

Contents

Contents	ii
List of Figures	iv
1 Mathematical Background	1
1.1 Introduction	1
1.2 Metric Spaces	1
1.3 Convergence and completeness	3
1.4 Vector spaces and Banach spaces	5
1.5 Inner product spaces and Hilbert spaces	6
1.6 Manifolds	9
1.7 Summary	9
2 Automatic Speech Recognition	11
2.1 Introduction	11
2.2 Fundamentals of the speech production mechanism	11
2.3 Design parameters of an ASR system	13
2.4 ASR system architecture	13
2.4.1 Computational formulation of ASR	15
2.4.2 Feature Extraction	16
2.4.3 The Language Model	23
2.4.4 The Acoustic Model	27
2.4.5 The Decoder	39
2.4.6 Evaluation	44
3 Manifold and Representation learning	45
3.1 Manifold learning	45
3.1.1 Mathematical formulation of manifold learning	46
3.2 Manifold learning algorithms	48
3.2.1 Dimensionality reduction	48

3.2.2	Algorithms	49
3.3	Discriminative manifold learning algorithms in ASR	55
3.3.1	Manifolds and speech data	55
3.3.2	Discriminative manifold learning	58
4	Deep Neural Networks	65
4.1	History and Motivation	65
4.2	Constructing DNNs: the Multilayer Perceptron	66
4.2.1	Architecture of the Deep Multilayer Perceptron	66
4.2.2	Training and related issues	69
4.3	Pretraining	77
4.4	Recurrent Neural Networks	86
4.4.1	State-Space formulation of the basic RNN	87
4.4.2	Training	88
4.4.3	RNNs with Long-Short-Term Memory cells	91

List of Figures

1.1	<i>The unit balls in the euclidean space defined using the euclidean norm (B_2), the sum norm (B_1) and the maximum norm (B_{\max}).[Kre78]</i>	3
1.2	<i>Illustration of a mapping T. [Kre78]</i>	4
1.3	<i>A one-dimensional manifold embedded in the three-dimensional space. [Cay05]</i>	10
1.4	<i>Overview of various spaces and their relations.(Image taken from:https://en.wiki2.org/wiki/mathematics)</i>	10
2.1	<i>The vocal organs.[JM09]</i>	12
2.2	<i>The noisy channel model. The decoder searches through all possible sentences and finds the one that after passing through the channel best matches the initial input signal[JM09].</i>	14
2.3	<i>The main building blocks of an ASR system[You02].</i>	14
2.4	<i>MFCC feature extraction from a quantized digital waveform[JM09].</i>	16
2.5	<i>Part of the spectrum of vowel [aa] before (a) and after (b) pre-emphasis[JM09].</i>	17
2.6	<i>Rectangular and Hamming windows and their effect on the signal[JM09].</i>	18
2.7	<i>Windowing process with a rectangular window.After a figure by Brian Pellom[JM09].</i>	18
2.8	<i>The mel filterbank after Davis and Mermelstein (1980)[JM09].</i>	19
2.9	<i>Effect of vocal tract on source signal.Tomi H. Kinnunen, Speech Technology Workshop</i>	20
2.10	<i>Cepstrum computation.Source: test.virtual-labs.ac.in</i>	20
2.11	<i>Cepstrum example: (a) magnitude spectrum, (b) log magnitude spectrum, (c) cepstrum.[JM09].</i>	20
2.12	<i>Feature extraction process.[You02].</i>	22
2.13	<i>Process of building a language model.[YEK⁺02].</i>	24
2.14	<i>Generate observation sequence from an HMM .[YD14]</i>	34
2.15	<i>The Forward algorithm for HMM probability evaluation.[JM09]</i>	36
2.16	<i>The Viterbi algorithm for HMM decoding.[YD14]</i>	36
2.17	<i>The Backward algorithm for HMM decoding.[DHS00]</i>	37

2.18	Three-state, left-to-right HMM model.[ST04]	39
2.19	Composition algorithm .[MPR01].	40
2.20	Determinization algorithm .[MPR01].	41
2.21	Transducer example: (a) Grammar G , (b) Lexicon \bar{L} [MPR01].	43
2.22	Transducer example: (c) $\bar{L} \circ G$, (d) $\bar{L} \circ G$ determinized, (e) $\min_{\text{tropical sem}} \det(\bar{L} \circ G)$ [MPR01].	44
3.1	Tangent plane and directions of variation on the manifold [BGC15].	47
3.2	Swiss roll (top) and the underlying manifold (bottom). [Cay05].	48
3.3	Tangent planes are tiled together to cover the manifold, forming a global coordinate system [BGC15].	50
3.4	Classical Multidimensional Scaling [Cay05].	51
3.5	ISOMAP [Cay05].	52
3.6	LLE [Cay05].	54
3.7	Laplacian Eigenmaps [Cay05].	55
4.1	A deep neural network with three hidden layers [YD14].	67
4.2	Computation of DNN output [YD14].	69
4.3	Back-propagation algorithm [YD14].	70
4.4	An example of a Restricted Boltzmann Machine [YD14].	79
4.5	Marginal probability density function of a Gaussian-Bernoulli RBM [YD14].	81
4.6	Contrastive Divergence algorithm for RBM training [YD14].	82
4.7	Different perspectives of the same RBM [YD14].	83
4.8	Discriminative pretraining [YD14].	85
4.9	Backpropagation-through-time [YD14].	89

Chapter 1

Mathematical Background

1.1 Introduction

An important part of every pattern recognition project handling large amounts of data, is the dimensionality reduction problem, that is, the transfer of the input data into a lower dimensional space, which will allow for more efficient processing and lower model error, as fewer parameters will have to be estimated. In order to develop algorithms for this task, one must be able to understand basic notions often arising in such context, e.g. Hilbert space, convergence, metric, manifolds. The purpose of this chapter is to offer a brief introduction to these concepts in connection with the broader topic of the project([Kre78]).

1.2 Metric Spaces

A *space* in the broader mathematical sense, is a set of elements X with some added structure.

Specifically, we will be working on *metric spaces*. Prior to defining a *metric space* we must first define what a *metric* is.

Definition 1. A *metric* d on a set X (or *distance function* on X) is a function defined on $X \times X$, such that for all $x, y, z \in X$, the following properties hold true:

- d is real-valued, finite and non-negative
- $d(x, y) = 0$ iff $x = y$
- $d(x, y) = d(y, x)$ (*Symmetry*)
- $d(x, y) \leq d(x, z) + d(z, y)$ (*Triangle Inequality*)

Definition 2. A metric space is a pair (X, d) where X is a set and d is a metric on X .

Set X is also called the *underlying set* of (X, d) . For the fixed points x, y the non-negative number $d(x, y)$ is called the distance from x to y .

Examples of metric spaces are the well known \mathbb{R} with the distance function

$$d(x, y) = |x - y|$$

, and the Euclidean space \mathbb{R}^2 with the Euclidean metric, defined by

$$d(x, y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2}$$

for elements $x = (x_1, x_2)$ and $y = (y_1, y_2)$.

Spaces and their structure play a vital role in dimensionality reduction. The goal of this process is to transfer the input space into a lower dimensional one, yet at the same time retain the original relations between input data. In order to do this, we need to make sure that the space we move to, has certain properties. We are especially interested in spaces where the triangular inequality is satisfied, so that we can take advantage of the convergence properties it offers to the space. This is the reason why we built on a metric space, adding properties and operations, to derive new spaces which offer the required structure.

We will now introduce some auxiliary concepts which will enable a smooth transition into further defining kinds of spaces.

Definition 3. Given a point $x_0 \in X$ and a real number $r > 0$, we define three types of sets:

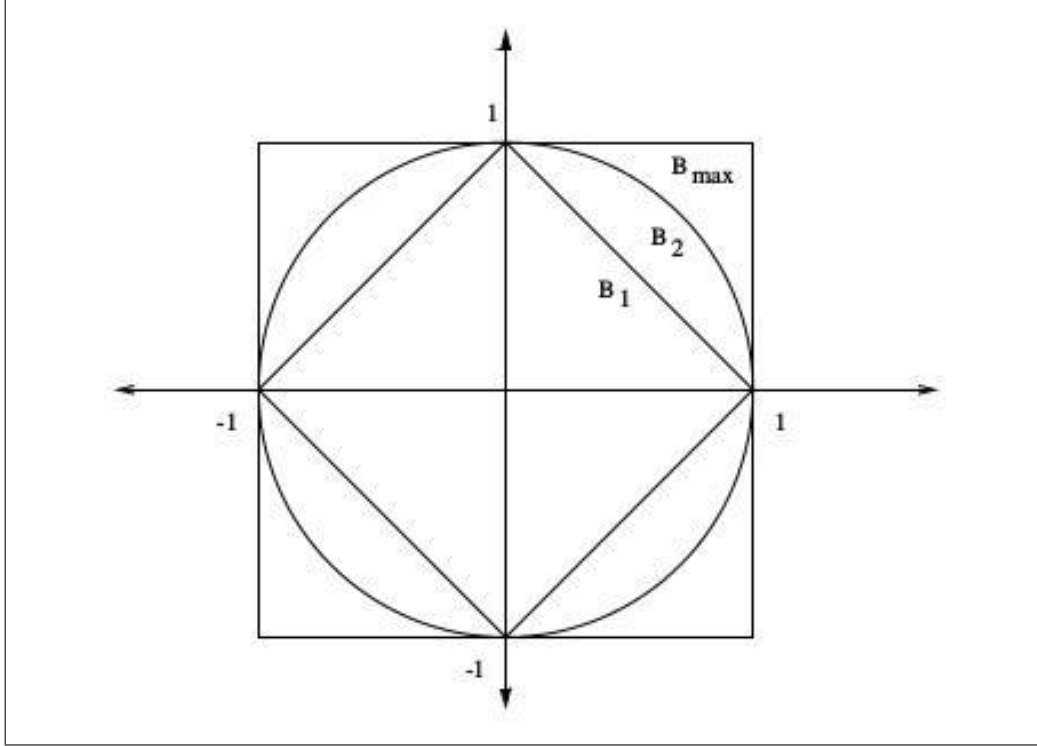
- $B(x_0; r) = \{x \in X \mid d(x, x_0) < r\}$ (Open Ball)
- $\bar{B}(x_0; r) = \{x \in X \mid d(x, x_0) \leq r\}$ (Closed Ball)
- $S(x_0; r) = \{x \in X \mid d(x, x_0) = r\}$ (Sphere)

where point x_0 is called the center and r the radius.

Given the definition of a *ball*, we can define the notion of a neighborhood of a point.

Definition 4. An open ball $B(x_0, \varepsilon)$, $\varepsilon > 0$, is called an ε -neighborhood of x_0 . A neighborhood of x_0 is any subset of X which contains an ε -neighborhood of x_0 .

Figure 1.1: The unit balls in the euclidean space defined using the euclidean norm (B_2), the sum norm (B_1) and the maximum norm (B_{\max}). [Kre78]



Since the process of dimensionality reduction is, in its essence, a mapping from one space to another, we proceed by defining the notion of a continuous mapping.

Definition 5. Let $X = (X, d)$ and $Y = (Y, \bar{d})$ be metric spaces. A mapping $T : X \rightarrow Y$ is said to be continuous at a point $x_0 \in X$ if for every $\varepsilon > 0$ there exists a $\delta > 0$ such that

$$\bar{d}(Tx, Tx_0) < \varepsilon \quad \forall x \text{ satisfying } d(x, x_0) < \delta$$

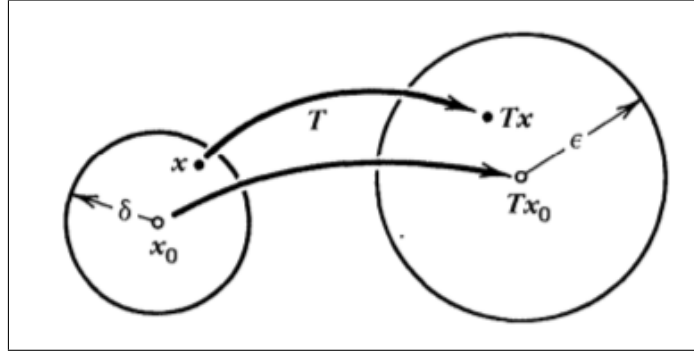
T is said to be continuous if it is continuous at every point of X .

Having presented the ideas of a metric space and some of its subsets, as well as defined a continuous mapping, we proceed to remind the notions of convergence and completeness which will be used to move on to further spaces.

1.3 Convergence and completeness

Definition 6. A sequence (x_n) in a metric space $X = (X, d)$ is said to converge or to be convergent if there is a $x \in X$ such that

$$\lim_{n \rightarrow \infty} d(x_n, x) = 0.$$

Figure 1.2: Illustration of a mapping T . [Kre78]

x is called the limit of x_n and we write

$$\lim_{n \rightarrow \infty} x_n = x \text{ or } x_n \rightarrow x.$$

It is obvious now why we need to be in a metric space to define the convergence of x_n : the metric d produces the sequence $a_n = d(x_n, x)$ whose convergence defines that of x_n . Furthermore, using the triangular inequality we can prove the following lemma:

Lemma 1. *Let $X=(X,d)$ be a metric space. Then:*

- *A convergent sequence in X is bounded and its limit is unique.*
- *If $x_n \rightarrow x$ and $y_n \rightarrow y \in X$, then $d(x_n, y_n) \rightarrow d(x, y)$.*

Convergence plays an important role in the definition of completeness which will help us later to define Banach and Hilbert spaces (where we mostly work).

Definition 7. *Let x_n be a sequence in \mathbb{R} or \mathbb{C} . x_n converges and we call it a Cauchy sequence if and only if it satisfies the Cauchy convergence criterion, that is, if and only if for every given $\varepsilon > 0$ there is a $N=N(\varepsilon)$ such that*

$$|x_m - x_n| < \varepsilon \quad \forall m, n > N.$$

However, this is the case only for \mathbb{R} and \mathbb{C} . Given that most pattern recognition tasks work on multidimensional spaces, we have to generalize this convergence property in such spaces.

Definition 8. *A sequence x_n in a metric space $X = (X, d)$ is said to be Cauchy (or fundamental), if for every $\varepsilon > 0$ there is a $N=N(\varepsilon)$ such that*

$$d(x_m, x_n) < \varepsilon \quad \forall m, n > N.$$

This generalization allows us to give the following definition of completeness:

Definition 9. A metric space $X = (X, d)$ is said to be complete if every Cauchy sequence in X converges (that is, it has a limit which is an element of X).

It is important to notice that convergence is not a property of the sequence by itself, but it also depends on the metric space that the sequence lies: the limit of the sequence must be *in the space*.

(EXAMPLES NEEDED??)

1.4 Vector spaces and Banach spaces

As we explore the structure of spaces so that we reach a suitable one for our problem, we come across the notion of *vector spaces*.

Definition 10. A vector space (or linear space) over a field K is a non-empty set X of elements x, y, \dots , which are called vectors, together with two algebraic operations : vector addition and multiplication of vectors by scalars, that is by elements of K .

Vector addition associates with every ordered pair (x, y) of vectors, a vector $x + y$, called the sum of $x + y$, in such a way that the following properties are satisfied:

$$x + y = y + x$$

$$x + (y + z) = (x + y) + z$$

Furthermore, there exists a vector θ , called the zero vector, and for every vector x there exists a vector $-x$, such that for all vectors we have:

$$x + \theta = x$$

$$x + (-x) = \theta$$

Vector multiplication by scalars associates with every vector x and scalar α a vector αx (or $x\alpha$) called the product of α and x in such a way that for all vectors x, y and scalars α, β the following hold:

$$\alpha(\beta x) = (\alpha\beta)x$$

$$1x = x$$

$$\alpha(x + y) = \alpha x + \alpha y$$

$$(\alpha + \beta)x = \alpha x + \beta x$$

A vector space X exactly as defined above, may or may not be a metric space. To make sure that a relation between the algebraic structure of X and the metric exists, and thus be able to combine algebraic and metric concepts, we have to define on X a metric d based on a *norm*.

Definition 11. A norm on a vector space X (over \mathbb{R} or \mathbb{C}) is a real-valued function on X whose value at a $x \in X$ is denoted by

$$\|x\|$$

and which has the properties

$$\|x\| \geq 0$$

$$\|x\| = 0 \Rightarrow x = 0$$

$$\|\alpha x\| = |\alpha| \|x\|$$

$$\|x + y\| \leq \|x\| + \|y\|$$

where $x, y \in X$ and $\alpha \in K$.

A norm on X defines a metric d on X which is given by

$$d(x, y) = \|x - y\| \quad x, y \in X$$

and is called the metric induced by the norm.

A vector space X with a norm defined on it, is called a normed space and is denoted by X or $(X, \|\cdot\|)$.

If a normed space is complete in the metric defined by the norm, we call it a *Banach space*.

1.5 Inner product spaces and Hilbert spaces

In a vector space, vectors can be added and multiplied giving the space its algebraic properties. By defining a norm on such a space, the concept of the length of a vector is generalized allowing us to define a metric and establish a relation between the algebraic and geometrical structure of the space.

To connect the above with pattern recognition, one should notice that when the problem is transferred to a lower dimensional space, it is often desirable to keep not only the length/magnitude relations between the data (as represented by the distance/norm for example) but also the angle between them. This reminds us of the inner product of a Euclidean space, which is what we will expand over normed spaces.

Definition 12. An inner product on a vector space X is a mapping of $X \times X$ into the scalar field K of X ; that is, with every pair of vectors x and y , there is associated a scalar which is written

$$\langle x, y \rangle$$

and is called the inner product of x and y , such that for all vectors x, y, z and scalar α the following properties hold:

$$\langle x + y, z \rangle = \langle x, z \rangle + \langle y, z \rangle$$

$$\langle \alpha x, y \rangle = \alpha \langle x, y \rangle$$

$$\langle x, y \rangle = \overline{\langle y, x \rangle}$$

$$\langle x, x \rangle \geq 0$$

$$\langle x, x \rangle \iff x = 0$$

An inner product on X defines a norm on X given by

$$\|x\| = \sqrt{\langle x, x \rangle}$$

and a metric induced by the above norm given by

$$d(x, y) = \|x - y\| = \sqrt{\langle x - y, x - y \rangle}$$

A vector space X with an inner product defined on it, is called an *inner product space* (or *pre-Hilbert space*). An inner product space which is complete in the metric defined by the inner product, is called a *Hilbert space*.

It is apparent that an inner product space is a normed space and a Hilbert space is a Banach space.

An important theorem of Banach spaces, and thus Hilbert spaces, is the *Fixed point theorem* or *Contraction theorem*. Before we present the theorem and its useful proof we have to give the following definitions.

Definition 13. A fixed point of a mapping $T: X \rightarrow X$ of a set X into itself, is a $x \in X$ which is mapped onto itself (it is kept “fixed” by T), that is,

$$Tx = x$$

the image Tx coincides with x .

Definition 14. Let $X = (X, d)$ be a metric space. A mapping $T: X \rightarrow X$ is called a contraction on X if there is a positive real number $\alpha < 1$ such that for all $x, y \in X$

$$d(Tx, Ty) \leq \alpha d(x, y)$$

that is, the images of any points x, y are closer together than the points x, y .

The Banach fixed point theorem is an existence and uniqueness theorem for fixed points of certain mappings, and its proof gives a constructive procedure for getting closer and closer to the fixed point starting from an initial approximation. Thinking this idea in connection with pattern recognition one should think of the fixed point as the pattern to be found and the initial approximation as the observation available.

Theorem 1. Banach fixed point theorem.

Consider a metric space $X = (X, d)$ where $X \neq \emptyset$. Suppose that X is complete and let $T: X \rightarrow X$ be a contraction on X . Then T has precisely one fixed point.

Proof Idea: We construct a sequence (x_n) and show that it is Cauchy so that it converges in the complete space X . Then we prove that its limit x is a fixed point of T and T has no further fixed points.

We choose any $x_0 \in X$ and define the iterative sequence (x_n) by

$$x_0, x_1 = Tx_0, x_2 = Tx_1 = T^2x_0, \dots, x^n = T^n x_0, \dots \quad (1.1)$$

which is the sequence of the images of x under repeated application of T . We now show that (x_n) is Cauchy. Since T is a contraction we have:

$$\begin{aligned} d(x_{m+1}, x_m) &= d(Tx_m, Tx_{m-1}) \\ &\leq \alpha d(x_m, x_{m-1}) = \alpha d(Tx_{m-1}, Tx_{m-2}) \\ &\leq \alpha^2 d(x_{m-1}, x_{m-2}) \\ &\dots \leq \alpha^m d(x_1, x_0). \end{aligned}$$

Hence, using the triangular inequality we obtain for $n > m$:

$$\begin{aligned} d(x_m, x_n) &\leq d(x_m, x_{m+1}) + \dots + d(x_{n-1}, x_n) \\ &\leq (\alpha^m + \alpha^{m+1} + \dots + \alpha^{n-1}) d(x_0, x_1) \\ &= \alpha^m \left(\frac{1 - \alpha^{n-m}}{1 - \alpha} \right) d(x_1, x_0) \end{aligned}$$

Since $0 < \alpha < 1$, in the numerator we have $1 - \alpha^{n-m} < 1$. Consequently,

$$d(x_m, x_n) \leq \left(\frac{\alpha^m}{1 - \alpha} \right) d(x_1, x_0)$$

On the right $0 < \alpha < 1$ and $d(x_0, x_1)$ is fixed, so that we can make the right-hand side as small as we please by taking m sufficiently large (and $n > m$). This proves that (x_m) is Cauchy. Since X is complete, (x_m) converges, say $x_m \rightarrow x$. We now have to show that this limit x is a fixed point of the mapping T .

From the triangle inequality and the definition of contraction we have:

$$\begin{aligned} d(x, Tx) &\leq d(x, x_m) + d(x_m, Tx) \\ &\leq d(x, x_m) + \alpha d(x_{m-1}, x) \end{aligned}$$

The sum on the second line can be made smaller than any preassigned $\varepsilon > 0$ because $x_m \rightarrow x$. We draw the conclusion that $d(x, Tx) = 0$ and consequently $x = Tx$, which means that x is a fixed point of T .

We now have to show that this fixed point is unique. Let there be a second fixed point \bar{x} . Then

$$d(x, \bar{x}) = d(Tx, T\bar{x}) \leq \alpha d(x, \bar{x})$$

which implies that $d(x, \bar{x}) = 0$ since $\alpha < 1$. Hence $x = \bar{x}$ and the proof of the theorem ends here.

1.6 Manifolds

An important notion that we will come across very often in the present work is the notion of a *manifold*.

An intuitive representation of a manifold is the following: suppose we have a set of data vectors in \mathbb{R}^n , and that a subset of their dimensions represents certain features of interest. If we focus on these dimensions only, “keeping” the rest “steady”, we notice that the data vectors move along a certain path in \mathbb{R}^n . This path, or curve formed in the space, is the *manifold* that the features lie on in \mathbb{R}^n .

For example, in figure 1.3 we see a curve in \mathcal{R}^3 , which has zero volume and zero area. It can therefore be parameterized by a single variable. Consequently, despite being in the three-dimensional space, the curve has an intrinsic dimensionality of one, since it locally resembles \mathcal{R}^1 . A manifold has a *locally* Euclidean geometry, in a neighborhood around each point, but on a global scale the relations between its elements are non metric.

We will give a formal definition of manifolds, when we talk about their use in representation learning.

1.7 Summary

In this chapter some basic ideas about spaces and their properties were presented. Starting from the metric and metric space, we moved on to the norm and vector space and based on these we defined the Hilbert space which is the closest generalization of the well-known Euclidean space. Hilbert spaces provide us with the tools we need in our work: convergence, completeness, contractions etc.

Figure 1.3: A one-dimensional manifold embedded in the three-dimensional space. [Cay05]

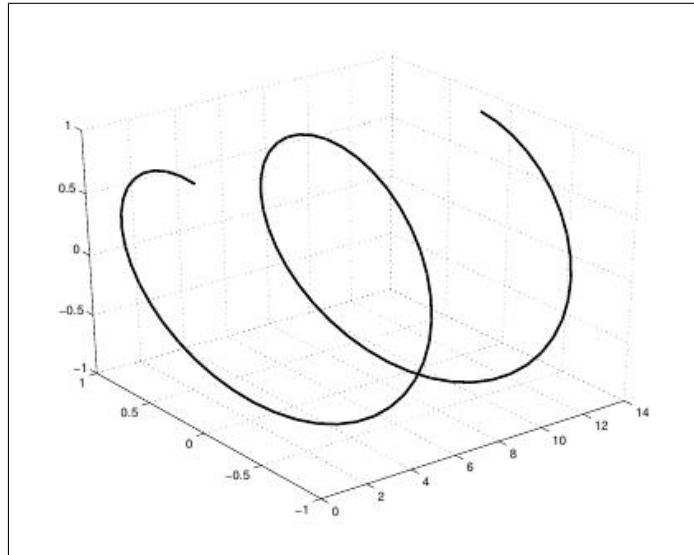
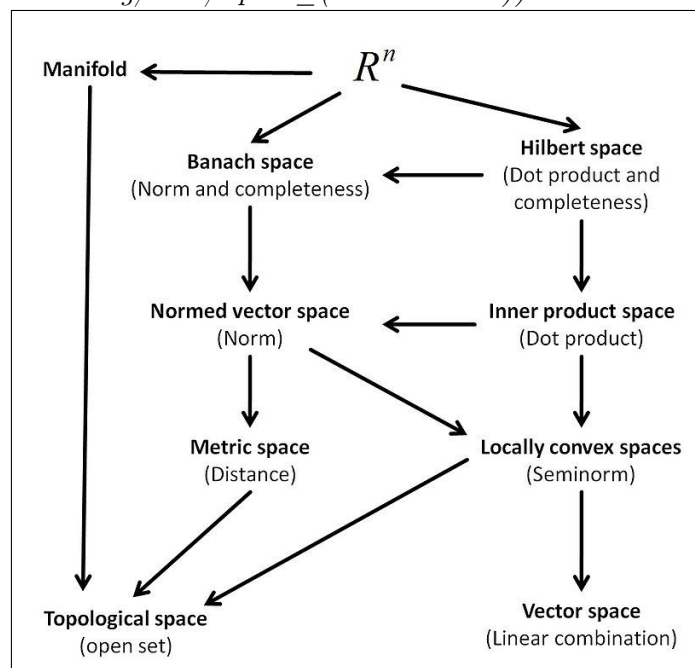


Figure 1.4: Overview of various spaces and their relations. (Image taken from: [https://en.wiki2.org/wiki/Space_\(mathematics\)](https://en.wiki2.org/wiki/Space_(mathematics)))



Chapter 2

Automatic Speech Recognition

2.1 Introduction

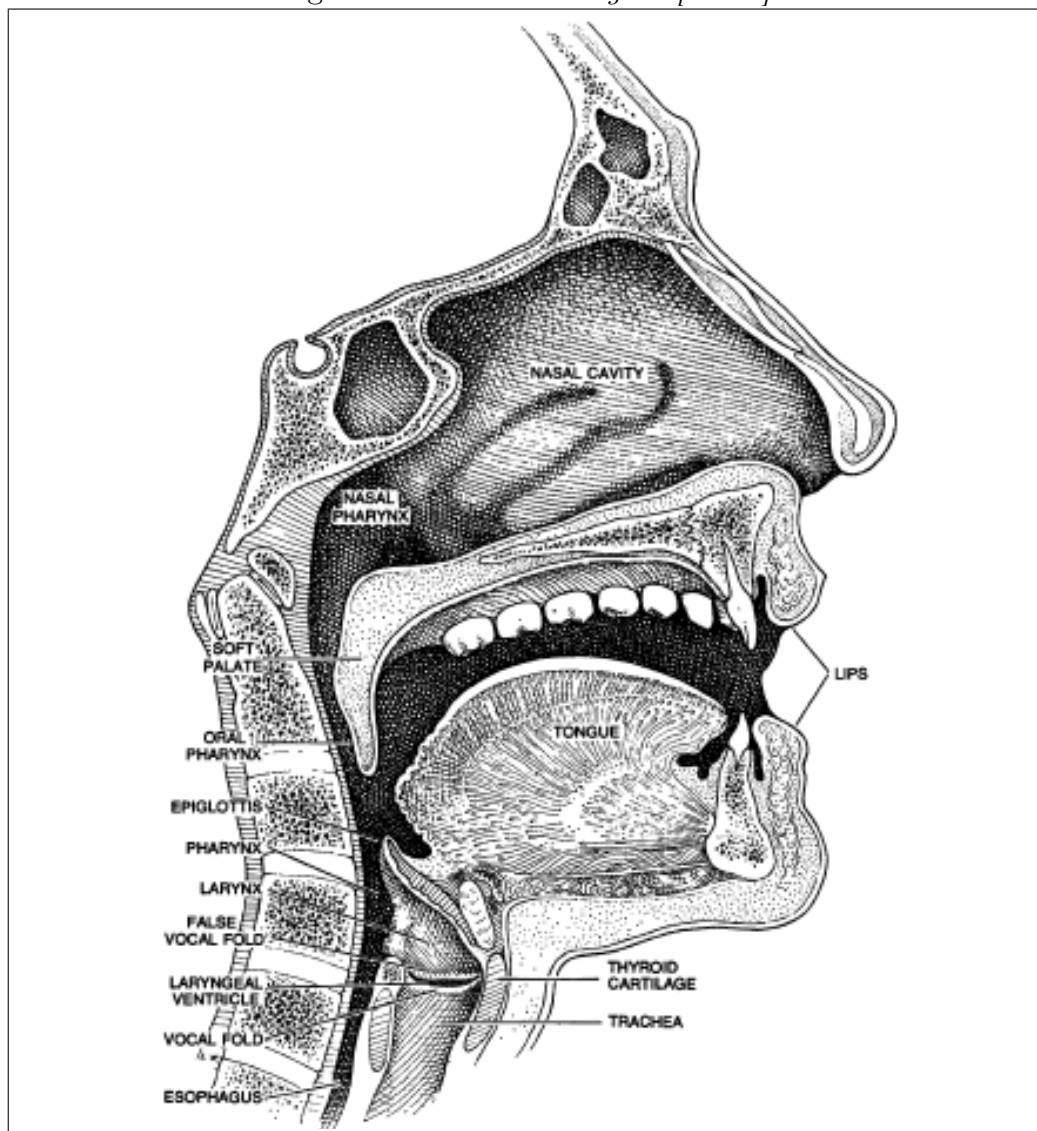
The aim of the work on Automatic Speech Recognition (*ASR*) is to build systems that recognize spoken speech, that is, they are able to map acoustic signals to strings of words. In contrast to natural language processing, ASR does not try to determine the meaning behind speech or find the multiple meanings of words; it merely tries to recognize *which* words were spoken.

