

Chapter 1

Speech Recognition in a Medical Environment

1.1 Background

Automatic speech recognition (ASR) has been a topic that interested many from an early age. Many consider it to have started in the 1950s with Bell Labs' *Audrey* [Moskvitch, 2017], which was able to do a single-speaker digit recognition. Seventy years later, the ASR technologies have grown so that they are present within our personal homes, with Amazon Alexa and Google Home devices [Fowler, 2018], as they have become more affordable. As the speech technologies advance, they approach the concept of ubiquitous computing (often also known as ambient intelligence), which is highly desirable for many industries. Although most of the ASR systems available for commerce are trained in normal language and not medical language, one industry that could greatly benefit from ASR technologies is the health care system, given that the ASR system is able to understand medical language.

Health care is a system which has a very high demand and the services are required to be as detailed as possible due to various reasons, one of which is the concept that it is paramount for a hospital to have a clear and rich track of a patient's medical history. Currently, in order to maintain a patient's medical history, a physician sees multiple people during his/her working hours, and only after their shift is over, does (s)he sit down to write the medical notes. The lag in between seeing a patient and taking notes may sometimes go up to 8-10 hours, and then those notes are often inaccurate. Having an ASR system in a medical environment could greatly help in keeping track of a patient's medical history, where a physician could easily dictate the notes. In a hospital environment, however, such as in an intensive care unit (ICU), a physician who is trying to dictate notes may find him or herself in trouble, as in the ICU, there are multiple background sounds from machines, multiple people speaking, and a great amount of white noise.

Current companies have been creating ASR systems that can understand medical language and allows physicians to dictate their medical notes. One major company that has been "dominating" a lot of the market is Nuance, with their Dragon Medical system. Unfortunately, their dictation system is not yet capable of inferring punctuation marks and markup language onto the text, thus, in order to dictate a segment such as: "37-year-old female presents complainint of urinary frequency, urgency and dysuria along with hematuria and low-grade fever." needs to be dictated as:

"37-year-old female presents complainint of urinary frequency **comma** urgency and dysuria along with hematuria and low-grade fever **period**"

Although the system has very high accuracy results, the system works best when one is in a quiet environment, which is not always a realistic scenario. The goal of this project is to create a transcription engine that can recognize medical language in a noisy environment.

1.2 Preliminary Data

1.2.1 Comparison of Different Commercial Engines

Compare commercial (+ CMU) Text- \rightarrow Voice- \rightarrow Text MIMIC II + 20kLeagues

1.2.2 Grade Reading Level of Various Texts

MIMIC II + 20kLeagues

1.3 Automatic Speech Recognition Schematic for a Medical Setting

In order to have a speech recognition system to be able to understand medical language, it is important to obtain a big picture of the project.

[INSERT FIGURE]

As one can see from the diagram above, the sound signals would come to a microphone or an array of microphones, which would then be passed through a system that performs blind source separation on the sound signal. It is very possible that, in a hospital environment, there may be multiple sources, such as the heart rate monitor, the ventilator, the healthcare providers talking, and others. A blind source separation algorithm would ensure to separate the audio channel to multiple sources. Once the sources have been separated, there would be a classification methodology to identify the physician from the patient, and ...

1.3.1 AI-Complete and ASR-Complete

Although there have been many advancements in the speech technologies, especially through black boxes such as neural networks, and there currently exists very accurate engines to do voice recognition, it still remains to be an unsolved problem. The speech recognition problem is ASR-complete, which falls under the umbrella of AI-complete, which, by analogy of NP-completeness, means that the problem is hypothesized to deal with multiple parts of the AI world, and it cannot be solved by a single algorithm. Current technologies are not able to solve AI-complete problems, and they often require human computation. Since speech recognition is an AI-complete problem, there is still a lot of room for improvement in the fundamental problems that have not yet been solved, and although deep learning has recently shown promising results, deep learning neural networks (DNN) do not necessarily solve fundamental problems, as no one has a truthful understanding of how they work and behave. Therefore, there are a lot of possibilities for the formulation design of the framework (i.e. pipeline).

1.4 Speech Enhancement

Speech enhancement is one of fundamental problems of speech that has not yet been fully solved yet, as noise is subjective, and there isn't necessarily one specific way of performing a fully accurate noise reduction for any type of noise. As previously mentioned before, a ventilator and heart rate monitor are examples of background noise that do not necessarily have a Gaussian distribution (i.e. it is not a white noise), yet, for ASR, they are both irrelevant. There are a number of methodologies that have been created to attempt to solve the problem of speech enhancement, and they may potentially be useful to be implemented onto the pipeline.

1.4.1 Recurrent Neural Networks for Noise Suppression

The work on recurrent neural networks (RNN) was originally worked sequentially by [Jordan, 1986] and [Barak, 1988], and it finally came to be coined by [Cleeremans et al., 1989]. They have shown to work very efficiently with time-series data as their cyclic nature allows them to adapt to incoming inputs in a sequential form, working often times better than convolutional neural networks (CNN), when dealing with time-series data. Other variants of RNNs have been created, which allow them to store "memory" for future cases, with Long Short-Term Memory (LSTM) [Hochreiter and Schmidhuber, 1997], which were then further expanded with gates for Gated Recurrent Units (GRU) [Gers et al., 2000] specifically to allow such units to forget information, and a visual representation of LSTMs and GRUs are shown on Figure 1.1

