



# 语音识别学习笔记

## 学习大法好

作者：杜沪

组织：BIT, SpeechOcean

时间：September 18, 2019

版本：0.1



*Victory won't come to us unless we go to it. — M. Moore*

# 目 录

<b>1</b>	<b>数学知识总结</b>	<b>1</b>
1.1	各类矩阵定义 . . . . .	1
1.2	瑞利商 . . . . .	4
1.3	EM 算法 . . . . .	8
1.3.1	Jensen's Inequality . . . . .	8
1.3.2	EM 算法推导 . . . . .	9
1.3.3	EM 算法收敛性证明 . . . . .	12
1.4	混合高斯分布 . . . . .	13
1.5	HMM 相关知识点总结 . . . . .	18
1.5.1	Markov Models . . . . .	18
1.5.2	Markov Model 的两个基本问题 . . . . .	19
1.5.3	Hidden Markov Model . . . . .	23
1.5.4	HMM 的三个基本问题 . . . . .	23
1.6	线性判别分析 . . . . .	33
1.7	最大似然线性变换 . . . . .	36
1.8	Beta 分布 . . . . .	36
1.9	MLE 和 MAP . . . . .	36
<b>2</b>	<b>微软 Edx 语音识别笔记</b>	<b>37</b>
2.1	Background and Fundamentals . . . . .	37
2.1.1	Phonetics . . . . .	37
2.1.2	Words and Syntax . . . . .	38
2.1.3	Measuring Performance . . . . .	38
2.1.4	Significance Testing . . . . .	39
2.1.5	Other Consideration . . . . .	40
2.1.6	The Fundamental Equation . . . . .	40
2.1.7	Lab 1: Create a speech recognition scoring program . . . . .	41
2.2	Speech Signal Processing . . . . .	44
2.2.1	Introduction . . . . .	44
2.2.2	Feature Extraction . . . . .	46
2.2.3	Mel Filtering . . . . .	46

2.2.4	Log Compression . . . . .	47
2.2.5	Feature Normalization . . . . .	48
2.2.6	Summary . . . . .	49
2.2.7	Lab 2: Feature extraction for speech recognition . . . . .	50
2.3	Acoustic Modeling . . . . .	52
2.3.1	Introduction . . . . .	52
2.3.2	Markov Chains . . . . .	52
2.3.3	Hidden Markov Models for Speech Recognition . . . . .	54
2.3.4	Deep Neural Network Acoustic Models . . . . .	57
2.3.5	Training Feedforward Deep Neural Networks . . . . .	58
2.3.6	Using a Sequence based Objective Function . . . . .	61
2.3.7	Decoding with Neural Network Acoustic Models . . . . .	62
2.3.8	Lab 3 . . . . .	62
2.4	Language Modeling . . . . .	66
2.4.1	Introduction . . . . .	66
2.4.2	N gram Models . . . . .	67
2.4.3	Language Model Evaluation . . . . .	69
2.4.4	Operations on Language Models . . . . .	69
2.4.5	Advanced LM Topics . . . . .	69
2.4.6	Lab 4 . . . . .	69
2.5	Speech Decoding . . . . .	69
2.5.1	Overview . . . . .	69
2.5.2	Weighted Finite State Transducers . . . . .	69
2.5.3	WFSTs and Acceptors . . . . .	69
2.5.4	Graph Composition . . . . .	69
2.5.5	Lab 5 . . . . .	69
2.6	Advanced Acoustic Modeling . . . . .	69
2.6.1	Imporved Objective Functions . . . . .	69
2.6.2	Sequential Objective Function . . . . .	69
2.6.3	Connectionist Temporal Classsification . . . . .	69
2.6.4	Sequence Discriminative Objective Functions . . . . .	69
2.6.5	Lab 6 . . . . .	69
2.7	补充知识点 . . . . .	69
2.7.1	傅里叶变换 . . . . .	69



2.7.2	Nyquist 定理	69
2.7.3	编辑距离	69
2.7.4	nonlinear infinite impulse response (IIR)	70
2.7.5	序列区分性训练中的目标函数详解	70
2.7.6	概率与似然度的的是是非非	70
<b>3</b>	<b>Kaldi 学习笔记</b>	<b>71</b>
3.1	kaldi 中的数据扰动	71
3.1.1	速度扰动	71
3.1.2	音量扰动	71
3.2	kaldi 中的 UBM	72
3.3	kaldi 通过 lattice 输出语音对齐音素和词	72
3.4	kaldi 中的数据准备	72
3.4.1	音频、文本、说话人准备	72
3.5	kaldi 中的决策树	73
<b>4</b>	<b>FFmpeg 和 sox</b>	<b>74</b>
4.1	FFmpeg	74
4.1.1	安装 FFmpeg	74
4.2	sox	75
<b>5</b>	<b>Linux 相关笔记</b>	<b>76</b>
5.1	linux 备忘录	76
5.2	Shell 指令笔记	78
5.2.1	简单的 shell 规则	78
5.2.2	shell 中的条件测试操作	80
5.2.3	find	81
5.2.4	grep	82
5.2.5	awk	83
5.2.6	sed	83
<b>6</b>	<b>Windows 相关</b>	<b>84</b>
6.1	win10 .net framework 3.5 安装报错 0x800F0954 问题	84
<b>7</b>	<b>Python 笔记</b>	<b>85</b>
7.1	一些小技巧	85

7.2	python 中的线程、进程、协程与并行、并发 . . . . .	86
7.2.1	进程和线程 . . . . .	87
7.3	客户端向服务端传送一个音频文件及信息 . . . . .	88
7.3.1	wave . . . . .	88
7.3.2	struct . . . . .	89
7.3.3	socket . . . . .	92
7.3.4	socketserver . . . . .	94
7.3.5	网络传输音频并保存 . . . . .	95
7.4	Python 中的正则表达式 . . . . .	98
7.5	替代 os 的 pathlib . . . . .	98
8	C++ 学习笔记 . . . . .	99
9	Docker . . . . .	100
9.1	安装 docker 和 nvidia-docker . . . . .	100
9.2	常用操作 . . . . .	100
9.3	实践要求 . . . . .	101
9.4	docker 安装 TensorFlow . . . . .	101
10	端到端语音识别汇总 . . . . .	102
10.1	CTC . . . . .	102
10.1.1	白话 CTC . . . . .	102
10.1.2	CTC 中的前后向算法 . . . . .	104
10.1.3	CTC 中的 loss 函数和梯度 . . . . .	108
10.1.4	CTC 的解码 . . . . .	111
10.2	RNN-Tranducer . . . . .	118
10.3	Attention . . . . .	118
10.4	Transformer . . . . .	118
10.5	CNNs . . . . .	118
10.6	Mixed Models . . . . .	118
10.6.1	Self-Attention Transducers for End-to-End Speech Recognition . . . . .	118
10.7	如何计算 WER? . . . . .	118
11	论文阅读笔记 . . . . .	119
11.1	Light Gated Recurrent Units for Speech Recogntion . . . . .	119
11.1.1	GRU 的介绍 . . . . .	119



11.1.2 Li-GRU 的学习 . . . . .	120
11.1.3 个人心得体会 . . . . .	123
<b>12 深度学习框架笔记</b>	<b>124</b>
12.1 paddlepaddle . . . . .	124
12.2 pytorch . . . . .	124
12.3 tensorflow . . . . .	124



# 插图目录

1.1	Jensen's Inequality 示例 . . . . .	9
1.2	二类 LDA 转换效果图 . . . . .	33
1.3	多类 LDA 的类间散度矩阵示意图 . . . . .	35
2.1	美式英语的音素和一般实现办法 . . . . .	37
2.2	完整音频波形图与部分波形图 . . . . .	45
2.3	mel filterbank . . . . .	47
2.4	提取 Fbank 特征（左）与原始信号的频谱图（右） . . . . .	48
2.5	原始信号的频谱图（左）、提取 Fbank 特征（右）和归一化后的 Fbank 特征（下） . . . . .	50
2.6	Lab 2 的期望输出 . . . . .	52
2.7	天气预报模型的马尔科夫链 . . . . .	53
2.8	天气预报模型的隐马尔科夫模型 . . . . .	54
2.9	音素"/uh/" 的 HMM . . . . .	55
2.10	单词"cup" 的 HMM . . . . .	55
2.11	单词"cup" 的三音素 HMM . . . . .	56
2.12	senones 的生成 . . . . .	57
2.13	MLF 文件示例 . . . . .	58
2.14	基于 DNN 的声学模型训练过程 . . . . .	59
2.15	基于 DNN 的声学模型训练过程 . . . . .	61
2.16	DNN 声学模型训练图 . . . . .	64
3.1	lexicon.txt 中存储格式 . . . . .	73
5.1	VScode 设置脚本格式为 Unix . . . . .	78
10.1	同一单词不一样的输出音频的图解 . . . . .	103
10.2	原 Alex 博士论文中链式求导部分的错误 . . . . .	111
10.3	Prefix Beam Search 原论文中算法错误地方 . . . . .	114
11.1	GRU 模型结构图 . . . . .	120
11.2	TIMIT 中更新门与重置门音频上的时域关联 . . . . .	121
11.3	Auto-correlation $C(z, z)$ 和 cross-correlation $C(z, r)$ . . . . .	122

# 表格目录

2.1	WER 计算公式中的三种错误实例演示 . . . . .	39
3.1	中英文 word 和对应 pronunciation 示例 . . . . .	73
5.1	常用的文件操作符及其意义 . . . . .	80
5.2	常用的字符串操作符及其意义 . . . . .	81
5.3	常用的整数操作符及其意义 . . . . .	81
5.4	常用的逻辑操作符及其意义 . . . . .	81
7.1	struct 中常用的数据类型、C 语言的对应类型和占用字节数 . . . . .	90
7.2	struct 中的字节对齐操作符号及其含义 . . . . .	90
11.1	音素分类及示例 . . . . .	123



# 第 1 章 数学知识总结

## 1.1 各类矩阵定义

### 定义 1.1. 转置矩阵

把矩阵  $A$  的行换乘同序数的列得到一个新矩阵，就叫做  $A$  的转置矩阵，记作  $A^T$ 。例如矩阵

$$A = \begin{bmatrix} 1 & 2 & 0 \\ 3 & -1 & 1 \end{bmatrix} \quad (1.1)$$

的转置矩阵为

$$A^T = \begin{bmatrix} 1 & 3 \\ 2 & -1 \\ 0 & 1 \end{bmatrix} \quad (1.2)$$



### 定义 1.2. 对称矩阵

设  $A$  为  $n$  阶方阵，如果满足  $A^T = A$ ，即：

$$a_{ij} = a_{ji} (i, j = 1, 2, \dots, n) \quad (1.3)$$

那么  $A$  称为对称矩阵，简称为对称阵。对称阵的特点是：它的元素以对角线为对称轴对应相等。



### 定义 1.3. 复共轭矩阵

设  $A \in C^{m \times n}$ ，用  $\bar{A}$  表示以  $A$  的元素的共轭复数为元素组成的矩阵，命：

$$A^H = (\bar{A})^T \quad (1.4)$$

则称  $A^H$  为  $A$  的复共轭转置矩阵。



### 定义 1.4. Hermitian 矩阵

设  $A \in R^{n \times n}$ ，若  $A^H = A$ ，则称  $A$  为 Hermitian 矩阵。若  $A^H = -A$ ，则称  $A$  为反 Hermitian 矩阵。



**定义 1.5. 正交矩阵**

如果  $n$  阶矩阵  $A$  满足

$$A^T A = E \quad (1.5)$$

即:

$$A^T = A^{-1} \quad (1.6)$$

则称  $A$  为正交矩阵, 简称正交阵。

**定义 1.6. 酉矩阵**

如果  $n$  阶复矩阵  $A$  满足

$$A^H A = A A^H = E \quad (1.7)$$

则称  $A$  为酉矩阵, 记作  $A \in U^{n \times n}$ 。

**定义 1.7. 奇异矩阵**

当  $|A| = 0$  时,  $A$  称为奇异矩阵, 否则称为非奇异矩阵。 $A$  是可逆矩阵的充分必要条件是  $|A| \neq 0$ , 即可逆矩阵就是非奇异矩阵。

**定义 1.8. 正规矩阵**

设  $A \in C^{n \times n}$ , 若:

$$A^H A = A A^H \quad (1.8)$$

则称  $A$  为正规矩阵,  $A \in R^{n \times n}$ , 显然有  $A^H = A^T$ , 上式就变成了:

$$A^T A = A A^T \quad (1.9)$$

则称  $A$  为实正规矩阵。

**定义 1.9. 幂等矩阵**

设  $A \in C^{n \times n}$ , 若:

$$A^2 = A \quad (1.10)$$

则称  $A$  是幂等矩阵。



### 定义 1.10. 正定矩阵

设  $A \in C^{n \times n}$ , 若  $A$  的所有特征值均为正数, 则称  $A$  为正定矩阵; 若  $A$  的特征值均为非负数, 则称  $A$  为半正定矩阵。

判断一个矩阵为正定矩阵的充要条件有:

1.  $A$  的所有特征值  $\lambda_i$  均为正数;
2.  $x^T A x \geq 0$  对所有非零向量  $x$  都成立;
3. 存在秩满矩阵  $R$ , 使得  $A = R^T R$ 。



### 定义 1.11. Jacobi 矩阵

假设某函数从  $f: R^n \rightarrow R^m$ , 从  $x \in R^n$  映射到向量  $f(x) \in R^m$ , 其 Jacobi 矩阵的维度是  $m \times n$ , 如下所示:

$$H = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_1} & \cdots & \frac{\partial f}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \quad (1.11)$$



### 定义 1.12. Hessian 矩阵

若实值函数  $f(x_1, x_2, \dots, x_n)$  的所有二阶偏导都存在并在定义域内连续, 那么函数  $f$  的 Hessian 矩阵为:

$$H = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix} \quad (1.12)$$

根据:

$$\frac{\partial^2 f}{\partial x_1 \partial x_2} = \frac{\partial^2 f}{\partial x_2 \partial x_1} \quad (1.13)$$

可知 Hessian 矩阵为对称阵。



## 1.2 瑞利商

对于一个Hermitian 矩阵 $M$  及非零向量  $x$ ，瑞利商(Rayleigh quotient) 的定义如公式1.14，其中  $x^H$  为  $x$  的共轭转置向量。

$$R(M, x) = \frac{x^H M x}{x^H x} \quad (1.14)$$

若  $M$  和  $x$  中元素均为实数，瑞利商可以写成公式1.15。

$$R(M, x) = \frac{x^T M x}{x^T x} \quad (1.15)$$

设  $M$  的特征值与特征向量分别为  $\lambda_1, \dots, \lambda_n$  和  $v_1, \dots, v_n$ ，且满足  $\lambda_{\min} = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n = \lambda_{\max}$ ，那么在  $M$  已知的情况下有：

$$\begin{aligned} \max_x R(M, x) &= \lambda_n \\ \min_x R(M, x) &= \lambda_1 \end{aligned} \quad (1.16)$$

以下为证明公式1.16的过程：

由于  $M$  是 Hermitian 矩阵，存在一个酉矩阵  $U$ ，满足公式1.17。

$$M = U A U^T \quad (1.17)$$

其中  $A = \text{diag}\{\lambda_1, \dots, \lambda_n\}$ 。

因此公式1.15可以转换如下：

$$\begin{aligned} R(M, x) &= \frac{x^T U A U^T x}{x^T x} \\ &= \frac{(U^T x)^T A (U^T x)}{x^T x} \end{aligned} \quad (1.18)$$

设  $P = U^T x$ ，则：

$$\begin{aligned} R(M, x) &= \frac{P^T A P}{P^T P} \\ &= \frac{\sum_{i=1}^n \lambda_i |P_i|^2}{\sum_{i=1}^n |P_i|^2} \end{aligned} \quad (1.19)$$

根据特征值的大小关系，我们可以得到不等式1.20。

$$\lambda_1 \sum_{i=1}^n |P_i|^2 \leq \sum_{i=1}^n \lambda_i |P_i|^2 \leq \lambda_n \sum_{i=1}^n |P_i|^2 \quad (1.20)$$

所以公式1.19的范围如下：

$$\lambda_1 \frac{\sum_{i=1}^n |P_i|^2}{\sum_{i=1}^n |x_i|^2} \leq R(M, x) \leq \lambda_n \frac{\sum_{i=1}^n |P_i|^2}{\sum_{i=1}^n |x_i|^2} \quad (1.21)$$

设  $U$  第  $i$  行第  $j$  列的元素为  $u_{ij}$ ，则  $U^T$  第  $i$  行第  $j$  列的元素为  $u_{ji}$ ，由  $P = U^T x$  和  $P^T = x^T U$  可得：

$$\begin{aligned} p_i &= \sum_{j=1}^n u_{ji} x_j \\ p_i^T &= \sum_{j=1}^n x_j u_{ij} \end{aligned} \quad (1.22)$$

则：

$$|p_i|^2 = p_i^T p_i = \sum_{j=1}^n \sum_{k=1}^n x_j u_{ij} u_{ki} x_k \quad (1.23)$$

于是：

$$\begin{aligned} \sum_{i=1}^n |p_i|^2 &= \sum_{i=1}^n p_i^T p_i \\ &= \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n x_j u_{ij} u_{ki} x_k \\ &= \sum_{j=1}^n \sum_{k=1}^n \left( \sum_{i=1}^n u_{ki} u_{ij} \right) x_j x_k \end{aligned} \quad (1.24)$$

因为  $U$  是酉矩阵，满足  $U^T U = I$ ，所以：

$$I_{jk} = \sum_{i=1}^n u_{ji} u_{ik} \quad (1.25)$$



其满足如下等式：

$$I_{jk} = \begin{cases} 1 & j = k \\ 0 & j \neq k \end{cases} \quad (1.26)$$

结合1.24和1.26可得如下等式：

$$\sum_{i=1}^n |p_i|^2 = \sum_{i=1}^n |x_i|^2 \quad (1.27)$$

代入公式1.21可得：

$$\lambda_1 \leq R(M, x) \leq \lambda_n \quad (1.28)$$

且：

$$R(M, x) = \begin{cases} \lambda_1 & x = v_1 \\ \lambda_n & x = v_n \end{cases} \quad (1.29)$$

如果用  $x' = cx$  代入公式1.15有：

$$\begin{aligned} R(M, x') &= \frac{x'^T M x'}{x'^T x'} \\ &= \frac{c^2 x^T M x}{c^2 x^T x} \\ &= \frac{x^T M x}{x^T x} \end{aligned} \quad (1.30)$$

由此可以看出对  $x$  进行缩放不影响瑞利商的值，即：

$$R(M, cx) = R(M, x) \quad (1.31)$$

因此我们可以限定  $x^T x = 1$ ，那么公式1.15可以简化为：

$$R(M, x) = x^T M x \quad (1.32)$$

那么  $R(M, x)$  的极值就可以转换成约束条件下的拉格朗日乘法，如公式1.33。

$$L(x, \lambda) = x^T M x - \lambda(x^T x - 1) \quad (1.33)$$

对  $x$  求导并置为 0 可得：

$$\nabla L(x, \lambda) = Mx - \lambda x = 0 \quad (1.34)$$

即  $M$  的特征值能使得瑞利商取极值，且有：

$$R(M, x) = \lambda \quad (1.35)$$

瑞利商可以推广至广义瑞利商 (Generalized Rayleigh Quotient)，其形式如公式1.37。

$$R(A, B, x) = \frac{x^H A x}{x^H B x} \quad (1.36)$$

其中  $A, B$  均为  $n \times n$  的 Hermitian 矩阵，且  $B$  为正定矩阵。

令  $x = B^{-\frac{1}{2}} x'$ ，广义瑞利商可以改写成：

$$\begin{aligned} R(A, B, x) &= \frac{(B^{-\frac{1}{2}} x')^H A (B^{-\frac{1}{2}} x')}{(B^{-\frac{1}{2}} x')^H B (B^{-\frac{1}{2}} x')} \\ &= \frac{x'^H (B^{-\frac{1}{2}})^H A B^{-\frac{1}{2}} x'}{x'^H (B^{-\frac{1}{2}})^H B B^{-\frac{1}{2}} x'} \\ &= \frac{x'^H (B^{-\frac{1}{2}})^H A B^{-\frac{1}{2}} x'}{x'^H x'} \end{aligned} \quad (1.37)$$

此时  $R(A, B, x)$  的最大特征值和最小特征值即为  $(B^{-\frac{1}{2}})^H A B^{-\frac{1}{2}}$  的最大和最小特征值。其实等价于当  $M = (B^{-\frac{1}{2}})^H A B^{-\frac{1}{2}}$  时的  $R(M, x')$ ， $x' = B^{\frac{1}{2}} x$ 。

为简单起见，我们可以令  $P = B^{-\frac{1}{2}}$ ，公式1.37可以写作：

$$\begin{aligned} R(A, B, x) &= \frac{x'^H (B^{-\frac{1}{2}})^H A B^{-\frac{1}{2}} x'}{x'^H x'} \\ &= \frac{x'^H P^H A P x'}{x'^H x'} \\ &= \frac{(P x')^H A P x'}{x'^H x'} \end{aligned} \quad (1.38)$$

类比上面提到的拉格朗日乘法，我们可以得到如下等式：

$$\nabla L(x, \lambda) = P^H A P x' - \lambda x' = 0 \quad (1.39)$$

代入  $x' = P^{-1} x$  有：

$$\nabla L(x, \lambda) = P^H A P P^{-1} x - \lambda P^{-1} x \quad (1.40)$$

解得：

$$PP^H Ax = \lambda x \quad (1.41)$$

又因为  $B^{-1} = PP^H$ ，所以最终求解特征值和特征向量可以依据：

$$B^{-1}Ax = \lambda x \quad (1.42)$$

## 1.3 EM 算法

本节来自于 Andrew Ng 的课堂讲义<sup>(12)</sup> 以及李航老师的《统计学习方法》<sup>(19)</sup>。

### 1.3.1 Jensen's Inequality

首先介绍一下 **Jensen 不等式**。

定义  $f$  为实域的函数，若  $f''(x) \geq 0$ ， $x \in R$ ，则  $f$  为凸函数。若自变量  $x$  为向量，则当  $f$  关于  $x$  的 Hessian 矩阵为半正定矩阵，即  $H \geq 0$  时， $f$  为凸函数；

如果对于所有的  $x \in R$ ，有  $f''(x) > 0$  或对于所有的向量  $x$ ，有  $H > 0$ ， $f$  为严格意义上的凸函数。那么 Jensen 不等式定义定理 1.1：

#### 定理 1.1. Jensen's Inequality

假定  $f$  是一个凸函数， $X$  是一个随机变量，那么：

$$E[f(X)] \geq f(EX) \quad (1.43) \heartsuit$$

且，如果  $f$  是严格凸函数，则当且仅当  $X = E[X]$  时定理 1.1 取等号。图 1.1 提供了一个理解 Jensen 不等式的例子。

图中实线部分  $f$  是一个凸函数，随机变量  $X$  分别以 0.5 的概率取值  $a$  和  $b$ ，因此随机变量  $X$  的期望为  $a$  和  $b$  的中间值，即：

$$EX = \frac{a+b}{2} \quad (1.44)$$

期望对应的函数值为  $f(EX)$ ，如图所示。而  $E[f(X)]$  等于：

$$E[f(X)] = \frac{f(a) + f(b)}{2} \quad (1.45)$$

从图上我们可以很明显的看出来二者的大小关系，又因为  $X$  是随机变量，当我们取  $X = EX$



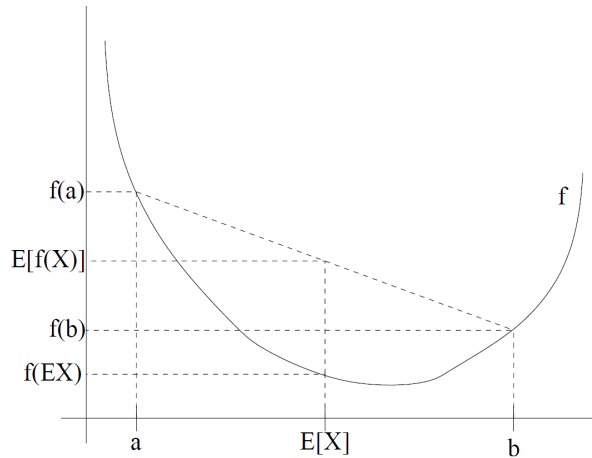


图 1.1: Jensen's Inequality 示例

的时候，函数的期望  $E[f(x)]$  和期望的函数  $f(EX)$  是相等的，所以有定理1.1。

### 1.3.2 EM 算法推导

一般地，用  $X$  表示观测随机变量的数据， $Z$  表示隐随机变量的数据， $X$  和  $Z$  连在一起称为完全数据。观测数据  $X$  又称为不完全数据。假设给定观测数据  $X$ ，其概率分布是  $P(X|\theta)$ ，其中  $\theta$  是需要估计的参数，那么不完全数据  $X$  的似然函数是  $P(X|\theta)$ ，对数似然函数是  $L(\theta) = \log P(X|\theta)$ ；假设  $X$  和  $Z$  的联合概率分布是  $P(X, Z|\theta)$ ，那么完全数据的对数似然函数是  $L(\theta) = \log P(X, Z|\theta)$ 。

假设有  $m$  个相互独立的训练样本  $\{x^{(1)}, \dots, x^{(m)}\}$ ，我们需要估计这  $m$  个样本概率分布的参数。若隐变量为  $z$ ，想要找到拟合模型  $p(x, z)$  的参数，其对数似然函数如公式1.46。

2005/06/2

由于隐变量  $z$  的存在，导致很难使用最大似然估计来得到参数  $\theta$ ，在这种情况下，EM 算法做的事情是给似然函数找个下界（E 步），然后最大化这个下界（M 步），如此反复迭代多次，最终收敛为止。

对于每一个样本  $x^{(i)}$ ，我们定义  $Q_i$  为隐随机变量  $z$  的分布，其满足：

$$\sum_z Q_i(z) = 1, Q_i(z) \geq 0 \quad (1.47)$$

下面，我们来推导一下 EM 算法的精髓过程，见公式1.48。

$$\begin{aligned}
 \sum_{i=1}^m (x_{(i)}; \theta) &= \sum_{i=1}^m \log \sum_{z^{(i)}} p(x^{(i)}, z; \theta) \\
 &= \sum_{i=1}^m \log \sum_{z^{(i)}} Q_i(z^{(i)}) \frac{p(x^{(i)}, z; \theta)}{Q_i(z^{(i)})} \\
 &\geq \sum_{i=1}^m \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z; \theta)}{Q_i(z^{(i)})}
 \end{aligned} \tag{1.48}$$

公式1.48的最后一步是通过 Jensen's Inequality 得到的。我们分析下这个推导。

首先我们知道对数函数为凹函数，因为对数函数的二阶导  $f''(x) = -1/x^2 < 0$ ，其中  $x \in R^+$ 。

而

$$\sum_{z^{(i)}} Q_i(z^{(i)}) \frac{p(x^{(i)}, z; \theta)}{Q_i(z^{(i)})}$$

是

$$\frac{p(x^{(i)}, z; \theta)}{Q_i(z^{(i)})}$$

关于服从分布  $Q_i$  的随机变量  $z_i$  期望。

由 Jensen's Inequality 可以得到公式1.49：

$$f\left(E_{z^{(i)} \sim Q_i} \left[ \frac{p(x^{(i)}, z; \theta)}{Q_i(z^{(i)})} \right]\right) \geq E_{z^{(i)} \sim Q_i} \left[ f\left( \frac{p(x^{(i)}, z; \theta)}{Q_i(z^{(i)})} \right) \right] \tag{1.49}$$

其中下角标  $z^{(i)} \sim Q_i$  表示这个期望是针对服从分布  $Q_i$  的随机变量  $z^{(i)}$ ，根据公式1.49就可以推出1.48的最后一步。

ok，关于公式1.49多说两句，这里面其实用到的公式就是 Jensen 不等式，但是这个式子和前面我们提到的定理1.1中的区别在于  $E[f(X)] \geq f(EX)$  中的  $X$  同样是个函数，其等于：

$$\frac{p(x^{(i)}, z; \theta)}{Q_i(z^{(i)})}$$

以上公式1.48给  $l\theta$  提供了一个下界，我们不断地最大化这个下界，就可以不断地增加  $l(\theta)$ ，说明每一次迭代模型都能找到拟合的更好的参数。