Research in ASR has come a long way in the last few years, which has allowed us to take advantage of it in multiple areas with very satisfying results: human-computer interaction (speech-only or multimodal interfaces), telephony (information passing, call routing) and dictation are examples where ASR systems can perform well irrespective of the speaker or their environment.

We will begin by providing a quick description of the way human speech is produced and move on to present the concept and mechanisms behind Automatic Speech Recognition.

2.2 Fundamentals of the speech production mechanism

[JM09] Sound is produced by the rapid movement of air from the lungs through the “windpipe”, also called the *trachea* and out of the mouth or nose. As it flows through the trachea, it passes through the *larynx* (or *Adam’s apple*) and there it affects the position of two small folds of muscle which are known as *vocal folds* or *vocal chords*. The possible movements of these muscles are either moving closer together or apart; the space between them is called the *glottis*. If the vocal folds are close together they will vibrate as air passes through the glottis and produce sounds

Figure 2.1: *The vocal organs.*[JM09]

that are known as *voiced* sounds. The vowels are examples of voiced sounds. On the other hand, if they are far apart they will not vibrate and the sounds produced are called *unvoiced*. Examples of unvoiced sounds are [p], [t] and [k].

After passing through the trachea and before exiting the body, air passes through the vocal tract, which consists of the *nasal* and *oral* tract. The vocal tract will act as a filter on the speech signal and the output of the filter will be the sound we will produce. The speech signal will vary according to what obstacles the air meets on its way out : the tongue, the lips or the teeth. These obstacles will define the vocal tract filter applied on the speech signal.

2.3 Design parameters of an ASR system

This section attempts to provide a short introduction into Automatic Speech Recognition and the most common architectures behind ASR systems[JM09].

However, before deciding on the architecture of a speech recognition system, one has to take into account its application domain. The most decisive parameters are:

- The vocabulary size, which is the number of distinct words that the system should be able to recognize. Few words imply relatively easy set-up and training of the system, whereas systems recognizing thousands of words, as in a broadcasting news vocabulary, are more complicated and harder to train.
- The fluency of the speech that the system will be asked to recognize. This includes whether the speech will be continuous or just isolated words as well as the speed and clarity of the speaker. Isolated word recognition systems, e.g. recognizing commands to a computer, are easier than ASR systems for continuous, conversational speech, e.g. a telephone conversation between humans.
- The environment where the recognition might have to take place. Systems designed to perform well in noisy environments with high distortion in the speech recorded are more demanding than systems that can clearly capture the speech for recognition in an isolated environment.
- The speaker variability. Speech recognition is easier if the system is expected to recognize the speech of a limited number of people. On the contrary, a general ASR system that should work with any speaker, regardless of sex, age, or accent is much more difficult to implement.

2.4 ASR system architecture

As mentioned in the introduction, the problem of ASR is, in principle, a structured sequence classification task, where a (relatively long) sequence of acoustic data is used to infer a (relatively short) sequence of the linguistic units such as words. Modern ASR systems use the model of a *noisy channel* to deal with the classification task. The idea behind this model is to think of the input signal as a distorted version of the corresponding words, which was produced as they passed through a noisy communications channel. If we manage to understand how the channel affected the signal, we can then match it to the original, noise-free set of words, by passing every acceptable -by the grammar of the language- sentence through the channel to get its distorted version and see which matches best the initial input signal.

Figure 2.2: *The noisy channel model. The decoder searches through all possible sentences and finds the one that after passing through the channel best matches the initial input signal*[JM09].

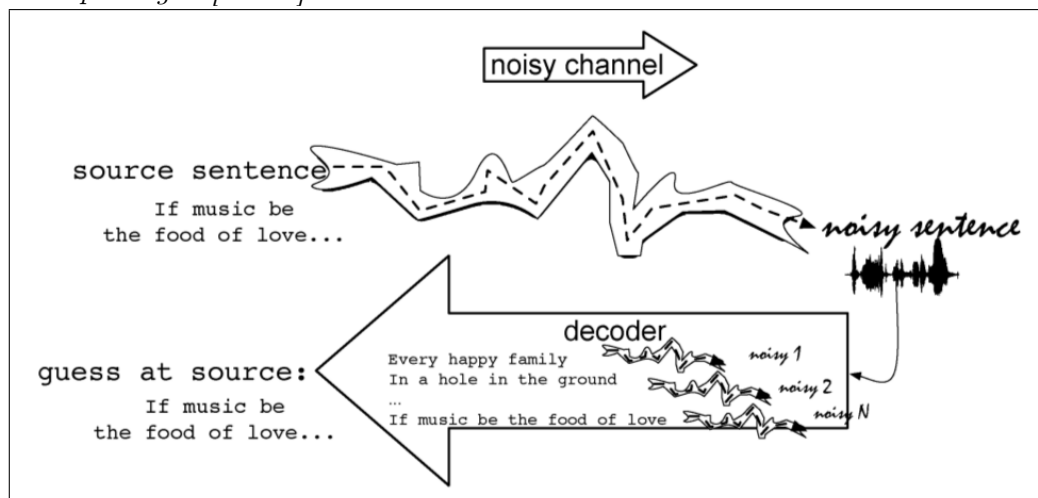
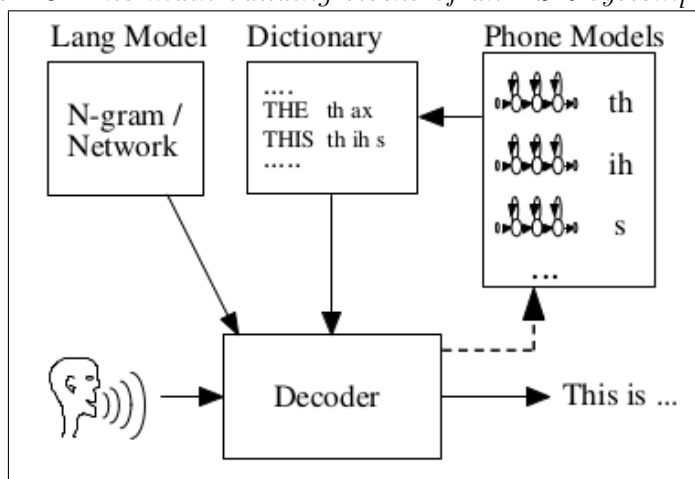


Figure 2.3: *The main building blocks of an ASR system*[You02].



In order to achieve the modelling of the noisy channel and the subsequent decoding of a new acoustic signal, we will need to have at our disposal the following tools: the prior probability of each sentence of the language, the probability of words being the concatenation of certain speech units and the probability of these speech units being realised as acoustic or spectral features, which are drawn from the input signal.

These tools define the main components of a modern automatic speech recognition system, which will be presented below, following the formulation of the computational/mathematical problem of Automatic Speech Recognition.

2.4.1 Computational formulation of ASR

We will now use mathematical notation and probability theory to answer the basic question in ASR: “What is the most likely sentence \hat{W} out of all sentences belonging to language \mathbb{L} given an acoustic signal O ?”.

The acoustic input signal O is a sequence of observations o_i ,

$$O = o_1, o_2, o_3, \dots o_t$$

each one representing features of a specific part of the input speech, which is usually split into overlapping parts of a certain duration (*frames*). In the same way, if we treat each sentence of the language as a string of words,

$$W = w_1, w_2, w_3, \dots w_n$$

we can write the answer to the ASR question in the following way:

$$\hat{W} = \operatorname{argmax}_{W \in \mathbb{L}} P(W|O)$$

, which we cannot compute directly. However, if we apply Bayes’ rule, the equation takes the following form:

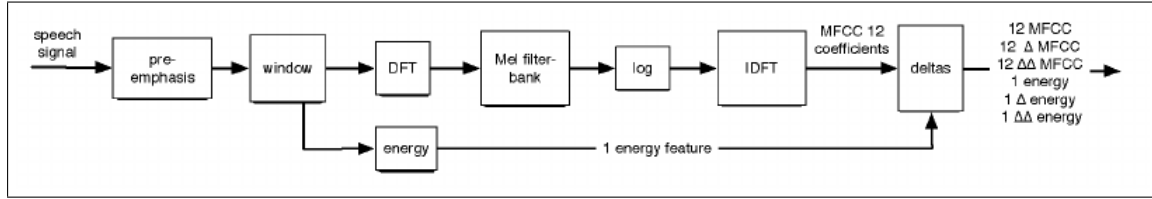
$$\hat{W} = \operatorname{argmax}_{W \in \mathbb{L}} \frac{P(O|W)P(W)}{P(O)}$$

We can simplify the computation even further, if we consider that the probability of the observation in the denominator, $P(O)$, does not affect the maximization with respect to W , since the observation signal stays the same as we search over the sentences space. Consequently, the answer to our problem can be computed by the following form:

$$\hat{W} = \operatorname{argmax}_{W \in \mathbb{L}} P(O|W)P(W)$$

The two probabilities on the right hand side of the equation, represent the tools we need to address the recognition problem, as we have mentioned above: $P(W)$, which is the prior probability of each sentence of the language, is computed by the *language model*, whereas $P(O|W)$, which includes the probability of words being the concatenation of certain speech units and the probability of these speech units being realized as certain features, is calculated by the *acoustic model*.

Having presented the computational formulation of the automatic speech recognition problem, we can move on to present the required steps to be taken and models to be constructed in order to build a speech recognition system.

Figure 2.4: *MFCC feature extraction from a quantized digital waveform*[JM09].

2.4.2 Feature Extraction

[JM09] The first issue we have to address is how and in what format do we “insert” the speech waveform into our system. This process, known as *feature extraction* results in the extraction from the speech waveform of a sequence of acoustic *feature vectors*, each of them representing the information included in a small time window of the signal. These feature vectors are further preprocessed and finally presented as the input to the ASR system.

Mel Frequency Cepstral Coefficients

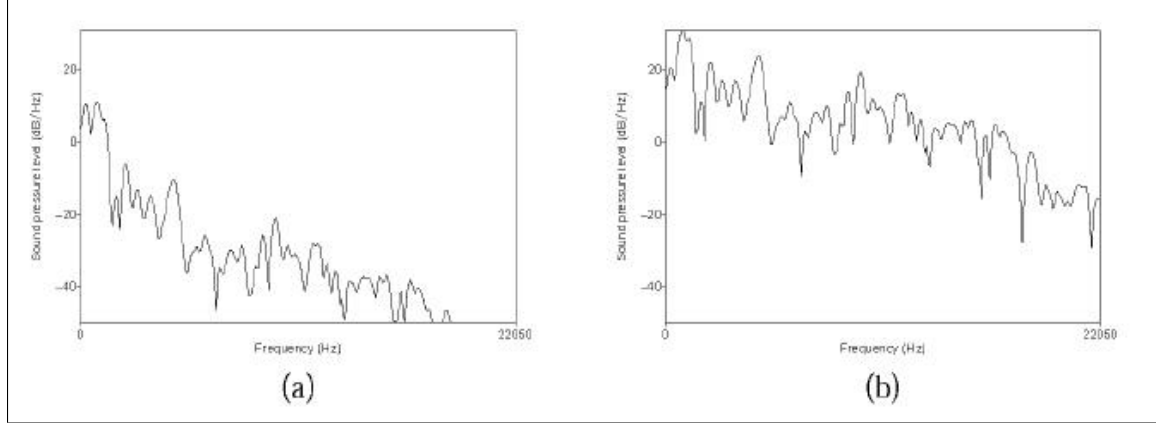
Although research in feature extraction moves towards using the raw waveform directly as input to the system, so far the most common feature representation in speech recognition systems is the Mel Frequency Cepstral Coefficients (*MFCC*). The steps involved in extracting the MFCC feature vectors follow the analog to digital conversion of the speech signal (sampling and quantization) and are outlined in the next paragraphs:

Pre-emphasis In the pre-emphasis step we want to amplify the amount of energy in the high frequencies of the signal. If we take a look at the spectrum of a vowel we will note a drop in energy as we move on to higher frequencies; this is known as *spectral tilt* and is due to the nature of the glottal pulse. **(EXPLAIN STH???)** Amplifying the energy of higher frequencies will improve phone detection accuracy as it will provide more information to the acoustic model coming from the boosted higher formants. In essence, pre-emphasis is applying to the signal a first order, high-pass filter whose equation is:

$$y[n] = x[n] - \alpha x[n - 1], \quad 0.9 \leq \alpha \leq 1.0$$

Windowing Given that the feature vectors we want to extract will be used to train phone classifiers, i.e. the acoustic model, we want them to be able to capture the spectral properties corresponding to these fundamental speech units. Consequently, since speech is a non-stationary signal - its statistical properties are not constant

Figure 2.5: Part of the spectrum of vowel [aa] before (a) and after (b) pre-emphasis[JM09].



across time - we extract the feature vectors from a small window of the speech signal that corresponds to a phone or subphone, where the signal can be considered stationary - its statistical properties stay constant across time.

The windowing process comes down to applying a filter to the signal that is non-zero inside some region and zero elsewhere, moving this filter along the speech signal and extracting segments of the signal (or *frames*). The window is characterized by its *width* (in milliseconds, also called *frame size*), the *overlap* between successive windows (usually a percentage of the width) and its *shape*, e.g. rectangular, Hamming etc. The Hamming window is usually preferred over the rectangular as it gradually reduces the signal values at the boundaries of the window towards zero, thus avoiding discontinuities which cause problems during the next step of feature extraction (Fourier analysis):

$$w_{\text{hamming}}[n] = \begin{cases} 0.54 - 0.46 \cos(\frac{2\pi n}{L}), & 0 \leq n \leq L - 1 \\ 0 & \text{otherwise} \end{cases}$$

The application of the filter is an element-wise multiplication of the signal values at each time step n with the values of the window:

$$y[n] = w[n]s[n]$$

where w is the window and s the signal.

Fourier Analysis The next step in the feature extraction process is to acquire the spectral information included in the windowed segments. The tool to extract such information, e.g. the amount of energy included in different discrete frequency bands of each segment, is the Discrete Fourier Transform. Given a part of a signal, the DFT will produce a complex number representing the magnitude and phase of

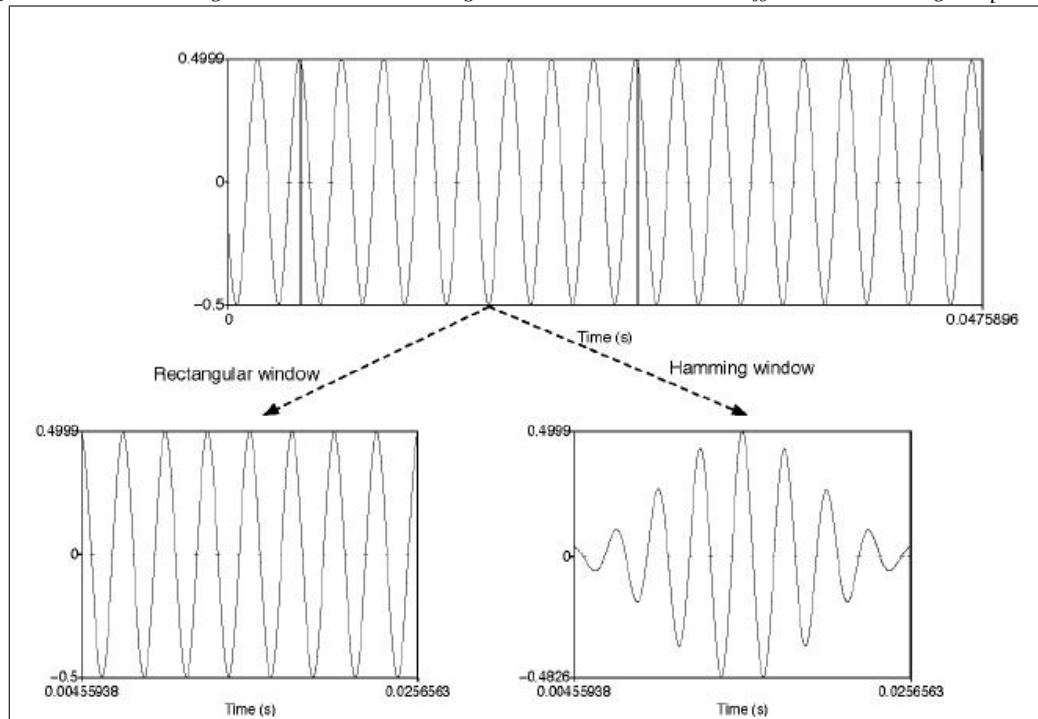
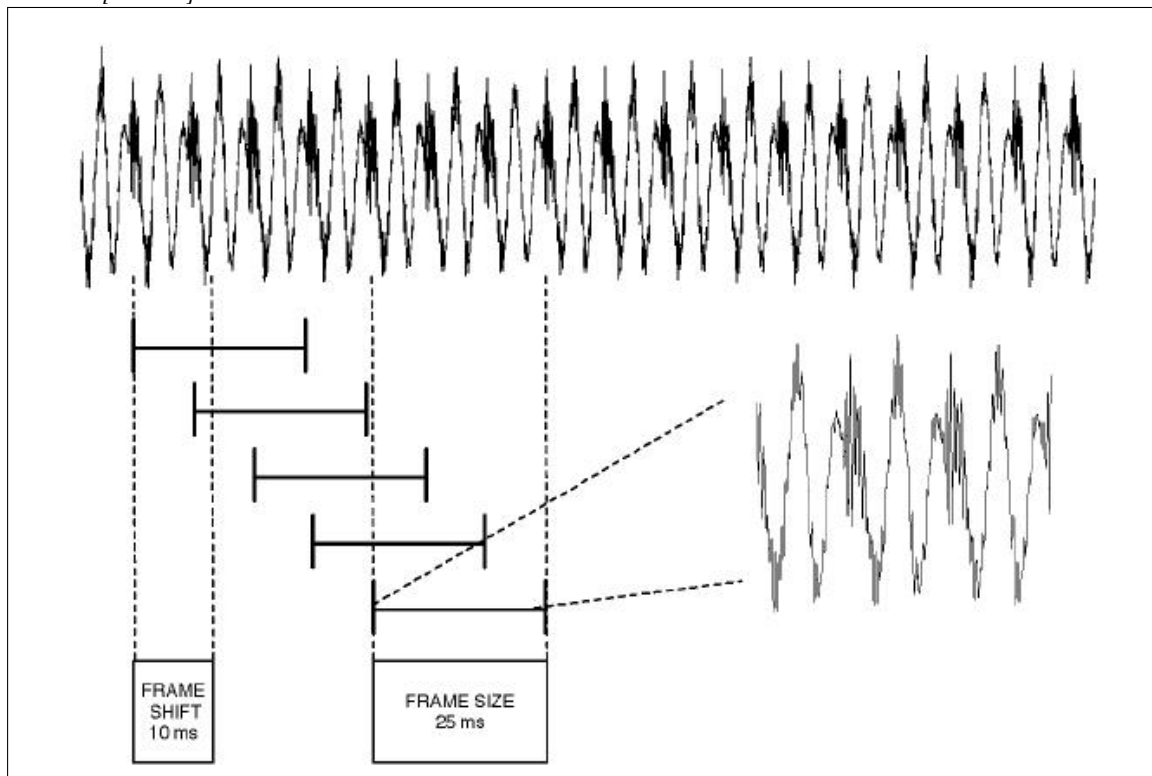
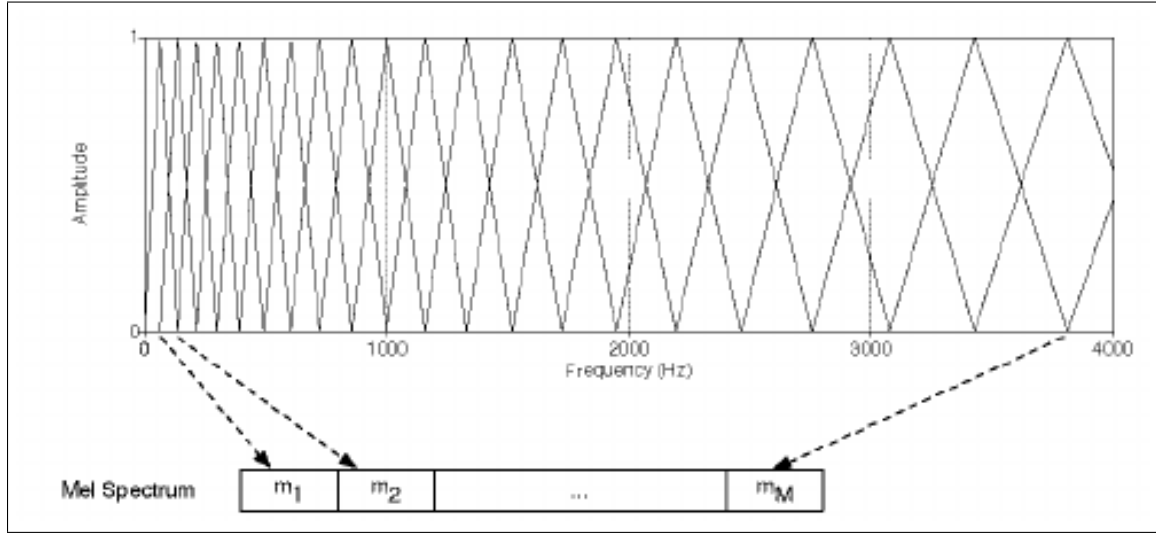
Figure 2.6: *Rectangular and Hamming windows and their effect on the signal*[JM09].Figure 2.7: *Windowing process with a rectangular window. After a figure by Brian Pellom*[JM09].

Figure 2.8: *The mel filterbank after Davis and Mermelstein (1980)[JM09].*

each frequency component of the corresponding segment.

Mel Filterbank Research has shown that the human hearing is not equally sensitive at all frequency bands - it is less sensitive at higher frequencies (above 1000 Hz). Moreover, humans are less sensitive to differences in amplitude at high amplitudes than at low amplitudes. Consequently, if we model this information and introduce it into the ASR system, we will improve its performance.

The way we take advantage of this information during the feature extraction process is by passing the DFT of the windowed signal parts through an array of triangular filters (a *filterbank*), which collect energy from each frequency band. These filters have their center frequencies spread on a *mel* scale, that is, 10 of them are spaced linearly below 1000 Hz and the remaining filters of the bank are spaced logarithmically above 1000 Hz. A *mel* is a unit of pitch defined so that pairs of sounds which are perceptually equidistant in pitch are separated by an equal number of mels[SVN37]. The mapping between frequency in Hz and the mel scale is described by the relation:

$$mel(f) = 1127 \ln\left(1 + \frac{f}{700}\right)$$

Finally, having passed the DFT of the signal through the filterbank, the final step is to take the logarithm of the mel spectrum values, which is a form of normalization to make the feature estimates less sensitive to variations in input.

The Cepstrum The next step in the MFCC feature extraction process is the computation of the *cepstrum*. As mentioned before, the speech waveform is created when a glottal source waveform is passed through the vocal tract which acts as a filter. As we have already mentioned, the shape of the vocal tract will determine the

Figure 2.9: *Effect of vocal tract on source signal.* Tomi H. Kinnunen, *Speech Technology Workshop*

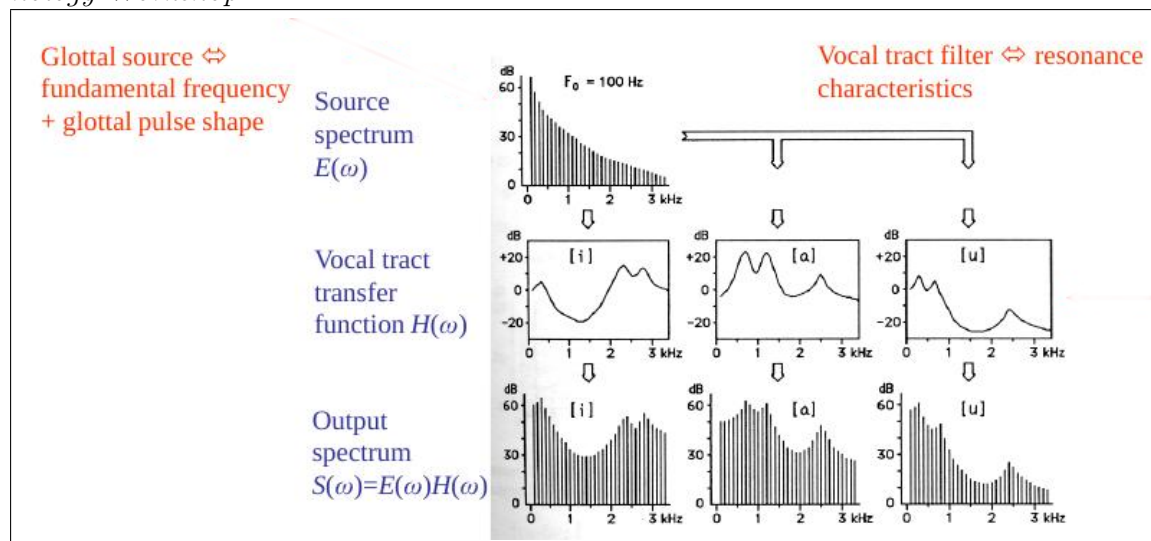
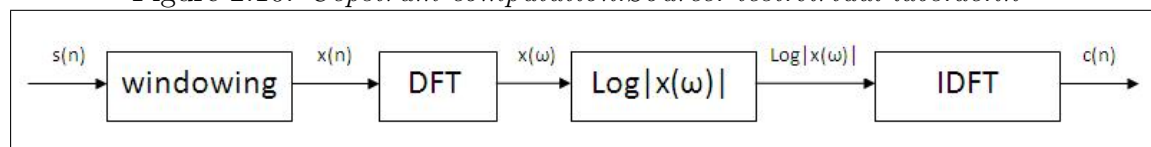


Figure 2.10: *Cepstrum computation.* Source: test.virtual-labs.ac.in

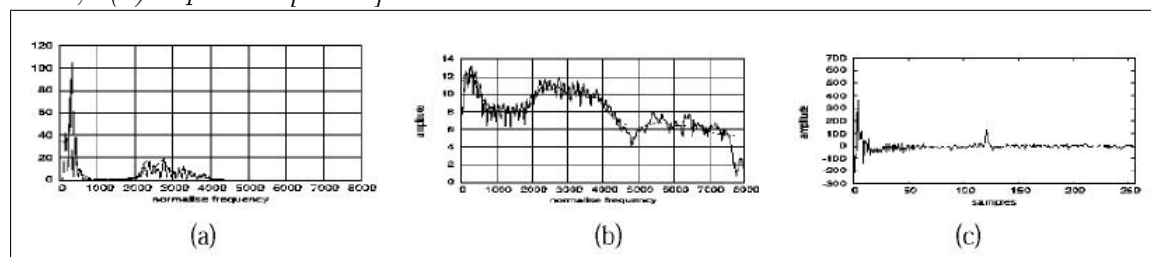


outcome of this filtering process, i.e. the sound, the phone that will be produced.

Therefore we aim to have some information about the vocal tract inserted into the feature vectors. The cepstrum provides us with a way of separating the glottal source from the vocal tract filter.

The peaks at lower values on the x-axis in the cepstrum of Figure 9 (c) correspond to the vocal tract characteristics whereas peaks at higher values correspond to the glottal source. Therefore, since we need information about the way a phone was produced, i.e. about the vocal tract shape, for the feature vectors we want to extract, we will keep a few of the first cepstral values (usually 12). Furthermore, an

Figure 2.11: *Cepstrum example:* (a) magnitude spectrum, (b) log magnitude spectrum, (c) cepstrum. [JM09].



important property of the cepstral coefficients is that their variance is uncorrelated, contrary to spectral coefficients, which are correlated at different frequency bands. This is extremely important for acoustic models based on Gaussian Mixture Models, as it allows us to keep the number of their parameters low.

Feature vectors The cepstral coefficients extracted from the previous process are just a part of the feature vectors. They are further enhanced by adding a few extra coefficients that provide more information helpful towards determining speech units such as phones.

The first extra piece of information we include by adding one more coefficient is the energy of the frame, which is defined as:

$$Energy = \sum_{t=t_1}^{t_2} s^2[t],$$

where s is the speech signal and t_1, t_2 are the boundaries of the frame.