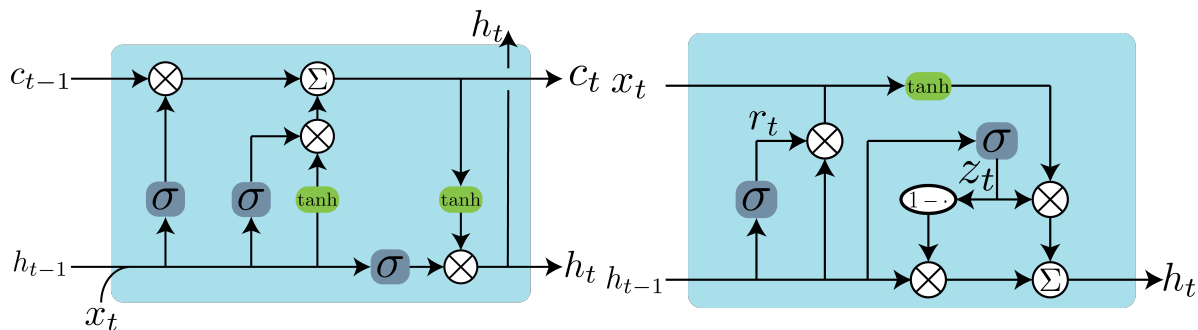


Figure 1.1: A visual representation of the LSTM cell (Left) and the GRU cell (Right). The LSTM contains three representations of its states: the cell state/memory c_t , the hidden states h_t , and the input state x_t , all of which occurs at time t . On the GRU, there are the input state x_t , the hidden state h_t , the update gate z_t and the reset gate r_t . The σ block is used to represent a sigmoid function, and the \tanh block is used to represent a hyperbolic tangent for both images. Please see A.1 for the mathematical expressions of the cells above

Because of the nature of the network being able to receive continuous sequential inputs, it becomes a very attractive model to potentially behave as an adaptive filter. Hence, members from the Xiph.Org Foundation, a non-profit organization, along with Mozilla have created a RNN architecture that is based on GRUs called RRNoise [Jean-Marc et al., 2017], which, after trained, is able to reduce a lot of the background noise of audio files, thus enhancing the signal. Below on Figure 1.2 is the topology of the RRNoise architecture, where it outputs voice activity detection (VAD) as well as the gains from the input features.

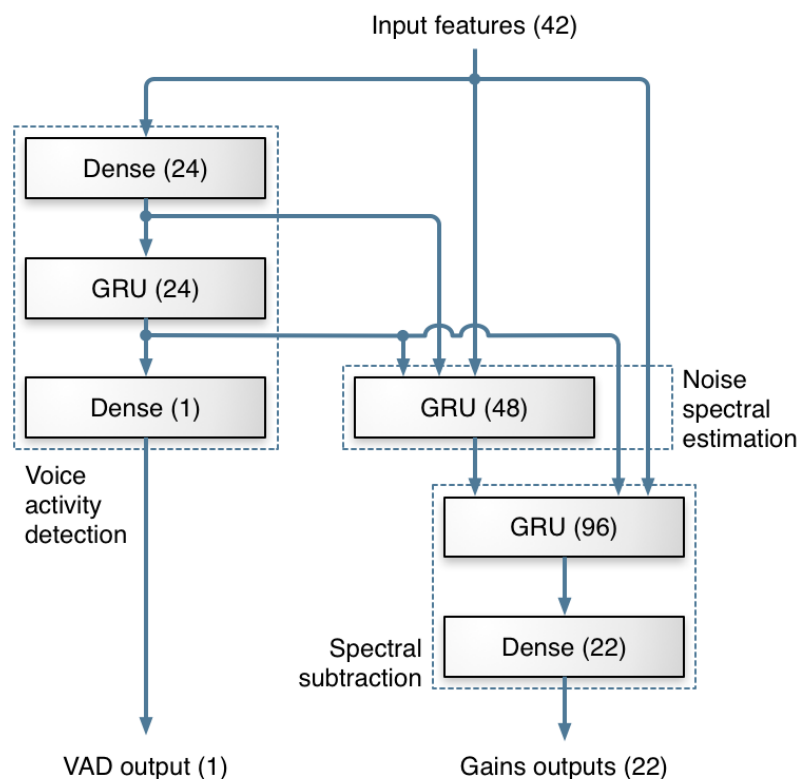


Figure 1.2: The topology of the RRNoise network used to perform noise reduction. This figure is taken from [Jean-Marc et al., 2017]

1.4.2 Non-Negative Matrix Factorization

Another very neat way of performing speech enhancement is via non-negative matrix factorization. Please skip to the Blind Source Separation's subsection on [Non-Negative Matrix Factorization](#) [Ng et al., 2001]

1.5 Blind Source Separation

Given that a recording seldom contains one and only one speaker speaking at a time, an ASR may find itself troubled in detecting who is the speaker and who should be taken as the background. This is namely known as the Cocktail Party Effect. Its name originates from the idea that, in a cocktail party, there are multiple people speaking to one another. While a human is able to have a conversation with another, focusing on his/her friend while having all of the other speakers (i.e. sources) as background, an ASR system is not able to do so. Hence it is important to be able to separate a mixed signal to then be able to focus on a single individual and this is called blind source separation (BSS). There are two methodologies that are proposed to perform BSS: independent component analysis, and non-negative matrix factorization.

1.5.1 Independent Component Analysis

1.5.2 Non-Negative Matrix Factorization (NMF) and Sparse NMF (SNMF)

Background

One of the major drawbacks from any methodology in ICA is that it requires N or more observations (i.e. microphone recordings) for N sources (e.g. speakers, noise), which, depending on the circumstance, may be limiting. The method of non-negative matrix factorization (NMF), however, is able to perform BSS with very high accuracies, while using a single-channel recording. NMF originated in the field of chemometrics, where chemists were trying to identify contents of solutions via positive matrix factorization [Paatero and Tapper, 1994] [Anttila et al., 1995]. It was later further developed in [Lee and Seung, 1999] and algorithms were then better formulated in [Lee and Seung, 2000].