好的，下一个问题是，我们如何最大化这个下界呢？假设  $t$  轮迭代后的参数为  $\theta^{(t)}$ ，我们希望通过这个  $\theta^{(t)}$  来得到  $Q_i(z^{(i)})$ ，即隐变量的分布。而  $t$  轮迭代的时候一定是满足下界最大化的，

而最大化很明显就是变不等号为等号的临界部分，即定理1.1中等号成立的条件。所以第  $t$  轮得到的  $\theta^{(t)}$  必然满足公式1.50。

$$E_{z^{(i)} \sim Q_i} \left[ \frac{p(x^{(i)}, z; \theta)}{Q_i(z^{(i)})} \right] = \frac{p(x^{(i)}, z; \theta)}{Q_i(z^{(i)})} \quad (1.50)$$

展开之后有：

$$\begin{aligned} \sum_z Q_i(z^{(i)}) \frac{p(x^{(i)}, z; \theta)}{Q_i(z^{(i)})} &= \sum_z p(x^{(i)}, z; \theta) \\ &= \frac{p(x^{(i)}, z; \theta)}{Q_i(z^{(i)})} \end{aligned} \quad (1.51)$$

所以我们有：

$$\begin{aligned} Q_i(z^{(i)}) &= \frac{p(x^{(i)}, z; \theta)}{\sum_z p(x^{(i)}, z; \theta)} \\ &= \frac{p(x^{(i)}, z; \theta)}{p(x^{(i)}; \theta)} \\ &= p(z|x^{(i)}; \theta) \end{aligned} \quad (1.52)$$

上面的推导过程其实就是 EM 算法中的 E 步。

总结下 EM 算法步骤：

1. E 步：对于每一个样本，作：

$$Q_i := p(z|x^{(i)}; \theta) \quad (1.53)$$

2. M 步：

$$\theta := \arg \max_{\theta} \sum_{i=1}^m \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z; \theta)}{Q_i(z^{(i)})} \quad (1.54)$$

在用 EM 算法求模型参数的时候，首先定义模型的初始化参数，得到初始化参数之后，我们就可以得到隐随机变量  $z$  的分布，得到这个分布之后，我们再通过 M 步去算最大似然概率，得到新的参数，如此反复迭代，最终到 EM 算法收敛为止。



### 1.3.3 EM 算法收敛性证明

如何证明 EM 算法最终一定会收敛呢？假设  $\theta^{(t)}$  和  $\theta^{(t+1)}$  为连续的两轮迭代后的参数，我们只需要证明最大似然随着迭代次数是单调递增即可，即证明不等式 1.55。

$$\theta^{(t+1)} \geq \theta^{(t)} \quad (1.55)$$

根据  $t$  轮迭代后的参数  $\theta^{(t)}$ ，我们可以得到

$$Q_i^{(t)} = p(z|x^{(i)}; \theta^{(t)})$$

那么  $t$  步的似然度为（此时公式 1.48 的最后一步取等号）：

$$l(\theta^{(t)}) = \sum_{i=1}^m \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{p(x^{(i)}, z; \theta^{(t)})}{Q_i^{(t)}(z^{(i)})}$$

参数  $\theta^{(t+1)}$  则是通过最大化上述公式得到的。因此我们有：

$$l(\theta^{(t+1)}) \geq \sum_{i=1}^m \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{p(x^{(i)}, z; \theta^{(t+1)})}{Q_i^{(t)}(z^{(i)})} \quad (1.56)$$

$$\geq \sum_{i=1}^m \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{p(x^{(i)}, z; \theta^{(t)})}{Q_i^{(t)}(z^{(i)})} \quad (1.57)$$

$$= l(\theta^{(t)}) \quad (1.58)$$

不等式 1.56 是由公式 1.48 推导出来的，其对于任意的  $Q_i$  和  $\theta$  都是成立的，自然对于  $Q_i = Q_i^{(t)}$  和  $\theta = \theta^{(t)}$  也成立。而不等式 1.57 是通过 EM 算法中的 M 步得到的，因为：

$$\theta^{(t+1)} = \arg \max_{\theta} \sum_{i=1}^m \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{p(x^{(i)}, z; \theta)}{Q_i^{(t)}(z^{(i)})} \quad (1.59)$$

公式 1.58 即为 Jensen 不等式处于临界条件时的结果。

由此可以得知  $l(\theta)$  是随着迭代次数的增加单调递增的，即似然度是单调收敛的。因此 EM 算法的终止条件是似然度收敛，即似然度不会再增加了为止。但是实际情况中一般是看两个相邻的迭代过程似然度的差是否低于一个预设的值，如果低于该值，则停止迭代。

定义公式 1.60：

$$J(Q, \theta) = \sum_{i=1}^m \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z; \theta)}{Q_i(z^{(i)})} \quad (1.60)$$

我们知道  $l(\theta) \geq J(Q, \theta)$ ，我们也可以把 EM 算法看成一个坐标上升的过程，在 E 步的时候，相当于固定  $\theta$  求  $Q$ ；在 M 步的时候，相当于固定了  $Q$  求  $\theta$ 。

**题外话：**这个固定一个变量，最大化另外一个变量，通过另外一个变量求这个变量这个思想感觉跟 GAN 很类似啊……GAN 里面也是这样的，先固定生成器，优化判别器；然后固定判别器，优化生成器。有时间把这两个联系到一起来想想，还挺有意思的……

## 1.4 混合高斯分布

### 定义 1.13. 方差

方差用于描述数据的离散或波动程度。假定变量为  $X$ ，均值为  $\bar{X}$ ， $N$  为总体样本数，方差计算公式如下：

$$\text{var}(X) = \frac{\sum_{i=1}^N (X_i - \bar{X})^2}{N - 1} \quad (1.61)$$

### 定义 1.14. 协方差

协方差表示了变量线性相关的方向，取值范围是  $[-\infty, \infty]$ ，一般来说协方差为正值，说明一个变量变大另一个变量也变大；取负值说明一个变量变大另一个变量变小，取 0 说明两个变量没有相关关系。

$$\text{cov}(X) = \frac{\sum_{i=1}^N (X_i - \bar{X})(Y_i - \bar{Y})}{N - 1} \quad (1.62)$$

### 定义 1.15. 相关系数

协方差可反映两个变量之间的相互关系及相关方向，但无法表达其相关的程度，皮尔逊相关系数不仅表示线性相关的方向，还表示线性相关的程度，取值  $[-1, 1]$ ，也就是说，相关系数为正值，说明一个变量变大另一个变量也变大；取负值说明一个变量变大另一个变量变小，取 0 说明两个变量没有相关关系，同时，相关系数的绝对值越接近 1，线性关系越显著。

$$\rho_{XY} = \frac{\text{cov}(X, Y)}{\sqrt{DX}\sqrt{DY}} \quad (1.63)$$

### 定义 1.16. 协方差矩阵

当  $X \in R^n$  为高维数据时，协方差矩阵可以很好的反映数据的性质，在协方差矩阵中，对角线元素反映了数据在各个维度上的离散程度，协方差矩阵为对角阵，非对角线元素反映

了数据各个维度的相关性，其形式如下：

$$\Sigma = \begin{bmatrix} \text{cov}(x_1, x_1) & \text{cov}(x_1, x_2) & \cdots & \text{cov}(x_1, x_n) \\ \text{cov}(x_2, x_1) & \text{cov}(x_2, x_2) & \cdots & \text{cov}(x_2, x_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(x_n, x_1) & \text{cov}(x_n, x_2) & \cdots & \text{cov}(x_n, x_n) \end{bmatrix} \quad (1.64)$$

单变量高斯分布公式如1.65，其中  $\mu$  和  $\sigma^2$  分别为均值和方差。

$$\mathcal{N}(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(x - \mu)^2}{2\sigma^2} \right\} \quad (1.65)$$

多变量高斯分布公式如1.66，其中  $\mu$  和  $\Sigma$  分别为均值和协方差矩阵。

$$\mathcal{N}(x; \mu, \Sigma) = \frac{1}{(2\pi)^{-\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\} \quad (1.66)$$

混合高斯模型 (Gaussian Mixture Model) 表示的是多个高斯分布叠加在一起的分布，其公式如1.67，其中  $K$  为高斯分量的个数， $\pi_k$  为各个分量的权重，其满足  $0 \leq \pi_k \leq 1$ ，且  $\sum_{k=1}^K \pi_k = 1$ 。 $p(x)$  表示的是多个高斯分量加权后的分布。

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \quad (1.67)$$

那么给定一堆训练数据，我们如何根据这些数据来得到 GMM 的参数呢？

假设训练集  $\{x^{(1)}, \dots, x^{(m)}\}$ ，由于估计 GMM 的算法为无监督学习算法，因此这些数据都是没有标签的。我们用一个联合分布模型来对这些数据进行建模，即：

$$p(x^{(i)}, z^{(i)}) = p(x^{(i)} | z^{(i)}) p(z^{(i)}) \quad (1.68)$$

其中， $z^{(i)} \sim \text{Multinomial}(\phi)$ ， $s_j \geq 0$ ，且  $\sum_{j=1}^k \phi_j = 1$ ， $\phi_j$  表示的是  $p(z^{(i)} = j)$ ，此外  $p(x^{(i)} | z^{(i)} = j) \sim \mathcal{N}(\mu_j, \Sigma_j)$ 。 $k$  表示的是  $z^{(i)}$  所能取的值的个数。

我们选取的模型假设每个  $x^{(i)}$  都由随机选取的  $z^{(i)}$  生成， $z^{(i)} \in \{1, \dots, k\}$ 。我们可以理解成 GMM 模型有  $k$  个高斯分量， $z$  就是表征每个分量权重的随机变量，因此  $z$  就是隐变量。 $p(z^{(i)}) = \phi_j$  就是  $x^{(i)}$  这个数据来源于分量  $j$  的概率。

综上，GMM 模型的参数为  $\mu, \Sigma, \phi$ 。现在的问题就是如何求解这些参数。我们可以得到似然

函数如公式1.69。

$$\begin{aligned}
 l(\mu, \Sigma, \phi) &= \sum_{i=1}^m \log p(x^{(i)}; \mu, \Sigma, \phi) \\
 &= \sum_{i=1}^m \log \sum_{z^{(i)}=1}^k p(x^{(i)}|z^{(i)}; \mu, \Sigma) p(z^{(i)}; \phi)
 \end{aligned} \tag{1.69}$$

如果我们直接对这个似然函数进行求导，由于隐变量的存在，我们无法得到关于参数的确切解。那么既然问题出在隐变量上，如果我们已经知道了每一个  $x^{(i)}$  所属的分量，公式1.69就可以写成公式1.70，需注意此时隐变量的分布我们还是不知道，我们只知道训练集中的样本分别属于哪一个高斯分量。

$$\begin{aligned}
 l(\mu, \Sigma, \phi) &= \sum_{i=1}^m \log p(x^{(i)}|z^{(i)}; \mu, \Sigma) p(z^{(i)}; \phi) \\
 &= \sum_{i=1}^m \left[ \log p(x^{(i)}|z^{(i)}; \mu, \Sigma) + \log p(z^{(i)}; \phi) \right]
 \end{aligned} \tag{1.70}$$

此时用最大似然概率求解的办法就可以得到这些参数，我们写细一点，求解一下看看。

首先，求解  $\phi_j$ ，在此之前，我们可以看出  $l$  中与  $\phi$  相关的项为：

$$\sum_{i=1}^m \log p(z^{(i)}; \phi)$$

我们知道  $\phi$  是服从多项式分布的， $\phi$  受限于  $\sum_{j=1}^k \phi_j = 1$ ，因此关于  $\phi$  的求解我们需要构建一个拉格朗日函数如下：

$$\mathcal{L}(\phi) = \sum_{i=1}^m \log p(z^{(i)}; \phi) + \beta \left( \sum_{j=1}^k \phi_j - 1 \right)$$

因此关于求解  $\phi$  的步骤如公式1.71：

$$\begin{aligned}
 \frac{\partial \mathcal{L}(\phi)}{\partial \phi_j} &= \frac{\partial \sum_{i=1}^m \log p(z^{(i)}; \phi) + \beta (\sum_{j=1}^k \phi_j - 1)}{\partial \phi_j} \\
 &= \sum_{i=1}^m \frac{\partial \log p(z^{(i)}; \phi)}{\partial \phi_j}
 \end{aligned} \tag{1.71}$$

因为不一定每个  $x^{(i)}$  都有  $j$  分量产生，我们定义  $1\{z^{(i)} = j\}$  如下，也就是说第  $i$  个样本由第

$j$  个样本产生的时候，我们取为 1。

$$1\{z^{(i)} = j\} = \begin{cases} 1 & \text{if } z^{(i)} = j \\ 0 & \text{otherwise} \end{cases} \quad (1.72)$$

公式1.71可以写成公式1.73：

$$\begin{aligned} \frac{\partial \mathcal{L}(\phi)}{\partial \phi_j} &= \frac{\partial \sum_{i=1}^m \log p(z^{(i)}; \phi) + \beta(\sum_{j=1}^k \phi_j - 1)}{\partial \phi_j} \\ &= \sum_{i=1}^m 1\{z^{(i)} = j\} \frac{\partial \log \phi_j + \beta \phi_j}{\partial \phi_j} \\ &= \sum_{i=1}^m 1\{z^{(i)} = j\} \left[ \frac{1}{\phi_j} + \beta \right] \end{aligned} \quad (1.73)$$

令公式1.73为 0，解得：

$$\phi_j = \frac{\sum_{i=1}^m 1\{z^{(i)} = j\}}{-\beta} \quad (1.74)$$

由此可知： $\phi_j \propto \sum_{i=1}^m 1\{z^{(i)} = j\}$ ，又因为  $\sum_{j=1}^k \phi_j = 1$ ，将  $\phi_j$  代入约束条件有：

$$\sum_{j=1}^k \phi_j = \frac{\sum_{j=1}^k \sum_{i=1}^m 1\{z^{(i)} = j\}}{-\beta} = 1 \quad (1.75)$$

即

$$-\beta = \sum_{j=1}^k \sum_{i=1}^m 1\{z^{(i)} = j\} \quad (1.76)$$

也就是说  $-\beta$  等于属于各个高斯分量的训练样本的和，也就是  $m$ 。至此，我们可以得到如下公式：

$$\phi_j = \frac{1}{m} \sum_{i=1}^m 1\{z^{(i)} = j\} \quad (1.77)$$



接着我们来求  $\mu$ ，如公式1.78。

$$\begin{aligned}
 \frac{\partial l(\mu, \Sigma, \phi)}{\partial \mu_j} &= \frac{\partial \sum_{i=1}^m \left[ \log p(x^{(i)}|z^{(i)}; \mu, \Sigma) + \log p(z^{(i)}; \phi) \right]}{\partial \mu_j} \\
 &= \frac{\partial \sum_{i=1}^m \log p(x^{(i)}|z^{(i)}; \mu, \Sigma)}{\partial \mu_j} \\
 &= \frac{\partial \sum_{i=1}^m 1\{z^{(i)} = j\} \log p(x^{(i)}|z^{(i)}; \mu, \Sigma)}{\partial \mu_j} \\
 &= \frac{\partial \sum_{i=1}^m 1\{z^{(i)} = j\} \log \mathcal{N}(x^{(i)}; \mu_j, \Sigma_j)}{\partial \mu_j} \\
 &= \frac{\partial \sum_{i=1}^m 1\{z^{(i)} = j\} \log \frac{1}{(2\pi)^{-\frac{n}{2}} |\Sigma_j|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \right\}}{\partial \mu_j} \quad (1.78) \\
 &= \frac{\partial \sum_{i=1}^m 1\{z^{(i)} = j\} \left[ -\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \right]}{\partial \mu_j} \\
 &= \frac{\partial \sum_{i=1}^m 1\{z^{(i)} = j\} \left[ -\frac{1}{2} (x^{(i)T} \Sigma_j^{-1} x^{(i)} - \mu_j^T \Sigma_j^{-1} x^{(i)} - x^{(i)T} \Sigma_j^{-1} \mu_j + \mu_j^T \mu_j) \right]}{\partial \mu_j} \\
 &= \frac{\partial \sum_{i=1}^m 1\{z^{(i)} = j\} \left[ -\frac{1}{2} (-\mu_j^T \Sigma_j^{-1} x^{(i)} - x^{(i)T} \Sigma_j^{-1} \mu_j + \mu_j^T \mu_j) \right]}{\partial \mu_j} \\
 &= \sum_{i=1}^m 1\{z^{(i)} = j\} \Sigma_j^{-1} (x^{(i)} - \mu_j)
 \end{aligned}$$

令公式1.78为零，则有：

$$\sum_{i=1}^m 1\{z^{(i)} = j\} \Sigma_j^{-1} (x^{(i)} - \mu_j) = 0 \quad (1.79)$$

解得：

$$\mu_j = \frac{\sum_{i=1}^m 1\{z^{(i)} = j\} x^{(i)}}{\sum_{i=1}^m 1\{z^{(i)} = j\}} \quad (1.80)$$

最后，根据方差的定义，我们可以计算出方差的值，如公式1.81

$$\Sigma_j = \frac{\sum_{i=1}^m 1\{z^{(i)} = j\} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^m 1\{z^{(i)} = j\}} \quad (1.81)$$

上面我们算了一下在已知每个样本所属的高斯分量的情况下参数  $\phi, \mu, \Sigma$  的估计值。回到最初那个问题，训练数据是没有标签的，那么怎么估计参数？

只要知道了每一个训练样本所属的高斯分量，我们就可以通过最大似然估计求出参数。既

然如此，我们可以先随机的将训练样本进行分配，初始化时，随机指定每一个训练样本所属的高斯分量，结合1.3中讲的 EM 算法，我们可以算出  $\phi, \mu, \Sigma$ ，再用这个参数来更新每一个样本所属的高斯分量，因为在得到参数后，我们将训练样本  $x^{(i)}$  代入到 GMM 中，可以得到该样本分别属于每一个高斯分量的概率值，即求取  $p(z^{(i)} = j | x^{(i)}; \phi, \mu, \Sigma), j \in \{1, \dots, K\}$ 。得到这些概率之后，我们认为每一个高斯分量以权重  $p(z^{(i)} = j | x^{(i)}; \phi, \mu, \Sigma)$  生成了样本  $x^{(i)}$ ，知道这些权重之后就可以继续更新参数，这样往往复复，不停的迭代，就会不停的朝着正确分类的方向走去，最终收敛即可。

下面介绍下 GMM 参数估计的 EM 算法：

1. E 步：对于  $i, j$ ，作：

$$w_j^{(i)} := p(z^{(i)} = j | x^{(i)}; \phi, \mu, \Sigma) \quad (1.82)$$

2. M 步：

$$\phi_j := \frac{1}{m} \sum_{i=1}^m w_j^{(i)} \quad (1.83)$$

$$\mu_j := \frac{\sum_{i=1}^m w_j^{(i)} x^{(i)}}{\sum_{i=1}^m w_j^{(i)}} \quad (1.84)$$

$$\Sigma := \frac{\sum_{i=1}^m w_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^m w_j^{(i)}} \quad (1.85)$$

## 1.5 HMM 相关知识总结

本节笔记来自于 CS229 课程中 Daniel Ramage 的 Hidden Markov Models Fundamentals<sup>(15)</sup>。

### 1.5.1 Markov Models

给定状态集合  $S = \{s_1, s_2, \dots, s_{|S|}\}$ ，我们可以观测到沿着时间线的序列  $\vec{z} \in S^T$ 。比方说一个天气系统  $S = \{sun, cloud, rain\}$ ，其  $|S| = 3$ 。假设连续五天观察到的天气序列为  $\{z_1 = s_{sun}, z_2 = s_{cloud}, z_3 = s_{cloud}, z_4 = s_{rain}, z_5 = s_{cloud}\}$ 。

上述天气系统中观测到的状态序列是时间线上的随机过程，如果不作进一步假设的话， $t$  时刻的状态  $s_j$  可能是任意变量的函数，不仅仅包括从时刻 1 到时刻  $t-1$  的状态，甚至还包括一些没有建模的变量，所以我们作出两个 **Markov Assumptions** 使模型更可行：

1. **Limited Horizon Assumption**：假设  $t$  时刻的处于某个状态的概率只取决于  $t-1$  时刻的状态。本假设的直观感受是： $t-1$  时刻的状态囊括了对过去时刻状态的总结，因而可以用于

预测  $t$  时刻的输出。即：

$$P(z_t | z_{t-1}, z_{t-2}, \dots, z_1) = P(z_t | z_{t-1}) \quad (1.86)$$

**2. Stationary Process Assumption:** 给定当前状态，下一个时刻的状态分布不随着时间变化而变化。即：

$$P(z_t | z_{t-1}) = P(z_2 | z_1) \quad (1.87)$$

其中  $z_{t-1} = z_1$ 。

按照惯例，我们假设存在初始状态和初始观测  $z_0 \equiv s_0$ ， $s_0$  表示的是在  $t = 0$  处的初始概率分布。这么表达初始状态方便我们计算第一个真实状态的先验概率为  $P(z_1 | z_0)$ ，且可以得到  $P(z_t | z_{t-1}, \dots, z_1) = P(z_t | z_{t-1}, \dots, z_1, z_0)$ 。此外，也有表达方式将这些先验置信度表示为  $\pi \in R^{|S|}$ 。

我们定义一个状态转移矩阵  $A \in R^{(|S|+1) \times (|S|+1)}$  来表示各个状态之间的转移，其中  $A_{ij}$  表示的是任意时刻从状态  $i$  转移到状态  $j$  的概率。对于上述提到的天气系统，其转移矩阵可能如下：

$$A = \begin{matrix} & \begin{matrix} s_0 & s_{sun} & s_{cloud} & s_{rain} \end{matrix} \\ \begin{matrix} s_0 \\ s_{sun} \\ s_{cloud} \\ s_{rain} \end{matrix} & \begin{bmatrix} 0 & 0.33 & 0.33 & 0.33 \\ 0 & 0.8 & 0.1 & 0.1 \\ 0 & 0.2 & 0.6 & 0.2 \\ 0 & 0.1 & 0.2 & 0.7 \end{bmatrix} \end{matrix} \quad (1.88)$$

上面这个编出来的矩阵说明了天气是自我关联的，如果今天是晴天，明天还是晴天的概率会比较大；如果今天多云，明天还是多云的概率就比较大。这个规律在很多 Markov Model 里都有，从转移矩阵的对角线上也可以看出来。在这个天气系统中，初始化天气是均匀分布的。

### 1.5.2 Markov Model 的两个基本问题

将上节中讲到的 **Markov Assumption** 和状态转移矩阵  $A$  结合起来，我们就可以回答 Markov Chain 中关于状态序列的两个基本问题了。

- 给定一个状态序列  $\vec{z}$ ，如何求这个序列的概率？
- 如何通过最大化一个观测序列  $\vec{z}$  的似然概率来估计模型的参数？

### 1.5.2.1 状态序列的概率

给定一个状态序列  $\vec{z}$ ，其长度为  $t$ 。我们通过概率的链式法则（Chain Rule）来求这个状态序列出现的概率，如公式1.89。

$$\begin{aligned}
 P(\vec{z}) &= P(z_t, z_{t-1}, \dots, z_1; A) \\
 &= P(z_t, z_{t-1}, \dots, z_1, z_0; A) \\
 &= P(z_t | z_{t-1}, \dots, z_1, z_0; A) P(z_{t-1} | z_{t-2}, \dots, z_1, z_0; A) \dots P(z_1 | z_0; A) \\
 &= P(z_t | z_{t-1}; A) P(z_{t-1} | z_{t-2}; A) \dots P(z_1 | z_0; A) \\
 &= \prod_{t'=1}^t P(z_{t'} | z_{t'-1}; A) \\
 &= \prod_{t'=1}^t A_{z_{t'-1} z_{t'}}
 \end{aligned} \tag{1.89}$$

根据公式1.89，我们可以求一下上节中提到的天气序列的概率：

$$\begin{aligned}
 &P(z_1 = s_{sun}, z_2 = s_{cloud}, z_3 = s_{cloud}, z_4 = s_{rain}, z_5 = s_{cloud}) \\
 &= P(s_{sun} | s_0) P(s_{cloud} | s_{sun}) P(s_{cloud} | s_{cloud}) P(s_{rain} | s_{cloud}) P(s_{cloud} | s_{rain}) \\
 &= 0.33 \times 0.1 \times 0.2 \times 0.7 \times 0.2 \\
 &= 0.000924
 \end{aligned} \tag{1.90}$$

### 1.5.2.2 Markov Model 的最大似然估计

出于学习的目的，我们希望找到参数  $A$  能够最大化长度为  $T$  的观测序列  $\vec{z}$  的最大对数似然。比方说求出天气系统中的从 sunny 转移到 cloudy 或者从 sunny 转移到 sunny 的概率等等，以使得观测序列最有可能出现。定义对数似然函数如公式1.91。

$$\begin{aligned}
 l(A) &= \log P(\vec{z}; A) \\
 &= \log \prod_{t=1}^T A_{z_{t-1} z_t} \\
 &= \sum_{t=1}^T \log A_{z_{t-1} z_t}
 \end{aligned} \tag{1.91}$$

我们要估计的是  $A$ ，所以我们需要将  $A$  的每一个元素  $A_{ij}$  都呈现在对数似然函数中，可是就从公式1.91中，我们没法呈现  $A_{ij}$ 。当然  $\vec{z}$  是已知的，即每个时刻的状态是确定的，但是我们没办法直接去指定  $\vec{z}$  中的每一个元素去估计  $A$ ，因为我们希望得到的是一般意义的最大似然函

数，即不论  $\vec{z}$  中都有哪些状态，公式1.91都可以表示针对样本序列  $\vec{z}$  的最大似然函数。因此我们可以定义一个 **indicator function** 如公式1.92（这个函数在 GMM 那一讲中也出现过）。

$$1\{z_t = j \wedge z_{t-1} = i\} = \begin{cases} 1 & \text{if } z_t = j \text{ and } z_{t-1} = i \\ 0 & \text{otherwise} \end{cases} \quad (1.92)$$

这个函数的意义是当  $t-1$  时刻状态为  $i$ ，并且  $t$  时刻状态为  $j$ ，其取值为 1。这样就和  $A_{ij}$  对应起来了，而且无需担心观测序列  $\vec{z}$  中是否有前一时刻为  $i$ 、当前时刻为  $j$  的存在。那么我们可以继续推导公式1.91，即枚举所有可能出现的序列，而对应的状态序列是否存在则由  $\vec{z}$  来决定，如果某个状态转移存在，那么其状态转移的概率就会被保留下来用于计算最大似然函数，如果不存在，那么其状态转移概率就会被扔掉。

$$\begin{aligned} l(A) &= \sum_{t=1}^T \log A_{z_{t-1}z_t} \\ &= \sum_{i=1}^{|S|} \sum_{j=1}^{|S|} \sum_{t=1}^T 1\{z_t = j \wedge z_{t-1} = i\} \log A_{ij} \end{aligned} \quad (1.93)$$

由于  $1\{z_t = j \wedge z_{t-1} = i\}$  的限制，公式1.93一定会迫使着最大似然函数是对应着观测序列  $\vec{z}$  的。

为了解决这个优化问题，我们需要保证最终求出来的转移矩阵  $A$  是有效的。那么针对每一个状态，从该状态跳转到下一个时刻状态的概率和应当等于 1，当然， $A$  中的每一个元素都必须是非负的。因此，我们用拉格朗日算子来解决这个问题，如公式1.94。

$$\begin{aligned} A &= \arg \max_A l(A) \\ s.t. \quad &\sum_{j=1}^{|S|} A_{ij} = 1, i \in \{1, 2, \dots, |S|\} \\ &A_{ij} \geq 0, i, j \in \{1, 2, \dots, |S|\} \end{aligned} \quad (1.94)$$

我们可以构造一个拉格朗日函数来解决这个受约束的最优化问题，如公式1.95。

$$\begin{aligned} \mathcal{L}(A) &= l(A) + \sum_{i=1}^{|S|} \alpha_i \left(1 - \sum_{j=1}^{|S|} A_{ij}\right) \\ &= \sum_{i=1}^{|S|} \sum_{j=1}^{|S|} \sum_{t=1}^T 1\{z_t = j \wedge z_{t-1} = i\} \log A_{ij} + \sum_{i=1}^{|S|} \alpha_i \left(1 - \sum_{j=1}^{|S|} A_{ij}\right) \end{aligned} \quad (1.95)$$

求偏微分并置 0，有：

$$\begin{aligned}\frac{\partial \mathcal{L}(A)}{\partial A_{ij}} &= \frac{\partial \sum_{i=1}^{|S|} \sum_{j=1}^{|S|} \sum_{t=1}^T 1\{z_t = j \wedge z_{t-1} = i\} \log A_{ij} + \sum_{i=1}^{|S|} \alpha_i (1 - \sum_{j=1}^{|S|} A_{ij})}{\partial A_{ij}} \\ &= \frac{1}{A_{ij}} \sum_{t=1}^T 1\{z_t = j \wedge z_{t-1} = i\} + \alpha_i \\ &= 0\end{aligned}\tag{1.96}$$

解得：

$$A_{ij} = \frac{1}{\alpha_i} \sum_{t=1}^T 1\{z_t = j \wedge z_{t-1} = i\}\tag{1.97}$$

又因为：

$$\sum_{j=1}^{|S|} A_{ij} = 1, i \in \{1, 2, \dots, |S|\}\tag{1.98}$$

$$\tag{1.99}$$

代入公式1.97求得的  $A_{ij}$  有：

$$\sum_{j=1}^{|S|} \frac{1}{\alpha_i} \sum_{t=1}^T 1\{z_t = j \wedge z_{t-1} = i\} = 1\tag{1.100}$$

解得：

$$\begin{aligned}\alpha_i &= \sum_{j=1}^{|S|} \sum_{t=1}^T 1\{z_t = j \wedge z_{t-1} = i\} \\ &= \sum_{t=1}^T 1\{z_{t-1} = i\}\end{aligned}\tag{1.101}$$

所以：

$$A_{ij} = \frac{\sum_{t=1}^T 1\{z_t = j \wedge z_{t-1} = i\}}{\sum_{t=1}^T 1\{z_{t-1} = i\}}\tag{1.102}$$

以上为 Markov Model 在已知观测序列的条件下估计状态转移矩阵  $A$  的推导。从公式1.102中，我们可以直观的看到  $A_{ij}$  就等于观测序列中，状态  $i$  跳转到状态  $j$  的次数除以位于状态  $i$  的个数。

### 1.5.3 Hidden Markov Model

Markov Model 是处理时序数据的一大利器，但是有些时候我们是无法观测状态序列的，在这种情况下，Markov Model 就无法使用了。在状态序列无法观测的情况下，我们知道这些状态的概率函数，那么如何利用这些状态序列的概率函数来对模型建模呢？比如说语音识别中的，我们能观测到的是根据词发出的声音，但是其对应的词我们是不知道的，如何根据这些声音的声学特征来估计对应的词呢？

为了介绍 HMM，我们借用下 Jason Eisner 论文中提出的 “Ice Cream Climatology”<sup>(4)</sup>，此处附上原文：

The situation: You are a climatologist in the year 2799, studying the history of global warming. You can't find any records of Baltimore weather, but you do find my (Jason Eisner's) diary, in which I assiduously recorded how much ice cream I ate each day. What can you figure out from this about the weather that summer?

