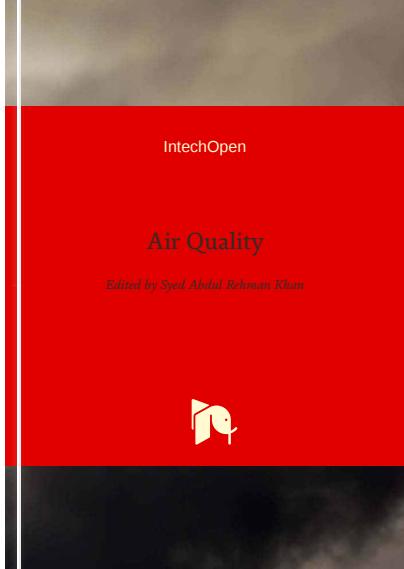
In addition to well-established academics in the field of scientific research, we also welcome and encourage the next up-and-coming generation of scientists looking to make their mark.

No matter where you are in your career, we would welcome you and encourage you to consider joining our community.



[VIEW ALL BOOKS OPEN FOR CHAPTER SUBMISSIONS](#)

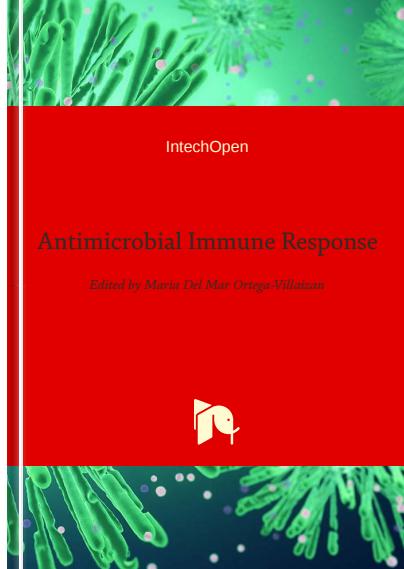
Books open for chapter submissions



Air Quality

Edited by Syed Abdul Rehman Khan

Open for chapter submissions



Antimicrobial Immune Response

Edited by Maria Del Mar Ortega-Villaizan

Open for chapter submissions



Book Subject Areas

Physical Sciences, Engineering and Technology

Chemistry (169)

Computer and Information Science (433)

Earth and Planetary Sciences (172)

Engineering (840)

Materials Science (279)

Mathematics (55)

Nanotechnology and Nanomaterials (114)

Physics (138)

Robotics (99)

Technology (111)

More >

Life Sciences

Agricultural and Biological Sciences (341)

Biochemistry, Genetics and Molecular Biology (273)

Environmental Sciences (171)

Immunology and Microbiology (64)

Neuroscience (59)

More >

Health Sciences

Medicine (1490)

Pharmacology, Toxicology and Pharmaceutical Science (76)

Veterinary Medicine and Science (28)

Social Sciences and Humanities

Business, Management and Economics (120)

Psychology (23)

Social Sciences (80)

>

Home

News

Contact

Careers

About

Our Authors and Editors

Scientific Advisors

Team

Events

Partnerships

Publish

About Open Access

How it works

OA Publishing Fees

Open Access funding

Peer Reviewing

Editorial Policies



Headquarters

IntechOpen Limited

5 Princes Gate Court,
London, SW7 2QJ,
UNITED KINGDOM

Phone: +44 20 8089 5702

Search



AUTHOR PANEL SIGN IN

[Terms and Conditions](#) | [Privacy Policy](#) | [Customer Complaints](#)

© 2020 IntechOpen. All rights reserved.