The concept behind NMF is that a mixed signal should be able to be decomposed into its unmixed signals, under the assumption that all of the superpositioned signals are non-negatively added to one another. In other words, they are positively added, but their value may also be zero. Conceptually, this makes absolute sense. If one see the work in [Lee and Seung, 1999], one will see that many of the components added are interpretations of existing physical features, such as eyebrows and mustache. It does not make sense to add a negative eyebrow to reconstruct one's face. Hence, in retrospect, it makes complete sense why this methodology was being used in chemometrics, as it does not make sense to add a negative sulfate ion to a solution.

One may ask how is this factorization different from other types of factorizations, such as principal component analysis (PCA). First, every factorization has different purposes; for instance, a Cholesky decomposition was specifically made to obtain lower computation complexity in order to invert matrices, where as PCA was made to obtain information about its principal orthogonal components. Second, when looking at the comparison between a signal reconstruction or separation from an NMF and from PCA, PCA poses a lot of problems. PCA only allows for the feature space to represent vectors that are orthogonal, which is not always the case. In the sense of speech, it is almost impossible for one speech signal to be orthogonal to another one. PCA also assumes that the probability distribution of the data follows a multivariate Gaussian, which again is not always the case. A NMF is very promising as it performs a matrix decomposition, which is very desirable for computational reconstruction in the last step of the BSS, and it is able to attain reconstructions just as detailed as some of the other reconstruction techniques, plus one is allowed to determine how many components to use.

Here, the notation follows the notations shown in [Schmidt and Olsson, 2006].

Consider the magnitude of a spectrogram of a mixed speech signal $\mathbf{Y} \in \mathbb{R}^{M \times N}$, which is a summation of R source signals (i.e. $\mathbf{Y} = \sum_{i=1}^R \mathbf{Y}_i$). The spectrogram can be sparsely represented in an overcomplete basis as

$$\mathbf{Y} = \mathbf{D}\mathbf{H}$$

where $\mathbf{D} \in \mathbb{R}^{M \times R}$ represents a compendium of dictionaries with R columns, which can be chosen by the user, and $\mathbf{H} \in \mathbb{R}^{R \times N}$ represents a code matrix, which is sparse and it contains code matrices associated with the

dictionaries in D . The matrix D and H can be represented as

$$D = \begin{bmatrix} | & | & & | \\ D_1 & D_2 & \cdots & D_R \\ | & | & & | \end{bmatrix}$$

$$H = \begin{bmatrix} --- & --- & --- & H_1 & --- & --- & --- \\ --- & --- & --- & H_1 & --- & --- & --- \\ & & & \vdots & & & \\ --- & --- & --- & H_R & --- & --- & --- \end{bmatrix}$$

A simple version of how a NMF would work is shown in Figure 1.3, where the figure below was inspired on the works from [Schmidt and Olsson, 2006].

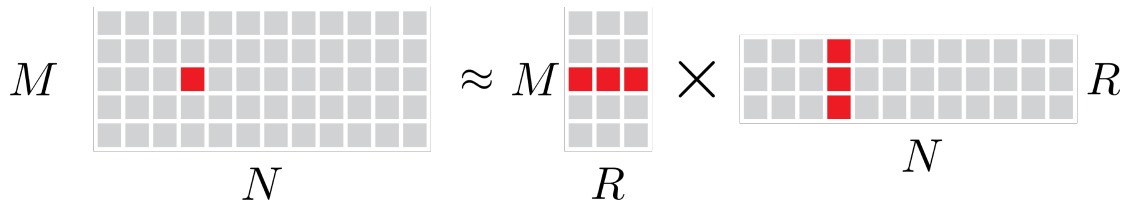


Figure 1.3: A very simple schematic of a non-negative matrix factorization, where the matrix on the left represents the mixed signal, the matrix on the center represents the dictionary matrix, and the matrix on the right represents the coding matrix

There are two different ways to obtain the desired decomposition, and that is based on what optimization approach one wishes to take. One approach is to minimize the squared Euclidean-based Frobenius norm, $\|Y - DH\|_F^2$, and another approach is to minimize the divergence, $\mathcal{D}(Y \| DH)$, both cases with respect to W and H , subject to the constraints that $D, H \geq 0$ [Lee and Seung, 2000]. Both of their definitions are shown below, where the subscript i represents the i^{th} row, and the j represents the j^{th} column. Notice that the divergence measurement below is not the Kullback-Leibler (KL) divergence, but it does converge to the KL divergence if $\sum_{ij} A_{ij} = \sum_{ij} B_{ij} = 1s$

$$\begin{cases} \text{Norm:} & \|A - B\|_F^2 = \sum_{i,j} (A_{ij} - B_{ij})^2 \\ \text{Divergence:} & \mathcal{D}(A \| B) = \sum_{i,j} \left(A_{ij} \log \left(\frac{A_{ij}}{B_{ij}} - A_{ij} + B_{ij} \right) \right) \end{cases}$$

Because this is an iterative algorithm, if one follows the norm minimization, the update rules are:

NMF Algorithm for Norm Minimization

1. Initialize $D, H \geq 0$
2. $H_{ij} \leftarrow H_{ij} \frac{(D^T Y)_{ij}}{(D^T D H)_{ij}}$
3. $D_{ij} \leftarrow D_{ij} \frac{(Y H^T)_{ij}}{(D H H^T)_{ij}}$
4. Repeat 2,3 until stability

If one follows the divergence minimization, the update rules are

NMF Algorithm for Divergence Minimization

1. Initialize $\mathbf{D}, \mathbf{H} \geq 0$
2. $\mathbf{H}_{ij} \leftarrow \mathbf{H}_{ij} \frac{\sum_k \mathbf{D}_{ki} \mathbf{Y}_{kj} / (\mathbf{D}\mathbf{H})_{kj}}{\sum_l \mathbf{D}_{li}}$
3. $\mathbf{D}_{ij} \leftarrow \mathbf{D}_{ij} \frac{\sum_k \mathbf{H}_{jk} \mathbf{Y}_{ik} / (\mathbf{D}\mathbf{H})_{ik}}{\sum_l \mathbf{H}_{jl}}$
4. Repeat 2,3 until stability

If one performs an NMF, a visual representation is shown on Figure 1.4, where the figure below was inspired on the works from [Schmidt and Olsson, 2006].