我们可以用 HMM 来解决这个问题。这个问题中模型的状态的序列（每一天的天气）是未知的，我们只能观测到状态的输出序列（每天吃多少冰淇淋）。

正式地，HMM 是一种 Markov Model，其有一系列的观测输出  $x = \{x_1, x_2, \dots, x_T\}, x_t \in V, t = 1, \dots, T$ ，其中  $V = \{v_1, v_2, \dots, v_{|V|}\}$  为观测输出集合。同样有状态序列  $z = z_1, z_2, \dots, z_T, z_t \in S$ ，其中  $S = \{s_1, s_2, \dots, s_{|S|}\}$  为状态的集合。在 HMM 中，状态是不可观测的。状态之间的转移概率仍然由状态转移矩阵  $A$  表示。

同时生成观测输出的概率为隐状态的函数。为了满足这个条件，我们假设 HMM 中所有的输出是相互独立的，且定义：

$$P(x_t = v_k | z_t = s_j) = P(x_t = v_k | x_1, \dots, x_T, z_1, \dots, z_T) \quad (1.103)$$

$$= B_{jk} \quad (1.104)$$

矩阵  $B$  表示的是在某个时刻处于状态  $s_j$  生成观测  $v_k$  的概率值。

回归到上面讲的天气例子，假设有连续四天的冰淇淋消耗量的日志  $\vec{x} = \{x_1 = v_3, x_2 = v_2, x_3 = v_1, x_4 = v_2\}$ ，对应的  $V = \{v_1 = 1icecream, v_2 = 2icecreams, v_3 = 3icecreams\}$ 。那么我们要解决 HMM 的哪些问题呢？

### 1.5.4 HMM 的三个基本问题

HMM 有三个基本问题，

- 某个特定观测序列的概率值为多少（我们有多大的可能性看到连续四天消耗冰淇淋的量为 3,2,1,2）？

- 生成某个特定观测序列的最有可能的状态序列是什么（那四天的冰淇淋消耗量对应的状态序列是多少）？
- 给定一些观测序列，怎么去估计 HMM 的参数  $A$  和  $B$ ？

### 1.5.4.1 概率计算问题

概率计算问题指的是已知观测序列，求的是 HMM 模型输出该观测序列的概率值，此时，状态转移矩阵  $A$  和发射矩阵  $B$  都是已知的。

在 HMM 中，数据生成的过程如下：

假定状态序列  $\vec{z}$  由一个 Markov Model 生成，其状态转移矩阵为  $A$ 。在每一个时间步  $t$ ，我们选定输出  $x_t$ ， $x_t$  是关于  $z_t$  的函数。为了得到观测序列的概率，给定每一个可能的观测序列，我们需要加上数据  $\vec{x}$  的似然度，见公式 1.105

$$\begin{aligned} P(\vec{x}; A, B) &= \sum_{\vec{z}} P(\vec{x}, \vec{z}; A, B) \\ &= \sum_{\vec{z}} P(\vec{x}|\vec{z}; A, B)P(\vec{z}; A; B) \end{aligned} \quad (1.105)$$

公式 1.105 对于任意概率分布都成立；然而 HMM 的假设简化这个表达式，如公式 1.106。

$$\begin{aligned} P(\vec{x}; A, B) &= \sum_{\vec{z}} P(\vec{x}|\vec{z}; A, B)P(\vec{z}; A; B) \\ &= \sum_{\vec{z}} \left( \prod_{t=1}^T P(x_t|z_t; B) \right) \left( \prod_{t=1}^T P(z_t|z_{t-1}; A) \right) \\ &= \sum_{\vec{z}} \left( \prod_{t=1}^T B_{z_t x_t} \right) \left( \prod_{t=1}^T A_{z_{t-1} z_t} \right) \end{aligned} \quad (1.106)$$

在参数已知的情况下，公式 1.106 是一个比较简单的表达式。HMM 有三大假设：

- output independent assumption;
- markov assumption;
- stationary process assumption.

但是在计算公式 1.106 的过程中，由于求和是针对每一个可能的  $\vec{z}$  的，而任意时刻的  $z_t$  可以取  $|S|$  个值，计算上述的求和需要  $O(|S|^T)$  次操作。一般我们使用一种动态规划算法：前后向算法来计算  $P(\vec{x}; A, B)$ 。

前向算法定义  $\alpha_i(t) = P(x_1, x_2, \dots, x_t, z_t = s_i; A, B)$ ， $\alpha_t$  表示的是直到  $t$  时刻的观测序列的概率



值，且此时的状态为  $s_i$ 。那么观测序列的概率值为：

$$\begin{aligned}
 P(\vec{x}; A, B) &= P(x_1, x_2, \dots, x_T; A, B) \\
 &= \sum_i^{|S|} P(x_1, x_2, \dots, x_T, z_T = s_i; A, B) \\
 &= \sum_i^{|S|} \alpha_i(T)
 \end{aligned} \tag{1.107}$$

算法1.1描述了前向算法的迭代过程，其提供了一个效率更高的方式来计算  $\alpha_i(t)$ 。由算法可知，在每个时间步只需要进行  $O(|S|)$  次运算，最终计算出  $P(\vec{x}; A, B)$  需要  $O(|S|T)$  次运算。

---

**Algorithm 1.1** HMM 中的前向算法
 

---

1. Base Case:  $\alpha_i(0) = A_{0i}, i = 1, 2, \dots, |S|$ ;
  2. Recursion:  $\alpha_j(t) = \sum_{i=1}^{|S|} \alpha_i(t-1) A_{ij} B_{jx_t}, j = 1, \dots, |S|, t = 1, \dots, T$
- 

类似的算法是后向算法，其定义后向概率  $\beta_i(t) = P(x_T, x_{T-1}, \dots, x_{t+1}, z_t = s_i; A, B)$ 。我们同样可以得到：

$$\begin{aligned}
 P(\vec{x}; A, B) &= P(x_1, x_2, \dots, x_T; A, B) \\
 &= \sum_i^{|S|} P(x_T, x_{T-1}, \dots, x_1, z_0 = s_i; A, B) \\
 &= \sum_i^{|S|} \beta_i(0)
 \end{aligned} \tag{1.108}$$

其复杂度与前向算法相同，如算法1.2。

---

**Algorithm 1.2** HMM 中的后向算法
 

---

1. Base Case:  $\beta_i(T) = 1, i = 1, 2, \dots, |S|$ ;
  2. Recursion:  $\beta_j(t) = \sum_{i=1}^{|S|} A_{ji} B_{jx_t} \beta_i(t+1), j = 1, \dots, |S|, t = 1, \dots, T$
-

### 1.5.4.2 预测问题

HMM 中最常见的一个问题是：给定观测序列  $\vec{x} \in V^T$ ，最有可能产生  $\vec{x}$  的状态序列  $\vec{z}$  是什么？正式地：

$$\arg \max_{\vec{z}} P(\vec{z}|\vec{x}; A, B) = \arg \max_{\vec{z}} \frac{P(\vec{z}, \vec{x}; A, B)}{P(\vec{x}; A, B)} \quad (1.109)$$

$$= \arg \max_{\vec{z}} \frac{P(\vec{x}, \vec{z}; A, B)}{\sum_{\vec{z}} P(\vec{x}, \vec{z}; A, B)} \quad (1.110)$$

$$= \arg \max_{\vec{z}} P(\vec{x}, \vec{z}; A, B) \quad (1.111)$$

公式1.109和公式1.110利用的是贝叶斯定理和全概率公式，由公式1.109可知，其分母与  $\vec{z}$  无关，因此式1.111成立。我们可以尝试所有可能的  $\vec{z}$ ，然后得到每一个  $\vec{z}$  的后验概率，选取其中最大值作为上述公式的解，然而这种解法同样需要  $|S|^T$  次操作。参考前面提到的前向算法，我们同样可以利用动态规划算法来求解  $\vec{z}$ 。由于我们需要求解的是每一次迭代的最大值，而不是对每一个时间步求和，因此将  $\sum_{\vec{z}}$  替换成  $\max$  即可。我们称之为 Viterbi 算法，如算法1.3。得到每一个时间步的最大概率状态序列，一直迭代到  $T$  为止，我们就得到了最有可能输出给定观测序列的状态序列。

---

#### Algorithm 1.3 HMM 中的解码算法

---

1. Base Case:  $\alpha_i(0) = \max(A_{0i}), i = 1, 2, \dots, |S|$ ;
  2. Recursion:  $\alpha_j(t) = \max(\alpha_i(t-1)A_{ij}B_{jx_t}), \text{ for } i \text{ in } \{1, 2, \dots, |S|\}, j = 1, \dots, |S|, t = 1, \dots, T$
- 

### 1.5.4.3 学习问题

关于 HMM，在给定一系列观测序列的条件下，如何去估计最有可能产生这些观测序列的状态转移矩阵  $A$  和发射概率矩阵  $B$ ？比如说语音识别中，已知语音训练集，我们可以通过这些语音数据去训练 HMM 模型，估计出最适合的参数，从而可以利用这些参数去算未知语音的状态序列的最大似然。换大白话说一下这个例子：我们有一堆语音数据，这些数据的音频特征是已知的，这就是观测值，而对应的建模单元（音素什么的）是关于这些音频特征的函数，也就是状态值，我们希望通过这个训练数据来训练 HMM 模型，得到 HMM 的参数值，知道参数值之后，就可以用于识别其他音频了，这时候就变成了 HMM 的预测问题了。

所以如何通过训练集估计出模型的参数就变成了一个很重要的问题。本小节我们介绍一种 EM 算法的变体用于估计 HMM 的参数。算法1.4。从算法中，我们可以看到 M-Step 有一些约束条件，这些约束条件保证了 HMM 的两个矩阵  $A$  和  $B$  的有效性。正如前面我们求解 Markov Model 的参数一样，对于这个这个优化问题需构建一个拉格朗日算子。可是不论是 E-Step 还是 M-Step，都是对  $S^T$  内所有可能的状态序列  $\vec{z}$  进行枚举，E、M 两步对于所有可能的标签序列（状

态序列)  $\vec{z}$  都有  $|S|^T$  的复杂度。为了计算效率, 同样使用前面提到的前后向算法来计算 E、M 两步的结果。

---

**Algorithm 1.4** HMM 中的 EM 算法 1
 

---

Repeat until convergence { (E-Step) For every possible labeling  $\vec{z} \in S^T$ , set

$$Q(\vec{z}) := p(\vec{z}|\vec{x}; A, B)$$

(M-Step) Set

$$\begin{aligned} A, B &:= \arg \max_{A, B} \sum_{\vec{z}} Q(\vec{z}) \log \frac{P(\vec{z}, \vec{x}; A, B)}{Q(\vec{z})} \\ s.t. \sum_{j=1}^{|S|} A_{ij} &= 1, i = 1, 2, \dots, |S|; A_{ij} \geq 0, i, j = 1, 2, \dots, |S| \\ \sum_{k=1}^{|V|} B_{ik} &= 1, i = 1, 2, \dots, |S|; B_{ik} \geq 0, i = 1, \dots, |S|, k = 1, \dots, |V| \end{aligned}$$

}

---

首先让我们利用 Markov Assumption 重写估计 HMM 参数的目标函数, 目标函数推导第二行是因为对数函数中的分母与所求的值无关, 舍去了; 第三行利用了 Markov Assumption, 而最后一行使用了 **indicator functions**, 这么表示之后, 我们可以直接利用矩阵的索引来表示  $A$  和  $B$ 。

$$\begin{aligned} A, B &= \arg \max_{A, B} \sum_{\vec{z}} Q(\vec{z}) \log \frac{P(\vec{z}, \vec{x}; A, B)}{Q(\vec{z})} \\ &= \arg \max_{A, B} \sum_{\vec{z}} Q(\vec{z}) \log P(\vec{z}, \vec{x}; A, B) \\ &= \arg \max_{A, B} \sum_{\vec{z}} Q(\vec{z}) \log \left[ \prod_{t=1}^T p(\vec{z}_t | \vec{z}_{t-1}; A) p(\vec{x}_t | \vec{z}_t; B) \right] \\ &= \arg \max_{A, B} \sum_{\vec{z}} Q(\vec{z}) \log \left[ (p(\vec{z}_t | \vec{z}_{t-1}; A)) \left( \prod_{t=1}^T p(\vec{x}_t | \vec{z}_t; B) \right) \right] \\ &= \arg \max_{A, B} \sum_{\vec{z}} Q(\vec{z}) \sum_{t=1}^T \left[ \log A_{z_{t-1} z_t} + \log B_{z_t x_t} \right] \\ &= \arg \max_{A, B} \sum_{\vec{z}} Q(\vec{z}) \sum_{i=1}^{|S|} \sum_{j=1}^{|S|} \sum_{k=1}^{|V|} \sum_{t=1}^T \left[ 1\{z_{t-1} = s_i \wedge z_t = s_j\} \log A_{ij} + 1\{z_t = s_j \wedge x_t = v_k\} \log B_{jk} \right] \end{aligned} \tag{1.112}$$

与前面求解 Markov Model 的最大似然估计一样, 根据 HMM 的约束条件, 我们构造一个拉

格朗日函数，如公式1.113，将约束条件包含到拉格朗日函数之后，这样我们就不需要担心这些不等式约束了，因为求解出来的值都是非负数。

$$\begin{aligned} \mathcal{L}(A, B, \delta, \epsilon) = & \sum_{\vec{z}} Q(\vec{z}) \sum_{i=1}^{|S|} \sum_{j=1}^{|S|} \sum_{k=1}^{|V|} \sum_{t=1}^T \left[ 1\{z_{t-1} = s_i \wedge z_t = s_j\} \log A_{ij} + 1\{z_t = s_j \wedge x_t = v_k\} \log B_{jk} \right] \\ & + \sum_{j=1}^{|S|} |S| \epsilon_j (1 - \sum_{k=1}^{|V|} B_{jk}) + \sum_{i=1}^{|S|} |S| \delta_i (1 - \sum_{j=1}^{|S|} A_{ij}) \end{aligned} \quad (1.113)$$

公式1.113分别对 HMM 的参数矩阵中的各个元素求偏导数并置 0。

对  $A_{ij}$  有：

$$\frac{\partial \mathcal{L}(A, B, \delta, \epsilon)}{\partial A_{ij}} = \sum_{\vec{z}} Q(\vec{z}) \frac{1}{A_{ij}} \sum_{t=1}^T 1\{z_{t-1} = s_i \wedge z_t = s_j\} - \delta_i \quad (1.114)$$

置 0 解得：

$$A_{ij} = \frac{1}{\delta_i} \sum_{\vec{z}} Q(\vec{z}) \sum_{t=1}^T 1\{z_{t-1} = s_i \wedge z_t = s_j\} \quad (1.115)$$

对  $B_{jk}$  有：

$$\frac{\partial \mathcal{L}(A, B, \delta, \epsilon)}{\partial B_{jk}} = \sum_{\vec{z}} Q(\vec{z}) \frac{1}{B_{jk}} \sum_{t=1}^T 1\{z_t = s_j \wedge x_t = v_k\} - \epsilon_j \quad (1.116)$$

置 0 解得：

$$B_{jk} = \frac{1}{\epsilon_j} \sum_{\vec{z}} Q(\vec{z}) \sum_{t=1}^T 1\{z_t = s_j \wedge x_t = v_k\} \quad (1.117)$$

再对拉格朗日算子中的参数求偏导，并代入上述求得的  $A_{ij}$  和  $B_{jk}$ 。

对  $\delta_i$  有

$$\begin{aligned} \frac{\partial \mathcal{L}(A, B, \delta, \epsilon)}{\partial \delta_i} &= 1 - \sum_{j=1}^{|S|} A_{ij} \\ &= 1 - \sum_{j=1}^{|S|} \frac{1}{\delta_i} \sum_{\vec{z}} Q(\vec{z}) \sum_{t=1}^T 1\{z_{t-1} = s_i \wedge z_t = s_j\} \end{aligned} \quad (1.118)$$



置 0 解得：

$$\begin{aligned}\delta_i &= \sum_{j=1}^{|S|} \sum_{\vec{z}} Q(\vec{z}) \sum_{t=1}^T 1\{z_{t-1} = s_i \wedge z_t = s_j\} \\ &= \sum_{\vec{z}} Q(\vec{z}) \sum_{t=1}^T 1\{z_{t-1} = s_i\}\end{aligned}\quad (1.119)$$

对  $\epsilon_i$  有

$$\begin{aligned}\frac{\partial \mathcal{L}(A, B, \delta, \epsilon)}{\partial \epsilon_i} &= 1 - \sum_{k=1}^{|V|} B_{jk} \\ &= 1 - \sum_{k=1}^{|V|} \frac{1}{\epsilon_j} \sum_{\vec{z}} Q(\vec{z}) \sum_{t=1}^T 1\{z_t = s_j \wedge x_t = v_k\}\end{aligned}\quad (1.120)$$

置 0 解得：

$$\begin{aligned}\epsilon_i &= \sum_{k=1}^{|V|} \sum_{\vec{z}} Q(\vec{z}) \sum_{t=1}^T 1\{z_t = s_j \wedge x_t = v_k\} \\ &= \sum_{\vec{z}} Q(\vec{z}) \sum_{t=1}^T 1\{z_t = s_j\}\end{aligned}\quad (1.121)$$

再将  $\delta_i$  和  $\epsilon_j$  代回到公式1.115和公式1.117有：

$$\begin{aligned}A_{ij} &= \frac{\sum_{\vec{z}} Q(\vec{z}) \sum_{t=1}^T 1\{z_{t-1} = s_i \wedge z_t = s_j\}}{\sum_{\vec{z}} Q(\vec{z}) \sum_{t=1}^T 1\{z_{t-1} = s_i\}} \\ B_{jk} &= \frac{\sum_{\vec{z}} Q(\vec{z}) \sum_{t=1}^T 1\{z_t = s_j \wedge x_t = v_k\}}{\sum_{\vec{z}} Q(\vec{z}) \sum_{t=1}^T 1\{z_t = s_j\}}\end{aligned}\quad (1.122)$$

公式1.122中的  $A_{ij}$  和  $B_{jk}$  即为每一轮迭代中的 M-Step 所求，但是从公式上来看，都是对所有可能的标签序列  $\vec{z} \in S^T$  的枚举求和。前面提到的 E-step 中的  $Q_{\vec{z}}$  的定义是  $P(\vec{z}|\vec{x}; A, B)$ 。我们

考虑下如何用前后向概率  $\alpha_i(t)$  和  $\beta_j(t)$  来表征参数  $A_{ij}$  中的分子：

$$\begin{aligned}
 & \sum_{\vec{z}} Q(\vec{z}) \sum_{t=1}^T 1\{z_{t-1} = s_i \wedge z_t = s_j\} \\
 &= \sum_{t=1}^T \sum_{\vec{z}} 1\{z_{t-1} = s_i \wedge z_t = s_j\} Q(\vec{z}) \\
 &= \sum_{t=1}^T \sum_{\vec{z}} 1\{z_{t-1} = s_i \wedge z_t = s_j\} P(\vec{z}|\vec{x}; A, B) \\
 &= \frac{1}{P(\vec{x}; A, B)} \sum_{t=1}^T \sum_{\vec{z}} 1\{z_{t-1} = s_i \wedge z_t = s_j\} P(\vec{z}, \vec{x}; A, B) \\
 &= \frac{1}{P(\vec{x}; A, B)} \sum_{t=1}^T \alpha_i(t) A_{ij} B_{jx_t} \beta_j(t+1)
 \end{aligned} \tag{1.123}$$

解释下上面这个牛逼的公式。第二行和第三行是重新整理了下公式并且代入了  $Q(\vec{z})$  的定义。第四步是根据贝叶斯定理。第五步就是我个人认为非常吊的一步了。我们仔细看下  $1\{z_{t-1} = s_i \wedge z_t = s_j\}$ ，这个式子的意思是  $t$  时刻处于状态  $s_i$ ， $t+1$  时刻处于  $s_j$ ，好了，其他时刻，从 1 到  $t-1$ ，从  $t+2$  到  $T$ ，这些个时候处于什么状态，我们不在乎。是啥都可以。那么我们先来考虑下从 1 到  $t$  这个时间段，再配合上  $\sum \vec{z} z_{t-1} = s_i$ 、 $z_t = s_j$ ，一个从到  $s_i$ ，一个是从  $s_j$  出发，前面就是终止状态为  $s_i$  的前向概率，后面就是初始状态为  $s_j$  的后向概率。

所以这个转换太厉害了。内心鼓掌……

……

同样分母可以转换成：

$$\begin{aligned}
 & \sum_{\vec{z}} Q(\vec{z}) \sum_{t=1}^T 1\{z_{t-1} = s_i\} \\
 &= \sum_{t=1}^T \sum_{\vec{z}} 1\{z_{t-1} = s_i\} Q(\vec{z}) \\
 &= \sum_{t=1}^T \sum_{\vec{z}} 1\{z_{t-1} = s_i\} P(\vec{z}|\vec{x}; A, B) \\
 &= \frac{1}{P(\vec{x}; A, B)} \sum_{t=1}^T \sum_{\vec{z}} 1\{z_{t-1} = s_i\} P(\vec{z}, \vec{x}; A, B) \\
 &= \frac{1}{P(\vec{x}; A, B)} \sum_{t=1}^T \alpha_i(t) A_{ij}
 \end{aligned} \tag{1.124}$$

需要时刻谨记的是上述分子分母中的  $A_{ij}$  和  $B_{jx_t}$  都是上一轮迭代出来的参数，用于更新本

轮参数，为了避免产生混淆，我们将本轮更新后的参数分别记作  $\hat{A}_{ij}$  和  $\hat{B}_{jx_t}$ 。因此：

$$\hat{A}_{ij} = \frac{\sum_{t=1}^T \alpha_i(t) A_{ij} B_{jx_t} \beta_j(t+1)}{\sum_{t=1}^T \alpha_i(t) A_{ij}} \quad (1.125)$$

另外，此处说明一下，分母的计算与参考资料<sup>(15)</sup>中有些不同，其计算过程为：

$$\begin{aligned} & \sum_{\vec{z}} Q(\vec{z}) \sum_{t=1}^T 1\{z_{t-1} = s_i\} \\ &= \sum_{j=1}^{|S|} \sum_{t=1}^T \sum_{\vec{z}} 1\{z_{t-1} = s_i \wedge z_t = s_j\} Q(\vec{z}) \\ &= \frac{1}{P(\vec{x}; A, B)} \sum_{j=1}^{|S|} \sum_{t=1}^T \alpha_i(t) A_{ij} B_{jx_t} \beta_j(t+1) \end{aligned} \quad (1.126)$$

最终求得的：

$$\hat{A}_{ij} = \frac{\sum_{t=1}^T \alpha_i(t) A_{ij} B_{jx_t} \beta_j(t+1)}{\sum_{j=1}^{|S|} \sum_{t=1}^T \alpha_i(t) A_{ij} B_{jx_t} \beta_j(t+1)} \quad (1.127)$$

同样地， $\hat{B}_{jk}$  的分子可以转换如下：

$$\begin{aligned} & \sum_{\vec{z}} Q(\vec{z}) \sum_{t=1}^T 1\{z_t = s_j \wedge x_t = v_k\} \\ &= \sum_{t=1}^T \sum_{\vec{z}} 1\{z_t = s_j \wedge x_t = v_k\} Q(\vec{z}) \\ &= \sum_{t=1}^T \sum_{\vec{z}} 1\{z_t = s_j \wedge x_t = v_k\} P(\vec{z} | \vec{x}; A, B) \\ &= \frac{1}{P(\vec{x}; A, B)} \sum_{t=1}^T \sum_{\vec{z}} 1\{z_t = s_j \wedge x_t = v_k\} P(\vec{z}; \vec{x}; A, B) \\ &= \frac{1}{P(\vec{x}; A, B)} \sum_{i=1}^{|S|} \sum_{t=1}^T \sum_{\vec{z}} 1\{z_{t-1} = s_i \wedge z_t = s_j \wedge x_t = v_k\} P(\vec{z}; \vec{x}; A, B) \\ &= \frac{1}{P(\vec{x}; A, B)} \sum_{i=1}^{|S|} \sum_{t=1}^T 1\{x_t = v_k\} \alpha_i(t) A_{ij} B_{jx_t} \beta_j(t+1) \end{aligned} \quad (1.128)$$

其分母可以转换为：

$$\begin{aligned}
 & \sum_{\vec{z}} Q(\vec{z}) \sum_{t=1}^T 1\{z_t = s_j\} \\
 &= \sum_{t=1}^T \sum_{\vec{z}} 1\{z_t = s_j\} Q(\vec{z}) \\
 &= \sum_{t=1}^T \sum_{\vec{z}} 1\{z_t = s_j\} P(\vec{z}|\vec{x}; A, B) \\
 &= \frac{1}{P(\vec{x}; A, B)} \sum_{i=1}^{|S|} \sum_{t=1}^T \sum_{\vec{z}} 1\{z_{t-1} = s_i \wedge z_t = s_j\} P(\vec{z}; \vec{x}; A, B) \\
 &= \frac{1}{P(\vec{x}; A, B)} \sum_{i=1}^{|S|} \sum_{t=1}^T \alpha_i(t) A_{ij} B_{jx_t} \beta_j(t+1)
 \end{aligned} \tag{1.129}$$