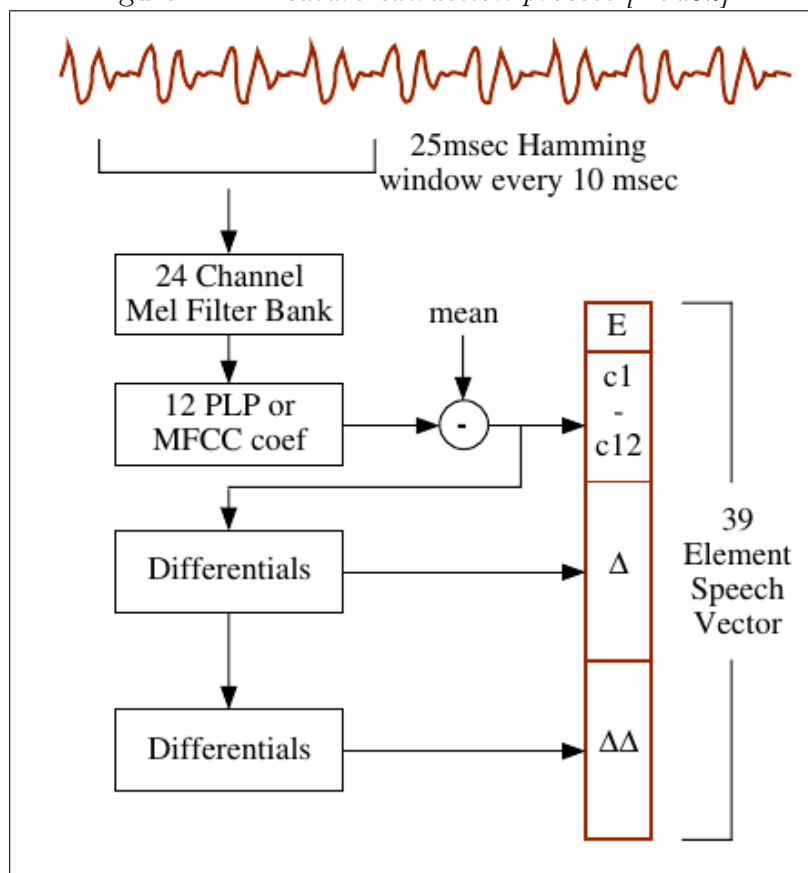
Energy is useful for phone detection since it correlates with phone identity: higher energy reveals the presence of e.g. a vowel whereas very low energy could identify a pause in speech.

Considering that speech properties change from frame to frame it is reasonable to expect that capturing these changes would provide more information about the nature of the speech signal. This motivated the use of *delta* and *delta-delta* coefficients for each one of the cepstral coefficients and the energy. The *delta* coefficients capture the change of the corresponding feature between successive frames whereas the *delta-delta* coefficients capture the change of the *delta* features between successive frames. The simplest way to compute these *delta* and *delta-delta* coefficients is by taking the difference of the corresponding features between successive frames:

$$d(t) = \frac{c(t+1) - c(t-1)}{2},$$

where c is the cepstral feature with delta coefficient d at time t .

This concludes the construction of the feature vectors. As mentioned before, we usually pick 12 cepstral coefficients plus an energy coefficient, thus we will have 39-dimensional feature vectors: 13 + 13 *delta* + 13 *delta-delta* coefficients. It is also common to concatenate the cepstral coefficients to produce higher dimensional vectors. This way we manage to include contextual information in the feature vectors. To deal with computational problems arising from the increase in the number of dimensions, various dimensionality reduction techniques are used, from simple ones such as Linear Discriminant Analysis or Principal Component Analysis, to more sophisticated such as techniques aiming to discover the lower dimensional manifold on which the feature vectors lie.

Figure 2.12: *Feature extraction process.*[You02].

2.4.3 The Language Model

The language model (*LM*) expresses how likely a given string of words is, taken into consideration certain linguistic constraints. In order to do this, we build on the idea of predicting the next word in a sequence of words, which is formalized with probabilistic models called *N-gram models*.

N-grams

An *N-gram* is a sequence of N words, e.g. a 2-gram or bigram is a sequence of two words, a 3-gram or trigram is a sequence of three words etc. An *N-gram model* is a probabilistic model which computes the N^{th} word of a sequence of N words given the previous N-1.

The power of N-grams becomes evident in areas such as speech or handwriting recognition, machine translation, spelling correction and natural language processing tasks. What all these areas have in common is that they might have to deal with noisy or ambiguous input. N-grams can deal with ambiguity by assigning a higher likelihood to word sequences that are valid according to the language constraints.

Since these models are capable of assigning a conditional probability to the next possible word in a group, we can exploit them to compute the joint probability of a sequence of words, i.e. a sentence, which is what we were aiming to from the beginning.

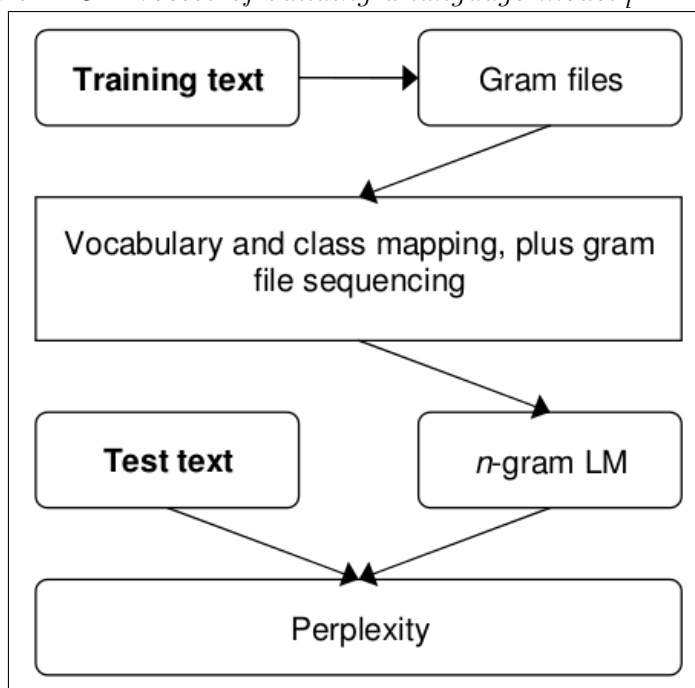
Suppose we have a sequence of n words $W = w_1, w_2, w_3, \dots, w_n$. Then, the probability $P(W)$ can be computed in the following way:

$$P(W) = P(w_1, w_2, w_3, \dots, w_n) = \prod_{i=1}^n P(w_i | w_1, w_2, w_3, \dots, w_{i-1})$$

As this computation for every possible word sequence in the language is very difficult if not impossible, we make the assumption that the i^{th} word depends only on the previous N-1 words (its *history*). At this point we can take advantage of the intuition behind the N-gram model:

$$P(W) = \prod_{i=1}^n P(w_i | w_1, w_2, w_3, \dots, w_{i-1}) \approx \prod_{i=1}^n P(w_i | w_{i-N+1}, \dots, w_{i-1})$$

to approximately compute the probability of the sequence. The expressions on the two sides would be exactly equal for sufficiently high n and if the language were ergodic, that is, the probability of any word could be estimated from sufficient history independent of the starting conditions.

Figure 2.13: *Process of building a language model.*[YEK⁺02].

Building a language model based on N-grams

[YEK⁺02]

Considering the intuition of N-grams and assuming that the probability of an N-gram occurring in an unknown text can be estimated from its frequency in a given training text, we can build language models based on N-grams.

The construction of such an LM can be broken down into three stages:

- Collect and store the N-grams of the training text (*corpus*)
- Possibly map some words into classes, e.g. out-of-vocabulary class mapping
????????????
- Count the N-grams and compute the N-gram probabilities

The last step of computing the probabilities is based on maximum likelihood estimation:

$$\hat{P}(w_i | w_{i-N+1}, \dots, w_{i-1}) = \frac{C(w_{i-N+1}, \dots, w_i)}{C(w_{i-N+1}, \dots, w_{i-1})}$$

where $C(.)$ is the count of a given word sequence extracted from the training text.

When building a language model based on N-grams, there are several factors one has to take into account before they decide on N. Resources constraints (e.g. storage) and size of the *vocabulary* (i.e. the set of distinct words in the language) will play a major role in deciding on N, as the number of parameters of the model grows

exponentially with $N : |V|^N$, where V is the vocabulary. However, storage needs because of parameters' size are lower than one would expect, because not all N -word combinations are acceptable/valid sequences in the language. What increases storage and computational needs however, is the large training sets required, so that our model estimates parameters with a minimum acceptable degree of confidence. Apart from processing, also acquiring these training sets might be difficult, especially for domain-specific applications, where training sets have to be specifically constructed.

Data sparsity and smoothing

[JM09] However, since the training set will always be finite, one can never have a sufficient number of N -grams for every valid N word sequence of the language. This issue faced when developing LMs is known as *data sparsity* and the technique used to deal with it is called *smoothing*. Smoothing aims to increase the robustness of the language model by redistributing the probability mass assigned by the maximum likelihood estimates: it removes some of it from higher counts of N -grams and assigns it to very low or zero counts, in order to ‘smooth’ the distribution.

Laplace Smoothing. *Laplace* or *add-one* smoothing is the simplest form of smoothing : we just add one to all the counts. Before we compute the ML probabilities, one has to take into account the extra $|V|$ ‘words’ that we added, in order to maintain the sum of all probabilities equal to 1. As this last step is necessary and because Laplace smoothing does not perform well, it is more convenient to use an *adjusted count*

$$c_i^* = (c_i + 1) \frac{K}{K + V}$$

where c_i is the original count and K is the number of word tokens.

Discounting. Instead of adding the same amount of probability mass to all N -gram probabilities, a different approach would be to remove some mass out of the higher counts and assign it to lower ones. We can therefore define a *relative discount* as the ratio of the new counts and the originals:

$$d_c = \frac{c^*}{c}$$

One algorithm applying discounting, is the *Good-Turing discounting*. The intuition behind it is that we use the MLE of N -grams occurring $c+1$ times in the training set, to define the MLE of N -grams occurring c times. The new smooth count c^* is thus defined as:

$$c^* = (c + 1) \frac{N_{c+1}}{N_c}$$

where N_i is the number of different N -grams occurring i times in the training set.

Back-off and interpolation. Discounting techniques allow us to distribute some probability mass equally to the unseen events. However, we can distribute it fairer if we take into consideration information from lower or higher order N-grams. This is the idea of *back-off* smoothing. In particular, *Katz back-off* always resorts to the (i-1)-gram if the i-gram has zero counts (i starting from our originally selected N). From a different point of view, we back-off to lower order N-grams only if we have zero evidence for a higher order one. On the other hand, *interpolation* deals with zero counts by summing estimates of all N-grams using weights (e.g. interpolate the estimates of unigrams, bigrams and trigrams).

Language model evaluation

Given the numerous applications of N-grams and LMs as well as their inherent drawbacks (static and finite vocabulary, finite training sets, more N-grams than can ever be collected and utilized), it is evident that we must have a way to compare LMs and evaluate their performance.

The most obvious way to compare two different language models would be to use them in our application and see which gives the best results. However, this way is expensive and time consuming as it is based on training and evaluating systems using huge speech datasets.

Another way to compare two LMs independent of the application, is to use the *perplexity* metric. The perplexity of an LM on a test set is a function of the probability that the LM assigns on it and is defined as:

$$PP(W) = P(w_1 w_2 w_3 \dots w_N)^{-\frac{1}{N}}$$

where $W = w_1 w_2 w_3 \dots w_N$ is the test set.

The best LM would be the one that has the minimum perplexity, since that would mean that it maximizes the test set probability, i.e. it better predicts the details of the test set. Another way to look at perplexity is as the weighted average branching factor of the language, that is, the number of best possible words following a word sequence. The smaller that number is, the better the work of the LM on coping with the ambiguity of the language.

Finally, when comparing language models using the perplexity metric, one should take care to use the *same vocabulary* for both LMs and evaluate them on the same test set which will be presented to the system for the first time during evaluation.

Recent advances in language modeling

In a paper presented recently [BDVJ03], a novel approach to building language models was presented which takes advantage of neural networks. The motivation

behind using NNs to build language models is based on the following problems associated with N-gram models:

- a huge amount of training data is needed to train LMs which will still have limited context capabilities (1-2 words)
- N-gram models ignore word similarity, which makes generalizing difficult

To deal with these issues becomes even more important when one considers that language models are probabilistic models using discrete random variables (words) which largely increases the amount of free parameters they need.

The writers suggest using distributed representations of words (i.e. real-valued feature vectors) which allows them 1) to identify similarities between words, since similar words will have similar feature vectors and 2) to exploit the smooth probability function modeled by a neural network in order to generalize: in this way, each vector representing a word will be able to provide information about a huge number of similar words, i.e. its “neighbors” in the feature space. In their work they present a neural network which simultaneously learns its parameters and the feature vectors associated with each word in the training set, and most importantly, the number of parameters it uses scales linearly with the vocabulary and context size. Due to the high computational cost of the training (higher than N-gram based models), they use parallel methods to efficiently train the model.

In recent years the use of neural networks for language modeling has included using recurrent neural networks which can take advantage of arbitrarily long contexts for each word (like humans do), something that was not possible with feedforward NNs [MKB⁺10].

2.4.4 The Acoustic Model

[You02]

According to the computational formulation of the ASR problem, we need the likelihood of the observed data (i.e. the acoustic signal) given the word sequence $P(O|W)$. However, it would be impractical and inefficient if we tried to compute this likelihood by building a separate model for each word in the language, since sub-word units are shared among different words. Instead, as mentioned before, the *acoustic model* calculates the probability of words being the concatenation of certain speech units and combines it with the probability of these speech units being realized as certain features, to produce the desired likelihood for a word.

We will first present some basic notions in acoustic modeling and then we will go into more details about this important part of the system.

Phones

The basic unit of speech analysis we use, is the *phone* which is the smallest identifiable unit we find in a stream of speech. The sequence of phones that constitutes each word in the training dataset is determined by a *pronouncing dictionary*. Using a phone sequence to represent each word makes it easy to add new words in the dataset just by adding them and their phone sequence to the dictionary.

Context-dependent phones

Given that there are thousands of words in a language, but just tens of phones (e.g. 44 in the English language) the computational and storage gain acquired from the use of phones as the basis of the acoustic model becomes immediately apparent. However, contextual effects like co-articulation cause large variations in the way that different sounds are produced even if in principle they correspond to the same phone. Hence, to achieve good phonetic discrimination, we build *context-dependent* phone models, with the most common being the *triphone*: for each phone there is a different model for every unique pair of left and right neighboring phones.

There are two dominant triphone models:

- *cross-word triphones*, which include phones of the previous and following words in the first and last triphones of the word of interest:

$$ten\ pots \rightarrow sil\ sil t_e\ t e_n\ e n_p\ n p_o\ p o_t\ o t_s\ t s_{sil}\ sil$$

The advantage of this approach is that they model co-articulation across word boundaries, but on the other hand, they complicate the decoding process since the phone models of each word depend on the following and preceding words as well.

- *word-internal triphones*, which explicitly encode word boundaries, thus making decoding easier:

$$ten\ pots \rightarrow sil\ sil t_e\ t e_n\ e n_- _p_o\ p o_t\ o t_s\ t s_{sil}\ sil$$

State of the art systems use mostly cross word triphones because of their ability to model contextual effects.

As a consequence, the number of distinct triphones greatly increases and the number of parameters for such systems can grow up to hundreds of millions, while at the same time we have too little training data in our disposal. In addition, we might

have unseen triphones appearing in evaluation tasks. To deal with these problems we have developed smoothing techniques, just as was the case with language models.

Smoothing techniques

Back-off and interpolation. When too little data is available for the training of a context-dependent model of a particular order, one can instead use a model of lower order at the expense of some inaccuracy in the modeling of the context: e.g. use a biphone or a context-independent phone (monophone) when we cannot use a triphone. In order to implement a more robust model one can use a weighted combination of models with various levels of context dependency (*interpolation*).

Parameter tying. An alternative that offers a greater degree of flexibility while maintaining the high level of context-dependency in the model, is the technique of *parameter-tying*, in which parameters of context-dependent phone models that are acoustically indistinguishable are tied together, to facilitate training in case there is little training data available. Before tying parameters together one has to apply some form of clustering to build the groups of phone models that will share their parameters. In practice, the most commonly used clustering technique is the *phonetic decision tree*, where a binary tree is built for each phone model and its leaves contain the parameters to be shared.

Acoustic Modeling with Gaussian Mixture and Hidden Markov models

State-of-the-art ASR systems use Hidden Markov Models to represent each phone in conjunction with Gaussian Mixture Models to determine the probability that an acoustic observation was produced by a certain phone. The high representational capabilities and ease of training of these models are what has made them prevalent in ASR. We will present the GMM/HMM acoustic model and in the next chapters we will examine their most recent competitor, that is, deep neural networks.

Gaussian Mixture Models.[You02][YD14]

Multivariate Gaussian random variables and Mixture Models.

Given that the feature vector corresponding to an acoustic observation is multi-dimensional (usually 39-dimensional as was presented above), we will have to treat it as a multivariate random variable and use a multivariate distribution to assign a probability to it. The reason why we choose the Gaussian distribution is not only its desirable computational properties but also its ability to approximate many real-world data (owing to the law of large numbers), such as speech features.

Supposing that Σ is the co-variance matrix, μ is the mean vector and d is the

number of dimensions of the feature vector, the multivariate Gaussian distribution is defined as:

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp(-1/2(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

However, because of the inherent multimodality of the speech features, a single Gaussian distribution is insufficient to describe them. Therefore, we use a *mixture* of Gaussian distributions:

$$p(\mathbf{x}) = \sum_{i=1}^M c_i \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

where c_i are the positive mixture weights and $\sum_{i=1}^M c_i = 1$. Usually, the number of the mixture components is chosen according to the nature of the problem and the information we have about the data. The variability and multimodality of the features may be due to multiple latent causes; provided we can identify these underlying causes we can match each one to the corresponding mixture component in the distribution.

As mentioned earlier, the Gaussian distribution is favorable both for its modeling and its computational properties. GMMs can model complex, multimodal distributions to any required level of accuracy and they can be trained using standard maximum likelihood approaches. Their attractiveness is also due to research into GMM training having come up with approaches to optimize the trade-off between their modeling effectiveness and the amount of training time and data needed. For example, we have the ability to reduce the number of free parameters (from $M \times d^2$ down to M) while still achieving high performance, if instead of using full co-variance matrices we opt for diagonal $\boldsymbol{\Sigma}$ or even use the same matrix for all mixture components. The use of diagonal co-variance matrices has been thought to impose uncorrelatedness among features, but, given that a mixture of Gaussians with diagonal $\boldsymbol{\Sigma}$ can at least effectively describe the correlations modeled by a single full co-variance Gaussian, this thought has been misleading.

Specifically for speech recognition, a number of ways has been proposed to improve recognition accuracy of a GMM system. We can discriminatively train the system after the generative maximum-likelihood training, so that we maximize the probability of generating the observed speech features in the training data, or augment the input speech features with bottleneck features acquired using neural networks (the latter will be examined later on in this project).

The set of free parameters to be estimated for a Gaussian-mixture distribution is denoted by $\boldsymbol{\Theta}$ and consists of : $\{c_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i\}$. In order to acquire the parameters

we rely on maximum likelihood methods and in particular on the Expectation-Maximization (EM) algorithm [DLR77]. The EM algorithm is used to find locally maximum likelihood parameter estimates of statistical models when the equations cannot be solved directly. It is especially useful for models involving latent variables apart from the parameters-to-be-estimated and the observable data. A GMM can be treated as such a model if we assume that each observable data point has a corresponding hidden data point specifying the component of the mixture that each point belongs to. Furthermore, the EM algorithm provides us with closed-form expressions for the computation of the estimates in the M-step:

$$c_i^{(j+1)} = \frac{1}{N} \sum_{t=1}^N h_i^{(j)}(t),$$

$$\boldsymbol{\mu}_i^{(j+1)} = \frac{\sum_{t=1}^N h_i^{(j)}(t) \mathbf{x}^{(t)}}{\sum_{t=1}^N h_i^{(j)}(t)},$$

$$\boldsymbol{\Sigma}_i^{j+1} = \frac{\sum_{t=1}^N h_i^{(j)}(t) [\mathbf{x}^{(t)} - \boldsymbol{\mu}_i^{(j)}][\mathbf{x}^{(t)} - \boldsymbol{\mu}_i^{(j)}]^T}{\sum_{t=1}^N h_i^{(j)}(t)}$$

where the posterior probabilities, i.e. the “latent” variables corresponding to the mixture components, computed in the E-step are:

$$h_i^{(j)}(t) = \frac{c_i^{(j)} \mathcal{N}(\mathbf{x}^{(t)} | \boldsymbol{\mu}_i^{(j)}, \boldsymbol{\Sigma}_i^{(j)})}{\sum_{m=1}^n c_m^{(j)} \mathcal{N}(\mathbf{x}^{(t)} | \boldsymbol{\mu}_m^{(j)}, \boldsymbol{\Sigma}_m^{(j)})}$$

The last equation computes the conditional probability for a given data point $\mathbf{x}^{(t)}$, $t = 1, \dots, N$ being generated from mixture component i using the current (denoted by j) parameter estimate.

Despite the ease of training of GMMs, they have two serious drawbacks when it comes to speech recognition systems. The first one is that they cannot model the sequence information contained in speech features. To balance their inability, we combine GMMs with more general models able to capture sequence information: the Hidden Markov Models, which will be presented next. The second disadvantage, is that, in spite of their huge modeling capabilities, GMMs are statistically inefficient for modeling data lying on or near a nonlinear manifold in the data space; this is the case however for speech features, despite their seemingly high dimensionality. To deal with this matter, there are a number of techniques we can apply to extract

the lower dimensional manifold of the features. We will go into more details about manifolds and speech features in the following chapter.

Hidden Markov Models and acoustic modeling.[YD14]

As we have already mentioned, mixture-of-Gaussian random variables (single- or multidimensional) lack a “temporal” dimension, which would make the length of the random vectors variable, in order to follow the length of the speech sequence we intent to model. Therefore, although Gaussian mixture models are appropriate for short-term sound patterns, we will need to introduce a new model appropriate for sequences of speech acoustic vectors.

Extending the notion of the random variable to the discrete-time random sequence will provide us with the necessary tool to model acoustic vector sequences of variable length. A discrete-time random sequence is a collection with variable length, consisting of random variables indexed by uniformly spaced discrete times. We will focus on the most commonly used class of random sequences which is the Markov sequences.

Markov sequences

The concept of *state* is a key point in Markov sequences. If we think of a system functioning as a Markov sequence generating random variables, then the configuration of the system at each time step is defined by a specific *state* of the sequence. If the state of the Markov sequence is confined to be discrete, then the Markov sequence is called a *Markov chain* and the possible values of each discrete state constitute the discrete state space. When each discrete state value is generalized to be a new random variable (either discrete or continuous) the Markov chain is generalized to the *Hidden Markov Sequence*, also called *Hidden Markov Model* when it characterizes statistical properties of real-world data sequences. The Hidden Markov Model is the tool that we will use to model speech units used, such as sub-phones.

Markov chains

A Markov chain is a discrete-time Markov sequence. Its state space is of discrete nature, finite and each element of the space is associated with a state in the chain:

$$q_t \in s^{(j)}, j = 1, 2, \dots, N$$

where q_t symbols a state.

A Markov chain denoted by $q_1^T = q_1, q_2, \dots, q_T$, is completely characterized by the initial state distribution probabilities (*priors*) and the transition probabilities defined by:

$$P(q_t = s^{(j)} | q_{t-1} = s^{(i)}) \doteq a_{ij}(t), \quad i, j = 1, 2, \dots, N$$

Given the transition probabilities of a Markov chain, the state-occupation probability

$$p_j \doteq P[q_t = s^{(j)}]$$

can be recursively computed by

$$p_i(t) = \sum_{j=1}^N a_{ij} p_j(t-1), \quad \forall i$$

Hidden Markov Models[RJ86][Rab89][DHS00]

If the states of a Markov chain are *emitting*, that is, they are able to generate observational output variables, then we call the chain an observable Markov sequence. However, since there is an one-to-one correspondence between the output of the chain and the states, the model is inadequate to describe real-world informational sources such as sequences of speech features. To overcome this limitation, we will add randomness to the Markov chain by associating each state with an observation probability distribution, thus creating the hidden Markov sequence. It is called hidden because the underlying Markov chain is no longer directly observable but it can be observed only through a separate random function characterized by the observation probability distributions, which overlap across the states.

A Hidden Markov Model (*HMM*) is characterized by:

- N , the number of states in the model
- K , the number of distinct observation symbols per state
- The transition probabilities, $\mathbf{A} = [\alpha_{ij}]$, $i, j = 1, 2, \dots, N$ where

$$\alpha_{ij} = P(q_t = j | q_{t-1} = i), \quad i, j = 1, 2, \dots, N$$

- The initial Markov chain prior probabilities

$$\pi = [\pi_i], \quad i = 1, 2, \dots, N$$

where $\pi_i = P(q_1 = i)$

- The observation probability distribution, $P(\mathbf{o}_t | s^{(i)})$, $i = 1, 2, \dots, N$.

If \mathbf{o}_t is discrete, the distribution associated with each state gives the probabilities of symbolic observation $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$:

$$b_i(k) = P[\mathbf{o}_t = \mathbf{v}_k | q_t = i], \quad i = 1, 2, \dots, N.$$

If the observation probability distribution is continuous, then the parameters Θ_i in the p.d.f. characterize state i in the HMM. The most common p.d.f. used

Figure 2.14: *Generate observation sequence from an HMM* .[YD14]

```

1: procedure DRAWFROMHMM( $A, \pi, P(\mathbf{o}_t|s^{(i)})$ )
     $\triangleright A$  is the transition probability
     $\triangleright \pi$  is the initial state occupation probability
     $\triangleright P(\mathbf{o}_t|s^{(i)})$  is the observation probability given a state (either Eq. 3.7 if discrete or Eq.
    3.8 if continuous)
2:   Select an initial state  $q_1 = s^{(i)}$  by drawing from the discrete distribution  $\pi$ 
3:   for  $t \leftarrow 1; t \leq T; t \leftarrow t + 1$  do
4:     Draw an observation  $\mathbf{o}_t$  based on  $P(\mathbf{o}_t|s^{(i)})$ 
5:     Make a Markov-chain transition from the current state  $q_t = s^{(i)}$  to a new state  $q_{t+1} = s^{(j)}$ 
        according to the transition probability  $a_{ij}$ , and assign  $i \leftarrow j$ .
6:   end for
7: end procedure

```

in speech processing is, as we have seen, the multivariate mixture of Gaussian distributions:

$$b_i(\mathbf{o}_t) = \sum_{m=1}^M c_{i,m} \mathcal{N}(\mathbf{o}_t | \boldsymbol{\mu}_{i,m}, \boldsymbol{\Sigma}_{i,m})$$

with $\boldsymbol{\Theta}_i = \{c_{i,m}, \boldsymbol{\mu}_{i,m}, \boldsymbol{\Sigma}_{i,m}\}$

Given these parameters one could consider the HMM as a generative model producing a sequence of observational data, \mathbf{o}_t , $t = 1, 2, \dots, T$. According to this perspective, the data at each time t is generated from the model according to:

$$\mathbf{o}_t = \boldsymbol{\mu}_i + \mathbf{r}_t(\boldsymbol{\Sigma}_i)$$

where state i at a given time t is determined by the evolution of the Markov chain characterized by α_{ij} and

$$\mathbf{r}_t(\boldsymbol{\Sigma}_i) = \mathcal{N}(0, \boldsymbol{\Sigma}_i)$$

is a zero-mean, independent and identically distributed (*IID*) residual sequence. Given that $\boldsymbol{\mu}_i$ is constant, the observation \mathbf{o}_t is also IID given the state. Consequently, the HMM would produce locally stationary sequences making it appropriate to model sub-phone units. A procedure to generate sequences of observations from an HMM is described in the figure 2.14.

The three basic problems for an HMM.

Given the HMM model as presented above, there are three main problems associated with it that apply to real-world problems:

- The evaluation problem: Suppose we have an HMM $(A_{ij}, b_{ik}, \boldsymbol{\Theta}_i)$. How do we determine the probability that a given sequence of observations was generated by that model?
- The decoding problem: Suppose we have an HMM $(A_{ij}, b_{ik}, \boldsymbol{\Theta}_i)$. How do we determine the most likely hidden state sequence that led to the generation of a given observation sequence?

- The parameter estimation problem: Given the basic structure of an HMM (number of states and number of distinct observation symbols) as well as a set of training observations, how do we determine the parameters $(A_{ij}, b_{ik}, \Theta_i)$?

The evaluation problem. Let \mathbf{q}_1^T be a finite length sequence of states in a Gaussian-mixture HMM and $P(\mathbf{o}_1^T, \mathbf{q}_1^T)$ be the joint likelihood of the observation sequence \mathbf{o}_1^T and the state sequence \mathbf{q}_1^T .