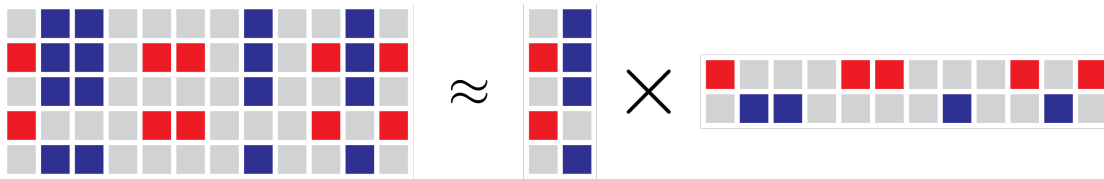


Figure 1.4: A visual example of an NMF result, with two sources being separated.

The algorithms for NMF were further improved by [Eggert and Korner, 2004], where the authors created a Sparse NMF. Their work was inspired by the basic NMF problems that, when one had overcomplete dictionaries, there was no well defined solution such that you would obtain sparse solutions. Thus, they wanted to make sure to obtain sparse solutions, especially on the coding matrix. To tune the sparsity measurement, the authors simply applied a L1 regularization based on the coding matrix coefficients on the cost functional of the norm, and called it a Sparse NMF. In other words, the cost (or energy) functional that originally was $\mathcal{E} = \|\mathbf{Y} - \bar{\mathbf{D}}\mathbf{H}\|_F^2$ became

$$\mathcal{E} = \|\mathbf{Y} - \bar{\mathbf{D}}\mathbf{H}\|_F^2 + \underbrace{\lambda \sum_{i,j} H_{ij}}_{L1 \text{ regularization}} \quad s.t. \quad \mathbf{D}, \mathbf{H} \geq 0$$

where $\bar{\mathbf{D}}$ still follows the notation of [Schmidt and Olsson, 2006], representing a column-wise normalized dictionary matrix, and λ is the parameter to control the degree of sparsity. To solve this optimization problem, if one was to set $\mathbf{R} = \mathbf{D}\mathbf{H}$, then one would obtain the following algorithm. Note, for this one case, that the dot sign \cdot and division operator $\frac{a}{b}$ represent point-wise multiplication and division

Sparse NMF Algorithm for Frobenius Norm L1 Regularized Minimization

1. Initialize $\mathbf{D}, \mathbf{H} \geq 0$
2. $\mathbf{H}_{ij} \leftarrow \mathbf{H}_{ij} \cdot \frac{\mathbf{Y}_i^T \bar{\mathbf{D}}_j}{\mathbf{R}_i^T \bar{\mathbf{D}}_j + \lambda} \in \mathbb{R}^1 \cdot \frac{\mathbb{R}^{1 \times M} \mathbb{R}^{M \times 1}}{\mathbb{R}^{1 \times M} \mathbb{R}^{M \times 1} + \mathbb{R}^1}$
3. $\mathbf{D}_{ij} \leftarrow \mathbf{D}_{ij} \cdot \frac{\sum_i \mathbf{H}_{ij} [\mathbf{Y}_i + (\mathbf{R}_i^T \bar{\mathbf{D}}_j) \bar{\mathbf{D}}_j]}{\sum_i \mathbf{H}_{ij} [\mathbf{R}_i + (\mathbf{Y}_i \bar{\mathbf{D}}_j) \bar{\mathbf{D}}_j]} \in \mathbb{R}^{M \times 1} \cdot \frac{\sum_i \mathbb{R}^1 + (\mathbb{R}^{1 \times M} \mathbb{R}^{M \times 1}) \mathbb{R}^{M \times 1}}{\sum_i \mathbb{R}^1 [\mathbb{R}^{M \times 1} + (\mathbb{R}^{1 \times M} \mathbb{R}^{M \times 1}) \mathbb{R}^{M \times 1}]}$
4. Repeat 2,3 until stability

All of the algorithms shown above are the steps necessary to train the algorithm, where \mathbf{Y} contains all of the training data. In order to actually perform a speech reconstruction, one should simply perform the algorithms above, but to keep the dictionary matrix \mathbf{D} fixed and update only the code matrix \mathbf{H} . The speech is then separated by computing the reconstruction of parts of the sparse decomposition. In other words:

1. Apply NMF or Sparse NMF to learn the dictionaries of individual speakers

2. To separate the mixtures, keep \mathbf{D} fixed, and only update \mathbf{H}
3. Reconstruct the signal by using the desired parts of the decomposition

Approaches to Learn Dictionaries

As shown in [Schmidt and Olsson, 2006], there are multiple approaches to perform BSS with NMF or SNMF: unsupervised approach and segmenting the training data and solving multiple NMF problems

In the unsupervised approach, one would compute the SNMF over a very large dataset matrix \mathbf{Y} of a single speaker to obtain a single dictionary. This, however, may take a lot of computation time. The second approach, used in [Schmidt and Olsson, 2006] involved segmenting the training data according to phoneme labels obtained by a speech recognition software with hidden-Markov models (HMM) for a single speaker, and creating a sparse dictionary for each phoneme, and finally, a final dictionary would be constructed by the concatenation of individual phoneme dictionaries. In order to create the dictionary of each phoneme, each type of phoneme would be concatenated together, to create that phoneme's dictionary, and then, the phonemes dictionaries would be concatenated together to create a dictionary that is representative of the speaker. Figure 1.5 shows an illustration on how to create the dictionary.