故而：

$$\hat{B}_{jk} = \frac{\sum_{i=1}^{|S|} \sum_{t=1}^T 1\{x_t = v_k\} \alpha_i(t) A_{ij} B_{jx_t} \beta_j(t+1)}{\sum_{i=1}^{|S|} \sum_{t=1}^T \alpha_i(t) A_{ij} B_{jx_t} \beta_j(t+1)} \tag{1.130}$$

算法1.5描述了利用前后向算法来进行参数学习的 Baum-Welch 算法：

---

**Algorithm 1.5** HMM 中的 Baum-Welch 算法

---

Initialization: Set  $A$  and  $B$  as random valid probability matrices where  $A_{i0} = 0$  and  $B_{0k} = 0$ , for  $i = 1, \dots, |S|$  and  $k = 1, \dots, |V|$ . repeat until convergence { (E-Step) Run the forward-backward algorithms to compute  $\alpha_i$  and  $\beta_j$  for  $i, j = 1, \dots, |S|$

$$\gamma_t(i, j) := \alpha_i(t) A_{ij} B_{jx_t} \beta_j(t+1) \tag{1.131}$$

(M-Step) Re-estimate the maximum likelihood parameters as:

$$\begin{aligned}
 \hat{A}_{ij} &= \frac{\sum_{t=1}^T \gamma_t(i, j)}{\sum_{j=1}^{|S|} \sum_{t=1}^T \gamma_t(i, j)} \\
 \hat{B}_{jk} &= \frac{\sum_{i=1}^{|S|} \sum_{t=1}^T 1\{x_t = v_k\} \gamma_t(i, j)}{\sum_{i=1}^{|S|} \sum_{t=1}^T \gamma_t(i, j)}
 \end{aligned} \tag{1.132}$$

}

---

仍然需要对算法1.5进行一些补充说明。与其计算所有可能的  $\vec{z}$  对应的  $Q(\vec{z})$ ，我们计算更有效率的  $\gamma_t(i, j) = \alpha_i(t) A_{ij} B_{jx_t} \beta_j(t+1)$ 。

$\gamma_t(i, j)$  表示的是所有的观测序列  $\vec{x}$ ，这些  $\vec{x}$  对应的  $t$  和  $t+1$  时刻状态分别为  $s_i$  和  $s_j$ 。再仔



细观察下 M-Step 的公式，我们可以更直观的感受  $A_{ij}$  等于从状态  $s_i$  跳转到  $s_j$  的次数除以出现状态  $s_i$  的次数， $B_{jk}$  等于状态  $s_j$  生成观测  $v_k$  除以状态  $s_j$  的次数。

正如其他很多 EM 算法的应用场景一样，HMM 的参数学习是一个非凸优化问题，其包含很多个局部最优点。不同的初始化参数会收敛到不同的最大似然，所以最好多跑几次。另外对参数矩阵的概率平滑也是很有必要的，这样可以防止参数矩阵中出现转移概率或者发射概率为 0 的情况。

## 1.6 线性判别分析

线性判别分析（Linear Discriminative Analysis, LDA）是一种有监督的降维学习方法，其不仅仅可以达到降维的目的，还可以对原始数据进行聚类，使得类间距变大，类内距变小。有监督意味着 LDA 中所有的数据都是有标签的，这也是和 PCA 的一个重要区别，PCA 无需样本类别，是一种无监督的降维方法。

LDA 概况起来就是“投影后类内方差最小，类间方差最大。”LDA 是对数据进行投影，将其投影到低维空间，投影后相同类别的样本距离更近，不同类别的类别中心更远。本节首先对二类 LDA 进行分析，再推广至多类 LDA。

二类 LDA 的形象表述如图 1.2 右。左边的图不同类之间有交叉，决策边界有重合，而右图既使得相同类更集中，也使得不同类的分类边界更清晰，这就是 LDA 达到的效果。

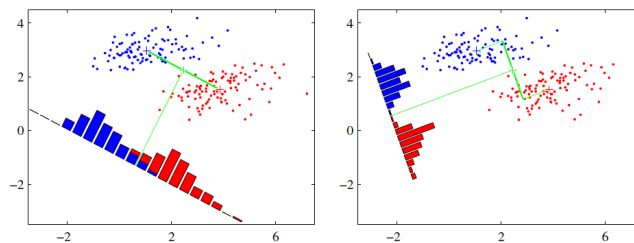


图 1.2: 二类 LDA 转换效果图

假定数据集  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$ ，其中  $x_i \in R^n$ ,  $y_i \in \{0, 1\}$ 。我们定义  $N_j (j = 0, 1)$  为第  $j$  类样本的个数， $X_j (j = 0, 1)$  为第  $j$  类样本的合集， $\mu_j (j = 0, 1)$  为第  $j$  类样本的均值向量， $\Sigma_j (j = 0, 1)$  为第  $j$  类样本缺少分母部分的协方差矩阵。那么  $\mu_j$  和  $\Sigma_j$  的表达式分别如公式 1.133 和公式 1.134 所示。

$$\mu_j = \frac{1}{N_j} \sum_{x \in j} x \quad (1.133)$$

$$\Sigma_j = \sum_{x \in j} (x - \mu_j)(x - \mu_j)^T \quad (1.134)$$

由于只有两类数据，所以只需要将这些数据投影到一条直线上就可以，假设投影向量为  $w$ ，

则对任意一个样本，其在直线上的投影为  $w^T x$ ，类别中心的投影分别为  $w^T \mu_0$  和  $w^T \mu_1$ ，LDA 要求不同类别之间的类别中心尽可能的远，所以需要最大化  $\|w^T \mu_0 - w^T \mu_1\|_2^2$ ，同时我们还希望同一类别尽可能接近，也就是样本投影之后的协方差尽可能的小，投影后的协方差如公式 1.135。

$$\begin{aligned}\Sigma'_j &= \sum_{x \in j} (w^T x - w^T \mu_j)(w^T x - w^T \mu_j)^T \\ &= \sum_{x \in j} w^T (x - \mu_j)(x - \mu_j)^T w \\ &= w^T \Sigma_j w\end{aligned}\quad (1.135)$$

所以我们希望最小化  $w^T \Sigma_0 w + w^T \Sigma_1 w$ ，由此我们可以得到需要优化的目标函数，如公式 1.136。

$$\begin{aligned}\arg \max_w J(w) &= \arg \max_w \frac{\|w^T \mu_0 - w^T \mu_1\|_2^2}{w^T \Sigma_0 w + w^T \Sigma_1 w} \\ &= \arg \max_w \frac{w^T (\mu_0 - \mu_1)(\mu_0 - \mu_1)^T w}{w^T (\Sigma_0 + \Sigma_1) w}\end{aligned}\quad (1.136)$$

类内散度矩阵  $S_w$  和类间散度矩阵  $S_b$  分别定义为公式 1.137 和公式 1.138。

$$S_w = \Sigma_0 + \Sigma_1 \quad (1.137)$$

$$S_b = (\mu_0 - \mu_1)(\mu_0 - \mu_1)^T \quad (1.138)$$

所以目标函数就变成了：

$$\arg \max_w J(w) = \arg \max_w \frac{w^T S_b w}{w^T S_w w} \quad (1.139)$$

也就是求解出  $w$  使得  $J(w)$  最大。根据 1.2 中的介绍，我们可以通过计算矩阵  $S_w^{-1} S_b$  的特征值和特征向量得到对应的  $w$ ，即求解公式 1.140。 $J(w)$  的最大值为  $S_w^{-1} S_b$  的最大特征值，最小值为  $S_w^{-1} S_b$  的最小特征值，而  $S_w$  和  $S_b$  均可由原始数据求解得出，因此很容易就可以求解出  $J(w)$  的最大值。

$$S_w^{-1} S_b w = \lambda w \quad (1.140)$$

接下来我们分析下多类 LDA 的原理。

假定数据集  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$ ，其中  $x_i \in R^n$ ， $y_i \in \{C_1, C_2, \dots, C_k\}$ 。我们定义  $N_j (j = 0, 1, \dots, k)$  为第  $j$  类样本的个数， $X_j (j = 0, 1, \dots, k)$  为第  $j$  类样本的合集， $\mu_j (j = 0, 1, \dots, k)$  为第  $j$  类样本的均值向量， $\Sigma_j (j = 0, 1, \dots, k)$  为第  $j$  类样本缺少分母部分的协方差矩阵。此时是多类分类，因此投影后的空间不再是一条直线，而是一个超平面。假设投影后的低维空间维度

为  $d$ ，对应的基向量为  $(w_1, w_2, \dots, w_d)$ ，基向量组成的矩阵为  $W \in R^{n \times d}$ 。

此时类内的散度矩阵  $S_W$  仍旧存在，如公式1.141。

$$S_W = \sum_{j=1}^k \Sigma_j \quad (1.141)$$

但是类间的散度矩阵就有所不同了。此时用每个类别的均值到全局均值的距离来衡量类间距如图1.3。

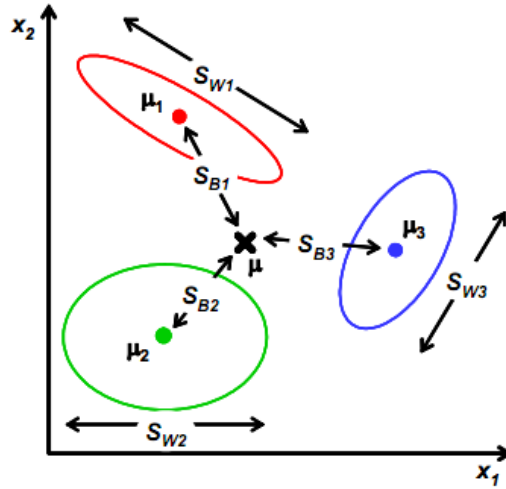


图 1.3: 多类 LDA 的类间散度矩阵示意图

其定义为公式1.142。

$$S_B = \sum_{j=1}^k N_j (\mu_j - \mu)(\mu_j - \mu)^T \quad (1.142)$$

其中：

$$\mu = \frac{1}{m} \sum_{i=1}^m x_i \quad (1.143)$$

$$\mu_j = \frac{1}{N_j} \sum_{x_j \in C_j} x_j \quad (1.144)$$

同样此时的优化目标为公式1.145。

$$\arg \max_W J(W) = \arg \max_W \frac{W^T S_B W}{W^T S_W W} \quad (1.145)$$

此时目标函数求解转换成了公式1.146:

$$S_W^{-1} S_B W = \lambda W \quad (1.146)$$

以上，可以总结出多类 LDA 的求解步骤：

1. 计算每个类别的均值向量和方差，以及全局均值向量；
2. 根据均值向量和方差，计算  $S_W$  和  $S_B$ ；
3. 对  $S_W^{-1} S_B W = \lambda W$  进行求解，求出  $S_W^{-1} S_B$  的特征值和特征向量；
4. 对特征向量进行排序，设定低维空间的维度  $d$ ，选取前  $d$  个特征值和特征向量，特征向量组合成投影矩阵  $W$ ；
5. 通过投影矩阵计算出投影后的输入数据  $x'_i = W^T x_i$ ；
6. 得到输出的新数据集： $\{(x'_1, y_1), (x'_2, y_2), \dots, (x'_m, y_m)\}$ 。

## 1.7 最大似然线性变换

最大似然线性变换（Maximum Likelihood Linear Transform）

在 HMM 系统中，协方差矩阵的选择可以是对角阵，分块对角阵或者全矩阵。相对于对角阵来说，全矩阵的优势在于对特征向量元素之间关系的建模，劣势在于参数量巨大。

## 1.8 Beta 分布

## 1.9 MLE 和 MAP

MLE vs MAP

## 第 2 章 微软 Edx 语音识别笔记

本章笔记主要是对微软 Edx 的课程 **Speech Recognition System** 的记录，首版主要是翻译，再加上自己翻阅其他资料综合起来的一些思考和总结。代码见 **Speech-Recognition**

### 2.1 Background and Fundamentals

#### 2.1.1 Phonetics

Phonetics（语音学）是 Linguistics（语言学）的一个分支，其研究的是人类语音发出的声音（sound）。语音学围绕着声音的产生（通过人类的发音器官）、声音的声学特性和感知。语音学有三个基本的分支，这三个分支都与 ASR 有关系。

1. Articulatory Phonetics（发音语音学）：通过发音器官、不同说话人而产生的声音；
2. Acoustic Phonetics（声学语音学）：声音从说话人到听者的传输；
3. Auditory Phonetics（听觉语音学）：听者对于声音的接收和感知。

声音的最小单元我们成为 **Phoneme**，即音素。序列中的词（Words）是由一个或多个音素组成的。一个音素的声学实现称为 **Phone**。图2.1展示了美式英语的音素和一般实现办法。

Phonemes	Word Examples	Description
ih	fill, hit, lid	front close unrounded (lax)
ae	at, carry, gas	front open unrounded (tense)
aa	father, ah, car	back open unrounded
ah	cut, bud, up	open-mid back unrounded
ao	dog, lawn, caught	open-mid back round
ay	tie, ice, bite	diphthong with quality: aa + ih
ax	ago, comply	central close mid (schwa)
ey	ate, day, tape	front close-mid unrounded (tense)
eh	pet, berry, ten	front open-mid unrounded
er	turn, fur, meter	central open-mid unrounded rhotic
ow	go, own, tone	back close-mid rounded
aw	fault, have, our	diphthong with quality: aa + uh
oy	toy, coin, oil	diphthong with quality: ao + ih
uh	book, pull, good	back close-mid unrounded (lax)
uw	tool, crew, moose	back close round
b	big, able, tab	voiced bilabial plosive
p	put, open, tap	voiceless bilabial plosive
d	dig, idea, wad	voiced alveolar plosive
t	talk, sat	voiceless alveolar plosive &
t	meter	alveolar flap
g	gut, angle, tag	voiced velar plosive
k	cut, ken, take	voiceless velar plosive
f	fork, after, if	voiceless labiodental fricative
v	vat, over, have	voiced labiodental fricative
s	sit, cast, toss	voiceless alveolar fricative
z	zap, lazy, haze	voiced alveolar fricative
th	thin, nothing, truth	voiceless dental fricative
dh	Then, father, scythe	voiced dental fricative
sh	she, cushion, wash	voiceless postalveolar fricative
zh	genre, azure	voiced postalveolar fricative
l	lid	alveolar lateral approximant
l	elbow, sail	velar lateral approximant
r	red, part, far	retroflex approximant
y	yacht, yard	palatal sonorant glide
w	with, away	labiodental sonorant glide
hh	help, ahead, hotel	voiceless glottal fricative
m	mat, amid, aim	bilabial nasal
n	no, end, pan	alveolar nasal
ng	sing, anger	velar nasal
ch	chin, archer, march	voiceless alveolar affricate: t + sh
jh	joy, agile, edge	voiced alveolar affricate: d + zh

图 2.1: 美式英语的音素和一般实现办法

一般我们将音素分为两类：元音（Vowel）和辅音（Consonants）。

1. Vowels: 元音有两个特点，一是元音都是发声的声音（voiced sound），这意味着从声带（vocal

chords) 到口腔 (mouth cavity) 的气流是由声带的某种基频的震动 (或者音高) 产生的。二是舌头在生产过程中不会以任何方式形成气流收缩。每个元音发声的时候舌头、嘴唇和下巴的造型都不一样。这些不同的方式形成了不同的共振态, 我们称之为共振峰。这些共振峰的共振态频率形成了不同的元音。

2. Consonants: 辅音是通过在口腔中或者空气中很明显的气流收缩形成的。有些辅音和元音一样是发声的, 有些是不发声的。不发声的音素不会激活声带, 因此也不存在基频或者音高。一些辅音音调成对出现, 只有在有声或无声的情况下才有所不同, 但在其他方面是相同的。比如说 /b/ 和 /p/ 这两个音素的发音方式是相同的, 因为你的嘴唇、下巴还有舌头的姿势是一样的。但是 /b/ 是发声的, /p/ 是不发声的。

音素的另外一个重要特性是 **其根据不同的上下文音素发音是会改变的**。我们称之为 Phonetic Context。之所以会这样, 是因为协同发音 (coarticulation)。这些声音连续起来发音会改变其原有的特征。由协同发音产生的音素我们称为音素变体 (allophones)。

所有当前的这些语音识别系统都使用了音素的语境相关的特性来建立处于不同 phonetic context 的音素模型。

### 2.1.2 Words and Syntax

Syllable 是一串声音, 是个序列, 由一个核心的音素, 可能有初始音素和终止音素, 这个核心音素一般是个元音或者一个音节辅音 (syllabic consonant), 是能够唱出来或者吼出来的声音。

举个例子, 英文单词 "bottle" 包含两个 syllable。第一个 syllable 有三个 phone, 在 Arpabet 音素描述代码里, 是 "b aa t"。这个 "aa" 就是核心音素, "b" 是发声的初始音素, "t" 是不发音的终止音素; 第二个 syllable 是只包含一个 syllabic consonant "l"。

一个 syllable 也可以组成一个词, 其本身就是一个单独的音素, 比如说, "Eye", "uh", 或者 "eau" (医: 水)。

语音识别里面, syllable 很少会考虑作为声学模型的建模单元, 而词一般是变成音素来建模。

Syntax (句法规则) 描述了给定词和定义了语法的规则下, 句子的形成。而 Semantics (语义学) 一般指代的是句子中的词或者短语是如何形成句意的。Syntax 和 Semantics 是 NLP 的重要组成部分, 但是在语音识别里面, 不起主要作用。

### 2.1.3 Measuring Performance

在语音识别系统的搭建和实验中, 如何来衡量一个系统的好坏呢? 由于语音识别是一个序列任务, 跟图像当中的分类不一样, 因此我们在衡量系统的性能时需要考虑到整个序列。

语音识别准确率衡量最常用的一个指标是词错误率 (word error rate, WER)。一般识别出来的结果可能会产生三种错误: 替换 (substitution)、删除 (delete) 和插入 (insert)。替换指的是一个词被识别成了另外一个词; 删除指的是原本有词, 但是没有识别出来; 插入指的是原本没

有词，多识别出来了词。WER 的计算方式如公式2.1。

$$WER = \frac{N_{sub} + N_{ins} + N_{del}}{N_{ref}}$$

(2.1)

其中  $N_{sub}$ 、 $N_{ins}$  和  $N_{del}$  分别是替换、插入和删除的数量，而  $N_{ref}$  是参考文本描述中词的个数。

WER 的计算用的是通过计算实际输出描述和参考文本描述之间的字符串编辑距离得到的。编辑距离的实现通过动态规划算法。因为长文本的编辑距离可能不可靠，所以我们通过逐句的计算累积的错误，这些错误最终整合到一起来计算测试集的 WER。编辑距离的原理见补充知识点中的2.7.3节。

表2.1呈现了实际输出和参考文献之间的不同，以及对应的三种错误。

1

Ref: however a little later we had a comfortable chat

2

Hyp: how never a little later he had comfortable chat

表 2.1: WER 计算公式中的三种错误实例演示

Reference	Hypothesis	Error
however	how	Substitution
	never	Insertion
a	a	
little	little	
later	later	
we	he	Substitution
had	had	
a		Deletion
comfortable	comfortable	
chat	chat	

在某些情况中，这三种错误的成本不对等，那么计算编辑距离的时候可以作相应的调整。

句错误率（Sentence Error Rate, SER）是另外一种衡量系统的标准，其计算方式是整句没有出现任何错误。SER 仅仅作为一个指标，来看下错误的句子占全部句子的比例。

2.1.4 Significance Testing

统计显著性检验（statistical significance testing）涉及测量两个实验（或算法）之间的差异在多大程度上归因于两个算法中的实际差异，或者仅仅是数据，实验设置或其他因素中的结果固有变异性。统计显著性是所有分类任务的基石，只是统计显著性检验的方法取决于任务的特性。大多数方法的核心是假设检验的概念中存在一个无效假设。问题在于你有多大的 confidence 能够说无效假设会被拒绝。





对于语音识别来说，比较两个实验或者算法最常用的方法是 Matched Pairs Sentence-Segment Word Error(MAPSSWE) 检验，简称为 Matched Pairs Test<sup>(5)</sup>。

在这个方法中，测试集被分为几份，假设这些子测试集中任意一个的错误都与其他子测试集统计独立。这个假设和语音识别的实验很贴合，因为测试数据都是一句一句的经由识别器输出结果。给定了每一个句子的 WER，就很容易构建一个 matched pairs<sup>(13)</sup>。

### 2.1.5 Other Consideration

除去准确率，对识别系统性能的影响还可能包括计算需求，处理速度或者延迟什么的。解码速度一般用实时因子（real-time factor, RTF）来衡量。RTF 为 1.0 指的是系统处理 10s 的数据需要花 10s 的时间。

RTF 高于 1.0 意味着系统要花更多的时间来解码，对于某些应用，也许是可以接受的，比如希望获得一个会议或者讲座的转写，相对于快速的得到转写结果，准确率可能更重要一些，因此多花一些时间也是可以接受的。

当 RTF 低于 1.0，系统会在当前数据达到前就处理好了之前的数据。当不止一个系统在同一个机器上运行的时候，这个就比较有用了。在这种情况下，我们可以用多线程来并行处理多路音频流。此外 RTF 低于 1.0 意味着系统能够满足在线实时解码的音频流应用。比如说，当我们在处理一个手机上远程音频需求的时候，网络阻塞可能会使得服务器接收音频产生间隙和延迟。如果语音识别器能够以比实时更快的速度来处理，那么它就可以在数据达到之后迅速跟进，追上最新的音频进度，以速度来掩盖网络的延迟。

一般来说，语音识别系统能够在准确率和速度之间调整，但是这种调整也是有限的，不可能无限好或者无限快。对于一个给定的模型和测试集，speed-accuracy 图有一条不可被突破的渐近线（asymptote），即便给予无限的算力。所以准确率是有个极限的，这个时候的错误率可以认为就是模型带来的错误。一旦根据模型搜索找到了最好的结果，进一步的处理也不会带来准确率上的提升。

### 2.1.6 The Fundamental Equation

语音识别可以看作一个优化任务。特别地，给定一个观测序列  $O = \{O_1, \dots, O_N\}$ ，我们找寻的是最有可能的词序列  $W = \{W_1, \dots, W_M\}$ ，也就是说我们要找到最大化后验概率  $P(W|O)$  的词序列，如公式2.2。

$$\hat{W} = \arg \max_W P(W|O) \quad (2.2)$$



利用贝叶斯规则，我们得到公式2.3。

$$P(W|O) = \frac{P(W)P(O|W)}{P(O)} \quad (2.3)$$

因为词序列并不依赖于观测序列的边缘概率分布  $P(O)$ ，我们可以忽略这个部分，综合上述两个公式，我们得到公式2.4。

$$\hat{W} = \arg \max_W P(W)P(O|W) \quad (2.4)$$

这就是语音识别的基本公式。语音识别问题就可以看作是在这个联合模型上的搜索。

公式中的  $P(O|W)$  叫做**声学模型 (acoustic model)**。这个模型描述了在给定词序列  $W$  的条件下，声学观测  $O$  的分布。声学模型表征的是词序列是如何转换成声学实现的，进而转换成 ASR 系统的声学观测的。

公式中的  $P(W)$  叫做**语言模型 (language model)**，其只取决于词序列  $W$ 。语言模型给每一个可能的词序列一个概率值。它是由日常使用的一些词序列训练成的。一个训练好的英语语言模型会给 "I like turtles" 高的概率值，给 "Turtles sing table" 低的概率值。语言模型促使着词序列的搜索沿着训练数据中的模式开展。语言模型也可以在一些纯文本的应用中见到，比如浏览器的自动补全等。

由于诸多原因，构建一个语音识别系统比这个简单地公式所能呈现的要复杂的多得多。

### 2.1.7 Lab 1: Create a speech recognition scoring program

本模块有一个实验作业，标题为 "Create a speech recognition scoring program"。

#### Required files:

- [wer.py](#)
- [M1\\_score.py](#)

#### Instructions:

In this lab, you will write a program in Python to compute the word error rate (WER) and sentence error rate (SER) for a test corpus. A set of hypothesized transcriptions from a speech recognition system and a set of reference transcriptions with the correct word sequences will be provided for you.

This lab assumes the transcriptions are in a format called the "trn" format, created by NIST. The format is as follows. The transcription is output on a single line followed by a single space and then the root name of the file, without any extension, in parentheses. For example, the audio file "tongue\_twister.wav" would have a transcription

sally sells seashells by the seashore (tongue\_twister)



Notice that the transcription does not have any punctuation or capitalization, nor any other formatting (e.g. converting "doctor" to "dr.", or "eight" to "8"). This formatting is called "Inverse Text Normalization" and is not part of this course.

The python code [M1\\_Score.py](#) and [wer.py](#) contain the scaffolding for the first lab. A main function parses the command line arguments and `string_edit_distance()` computes the string edit distance between two strings.

Add code to read the trn files for the hypothesis and reference transcriptions, to compute the edit distance on each, and to aggregate the error counts. Your code should report:

- Total number of reference sentences in the test set
- Number of sentences with an error
- Sentence error rate as a percentage
- Total number of reference words
- Total number of word errors
- Total number of word substitutions, insertions, and deletions
- The percentage of total errors (WER) and percentage of substitutions, insertions, and deletions

The specific format for outputting this information is up to you. Note that you should not assume that the order of sentences in the reference and hypothesis trn files is consistent. You should use the utterance name as the key between the two transcriptions.

When you believe your code is working, use it to process `hyp.trn` and `ref.trn` in the `misc` directory, and compare your answers to the solution.