Then, $P(\mathbf{o}_1^T | \mathbf{q}_1^T)$ denotes the likelihood that the observation sequence \mathbf{o}_1^T is generated by the model conditioned on the state sequence \mathbf{q}_1^T and is in the form of:

$$\prod_{i=1}^T b_i(\mathbf{o}_t)$$

whereas the probability of state sequence \mathbf{q}_1^T is the product of transition probabilities:

$$P(\mathbf{q}_1^T) = \pi_{q_1} \prod_{t=1}^{T-1} a_{q_t q_{t+1}}$$

The joint likelihood $P(\mathbf{o}_1^T, \mathbf{q}_1^T)$ can be obtained as:

$$P(\mathbf{o}_1^T, \mathbf{q}_1^T) = P(\mathbf{o}_1^T | \mathbf{q}_1^T) P(\mathbf{q}_1^T)$$

Since the hidden state sequence \mathbf{q}_1^T is not known, we will have to sum over all possible state sequences in order to compute the desired probability:

$$P(\mathbf{o}_1^T) = \sum_{\mathbf{q}_1^T} P(\mathbf{o}_1^T, \mathbf{q}_1^T)$$

The amount of this computation though, is exponential in the length T of the observation sequence. However, an efficient algorithm (linear complexity in T) to evaluate the above expression has been found, based on the principle of optimality (dynamic programming) [Bel03]. The algorithm, known as *Forward algorithm* is described in figure 2.15

The decoding problem. The decoding problem consists of finding the most probable sequence of HMM hidden states given a sequence of observations. It is essentially a path-finding optimization problem that will be dealt with using again the dynamic programming paradigm. In fact, the decoding algorithm, also known as the *Viterbi algorithm* (figure 2.16) is very similar to the *Forward algorithm* presented above.

The Viterbi algorithm returns the maximum joint likelihood of the observation and state sequence as well as the corresponding state transition path. The optimal path for a left-to-right HMM, i.e. an HMM where transitions are only allowed in the forward direction, is equivalent to the information required to determine the optimal segmentation of the HMM states to match the observation sequence.

Figure 2.15: *The Forward algorithm for HMM probability evaluation.*[JM09]

```

function FORWARD(observations of len  $T$ , state-graph of len  $N$ ) returns forward-prob

  create a probability matrix forward[ $N+2,T$ ]
  for each state  $s$  from 1 to  $N$  do                                ; initialization step
    forward[ $s,1$ ]  $\leftarrow a_{0,s} * b_s(o_1)$ 
  for each time step  $t$  from 2 to  $T$  do                            ; recursion step
    for each state  $s$  from 1 to  $N$  do
      
$$forward[s,t] \leftarrow \sum_{s'=1}^N forward[s',t-1] * a_{s',s} * b_s(o_t)$$

    
$$forward[q_F,T] \leftarrow \sum_{s=1}^N forward[s,T] * a_{s,q_F}$$
      ; termination step
  return forward[ $q_F,T$ ]

```

Figure 2.16: *The Viterbi algorithm for HMM decoding.*[YD14]

```

function VITERBI(observations of len  $T$ , state-graph of len  $N$ ) returns best-path

  create a path probability matrix viterbi[ $N+2,T$ ]
  for each state  $s$  from 1 to  $N$  do                                ; initialization step
    viterbi[ $s,1$ ]  $\leftarrow a_{0,s} * b_s(o_1)$ 
    backpointer[ $s,1$ ]  $\leftarrow 0$ 
  for each time step  $t$  from 2 to  $T$  do                            ; recursion step
    for each state  $s$  from 1 to  $N$  do
      
$$viterbi[s,t] \leftarrow \max_{s'=1}^N viterbi[s',t-1] * a_{s',s} * b_s(o_t)$$

      
$$backpointer[s,t] \leftarrow \operatorname{argmax}_{s'=1}^N viterbi[s',t-1] * a_{s',s}$$

    
$$viterbi[q_F,T] \leftarrow \max_{s=1}^N viterbi[s,T] * a_{s,q_F}$$
      ; termination step
    
$$backpointer[q_F,T] \leftarrow \operatorname{argmax}_{s=1}^N viterbi[s,T] * a_{s,q_F}$$
 ; termination step
  return the backtrace path by following backpointers to states back in
    time from backpointer[ $q_F,T$ ]

```


Figure 2.17: *The Backward algorithm for HMM decoding.*[DHS00]

```

1 initialize  $\omega(T), t = T, a_{ij}, b_{jk}$ , visible sequence  $V^T$ 
2 for  $t \leftarrow t - 1$ ;
4      $\beta_j(t) \leftarrow \sum_{i=1}^c \beta_i(t+1) a_{ij} b_{jk} v(t+1)$ 
5 until  $t = 1$ 
7 return  $P(V^T) \leftarrow \beta_i(0)$  for the known initial state
8 end

```

The parameter estimation problem. The goal in HMM training, is to extract the model parameters so as to minimize the empirical risk with respect to the joint likelihood loss, involving a sequence of acoustic data and their corresponding linguistic labels. To estimate the parameters of an HMM model given training data, we will apply the Expectation-Maximization algorithm, also known as *Baum-Welch algorithm* in the context of HMM training. First however, we will introduce the *Backward algorithm*, which is a part of the EM computation for HMMs.

The *Backward algorithm* (figure 2.17) is very similar to the *Forward* but now we are moving backwards, that is, the algorithm computes

$$\beta_t(i) = P(\mathbf{o}_{t+1}^T | q_t = i), \quad t = 1, \dots, T-1$$

The EM algorithm uses both the Forward and Backward algorithms in the expectation E-step in order to obtain:

- the posterior state transition probabilities in the HMM

$$\xi_t(i, j) = \frac{\alpha_t(i) \beta_{t+1}(j) a_{ij} \exp N_{t+1}(j)}{P(\mathbf{o}_1^T | \theta_0)}, \quad t = 1, \dots, T-1$$

where $N_t(i)$ is the logarithm of the Gaussian p.d.f. associated with state i ,

- the posterior state occupancy probabilities

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j)$$

In the maximization M-step the parameters are computed using the current estimates of ξ and γ :

$$\hat{\alpha}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$

$$\hat{\Sigma}_i = \frac{\sum_{t=1}^T \gamma_t(i) (\mathbf{o}_t - \hat{\boldsymbol{\mu}}_i)(\mathbf{o}_t - \hat{\boldsymbol{\mu}}_i)^T}{\sum_{t=1}^T \gamma_t(i)}$$

$$\hat{\boldsymbol{\mu}}_i = \frac{\sum_{t=1}^T \gamma_t(i) \mathbf{o}_t}{\sum_{t=1}^T \gamma_t(i)}$$

for each state i .

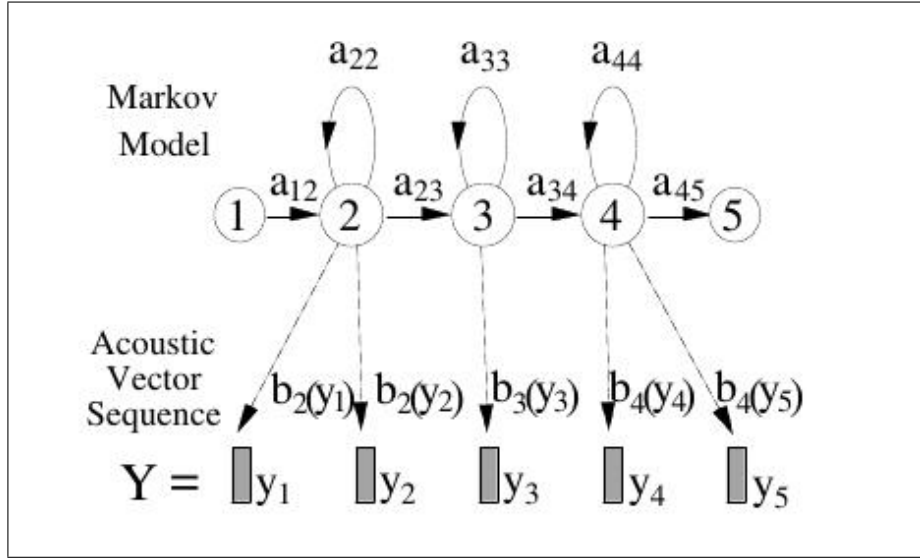
HMMs in speech modeling and recognition.[YD14][ST04] We have already seen that HMMs can be used as generative models to produce sequences of observations. They are able to produce sequences of variable length, which is of utmost importance for speech modeling and recognition, and they have proven to be good models for the statistical distribution of sequence data of speech acoustics. Consequently, HMMs have become very popular in the ASR community.

GMM/HMMs in ASR. As we have seen, a GMM/HMM is a statistical model that describes two dependent random processes, an observable and a hidden, where the observable is assumed to be generated by a hidden state according to a Gaussian mixture p.d.f.

In the context of speech, we think of a single HMM as a system generating acoustic features of a modeled speech unit, which can be a word, a syllable, a single phone, or, usually a context-dependent phone. The motivation behind choosing context-dependent phones, and therefore states, stems from the effort to reduce output variability of speech feature vectors associated with each state, leading to more detailed generative modeling. However this leads to an expansion of the state space, which we have seen how it is dealt with, at the beginning of the chapter.

The most common HMM model is the three-state left-to-right HMM (figure 2.18). The number of states is chosen based on the behavior of the vocal tract. It goes through three states when uttering a phone : changing from the previous phone, steady pronunciation of the current phone and changing to the next phone.

However, despite their advantages, HMMs have been found to have several weaknesses. The temporal independence of speech data conditioned on the HMM states and the lack of proven correlation between acoustic features and ways in which speech sounds are produced (e.g. speaking rate and style) have motivated the replacement of GMMs associated with each state by more realistic, temporally correlated dynamic systems containing hidden, continuous-valued dynamic structure (e.g. [Bil03]).

Figure 2.18: *Three-state, left-to-right HMM model.[ST04]*

2.4.5 The Decoder

[MPR01] The components described above are finally combined in the decoder which will find the most likely word sequence given a sequence of feature vectors. In order to define a common framework for the representation and use of the aforementioned models in LVCSR, we represent them by Weighted Finite State Transducers; an approach which provides significant algorithmic and engineering benefits.

A *finite-state transducer* is a finite automaton whose state transitions are labeled with both input and output symbols. Consequently, a path through the transducer encodes a mapping from an input to an output symbol sequence. *Weighted finite-state transducers (WFSTs)*, in addition to input/output symbols, have weights on the transitions which accumulate along paths to compute the total cost of a mapping from an input to an output symbol sequence. Thus, seeing the components of an ASR system as WFSTs and mathematical operations on them, allows for generalizing and efficiently implementing many of the common processing methods in speech recognition. We will briefly present the WFST framework in speech recognition, yet first we will introduce some necessary notation and algorithms.

Notation and Algorithms

- A *semi-ring* $(\mathbb{K}, \oplus, \odot, \bar{0}, \bar{1})$ is defined by a set of values \mathbb{K} , two binary operations of addition (\oplus) and multiplication (\odot) and two designated values $\bar{0}$ and $\bar{1}$. The addition operation is associative, commutative and has $\bar{0}$ as the identity element. The multiplication operation is associative, has $\bar{1}$ as the identity element, is distributive with respect to addition and has $\bar{0}$ as the annihilator

Figure 2.19: *Composition algorithm* .[MPR01].

```

WEIGHTED-COMPOSITION( $T_1, T_2$ )
1   $Q \leftarrow I_1 \times I_2$ 
2   $S \leftarrow I_1 \times I_2$ 
3  while  $S \neq \emptyset$  do
4       $(q_1, q_2) \leftarrow \text{HEAD}(S)$ 
5       $\text{DEQUEUE}(S)$ 
6      if  $(q_1, q_2) \in I_1 \times I_2$  then
7           $I \leftarrow I \cup \{(q_1, q_2)\}$ 
8           $\lambda(q_1, q_2) \leftarrow \lambda_1(q_1) \otimes \lambda_2(q_2)$ 
9      if  $(q_1, q_2) \in F_1 \times F_2$  then
10          $F \leftarrow F \cup \{(q_1, q_2)\}$ 
11          $\rho(q_1, q_2) \leftarrow \rho_1(q_1) \otimes \rho_2(q_2)$ 
12         for each  $(e_1, e_2) \in E[q_1] \times E[q_2]$  such that  $o[e_1] = i[e_2]$  do
13             if  $(n[e_1], n[e_2]) \notin Q$  then
14                  $Q \leftarrow Q \cup \{(n[e_1], n[e_2])\}$ 
15                  $\text{ENQUEUE}(S, (n[e_1], n[e_2]))$ 
16              $E \leftarrow E \cup \{((q_1, q_2), i[e_1], o[e_2], w[e_1] \otimes w[e_2], (n[e_1], n[e_2]))\}$ 
17 return  $T$ 

```

element: $\forall a \in \mathbb{K}, a \odot \bar{0} = \bar{0} \odot a = \bar{0}$. If \odot is also commutative the semi-ring is called *commutative*, which will be the case for all the semi-rings mentioned later.

- A WFST $T = (\mathcal{A}, \mathcal{B}, \mathcal{Q}, \mathcal{I}, \mathcal{F}, \mathcal{E}, \lambda, \rho)$ over a semi-ring \mathbb{K} is specified by a finite input alphabet \mathcal{A} , a finite output alphabet \mathcal{B} , a finite set of states \mathcal{Q} , a set of initial states $\mathcal{I} \subseteq \mathcal{Q}$, a set of final states $\mathcal{F} \subseteq \mathcal{Q}$, a finite set of transitions $\mathcal{E} \subseteq \mathcal{Q} \times (\mathcal{A} \cup \epsilon) \times (\mathcal{B} \cup \epsilon) \times \mathcal{K} \times \mathcal{Q}$, an initial state weight assignment $\lambda : \mathcal{I} \rightarrow \mathcal{K}$ and a final state weight assignment $\rho : \mathcal{F} \rightarrow \mathcal{K}$. $\mathcal{E}[q]$ denotes the sum of the number of states and transitions of \mathcal{T} .

Based on the notation just introduced we will present the operations on WFSTs that will be used in speech recognition applications.

- *Composition* is the basic operation that allows us to create complex WFSTs from simpler ones, thus putting together all the fundamental components of an ASR system.
- *Determinization* removes non-determinacy from the WFST by ensuring that each state has no more than a single output transition for a given input label. Not every WFST is determinizable, however, there is a *pre-determinization* algorithm that can be used to make determinizable an arbitrary WFST over the tropical semi-ring $((\mathcal{R} \cup \{-\infty, +\infty\}, \min, +, +\infty, 0))$ by inserting transitions labeled with special symbols. The determinization operation is particularly important in ASR considering the redundancy found in e.g. the WFST

Figure 2.20: *Determinization algorithm* .[MPR01].

```

WEIGHTED-DETERMINIZATION( $A$ )
1   $i' \leftarrow \{(i, \lambda(i)) : i \in I\}$ 
2   $\lambda'(i') \leftarrow \bar{1}$ 
3   $S \leftarrow \{i'\}$ 
4  while  $S \neq \emptyset$  do
5       $p' \leftarrow \text{HEAD}(S)$ 
6       $\text{DEQUEUE}(S)$ 
7      for each  $x \in i[E[Q[p']]]$  do
8           $w' \leftarrow \bigoplus \{v \otimes w : (p, v) \in p', (p, x, w, q) \in E\}$ 
9           $q' \leftarrow \{(q, \bigoplus \{w'^{-1} \otimes (v \otimes w) : (p, v) \in p', (p, x, w, q) \in E\}) :$ 
            $q = n[e], i[e] = x, e \in E[Q[p']]\}$ 
10          $E' \leftarrow E' \cup \{(p', x, w', q')\}$ 
11         if  $q' \notin Q'$  then
12              $Q' \leftarrow Q' \cup \{q'\}$ 
13             if  $Q[q'] \cap F \neq \emptyset$  then
14                  $F' \leftarrow F' \cup \{q'\}$ 
15                  $\rho'(q') \leftarrow \bigoplus \{v \otimes \rho(q) : (q, v) \in q', q \in F\}$ 
16              $\text{ENQUEUE}(S, q')$ 
17 return  $T'$ 

```

representing the pronunciation lexicon. A deterministic lexicon WFST will contain at most one path for any input string, thus less time and space will be needed to process the string (figure 2.20).

- *Minimization* transforms a WFST to an equivalent one with the fewest possible states and transitions which saves both space and time during its processing. Before minimising the transducer, a form of re-weighting (*weight-pushing*) is performed to redistribute weight among transitions as well as to improve search operations.

WFST models in speech recognition

There are four principal models of weighted finite-state transducers that are used in speech recognition:

- G : the word level grammar
- L : the pronunciation lexicon
- C : the context-dependency transducer
- H : the HMM transducer

The word level *grammar transducer* G has a state w_i for each word and transitions are added according to the N-gram model used in the grammar. For example,

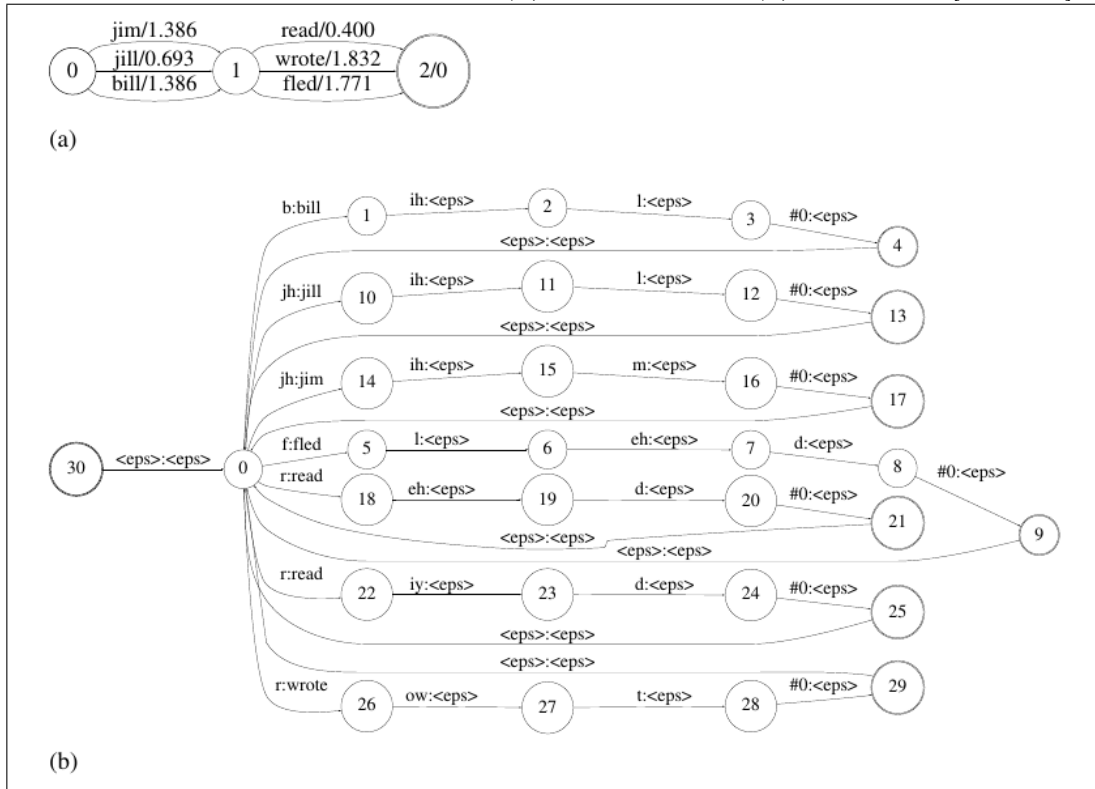
a bigram grammar has a transition from state w_1 to w_2 for every bigram w_1w_2 seen in the training corpus. The label of the transition is then w_2 and the weight is the negative logarithm of the transition probability ($-\log(\hat{p}(w_2|w_1))$). To deal with unseen N-grams while keeping the complexity of constructing the WFST low, we introduce a back-off state b . An unseen bigram w_1w_3 is then represented as two transitions: an ϵ transition w_1b with weight $-\log(\beta(w_1))$ and a transition bw_3 with weight $-\log(\hat{p}(w_3))$. Because ϵ transitions introduce non-determinism in the WFST, we can treat ϵ labels as normal symbols during determinization, thus keeping the number of transitions low - otherwise transitions become quadratic with respect to the size of the vocabulary after determinization.

The *pronunciation lexicon transducer* L is the Kleene closure of the union of individual word pronunciations. It is easy to see that L is, in general, not determinizable, considering the existence of homophones and the fact that the first word of the output string might be impossible to determine before the entire phone string is scanned. To make L determinizable we add a number of disambiguation symbols ($\#_i$) as well as a symbol ($\#_0$) to mark the end of the phonetic transcription of each word. The new lexicon transducer after the addition of these symbols is determinizable and is denoted by \bar{L} .

The *context dependency transducer* represents a mapping from context independent phones to context dependent units. The transducer has a state for every pair of phones with label (a, b) , where a is the past and b is the future phone, and transitions marked as $a:\text{phone}/\text{left context_right context}$. To apply the context dependent triphone models often used in ASR in the WFST framework, we need to be able to compose a context dependency transducer with the lexicon transducer introduced earlier. In order to make this composition feasible we first invert the context dependency transducer (interchange input and output labels) and create the transducer C which maps from context dependent triphones to context independent phones.

The final transducer that is used in the decoding process is the *Hidden Markov Models transducer* H which is the closure of the union of the individual HMMs of the acoustic model.

Applying the composition operation on the transducers presented here will output the decoding transducer which we will further optimize to help decoding and make it as efficient as possible. First, composing the lexicon and grammar transducers gives a new transducer that maps from phones to word strings restricted to the grammar ($L \circ G$). The resulting transducer is then composed with the context dependency transducer C and the resulting transducer maps from context dependent phones to word strings restricted to the grammar ($C \circ L \circ G$). Finally, $C \circ L \circ G$ is composed with the HMM transducer H and the outcome is a transducer that maps

Figure 2.21: *Transducer example: (a) Grammar G , (b) Lexicon \bar{L} [MPR01].*

from the identifiers of context-dependent HMM states to word strings restricted to G ($H \circ C \circ L \circ G$).

After the construction of the final transducer we have to determinize and minimize it so that we have the optimal transducer for recognition. By determinizing it we eliminate redundant paths which reduces recognition time. Moreover, if the determinization is applied after each composition stage during the construction of the graph, the composition operations that follow are performed more efficiently and the total size of the transducer is reduced. Once we have determinized the transducer we can further optimize it by minimization and the weight pushing process that precedes it. As far as weight pushing is concerned, provided we are using the $-\log$ semi-ring, we can have large efficiency gains during the Viterbi beam search, and make sure that all weights leaving a state sum up to one (which is desirable for a language model).

Given the $HCLG$ transducer and an utterance of N frames, how do we decode it, namely how do we find the most likely word sequence and its corresponding state-level alignment?

The first step is to construct an acceptor U of the utterance, which is a WFST with identical input and output symbols. The acceptor has $N+1$ states with an arc for each (time, context-dependent HMM state) combination and weights on these

Chapter 3

Manifold and Representation learning

In the first chapter we briefly saw an intuitive definition of manifolds. In this chapter we will give a mathematical definition of manifolds and examine the role that they play in representation learning.

Representation learning is the process of learning good representations of the data, that make it easier to extract information that will be useful when we build classifiers or other predictors. A good representation is one that captures the posterior distribution of the latent explanatory factors that led to the observed data and as such, it is particularly effective as input to a supervised predictor [BCV13].

3.1 Manifold learning

[BGC15]

Manifold learning is an approach to representation learning that builds on the *manifold hypothesis* ([Cay05]). According to this hypothesis, real-world data presented in high dimensional spaces is expected to concentrate on the vicinity of a manifold \mathcal{M} of much lower dimensionality $d_{\mathcal{M}}$, embedded in high dimensional input space \mathcal{R}^{d_x} .

This hypothesis is a kind of prior assumption about the data generating distribution that seems particularly fit to artificial intelligence tasks involving data, such as speech, text, music and images. The common thing about such data, is that if we choose configurations of the observed variables at random, according to a factored distribution (e.g. a uniform distribution), it is very unlikely to generate the kind of observations we want to model. For example, if we uniformly pick values for an acoustic signal, we will most likely create an unnatural speech sample. Consequently, manifold learning algorithms try to discover where probability concentrates in the

input space, as it will be a very small region of the total space of configurations.

Furthermore, making very small changes to the input data, e.g the values of the pixels of an image, will produce natural-looking data similar to the original input. This is another aspect of the manifold hypothesis, that is, probable configurations are likely to be surrounded by other possible configurations.

3.1.1 Mathematical formulation of manifold learning

[Cay05]

We will now give the mathematical definition of manifolds and mathematically formulate manifold learning.

Definition 15. *A homeomorphism is a continuous function whose inverse is also continuous.*

Definition 16. *A d -dimensional manifold \mathcal{M} is a set that is locally homeomorphic with \mathcal{R}^d . That is, $\forall x \in \mathcal{M}$ there is an open neighborhood around x , \mathcal{N}_x , and a homeomorphism $f : \mathcal{N}_x \rightarrow \mathcal{R}^d$. These neighborhoods are known as coordinate patches and the map as a coordinate chart. The image of the coordinate charts is referred to as the parameter space.*

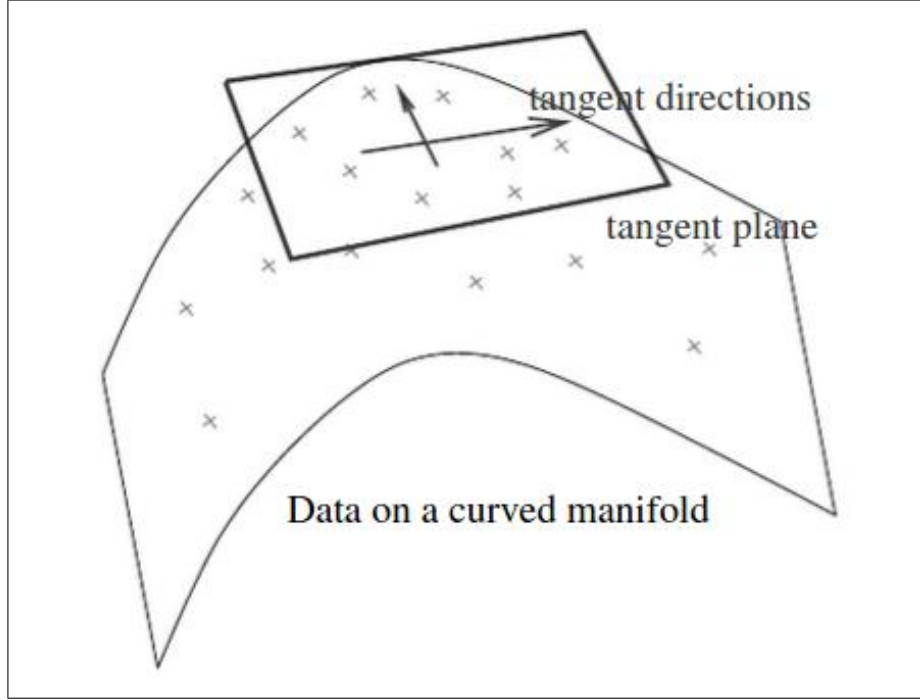
Intuitively, the dimension d of the manifold indicates the number of independent ways by which a probable configuration can be locally transformed into another probable configuration (remember the example of tiny changes in pixel values of an image). Another concept associated with this intuition is the set of *tangent planes* of a manifold. A tangent plane at a point x on a d -dimensional manifold is given by d basis vectors which cover the local dimensions of variation on the manifold, i.e. they show how x can be changed while staying on the manifold (figure 3.1).

We will only deal with manifolds which are subsets of \mathcal{R}^D where $D \gg d$. That is, the manifold will lie in \mathcal{R}^D but will be homeomorphic to a lower dimensional space \mathcal{R}^d .

Before moving on to manifold learning we will further constrain the manifolds of interest.

Definition 17. *A smooth or differentiable manifold is a manifold such that each coordinate chart is differentiable with a differentiable inverse. We say that such a chart is a diffeomorphism.*

Definition 18. *An embedding of a manifold \mathcal{M} into \mathcal{R}^D is a smooth homeomorphism from \mathcal{M} to a subset of \mathcal{R}^D .*

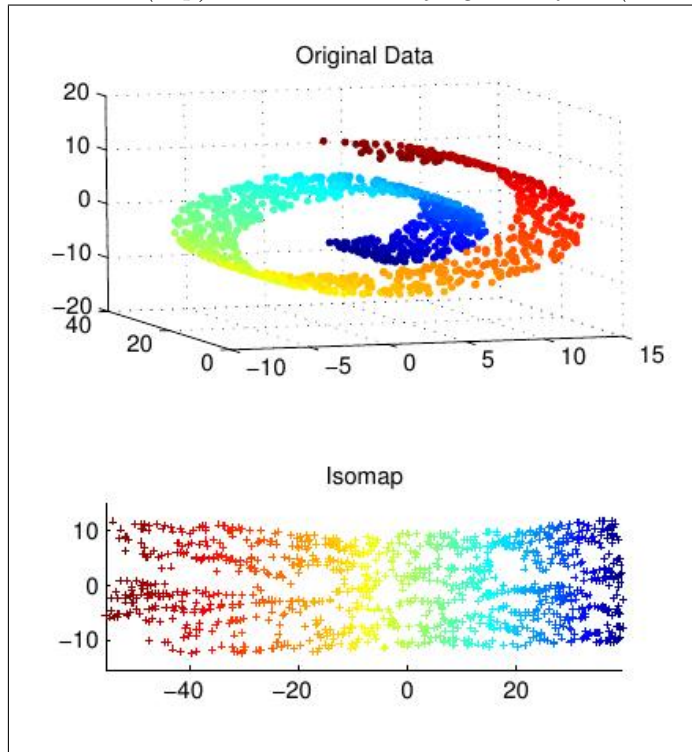
Figure 3.1: *Tangent plane and directions of variation on the manifold [BGC15].*

An embedding is an important concept that characterizes a manifold \mathcal{M} and is a representation of the data points on or projected on the manifold. It is usually given by a low-dimensional vector, with less dimensions than the surrounding space, in which the manifold lies. The learning algorithms that will be presented later will learn either directly an embedding for each training data point, or a more general mapping, also called encoder or representation function, that maps any point in the high-dimensional input space to its embedding.