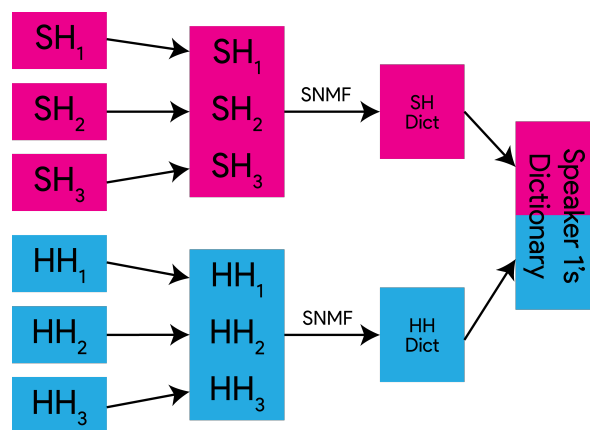


Figure 1.5: An illustration on how to create a speaker's dictionary \mathbf{D}_i

NMF for Speech Enhancement

The works from [Li and Xiao,] have shown that it is also possible to perform speech enhancement on a speech signal via non-negative matrix factorization. An illustration on Figure 1.6, inspired by the images on [Li and Xiao,], shows the procedure. It is possible to generate the dictionary matrix \mathbf{D} and the coding matrix \mathbf{H} from a noisy signal, however one would not know which of the subdictionaries \mathbf{D}_i or the code submatrices \mathbf{H}_i represent the ones associated with the speech and with the noise. Through spectral clustering (for technical details, see A.3), it is possible to group the subparts of the dictionary matrix and the code matrix so that it is differentiated between speech and noise. Given that the reconstruction process is based on linear algebra rules, it is then possible to simply multiply the dictionary matrix related to speech \mathbf{D}_{speech} and the code matrix related to speech \mathbf{H}_{speech} to obtain the cleaned speech.

Of course this solution is not guaranteed to remove every single noise source due to the possible stochasticity that is propagated to the eigenvectors of the Laplacian, and the fact that the solution is not unique, which is based on the initialization of the centroids of the k-means algorithm, but it could be a good fundamental method to implement.

1.6 Speech Recognition

1.6.1 Hidden-Markov Models and Gaussian Mixture Models

A well documented work that I've found on Gaussian mixture models was [Figueiredo and Jain, 2002], though works on the Expectation-Maximization (EM) algorithm have been shown as early as in [Duda and Hart, 1973].

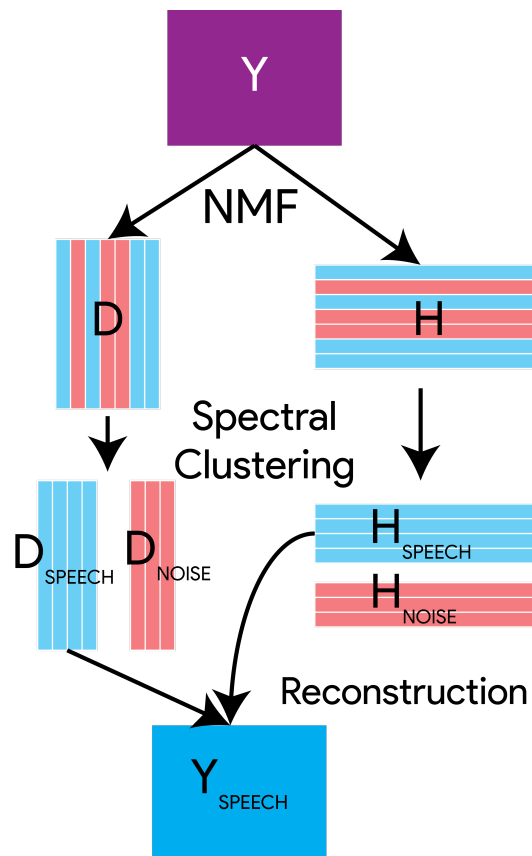


Figure 1.6: An illustration on how to use non-negative matrix factorization to denoise a signal, where Y represents the noisy signal, D represents the dictionary matrix obtained through NMF, and H represents the code matrix, also attained through the NMF.

Such models are often used for unsupervised clustering classification problems, where one assumes the probability distribution of different classes to follow a Gaussian. As mentioned before, it is critical that the data has compactness in order to make the assumption that the data distributions follow Gaussians, and if not, spectral clustering (see A.3) may be used for clustering problems. Nonetheless, for speech, depending on the features that one may be using, it is very safe to assume that they follow Gaussian distributions, which has led a lot of work in speech to be dependent on derivatives of Gaussian mixture models. For more detailed information on GMMs, please see A.2

One of the most important works for speech are the hidden Markov models (HMM), introduced by [Rabiner and Juang, 1986] (which happens to be yet one of the best tutorials on HMMs), and shown to work greatly on speech [Rabiner, 1989]. The Viterbi algorithm, which was first proposed in [VITERBI, 1967], and later formalized in [Forney, 1973], is a great way to solve HMM problems, via dynamic programming. Many consider HMMs to be dynamic forms of GMMs, where a single point in time of a HMMs is essentially a GMM (assuming a Gaussian distribution between the classifications), but this fails if the prior probabilities are not Gaussian. Nonetheless, it became very common to use HMMs for speech recognition. For instance, the work from [Schmidt and Olsson, 2006] could not be completed if the authors did not have an HMM module to output phonemes to segment their signals in order to create dictionaries for SNMF.