### Solutions:

我们将课程提供的代码补充完整，加上对应的代码注释和结果展示。首先说下测试文本的格式。测试文本有两个

- **hyp.trn**: 模型输出结果;
- **ref.trn**: 参考文本。

hyp.trn 的内容如下:

```
1 apple coconut date eggplant elephant fig (0000-000000-0000)
2 one tiger three flamingo five six (0000-000000-0001)
3 delaware cat georgia dog mouse massachusetts (0000-000000-0002)
```

ref.trn 的内容如下:

```
1 apple banana coconut date eggplant fig (0000-000000-0000)
2 one two three four five six (0000-000000-0001)
3 delaware pennsylvania new_jersey georgia connecticut massachusetts
```



(0000-00000-0002)

可以很明显的看到，hyp 和 ref 是一一对应的，格式相同，都是 trn 文件，其内部都是"text (token)" 的形式。

本节的代码逻辑是构建一个函数 **string\_edit\_distance**，其接受 ref 和 hyp 两个变量，返回的是一个 tuple，依次为 token、edits、deletions、insertions 和 substitutions。该函数存储于 **wer.py** 中，以供代码 **M1\_score.py** 调用。

我们先来看 **wer.py**，理解了函数 **string\_edit\_distance** 自然就会调用了。代码和注释如下：

```

1 import numpy as np
2
3
4 def string_edit_distance(ref=None, hyp=None):
5
6     if ref is None or hyp is None:
7         RuntimeError("ref and hyp are required, cannot be None")
8
9     x = ref
10    y = hyp
11    tokens = len(x)
12    if (len(hyp)==0):
13        return (tokens, tokens, tokens, 0, 0)
14
15    # p[ix,iy] consumed ix tokens from x, iy tokens from y
16    p = np.PINF * np.ones((len(x) + 1, len(y) + 1)) # track total errors
17    e = np.zeros((len(x)+1, len(y) + 1, 3), dtype=np.int) # track deletions
18    , insertions, substitutions
19    p[0] = 0
20    for ix in range(len(x) + 1):
21        for iy in range(len(y) + 1):
22            cst = np.PINF*np.ones([3])
23            s = 0
24            if ix > 0:
25                cst[0] = p[ix - 1, iy] + 1 # deletion cost
26            if iy > 0:
27                cst[1] = p[ix, iy - 1] + 1 # insertion cost

```

```

27         if ix > 0 and iy > 0:
28             s = (1 if x[ix - 1] != y[iy - 1] else 0)
29             cst[2] = p[ix - 1, iy - 1] + s # substitution cost
30         if ix > 0 or iy > 0:
31             idx = np.argmin(cst) # if tied, one that occurs first wins
32             p[ix, iy] = cst[idx]
33
34         if (idx==0): # deletion
35             e[ix, iy, :] = e[ix - 1, iy, :]
36             e[ix, iy, 0] += 1
37         elif (idx==1): # insertion
38             e[ix, iy, :] = e[ix, iy - 1, :]
39             e[ix, iy, 1] += 1
40         elif (idx==2): # substitution
41             e[ix, iy, :] = e[ix - 1, iy - 1, :]
42             e[ix, iy, 2] += s
43
44     edits = int(p[-1,-1])
45     deletions, insertions, substitutions = e[-1, -1, :]
46     return (tokens, edits, deletions, insertions, substitutions)

```

## 2.2 Speech Signal Processing

### 2.2.1 Introduction

通过空气传播的音频波形通过麦克风的捕捉，将这些压力波转换成可捕捉的电信号活动。对这些电信号活动进行采样，这样就得到了一系列采样波形，我们用这些波形来描述信号。音乐信号一般采样率为 44100 Hz，即一秒钟会有 44100 个采样点。根据奈奎斯特定理（Nyquist theorem），只有频率低于 22050 Hz 的音频才可以被捕捉到。如果信号的高频部分比较少，最高频率为 8000Hz 的话，采样率一般就是 16000Hz。传统的电话和大部分的手机带限是 3400 Hz，所以 8000 Hz 的采样率就足够了。所以电话语音的采样率一般就是 8000Hz。

一个典型的音频波形图如2.2（左），其说的句子是"speech recognition is cool stuff"。

回忆上一节中讨论的发声音素和不发声音素，我们来瞅一下最后一个单词"stuff"的波形图，如图2.2（右）。从图上我们看到，这个单词的发音有三个不同的部分，初始不发声语音"st"，中

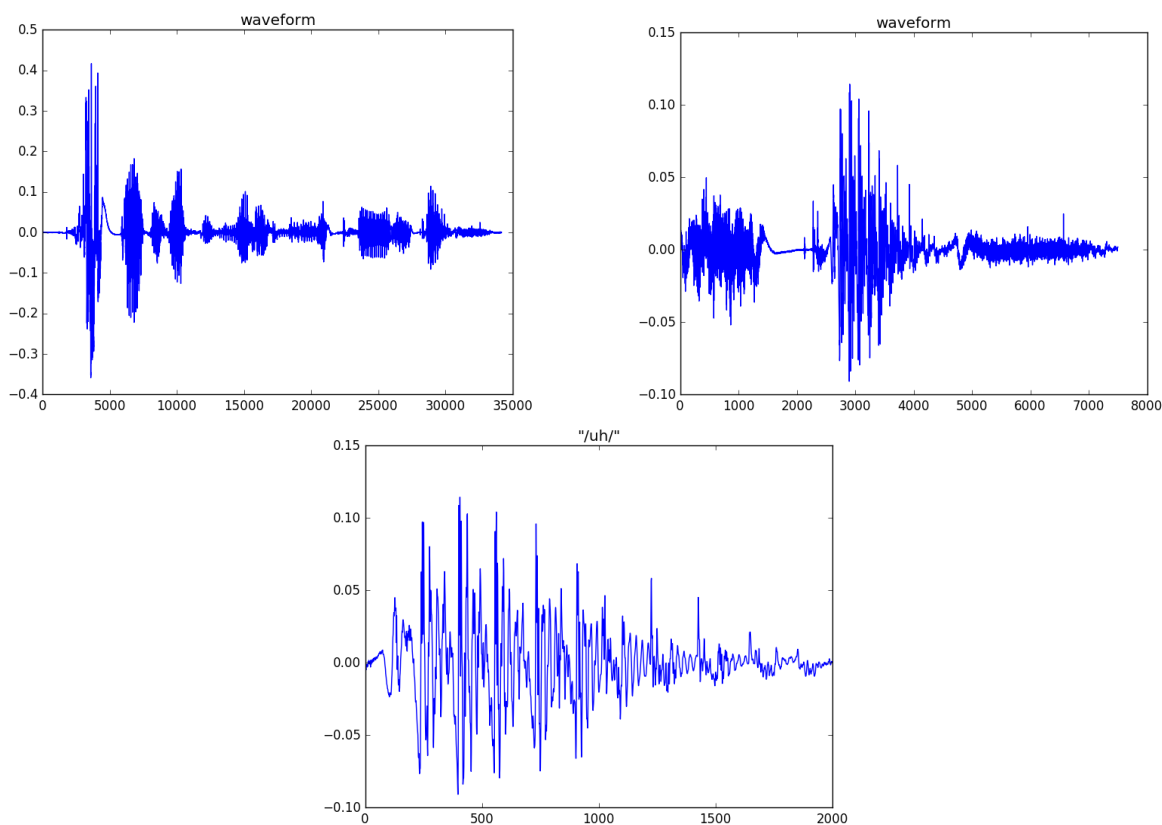


图 2.2: 完整音频波形图与部分波形图

间发声语音"uh", 终止不发声语音"f"。不发声的语音部分看上去像噪音, 具有随机性。而发声部分的语音由于声带的振动而具有周期性。

在拉近一点, 我们看看发声的这个元音"/uh/", 如图2.2 (下) 可以更明显的看到这个元音的周期性。

从这些波形图中, 我们可以看到波形的特征形成有两方面的因素: (1) 声带的应激反应 (excitation) 趋势着空气从声道和嘴巴流出; (2) 发出某种特定声音时, 声道本身的形状。

举例来说, 图2.2 (右) 中"st" 和"f" 看上去都像噪声, 但是他们的形状却不相同, 就是因为他们是不一样的声音。而"uh" 的声音更具周期性, 因为发声的应激反应, 其由于发声的时候声道的作用有独特的形状。所以对于同一个说话人来说, 不同的元音可能会有着相近的周期, 但是波形的整体形状是不一样的, 就因为这些波形都是由相同的声带产生的, 但是发不同的声音声道是不一样的。

在信号处理中, 一般用源滤波器模型 (source-filter model) 来对这个语音产生的过程进行建模。声源是由通过声道的声带产生的激励信号, 我们将其建模成为时变线性滤波器。源滤波器在语音识别中有很多应用, 比如语义分析和编码。而且有很多种办法来评估源信号和滤波器的参数, 比如很有名的线性预测编码 (Linear Predictive Coding, LPC)。

对于语音识别来说，音素分类在很大程度上取决于声道的形状，也就是说取决于源滤波器模型的滤波器部分。激励信号或者原信号大都被忽略或者舍弃了。所以语音识别的特征提取过程一般设计成捕捉话语过程的时变滤波器形状。

## 2.2.2 Feature Extraction

从波形图中，很明显语音是非平稳信号（non-stationary signal），这就意味着语音信号的统计特性会随着时间变化而变化。所以为了分析语音信号，我们需要将信号分成一个又一个 chunk（也成为窗或者帧），这些 chunk 短到可以认为它们是平稳信号。这样我们就可以去分析一系列短时的有重叠的语音帧。在语音识别中，我们一般选窗长为 25ms，窗移位 10ms，也就是说一秒钟会被分成 100 帧。

因为我们是从一个长音频提取的 chunk，所以对于每一个 chunk 的边缘，我们要进行一些处理。一般对每一帧数据加个窗函数，常用的是汉明窗（Hamming Windows），当然也有用其他窗函数的。定义  $m$  为某一帧的索引， $n$  为采样点的索引， $L$  是这一帧采样点的个数， $N$  是采样中的偏移量。那么从原始信号中提取出来的每一帧计算公式见 2.5。

$$x_m[n] = w[n]x[mN + n], n = 0, 1, \dots, L - 1 \quad (2.5)$$

其中  $w[n]$  是窗函数。

然后我们利用离散傅里叶变换将每一帧的数据转换到频域，如公式 2.6。所有现代软件中都可以有效的计算快速傅里叶变换。

$$X_m[k] = \sum_{n=0}^{N-1} x_m[n]e^{-j2\pi knN} \quad (2.6)$$

傅里叶表征  $X_M[k]$  是一个很复杂的数，因为它包含了每一帧和每一个频率的频谱幅值（绝对幅值）和相位信息。为了提取特征，我们去掉了相位信息，只考虑幅值  $|X_m[k]|$ 。

频谱图描述了对语音信号进行 FFT 操作得到的 log 幅值（或者 log-power），如图 2.4（右）。横轴是帧索引（单位为 10ms），纵轴是频率，其范围是 0Hz 到采样率的一半，也就是对应的 Nyquist 频率。图中呈现的是 "speech recognition is cool stuff"。在频谱图中，黄色和红色区域表示该区域能量高。

## 2.2.3 Mel Filtering

从频谱图中可以看出高频的高能量区域大致对应着不发声的辅音，低频的高能量区域大致对应着发声的元音。频谱图中，发声区域的水平线（horizontal lines）呈现的是语音的谐波结构（harmonic structure）。

由于发声区域的谐波结构和不发声区域的随机噪声,频谱中存在着变数(variability)。为了移除这些变数,我们对幅度谱(magnitude spectrum)进行频谱光滑操作。受听觉系统处理语音信号的启发,我们对频谱图进行滤波器组操作(filterbank),该滤波器组对频率轴进行了 approximately logarithmic scale。也就是说随着频率的升高,滤波器也会变得更宽间隔更大。最常用于特征提取的 filterbank 是 **mel filterbank**。一个 mel filterbank 包含 40 个滤波器,如图2.3。每一个滤波器会对不同频率区间的能量谱求平均。

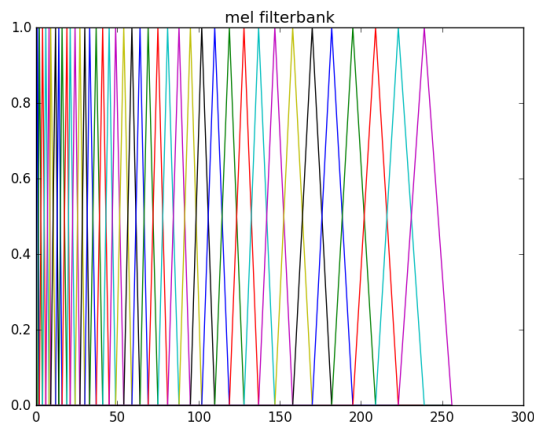


图 2.3: mel filterbank

mel filterbank 图在左边的滤波器很密集,在右边的滤波器间隔就比较远了。这是符合人耳听觉系统的,因为人耳对低频的信号更加敏感,对高频的信号不太敏感,所以低频的信息是更重要的,那么就需要多一些滤波器,从而提取更多有效特征。

P 维的 mel filterbank 的系数计算公式如2.7。一个 mel filterbank 一般会计算出 40 个系数,虽然现代系统有的会多一些或者少一些。平滑过多,系数就少一些,反之则反之。

$$X_{\text{mel}}[p] = \sum_k M[p, k] |X_m[k]|, \quad p = 0, 1, \dots, P-1 \quad (2.7)$$

### 2.2.4 Log Compression

特征提取的最后一步是对经过滤波器组得到的系数进行对数压缩。这个操作有助于压缩信号的动态范围,还能模拟听觉系统对声音的非线性压缩效果。我们把对数压缩后的输出称为“filterbank”系数。

将提取的特征以频谱图式的方式呈现出来之后,如图2.4(左),与原始信号频谱图(右)进行比较,我们可以看出沿着频率轴的 Fbank 系数要平滑得多,这是因为高频噪声和 pitch/谐波结构都被移除了。

除了上面的这些操作,在提取特征的过程中可能还会有一些其他的操作。其中有:



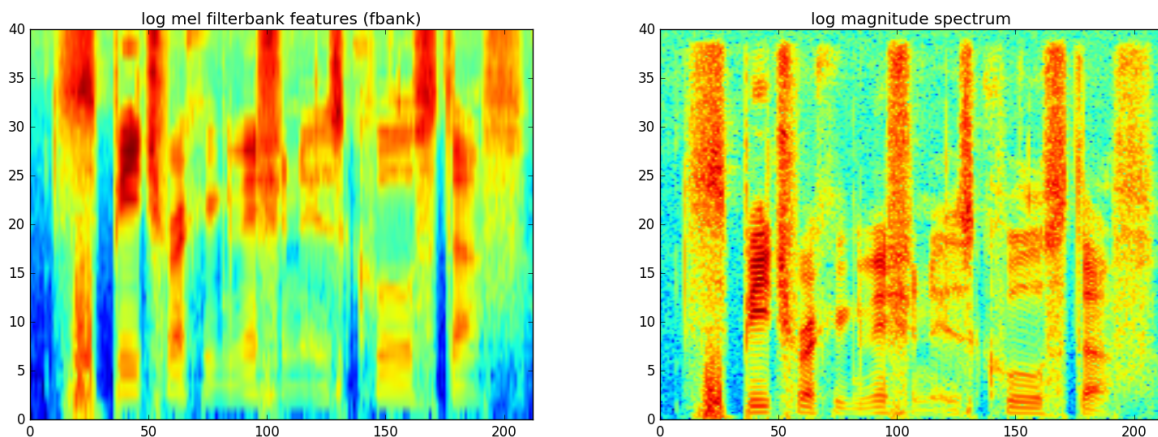


图 2.4: 提取 Fbank 特征 (左) 与原始信号的频谱图 (右)

1. **Dithering** (抖动): 在原始音频信号中加入一个很小的噪声, 为了防止在提取特征的时候出现数学问题, 尤其是出现  $\log 0$ ;
2. **DC-removal** (直流常数移除): 在提取特征之前, 去除音频中的常数偏置;
3. **Pre-emphasis** (预加重): 在提取特征之前用一个高通滤波器处理信号, 因为发声的语音部分低频能量比不发声的语音部分高频能量要大得多, 用一个高通滤波器来抵消下这个问题。实际操作的时候就是用了个简单地线性滤波器, 见公式 2.8, 其中  $\alpha = 0.97$ 。

$$y[n] = x[n] - \alpha x[n-1] \quad (2.8)$$

4. **DCT** (离散余弦变换): Fbank 系数是强相关的, 为了减弱这种相关性, 将上述步骤提取出来的 Fbank 特征值再经过 DCT。经过 DCT 之后的特征, 一般取前 13 个, 余下的由于所含信息不足舍弃了。DCT 公式如 2.9。

$$d_t = \frac{\sum_{n=1}^N n(c_{t+n} - c_{t-n})}{2 \sum_{n=1}^N n^2} \quad (2.9)$$

### 2.2.5 Feature Normalization

通信信道可能会对捕获的语音信号引入一些偏差 (恒定滤波)。比如说, 麦克风的频率响应不平稳, 此外, 即使相同语音的基础信号, 其信号增益的变化也可能导致计算的滤波器组系数的差异。这些信道的影响可以用时域上的卷积进行建模, 等价于频域表征的信号进行点乘 (elementwise multiplication)。

因此信道的影响可以用恒定滤波来建模 (constant filter), 如公式 2.10。

$$X_{t,\text{obs}}[k] = H[k]X_t[k] \quad (2.10)$$



其观测幅度为：

$$|X_{t,\text{obs}}[k]| = |H[k]| |X_t[k]| \quad (2.11)$$

如果我们对公式2.11两边取对数，并计算句子中所有帧的均值，则我们有公式2.12。

$$\begin{aligned} \mu_{\text{obs}} &= \frac{1}{T} \sum_t \log(|X_{t,\text{obs}}[k]|) \\ &= \frac{1}{T} \sum_t \log(|H[k]| |X_t[k]|) \\ &= \frac{1}{T} \sum_t \log(|H[k]|) + \frac{1}{T} \sum_t \log(|X_t[k]|) \end{aligned} \quad (2.12)$$

假设滤波器在时间轴上是常数，且语音信号的对数幅值均值为 0，那么公式2.12可以简化为2.13。

$$\mu_{t\text{to}bs} = \log(|H[k]|) \quad (2.13)$$

以上，如果我们计算出句子的对数幅值的均值，并且对句子中的每一帧都减去这个均值，这样我们就可以除去信号中所有的恒定信道效应。

为了简便，我们直接对取了  $\log$  之后 **fbank** 特征进行归一化（normalization）。为了对比一系列操作的结构，图2.5展现了原始信号的频谱图（左）、**fbank** 特征的频谱图（右）和对 **fbank** 特征归一化后的频谱图（下）。

## 2.2.6 Summary

为了从语音信号中提取语音识别所需的特征，我们希望提取到与声道形状有关的时变频谱信息，其由源滤波器模型中的一个滤波器建模，计算语句中的特征步骤如下：

1. 预处理信号，包括预加重和 **dithering**；
2. 将信号切割成有重叠部分的帧，一般帧长为 25ms，帧移为 10ms；
3. 对于每一帧：
  - 用汉明窗处理信号；
  - 使用 **FFT** 进行傅里叶变换；
  - 计算频谱的幅值；
  - 应用 **mel filterbank**；
  - 进行对数操作；
4. 如果需要进行信道补偿，则对每一帧的 **fbank** 系数进行均值归一化。

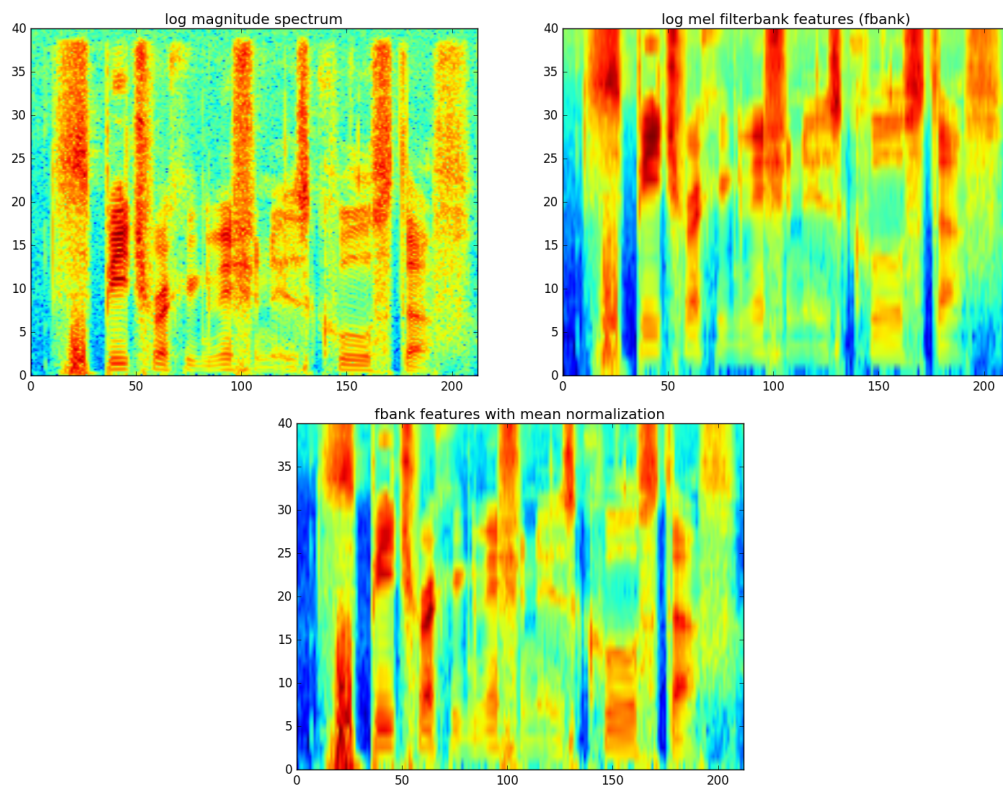


图 2.5: 原始信号的频谱图（左）、提取 Fbank 特征（右）和归一化后的 Fbank 特征（下）

### 2.2.7 Lab 2: Feature extraction for speech recognition

本模块有一个实验作业，标题为"Feature extraction for speech recognition"。

#### Required files:

- [M2\\_Wav2Feat\\_Single.py](#)
- [M2\\_Wav2Feat\\_Batch.py](#)
- [speech\\_sigproc.py](#)
- [htk\\_featio.py](#)

#### Instructions:

In this lab, you will write the core functions necessary to perform feature extraction on audio waveforms. Your program will convert an audio file to a sequence of log mel frequency filterbank ("FBANK") coefficients.

The basic steps in features extraction are

1. Pre-emphasis of the waveform
2. Dividing the signal into overlapping segments or frames
3. For each frame of audio:
  - Windowing the frame

- Computing the magnitude spectrum of the frame
- Applying the mel filterbank to the spectrum to create mel filterbank coefficients
- Applying a logarithm operation to the mel filterbank coefficient

In the lab, you will be supplied with python file called **speech\_sigproc.py**. This file contains a partially completed python class called **FrontEnd** that performs feature extraction, using methods that perform the steps listed above. The methods for dividing the signal into frames (step 2) will be provided for you, as will the code for generating the coefficients of the mel filterbank that is used in step 3c. You are responsible for filling in the code in all the remaining methods.

There are two top-level python scripts that call this class. The first is called **M2\_Wav2Feat\_Single.py**. This function reads a single pre-specified audio file, computes the features, and writes them to a feature file in HTK format.

In the first part of this lab, you are to complete the missing code in the **FrontEnd** class and then modify **M2\_Wav2Feat\_Single.py** to plot the following items:

1. Waveform
2. Mel frequency filterbank
3. Log mel filterbank coefficients

You can compare the figures to the figures below. Once the code is verified to be working, the feature extraction program should be used to create feature vector files for the training, development, and test sets. This will be done using **M2\_Wav2Feat\_Batch.py**. This program takes a command line argument **-set** (or **-s**) which takes as an argument either **train**, **dev**, or **test**. For example

```
$ python M2_Wav2Feat_Batch.py -set train
```

This program will use the code you write in the **FrontEnd** class to compute feature extraction for all the files in the LibriSpeech corpus. You need to call this program 3 times, once each for train, dev, and test sets.

When the training set features are computed (**-set train**) the code will also generate the global mean and precision (inverse standard deviation) of the features in the training set. These quantities will be stored in two ASCII files in the **am** direction for use by CNTK during acoustic model training in the next module.

Here are the outputs you should get from plotting:



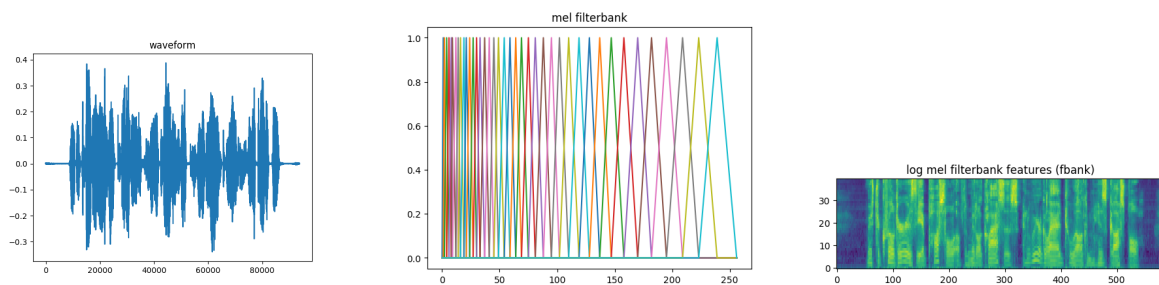


图 2.6: Lab 2 的期望输出

## 2.3 Acoustic Modeling

### 2.3.1 Introduction

本节讨论的是语音识别器中的声学模型。声学模型是一个混合模型，通过 DNN 得到逐帧的预测标签，再通过 HMM 将这些预测的音素转换成序列预测。HMM 常用于对离散时间序列事件的建模。HMM 的基本概念可以追溯到数十年前，而且 HMM 有很多应用。

### 2.3.2 Markov Chains

了解一点马尔科夫链（Markov Chains）对学习 HMM 大有帮助。马尔科夫链是一种对随机过程进行建模的方法。在马尔科夫链中，用一系列状态（states）来对离散时间进行建模。状态之间的运动由随机过程控制。

举例说明，在一个天气预测的应用中，状态为"Sunny"、"Partly Cloud"、"Cloudy"、和"Raining"。我们考虑某个连续五天的特定天气概率，比如  $P(p, p, c, r, s)$ ，我们可以使用贝叶斯规则将这个联合概率分布打散成一系列条件概率的乘积，如公式 2.14。

$$p(X_1, X_2, X_3, X_4, X_5) = p(X_5|X_4, X_3, X_2, X_1)p(X_4|X_3, X_2, X_1)p(X_3|X_2, X_1)p(X_2|X_1)p(X_1) \quad (2.14)$$

假设天气模型满足一阶马尔科夫假设，即满足公式 2.15。

$$p(X_i|X_1, \dots, X_{i-1}) = p(X_i|X_{i-1}) \quad (2.15)$$

那么连续五天天气的联合概率分布可以简化为公式 2.16。

$$\begin{aligned} p(X_1, X_2, X_3, X_4, X_5) &= p(X_5|X_4)p(X_4|X_3)p(X_3|X_2)p(X_2|X_1)p(X_1) \\ &= p(X_1) \prod_{i=2}^5 p(x_i|x_{i-1}) \end{aligned} \quad (2.16)$$

一个马尔科夫链的核心元素有**状态的定义**（此处为天气预测）和**转移概率**  $p(X_i|X_{i-1})$ ，转移概率描述的是从一个状态移动到另一个状态的概率值（也包括转移到自身状态）。

比如说，天气预报的一个完整的（大体完整的）马尔科夫链可以用图2.7来表示。

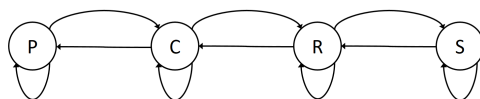


图 2.7: 天气预报模型的马尔科夫链

需要注意的是除了上面说的转移概率  $p(X_i|X_{i-1})$ ，我们还需要知道这个序列的第一个元素，即第一天的某种天气的概率值  $p(X_1)$ 。

所以除了状态清单和状态转移概率，我们还需要知道从马尔科夫链每一个状态开始的初始概率值。假设先验概率（每个状态的初始概率值）如公式2.17。

$$\begin{aligned} p(p) &= \pi_p \\ p(c) &= \pi_c \\ p(s) &= \pi_s \\ p(r) &= \pi_r \end{aligned} \tag{2.17}$$

现在我们回到这个例子，公式2.18说明了如何求  $P(p, p, c, r, s)$ 。

$$\begin{aligned} p(p, p, c, r, s) &= p(s|p, p, c, r)p(r|p, p, c)p(c|p, p)p(p|p)p(p) \\ &= p(s|r)p(r|c)p(c|p)p(p|p)p(p) \end{aligned} \tag{2.18}$$

上述 Markov Models 也可以称为可观测的马尔科夫模型（observable markov models）。因为当某个时刻所处状态已知，其输出是显性可见的，比如说今天会下雨。而 HMM 中的每个状态的输出并不是确定性事件或者确定性观察，而是事件或者观察的概率分布。这就使得 HMM 是双随机的（double stochastic）：状态之间的转移是随机的，状态的观测值也是随机的。

同样以天气举例子，我们可以将天气系统的 Markov model 转变成 HMM，将原先的状态"Sunny"、"Partly Cloud"、"Cloudy"、和"Raining" 替换成"Hot"、"Chilly"、"Cold" 和"Stormy"，这些状态以不同的概率值对应着不同的天气，如图2.8，当状态为"Hot" 的时候，观测值有 0.9 的概率为"Sunny"，有 0.1 的概率为"Partly Cloudy"。

因此一个包含了  $N$  个状态的 HMM 模型，其参数为：

- 转移矩阵  $A$ ，其元素为  $a_{ij}$ ，表征了从一个状态跳转到下一个状态的转移概率；
- 发射矩阵  $B$ ，其元素为  $b_{ix}$ ，表征了状态  $i$  的观测值为  $x$  的概率值，其中  $i = 1, 2, \dots, N$ ；
- 初始化状态的先验概率，即初始状态为某个状态的概率值， $\pi = \{\pi_1, \pi_2, \dots, \pi_N\}$ 。