Definition 19. Suppose we are given a set of \mathcal{N} points $x_1, x_2, \dots, x_N \in \mathcal{R}^D$ that lie on a d -dimensional manifold M which can be described by a single coordinate chart $f : \mathcal{M} \rightarrow \mathcal{R}^d$. Manifold learning is the process of finding $y_1, y_2, \dots, y_N \in \mathcal{R}^d$, where $y_i \doteq f(x_i)$.

In other words, manifold learning is the process of trying to learn the manifold of a set of available points, or find an embedding of the underlying manifold. The most frequently used example in the manifold literature is the Swiss roll, which is a two-dimensional manifold embedded in the three-dimensional space (\mathcal{R}^3). What is important to note is the fact that the distances between points in the original dataset are maintained in the two-dimensional dataset. This is due to the chart f between the two datasets being a homeomorphism.

The fact that the point distances between the datasets are preserved is a basic characteristic of the manifold learning algorithms that will be presented in the next

Figure 3.2: *Swiss roll (top) and the underlying manifold (bottom).* [Cay05].

section.

3.2 Manifold learning algorithms

3.2.1 Dimensionality reduction

[Ver08] We know that most learning algorithms performing well in low-dimensional datasets, perform poorly when applied to high-dimensional datasets. This phenomenon is known as the *curse of dimensionality*. In order to effectively apply learning algorithms in high-dimensional data, we would like first to map it to a low-dimensional space, while preserving much of the important information, and then run the algorithms in the projected space. This process is known as *dimensionality reduction*.

One way to verify the quality of a dimensionality reduction technique is to test how well the mapping preserves pairwise distances. This idea is based on the assumption that the distances between points in space relate to the dissimilarity between the corresponding observations. It is evident that this distance preservation would be much easier if we could find the underlying manifold structure of the data; we would then limit our work on points lying on the manifold and not in the whole ambient space.

A question that arises here, is, whether we can project any dataset on a lower-dimensional space, which maintains structure. The *Johnson-Lindenstrauss lemma* states that any n points in high dimensional euclidean space can be mapped onto k dimensions where $k \geq \mathcal{O}(\frac{\log n}{\epsilon^2})$ without distorting the euclidean distance between any two points more than a factor of $1 \pm \epsilon$ [Mah09], [DG03].

Lemma 2. *For any $0 < \epsilon < 1$ and any integer n , let k be a positive integer such that $k \geq 4(\frac{\epsilon^2}{2} - \frac{\epsilon^3}{3})^{-1} \ln n$.*

Then for any set V of n points in \mathcal{R}^d there is a map $f : \mathcal{R}^d \rightarrow \mathcal{R}^k$ such that $\forall u, v \in V$:

$$(1 - \epsilon)\|u - v\|^2 \leq \|f(u) - f(v)\|^2 \leq (1 + \epsilon)\|u - v\|^2$$

This map can be found in randomized polynomial time.

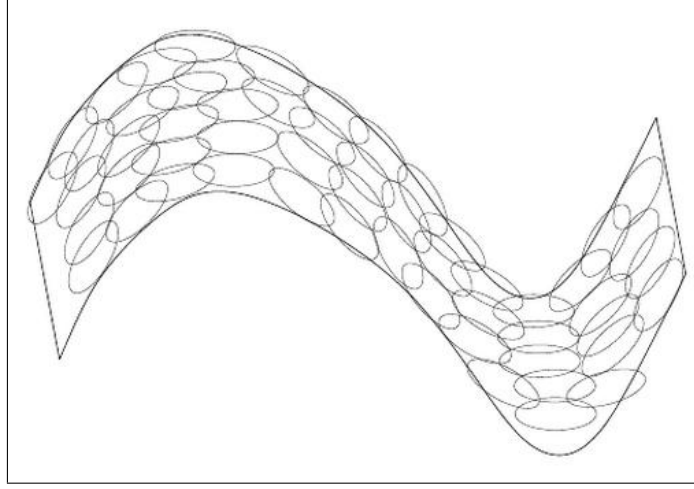
A general outline of the lemma is the following:

- Construct a random projection over k -dimensional subspaces.
- Prove that the expected value of the euclidean distance of the random projection is equal to the euclidean distance of the original subspace.
- Prove that the variance of the euclidean distance is greater than the specified error factor only with a probability $\frac{2}{n^2}$, such that the union bound of this probability across all points is less than $1 - \frac{1}{n}$.

3.2.2 Algorithms

So far, most manifold learning algorithms are unsupervised learning procedures attempting to uncover the underlying manifold structure. They use a nearest-neighbor graph which has one node per training example and edges connecting near neighbors. Each node is associated with a tangent plane which spans the directions of variations associated with a neighborhood of the graph and a coordinate system that associates each training example with an embedding. This coordinate system can be thought of as a local Euclidean system or a locally flat Gaussian, with very small variance in the directions orthogonal to the plane, and a very large variance in the directions defining the local coordinate system (figure 3.3). With this approach, a generalization is possible to new examples by a form of interpolation between neighbors and a global coordinate system can be obtained through an optimization or solving a linear system. Furthermore, by formulating a problem in terms of graph structures one avoids making any assumptions about the data distribution ([YXZ⁺07]).

Figure 3.3: *Tangent planes are tiled together to cover the manifold, forming a global coordinate system [BGC15].*



However, the methods based on nearest-neighbor graph work well provided we have a very large number of training examples to cover all the curves and twists of the underlying manifold. This need for huge amounts of data, as well as the fact that manifolds of interest in AI, such as speech and images, have many curves and twists, has motivated the use of deep learning methods and distributed representations to uncover the manifold structure. These methods try to learn a coordinate system for the main latent factors that explain the structure of the underlying generating distribution.

Furthermore, such algorithms work in batch mode and they do not provide a way of mapping new points from the high to the low dimensional space without re-running the algorithm from scratch; to tackle this issue, a number of approaches have been suggested ([AE], [SR03],[BPV03]).

We will now introduce some basic manifold learning algorithms.

Isomap

Isometric feature mapping (*ISOMAP*) was one of the first manifold learning algorithms. Isomap contains the following three steps:

- Construct the k nearest-neighbor graph. The number of neighbors k is defined after trying with various values and comparing the results. It is also possible to find the neighbors with an approximate nearest-neighbor procedure, e.g. use neighborhoods with a certain ϵ radius.
- Estimate the distances along the manifold (geodesic distances) between input points using shortest-path distances on the neighbor graph constructed in the

Figure 3.4: *Classical Multidimensional Scaling [Cay05].*

classical Multidimensional Scaling
input: $D \in \mathbb{R}^{n \times n}$ ($D_{ii} = 0$, $D_{ij} \geq 0$), $d \in \{1, \dots, n\}$

1. Set $B := -\frac{1}{2}HDH$, where $H = I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top$ is the centering matrix.
2. Compute the spectral decomposition of B : $B = U\Lambda U^\top$.
3. Form Λ_+ by setting $[\Lambda_+]_{ij} := \max\{\Lambda_{ij}, 0\}$.
4. Set $X := U\Lambda_+^{1/2}$.
5. Return $[X]_{n \times d}$.

previous step.

Isomap assumes that there is an isometric chart mapping the input points to the low-dimensional space, i.e. it preserves the distances between points. If the manifold is smooth enough then the geodesic distance between nearby points is almost linear. Given that the manifold locally resembles the Euclidean space, the Euclidean distance can be used to estimate distances. On the other hand, for points that are far apart from each other, we cannot estimate distances in the same way, because the manifold structure is not the same globally as it is locally. To compute such distances we use shortest-paths algorithms on the neighbors graph we have constructed.

- Using Multidimensional Scaling find points in lower-dimensional Euclidean space whose interpoint distances match the geodesic distances found in the previous step.

In the third step Isomap finds points whose Euclidean distances equal the geodesic distances found earlier. Supposing that the manifold is isometrically embedded, we are certain that such points exist and are unique to translation and rotation. To find such points we apply a technique called *Multidimensional scaling* (figure 3.4), which, given a matrix $D \in \mathbb{R}^{n \times v}$ of dissimilarities constructs a set of points whose interpoint Euclidean distances match closely those in D .

Furthermore, Isomap has the following two important properties. First, it automatically estimates the dimensionality of the underlying manifold: the number of non-zero eigenvalues found by MDS is the underlying dimensionality.

Second, under the following assumptions, Isomap is guaranteed to recover the parameterization of a manifold.

- The manifold is isometrically embedded into \mathbb{R}^D .

Figure 3.5: *ISOMAP* [Cay05].

Isomap
input: $x_1, \dots, x_n \in \mathbb{R}^D, k$

1. Form the k -nearest neighbor graph with edge weights $W_{ij} := \|x_i - x_j\|$ for neighboring points x_i, x_j .
2. Compute the shortest path distances between all pairs of points using Dijkstra's or Floyd's algorithm. Store the squares of these distances in D .
3. Return $Y := \text{cMDS}(D)$.

- The underlying parameter space is convex, or intuitively, the parameter space of the manifold cannot contain any holes.
- The manifold is compact and well-sampled everywhere.

Locally Linear Embedding

The intuition behind Locally Linear Embedding (*LLE*) stems from thinking of the manifold as a collection of overlapping coordinate patches. Then, if the neighborhood sizes are small and the manifold does not have many curves or twists, then these patches will be approximately linear. The concept is to find these patches, discover their geometry and find a mapping from the manifold to \mathcal{R}^d that preserves the local geometry and is approximately linear. Since these patches are overlapping, the local reconstructions will combine into a global one.

The first step of LLE tries to model the manifold as a collection of linear patches and uncover their geometry. In order to do so, it represents each point x_i as a weighted, convex combination of its nearest neighbors. The weights are chosen based on minimizing the squared error:

$$\|x_i - \sum_{j \in N(x_i)} W_{ij} x_j\|^2$$

where $N(i)$ is the set of nearest-neighbors of x_i .

The weight matrix W reveals the layout of the neighbors around each point; thus it is W that captures the local geometry of each patch. Apart from the convexity, we impose the following constraint on the weights:

$$W_{ij} = 0 \text{ if } j \notin N(i)$$

These two constraints have an important physical meaning: the convexity makes the weights invariant to global rotations, translations and scalings, e.g.

$$\|x_i + a - \sum_{j \in N(x_i)} W_{ij}(x_j + a)\|^2 = \|x_i - \sum_{j \in N(x_i)} W_{ij}x_j\|^2$$

whereas the second enforces locality on the patches.

The above minimization problem has a closed-form solution for the weight matrix which, for each x_i is given by:

$$\hat{W}_i = \frac{\sum_k C_{jk}^{-1}}{\sum_{lm} C_{lm}^{-1}}$$

where C is the local covariance matrix

$$C_{jk} \doteq (x_i - \eta_j)^T (x_i - \eta_k)$$

and η_j, η_k are neighbors of x_i . $\hat{W} \in \mathcal{R}^{n \times k}$ is then transformed into the sparse $W \in \mathcal{R}^{n \times n}$ by setting $W_{ij} = \hat{W}_{il}$ if x_j is the l^{th} neighbor of x_i and $W_{ij} = 0$ if $j \notin N(i)$.

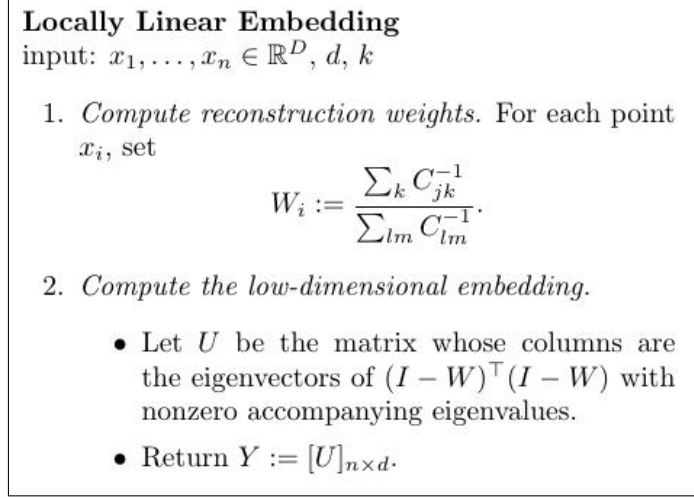
The second step in LLE is to find a configuration of points in the d -dimensional parameter space whose local geometry is described well by W . Contrary to Isomap LLE cannot discover the number of dimensions d ; we have to provide it based on prior knowledge about the parameter space or estimate it using other techniques (e.g. Conformal Eigenmaps, [SS05]). Such a configuration is found by minimizing the cost function:

$$\sum_i \|y_i - \sum_j W_{ij}y_j\|^2 = Y^T[(I - W)^T(I - W)]Y = Y^TMY$$

with respect to $y_1, y_2, \dots, y_n \in \mathcal{R}^d$ and under constraints: $Y^TY = I$ and $\sum_i Y_i = \mathbf{0}$. The first constraint forces the solution to be of rank d whereas the second centers the embedding around the beginning of the axes. Minimizing the cost function corresponds to setting the columns of Y to the eigenvectors that correspond to the minimum d eigenvalues of \mathcal{M} . However, the minimum (eigenvector, eigenvalue) pair is $(\mathbf{1}, 0)$. Therefore, to avoid the last coordinate being identical for all points, we pick the minimum d eigenvectors that are not constant.

Laplacian Eigenmaps

The Laplacian Eigenmaps algorithm is based on the concept of the Laplacian of a graph: Given a graph and the corresponding matrix containing edge weights, the Laplacian of the graph is defined as $L \doteq D - W$, where D is defined as the diagonal

Figure 3.6: *LLE* [Cay05].

matrix defined as $D_{ii} = \sum_j W_{ij}$. The eigenvalues of L reveal information about the structure of the graph, such as whether it is complete or connected. Since in our case the weight matrix contains information about the neighborhoods of the points, its Laplacian will help us discover local information about the underlying manifold.

In Laplacian Eigenmaps we have two choices for constructing for the weight matrix W . Either use the simple scheme:

$$W_{ij} = \begin{cases} 1 & \text{if } j \in N(i) \\ 0 & \text{if } j \notin N(i) \end{cases}$$

or use the Gaussian heat kernel:

$$W_{ij} = \begin{cases} e^{\frac{-\|x_i - x_j\|^2}{2\sigma^2}} & \text{if } j \in N(i) \\ 0 & \text{if } j \notin N(i) \end{cases}$$

Like LLE we now use W to find points in the d -dimensional parameter space (d is set or estimated in advance) that preserve the relations of their original counterparts.

The cost function we are aiming to minimize now is based on the fact that we want to maintain the distance between points:

$$\sum_{ij} W_{ij} \|y_i - y_j\|^2 = \text{tr}(Y^T L Y)$$

To avoid getting to a solution with less than d dimensions, we enforce the constraint $Y^T D Y = I$. In this way we will ensure that the rank of Y is d . We now have to solve the generalized eigenvalue problem $L\mathbf{y} = \lambda D\mathbf{y}$ and populate Y with the eigenvectors that correspond to the d smallest, non-zero eigenvalues.

Figure 3.7: *Laplacian Eigenmaps [Cay05]*.

Laplacian Eigenmaps
input: $x_1 \dots x_n \in \mathbb{R}^D$, d , k , σ .

1. Set $W_{ij} = \begin{cases} e^{-\|x_i - x_j\|^2 / 2\sigma^2} & \text{if } x_j \in N(i) \\ 0 & \text{otherwise.} \end{cases}$
2. Let U be the matrix whose columns are the eigenvectors of $Ly = \lambda Dy$ with nonzero accompanying eigenvalues.
3. Return $Y := [U]_{n \times d}$.

3.3 Discriminative manifold learning algorithms in ASR

The reason why we presented some basic ideas and algorithms coming from the area of manifold learning is the fact that speech sounds are suspected to lie on a low-dimensional manifold, which we can exploit to reduce the dimensionality of speech feature vectors. Approaches to automatic speech recognition using manifold learning is an active field of research that has so far produced promising results.

3.3.1 Manifolds and speech data

We had mentioned earlier that natural data is usually high dimensional, yet the latent factors that generate it may lie in lower-dimensional spaces. Intuitively one can see why the manifold assumption is true for certain classes of speech: given that there is only a small finite number of articulatory systems involved (tracheia, glottis, vocal tract, tongue, lips, teeth), it is reasonable to expect the existence of a low-dimensional structure in speech data.

In their work, Jansen and Niyogi ([JN05]) consider wave propagation in acoustic tube models and present a derivation of the sounds generated by it. They represent steady-state sounds produced by such models with the magnitude of their Fourier coefficients and thus, as a point in a high dimensional space. They prove that certain classes of sounds generated, as the configuration of the articulatory system varies, lie on a low dimensional manifold embedded in the ambient high-dimensional space.

We will outline the analysis in their paper ([JN05]) when a single tube is modeling the vocal tract.

In order to be able to arrive at a tractable problem they introduce the following

approximations: the vocal tract walls are rigid, their impedance is much greater than that of air and the transverse dimension of the vocal tract is much smaller than the signal's wavelength. Based on these assumptions the problem of acoustically analysing the vocal tract resonator comes down to the analysis of planar waves in one spatial dimension.

The notation used is:

- p is the air pressure
- ρ is the density
- u is the velocity
- $A(x)$ is the cross-section area of the tube as a function of the position
- $U = Au$ is the volume velocity of the air perturbations

Starting from the continuity equation for compressible fluid flow and the preservation of momentum, we can arrive at a partial differential equation with respect to the volume velocity U as a function of space and time, which can be further reduced to the standard wave equation in free space:

$$\frac{\partial^2 U}{\partial x^2} = \frac{1}{c^2} \frac{\partial^2 U}{\partial t^2}$$

or, using pressure p

$$\frac{\partial^2 p}{\partial x^2} = \frac{1}{c^2} \frac{\partial^2 p}{\partial t^2}$$

with the boundary conditions

$$U(0, t) = s(t)$$

where $s(t)$ represents the glottal vibrations as a function of time, and

$$\hat{U}(L, \omega) = \frac{\hat{p}(L, \omega)}{\mathbf{Z}_r(\omega)}$$

where the second condition is in the frequency domain and $\mathbf{Z}_r(\omega)$ is the acoustic impedance at the open end of the tube ($x = L$).

The solution $U(x, t)$ to the above equation can be expressed using its Fourier transform $\hat{U}(x, \omega)$:

$$U(x, t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \hat{U}(x, \omega) e^{i\omega t} d\omega$$

Substituting in the wave equation and solving we get the general solution for the volume and (after working in a similar way) pressure respectively:

$$\hat{U}(x, \omega) = U_1(\omega) e^{-ikx} + U_2(\omega) e^{ikx}$$

$$\hat{p}(x, \omega) = p_1(\omega)e^{-ikx} + p_2(\omega)e^{ikx}$$

The detailed computations for the magnitudes can be found in [JN05]. Eventually we arrive at the general solutions for U and p :

$$\hat{U}(x, \omega) = \hat{s}(\omega) \left[\frac{e^{-ikx}}{1 + \frac{B+1}{B-1}e^{-i2kL}} + \frac{e^{ikx}}{1 + \frac{B-1}{B+1}e^{i2kL}} \right]$$

$$\hat{p}(x, \omega) = \hat{s}(\omega) \mathbf{Z}_r(\omega) \left[\frac{e^{-ikx}}{1 + \frac{B+1}{B-1}e^{-i2kL}} + \frac{e^{ikx}}{1 + \frac{B-1}{B+1}e^{i2kL}} \right]$$

where $B = \frac{\rho_0 c}{A \mathbf{Z}_r(\omega)}$, ρ_0 is the equilibrium density air value.

Since any odd periodic function can be expressed as a Fourier series of sinusoidal sources, we pick a single-frequency sinusoidal as the source function:

$$s(t) = U_0 \sin \omega_0 t$$

If we compute its Fourier transform, substitute in the previous solutions and take the inverse Fourier, we arrive at the final form of the solution in the time-domain:

$$p(L, t) = \frac{U_0}{2i} \left[f(\omega_0, L, A) e^{i\omega_0 t} - f(-\omega_0, L, A) e^{-i\omega_0 t} \right] = \text{Im}(U_0 f(\omega_0, L, A) e^{i\omega_0 t})$$

where $f(\omega, L, A) = \mathbf{Z}_r(\omega) \left[\cos kL - i \frac{A \mathbf{Z}_r(\omega)}{\rho_0 c} \sin kL \right]^{-1}$ and we arrived at the second form given that $\mathbf{Z}_r(\omega) = \mathbf{Z}_r^*(-\omega)$ and consequently $f^*(\omega, L, A) = f(-\omega, L, A)$.

If we use a linear combination of H harmonics as the source function

$$s(t) = \sum_{n=1}^H a_n \sin(n\omega_0 t)$$

then at the output we expect a solution of the form

$$p(L, t) = \sum_{n=1}^H \beta_n \sin(n\omega_0 t + \phi_n)$$

where for each n

$$\beta_n \sin(n\omega_0 t + \phi_n) = \text{Im}(\alpha_n f(n\omega_0, L, A) e^{in\omega_0 t})$$

Consequently, the output Fourier coefficients β_n are given by

$$\beta_n(L, A) = \alpha_n |f(n\omega_0, L, A)|$$

If we now define a subset of \mathcal{R}^H for a given set $\alpha_i | i = 1, \dots, H$ by

$$\mathcal{M}_1(L_1, L_2) = (\beta_1, \beta_2, \dots, \beta_H) | L \in (L_1, L_2)$$

then \mathcal{M} traces-out a one-dimensional curve in the ambient Fourier space \mathcal{R}^H with the following properties:

- There exists a diffeomorphism $\phi : (L_1, L_2) \in \mathcal{R}^1 \rightarrow \mathcal{M}_1(L_1, L_2) \in \mathcal{R}^H$ for L_1, L_2 in the range of human vocal tract lengths.
- The diffeomorphism $\phi^{-1} : \mathcal{M}_1(L_1, L_2) \rightarrow \mathcal{R}^1$ is a coordinate chart on the set $\mathcal{M}_1(L_1, L_2)$.
- The set $\mathcal{M}_1(L_1, L_2)$ is open.

These properties tell us that the set $\mathcal{M}_1(L_1, L_2)$ is a smooth and open one dimensional manifold embedded in \mathcal{R}^H , which is not a straight line and spans the whole ambient space.

The manifold structure of data can be exploited in several ways: improve supervised learning techniques by applying regularization based on the manifold of speech data (*manifold regularization*) and find new methods of dimensionality reduction. One such method will be presented next.

3.3.2 Discriminative manifold learning

[TR14][TR13a]

In the context of automatic speech recognition we usually want to classify our data into the corresponding phoneme categories, e.g. biphones or triphones. It has recently been established ([SR03]) that the geometric and local structure of the space the data lie on, is important for classification tasks. Consequently, if we could identify the underlying manifold of the available data and apply a discriminative feature transformation, we should see an improvement in our classification efforts.

However, discriminative transformations alone are incapable of capturing the geometric and local distributional structure of the data space, whereas on the other hand, manifold learning algorithms that do uncover the geometric properties of the space, are unsupervised and non-discriminative. Rose and Tomar propose a new framework ([TR14]) that introduces a discriminative element in manifold learning techniques: they aim to maximize the separability between different classes while at the same time preserving the manifold-constraint relationships between data points belonging to the same class.

Their discriminative manifold learning framework starts -similarly to the already presented techniques- by embedding feature vectors into two high-dimensional graphs connecting neighborhoods of vectors, followed by an optimization of their structure according to certain constraints. These weighted, undirected graphs characterize the class-specific manifolds the feature vectors lie on. One of them, called the *intrinsic* graph, defines the relationships between same-class feature vectors

($\mathcal{G}_{int} = \{\mathbf{X}, \mathbf{W}_{int}\}$), whereas the other, called *penalty* graph, defines the relationships between feature vectors belonging to different classes ($\mathcal{G}_{pen} = \{\mathbf{X}, \mathbf{W}_{pen}\}$).

The characteristics of the two graphs, such as structure, connectivity and compactness, are mostly influenced by the weights on their edges, which are calculated according to the distance between points. Therefore, intuitively we understand that the graphs express the geometry of the data space.

Based on the metrics used to calculate distance, the writers propose two different approaches, both resulting in the estimation of a projection matrix $\mathbf{P} \in R^{d \times m}$, $m \ll d$, which maximizes the sub-manifold class discrimination in the projected space while at the same time it preserves the within-sub-manifold intrinsic data relations.

Locality Preserving Discriminant Analysis

[TR12a]

In Locality Preserving Discriminant Analysis (*LPDA*) the distance metric used is the Euclidean distance. The elements of the weight matrices \mathbf{W}_{int} and \mathbf{W}_{pen} are calculated using the Gaussian heat kernel:

$$w_{ij}^{int} = \begin{cases} e^{\frac{-\|x_i - x_j\|^2}{\rho}} & \text{if } C(x_i) = C(x_j), e(x_i, x_j) = 1 \\ 0 & \text{otherwise} \end{cases}$$

$$w_{ij}^{pen} = \begin{cases} e^{\frac{-\|x_i - x_j\|^2}{\rho}} & \text{if } C(x_i) \neq C(x_j), e(x_i, x_j) = 1 \\ 0 & \text{otherwise} \end{cases}$$

where ρ is the heat kernel scale parameter, $C(x_i)$ refers to the class of vector x_i and $e(x_i, x_j) = 1$ indicates that x_i is in the near neighborhood of x_j . The neighbors of vector x_i can be discovered either by an exact or by an approximate k -nearest neighbors method, such as a range search for points lying in a ball of radius r around x_i .

In essence, the meaning of the aforementioned weights assignment is that in graph \mathcal{G}_{int} a node x_i is connected to its k_{int} -nearest neighbors having the same label, whereas in \mathcal{G}_{pen} a node x_i is connected to its k_{pen} -nearest neighbors that have a different label than itself. The values of k_{int} and k_{pen} are determined empirically.

We can define a scatter measure for a graph G as follows:

$$F_G(\mathbf{P}) = \sum_{i \neq j} \|y_i - y_j\|^2 w_{ij} = 2\mathbf{P}^T \mathbf{X}(\mathbf{D} - \mathbf{W})\mathbf{X}^T \mathbf{P}$$

where \mathbf{D} is a diagonal matrix defined as $D_{ii} = \sum_j w_{ij}$. Depending on whether we want to retain or discard graph properties, the optimal projection matrix \mathbf{P} can be obtained by minimizing or maximizing F_G .

Specifically for LPDA, one would want to reinforce the fact that points belonging to different classes are conceptually far away from each other and on the other hand strengthen the relations between points of the same class. Intuitively it follows that one should aim to maximize the scatter of the penalty graph $F_{pen}(\mathbf{P})$ and minimize the scatter of the intrinsic graph $F_{int}(\mathbf{P})$. To combine the two scatters into a single measure the writers define their ratio as a measure of class-separability and graph preservation:

$$F(\mathbf{P}) = \frac{F_p(\mathbf{P})}{F_i(\mathbf{P})} = \frac{\mathbf{P}^T \mathbf{X} (\mathbf{D}_p - \mathbf{W}_p) \mathbf{X}^T \mathbf{P}}{\mathbf{P}^T \mathbf{X} (\mathbf{D}_i - \mathbf{W}_i) \mathbf{X}^T \mathbf{P}}$$

where i denotes *intrinsic* and p *penalty*.

The optimal \mathbf{P} is the one to maximize $F(\mathbf{P})$:

$$\mathbf{P}_{LPDA} = \underset{\mathbf{P}}{\operatorname{argmax}} F(\mathbf{P}) = \underset{\mathbf{P}}{\operatorname{argmax}} \left\{ \operatorname{tr} \left((\mathbf{X} (\mathbf{D}_i - \mathbf{W}_i) \mathbf{X}^T \mathbf{P})^{-1} (\mathbf{P}^T \mathbf{X} (\mathbf{D}_p - \mathbf{W}_p) \mathbf{X}^T \mathbf{P}) \right) \right\}$$

This maximization problem can be brought down to solving the following generalized eigenvalue problem:

$$(\mathbf{X} (\mathbf{D}_p - \mathbf{W}_p) \mathbf{X}^T) \mathbf{p}_{lpda}^j = \lambda_j (\mathbf{X} (\mathbf{D}_i - \mathbf{W}_i) \mathbf{X}^T) \mathbf{p}_{lpda}^j$$

where \mathbf{p}_{lpda}^j is the j^{th} column of the transformation matrix $\mathbf{P}_{lpda} \in R^{d \times m}$ and is the eigenvector associated with the j^{th} largest eigenvalue. Therefore $\mathbf{P}_{lpda}^{optimal}$ contains the m eigenvectors associated with the m largest eigenvalues.