As an example, Figure 1.7 shows an example of a HMM, where the states are the phonemes of the word "shoot". As it can be seen, there are multiple entry points to the model, where they can start at SH, UW or T, and as time progresses, which is demonstrated as the repetition of the three states to the right, it is possible that the SH will progress to SH, or any of the other phonemes, and this propagates throughout the entire signal. A set of features is computed over windows of a signal, and a set of features are extracted at each window. The Baum-Welch algorithm is then used to compute the statistics (i.e. "train") for the nodes and to compute the transition probabilities. Finally, when testing, the windows are classified over which states does it most likely

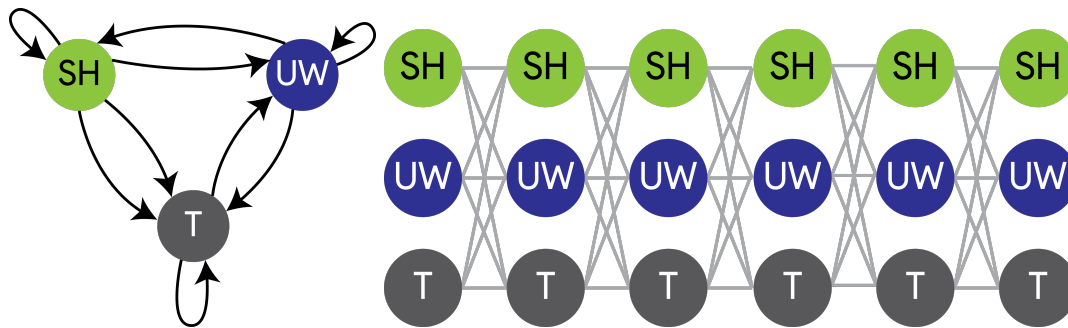


Figure 1.7: A simple HMM with the states being the phonemes (shown as ARPABET). The image on the left shows the HMM in its compact form, and the image on the right shows the HMM unrolled over time.

belong to, such as with a GMM. The Viterbi-Trellis algorithm is then used to compute the optimal path over the trellis, thus defining all of the states over time. Because the algorithm is computed over a successive multiplication of the probabilities and transition probabilities, it is very likely that the numbers quickly decrease. Therefore, it is optimal to utilize log computations, especially for computers to be able to maintain float accuracies. Once done so, an optimal path is computed, which is depicted on Figure 1.8 as a black line.

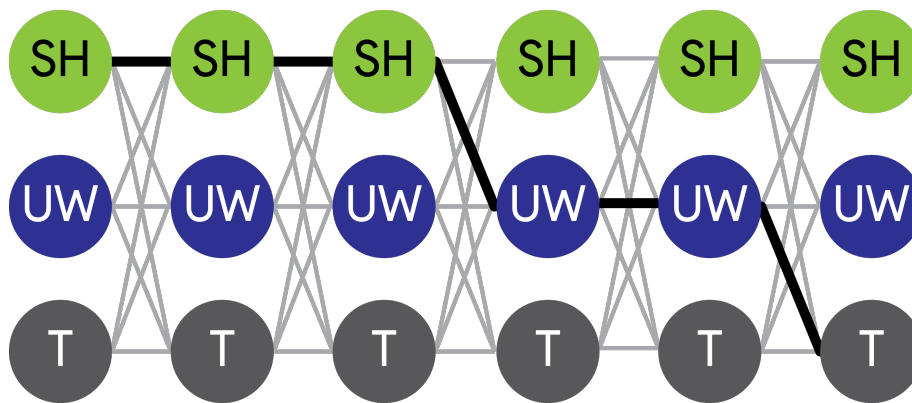


Figure 1.8: The theoretical Viterbi algorithm solution of a signal.

HMMs are often formulated where states have substates. Building on the example above, the phonemes could be broken down to substates, which is shown on Figure 1.9. On the illustration, dashed lines show the intra-state transition within each phonemes, and the solid lines represent the transition from one phoneme to another. Notice how, in a phoneme X , X_1 can only transition to itself or to a X_2 . Likewise, X_2 can only transition to itself or to X_3 , and X_3 can only transition to itself or to a new initial phoneme state Y_1 .

Because of the powerful formulation which allows classifications to continuously repeat themselves over the analysis of an entire frame, HMM are one the more powerful techniques for time-series analysis or dynamic datastreams, while also having the ability to perform encoding over a signal.

1.6.2 Deep Learning Architectures

Another methodology that some consider it to be "state-of-the-art" is with deep learning, by using recurrent neural networks, which were previously described in 1.4.1, there have been promising results using CNNs [Zhang et al., 2016], as well as RNNs [Graves et al., 2013] [Graves and Jaitly, 2014], and LSTMs [Han et al., 2018]. The logic behind utilizing CNNs is the fact that a neuron would be a windowing function and classifier, which captures the content/interpretation of a signal as it slides over time, and thus it activates its successive neurons. However, because the authors provide no incentive of understanding the nature of their experiments, nor do they even include a discussion section over their work, and when they do, they simply state already known concepts in