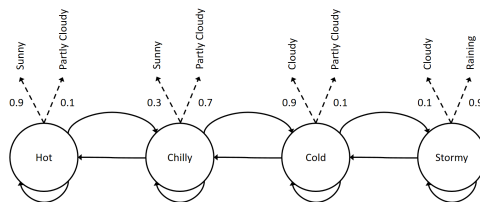


图 2.8: 天气预报模型的隐马尔科夫模型

我们将这些参数记成  $\Phi = \{A, B, \pi\}$ 。

HMM 有三个基本问题，每一个问题都有着很牛逼的解决办法，此处仅做一个简单介绍，详细的过程参考 1.5.3 节。

### The Evaluation Problem

给定一个模型和观测序列，这些观测序列由这个模型生成的概率有多大？

这个问题可以通过将所有可能的状态序列的概率加起来求得，这些状态序列都可以产生观测序列，通过转移概率和发射概率就可以求得。但是直接这么干是不划算的，因为其时间复杂度为  $O(N^T)$ ，其中  $T$  为观测序列的长度， $N$  为所有可能的观测序列。

前向算法是一种动态规划算法，相比较而言要有效得多。前向算法是在每个时间步都处理下子序列的前向概率。其在每个时间步都会存储  $N$  个值，时间复杂度为  $O(N^2T)$ 。

### The Decoding Problem

给定一个模型和观测序列，最有可能产生这个观测序列的状态序列是什么？

这个问题可以用 Viterbi 算法解决，Viterbi 算法被广泛的应用于语音识别的解码上，后续我们会介绍其是如何应用于语音识别以及如何将 Viterbi 算法并入训练准则之中。

### The Training Problem

给定一个模型和观测序列，如何去估计模型的参数  $\Phi$ ？

这个问题可以用 Baum-Welch 算法来解决，这个算法中就包含了前后向算法。本处的前向算法计算了  $t$  时刻处于状态  $i$  的所有可能的子序列的前向概率，后向算法计算了  $t$  时刻处于状态  $i$  为出发点直到  $T$  的所有可能的子序列的后向概率，此时时间步的初始值为  $t+1$ ，因此将前后向算法合并起来就是整个  $T$  时间步，所有在  $t$  时刻处于状态  $i$  的路径概率之和。

我们知道了每一个时刻各个状态的后验概率之后，Baum-Welch 算法将这些后验概率作为隐藏状态序列的直接观测值，并且更新模型的参数来优化目标函数。

## 2.3.3 Hidden Markov Models for Speech Recognition

在语音识别中 HMM 用于对声学观测（声学特征向量）建模，其建模单元为子词，比如说音素。

典型的做法是一个音素对应三个 HMM 的状态，用于对这个音素的开始、中间和结尾建模，每个状态有个自转移概率和转移到下一个状态的概率值，如图2.9。

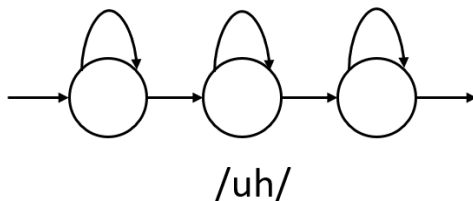


图 2.9: 音素"/uh/"的 HMM

词的 HMM 可以通过将其连续的音素 HMM 结合起来形成，比如说单词"cup"的 HMM 模型如图2.10。

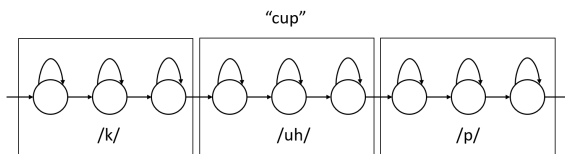


图 2.10: 单词"cup"的 HMM

因此，一个高质量的、包含每个单词音素表示的发音词典对于声学模型来说至关重要！

在 HMM-GMM 模型中，HMM 的状态都有一个概率分布，由 GMM 建模，其定义如公式2.19。

$$p(x|s) = \sum_m w_m N(x; \mu_m, \Sigma_m) \quad (2.19)$$

其中  $N(x; \mu_m, \Sigma_m)$  是高斯分布， $w_m$  是混合权重，其满足  $\sum_m w_m = 1$ 。因此模型每个状态都有自己的 GMM。Baum-Welch 算法会估计出所有的状态转移概率，以及 GMM 中的均值、方差和混合权重。

如今语音识别系统中不再使用 GMM 来对声学观测值建模了，DNN 取代了 GMM，DNN 的输出标签表示的是所有音素的所有状态。比如说有 40 个音素，每个音素是 3-state 的 HMM，那么 DNN 的输出就会有  $40 \times 3 = 120$  个标签。

这样的声学模型被称作混合模型，即 DNN-HMM，其取代了 GMM 来对声学进行建模。但是 HMM 其他的部分，尤其是 HMM 的拓扑结构和转移概率仍然会被用到。

### 2.3.3.1 Choice of Subword Units

我们现在已经知道了如何通过发音词典，链式的绑定音素来形成词 HMM。这些音素通常记作"Context Independent" Phones，即 CI phones。而事实上单个音素的发音很大程度上依赖于其前面的音素和后面的音素。比如说对于音素"/ah/"，在"bat"和"cap"中的发音就是不一样的。



有鉴于此, 如果使用"Context-Dependt"(CD) phones, 识别的准确率会有很明显的提升。因此对"bat"中, 我们用"/b-ah+t/"来表示"/ah/", 这样可以得到对应音素的上下文信息; 在"cap"中, 我们用"/k-ah+p/"来表示"/ah/"。比如说"cup"在句子"a cup of coffee"中, 即可用图2.11来对单词"cup"建模。

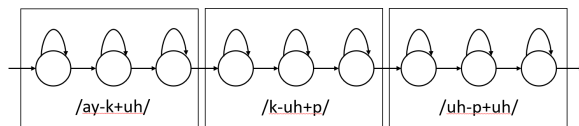


图 2.11: 单词"cup"的三音素 HMM

因为 CD phones 由三个连续的音素来表示的, 记作 **triphone**。有的系统用更长的上下文关系来建模, 比如"quinphones"就是 5 个连续的音素表示的, 不过这种的不常见。

当 ASR 使用 CI phones 来作为声学模型的建模单元时, 状态数并不多:  $N$  个音素乘以每个音素的状态数  $P$ 。美式英语一般有 40 个音素, 每个音素三个状态, 那么就有 120 个 CI states。那 triphone 作为建模单元的话, 会有多少个状态呢? triphone 有  $N^3$  个, 其状态数就有  $40^3 \times 3 = 192000$  个, 这会产生两个很严重的问题:

1. 用于训练每一个 triphone 的数据太少;
2. 可能有些训练集中没有的 triphone, 在测试集中出现了。

广为流传的一种办法是对多个 CD states 进行池化 (pooling), 这些 CD states 具备相似的性质, 将它们绑定到一起以形成单个"tied"或者"shared"HMM 状态。这些 tied states 叫做 **senone**。这些被绑定的原始状态都通过计算 senone 的声学分数来得到。

使用决策树 (decision tree) 将一个集合内的 CD triphones states 聚类到 senones 中。对每一个 CI phone 的每一个状态都会构建一个决策树。

决策树的聚类过程如下:

1. 将某个特定状态中, 共有中间音素的所有 triphone 放到一起形成根节点, 比如所有满足格式"/\*-p+\*/"的三音素的第二个状态;
2. 通过问一系列关于这个三音素的左边 context 或者右边 context 的语言学问题使得决策树生长, 这些问题都是二选一的问题。比如说: "is the left context phone a back vowel?" 或者 "is the right context phone voiced?". 在每一个节点, 选择训练数据似然度增加最大的问题;
3. 继续生长, 直到节点数达到预设值或者继续分割下去似然度的增加值小于某个阈值;
4. 决策树的叶子就定义这个 CD phone state 的 senones。

通过决策树聚类生成 senones 就可以解决上述两大问题, 首先, 映射成同一个 senones 的 triphone 之间可以共享数据, 所以其数据的估计也会比较稳定; 其次如果测试的时候有个 triphone 没有出现过, 其对应的 senone 就是顺着决策树搜索, 最终找到合适的输出。

几乎所有基于 phone 的语音识别系统都会使用 senone 作为包含上下文信息的建模单元。一



般工业级的语音识别器大概有 10000 个 senones。这虽然比 120 个 CI state 多不少，但是比 192000 少得多了。

关于决策树的生成，参考 Ross Young 的博文<sup>(16)</sup>，如图 2.12。

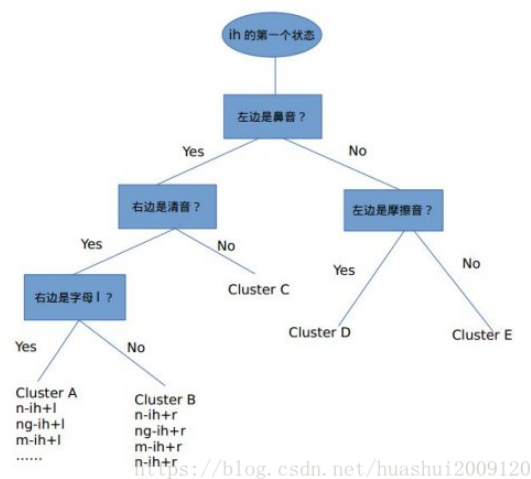


图 2.12: senones 的生成

此外，kaldi 中的决策树规则是由算法生成的，具体看 3.5。

### 2.3.4 Deep Neural Network Acoustic Models

这几年语音识别领域出现的最重要的改进就是使用 DNN 声学模型。前面也提到过，混合 DNN 系统使用单个 DNN 替换了原先的 GMMs（每个 senone 都有一个 GMM），DNN 的输出标签对应着 senones。

训练分类神经网络模型最常用的目标函数是交叉熵（cross entropy），对于一个  $M$  类的多分类任务，比如说 senones 分类，单个样本的目标函数如公式 2.20。

$$E = - \sum_{m=1}^M t_m \log y_m \quad (2.20)$$

其中  $t_m$  是标签值（如果这个样本是第  $m$  类， $t_m$  为 1，否则为 0）， $y_m$  是网络的实际输出值，神经网络最后一层的输出经过 softmax 函数映射到概率空间。因此对于每一帧来说，有个  $M$  维的 one-hot 向量，其对应的 groundtruth，即本帧真正正对应的 senone。

所以，每一帧我们都得给个标签。

为了给所有训练数据的每一帧加上对应 senones 的标签，我们需要对所有的数据进行强制对齐。强制对齐的操作基本上就是使用 HMM 模型对每一条语音数据进行解码，限制其解码出来的路径会产生正确的参考文本。强制对齐产生单条最有可能的路径，因此就得到了语句中每一帧对应的 senone 的标签。

强制对齐操作需要一个语音识别系统，可以是一个基于 GMM 的初始模型，也可以是一个训练好的 DNN 模型，但是前提是待对齐的和训练好的模型他们的 `senones` 是一样的。

强制对齐的输出文件包含了每一句的起始帧和终止帧，以及对应的 `senone` 标签的。不同的工具，这个文件的格式也不一样，以 HTK 为例，输出是一个 MLF 格式的文件，如图 2.13，该文件每一列表示的分别是：

1. 开始时间，单位是 100ns；
2. 结束时间，单位是 100ns；
3. Senone 的 ID；
4. Senone 片段的声学模型得分；
5. CD triphone HMM 模型；
6. triphone HMM 模型的声学模型得分；
7. Senone 对应文本中的词

```
#!MLF!#
"test utterance.lab"
0      17700000 sil[2]      0.031854 sil 719.849243 <s>
17700000 27100000 sil[3]      140.979706
27100000 34800000 sil[4]      152.653412
34800000 35500000 g_s2_6      8.032190 sil-g+uh 27.682287 good
35500000 35800000 g_s3_25      6.538751
35800000 36100000 g_s4_11      13.111345
36100000 36200000 uh_s2_7      3.356279 g-uh+d 11.706375
36200000 36300000 uh_s3_15      4.878324
36300000 36400000 uh_s4_15      3.471773
36400000 36500000 d_s2_51      1.766836 uh-d+m 16.887865
36500000 36800000 d_s3_44      11.397981
36800000 36900000 d_s4_26      3.723048
36900000 37200000 m_s2_7      13.519948 d-m+ao 38.862099 morning
37200000 37500000 m_s3_21      20.174339
37500000 37600000 m_s4_79      5.167810
37600000 37700000 ao_s2_51      3.014942 m-ao+r 11.643064
37700000 37900000 ao_s3_56      5.589801
37900000 38000000 ao_s4_46      3.038320
38000000 38100000 r_s2_178      5.339163 ao-r+n 26.518557
38100000 38300000 r_s3_100      14.487432
38300000 38400000 r_s4_61      6.691961
38400000 38600000 n_s2_28      17.327641 r-n+ih 35.669243
38600000 38700000 n_s3_145      9.536380
38700000 38800000 n_s4_112      8.805222
38800000 38900000 ih_s2_15      5.309596 n-ih+ng 25.129721
38900000 39000000 ih_s3_160      5.600624
39000000 39200000 ih_s4_123      14.219500
39200000 39300000 ng_s2_8      4.993470 ih-ng+m 7.568606
39300000 39400000 ng_s3_14      2.377748
39400000 39500000 ng_s4_4      0.197388
```

图 2.13: MLF 文件示例

有这样的文件或者类似的文件之后，我们就能得到训练 DNN 所需的标签了。

### 2.3.5 Training Feedforward Deep Neural Networks

构建声学模型的神经网络中，最简单也是最常见的就是前馈神经网络（feed-forward neural network）。前馈神经网络的相关知识网上到处都是，本处咱们只聊基于 DNN 的声学模型的关键信息。

尽管我们训练一个 DNN 是为了预测每一个输入帧的标签，如果我们提供一个 context window 内的帧作为输入的话，而不是使用单独的一帧作为输入，分类的效果会好很多。具体来讲呢，对于  $t$  时刻的某一帧，DNN 的输入是本帧的前  $N$  帧、本帧和本帧的后  $N$  帧，这个 context window 是对称的。因此假设  $x_t$  是  $t$  时刻的特征向量，那么对应该时刻 DNN 的输入为：

$$X_t = [x_{t-N}, x_{t-N+1}, \dots, x_t, \dots, x_{t+N-1}, x_{t+N}] \quad (2.21)$$

一般呢， $N$  的取值范围在 5 到 11 之间，这个取决于训练集的体量。越大的 context window 提供了越多的信息，但是同时会增加输入层的参数矩阵的大小，如果没有足够的数据的话，训练起来会比较困难。

通常还会做一些数据增强，通过加入这些特征的时间导数，记作 delta 特征。这些特征可以通过简单地差分或者复杂的回归公式计算得到，比如说：

$$\begin{aligned} \Delta x_t &= x_{t+2} - x_{t-2} \\ \Delta^2 x_t &= \Delta x_{t+2} - \Delta x_{t-2} \end{aligned} \quad (2.22)$$

这样的话，DNN 中每一帧的输入就是一个 context window 堆叠（stack）的特征，其包含原始的特征向量，delta features 和 delta-delta features，如下：

$$x_t, \Delta x_t, \Delta^2 x_t \quad (2.23)$$

上述的特征经过一系列全连接层，最终通过 softmax 函数映射成每个 senone 的概率值，以交叉熵为目标函数，通过 BP 算法来进行参数的迭代，直到训练结束。整个过程如图 2.14。

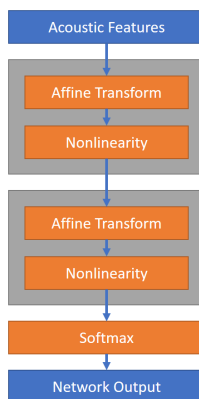


图 2.14: 基于 DNN 的声学模型训练过程

除了基于前馈神经网络的声学模型外，我们经常使用的还有 RNN。与 DNN 不同的是，RNN

处理的是序列数据，其权重具备时域依赖性。RNN 隐含层的输出如公式2.24。

$$h_t^i = f(W^i h_t^{i-1} + U^i h_{t-1}^i + c^i) \quad (2.24)$$

其中  $f(\cdot)$  是一个非线性函数，比如 sigmoid 函数或者 relu 函数， $i$  是网络的层数， $t$  表示的是帧数或者时间索引，输入  $h_t^{i-1}$  等于第  $i-1$  层的输出，其实就是上一层的输出，第一层的 RNN 的输入就等于  $x_t$ 。

与 DNN 相比，循环层的输出既依赖于当前输入，又依赖于上一个时刻的输出。如果你比较熟悉信号处理中的滤波操作，RNN 层就可以看作是一个非线性无限脉冲响应（nonlinear infinite impulse response, IIR）滤波器。

对于离线操作的应用来说，延迟并不是很重要的一个指标，所以可以利用双向 RNN（bi-RNN）来提升模型的性能。这种情况下，每一层都会有一些参数处理前向序列，同时还会有另一些参数处理后向序列。这样就会有二个输出，将它们拼接（concatenate）到一起作为下一层的输入，如公式2.25，其中下角标  $f$  和  $b$  分别指前向和后向。

$$\begin{aligned} \vec{h}_t^i &= f(W_f^i h_t^{i-1} + U_f^i h_{t-1}^i + c_f^i) \\ \overleftarrow{h}_t^i &= f(W_b^i h_t^{i-1} + U_b^i h_{t+1}^i + c_b^i) \\ h_t^i &= [\vec{h}_t^i, \overleftarrow{h}_t^i] \end{aligned} \quad (2.25)$$

由于 RNN 可以学习到特征向量序列之间的时域模式，因此很适合用于声学建模。为了训练 RNN，训练序列的序列特性就会保留下来。因此与其用 DNN 中基于单帧随机化的方式来训练模型，不如采用基于语句随机化的方式。这样的情况下，语句被随机化打乱了，但是语句中的序列特性还是会保留下来。

因为 RNN 本身的特性决定了模型会学习到数据沿着时间线上的关联关系，就不需要再用一个很宽的 context window 来捕获每一帧的上下文关系了。对于单向 RNN 来说，提供一些 future context 帧是有积极作用的，但是这个数量比 DNN 要少很多。在双向 RNN 中，提供 context window 就失去了意义，因为处理每一帧，网络都可以看到整句的信息，自然会吸收上下文的关联关系。

训练 RNN 同样使用交叉熵作为目标函数，只是在计算梯度的时候有些微不同。基于这个模型的时域特性，使用 BP 的变体 BPTT 算法来更新网络的参数。我们想想看单向 RNN 的特点，当前时刻的输出不仅依赖于当前输入，还依赖于之前所有时刻的输入。

正如标准的 BP 算法，BPTT 使用梯度下降来优化这个模型，目标函数关于模型参数的梯度是通过链式规则计算得到的，这样导数中就包含了很多梯度的乘积（取决于前面已经处理了多少个时间步）。这些梯度都不存在任何限制，很有可能就会出现梯度接近于 0（网络不再学习了）或者接近于无穷大（网络训练就不稳定，也无法收敛）。这个问题就是梯度消失或爆炸。

为了抑制梯度消失的问题，我们一般有两种办法：

1. 采用 RNN 的变体结构，比如 LSTM，这个我们一会再说；
2. 限制 BPTT 回溯的长度，就是说只允许 BPTT 利用限定长度的历史信息，这样就会限制连乘梯度的数量。

为了抑制梯度爆炸的问题，可以用 **gradient clipping**（梯度裁剪）的方法。设定一个梯度的上限，对于模型中所有参数的梯度，绝对值超过这个上限的，就让其梯度等于这个上限。

前面提到了 LSTM 可以抑制梯度消失或爆炸的问题以更好的学习训练集中的长时关系。LSTM 中有个概念叫 **cell**，这个东西就是用来存储状态信息的。**cell** 中的信息可以随着时间变化一直被保存着，也可以擦除过去的信息，存入当前的信息。这些操作都依赖于 **gate** 来实现。

**Gate** 接近于 0，信息就无法通过；**Gate** 接近于 1，信息就可以通过。输入门决定了是否将当前时间步的信息存储到 **cell** 中；遗忘门决定了是保存 **cell** 中信息还是擦掉重写入当前的信息；输出门决定了是否将 **cell** 中的信息传出去。LSTM 单个神经元的内部如图 2.15。更多关于 LSTM 的资料参考<sup>(3)</sup>。

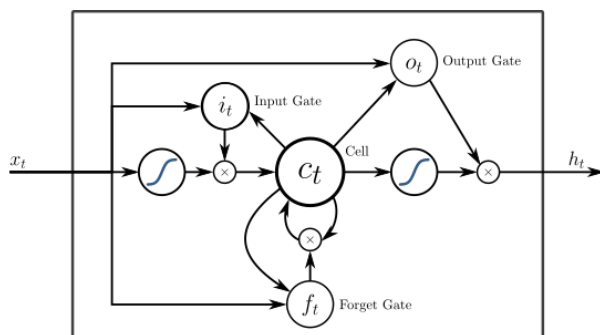


图 2.15: 基于 DNN 的声学模型训练过程

其他的 LSTM 的变体，比如说 GRU，就是 LSTM 的一个简化版，GRU 与 LSTM 的表现相当，但是参数更少，见 11.1.1。

### 2.3.6 Using a Sequence based Objective Function

RNN 是一个序列声学模型，因为它们可以对声学特征向量作为一个时域上连续的序列进行建模。但是其目标函数仍然是帧独立的，就是说在计算目标函数的时候，每一帧的损失是独立计算的。然而因为语音识别是一个序列分类任务，如果我们能搞个基于序列的目标函数，模型的性能应该会有所提升。

GMM-HMM 中就有基于序列的目标函数，如今也可以应用于 DNN-HMM 中。基于帧的交叉熵和序列判别性目标函数的区别就在于基于序列的目标函数更好的对解码过程进行了建模。

具体来说，序列区分性训练中应用了语言模型，来找到那些与参考文本很相似的输出，这些和 **groundtruth** 很接近的序列正是网络需要区分开的。相对而言，基于帧的交叉熵作为目标函数，所有不正确的分类都会受到惩罚，即便这些不正确的分类在解码的时候根本不会出现，当

然解码的时候需要用的 HMM 拓扑结构或者语言模型。在序列区分性训练中，正确分类的这些 competitors 是由解码训练数据得到的。

最常见的一个序列区分性训练目标函数是最大互信息 (maximum mutual information, MMI), 如公式2.26。

$$F_{MMI} = \sum_u \frac{\log p(X_u|S_u)p(W_u)}{\sum_{w'} \log p(X_u|S_{w'})p(W')} \quad (2.26)$$

其中  $u$  是训练集中句子的索引,  $W_u$  是语句  $u$  的参考词序列,  $S_u$  是对应的状态序列, 分母呢是所有可能的词序列的概率和, 分母通常由 word lattice 近似得到, 解码的时候 word lattice 能得到可能输出的词序列。

使用这个目标函数, 让声学模型的训练复杂了很多, 但是确实能提高识别性能。同样有很多 MMI 的变体, 比如 Minimum Phone Error (MPE), 和 state-level 的 Minimum Bayes Risk (SMBRa), 见2.7.5。

### 2.3.7 Decoding with Neural Network Acoustic Models

DNN 声学模型计算的是在已知声学特征的条件下, 某个 senone 的条件概率, 即  $p(s|x_t)$ , 这些 state-level 的后验概率必须转换成状态的似然度  $p(x_t|s)$ , 只有这样我们才能够使用 HMM 来解码。如何得到  $p(x_t|s)$  呢?

根据贝叶斯规则:

$$p(x_t|s) = \frac{p(s|x_t)p(x_t)}{p(s)} \propto \frac{p(s|x_t)}{p(s)} \quad (2.27)$$

公式2.27之所以成立是因为对于所有的 senones 观测鲜艳  $p(x_t)$  是个常数, 其只是说给似然度增加了个系数而已, 因此舍掉它。因此似然度  $p(x_t|s)$  等于 DNN 的输出除以 senone 的先验  $p(s)$ 。而 senone 的先验可以通过简单地统计训练集中对应 senone 出现的次数来估计。

上述似然度叫做 **scaled likelihood**, 指代的是其是通过 senone 的先验 scale 其后验概率得到的。

### 2.3.8 Lab 3

**Required files:**

- [M3\\_Train\\_AM.py](#)
- [M3\\_Plot\\_Training.py](#)

**Instructions:**

In this lab, we will use the features generated in the previous lab along with the phoneme state

alignments provided in the course materials to train two different neural network acoustic models, a DNN and an RNN.

The inputs to the training program are:

- **lists/feat\_train.rscp, lists/feat\_dev.rscp** - List of training and dev feature files, stored in a format called RSCP. This standard for relative SCP file, where SCP is HTK-shorthand for script file. It is simply a list of files in the two sets. The dev set is used in training to monitor overfitting and perform early stopping. These files should have been generated as part of completing lab 2.
- **am/feat\_mean.ascii, am/feat\_invstddev.ascii** - The global mean and precision (inverse standard deviation) of the training features, also computed in lab 2
- **am/labels\_all.cimlf**- The phoneme-state alignments that have been generated as a result of forced alignment of the data to an initial acoustic model. Generating this file requires the construction of a GMM-HMM acoustic model which is outside the scope of this course, so we are providing it to you. The labels for both the training and dev data are in this file.
- **am/labels\_ciphones**- The list of phoneme state symbols which correspond to the output labels of the neural network acoustic model
- **am/labels\_ciprior.ascii** - The prior probabilities of the phoneme state symbols, obtained by counting the occurrences of these labels in the training data.

The training, dev, and test RSCP files and the training set global mean and precision files were generated by the lab in Module 2. The remaining files have been provided for you and are in the **am** directory.

### PART1: TRAINING A FEEDFORWARD DNN

We have provided a python program called **M3\_Train\_AM.py** which will train a feed-forward deep network acoustic model using the files described above. The program is currently configured to train a network with the following hyperparameters:

- 4 hidden layers of size 512 hidden units per layer.
- 120 output units corresponding to the phoneme states
- input context window of 23 frames, which means the input to the network for a given frame is the current frame plus 11 frames in the past and 11 frames in the future
- minibatch size of 256
- Learning is performed with Momentum SGD with a learning rate of  $1e - 04$  per sample with momentum as a time constant of 2500
- One epoch is defined as a complete pass of the training data and training will run for 100 epochs
- The development set will be evaluated every 5 epochs.

This can be executed by running:





```
1 python M3_Train_AM.py
```

On a GTX 965M GPU running on a laptop, the network trained at a rate of 63,000 samples/sec or about 20 seconds per epoch. Thus 100 epochs will run in 2000 seconds or about 30 minutes.

After 100 epochs, the result of training, obtained from the end of the log file, was

```
1 Finished Epoch[100 of 100]: [CE_Training] loss = 1.036854 * 1257104,
    metric = 32.74% * 1257104 17.146s (73317.6 samples/s);
2 Finished Evaluation [20]: Minibatch[1-11573]: metric = 44.26% * 370331;
```

Thus, the training set has a cross entropy of 1.04 per sample, and a 32.74% frame error rate, while the held-out dev set has a frame error rate of 44.3

After training is complete, you can visualize the training progress using **M3\_Plot\_Training.py**. It takes a CNTK log file as input and will plot epoch vs. cross-entropy of the training set on one figure and epoch vs. frame error rate of the training and development sets on another figure.

```
1 python M3_Plot_Training.py --log <logfile>
```

For this experiment, **<logfile>** would be **../am/dnn/log**

Here is an example of the figure produced by this script.

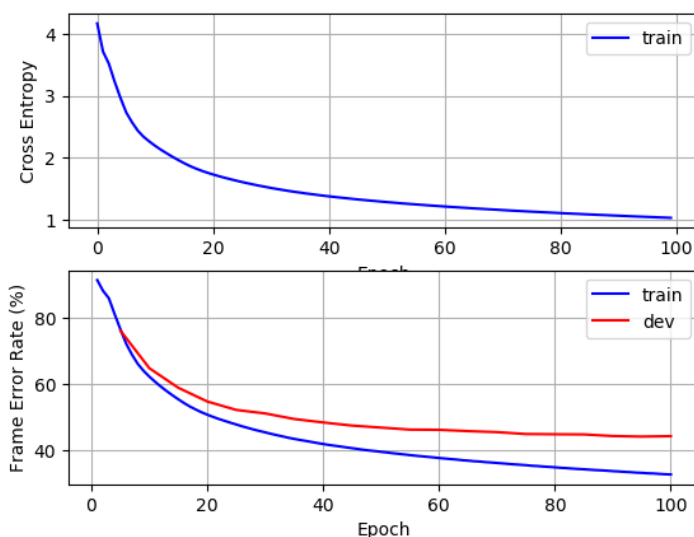


图 2.16: DNN 声学模型训练图

As you can see from the figure, overfitting has not yet occurred as the development set performance is still the best in the final epoch. It is possible that small additional improvements can be obtained with additional training iterations.