Correlation Preserving Discriminant Analysis

[TR13b][TR12b]

The second approach, Correlation Preserving Discriminant Analysis (CPDA), uses the cosine-correlation, i.e. normalized inner-product, as distance metric. This is motivated by the fact that models based on Euclidean distance are susceptible to noise, whereas cosine-correlation based models are expected to be less influenced by noise, given the robustness of the angles between cepstrum vectors.

Similarly to LPDA, CPDA uses two undirected graphs, $(\mathcal{G}_{int} = \{\mathbf{X}, \mathbf{W}_{int}\})$, $(\mathcal{G}_{pen} = \{\mathbf{X}, \mathbf{W}_{pen}\})$, to embed the feature vectors, yet this time the similarity between nodes is calculated based on the cosine-correlation distance:

$$w_{ij}^{int} = \begin{cases} e^{\frac{\langle x_i, x_j \rangle - 1}{\rho}} & \text{if } C(x_i) = C(x_j), e(x_i, x_j) = 1 \\ 0 & \text{otherwise} \end{cases}$$

$$w_{ij}^{pen} = \begin{cases} e^{\frac{\langle x_i, x_j \rangle - 1}{\rho}} & \text{if } C(x_i) \neq C(x_j), e(x_i, x_j) = 1 \\ 0 & \text{otherwise} \end{cases}$$

where the notation is the same as above. The scatter measure used in CPDA is the same as in LDPA, yet CPDA first projects the data points onto the surface of a unit hypersphere, thus discarding magnitude information and keeping correlation-based information about data points. Thus, the projection matrix $\mathbf{P} \in R^{d \times m}$, $m \ll d$, eventually estimated, will non-linearly project the features from the original d -dimensional hypersphere onto the target m -dimensional hypersphere.

The scatter for CPDA is:

$$F_G(\mathbf{P}) = \sum_{i \neq j} \|y_i - y_j\|^2 w_{ij} \stackrel{\mathbf{y} = \mathbf{P}^T \mathbf{x}}{=} \sum_{i \neq j} \left\| \frac{\mathbf{P}^T \mathbf{x}_i}{\|\mathbf{P}^T \mathbf{x}_i\|} - \frac{\mathbf{P}^T \mathbf{x}_j}{\|\mathbf{P}^T \mathbf{x}_j\|} \right\|^2 w_{ij} = 2 \sum_{i \neq j} \left(1 - \frac{f_{ij}}{f_i f_j} \right) w_{ij}$$

where for two arbitrary vectors x_u, x_v :

$$f_u = \sqrt{x_u^T \mathbf{P} \mathbf{P}^T x_u}$$

and

$$f_{uv} = \sqrt{x_u^T \mathbf{P} \mathbf{P}^T x_v}$$

As previously, we want to maximize the scatter of the penalty graph and minimize the scatter of the intrinsic graph. However, in CPDA the writers have chosen the difference of the two scatter measures as a measure of manifold separability and graph preservation:

$$F(\mathbf{P}) = F_{pen}(\mathbf{P}) - F_{int}(\mathbf{P}) = 2 \sum_{i \neq j} \left(1 - \frac{f_{ij}}{f_i f_j} \right) w_{ij}^{p-i}$$

where $w_{ij}^{p-i} = w_{ij}^{pen} - w_{ij}^{int}$. The optimal projection matrix \mathbf{P}_{cpda}^{opt} is obtained by maximizing $F(\mathbf{P})$:

$$\mathbf{P}_{CPDA} = \underset{\mathbf{P}}{\operatorname{argmax}} F(\mathbf{P})$$

Because of the projection of the feature vectors onto the unit hypersphere, the maximization problem cannot be solved by solving an equivalent generalized eigenvalue problem. Therefore, an iterative procedure is chosen, the well known gradient ascent:

$$\mathbf{P} := \mathbf{P} + \alpha \nabla_{\mathbf{P}} F$$

where

$$\nabla_{\mathbf{P}} F = 2 \sum_{i \neq j} \left[\frac{f_{ij} x_i x_i^T}{f_i^3 f_j} + \frac{f_{ij} x_j x_j^T}{f_i f_j^3} - \frac{x_i x_j^T + x_j x_i^T}{f_i f_j} \right] \mathbf{P} w_{ij}^{p-i}$$

where α is the gradient scaling factor and $\nabla_{\mathbf{P}} F$ represents the gradient of the scatter measure F with respect to \mathbf{P} . Given that F is non-linear and non-convex, a good initialization is important to avoid reaching just a local maximum. Such an

initialization can be obtained by neglecting the projection on the unit hypersphere, approximate a linear transformation:

$$\mathbf{y}_i = \mathbf{P}^T \mathbf{x}_i$$

and eventually acquire \mathbf{P}_{init} by solving the generalized eigenvalue problem:

$$(\mathbf{X}(\mathbf{D}_p - \mathbf{W}_p)\mathbf{X}^T)\mathbf{p}^j = \lambda_j(\mathbf{X}(\mathbf{D}_i - \mathbf{W}_i)\mathbf{X}^T)\mathbf{p}^j$$

where \mathbf{D} is a diagonal matrix with elements $D_{ii} = \sum_j w_{ij}$ and \mathbf{X} contains the normalized feature vectors. As before, \mathbf{P}_{cpda}^{init} is composed of the eigenvectors corresponding to the m largest eigenvalues.

Noise Aware Manifold Learning

[TR13c]

Furthermore, in their work Rose and Tomar explore a new method to increase noise robustness of manifold learning methods. They point out that since edge weights, which depend on the Gaussian heat kernel size and shape, characterize local neighborhoods, there is a strong relation between kernel size and the robustness to background noise.

The kernel size governs the compactness of the neighborhood and the smoothness of the manifolds, and in turn, it is governed by the scale parameter ρ . Therefore, the choice of ρ is important to the definition of local neighborhoods and thus to the characteristics of the transformation. If ρ had a value that is too large, then the kernel would tend to flatten, and all data pairs on the graph would have the same importance; on the other hand, if it had a too small value, the manifold would lack smoothing and thus the kernel would be too sensitive to noise. Therefore, the writers claim that the optimal value for ρ is dependent on the SNR level of the speech signal, and support their claims with successful experimental results.

They propose that multiple scale parameters be used, each specific to a noise condition. First, during training, they gather multiple SNR dependent LPDA or CPDA projections and Continuous Density HMMs. This procedure is broken down into three steps:

- Based on a set of ρ values, determine the set of the corresponding projection matrices
- Heuristically identify the optimal value of ρ and the corresponding transformation that maximizes ASR performance for a given SNR level

- Train CDHMM models from the features obtained using the set of projection matrices previously acquired

Having this set of SNR-dependent LPDA/CPDA transformations, SNR is estimated for each test utterance and the feature space projection is performed using the corresponding LPDA/CPDA projection matrix. Eventually, the corresponding CDHMM model is used in the ASR application.

This procedure is referred to as Noise-Aware LPDA/CPDA (*N-LPDA* or *N-CPDA*).

Results

All of the proposed techniques are verified by experiments the writers conducted using the Aurora-2 and Aurora-4 datasets.

LPDA and CPDA show a relative WER improvement ranging from 6 to 30% compared to conventional techniques, such as Linear Discriminant Analysis or Locality Preserving Projections. This verifies the assertion that manifold and discriminative learning result in a well-behaved and robust feature-space transformation.

Furthermore, CPDA performs better than LPDA under high noise conditions, supporting the noise robustness of cosine-correlation distance compared to Euclidean. As far as Noise-Aware manifold learning is concerned, N-LPDA performs slightly better on average as compared to LPDA for most noise conditions (0.10-0.20 reduced WER).

At this point it is important to point out the huge computational complexity involved in computing the weight matrices, given the high dimensionality of the feature space. To tackle this issue the writers propose using novel approximate nearest-neighbors techniques such as Locality Sensitive Hashing ([TR13b]).

Chapter 4

Deep Neural Networks

In this chapter we will give a general introduction to deep neural networks (*DNNs*). We will begin by presenting the ideas and motivation behind deep architectures and then move on to describe a few of these models as well as practical issues concerning building and training such models.

4.1 History and Motivation

[DY14]

Until recently, shallow-structured architectures were prevalent in machine learning and signal processing. These include Gaussian Mixture Models, Conditional Random Fields, Support Vector Machines, logistic regression and the well-known multilayer perceptron, all of which have been successful in solving well-constrained problems, yet their limited representational power does not allow them to effectively deal with applications involving natural signals such as human speech, natural sound and image.

In order to deal with such tasks efficiently, researchers had to improve the information processing mechanisms involved in the systems and they turned to the human brain for inspiration. Human information processing mechanisms (such as vision and audition) suggest the need for “deep”, hierarchical architectures ([SKK⁺07]) to capture non-linear, complex structure and build internal distributed representations from sensory inputs. Researchers however, recently have been stressing out that deep learning, i.e. the process of using deep architectures to extract information from data, is only *inspired* by the way the brain works, rather than mimicking it. For one thing, the human brain is full of different types of neurons each one serving a different purpose, whereas deep architectures have a single type of neuron performing the same task. Moreover, humans learn mostly in an unsupervised manner, contrary to deep models which usually need labeled examples to learn from.

Although deep architectures recently gained in popularity, the main idea behind them, that is to learn lower-level representations from raw input data, is not new. As early as 1998, Yann LeCun et al. ([LBBH98]) developed a Convolution Neural Network to recognize handwritten digits on banking checks. The novelty of LeNet-5 (the model proposed in the paper) was that it learned features from raw pixel images, as opposed to receiving hand-crafted features in the input end. The power of deep models lies exactly at that point: they can extract features and do discrimination jointly ([Ben09]).

However, at that time, it was not easy to implement a deep architecture onto , for example, a feed-forward neural network. The back-propagation algorithm, which had been very popular until then, did not work well in practice when more than a small number of hidden layers were involved. Given the presence of too many local optima in the non-convex objective function of such models, as well as the problems of diminishing or exploding gradients and the computational issues arising when the parameters' number rose too much, most research turned to shallow models with convex loss functions (e.g. SVMs, CRFs). It was not until recently that new theoretical breakthroughs appeared (e.g. RBM pretraining, primal-dual update algorithm, new non-linear activation functions) and together with the boom of computationally efficient hardware usage (GPUs, CPU/GPU clusters) allowed researchers to turn to deep (many hidden layers) and wide (many neurons) models.

4.2 Constructing DNNs: the Multilayer Perceptron

[YD14]

Although initially the term “deep neural network” was used to identify a multilayer perceptron with many hidden layers, it has now come to name a whole class of neural network models having a “deep” architecture, i.e. having at least three hidden layers. However, we will use the conventional multilayer perceptron to present the basic deep architecture and cover the training issues associated with it.

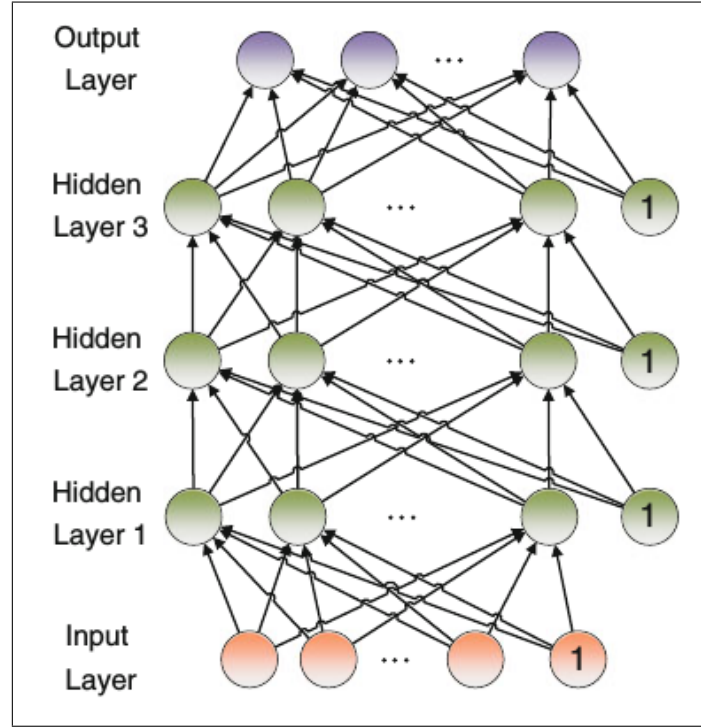
4.2.1 Architecture of the Deep Multilayer Perceptron

As already mentioned, the simplest deep neural network is a multilayer perceptron with at least three hidden layers, that is, excluding the input and output layers. A DNN with a total of five layers can be seen in figure 4.1

We can denote the input layer as layer 0 and the output layer as layer L for an $L + 1$ -layer DNN.

The neurons in the first layer correspond to each feature (or dimension) in the

Figure 4.1: A deep neural network with three hidden layers [YD14].



input vector, i.e.

$$\mathbf{v}^0 = \mathbf{x}_{input}$$

whereas for the first L layers the output is calculated as follows:

$$\mathbf{v}^l = f(\mathbf{z}^l) = f(\mathbf{W}^l \mathbf{v}^{l-1} + \mathbf{b}^l), \quad 0 < l < L$$

where

$$\mathbf{z}^l = \mathbf{W}^l \mathbf{v}^{l-1} + \mathbf{b}^l \in \mathcal{R}^{N_l \times 1}$$

$$\mathbf{v}^l \in \mathcal{R}^{N_l \times 1}$$

$$\mathbf{W}^l \in \mathcal{R}^{N_l \times N_{l-1}}$$

$$\mathbf{b}^l \in \mathcal{R}^{N_l \times 1}$$

and $N_l \in \mathcal{R}$, are respectively the excitation vector, the activation vector, the weight matrix, the bias vector and the number of neurons at layer l .

The function $f(\cdot) : \mathcal{R}^{N_l \times 1} \rightarrow \mathcal{R}^{N_l \times 1}$ is the activation function of the corresponding neuron that is applied element-wise to the excitation vector. The most common activation function is the sigmoid function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

followed by the second most common, the hyperbolic tangent

$$\tanh(z) = \frac{e^z + e^{-z}}{e^z - e^{-z}}$$

which is a rescaled version of the former. Both have the same modeling power, but the output range of the hyperbolic tangent is $(-1, +1)$, therefore the activation values are symmetric and this was believed to help during training. On the other hand the sigmoid function has an output range of $(0, 1)$, which stimulates sparse but asymmetric activation values. Recently, another type of activation function, the rectified linear unit (ReLU) has been suggested:

$$ReLU(z) = \max(0, z)$$

ReLU has been gaining in popularity since it enforces sparse activation values and has a simple gradient:

$$\nabla_{ReLU}(z) = \max(0, \text{sgn}(z))$$

As far as the output layer is concerned, its form depends on the task of interest.

In case of regression tasks, a linear layer is used to generate the output vector $\mathbf{v}^L \in \mathcal{R}^{N_l}$, where N_l is the output dimension:

$$\mathbf{v}^L = \mathbf{z}^L = \mathbf{W}^L \mathbf{v}^{L-1} + \mathbf{b}^L$$

In case of classification tasks, each output neuron represents a class $i \in 1, \dots, C$, where $C = N_L$ is the number of classes. In this case, the output of the DNN will be a probability distribution over the possible classes. The value of the i^{th} output neuron represents the probability $P_{dnn}(i|\mathbf{x}_i)$ that the input vector belongs to class i .

Since we want the output vector to be a valid probability distribution we impose the following constraints on it

$$v_i^L \geq 0$$

and

$$\sum_{i=1}^C v_i^L = 1$$

and use as output function the softmax function:

$$v_i^L = P_{dnn}(i|\mathbf{x}_i) = \text{softmax}(\mathbf{z}^L) = \frac{e^{z_i^L}}{\sum_{j=1}^C e^{z_j^L}}$$

where z_i^L is the i^{th} element of the excitation vector \mathbf{z}^L .

Given an input matrix with observations \mathbf{O} the output of the DNN can be computed as described in figure 4.2

Figure 4.2: *Computation of DNN output [YD14].*

```

1: procedure FORWARDCOMPUTATION(O)
2:    $\mathbf{V}^0 \leftarrow \mathbf{O}$ 
3:   for  $\ell \leftarrow 1; \ell < L; \ell \leftarrow \ell + 1$  do
4:      $\mathbf{Z}^\ell \leftarrow \mathbf{W}^\ell \mathbf{V}^{\ell-1} + \mathbf{B}^\ell$ 
5:      $\mathbf{V}^\ell \leftarrow f(\mathbf{Z}^\ell)$ 
6:   end for
7:    $\mathbf{Z}^L \leftarrow \mathbf{W}^L \mathbf{V}^{L-1} + \mathbf{B}^L$ 
8:   if regression then
9:      $\mathbf{V}^L \leftarrow \mathbf{Z}^L$ 
10:  else
11:     $\mathbf{V}^L \leftarrow \text{softmax}(\mathbf{Z}^L)$ 
12:  end if
13:  Return  $\mathbf{V}^L$ 
14: end procedure

```

▷ Each column of \mathbf{O} is an observation vector
 ▷ L is the total number of layers
 ▷ Each column of \mathbf{B}^ℓ is \mathbf{b}^ℓ
 ▷ $f(\cdot)$ can be sigmoid, tanh, ReLU, or other functions
 ▷ regression task
 ▷ classification task
 ▷ Apply softmax column-wise

4.2.2 Training and related issues

Given that a multilayer perceptron with a sufficiently large hidden layer can approximate any function ([HSW89]), it is evident that a multilayer perceptron with more than one hidden layers can serve as a universal approximator as well, that is, compute any mapping

$$g : \mathcal{R}^D \rightarrow \mathcal{R}^C$$

from the input space \mathcal{R}^D to the output space \mathcal{R}^C . Given a set of training samples we can estimate the network parameters $\{\mathbf{W}, \mathbf{b}\}$ with a training process, characterized by a training criterion and a learning algorithm.

Training criterion

The training criterion is chosen so that an improvement in it reflects an improvement in the final evaluation score. Therefore, it should somehow express the final goal of the task, yet at the same time be easy to evaluate.

Ideally, we would like to estimate the DNN parameters that minimize the expected cost:

$$J_E = \mathbb{E}(J(\mathbf{W}, \mathbf{b}; \mathbf{x}_i, \mathbf{y}))$$

where J is the cost (or loss) function given the model parameters $\{\mathbf{W}, \mathbf{b}\}$ together with the input vector \mathbf{x}_i and the corresponding output \mathbf{y} . However, since we cannot know in advance the output for unseen samples, we can only optimize empirical criteria based on the training set.

The two most popular empirical criteria used in DNN training are the mean square error (MSE) for regression tasks and the cross-entropy (CE) for classification

Figure 4.3: *Back-propagation algorithm [YD14].*

```

1: procedure BACKPROPAGATION( $\mathbb{S} = \{(\mathbf{o}^m, \mathbf{y}^m) \mid 0 \leq m < M\}$ )
     $\triangleright \mathbb{S}$  is the training set with  $M$  samples
2:   Randomly initialize  $\{\mathbf{W}_0^\ell, \mathbf{b}_0^\ell\}, 0 < \ell \leq L$ 
     $\triangleright L$  is the total number of layers
3:   while Stopping Criterion Not Met do
     $\triangleright$  Stop if reached max iterations or the training criterion improvement is small
4:     Randomly select a minibatch  $\mathbf{O}, \mathbf{Y}$  with  $M_b$  samples.
5:     Call ForwardComputation( $\mathbf{O}$ )
6:      $\mathbf{E}_t^L \leftarrow \mathbf{V}_t^L - \mathbf{Y}$ 
     $\triangleright$  Each column of  $\mathbf{E}_t^L$  is  $\mathbf{e}_t^L$ 
7:      $\mathbf{G}_t^L \leftarrow \mathbf{E}_t^L$ 
8:     for  $\ell \leftarrow L; \ell > 0; \ell \leftarrow \ell - 1$  do
9:        $\nabla_{\mathbf{W}_t^\ell} \leftarrow \mathbf{G}_t^\ell (\mathbf{v}_t^{\ell-1})^T$ 
10:       $\nabla_{\mathbf{b}_t^\ell} \leftarrow \mathbf{G}_t^\ell$ 
11:       $\mathbf{W}_{t+1}^\ell \leftarrow \mathbf{W}_t^\ell - \frac{\epsilon}{M_b} \nabla_{\mathbf{W}_t^\ell}$ 
     $\triangleright$  Update  $\mathbf{W}$ 
12:       $\mathbf{b}_{t+1}^\ell \leftarrow \mathbf{b}_t^\ell - \frac{\epsilon}{M_b} \nabla_{\mathbf{b}_t^\ell}$ 
     $\triangleright$  Update  $\mathbf{b}$ 
13:       $\mathbf{E}_t^{\ell-1} \leftarrow (\mathbf{W}_t^\ell)^T \mathbf{G}_t^\ell$ 
     $\triangleright$  Error backpropagation
14:      if  $\ell > 1$  then
15:         $\mathbf{G}_t^{\ell-1} \leftarrow f'(\mathbf{Z}_t^{\ell-1}) \bullet \mathbf{E}_t^{\ell-1}$ 
16:      end if
17:    end for
18:  end while
19:  Return  $dnn = \{\mathbf{W}^\ell, \mathbf{b}^\ell\}, 0 < \ell \leq L$ 
20: end procedure

```

tasks, which are defined as follows:

$$J_{MSE}(\mathbf{W}, \mathbf{b}; \mathbf{x}_{i,train}, \mathbf{y}_{o,train}) = \frac{1}{2M} \sum_{m=1}^M \|\mathbf{v}^L - \mathbf{y}\|^2$$

and

$$J_{CE}(\mathbf{W}, \mathbf{b}; \mathbf{x}_{i,train}, \mathbf{y}_{o,train}) = \frac{1}{M} \sum_{m=1}^M \left(- \sum_{i=1}^C y_i \log(v_i^L) \right)$$

where y_i and v_i in J_{CE} are respectively, the probability observed in the training set that \mathbf{x} belongs to class i ($y_i = P_{empirical}(i|\mathbf{x})$) and the same probability estimated from the DNN $v_i^L = P_{DNN}(i|\mathbf{x})$.

The minimization of the cross-entropy criterion is equivalent to minimizing the the Kullback-Leibler divergence between the empirical probability distribution and the probability distribution estimated from the deep neural network.

Training algorithm and practical considerations

Having selected an appropriate training criterion, we can train the network using the well known back-propagation algorithm which is based on the chain rule for gradient computation. The algorithm is summarized in figure 4.3.

Despite the simplicity of the algorithm, its application to DNNs gives rise to many problems, the most important ones being the speed of convergence and the vanishing or exploding gradients that get in the way of the training. It is therefore

crucial to address these issues by specifically tuning the DNN and the training set to avoid such drawbacks.

Data preprocessing Data preprocessing is an important part of the pre-training phase in the DNN construction. There are two principal ways of preprocessing each one attacking a different problem that might arise during training.

The first preprocessing technique is referred to as per-sample feature normalization and it results in reducing the variability in the final feature presented to the DNN. This is desired since we want to drop any variation that is not important to the final decision made by the model and would otherwise add extra difficulties both at the output stage as well as during training, as the network would try to learn too much redundant information. For example, in image recognition, it is desired to reduce the variability introduced by the brightness, or in speech recognition we want to avoid acoustic channel distortions. In this last case we can achieve our goal by calculating the per-utterance mean for each feature and subtracting it from all frames in the utterance (cepstral mean normalization).

The second preprocessing technique is the global feature standardization and it aims to scale the data along each dimension with a global transformation so that the final data vectors lie in a similar range of numerical values. The global transformation is estimated from the training data and is then applied to both training and test sets. Global feature standardization will improve performance in later steps of the training process. As far as DNN training is concerned, it allows us to use the same learning rate across all weight dimensions and still achieve good performance.

A common global transformation in speech recognition is to standardize each dimension of the feature vector to have zero mean and unit variance. If we use MFCC vectors without standardization, the energy component will dominate the learning procedure since its values are much greater than the rest of the coefficients.

Initialization of model parameters The back-propagation algorithm starts given a set of initial network parameters. Since the DNN is highly non-linear and the training criterion with respect to the model parameters is non-convex, initialization of the model parameters can greatly affect the end model. There are many suggestions for the initialization of a DNN model, all of which are based on two assumptions.

First, the weights should be initialized so that each neuron operates in the linear range of the activation function (e.g. sigmoid) at the start of the learning. If the weights were very large, neurons would saturate, i.e. take an extreme value (close to zero or one in the case of a sigmoid function) leading to vanishing gradients. This

becomes evident in the case of the sigmoid as we can observe from its derivative:

$$\sigma'(\mathbf{z}_t^l) = (1 - \sigma(\mathbf{z}_t^l))\sigma(\mathbf{z}_t^l) = (1 - \mathbf{v}_t^l)\mathbf{v}_t^l$$

On the other hand, if the neurons operate in the linear range, the gradients are large enough and learning can proceed smoothly. It is important to point out that the excitation vector depends on both the weights and the input values. If we have standardized the input values in the preprocessing stage, then the weight initialization can become easier.

Second, the weights should be initialized at random. Since neurons in DNNs are symmetric and interchangeable, random initialization serves the purpose of symmetry breaking. If all neurons were initialized with the same values then they would have the same output thus detect the same patterns in the lower layers.

Bengio et. al. ([GB10]) have suggested that the weight values should be uniformly sampled from an interval depending on the activation function. For the sigmoid function they suggested the interval

$$\left[-4\sqrt{\frac{6}{fan_{in} + fan_{out}}}, 4\sqrt{\frac{6}{fan_{in} + fan_{out}}} \right]$$

where fan_{in} is the number of units in the $(i-1)^{th}$ layer and fan_{out} is the number of units in the i^{th} layer. This approach is based on the idea that neurons with more inputs should have smaller weights ([Ben12]).

LeCun et. al. ([LBOM98]) suggested to draw the initial weight values for layer l from a zero-mean Gaussian distribution with standard deviation $\sigma_{\mathbf{W}^{l+1}} = \frac{1}{\sqrt{N_l}}$, where N_l is the number of connections feeding into the node.

When it comes to speech recognition applications, usually each DNN layer has 1000-2000 neurons and initializing the weights either from a Gaussian distribution $\mathcal{N}(w; 0, 0.05)$ or from a uniform distribution in the range of $[-0.05, 0.05]$ we can achieve good results.

As far as the bias vectors \mathbf{b}^l are concerned, they can be initialized to zero.

Regularization Taking into consideration the huge number of parameters in a DNN it is easy to understand the risk of overfitting the model to the training data. Overfitting occurs because although we aim to minimize the expected loss, we actually minimize the empirical loss defined on the training set. To control overfitting we try to affect the model parameters through the training function so that they do not fit the training data too well. This process is known as regularization or weight decay. Regularization aims to help the model learn patterns often seen in the training data yet not learn every noise peculiarity present in the training set ([Nie15]).

The most common forms of regularization are based on the L_1 and L_2 norms respectively:

$$R_1(\mathbf{W}) = \|\text{vec}(\mathbf{W})\|_1 = \sum_{l=1}^L \|\text{vec}(\mathbf{W}^l)\|_1$$

$$R_2(\mathbf{W}) = \|\text{vec}(\mathbf{W})\|_2^2 = \sum_{l=1}^L \|\text{vec}(\mathbf{W}^l)\|_2^2$$

Including the regularization terms the training criterion becomes

$$J_R(\mathbf{W}, \mathbf{b}; \mathcal{S}_{train}) = J(\mathbf{W}, \mathbf{b}; \mathcal{S}_{train}) + \lambda R(\mathbf{W})$$

where λ is an interpolation weight called regularization weight. Regularization is particularly helpful when the training set size is small compared to the number of parameters of the DNN. For DNNs used in speech recognition, there are usually millions of parameters in the model and the regularization weight should be small ($\approx 10^{-4}$) or zero when the training set is large.

In Bayesian context, the regularization term is the negative log-prior $-\log P(\theta)$ on the parameters θ . Adding the regularization term makes the training criterion correspond to the negative joint likelihood of training data and parameters

$$-\log P(\text{data}, \theta) = -\log P(\text{data}|\theta) - \log P(\theta)$$

with the loss function $L(z, \theta)$ being interpreted as $-\log P(z|\theta)$ and $-\log P(\text{data}|\theta) = -\sum_{t=1}^T L(z_t, \theta)$, where z_t , $t = 1, \dots, T$ are the training data. In case we are training our model with a stochastic gradient based method, as is usually the case, we want to use an unbiased estimation of the gradient of the total training criterion including both the loss function and the regularization term. In a mini-batch or online learning approach it is not trivial to include the regularization in the above sum, since it is different from the sum of the loss function on all examples and the regularization term. In these cases, the regularization term should be properly weighted not only with λ but also with the inverse of the number of updates needed to see all the training data ([Ben12]).

Another popular regularization method is dropout. The basic idea behind dropout is to randomly omit a certain percentage (e.g. p) of the neurons in each hidden layer for each presentation of the samples during training. This would mean that each random combination of the remaining neurons needs to perform well during training, which implies that each neuron depends less on other neurons to detect patterns.

One way to look at dropout is as a technique to add random noise to the training data, since each neuron takes input from a random collection of the neurons in the previous layer. Consequently, the excitation value will be different even if the same

input is fed to the DNN. The improvement in generalization comes from the fact that the capacity of the DNN is reduced since some weights must be dedicated to removing the noise inserted by dropout.

A different perspective of dropout is as a *bagging* technique. *Bagging* refers to the combination of different classifiers, each one performing better than the rest at identifying a particular pattern in the data. The resemblance of dropout to bagging comes from the fact that omitting neurons creates in essence a different layer, thus a different classifier, during each feed-forward pass. After multiple passes it is as if we are averaging the outputs of these “different” classifiers to produce the final output of our model.