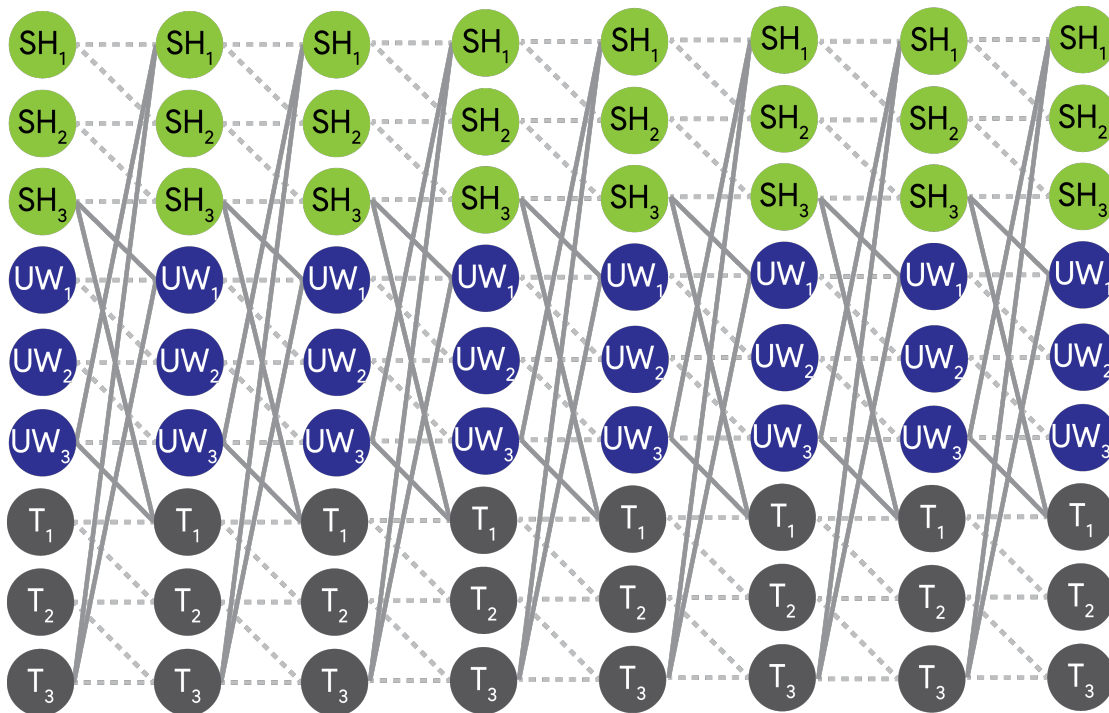


Figure 1.9: An extended version of the HMM for three phonemes, in which each phoneme contains three substates.

the field or regurgitate the aspects of their architecture, concluding the paper by stating that their future works will be to test on a larger dataset or try different architectures, I will move on with our discussion to more fruitful topics which contain explicit discussions with actual scientific approaches to solving problems.

1.7 Speaker Verification

When using an ASR system in a medical setting, privacy and accuracy is paramount. In other words, it is extremely important that medical records do not get leaked for obvious reasons, and to make sure that not anyone is logging information into a patient's medical records is also of extreme importance. The concept of a speaker verification system becomes an attractive mechanism for an ASR, such that it ensures that important information is coming specifically from the physician, and not anyone else.

It is important to state that, in contrary to popular belief, speaker verification is not the same as speaker diarization (sometimes spelled as *diarisation*). Speaker verification is a methodology used to determine whether the speaker is the target speaker or not, and this technology was made for security purposes, which in the end of the pipeline, it terminates with a yes/no answer. Speaker diarization, on the other hand, is used to answer the question: "Who spoke when?"

One technique that has been considered state-of-the-art in the past few years is the i-vector methodology. In order to understand i-vectors, it is important to first understand Joint Factor Analysis (JFA), since i-vectors were based on the concepts from JFA. Prior to going into JFA, it is important to understand the difference between two types of session variations:

- **Inter-Speaker Variation:** This is a variation originated from two utterances from different speakers.
- **Inter-Session Variation:** This is a variation originated from utterances from the same speaker. This may be due to:
 - *Channel effects:* when utterances are recorded from different channels (e.g. microphones, environment)
 - *Intra-Speaker Variation:* when utterances vary due to the speaker's health or emotional state

1.7.1 Gaussian Mixture Model-Universal Background Model and Supervectors

The Gaussian Mixture Model - Universal Background Model (GMM-UBM) approach

One great advantage of using the UBM model is that it has a very nice performance, even if the speaker-dependent data is small; however, the disadvantage is that it requires a gender-balanced large speakers set in order to train it efficiently [Tiwari and Lotia, 2012]

1.7.2 Joint Factor Analysis

1.7.3 *i*-vectors

1.8 Speaker Diarization

1.9 Features from Speech

Appendices

Appendix A

Methodologies

A.1 Long Short-Term Memory and Gated Recurrent Unit

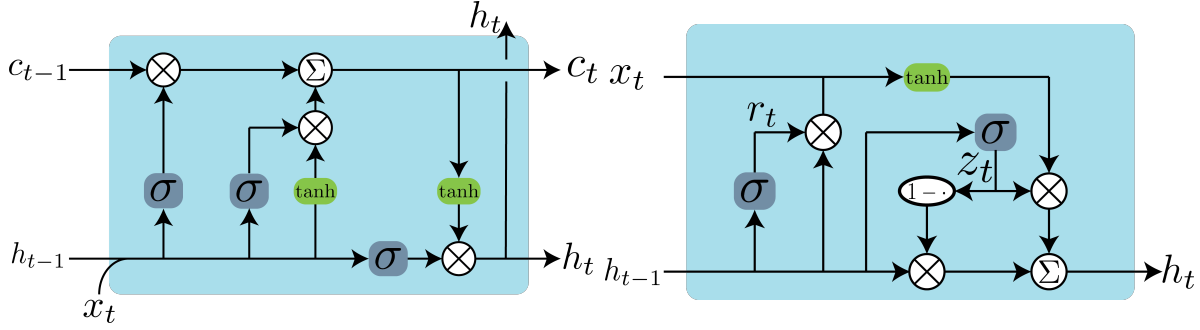


Figure A.1: A visual representation of the LSTM cell and the GRU cell

The LSTM and the GRU cell units are shown in Figure A.1, where σ depicts the sigmoid function. The equations that govern the LSTM are

$$\begin{cases} \begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \\ c_t f \cdot c_{t-1} + i \cdot g \\ h_t = o \cdot \tanh(c_t) \end{cases}$$

The equations that govern the GRU are

$$\begin{cases} \text{Update Gate:} & z_t = \sigma(W_z[h_{t-1}, x_t]) \\ \text{Reset Gate:} & r_t = \sigma(W_r[h_{t-1}, x_t]) \\ \text{Current Memory Content:} & \tilde{h}_t = \tanh(W \cdot [r_t \cdot h_{t-1}, x_t]) \\ \text{Final Memory:} & h_t = (1 - z_t) \cdot h_{t-1} + z_t \cdot \tilde{h}_t \end{cases}$$

Although sigmoid functions are often frowned upon in the deep learning community due to the inevitable fact that often times, the gradients become zero at very high or low values, the sigmoids behave as switch knobs for the acceptance or rejection of previous information on the cell memory c_t or on the hidden states h_t , which are coming from previous iterations. The vanishing gradient problem is avoided in LSTM and GRU or any network with forget gates through the concept of linear carousel.