You can now experiment with this neural network training script. You can modify the various hyperparameters to see if the performance can be further improved. For example, you can vary the

- Number of layers
- Number of hidden units in each layer
- Learning rate
- Minibatch size
- Number of epochs
- Learning algorithm (see the CNTK documentation for details on using other learners, such as Adam or AdaGrad)

## PART 2: TRAINING A RECURRENT NEURAL NETWORK

In the second part of this lab, you will modify the code to train a Bidirectional LSTM (BLSTM) network, a type of recurrent neural network.

To train an BLSTM, there are several changes in the code that you should be aware of.

1. In DNN training, all frames (samples) are processed independently, so the frames in the training data are randomized across all utterances. In RNN training, the network is trying to learn temporal patterns in the speech sequence, so the order of the utterances can be randomized but the utterances themselves must be kept intact. Thus, we set **frame\_mode=False** in the **MinibatchSource** instantiated by **create\_mb\_source()**.
2. Change the network creation to create a BLSTM

- In **create\_network()**, we've created a function called **MyBLSTMLayer** as specified below. This function uses the **Optimized\_RNN Stack** functionality in CNTK. A complete description and additional examples can be found in the CNTK documentation. One thing to be aware of is that with a BLSTM, the size of the hidden layer is actually applied to both directions. Thus, setting the number of hidden units to 512 means that both the forward and backward layers consist of 512 cells. The outputs of the forward and backward layer are then concatenated forming an output of 1024 units. This is then projected back to 512 using the weight matrix **W**.

```

1 def MyBLSTMLayer(hidden_size=128, num_layers=2):
2     W = C.Parameter((C.InferredDimension, hidden_size), init=C.
3         he_normal(1.0), name='rnn_parameters')
4     def _func(operand):
5         return C.optimized_rnnstack(operand, weights=W, hidden_size=
6             hidden_size, num_layers=num_layers, bidirectional=True,
7             recurrent_op='lstm')
```

```
5 return _func
```

- The code calls **MyBLSTMLayer** when the **model\_type** is **BLSTM**. We've reduced the number of hidden layers to 2, since the BLSTM layers have more total parameters than the DNN layers.
3. For utterance based processing, entire utterance needs to be processed during training. Thus the minibatch size specifies the total number of frames to process but will pack multiple utterances together if possible. Setting the minibatch size to a larger number will allow for efficient processing with multiple utterances in each minibatch size. We have set the minibatch size to 4096.

The training the BLSTM model, you can execute the following command.

```
1 python M3_Train_AM.py --type BLSTM
```

Because of the sequential nature of the BLSTM processing, they are inherently less parallelizable, and thus, train much slower than DNNs. On a GTX 965M GPU running on a laptop, the network trained at a rate of 440 seconds per epoch, or 20 times slower than the DNN. Thus, we will only train for 10 epochs to keep processing time reasonable.

Here too, you can use **M3\_Plot\_Training.py** to inspect the learning schedule in training. And again, if you are interested, you can vary the hyperparameters to try to find a better solution.

## 2.4 Language Modeling

### 2.4.1 Introduction

本模块介绍语言模型的基础知识。语言模型是语音识别的一个组成部分，其估计的是可能的话语的先验概率  $P(W)$ 。回忆一下前面说过的语音识别的基本公式，这个概率是和声学模型的似然度结合起来以获得最好的输出假设的，如公式2.28。

$$\hat{W} = \arg \max_W P(O|W)P(W) \quad (2.28)$$

因此 **LM** 表征的是识别器所具备的知识，这个知识表达了可能的词序列是什么，即使这个识别器还没有听到任何实际的音频，因为 **LM** 表达的是词序列的先验知识。**LM** 应当给可能的句子高的概率值，不太可能的句子低的概率值，而不是硬性的规定语法规则来允许一些句子的存在，不允许另一些句子的存在。**LM** 不应当定下一个绝对的规则是因为没有人知道说话人会说什么。**LM** 赋予某个句子的概率值也不是通过语言学或者什么规则得到的，它和声学模型一样，来源于数据。因此我们让实际数据的统计学特性来决定在一种语言中，在一个场景中或者在一个应用中，什么词是比较可能出现的。

注意下关于术语的使用，在语言模型中，我们一般把整条音频对应的句子（sentence）称作词序列（word sequence），这里面没有任何暗示表明这个词序列就是一个服从传统语法的、正确的、完整的句子。事实上一个句子的 LM 针对的是某个场景下说话人可能说出的任何句子。

## 2.4.2 N gram Models

### 2.4.2.1 Vocabulary

我们需要给每一个可能的词序列赋予一个概率值：

$$W = w_1 w_2 \dots w_n \quad (2.29)$$

其中  $n$  是词的个数，原则上来说这个数量是没有上限的。

首先，我们将问题简化下，用一个有限集合来限制词的可选择空间，即 LM 的 vocabulary。注意 LM 的 vocabulary 也就是识别器的 vocabulary，我们是没办法识别一个不被 LM 包括在内的词的，其概率值无限接近于 0 的。

vocabulary 之外的词呢一般称为是 out-of-vocabulary 的词，简称 OOV。在输入数据中如果存在一个 OOV，那么识别器就至少会产生一个错误，所以选择的词典应当最小化 OOV 出现的可能性。常用的策略是选择出现在训练集中的最频繁的那些词，比如说我们选择训练集中前  $N$  个出现的最频繁的词或者所有出现次数大于  $K$  的词，这个  $N$  或者  $K$  要选择适当。一般呢最优的词典其大小能够在识别速度和准确率之间找到一个平衡，因为词典越大意味着解码的时候效率越低，而减小 OOVs 又能提高准确率。当然添加那种基本上不咋出现的词对准确率的影响可以忽略不计，甚至可能降低准确率，因为解码搜索的时候会产生一些错误也会造成声学的模糊性。

### 2.4.2.2 Markov factorization and N-grams

即便有了一个有限数量的词典，我们同样有无限个词序列，很显然，列举出所有可能的句子是不现实的。就算说我们概念上可以这么干，这样做也没办法去评估 LM 的参数，因为大部分的句子出现的可能性很小，可能性越小的概率值就得需要更多的数据来评估，这样得到的概率值才是可靠的。

为了解决这两个问题，我们使用一种技巧：使用链式规则将句子的概率值分解成一堆词的概率连乘，如公式 2.30，然后应用马尔科夫假设限制状态和参数的数量。

$$P(W) = P(w_1) \times P(w_2|w_1) \times P(w_3|w_1 w_2) \times \dots \times P(w_n|w_1 \dots w_{n-1}) \quad (2.30)$$

如果我们假设是 LM 是一阶马尔科夫模型，那么我们就可以改写  $P(W)$  如公式 2.31，这就是 bigram 模型，因为模型只利用了相邻的两个词的统计特性，这样每个词靠前面出现的那个词就

可以预测出来。类似的，二阶马尔科夫的 LM 就称为 **trigram** 模型，当前词是靠前面出现的两个词预测出来的。

$$P(W) = P(w_1) \times P(w_2|w_1) \times P(w_3|w_2) \times \dots \times P(w_n|w_{n-1}) \quad (2.31)$$

这种技巧泛化之后我们就称之为 **N-gram** 模型，即当前词依赖于前面出现的  $N - 1$  个词。事实表明 **trigram** 比 **bigram** 效果要好不少，但是随着  $N$  的增加，性能提升的就越少。因此实际中我们很少使用超过 **4-gram** 或者 **5-gram** 以上的 LM。实验室一般用的 **trigram** 居多，后续讲解我们都以 **bigram** 来说明，但是要切记，这些 **bigram** 的特点或者 **bigram** 的应用都是可以泛化到 **N-gram** 的。



### 2.4.2.3 Sentence start and end

### 2.4.2.4 N-gram probability estimation

### 2.4.2.5 N-gram smoothing and discounting

### 2.4.2.6 Back-off in N-gram models

## 2.4.3 Language Model Evaluation

## 2.4.4 Operations on Language Models

## 2.4.5 Advanced LM Topics

## 2.4.6 Lab 4

# 2.5 Speech Decoding

## 2.5.1 Overview

## 2.5.2 Weighted Finite State Transducers

## 2.5.3 WFSTs and Acceptors

## 2.5.4 Graph Composition

## 2.5.5 Lab 5

# 2.6 Advanced Acoustic Modeling

## 2.6.1 Improved Objective Functions

## 2.6.2 Sequential Objective Function

## 2.6.3 Connectionist Temporal Classification

## 2.6.4 Sequence Discriminative Objective Functions

## 2.6.5 Lab 6

# 2.7 补充知识点

## 2.7.1 傅里叶变换

## 2.7.2 Nyquist 定理

## 2.7.3 编辑距离

给定 2 个字符串  $a, b$ , 编辑距离是将  $a$  转换为  $b$  的最少操作次数, 其操作只允许以下三种:

1. 插入一个字符 (Insert): 例如  $fj \rightarrow fxj$  *Delete*  $fxj \rightarrow fj$
2. 替换一个字符 (Substitute): 例如  $jxj \rightarrow fxj$

**Solutions:**

采用分而治之的思想来解决这个问题, 将复杂的问题分解成相似的子问题。

假设字符串  $a$ , 共有  $m$  位, 从  $a[1]$  到  $a[m]$ ; 字符串  $b$ , 共有  $n$  位, 从  $b[1]$  到  $b[n]$ 。

设  $d[i][j]$  表示字符串  $a[1]$  到  $a[i]$  转换为  $b[1]$  到  $b[j]$  的编辑距离。那么有如下递归规律, 其中  $a[i]$  和  $b[j]$  分别为

#### 2.7.4 nonlinear infinite impulse response (IIR)

#### 2.7.5 序列区分性训练中的目标函数详解

#### 2.7.6 概率与似然度的是是非非

## 第 3 章 Kaldi 学习笔记

### 3.1 kaldi 中的数据扰动

kaldi 程序中对原始数据进行扰动以达到数据增强的效果，一般是在单音素对齐，三音素对齐之后，在生成 ivector 的时候进行扰动处理，扰动有两种方式：速度扰动和音量扰动。如果存在 segment 文件的话，那么对应的起始和终止时间点会存在 segment 文件中，提取特征时会根据这个 segment 中存储的时间节点进行操作；如果不存在的话，相当于整段 wav 音频都是有效的，那么起始时间点和 wav 文件相同。扰动也会根据 segment 文件的有无进行对应的操作。

#### 3.1.1 速度扰动

速度扰动一般是对音频进行加速和减速，根据 Povey 大佬的论文《Audio Augmentation for Speech Recognition》中的第二部分 Audio Perturbation，对 mel 频谱进行一个偏移就能得到类似加速和减速的效果。首先定义一个扰动因子  $\alpha$ ，假定 segment 中某一段音频的起始时间和终止时间为  $t_1$  和  $t_2$ ，那么新的音频起始时间和终止时间计算方式如公式 3.1。

$$\begin{aligned} t'_1 &= \frac{t_1}{\alpha} \\ t'_2 &= \frac{t_2}{\alpha} \end{aligned} \tag{3.1}$$

kaldi 一般取  $\alpha$  为 0.9 和 1.1 以达到加速和减速的目的。得到了 segment 文件之后，在 wav.scp 文件中存储原始音频的位置，加速后音频 sox 指令和减速后音频 sox 指令。其详细的脚本指令见代码 utils/data/perturb\_data\_dir\_speed.sh 第 74 行。最终重新提取特征存于代码根目录下 mfcc\_perturbed 文件夹中。由于速度扰动对音频时间轴有改动，因此此时需要对音频进行重新对齐的操作。

#### 3.1.2 音量扰动

音量扰动一般是对音频进行增大音量和减小音量。音量增加或者减小的幅度默认取 [0.125, 2] 之间的正态分布值。使用 sox 工具中的 "sox - -vol volume" 来进行实际操作。其详细的脚本指令见代码 utils/data/perturb\_data\_dir\_volume.sh 第 71 行，其 sox 操作见代码 utils/data/internal/perturb\_volume.py。其重新提取的特征位于代码根目录下 mfcc\_hires 文件夹中。此时由于仅仅对音频进行音量大小的扰动，并没有对时间维度进行操作，因此无需再进行一遍对齐操作，其标签对齐直接采用上一步即速度扰动后生成的对齐结果。此时重新提取特征时，MFCC 特征的维度是 40d。其原因是 40d 的 MFCC 和 40d 的 Fbank 维度相同，保存的信息量相似，同时 MFCC 由于

其相关性较弱 (DCT 去相关), 所以能更好的压缩特征, 因此 Kaldi 一般都是采用 40d 的 MFCC 作为神经网络的输入特征 (见 kaldi 的各个 egs 里 conf 下 mfcc\_hires.conf)。

"Config for high-resolution MFCC features, intended for neural network training. Note: we keep all cepstra, so it has the same info as filterbank features, but MFCC is more easily compressible (because less correlated) which is why we prefer this method. "

## 3.2 kaldi 中的 UBM

通用背景模型 **UBM**(Universal Background Model)

## 3.3 kaldi 通过 lattice 输出语音对齐音素和词

我们通过解码之后得到一堆的 lat.\*.gz 文件, 这些文件是 lattice 对齐后生成的对齐文件, 其为压缩格式, 所以首先通过 gunzip 指令对齐解压, 我们以 lat.1.gz 为例来讲这一节的知识。

```
gunzip exp/chain/tdnn7q_sp_online/decode_data_tgsmall/lat.1.gz
```

这样会在 exp/chain/tdnn7q\_sp\_online/decode\_data\_tgsmall/ 下生成一个 lat.1 文件, 原先的 lat.1.gz 消失不见了……lat.1 文件是二进制的格式, 其由指令 online2-wav-nnet3-latgen-fatser 生成。

## 3.4 kaldi 中的数据准备

### 3.4.1 音频、文本、说话人准备

准备词典 词典的作用是将词映射成音素, 不论是基于 HMM-GMM 的对齐, 还是 HMM-DNN 的声学模型训练, 其建模单元为音素, 但是在实际交流中, 我们看到的都是词或者词组成的句子, 因此我们需要这么一个文件, 在对齐或者训练声学的时候, 将词映射成发音; 在解码的时候我们看到的是词而不是音素, 因此需要训练语言模型将音素转换成最有可能的词。因此, 词典是非常非常重要的。

首先 lexicon.txt 的格式是 "<word> <pronunciation>". 对应的 "<pronunciation>" 就是词的音素。表 3.1 分别展示了中英文的 lexicon.txt 的示例。

在准备词典的时候, 其对应的路径为 **data/local/lm**, 其文本内容类似于图 3.1。词和发音之间以 Tab 隔开, 音素之间以空格隔开。



表 3.1: 中英文 word 和对应 pronunciation 示例

Word	Pronouciation
七叶树	qi iH iH yi ehH ehL shu uH uL
episode	EH1 P IH0 S OW2 D

```
一八三零年 yi iH iL b aH aH s aH nnH li eLL ngH ni aL nnH
一八九七年 yi iH iL b aH aH ji oL uL qi iH iH ni aL nnH
一八九九年 yi iH iL b aH aH ji oL uL ji oL uL ni aL nnH
一八九八年 yi iH iL b aH aH ji oL uL b aH aH ni aL nnH
一八五七年 yi iH iL b aH aH wu uL uL qi iH iH ni aL nnH
一八五九年 yi iH iL b aH aH wu uL uL ji oL uL ni aL nnH
一八五二年 yi iH iL b aH aH wu uL uL ge erH erL ni aL nnH
一八五八年 yi iH iL b aH aH wu uL uL b aH aH ni aL nnH
一八八七年 yi iH iL b aH aH b aH aH qi iH iH ni aL nnH
一八四一年 yi iH iL b aH aH s ifH ifL yi iH iL ni aL nnH
一八四八年 yi iH iL b aH aH s ifH ifL b aH aH ni aL nnH
一八四零年 yi iH iL b aH aH s ifH ifL li eLL ngH ni aL nnH
一六二九年 yi iL iH li oH uL ge erH erL ji oL uL ni aL nnH
一六四零年 yi iL iH li oH uL s ifH ifL li eLL ngH ni aL nnH
```

图 3.1: lexicon.txt 中存储格式

另外一个很重要的文件是 vocab.txt，其路径为 **data/local/lm/vocab.txt**。其存储了待建模的数据库中包含的所有的词。如果是中文的话，就需要对中文语句进行分词，如果对应的数据库已经分好了词，那么

3.5 kaldi 中的决策树



## 第 4 章 FFmpeg 和 sox

### 4.1 FFmpeg

#### 4.1.1 安装 FFmpeg

在 Centos 中利用 yum 安装 FFmpeg 步骤如下：

##### 1. 升级系统

```
1 sudo yum install epel-release -y
2 sudo yum update -y
3 sudo shutdown -r now
```

##### 2. 安装 Nux Dextop Yum 源

由于 CentOS 没有官方 FFmpeg rpm 软件包。但是，我们可以使用第三方 YUM 源（Nux Dextop）完成此工作。

###### • CentOS 7

```
1 sudo rpm --import http://li.nux.ro/download/nux/RPM-GPG-KEY-nux.
  ro
2 sudo rpm -Uvh http://li.nux.ro/download/nux/dextop/el7/x86_64/nux-
  dextop-release-0-5.el7.nux.noarch.rpm
```

###### • CentOS 6

```
1 sudo rpm --import http://li.nux.ro/download/nux/RPM-GPG-KEY-nux.
  ro
2 sudo rpm -Uvh http://li.nux.ro/download/nux/dextop/el6/x86_64/nux-
  dextop-release-0-2.el6.nux.noarch.rpm
```

##### 3. 安装 FFmpeg 和 FFmpeg 开发包

```
1 sudo yum install ffmpeg ffmpeg-devel -y
```

##### 4. 测试是否安装成功

```
1 ffmpeg
```

备注：

(1) 查看机器的 centos 版本：cat /etc/redhat-release

(2) 在执行安装 FFMpeg 和 FFmpeg 包的时候，可能会出现一些错误，因为有些依赖包没有安装，看下未安装的包有哪些，然后去pkgs官网上去搜索下载对应的版本即可。

(3) FFMpeg的使用参考资料

- reach296的博客;
- feixiao 的 GitHub 库;
- ffprobe,ffplay ffmpeg 常用的命令行命令

## 4.2 sox

sox 的参考资料:

- SoX 一音频处理工具里的瑞士军刀

## 第 5 章 Linux 相关笔记

### 5.1 linux 备忘录

以下零零散散的一些笔记用于记录不太熟的 Linux 指令，不断扩充中……

1. 永久修改 ip 地址。首先 ifconfig 找到对应的网卡,其次编辑 vi /etc/sysconfig/network-scripts/ifcfg-lo, 此处假设网卡是 lo。
2. 当我们输入 history 的时候, 会出现历史命令。这些命令从 1 开始到 history 这个指令的索引数, 可以用 !+index 的方式来调用。比如 history 中显示的第五条命令是 ls, 那么当我们在终端输入 !5 的时候会自动调用 ls 指令。而且 ! 也可以接指令的部分内容, 其会自动找寻最近的一条与该内容相关的指令并执行。比如说历史中有两条 ls 指令。一条是 ls /; 一条是 ls /home, 假设第二条是最近的指令, 那么我们执行 ! ls, 那么系统就会执行 ls /home 这条指令。
3. alias short\_cmd='real\_cmd', 使用这个指令可以将 real\_cmd 用 short\_cmd 代替, 减少常用命令的输出时间。使用 unalias short\_cmd 可以取消对应的映射。alias 单独运行可以查看当前有哪些 short\_cmd。而且我们可以将这个指令直接添加到 ~/.bashrc 中, 这样就万事大吉, 省心了, 重启的时候也不会消失。
4. > » 2> 和 2> 还有 > bash run.sh 1»result 2>1。正确错误的都会传到 result 里面。
5. free -m 以 M 为单位显示磁盘存储空间
6. 修改 Ubuntu 下载源可以修改 /etc/apt/sources.list 文件
7. 查看 Ubuntu 版本可以使用 cat /etc/issue
8. Ubuntu 下解决 ifconfig command not found 的办法: sudo apt-get install net-tools
9. 当拥有多个用户, 想把某个文件或者文件夹的权限下放到某个用户或者限制某个用户对某个文件或者文件夹的访问的时候, 可以使用 setfacl 这个指令。具体的下放权限指令操作为: setfacl -m u:user1:rw test.txt, 清空对于某个文件或者文件夹设置的权限时指令操作为: setfacl -b test.txt。查看某个文件或者文件夹更细化的权限信息可以使用 getfacl 指令。对目录以及子目录设置 acl 权限时, 指令为: setfacl -m u:user1:rw -R /mnt。如果后续当前目录有生成新的子目录或者文件, 想要继承现有权限的时候, 需要在执行上述指令之后, 再执行一次 setfacl -m d:u:user1:rw -R /mnt
10. 设置用户对于某个命令的执行权限, 使用指令 visudo (需要在 root 下执行), 举例: 首先执行 visudo, 会打开一个文件, 要添加某个用户对于某个指令的权限时, 需要给出对应命令的绝对路径, 假设给予 user4 添加新用户的权限, 则在 visudo 打开的文件添加一句 user4 localhost=/usr/sbin/useradd, 之后使用 sudo /usr/sbin/useradd user5 之后再输入 user4 的密码,

就可以创建 user5。多个命令使用 “,” 隔开。

11. 分配无密码的 sudo 命令: 同样使用 visudo, 然后在文本中添加一句: user4 ALL=NOPASSWD: /usr/sbin/useradd, /usr/sbin/userdel, 之后再调用 sudo /usr/sbin/useradd user5 的时候就不需要 user4 的密码了。也可以使用 user4 localhost=NOPASSWD: /usr/sbin/useradd, /usr/sbin/userdel, 这就意味着只有在 user4 下的时候才可以不用输入密码, 而 ALL 表示在所有用户下使用 sudo 都不用输入密码。
12. 任务计划 crontab -e 创建一个新的任务计划。30 17 \* \* 5 task, 其中前面的表示周五的 17:30, 在这个时间点执行任务 task。可以用 crontab -l 查看任务计划的内容。
13. ls -R 递归式的显示当前目录所有的东西, 包括文件, 子目录, 子目录中的所有东西
14. 创建一个 ftp 服务站代码如下, 这个时候在 /var 目录下会生成一个 ftp 的目录, ftp 下还有个 pub 目录, 这个目录, 使得我们可以在其他操作系统访问。在 Windows 下, 打开文件夹, 在文件栏输入 ftp://your\_ip 即可访问 ftp 的 pub 目录, ftp 启动指令为 service vsftpd start。

```
1 yum -y install vsftpd*
```

15. 查看 linux 某个端口是否被占用指令: netstat -an|grep \$port
16. 当我们不需要显示输出结果, 但是 Linux 指令默认自动输出指令的时候, 我们可以在指令后面加上 &>/dev/null, 这样所有输出就重定向到一个黑洞里, 全部都消失了。
17. 解决 Linux 中文乱码问题: 以 docker 安装的 Ubuntu 为例。

```
1 >>locale -a
2 C
3 C.UTF-8
4 POSIX
5 >>export LC_ALL='C.UTF-8' #这是临时修改的
6 >>export LC_ALL='C.UTF-8' >> /etc/bash.bashrc | source /etc/bash.
   bashrc #这是永久修改。
7 >> docker restart <container-ip> #退出container, 重启container即可
```

18. docker 启动命令: systemctl start docker
19. bad interpreter: No such file or directory: 出现这个错误的原因是我用 winscp 连接服务器, 然后用 Windows 系统下的 VScode 写代码。因为 Windows 下文本默认的格式是 DOS, 而 Linux 是 unix, 所以在服务器上运行的时候会出现这个错误。修改的办法有两个: 第一个是在 Linux 里面用 vim 打开这个脚本, 输入 set ff 的时候会出现个 DOS, 所以我们要设置下属性, 如下所示。另外我们还可以在 VScode 中直接修改, 写好脚本之后, 点击右下角的 CRLF, 上边栏会出现一个选项, 选择 LF 即可, 如图5.1。

```
1 :set ff=unix
```

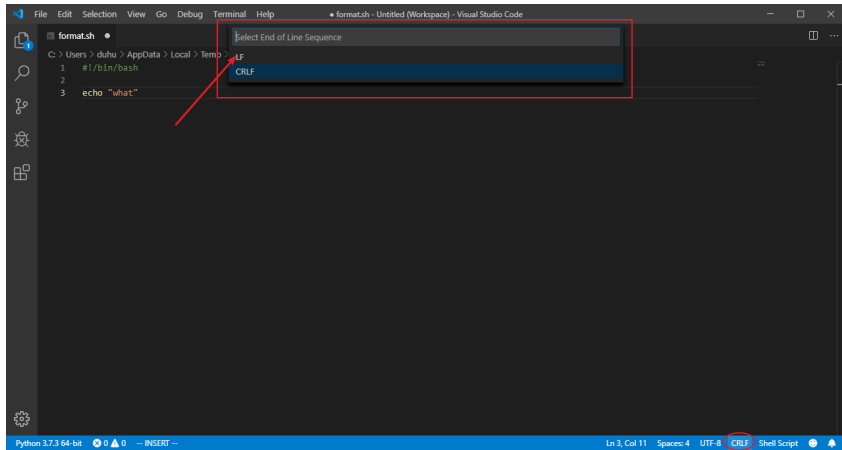


图 5.1: VScode 设置脚本格式为 Unix

20. 我们在 Linux 上安装软件时，经常需要选择机器所对应的版本，即机器是多少位的。查看指令有：

- `uname -a`
- `file /bin/l`
- `cat /proc/vision`
- `getconf LONG_BIT`

21. 删除文本中的空行：`grep -v '' file`

22. 如何对文本指定列进行排序：

```
1 sort -n -k2 file #numerical sort;k2表示的是某一列，从小到大排序
2 sort -g -k2 file #general-numerical sort
3 sort -nr -k2 file #从大到小排序
```

23. 出现错误：“[: too many arguments”：这个是在字符串比较中出现的，表示的是比较的对象数量不一致，也就是说可能出现多个字符串，这个时候给这个变量添加上双引号就可以了。

## 5.2 Shell 指令笔记

### 5.2.1 简单的 shell 规则

此处记载一些简单的 shell 指令和编程规则。不断补充……

1. 将键盘输入的内容赋值给变量：`read [-p "some message"] 变量名。`

2. 双引号 (") 允许通过 \$ 符号引用其他变量值；单引号 (') 禁止引用其他变量值，\$ 视为普通字符；反撇号 (`) 将命令执行的结构输出给变量。
3. 关于 shell 中外部传输变量的一些操作：
  - \$#: 命令行中位置参数的个数；
  - \$\*: 所有位置参数的内容，当格式为 "\$\*" 时，参数作为一个整体传递给程序；
  - \$? : 上一条命令执行后返回的状态，当返回状态值为 0 时，表示执行正常；非 0 则表示执行异常；
  - \$0: 当前执行的进程/程序名；
  - \$n:  $n \in \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$ ，表示外部传入的第 n 个参数。
  - \$@: 传递所有参数的内容，当格式为 "\$@" 的时候，表示将参数分开传递给程序；
4. shell 中的数学运算指令: `expr`。注意乘法是 `\*`
5. 解析字符串中的转义字符，可以使用 `echo -e`，比如 `echo -e "test\ntest"`，`echo -e` 还可以修改输出的颜色和背景色，指令是 `\[033[前景颜色; 背景颜色 m`。`\[033[0m` 表示恢复到系统默认颜色。其前景颜色的数值为：默认 =0，黑色 =30，红色 =31，绿色 =32，黄色 =33，蓝色 =34，紫色 =35，天蓝色 =36，白色 =3；背景颜色的数值为：默认 =0，黑色 =30，红色 =31，绿色 =32，黄色 =33，蓝色 =34，紫色 =35，天蓝色 =36，白色 =3。
6. shell 默认在输出之后换行的，如果想要不换行，可以选择参数 `-n`。例如 `echo -n "Enter your name: "`，这样在终端显示的时候就不会换行。只有一个 `echo` 的时候，会输出一个换行。
7. `cat` 的另外一种用法，比如可以用于制作菜单，如下程序所示，这个程序在执行的时候会保留原样输出两个 x 之间的内容，包括换行和空格等。

```
1 cat<<x
2   please input your name:
3       1) user1
4       2) user2
5       3) user3
6 x
```