It should be noted that regularization affects only the weights and not the biases of the layers. This is because having a large bias does not make the network as sensitive to its input as having large weights would.

Batch size The term *batch size* refers to the number of samples we observe before we make an update to the model parameters. This number will affect the convergence speed as well as the end model.

The simplest choice of batch size is the whole training set. This approach, also known as batch training, has several advantages: it has good convergence properties, there are many techniques that can speed it up like conjugate gradient ([HS52]) and L-BFGS ([LN89]) and it allows for easy parallelization. On the other hand, seeing the whole dataset before a parameter update, is prohibitive for large datasets, which are often met in speech processing tasks.

Another approach to setting the batch size is to set it equal to one. In this approach, commonly referred to as stochastic gradient descent (*SGD*) or online learning, we update the model parameters after seeing just a single sample from the training set. If the sample is independently and identically distributed, it can be shown that the gradient of the training criterion estimated from it is an unbiased estimation of the gradient on the whole training set, despite the fact that it is noisy. This noise is the advantage of the SGD over the batch training approach. Since deep neural networks are highly non-linear and non-convex, the loss function contains many local minima, many of which are poor. Batch learning will converge to the minimum of whatever basin the initial model parameters are in and therefore it makes the model highly dependent on its initialization. On the other hand, the noise present in the gradient estimation of the SGD approach can help the algorithm move away from a poor local minimum and go closer to a better one. This property is reminiscent of simulated annealing, which allows the model parameters to move in a direction locally poorer but globally better ([KGV83]).

Moreover, SGD often converges much faster than batch learning, especially in

large datasets. This is because first, it does not waste time in redundant samples in the dataset as batch learning does, and second because in each update the new gradient is estimated based on the new model rather than the old. This means that SGD moves faster towards the best model.

However, SGD cannot be easily parallelized and it will fluctuate around the local minimum due to the noise present in the gradient estimation, without fully converging to it. This fluctuation, can be sometimes useful, as it reduces overfitting.

The third approach to batch size is to set it to a number, usually small, between one and the training set size. This approach, called minibatch training, estimates the gradient based on a small number (a batch) of *randomly* drawn training samples. It must be noted that drawing the samples at random is crucial to the minibatch and the SGD approaches to training.

Similar to SGD, the minibatch-estimated gradient is also unbiased, yet the variance of the estimation is smaller than that in the SGD algorithm. Furthermore, we can easily parallelize the computations within the minibatch, which makes this approach converge faster than SGD.

Care should be taken as to how to choose the appropriate minibatch size, yet fortunately it can be chosen independently of the rest of the learning parameters. It can be determined automatically considering the variance of the gradient estimation, or it can be empirically determined by searching on a small subset of samples in each epoch. As far as speech recognition tasks are concerned, a minibatch size of 64-256 is good for the early stages of the training, which can be expanded to 1024-8096 in later stages.

Momentum In neural network training adding a momentum term to the update of the model parameters is known to improve the convergence speed. This is due to the fact that the parameters' updates will be based on all the previous gradients instead of only the current one and the oscillation problems commonly seen in the back-propagation algorithm will be reduced.

Learning rate and stopping criterion One of the most difficult to choose parameters in DNN training, yet maybe the most crucial, is the learning rate, as it can determine whether the training will convergence to an optimum, as well as the convergence speed. Most learning rate strategies choose an appropriate initial learning rate and modify it during training. The modifications are based on a strategy usually defined empirically.

A common practice to determine the initial optimal learning rate is to set it to the largest value that does not make the training criterion diverge. Consequently, one can start with a large value, observe the development of the training criterion and if it diverges, try setting the learning rate to a value three times smaller than

the initial and proceed likewise until there is no divergence ([Ben12]).

In the case of batch learning, an empirical suggestion for a learning strategy is to decrease the learning rate when the weight vectors oscillate and increase it when the weights follow a relatively steady course. However, this approach cannot be applied to SGD learning since the weights continuously fluctuate ([LBOM98]).

It has been proven that when the learning rate is set to

$$\epsilon = \frac{c}{t}$$

where t is the number of samples seen by the network and c is a constant, or generally

$$\epsilon_t = \frac{\epsilon_0 \tau}{\max(t, \tau)}$$

which maintains a constant learning rate for the first τ steps and then decreases it in $\mathcal{O}(\frac{1}{t})$, SGD converges asymptotically, although the convergence will be slow since ϵ will quickly become too small ([YD14],[LBOM98]).

Network Architecture Network architecture, meaning the number of hidden layers and the number of neurons in each layer, is an important model parameter that can affect the computational complexity and consequently the training speed, as well as the performance of the end model itself.

Each layer in the network works in essence as a feature extractor of the previous layer. Consequently, one must make sure that the number of neurons in it is large enough to capture the patterns of interest of the input vector. This becomes especially important in the case of the lower layers, since lower layer features are more variable, and more neurons are required to model the hidden patterns. However, having too many neurons encloses the danger of overfitting. As a rule of thumb, the wide, shallow networks are easier to overfit whereas the narrow, deep networks are easier to underfit to the training data.

Taking these into consideration, it is often the case that a different number of neurons is used in each layer, with wide layers with many neurons being placed close to the network input and narrow layers with less neurons being placed close to the output. If however all layers have the same number of neurons, simply adding more layers in some cases might cause the network to underfit since the extra layers add extra constraints to the model parameters.

On the other hand, Bengio ([Ben12]) reports that networks using same-size layers perform better or the same as using decreasing or increasing size layers. He mentions though, that this could depend on the data used.

Regarding speech recognition tasks, it has been found that 5-7 layer DNNs with 1000-3000 neurons per layer work well.

4.3 Pretraining

Considering the above discussion and the difficulties in DNN training, it must have already been clear that the initialization of the model parameters is of utmost importance to the performance of the network. We will now briefly examine some pretraining techniques which aim at properly initializing a DNN to facilitate training and improve performance.

Generative pretraining

The idea of generative pretraining stems from the fact that generative models are able to identify more complex relationships between data and target labels than discriminative models. Before describing this pretraining method it makes sense to introduce the generative models that it is based on.

Energy-based models[Ben09] Energy-based models are models basing their learning procedure on the idea of *energy*. Energy is a scalar function of each configuration of the variables of interest in our model. Such models learn by modifying the shape of the energy function so that desirable configurations of their variables have low energy.

The probability distribution of energy-based models assumes the following form:

$$P(x) = \frac{e^{-Energy(x)}}{Z}$$

where Z , called the *partition function*, is

$$Z = \sum_x e^{-Energy(x)}$$

If the energy function is a sum of functions f_i

$$Energy(x) = \sum_i f_i(x)$$

then the probability distribution becomes

$$P(x) \propto \prod_i P_i(x) \prod_i e^{-f_i(x)}$$

In that case, if we assume that f_i is an “expert” imposing a constraint on x , then $P(x)$ is a “product of experts” and all f_i form a distributed representation of the configuration space, where more than one expert can contribute to each region in space. In contrast, if a different energy form leads to a weighted sum of experts, i.e. $P(x)$ is a “mixture of experts”, then the experts do not create a distributed representation of the configuration space, since belonging to a region excludes all the rest.

Distributed representations of data, are internal representations of the observed data in such way that they are modeled as being explained by the interactions of many latent factors ([DY14]). Such representations provide robustness in representing the internal structure of data and are capable to generalize concepts and relations.

In case there are unobservable components in our data, or we want to introduce hidden components in order to enhance the expressive power of the model, we can inset a hidden part h in the energy function. The probability distribution then becomes

$$P(v, h) = \frac{e^{-Energy(v, h)}}{Z}$$

where

$$Z = \sum_{v, h} e^{-Energy(v, h)}$$

and

$$P(v) = \sum_h \frac{e^{-Energy(v, h)}}{Z}$$

where v is the visible component of the data. If we define the quantity of “free energy” as

$$F_e(v) = -\log\left(\sum_h e^{-Energy(v, h)}\right)$$

then the probability distribution function of the visible part of the data takes the form of an energy-based model p.d.f. without latent components:

$$P(v) = \frac{e^{-F_e(v)}}{Z}$$

and

$$Z = \sum_v e^{-F_e(v)}$$

Boltzmann Machines A Boltzmann Machine is a network of symmetrically connected units that make stochastic decisions about whether to be on or off ([DY14]).

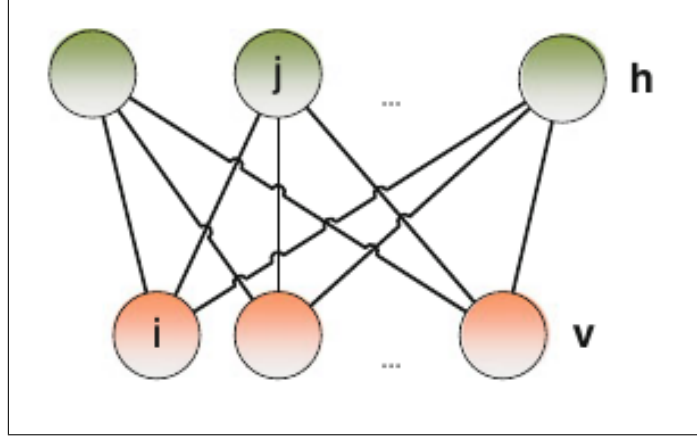
The energy function in a Boltzmann Machine is a general second-order polynomial:

$$Energy(\mathbf{v}, \mathbf{h}) = -\mathbf{a}'\mathbf{v} - \mathbf{b}'\mathbf{h} - \mathbf{h}'\mathbf{W}\mathbf{v} - \mathbf{v}'\mathbf{D}\mathbf{v} - \mathbf{h}'\mathbf{R}\mathbf{h}$$

where the offsets $\mathbf{a}_i, \mathbf{b}_i$, each associated with a single element of vector \mathbf{x} or \mathbf{h} , and the weight matrices $\mathbf{W}, \mathbf{D}, \mathbf{R}$, each associated with a pair of units are the model parameters ([DY14],[Ben09]).

Restricted Boltzmann Machines A special type of Boltzmann Machine is the Restricted Boltzmann Machine (*RBM*). It is an undirected graphical model built by a layer of stochastic visible neurons and a layer of stochastic hidden neurons, which

Figure 4.4: An example of a Restricted Boltzmann Machine [YD14].



together form a bipartite graph with no visible-visible or hidden-hidden connections (figure 4.4).

The energy function of the RBM has the same form of the energy of a Boltzmann Machine, only now $\mathbf{D} = \mathbf{R} = \mathbf{0}$ given the special architecture of the RBM.

For the Bernoulli-Bernoulli RBM, in which $\mathbf{v} \in \{0, 1\}^{N_v \times 1}$ and $\mathbf{h} \in \{0, 1\}^{N_h \times 1}$ the energy is

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{a}^T \mathbf{v} - \mathbf{b}^T \mathbf{h} - \mathbf{h}^T \mathbf{W} \mathbf{v}$$

where N_v, N_h are the number of visible and hidden neurons respectively, $\mathbf{W} \in \mathcal{R}^{N_h \times N_v}$ is the weight matrix connecting visible and hidden neurons and $\mathbf{a} \in \mathcal{R}^{N_v \times 1}, \mathbf{b} \in \mathcal{R}^{N_h \times 1}$ are the bias vectors for the visible and hidden layers respectively.

If the visible neurons take real values, i.e. $\mathbf{v} \in \mathcal{R}^{N_v \times 1}$, the RBM formed is called Gaussian-Bernoulli RBM and its energy function is

$$E(\mathbf{v}, \mathbf{h}) = -\frac{1}{2}(\mathbf{v} - \mathbf{a})^T(\mathbf{v} - \mathbf{a}) - \mathbf{b}^T \mathbf{h} - \mathbf{h}^T \mathbf{W} \mathbf{v}$$

Since RBM is an energy model, each configuration of its units (\mathbf{v}, \mathbf{h}) is associated with a probability

$$P(\mathbf{v}, \mathbf{h}) = \frac{e^{-E(\mathbf{v}, \mathbf{h})}}{Z}$$

where

$$Z = \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}$$

Because there are no connections between units of the same layer the posterior probabilities $P(\mathbf{v}|\mathbf{h})$ and $P(\mathbf{h}|\mathbf{v})$ can be efficiently computed. For the Bernoulli-Bernoulli RBM we can easily arrive at the expression for the posterior

$$P(\mathbf{h}|\mathbf{v}) = \prod_i \frac{e^{b_i h_i + h_i \mathbf{W}_{i,*} \mathbf{v}}}{\sum_{\bar{h}_i} e^{b_i \bar{h}_i + \bar{h}_i \mathbf{W}_{i,*} \mathbf{v}}} = \prod_i P(h_i|\mathbf{v})$$

where $\mathbf{W}_{i,*}$ is the i^{th} row of \mathbf{W} . The above expression reveals that the hidden units are conditionally independent given the visible vector. Since $h_i \in \{0, 1\}$

$$P(\mathbf{h} = \mathbf{1}|\mathbf{v}) = \sigma(\mathbf{W}\mathbf{v} + \mathbf{b})$$

where $\sigma(\cdot)$ is the element-wise logistic sigmoid function, and for the binary visible neurons we symmetrically obtain

$$P(\mathbf{v} = \mathbf{1}|\mathbf{h}) = \sigma(\mathbf{W}^T\mathbf{h} + \mathbf{a})$$

In the case of the Gaussian-Bernoulli RBM, the $P(\mathbf{h} = \mathbf{1}|\mathbf{v})$ is the same as above and

$$P(\mathbf{v}|\mathbf{h}) = \mathcal{N}(\mathbf{v}; \mathbf{W}^T\mathbf{h} + \mathbf{a}, I)$$

where I is the identity covariance matrix.

At this point it is important to note that the posterior probability of the hidden units given the visible, independently of the input values, has the same form as the activation function of the deep multilayer perceptron (DNN) with sigmoidal hidden units. This is where the idea of generative pretraining is based, since we can equate the inference for RBM hidden units with forward computation in a DNN and use the weights of an RBM to initialize a feed-forward DNN with sigmoidal hidden units.

Using the free energy defined earlier we can express the marginal probability $P(\mathbf{v})$ as

$$P(\mathbf{v}) = \sum_{\mathbf{h}} P(\mathbf{v}, \mathbf{h}) = \frac{e^{-F(\mathbf{v})}}{\sum_{\mathbf{v}} e^{-F(\mathbf{v})}}$$

For a Gaussian-Bernoulli RBM with no hidden neurons, the marginal probability density function is

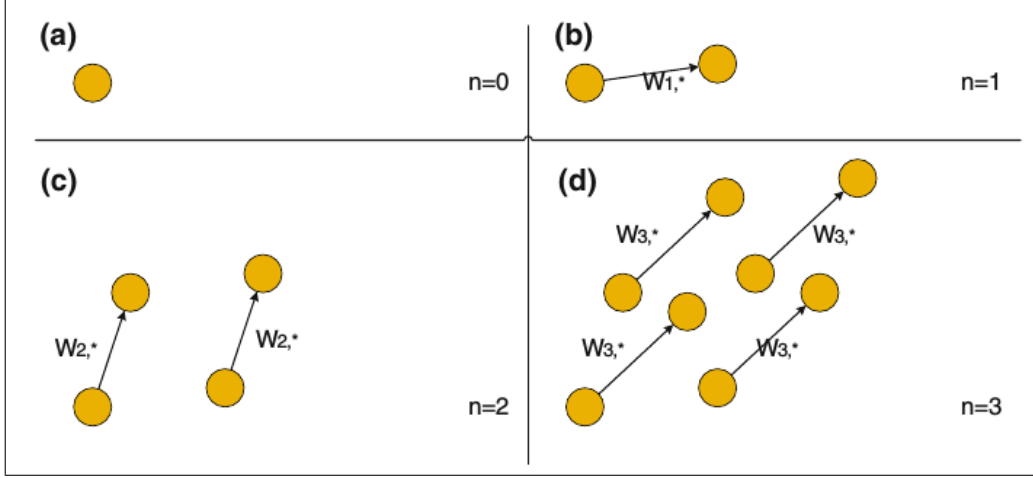
$$p_0(\mathbf{v}) = \frac{e^{-\frac{1}{2}(\mathbf{v}-\mathbf{a})^T(\mathbf{v}-\mathbf{a})}}{Z_0}$$

which is a unit-variance Gaussian distribution centered at vector \mathbf{a} (figure 4.5, part (a)).

It is easy to show ([YD14]) that when a new hidden neuron is added and the rest of the model parameters are fixed, the initial distribution is scaled and a copy of it is placed along the direction determined by $\mathbf{W}_{n,*}$ (figure 4.5, parts (b)-(d)). This means that RBMs represent their visible inputs as a Gaussian Mixture Model with an exponential number of unit-variance Gaussians. Consequently, RBMs can be used in generative models replacing GMMs, since, for example, a Gaussian-Bernoulli RBM can represent real-valued data distributions in a similar way to GMMs.

Another useful property of RBMs is that if we present the training data to the network, it can learn the correlation between different dimensions of the feature

Figure 4.5: *Marginal probability density function of a Gaussian-Bernoulli RBM [YD14].*



vectors, i.e. represent inter-visible connections through the hidden neurons they connect to. However, this also implies that the hidden layers can be used to learn a different representation of the raw input.

When it comes to learning the RBM parameters, one can perform stochastic gradient descent using as training criterion the minimization of the negative log-likelihood:

$$J_{NLL}(\mathbf{W}, \mathbf{a}, \mathbf{b}; \mathbf{v}) = -\log P(\mathbf{v}) = F(\mathbf{v}) + \log \sum_v e^{F(v)}$$

However, the computation of the gradient of the log-likelihood of the data is infeasible to compute exactly:

$$\nabla_{\theta} J_{NLL}(\mathbf{W}, \mathbf{a}, \mathbf{b}; \mathbf{v}) = - \left[\left\langle \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \theta} \right\rangle_{data} - \left\langle \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \theta} \right\rangle_{model} \right]$$

where θ is some model parameter and $\langle \cdot \rangle_{data}$ and $\langle \cdot \rangle_{model}$ are the expectation of the argument estimated from the training data and the model respectively. The first expectation $\langle \cdot \rangle_{data}$ can be computed from the training set, but $\langle \cdot \rangle_{model}$ takes exponential time to compute exactly when the hidden values are unknown. Consequently, we have to use approximated methods for RBM training, with the most widely used being the Contrastive Divergence learning algorithm ([Hin02]).

The Contrastive Divergence algorithm aims to locally approximate the gradient of the training criterion around a training point. Its goal is to discover a decision surface to separate high- from low-probability regions by comparing training samples and samples generated by the model. The term “contrastive” comes exactly from the fact that it builds on the contrast between these two classes of samples.

The algorithm starts by initializing a sampling process (*Gibbs sampling*) at a training data sample. It then generates a hidden sample from the visible sample

Figure 4.6: *Contrastive Divergence algorithm for RBM training [YD14].*

```

1: procedure TRAINRBMWITHCD( $\mathcal{S} = \{\mathbf{o}^m | 0 \leq m < M\}$ ,  $N$ )
    $\triangleright \mathcal{S}$  is the training set with  $M$  samples,  $N$  is the CD steps
2:   Randomly initialize  $\{\mathbf{W}_0, \mathbf{a}_0, \mathbf{b}_0\}$ 
3:   while Stopping Criterion Not Met do
    $\triangleright$  Stop if reached max iterations or the training criterion improvement is small
4:     Randomly select a minibatch  $\mathbf{O}$  with  $M_b$  samples.
5:      $\mathbf{V}^0 \leftarrow \mathbf{O}$   $\triangleright$  Positive phase
6:      $\mathbf{H}^0 \leftarrow P(\mathbf{H}|\mathbf{V}^0)$   $\triangleright$  Applied column-wise
7:      $\nabla_{\mathbf{W}} J \leftarrow \mathbf{H}^0 (\mathbf{V}^0)^T$ 
8:      $\nabla_{\mathbf{a}} J \leftarrow \text{sumrow}(\mathbf{V}^0)$   $\triangleright$  Sum along rows
9:      $\nabla_{\mathbf{b}} J \leftarrow \text{sumrow}(\mathbf{H}^0)$ 
10:    for  $n \leftarrow 0; n < N; n \leftarrow n + 1$  do  $\triangleright$  Negative phase
11:       $\mathbf{H}^n \leftarrow \mathbb{I}(\mathbf{H}^n > \text{rand}(0, 1))$   $\triangleright$  Sampling,  $\mathbb{I}(\bullet)$  is the indicator function
12:       $\mathbf{V}^{n+1} \leftarrow P(\mathbf{V}|\mathbf{H}^n)$ 
13:       $\mathbf{H}^{n+1} \leftarrow P(\mathbf{H}|\mathbf{V}^{n+1})$ 
14:    end for
15:     $\nabla_{\mathbf{W}} J \leftarrow \nabla_{\mathbf{W}} J - \mathbf{H}^N (\mathbf{V}^N)^T$   $\triangleright$  Subtract negative statistics
16:     $\nabla_{\mathbf{a}} J \leftarrow \nabla_{\mathbf{a}} J - \text{sumrow}(\mathbf{V}^N)$ 
17:     $\nabla_{\mathbf{b}} J \leftarrow \nabla_{\mathbf{b}} J - \text{sumrow}(\mathbf{H}^N)$ 
18:     $\mathbf{W}_{t+1} \leftarrow \mathbf{W}_t + \frac{\epsilon}{M_b} \Delta \mathbf{W}_t$   $\triangleright$  Update  $\mathbf{W}$ 
19:     $\mathbf{a}_{t+1} \leftarrow \mathbf{a}_t + \frac{\epsilon}{M_b} \Delta \mathbf{a}_t$   $\triangleright$  Update  $\mathbf{a}$ 
20:     $\mathbf{b}_{t+1} \leftarrow \mathbf{b}_t + \frac{\epsilon}{M_b} \Delta \mathbf{b}_t$   $\triangleright$  Update  $\mathbf{b}$ 
21:  end while
22:  Return  $\text{rbm} = \{\mathbf{W}, \mathbf{a}, \mathbf{b}\}$ 
23: end procedure

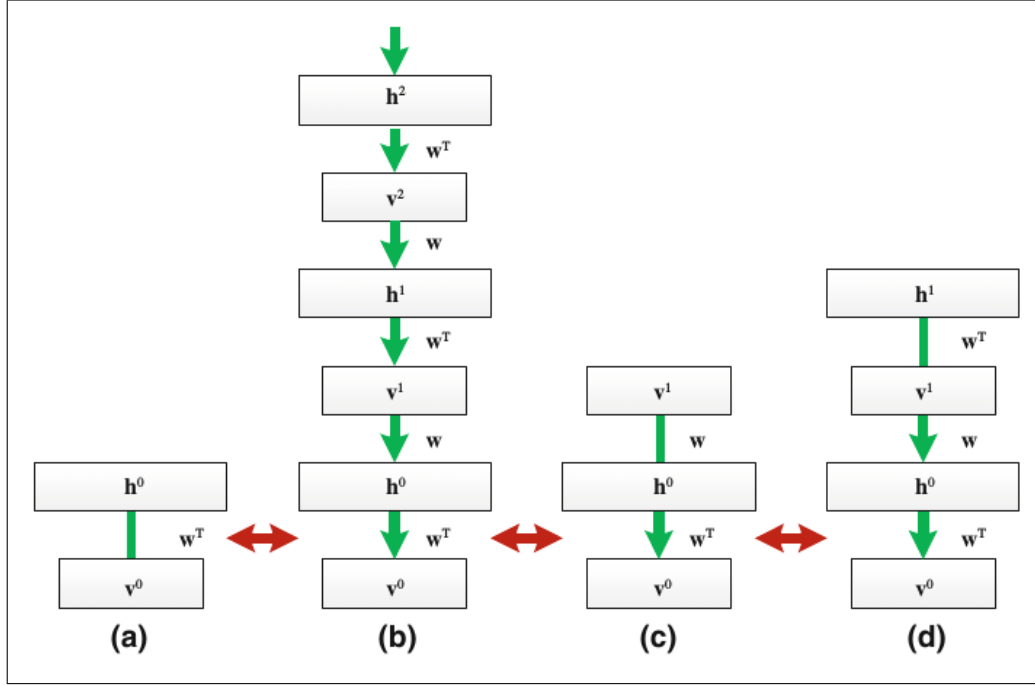
```

based on the posterior probability $P(\mathbf{h}|\mathbf{v})$, defined accordingly to the RBM model used (Gaussian-Bernoulli or Bernoulli-Bernoulli). This hidden sample is further used to generate a visible sample based on the posterior probability $P(\mathbf{v}|\mathbf{h})$. This process may continue for many steps. If it continues for an infinite number of steps the true expectation $\langle \cdot \rangle_{\text{model}}$ can be estimated. However, it has been found that even one step of the algorithm is enough to provide us with a good estimate of $\langle \cdot \rangle_{\text{model}}$ and consequently of the gradient of the training criterion ([Hin12]).

Deep Belief Networks We can view an RBM as a generative model with infinite number of layers all of which share the same weight matrix (figure 4.7, parts (a)-(b)).

If we separate the bottom layer from the deep model of figure 4.7(b) the remaining layers form another generative model with infinite number of layers sharing the same weight matrices, i.e. they are equivalent to another RBM whose visible layer and hidden layer are switched (figure 4.7, part (c)). This model is a special generative model called Deep Belief Network (*DBN*), where the bottom layers form a directed generative model and the top layer is an undirected RBM. Thinking in the same way we can derive the equivalent model shown in figure 4.7(d).

Because the DBN is related closely to the RBM, we can apply a layer-wise process

Figure 4.7: *Different perspectives of the same RBM [YD14].*

to train deep generative models. We start by training an RBM, which, after the training, can discover a different representation of the feature vectors. Hence, for each vector \mathbf{v} we compute a vector of expected hidden neuron activations \mathbf{h} . These hidden expectations can be fed to a new RBM as training data. Continuing in this way we can use each set of RBM weights to extract features from the output of the previous layer. When we do not need to train more RBMs, we can use the weight matrices to initialize the hidden layers of a DBN with as many hidden layers as the RBMs we have trained. The final DBN can be further fine-tuned.

This procedure can be used to stack RBMs with different or the same number of dimensions which allows us to improve the flexibility and performance of the DBN.

Weight initialization of a DNN The DBN weights can be used to initialize the weights of a sigmoidal DNN. This is due to the fact that $P(\mathbf{h}|\mathbf{v})$ in the RBM has the same form as that in the DNN provided the layers use the sigmoidal activation function. The DNN can be viewed as a statistical graphical model where each hidden layer $0 < l < L$ models posterior probabilities of conditionally independent hidden binary neurons \mathbf{h}^l given input vectors \mathbf{v}^{l-1} as Bernoulli distribution

$$P(\mathbf{h}^l|\mathbf{v}^{l-1}) = \sigma(\mathbf{z}^l) = \sigma(\mathbf{W}^l \mathbf{v}^{l-1} + \mathbf{b}^l)$$

and the output layer approximates the label \mathbf{y} as

$$P(\mathbf{y}|\mathbf{v}^{L-1}) = \text{softmax}(\mathbf{z}^L) = \text{softmax}(\mathbf{W}^L \mathbf{v}^{L-1} + \mathbf{b}^L)$$

Given the observation vector \mathbf{o} and its label \mathbf{y} , the precise modeling of $P(\mathbf{y}|\mathbf{o})$ requires integration over all possible values of \mathbf{h} across all layers. However, because this is infeasible, we replace the marginalization with the mean-field approximation:

$$\mathbf{v}^l = \mathbb{E}(\mathbf{h}^l|\mathbf{v}^{l-1}) = P(\mathbf{h}^l|\mathbf{v}^{l-1}) = \sigma(\mathbf{W}^l\mathbf{v}^{l-1} + \mathbf{b}^l)$$

which is the non-stochastic description of the DNN previously presented.

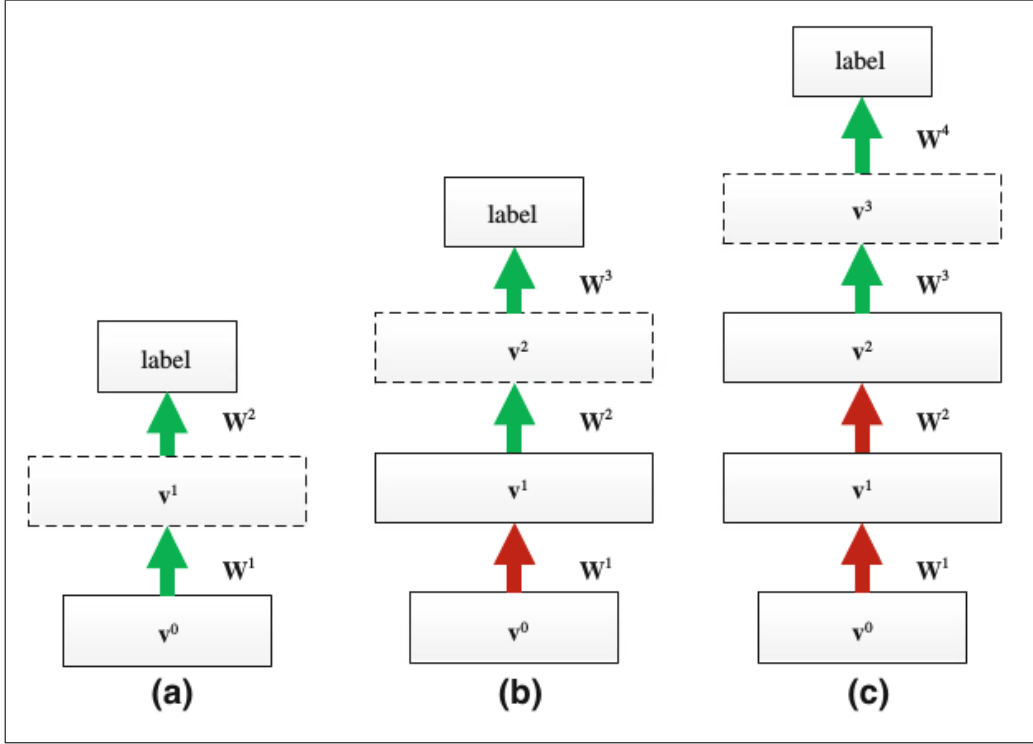
Weight initialization via generative pretraining can improve the DNN performance on the test set. This is due to three main reasons. First, the initial point of the learning algorithm can greatly affect the end model, especially in batch training. Second, since we only need labels for the fine-tuning part of the training, we can utilize a large amount of unlabeled data, which would be useless without pretraining. Third, the generative criterion used in pretraining is different from the discriminative criterion used in fine-tuning, which means that pretraining can act as a data-dependent regularizer.