A.2 Gaussian Mixture Models

A.3 Spectral Clustering

A.4 Mel-Frequency Cepstrum Coefficients

A.5 Linear Prediction Coefficients and its Derivatives

Bibliography

- [Anttila et al., 1995] Anttila, P., Paatero, P., Tapper, U., and Järvinen, O. (1995). Source identification of bulk wet deposition in finland by positive matrix factorization. *Atmospheric Environment*, 29(14):1705 – 1718.
- [Barak, 1988] Barak, A. P. (1988). Learning state space trajectories in recurrent neural networks. *Neural Computation*, 1:263–269.
- [Cleeremans et al., 1989] Cleeremans, A., Servan-Schreiber, D., and McClelland, J. L. (1989). Finite-state automata and simple recurrent networks. *Neural Computation*, 1(3):372–381.
- [Duda and Hart, 1973] Duda, R. O. and Hart, P. E. (1973). *Pattern Classification and Scene Analysis*. John Wiley & Sons, New York.
- [Eggert and Korner, 2004] Eggert, J. and Korner, E. (2004). Sparse coding and nmf. In *2004 IEEE International Joint Conference on Neural Networks (IEEE Cat. No.04CH37541)*, volume 4, pages 2529–2533 vol.4.
- [Figueiredo and Jain, 2002] Figueiredo, M. A. T. and Jain, A. K. (2002). Unsupervised learning of finite mixture models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(3):381–396.
- [Forney, 1973] Forney, G. D. (1973). The viterbi algorithm. *Proceedings of the IEEE*, 61(3):268–278.
- [Fowler, 2018] Fowler, G. (2018). I live with alexa, google assistant and siri. here’s which one you should pick. *The Washington Post: The Switch Review*.
- [Gers et al., 2000] Gers, F. A., Schmidhuber, J., and Cummins, F. A. (2000). Learning to forget: Continual prediction with lstm. *Neural Computation*, 12:2451–2471.
- [Graves and Jaitly, 2014] Graves, A. and Jaitly, N. (2014). Towards end-to-end speech recognition with recurrent neural networks. In *Proceedings of the 31st International Conference on International Conference on Machine Learning - Volume 32, ICML’14*, pages II–1764–II–1772. JMLR.org.
- [Graves et al., 2013] Graves, A., Mohamed, A., and Hinton, G. E. (2013). Speech recognition with deep recurrent neural networks. *CoRR*, abs/1303.5778.
- [Han et al., 2018] Han, K. J., Chandrashekar, A., Kim, J., and Lane, I. R. (2018). The CAPIO 2017 conversational speech recognition system. *CoRR*, abs/1801.00059.
- [Hochreiter and Schmidhuber, 1997] Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8):1735–1780.
- [Jean-Marc et al., 2017] Jean-Marc, Xiph.Org, and Mozilla (2017). Rnoise.
- [Jordan, 1986] Jordan, M. I. (1986). Attractor dynamics and parallelism in a connectionist sequential machine. *Cognitive Science Conference*, pages 531–546.
- [Lee and Seung, 1999] Lee, D. D. and Seung, H. S. (1999). Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791.
- [Lee and Seung, 2000] Lee, D. D. and Seung, H. S. (2000). Algorithms for non-negative matrix factorization. In *Proceedings of the 13th International Conference on Neural Information Processing Systems, NIPS’00*, pages 535–541, Cambridge, MA, USA. MIT Press.

- [Li and Xiao,] Li, P. and Xiao, Y. Matrix factorization for speech enhancement.
- [Moskvitch, 2017] Moskvitch, K. (2017). The machines that learned to listen. *BBC Future*.
- [Ng et al., 2001] Ng, A. Y., Jordan, M. I., and Weiss, Y. (2001). On spectral clustering: Analysis and an algorithm. In *ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS*, pages 849–856. MIT Press.
- [Paatero and Tapper, 1994] Paatero, P. and Tapper, U. (1994). Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values. *Environmetrics*, 5(2):111–126.
- [Rabiner, 1989] Rabiner, L. R. (1989). A tutorial on hidden markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286.
- [Rabiner and Juang, 1986] Rabiner, L. R. and Juang, B.-H. (1986). An introduction to hidden markov models. *IEEE ASSP Magazine*, 3:4–16.
- [Schmidt and Olsson, 2006] Schmidt, M. N. and Olsson, R. K. (2006). Single-channel speech separation using sparse non-negative matrix factorization. In *INTERSPEECH*.
- [Tiwari and Lotia, 2012] Tiwari, M. S. and Lotia, P. (2012). Speaker verification system using gaussian mixture model ubm. *International Journal of Digital Application Contemporary Research*, 1(3).
- [VITERBI, 1967] VITERBI, A. J. (1967). Error bounds for convolutional codes and an asymptotically optimal decoding algorithm. *IEEE Transactions on Information Theory*, 13:260–269.
- [Zhang et al., 2016] Zhang, Y., Pezeshki, M., Brakel, P., Zhang, S., Laurent, C., Bengio, Y., and Courville, A. (2016). Towards end-to-end speech recognition with deep convolutional neural networks. In *Interspeech 2016*, pages 410–414.