Bengio in [Ben09] debates whether unsupervised (generative) pretraining has in principle a regularization or optimization effect.

Given that generative pretraining is equivalent to imposing a constraint on where in the parameter space a solution is allowed, i.e. near solutions of the unsupervised training criterion that capture important statistical structure in the data, it can be seen as a form of regularization. Furthermore, generative pretraining plays a more important role in the initialization of the lower layers. Experiments have shown that when there are no constraints on the number of hidden neurons at the top two layers of the network, the representation computed by the lower layers is unimportant: the network can still be trained to minimize the training error as much as desired - despite the fact that it may generalize poorly. However, if we have to use smaller top layers, then pretraining of the lower layers becomes necessary to get a low training error and better generalization.

The optimization perspective of unsupervised pretraining requires to focus on tuning the lower layers while the top two layers are kept small in terms of the number of neurons or of the magnitude of their weights. Bengio suggests that if unsupervised pretraining worked purely as a regularizer, then if we had an infinite stream of training data and performed online learning - in which case we would in essence minimize the generalization error - then with or without pretraining the error would converge to the same level. However, after emulating such a setting, he discovered that the pretrained network achieved a lower minimum, which suggests that there is an optimization component in generative pretraining.

In any case, there has been no evidence so far of generative pretraining hurting the DNN training.

Figure 4.8: *Discriminative pretraining [YD14].*

Apparently, generative pretraining is highly dependent on the task. It also works best with two hidden layers and is of trivial importance when the network has only one hidden layer. If we add more hidden layers effectiveness often decreases, since modeling errors introduced by the mean-field approximation and the contrastive divergence algorithm accumulate ([YD14]).

Discriminative pretraining

A different, “incremental”, approach to pretraining is to discriminatively train each layer using back-propagation. Pretraining begins by training a one-hidden-layer network to convergence, using labels. Then, we insert another hidden layer before the output layer, randomly initialized, and retrain the whole network to convergence. We continue in the same way until the network reaches the desired number of hidden layers (figure 4.8)

It is important to note that all hidden layers are updated during training and not only the layer added last. Because of this, if we are planning to add more hidden layers, the network should not be trained to full convergence, to avoid having some hidden neurons operating in the saturation region and thus being unable to further update. A common heuristic is to go through the training set $\frac{1}{L}$ times of the total number of the data passes needed to convergence, where L is total number of layers

in the final model.

Discriminative pretraining aims to bring the weights close to a good local optimum and it does not have the regularization effect of generative pretraining. Therefore, it is best used when large amounts of training data are available.

Hybrid pretraining

Both pretraining techniques presented work well, however they both have some drawbacks.

Generative pretraining is not directly related to the task-specific objective function, hence it is not guaranteed to help the discriminative fine-tuning phase. However, it does help to reduce overfitting. On the other hand, discriminative pretraining minimizes the training criterion, yet it entails the danger of over-tuning the lower layers to the final objective thus making them unable to learn when new hidden layer are added.

To tackle these issues, it has been suggested to optimize a weighed sum of the generative and discriminative pretraining criteria. We thus have a hybrid criterion to optimize:

$$J_{HYB}(\mathbf{W}, \mathbf{b}; \mathcal{D}_{train}) = J_{DISC}(\mathbf{W}, \mathbf{b}; \mathcal{D}_{train}) + \alpha J_{GEN}(\mathbf{W}, \mathbf{b}; \mathcal{D}_{train})$$

The discriminative criterion can be cross-entropy or mean-square error whereas the generative can be negative log-likelihood. Intuitively, the generative criterion acts as a data-dependent regularizer for the discriminative.

Hybrid pretraining can outperform both the generative and the discriminative pretraining approaches. Despite the fact the the importance of pretraining diminishes as the amount of training data increases, pretraining can still help to make the training procedure more robust.

Dropout pretraining

We have already seen that dropout can be considered a bagging technique, thus generate a smoother objective surface. This implies that dropout could be used to pretrain a DNN to easily find a good starting point on the smoothed objective surface and then fine-tune the model without dropout.

Dropout requires labeled data and achieves similar performance to the discriminative pretraining, as well as being easier to implement and control.

4.4 Recurrent Neural Networks

[YD14]

Research in ASR has been working towards building models that can adequately emulate the human speech production and perception system. However, as we have seen, GMM/HMM models fail to model the dynamic and hierarchical structure of the human speech system and DNN/HMM models exhibit the same type of limitations. A Recurrent Neural Network (*RNN*) is a class of neural network models where many connections among its neurons form a directed cycle (hence the name *recurrent*), thus providing the network with a structure of internal states, or memory, which helps it exhibit a dynamic temporal behavior missing from the DNN.

The internal representation of dynamic speech features in the RNN is discriminatively formed by feeding the low-level features into the hidden layers, together with the hidden features from the past history. The RNN implements a time-delay operation over the temporal dimension which allows its internal states to function as memory and the network to exhibit a dynamic temporal behavior.

RNNs are considered deep models since if we unfold the network in time, we will create as many layers in the network as the length of the input speech utterance.

In this section we will present the fundamental RNN model and the two basic ways to train it, as well as a model which has received a lot of attention recently, the *Long-Short-Term Memory*.

4.4.1 State-Space formulation of the basic RNN

The RNN differs from the DNN in that it operates not only on its inputs but also on its internal states, which encode the past information that has already been processed.

We will express the one-hidden layer RNN in terms of the noise-free nonlinear state-space model often used in signal processing, which will enable the learning of sequentially extended dependencies over a long time span. If \mathbf{x}_t is the $K \times 1$ input vector, \mathbf{h}_t is the $N \times 1$ vector of hidden state values and \mathbf{y}_t is the $L \times 1$ output vector, the one-hidden layer RNN can be described as

$$\mathbf{h}_t = f(\mathbf{W}_{xh}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1}) \text{ (state equation)}$$

$$\mathbf{y}_t = g(\mathbf{W}_{hy}\mathbf{h}_t) \text{ (observation equation)}$$

where \mathbf{W}_{hy} is the $L \times N$ weight matrix connecting the N hidden units to the L outputs, \mathbf{W}_{xh} is the $N \times K$ weight matrix connecting the K inputs to the N hidden units, and \mathbf{W}_{hh} is the $N \times N$ weight matrix connecting the N hidden units from time $t - 1$ to time t .

Let us further define

$$\mathbf{u}_t = \mathbf{W}_{xh}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1}$$

as the $N \times 1$ vector of hidden layer potentials and

$$\mathbf{v}_t = \mathbf{W}_{hy} \mathbf{h}_t$$

as the $L \times 1$ vector of output layer potentials.

Then, $f(\mathbf{u}_t)$ is the hidden layer activation function and $g(\mathbf{v}_t)$ is the output layer activation function. Commonly used activation functions in the hidden layers are the *sigmoid*, *tanh* and rectified linear units (*ReLU*), whereas typical output layer activation functions are the *linear* and *softmax* functions.

In case we want to use the output from previous time frames to update the state vector the state equation becomes

$$\mathbf{h}_t = f(\mathbf{W}_{xh} \mathbf{x}_t + \mathbf{W}_{hh} \mathbf{h}_{t-1} + \mathbf{W}_{yh} \mathbf{y}_{t-1})$$

where \mathbf{W}_{yh} is the weight matrix connecting the output layer to the hidden layer.

4.4.2 Training

Back-propagation-through-time algorithm

The back-propagation-through-time (*BPTT*) method ([Bod01]) manages to learn the weight matrices of the RNN by first unfolding the network in time and then propagating errors backwards through time. The algorithm is an extension of the original back-propagation algorithm for feed-forward neural networks, yet now the stacked hidden layers for the same training frame t have been replaced by the T same single hidden layers across time $t = 1, 2, \dots, T$.

The training criterion of the BPTT algorithm is the well-known sum-square error

$$E = \frac{1}{2} \sum_{t=1}^T \|\mathbf{l}_t - \mathbf{y}_t\|^2$$

where \mathbf{l}_t is the target vector and \mathbf{y}_t is the output vector over all time frames. Our goal is to minimize this cost with respect to the weights and for that we will use the gradient descent rule. The algorithm is presented in figure 4.9

It should be noted that contrary to the DNN backpropagation algorithm, the RNN weight matrices are spatially duplicated for an arbitrary number of time steps, i.e. they are *tied*. That is the reason why in the update step of the algorithm we sum over all time frames.

The computational complexity of the BPTT is $O(M^2)$ per time step, where $M = LN + NK + N^2$ is the total number of the model parameters. It converges slower than the original back-propagation algorithm due to dependencies between frames, and is highly likely to converge to a poor local optimum due to exploding or

Figure 4.9: *Backpropagation-through-time [YD14].*

```

1: procedure BPTT( $\{\mathbf{x}_t, \mathbf{I}_t\} \ 1 \leq t \leq T$ )
     $\triangleright \mathbf{x}_t$  is the input feature sequence
     $\triangleright \mathbf{I}_t$  is the label sequence
     $\triangleright$  forward computation

2:   for  $t \leftarrow 1; t \leq T; t \leftarrow t + 1$  do
3:      $\mathbf{u}_t \leftarrow \mathbf{W}_{xh}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1}$ 
4:      $\mathbf{h}_t \leftarrow f(\mathbf{u}_t)$ 
5:      $\mathbf{v}_t \leftarrow \mathbf{W}_{hy}\mathbf{h}_t$ 
6:      $\mathbf{y}_t \leftarrow g(\mathbf{v}_t)$ 
7:   end for

8:    $\delta_T^y \leftarrow (\mathbf{I}_T - \mathbf{y}_T) \bullet g'(\mathbf{v}_T)$ 
9:    $\delta_T^h \leftarrow \mathbf{W}_{hy}^T \delta_T^y \bullet f'(\mathbf{u}_T)$ 
10:  for  $t \leftarrow T - 1; t \geq 1; t \leftarrow t - 1$  do
11:     $\delta_t^y \leftarrow (\mathbf{I}_t - \mathbf{y}_t) \bullet g'(\mathbf{v}_t)$ 
12:     $\delta_t^h \leftarrow [\mathbf{W}_{hh}^T \delta_{t+1}^h + \mathbf{W}_{hy}^T \delta_{t+1}^y] \bullet f'(\mathbf{u}_t)$ 
13:  end for
14:   $\mathbf{W}_{hy} \leftarrow \mathbf{W}_{hy} + \gamma \sum_{t=1}^T \delta_t^y \mathbf{h}_t^T$ 
15:   $\mathbf{W}_{hh} \leftarrow \mathbf{W}_{hh} + \gamma \sum_{t=1}^T \delta_t^h \mathbf{h}_{t-1}^T$ 
16: end procedure

```

\triangleright backpropagation through time
 $\triangleright \bullet$: element-wise multiplication
 \triangleright propagate from δ_t^y and δ_{t+1}^h
 \triangleright model update

vanishing gradients. Although we can limit the past history to the last p time steps and thus improve speed, it is still difficult to achieve good results without much experimentation and tuning.

Primal-dual learning algorithm

The well-known problems of the back-propagation algorithm in conjunction with gradient computation are even more evident in the case of RNNs. Learning RNNs is particularly difficult because of the exploding (too big) and vanishing (too small) gradients.

A sufficient condition for vanishing gradients to occur is $\|\mathbf{W}_{hh}\| < d$, where d is a constant dependent on the activation function of the hidden units ($d = 4$ for sigmoidal units and $d = 1$ for linear) and $\|\mathbf{W}_{hh}\|$ is the largest singular value of the weight matrix \mathbf{W}_{hh} .

Similarly, a necessary condition for exploding gradients is $\|\mathbf{W}_{hh}\| > d$.

Therefore, in order to effectively learn an RNN, we have to impose certain constraints on the weight matrix $\|\mathbf{W}_{hh}\| > d$. For example, to avoid vanishing gradients, we could add a regularization term to increase the gradient or exploit curvature-related information about the regularization function; to avoid exploding gradients we could clip the gradient at an empirically-determined threshold.

A more rigorous approach to RNN learning is would be to directly exploit the conditions described above regarding \mathbf{W}_{hh} , and this is the focus of the primal-dual

algorithm for training RNNs. This algorithm tries to preserve the *echo state property* which is closely related to the previous constraints and states that *if the network has been run for a very long time, the current network state is uniquely determined by the history of the input and the teacher-forced output*. Mathematically, a sufficient condition for the network to satisfy the echo state property is

$$\|\mathbf{W}_{hh}\|_{\infty} < d$$

where $\|\mathbf{W}_{hh}\|_{\infty}$ is the maximum absolute sum of \mathbf{W}_{hh} , $d = 4$ for sigmoidal units and $d = 1$ for hyperbolic tangent units.

Given that we want to preserve the echo state property, RNN learning can be formulated as a constraint optimization problem:

$$\min_{\Theta} E(\Theta) = \min_{\Theta} E(\mathbf{W}_{hh}, \mathbf{W}_{xh}, \mathbf{W}_{hy})$$

subject to

$$\|\mathbf{W}_{hh}\|_{\infty} \leq d$$

The Lagrangian of this problem can be formulated as:

$$L(\Theta, \boldsymbol{\lambda}) = E(\mathbf{W}_{hh}, \mathbf{W}_{xh}, \mathbf{W}_{hy}) + \sum_{i=1}^N \lambda_i \left(\sum_{j=1}^N |W_{ij}| - d \right)$$

where $\lambda_i \geq 0$ is the i th entry of the Lagrange vector $\boldsymbol{\lambda}$.

Then, the dual function $q(\boldsymbol{\lambda})$ is the solution to the following unconstrained optimization problem

$$q(\boldsymbol{\lambda}) = \min_{\Theta} L(\Theta, \boldsymbol{\lambda})$$

and is a lower bound of the original constrained optimization problem, i.e.

$$q(\boldsymbol{\lambda}) \leq E(\Theta^*)$$

Maximizing $q(\boldsymbol{\lambda})$ will be the best lower bound that can be obtained for $E(\Theta)$. This is called the *dual problem* of the original optimization problem:

$$\max_{\boldsymbol{\lambda}} q(\boldsymbol{\lambda}), \quad \lambda_i \geq 0, \quad i = 1, \dots, N$$

and is a convex optimization problem since we aim to maximize a concave objective with linear equality constraints.

Having solved the dual problem for $\boldsymbol{\lambda}^*$ we substitute the optimal dual variables into the Lagrangian and solve for the set of parameters Θ that minimize $L(\Theta, \boldsymbol{\lambda}^*)$:

$$\Theta^0 = \operatorname{argmin}_{\Theta} L(\Theta, \boldsymbol{\lambda}^*)$$

This solution will be an approximation to an optimal solution to the original constrained optimization problem. However, in general non-convex problems where acquiring a globally optimal solution is infeasible, this approximation will be good enough.

Following the two-stepped solution, the updates of the RNN parameters happen in two steps:

- primal update: minimize $L(\Theta, \lambda^*)$ w.r.t. Θ
- dual update: maximize $L(\Theta^*, \lambda)$ w.r.t. λ

The standard gradient descent algorithm can be applied to the update rules with some improvements being possible due to the structure in the objective function.

4.4.3 RNNs with Long-Short-Term Memory cells

Motivation

The basic RNN presented so far does not have the necessary structure to model complex temporal dynamics and is thus incapable of looking far back into the past in many types of input sequences.

One way to solve this problem is to introduce a memory structure into the RNN which leads to the incorporation of *long-short-term memory cells* (LSTM) ([HS97]). Due to its structure, the LSTM-RNN can recognize temporally extended patterns in noisy input sequences and temporal order of widely separated events in noisy input streams. Like the basic RNN, the LSTM-RNN is a universal computing machine, provided it has enough network units and receives proper training, and, in addition, it can learn from input sequence data to classify, process and predict time series with very long time lags of unknown lengths between important events.

The LSTM-RNN has been shown to perform very well in handwriting and phone recognition, keyword spotting, reinforcement learning for robot localization and control, online learning for protein structure prediction, learning music composition and grammars, language identification, speech synthesis and environment-robust speech recognition.

Architecture of LSTM cells

The basic notion behind LSTM cells in the RNN is to use various types of element-wise multiplication (or *gating*) structure to control the information flow in the network. Such a cell can be viewed as a complex and smart network unit able to remember information for a long time. This is the job of the gating structure which

is in charge of determining when the input is important enough to remember, when the cell should continue to remember or forget the information and when it should output the information.

An LSTM cell is mathematically described by the following forward operations iteratively over time $t = 1, 2, \dots, T$:

$$\begin{aligned} \mathbf{i}_t &= \sigma\left(\mathbf{W}^{(xi)}\mathbf{x}_t + \mathbf{W}^{(hi)}\mathbf{h}_{t-1} + \mathbf{W}^{(ci)}\mathbf{c}_{t-1} + \mathbf{b}^{(i)}\right) \\ \mathbf{f}_t &= \sigma\left(\mathbf{W}^{(xf)}\mathbf{x}_t + \mathbf{W}^{(hf)}\mathbf{h}_{t-1} + \mathbf{W}^{(cf)}\mathbf{c}_{t-1} + \mathbf{b}^{(f)}\right) \\ \mathbf{c}_t &= \mathbf{f}_t \bullet \mathbf{c}_{t-1} + \mathbf{i}_t \bullet \tanh\left(\mathbf{W}^{(xc)}\mathbf{x}_t + \mathbf{W}^{(hc)}\mathbf{h}_{t-1} + \mathbf{b}^{(c)}\right) \\ \mathbf{o}_t &= \sigma\left(\mathbf{W}^{(xo)}\mathbf{x}_t + \mathbf{W}^{(ho)}\mathbf{h}_{t-1} + \mathbf{W}^{(co)}\mathbf{c}_t + \mathbf{b}^{(o)}\right) \\ \mathbf{h}_t &= \mathbf{o}_t \bullet \tanh(\mathbf{c}_t) \end{aligned}$$

where \mathbf{i}_t , \mathbf{f}_t , \mathbf{c}_t , \mathbf{o}_t and \mathbf{h}_t are vectors, all with the same dimensionality, which represent five different types of information at time t of the input gate, forget gate, cell activation, output gate and hidden layer respectively; $\sigma(\cdot)$ is the logistic sigmoid function, \mathbf{W} 's are the weight matrices connecting different gates and \mathbf{b} 's are the corresponding bias vectors. The weight matrices are full except for $\mathbf{W}^{(ci)}$ that is diagonal. In addition to the above set, an additional output layer has to be provided on top of the LSTM-RNN's hidden layer.

Training

The Back-Propagation-Through-Time algorithm introduced for RNN training can be also used in the case of LSTM-RNNs in order to compute error derivatives (although the computation is more complex) and then apply stochastic gradient descent methods to update the weight parameters.

However, the problems of vanishing and exploding gradients encountered in RNN training are less evident here. This is because the error signals back-propagated from the output become trapped in the memory part of the LSTM cells. In this way, meaningful error signals are constantly fed back to each of the gates until the RNN parameters are well trained.

Bibliography

- [AE] John McKenna Andrew Errity. An investigation of manifold learning for speech analysis.
- [BCV13] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(8):1798–1828, August 2013.
- [BDVJ03] Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Janvin. A neural probabilistic language model. *J. Mach. Learn. Res.*, 3:1137–1155, March 2003.
- [Bel03] Richard Ernest Bellman. *Dynamic Programming*. Dover Publications, Incorporated, 2003.
- [Ben09] Yoshua Bengio. Learning deep architectures for ai. *Found. Trends Mach. Learn.*, 2(1):1–127, January 2009.
- [Ben12] Yoshua Bengio. Practical recommendations for gradient-based training of deep architectures. *CoRR*, abs/1206.5533, 2012.
- [BGC15] Yoshua Bengio, Ian J. Goodfellow, and Aaron Courville. Deep learning. Book in preparation for MIT Press, 2015.
- [Bil03] Jeff A. Bilmes. Buried markov models: a graphical-modeling approach to automatic speech recognition. *Computer Speech and Language*, 17(2.3):213 – 231, 2003. New Computational Paradigms for Acoustic Modeling in Speech Recognition.
- [Bod01] Mikael Boden. A guide to recurrent neural networks and backpropagation, 2001.
- [BPV03] Yoshua Bengio, Jean-Francois Paiement, and Pascal Vincent. Out-of-sample extensions for lle, isomap, mds, eigenmaps, and spectral clustering. In *In Advances in Neural Information Processing Systems*, pages 177–184. MIT Press, 2003.

- [Cay05] L. Cayton. Algorithms for manifold learning. *Univ. of California at San Diego Tech. Rep*, 2005.
- [DG03] Sanjoy Dasgupta and Anupam Gupta. An elementary proof of a theorem of johnson and lindenstrauss. *Random Struct. Algorithms*, 22(1):60–65, January 2003.
- [DHS00] Richard O. Duda, Peter E. Hart, and David G. Stork. *Pattern Classification (2Nd Edition)*. Wiley-Interscience, 2000.
- [DLR77] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *JOURNAL OF THE ROYAL STATISTICAL SOCIETY, SERIES B*, 39(1):1–38, 1977.
- [DY14] Li Deng and Dong Yu. Deep learning: Methods and applications. *Found. Trends Signal Process.*, 7(3–4):197–387, June 2014.
- [GB10] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *In Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS10). Society for Artificial Intelligence and Statistics*, 2010.
- [Hin02] Geoffrey E. Hinton. Training products of experts by minimizing contrastive divergence. *Neural Comput.*, 14(8):1771–1800, August 2002.
- [Hin12] Geoffrey E. Hinton. A practical guide to training restricted boltzmann machines. In Gregoire Montavon, Genevieve B. Orr, and Klaus-Robert MG’Oller, editors, *Neural Networks: Tricks of the Trade (2nd ed.)*, volume 7700 of *Lecture Notes in Computer Science*, pages 599–619. Springer, 2012.
- [HS52] M. R. Hestenes and E. Stiefel. Methods of conjugate gradients for solving linear systems. *Journal of research of the National Bureau of Standards*, 49:409–436, 1952.
- [HS97] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–1780, November 1997.
- [HSW89] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5):359 – 366, 1989.
- [JM09] Daniel Jurafsky and James H. Martin. *Speech and Language Processing (2Nd Edition)*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 2009.

- [JN05] A. Jansen and P. Niyogi. A geometric perspective on speech sounds. Technical report, 2005.
- [KGV83] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi. Optimization by simulated annealing. *SCIENCE*, 220(4598):671–680, 1983.
- [Kre78] E. Kreyszig. *Introductory Functional Analysis With Applications*. Wiley Classics Library. John Wiley & Sons, 1978.
- [LBBH98] Yann Lecun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. In *Proceedings of the IEEE*, pages 2278–2324, 1998.
- [LBOM98] Yann LeCun, Léon Bottou, Genevieve B. Orr, and Klaus-Robert Müller. Efficient backprop. In *Neural Networks: Tricks of the Trade, This Book is an Outgrowth of a 1996 NIPS Workshop*, pages 9–50, London, UK, UK, 1998. Springer-Verlag.
- [LN89] D. C. Liu and J. Nocedal. On the limited memory bfgs method for large scale optimization. *Math. Program.*, 45(3):503–528, December 1989.
- [Mah09] Michael Mahoney. Algorithms for massive datasets analysis, lecture notes, 2009.
- [MKB⁺10] Tomas Mikolov, Martin Karafiat, Lukas Burget, Jan Cernocky, and Sanjeev Khudanpur. Recurrent neural network based language model. In Takao Kobayashi, Keikichi Hirose, and Satoshi Nakamura, editors, *INTERSPEECH*, pages 1045–1048. ISCA, 2010.
- [MPR01] Mehryar Mohri, Fernando Pereira, and Michael Riley. Weighted finite-state transducers in speech recognition, 2001.
- [Nie15] Michael Nielsen. *Neural Networks and deep learning*, 2015.
- [PHB⁺12] Daniel Povey, Mirko Hannemann, Gilles Boulianne, Lukas Burget, Arnab Ghoshal, Milos Janda, Martin Karafiat, Stefan Kombrink, Petr Motlicek, Yanmin Qian, Korbinnian Riedhammer, Karel Vesely, and Ngoc Thang Vu. Generating exact lattices in the wfst framework. In *ICASSP*, pages 4213–4216. IEEE, 2012.
- [Rab89] Lawrence R. Rabiner. A tutorial on hidden markov models and selected applications in speech recognition. In *PROCEEDINGS OF THE IEEE*, pages 257–286, 1989.

- [RJ86] L. Rabiner and B. Juang. An introduction to hidden markov models. *IEEE Acoutics, Speech and Signal Processing Magazine*, 3:4–16, 1986.
- [SKK⁺07] Thomas Serre, Gabriel Kreiman, Minjoon Kouh, Charles Cadieu, Ulf Knoblich, and Tomaso Poggio. A quantitative theory of immediate visual recognition. *PROG BRAIN RES*, pages 33–56, 2007.
- [SR03] Lawrence K. Saul and Sam T. Roweis. Think globally, fit locally: Un-supervised learning of low dimensional manifolds. *J. Mach. Learn. Res.*, 4:119–155, December 2003.
- [SS05] Fei Sha and Lawrence K. Saul. Analysis and extension of spectral methods for nonlinear dimensionality reduction. In *Proceedings of the 22Nd International Conference on Machine Learning, ICML '05*, pages 784–791, New York, NY, USA, 2005. ACM.
- [ST04] Darius Silingas and Laimutis Telksnys. Specifics of hidden markov model modifications for large vocabulary continuous speech recognition. *Informatika, Lith. Acad. Sci.*, 15(1):93–110, 2004.
- [SVN37] S. S. Stevens, J. Volkmann, and E. B. Newman. A scale for the measurement of the psychological magnitude pitch. *The Journal of the Acoustical Society of America*, 8(3):185–190, 1937.
- [TR12a] Vikrant Singh Tomar and Richard C. Rose. Application of a locality preserving discriminant analysis approach to ASR. In *11th International Conference on Information Science, Signal Processing and their Applications, ISSPA 2012, Montreal, QC, Canada, July 2-5, 2012*, pages 103–107, 2012.
- [TR12b] Vikrant Singh Tomar and Richard C. Rose. A correlational discriminant approach to feature extraction for robust speech recognition. In *INTERSPEECH 2012, 13th Annual Conference of the International Speech Communication Association, Portland, Oregon, USA, September 9-13, 2012*, pages 555–558, 2012.
- [TR13a] Vikrant Singh Tomar and Richard C. Rose. Efficient manifold learning for speech recognition using locality sensitive hashing. In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2013, Vancouver, BC, Canada, May 26-31, 2013*, pages 6995–6999, 2013.
- [TR13b] Vikrant Singh Tomar and Richard C. Rose. Locality sensitive hashing for fast computation of correlational manifold learning based feature space

- transformations. In *INTERSPEECH 2013, 14th Annual Conference of the International Speech Communication Association, Lyon, France, August 25-29, 2013*, pages 1776–1780, 2013.
- [TR13c] Vikrant Singh Tomar and Richard C. Rose. Noise aware manifold learning for robust speech recognition. In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2013, Vancouver, BC, Canada, May 26-31, 2013*, pages 7087–7091, 2013.
- [TR14] Vikrant Singh Tomar and Richard C. Rose. A family of discriminative manifold learning algorithms and their application to speech recognition. *IEEE/ACM Transactions on Audio, Speech & Language Processing*, 22(1):161–171, 2014.
- [Ver08] Nakul Verma. Mathematical advances in manifold learning, 2008.
- [YD14] Dong Yu and Li Deng. *Automatic Speech Recognition: A Deep Learning Approach*. Springer Publishing Company, Incorporated, 2014.
- [YEK⁺02] S. Young, G. Evermann, D. Kershaw, G. Moore, J. Odell, D. Ollason, V. Valtchev, and P. Woodland. The htk book. *Cambridge University Engineering Department*, 3, 2002.
- [You02] Steve Young. Acoustic modelling for large vocabulary continuous speech recognition. *Cambridge University Engineering Department*, 2002.
- [YXZ⁺07] Shuicheng Yan, Dong Xu, Benyu Zhang, Hong-Jiang Zhang, Qiang Yang, and Stephen Lin. Graph embedding and extensions: A general framework for dimensionality reduction. *IEEE Trans. Pattern Anal. Mach. Intell.*, 29(1):40–51, January 2007.