8. `nl` 指令：给输出标上行，比如 `cat test.txt | nl`，这条指令会给 `test.txt` 中的每一行进行顺序标号；
9. `tee` 指令：在程序执行输出的时候额外保留一份输出到文件中，比如 `./test.sh | tee test.txt`。能够起到一个备份的作用。
10. `shift` 指令：用于迁移位置变量，将 \$1 到 \$9 依次向左传递。
11. 两个文件的交集、并集和去重：

```
1 cat file1 file2 | sort | uniq > result #求两个文件的并集，如果有重复的行
```

```
只保留一行。
2 cat file1 file2 | sort | uniq -d > result #求两个文件的交集，即两个文件
  中都有有的行。
3 cat file1 file2 | sort | uniq -u > result #求两个文件的差集，即只有一个
  文件中有的行。
4 cat file1 file2 > result #追击的方式，如果file1有n行，file2有m行，result
  为n+m行
5 paste file1 file2 > result #一个文件的内容在左边，一个文件命令在右边。
6 sort file |uniq #将重复的多行变为一行
7 sort file |uniq -u #sort file |uniq -u
```

5.2.2 shell 中的条件测试操作

1. test 命令:

- (a). 用途: 测试特定的表达式是否成立，成立返回 0，不成立则返回非 0;
- (b). 格式: test 条件表达式 [ 条件表达式]

2. 常见的测试类型

- (a). 测试文件状态:
  - 格式: [ 操作符文件或者目录]
  - 常用的文件操作符，如表5.1。

表 5.1: 常用的文件操作符及其意义

文件操作符	作用
-d	是否为目录 (Directory)
-f	是否为文件 (File)
-e	是否存在文件或者目录 (Exist)
-r	当前用户是否可读 (Read)
-w	当前用户是否可写 (Write)
-x	当前用户是否可操作 (Excute)
-L	文件是否为符号链接文件 (Link)

(b). 字符串比较:

- 格式: [ 字符串 1 = 字符串 2 ], [ 字符串 1 != 字符串 2 ], [ -z 字符串]
- 常用的字符串操作符，如表5.2。

(c). 整数值比较:

- 格式: [ 整数 1 操作符整数 2 ]
- 常用的整数值比较操作符，如表5.3。





表 5.2: 常用的字符串操作符及其意义

字符串操作符	作用
=	字符串内容相同
!=	字符串内容不相同
-z	字符串内容为空

表 5.3: 常用的整数操作符及其意义

整数操作符	作用
-eq	等于 (Equal)
-ne	不等于 (Not Equal)
-gt	大于 (Greater Than)
-lt	小于 (Less Than)
-ge	大于等于 (Greater or Equal)
-le	小于等于 (Less or Equal)

- (d). 逻辑测试: 此处需要注意的是与操作和或操作有的时候可以起到一个开关的作用, 比如 A&&B, 假设 B 是一条指令的话, 那么只有 A 为真的情况下才会进行下一步操作; 类似的, 如果是或操作, 只有前面为假才会执行后面的操作。
- 格式: [ 表达式 1 ] 操作符 [ 表达式 2 ] ...
  - 常用的逻辑操作符操作, 如表5.4。

表 5.4: 常用的逻辑操作符及其意义

逻辑操作符	作用
-a/∧∧	逻辑与
-o/∥	逻辑或
-!	逻辑否

5.2.3 find

基本使用 find <directory> -args contents, 比如说 find . -name fun。find 的指令和例子如下代码5.1所示。

Listing 5.1: find 指令详解

```
1 find . -name "[a-z]*" #寻找所有a-z开头的文件或者文件夹
2 find . -name "[0-9]*" #寻找所有0-9开头的文件或者文件夹
3 find . -perm 775 #根据权限寻找文件
4 find . -user root #根据文件创建者来寻找文件
5 find . -mtime -5 #寻找更改时间为5天以内的文件
```



```
6 find . -mtime +3 #寻找更改时间为3天以前的文件
7 find . -type d[f;I] #寻找类型为文件夹[文件；链接]的文件
8 find . -size +1000000c #寻找文件大小大于1M的文件
9 find . -perm 700 | xargs chmod 777 #找到权限为700的改为777
10 find . -type f | xargs ls -l #找到所有文件并展示出来
11 find . \( -name *.gt -o -name *.pre\) #同时找到不同类型的多个文件，每添加一个
    类型，都需要加上 -o -name
```

### 5.2.4 grep

本小节分为两块一个是 grep 指令的一般性操作，如代码5.2。

Listing 5.2: grep 指令的一般性操作

```
1 grep "<contents>" #查找内容
2 grep -c "<contents>" file #文件中有多少行匹配到查找的内容
3 grep -n "<contents>" file #文件中有多少行匹配到找到的文件，并显示行号
4 grep -i "<contents>" file #查找内容，并忽略大小写
5 grep -v "<contents>" file #过滤掉内容
```

另外一部分就是配合正则表达式进行文本的查找，先说下一些正则表达式的表示，如下所示：

1. '^linux': 以 linux 开头的行；
2. 'linux\$': 以 linux 结尾的行；
3. '.': 匹配任意单字符；
4. '.+': 匹配任意多个字符；
5. '\*': 匹配 0 个或多个字符；
6. '[0-9a-z]': 匹配 [] 内任意一个字符；
7. '(linux)+': 出现多次 linux 单词；
8. '(linux){2}': 出现两次 linux；
9. '\\': 转义字符；
10. '\$': 空行。
11. '[^linux]': 查找不以 linux 开头的行，[] 内的 '^' 表示取反。

### 5.2.5 awk

awk 是非常强的一个文本编辑的命令，主要按照列来对文本进行分割处理。一般指令为：awk -F: '{print \$1}'，这段指令的意思是以':' 为分隔符，输出第一列的数据。awk 不太好总结，因为用到的还不是很多，因此，以逐条总结的方式，将一些不太常用，但是关键时刻很管用的记录下来<sup>(20)</sup>，等以后更熟悉了再做全面的总结。代码如5.3。

Listing 5.3: awk 指令的一般性操作

```

1 >>cat log.txt
2 2 this is a test
3 3 Are you like awk
4 This's a test
5 10 There are orange,apple,mongo
6
7 >>awk -v #设置变量
8 >>awk -va=1 '{print $1,$1+a}' log.txt #此时a就是1，所以如果$1是数字，那么数字
   加一，如果不是，那么输出就是1
9 >>awk -vb=s '{print $1,$1+a}' log.txt #b为s，所以在第一列的所有元素后面加个s
10 >>awk '$1>2' log.txt #输出第一列中大于2的行
11 >>awk '$1==2 {print $1,$3}' log.txt #输出第一列中等于2的行的第一列和第三列
12 >>awk '$1>2 && $2="Are" {print $1, $2, $3}' log.txt #输出第一列大于2的行的第
   一、二和三列

```

awk 用于求和、求最大值、求最小值和求平均值，见代码5.4。

Listing 5.4: awk 求和、求平均、求最大最小值

```

1 cat data|awk '{sum+=$1} END {print "Sum = ", sum}' #求和
2 cat data|awk '{sum+=$1} END {print "Average = ", sum/NR}' #求平均值
3 cat data|awk 'BEGIN {max = 0} {if ($1>max) max=$1 fi} END {print "Max=",
   max}' #求最大值
4 awk 'BEGIN {min = 1999999} {if ($1<min) min=$1 fi} END {print "Min=", min
   }' #求最小值

```

### 5.2.6 sed

sed 与 awk 相辅相成，其主要对行进行操作，主要细节参考菜鸟教程<sup>(21)</sup>。

## 第 6 章 Windows 相关

本章用于记录一些 Windows 使用的一些问题，方便再遇到的时候查询。

### 6.1 win10 .net framework 3.5 安装报错 0x800F0954 问题

本问题解决参考[tOneDay](#)的博客，其他的一些博客都没有解决问题。以此为准：

1. 打开注册表：cmd+r 输入 regedit，确定；
2. 找到路径 HKEY\_LOCAL\_MACHINE-SOFTWARE-Policies-Microsoft-Windows-WindowsUpdate-AU，其中 UseWUSever 默认值为 1，改成 0；
3. 打开服务列表，重启 Windows Update service；
4. 此时可以正常安装.net framework 3.5；
5. 将第二步的修改还原，并重启 Windows Update service。

## 第 7 章 Python 笔记

### 7.1 一些小技巧

此处记录一些常用到小技巧，省时省力省心还漂亮……不断补充中……

#### 1. 按照 Value 对字典进行排序

```
1 xs = {'a':4, 'b':3, 'c':2, 'd':1}
2 sorted(xs.items(), key=lambda x:x[1])
3 import operator
4 sorted(xs.items(), key=operator.itemgetter(1))
```

#### 2. 假设 B 是列表 A 的子集，想要求 B 的补集，代码如下：

```
1 x = [1,2,3,4,5,6,7]
2 y = [1,2,3]
3 z = list(set(x+y))
```

#### 3. 假设列表 B 和列表 A 有重复元素，目的去除重复元素，去除的代码为函数 `remove_same(A, B)`。列表长度的判断是为了减少循环次数，为了测试时间写了一个简单的脚本，结果显示循环短列表大概块 1 秒钟，这个取决于短列表有多短，只是聊胜于无的减少了些代码运行时间。

```
1 from time import time
2 def remove_same(A, B):
3     a, b = A.copy(), B.copy()
4     for i in A:
5         if i in B:
6             a.remove(i)
7             b.remove(i)
8     return a, b
9 A = []
10 B = []
11 for i in range(100000):
12     if i%2 == 0:
13         A.append(i)
14     if i%33 == 0:
```

```

15         B.append(i)
16 print("lenA:{}\t\tlenB:{}".format(len(A), len(B)))
17 st = time()
18 b, a = remove_same(B, A)
19 mt = time()
20 a, b = remove_same(A, B)
21 et = time()
22 if len(A) > len(B):
23     b, a = remove_same(B, A)
24 else:
25     a, b = remove_same(A, B)
26 et2 = time()
27 print("CycleB:", mt-st)
28 print("CycleA:", et-mt)
29 print("CycleSmallerOne:", et2-et)

```

4. TextGrid 格式的文本处理参考[python 的 textgrid 库调研小结](#)

5. 进度条效果:

```

1 sys.stdout.write('.')
2 sys.stdout.flush()

```

## 7.2 python 中的线程、进程、协程与并行、并发

我们先介绍下设计到这几个东西的概念吧。

1. GIL: 在 Cpython 解释器中, 同一个进程下的多个线程, 同一时刻只能有一个线程执行, 无法利用多核优势。GIL 本质就是一把互斥锁, 即将并发运行变成串行, 以此来控制同一时间内共享数据只能被一个任务进行修改, 从而保证数据的安全性保护不同的数据时, 应该加不同的锁, GIL 是解释器级别的锁, 又叫做全局解释器锁 CPython 加入 GIL 主要的原因是为了降低程序的开发复杂度, 让你不需要关心内存回收的问题, 你可以理解为 Python 解释器里有一个独立的线程, 每过一段时间它起 wake up 做一次全局轮询看看哪些内存数据是可以被清空的, 此时你自己的程序里的线程和 Python 解释器自己的线程是并发运行的, 假设你的线程删除了一个变量, py 解释器的垃圾回收线程在清空这个变量的过程中的 clearing 时刻, 可能一个其它线程正好又重新给这个还没来得及得清空的内存空间赋值了, 结果就有可能新赋值的数据被删除了, 为了解决类似的问题, Python 解释器简单粗暴的加了

锁, 即当一个线程运行时, 其它人都不能动, 这样就解决了上述的问题, 这可以说是 Python 早期版本的遗留问题. 毕竟 Python 出来的时候, 多核处理还没出来呢, 所以并没有考虑多核问题, 以上就可以说明, Python 多线程不适合 CPU 密集型应用, 但适用于 IO 密集型应用

'''

In CPython, the global interpreter lock, or GIL, is a mutex that prevents multiple native threads from executing Python bytecodes at once. This lock is necessary mainly because CPython's memory management is not thread-safe. (However, since the GIL exists, other features have grown to depend on the guarantees that it enforces.)

'''

2. 进程
3. 线程
4. 协程
5. 并行
6. 并发

## 7.2.1 进程和线程

进程和线程操作系统中的概念, 这也是操作系统中的核心概念。

### 7.2.1.1 进程

进程是对正在运行程序的一个抽象, 即一个进程就是一个正在执行程序实例。从概念上说每个进程拥有它自己的虚拟 CPU, 当然, 实际上是真正的 CPU 在各个进程之间来回切换, 这种快速切换就是多道程序设计, 但是某一瞬间, 一个 CPU 只能运行一个进程, 但是在 1 秒钟期间, 它可能运行多个进程, 就是 CPU 在进行快速的切换, 有时人们所说的伪并行就是指这种情况。

#### 1. 创建进程

操作系统中有四种事件会导致进程的创建:

- 系统初始化, 启动操作系统时, 通常会创建若干个进程, 分为前台进程和后台进程;
- 执行了正在运行的进程所调用的进程创建系统调用;
- 用户请求创建一个新的进程;
- 一个批处理作业的初始化。

从技术上来看, 在所有这些情况中, 新进程都是由一个已经存在的进程执行了一个用于创建进程的系统调用而创建的。这个进程可以是一个运行的用户过程, 一个由键盘或者鼠标启动的系统进程或者一个批处理管理进程。这个进程所做的工作是执行一个用来创建新进



程的系统调用。在 Linux/Unix 系统中提供了一个 `fork()`，用来创建进程的子进程，在 python 的 `os` 模块中封装了常见的系统调用。代码如下：

```
1 import os
2 # os.getpid()获取父进程的ID
3 print("Process%sstart..." % os.getpid())
4 # fock()调用一次会返回两次
5 pid = os.fork()
6 # 子进程返回0
7 if pid == 0:
8     print("Iamchildprocess%sandmyparentis%s"%(os.getpid(), os.
9         getppid()))
9 # 父进程返回子进程的ID
10 else:
11     print("I%sjustcreatedachildprocess%s"%(os.getpid(), pid))
```

## 7.3 客户端向服务端传送一个音频文件及信息

假定我们有一个客户端程序，一个服务端程序，要求从客户端发送一个音频到服务端，包括音频的名字、时长、采样率、采样宽度和通道数。应该怎么去做？这里面主要有四个库：`socketserver`，用于服务端部署，`socket` 用于客户端发送数据，`struct` 用于将音频的额外信息打包，`wave` 用于读取客户端音频及其信息，生成二进制采样点，等音频的这些数据发送到服务端的时候，再通过传送过来的音频数据和音频信息生成完全相同的音频。

之所以要这么一个奇怪的需求，是因为我们可以在服务端部署一个在线语音识别引擎，这样服务端的引擎一直在运行，当接收到客户端传送过来的数据的时候就开始在客户端也生成一个同样的音频，再调用识别模块对音频进行解析，最终生成文本再传送回客户端。这样我们就可以完成在线语音识别的任务了。那么首先遇到的一个问题就是……

……………怎么传文件……………

首先我们挨个库介绍下吧……

### 7.3.1 wave

`wave` 是 python 中专门用于读取音频信息的一个库，可读可写，都是二进制的操作。音频有一些重要的信息包括采样率，采样宽度，时长和通道数，以及音频各个采样点的值都可以通过 `wave` 获得。通过下面的代码，我们就可以完成一个 `wave` 读取一个音频，再重新创建一个音频



文件，写入读取的音频信息生成一个完全相同的音频的过程。有点类似于复制的感觉……不解释太多，一切尽在代码7.1中。

**Listing 7.1:** wave 库读取和写入音频

```
1 import sys
2 import wave
3 def read_wav(audio_path)
4     f = wave.open(audio_path, 'rb')
5     nchannels, samplewidth, framerate, nframes = f.getparams()[:4]
6     audio_contents = f.readframes(nframes) #f.readframes(n)里面的n是帧数，
7                                           #通过这个参数我们可以对音频进行裁剪
8     f.close()
9     return audio_contents, nchannels, samplewidth, framerate, nframes
10 def write_wav(audio_contents, nchannels, samplewidth, framerate,
11               audio_path)
12     f = wave.open(audio_path, 'wb')
13     f.setnchannels(nchannels)
14     f.setsampwidth(samplewidth)
15     f.setframerate(framerate)
16     f.writeframes(audio_contents)
17     f.close()
18 audio_path, copied_audio = sys.argv[1:] #外部给原始地址和存储地址
19 audio_contents, nchannels, samplewidth, framerate, nframes = read_wav(
20     audio_path)
21 write_wav(audio_contents, nchannels, samplewidth, framerate, copied_audio)
```

### 7.3.2 struct

**struct**是用来处理二进制数据的。有三个函数是比较重要的：**pack**、**unpack** 和 **calsize**。当我们调用 **pack** 这个函数的时候，格式是 **struct.pack("<fmt>", data)**，其中 **<fmt>** 指的是打包数据的格式，因为不同的格式在计算机中占据的存储空间不同，因此需要指定下，同样，当我们在调用 **unpack** 这个函数的时候，格式是 **struct.unpack("<fmt>", data)**，其格式跟 **pack** 函数差不多，因为解包的时候也一样得知道这个压缩的数据包中都是什么样的数据，这样可以根据这个 **<fmt>** 得到原始的数据，其返回的是一个 **tuple**。而 **calsize** 这个函数就是用来计算如果以格式"**<fmt>**" 打包或者解包所需要的存储空间，其格式是 **struct.calsize("<fmt>")**。而不同类型的数据占据的字节

数不同，表7.1列出了一般的数据类型、其表示和占据字节数。

表 7.1: struct 中常用的数据类型、C 语言的对应类型和占用字节数

Format	C Type	Python	字节数
x	pad byte	no value	1
c	char	string of length 1	1
b	signed char	integer	1
B	unsigned char	integer	1
?	<i>bool</i>	bool	1
h	short	integer	2
H	unsigned short	integer	2
i	int	integer	4
I	unsigned int	integer or long	4
l	long	integer	4
L	unsigned long	long	4
q	long long	long	8
Q	unsigned long long	long	8
f	float	float	4
d	double	float	8
s	char[]	string	1
p	char[]	string	1
P	void *	long	unclear

这里面还有一些补充的信息如下：

- 每个格式前可以有一个数字，表示个数，比如"5i" 表示有五个整型数据，则占用 20 个字节；
- s格式表示一定长度的字符串，"4s" 表示长度为 4 的字符串，而 p 表示的是 pascal 字符串；
- q 和 Q 只在机器 64 位操作时有意义；
- P 用来转换一个指针，其长度和机器字长相关；

为了同 C 中的结构体交换数据，还要考虑有的 c 或者 C++ 编译器使用了字节对齐，通常是以 4 个字节为单位的 32 位系统，故而 struct 根据本地机器字节顺序转换，可以用格式中的第一个字符来改变对齐方式。这个字符的类型和定义如表7.2。

表 7.2: struct 中的字节对齐操作符号及其含义

Character	Bytes Order	Size and alignment
@	native	native 凑够 4 个字节
=	native	standard 按原字节数
<	little-endian	standard 按原字节数
>	big-endian	standard 按原字节数
!	network(=big-endian)	standard 按原字节数

在讲到 struct 在我们这个需求中应用之前，我们先看一些小例子，将不同类型的数据打包



起来，再按照原始格式给解包。需要注意的是字符串必须先转换成二进制，先编码再转换；解包之后也需要进行解码才可以得到原始的字符串。一切尽在代码7.2中……

**Listing 7.2:** struct 打包和解包不同类型的数据

```
1 import struct
2 a,b,c,d,e = 1,2,3,'this', 'good'
3 d,e = bytes(d.encode('utf-8')), bytes(e.encode('utf-8'))
4 y = struct.pack("3i4s4s", a,b,c,d,e)
5 p,q,r,s,t = struct.unpack("3i4s4s", y)
6 s,t = s.decode('utf-8'), t.decode('utf-8')
```

如果我们需要通过 socket 传输一个音频，并且希望可以在服务器端可以完全重建音频，那么我们需要打包的音频信息有哪些呢？

- 音频长度：因为 socket 传输的时候，是根据服务端来确定接收多少个字节的，一个音频比较长，需要分成很多段去接收，那么什么时候算接收完了呢？就需要这个音频长度去确定了。
- 音频的基本信息：采样率，采样宽度，通道数。在服务器端重建时，wave 这个库需要用到这些信息；
- 音频的名字：我们需要在服务器端生成一个完全一样的同名 wav 文件，那么文件名是必须传输的。

那么我们在客户端打包的时候格式很好确定，音频长度、采样率、采样宽度和通道数共四个整数；那我们需要确定文件名的长度啊，不然服务端怎么解包，所以再加一个整数：文件名的长度；最后还有文件名。那么如果想要在服务器端直接利用这些信息的格式来解包，那么我们就需要将这个格式传输过去，因为这个格式还是字符串，所以我们还需要传输一个格式的长度，这个是整数。因此这个包就分为了两块：第一块是存储格式的长度和对应的格式；第二块是存储上面说的音频信息。所以结合上面7.1，我们可以这样处理，代码7.3中仅写了打包发送音频信息的部分，其他部分见后面 socket。

**Listing 7.3:** struct 打包音频信息和数据

```
1 import os
2 import struct
3 def encode(str):
4     return bytes(str.encode('utf-8'))
5 cons, nchan, sampwid, frate, nfr = read_wav(audio_path)
6 filename = os.path.basename(audio_path)
7 #-----client-----
```

```

8  fmt = ">4i%is"%(len(filename))
9  info = struct.pack(">i%is4i%is"%(len(fmt), len(filename)), \
10                      len(fmt), encode(fmt), \
11                      len(cons), nchan, sampwid, \
12                      frate, encode(filename))
13  #-----server-----
14
15  fmt_len = struct.unpack(">i", info[:4])
16  fmt = struct.unpack("%i%is"%fmt_len, info[4:fmt_len])
17  fmt_size = struct.calcsize(fmt)
18  info = info[(4+fmt_len):]
19  conlen, nchan, sampwid, frate, fname = \
20      struct.unpack(fmt, info[:fmt_size])

```

### 7.3.3 socket

**socket**是个用来进行网络传输的库，爬虫的时候经常能看见这玩意，咱们介绍下 **socket** 常用的一些操作，并举一些例子，需要注意的是本需求中，我们只在客户端用 **socket** 库。

**socket** 常用的函数如下：

1. **socket(socket.AF\_INET, socket.SOCK\_STREAM)**: 创建了一个 **socket** 连接的实例，括号内的参数可变，一般常用的就是这两个，其他的不甚了解，以后再补上；
2. **bind((host, port))**: 绑定 **ip** 和端口，**host** 是要连接的服务器端的 **ip** 地址，**port** 是服务器端的端口；
3. **connect((host, port))**: 客户端连接服务端的 **ip** 和端口；
4. **listen(n)**: 服务器端开始监听，其中 **n** 表示监听的队列的个数；
5. **accept()**: 接收客户端传来的数据，返回两个值：(**conn**, **address**)，**conn** 是新的套接字对象，用来接收数据和返回数据，**address** 是客户端的地址和 **ip**，类型是 **tuple**，**conn** 是和 **address** 绑定的；
6. **recv(byte)**: 接收数据，其中的 **byte** 表示一次接收多少个字节；
7. **send(string)**: 发送二进制字符串，并返回发送的字节大小，就发送一次。这个字节长度可能是小于实际要发送的字节的长度的，假设要发送的字节数是 1025，服务端一次接收 1024 个字节，那么最后那个字节就丢了……所以如果用这个函数的话，可能就需要写一个多次发送的循环；
8. **sendall(string)**: 发送二进制字符串，发送成功返回 **None**，失败则出错。这个就是整个数据

都发送，发完为止：

基本的函数说完了，咱们就来分别写一个服务端和客户端的小 demo。见代码7.4和7.5，这里面用 struct 打包了要发送的数据的长度，为避免数据传输不完整。

**Listing 7.4:** socket 服务端代码

```
1 import socket
2 import struct
3 def server(host, port):
4     sock = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
5     sock.bind((host, port))
6     sock.listen(5)
7     while True:
8         conn, addr = sock.accept()
9         print("get msg from", addr)
10        data = conn.recv(1024)
11        len_msg = struct.unpack(">i", data[:4])[0]
12        print(len_msg)
13        msg = data[4:]
14        if data:
15            while len(msg) < len_msg:
16                data = conn.recv(1024)
17                msg += data
18            print(msg)
19            conn.sendall(msg)
20        else:
21            continue
22 if __name__ == "__main__":
23     host = 'localhost'
24     port = 8888
25     server(host, port)
```

**Listing 7.5:** socket 客户端代码

```
1 import socket
2 import struct
```

```
3 def client(host, port):
4     sock = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
5     sock.connect((host, port))
6     msg = b"hello"
7     sock.send(struct.pack('>i', len(msg))+msg)
8     data = sock.recv(1024)
9     print(data)
10 if __name__ == "__main__":
11     host = 'localhost'
12     port = 8888
13     client(host, port)
```

### 7.3.4 socketserver

**socketserver**是一个搭建网络服务器的库，基于 **socket** 扩展的。将其用于服务端的搭建更省心一些。其有一些基类，再对这些类中的一些方法进行复写。其定义的类有不少，介绍四个，以后再补充。如下所示：

1. **socketserver.TCPServer**: 负责处理 TCP 协议的类，网络传输数据的时候我们用这个比较多，因为比较稳定，这个是有链接的，客户端和客户端的交流只能是通过服务端来搞定；
2. **socketserver.UDPServer**: 负责处理 UDP 协议的类，这个是没有中心链接的，客户端之间可以直接交流，可以用来写聊天器（存疑 ing）；
3. **socketserver.BaseRequestHandler**: 开发者自定义的处理 request 的类，用来接收客户端的连接和数据传输等；
4. **socketserver.ThreadingMixIn**: 用来处理多线程连接的类。

同样我们是根据一些小 demo 来理解这个库，我们就改写下 **socket** 中的服务器端的代码，见代码7.6。

**Listing 7.6:** socketserver 构建服务器端代码

```
1 import socketserver
2 import struct
3 class myTCPServer(socketserver.ThreadingMixIn, socketserver.TCPServer):
4     def __init__(self, address, handler):
5         socketserver.TCPServer.__init__(self, address, handler)
6
7 class myTCPRequestHandler(socketserver.BaseRequestHandler):
```

```
8     def handle(self):
9         data = self.request.recv(1024)
10        len_msg = struct.unpack(">i", data[:4])[0]
11        msg = data[4:]
12        while len(msg) < len_msg:
13            data = self.request.recv(1024)
14            msg += data
15        print(msg)
16        self.request.sendall(msg)
17    if __name__ == "__main__":
18        host = 'localhost'
19        port = 8888
20        server = myTCPServer((host,port), myTCPRequestHandler)
21        server.serve_forever()
```

### 7.3.5 网络传输音频并保存

okokok, 累死我了。讲到这儿, 我觉得把上面讲的这些串起来, 写一个客户端发送音频, 服务端接收音频并原样保存的代码并不困难了, 那么直接把代码放上来, 不做太多解释了, 见代码7.7和7.8。

**Listing 7.7:** 音频传输的 client 端代码

```
1  import os
2  import sys
3  import wave
4  import struct
5  import socket
6  def callback():
7      if sys.argv[1] is not None:
8          audio_path = sys.argv[1]
9          filename = os.path.basename(audio_path)
10         cons, nchannels, samplewidth, framerate, _ = read_wav(audio_path)
11         sock = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
12         sock.connect((HOST, PORT))
13         # sock.sendall(struct.pack('>i', len(cons))+cons)
```

```

14     fmt = ">4i%is"%len(filename)
15     info = struct.pack('>i%is4i%is'%(len(fmt), len(filename)), \
16                        len(fmt), structcode(fmt), \
17                        len(cons), nchannels, \
18                        samplewidth, framerate, \
19                        structcode(filename))
20     sock.sendall(info+cons)
21     recived = sock.recv(1024)
22     print("FILENAME: ", recived)
23 def read_wav(path):
24     f = wave.open(path)
25     nchannels, samplewidth, framerate, nframes = f.getparams()[ :4]
26     cons = f.readframes(nframes)
27     return cons,nchannels, samplewidth, framerate, nframes
28 def structcode(_str):
29     return bytes(_str.encode('utf-8'))
30 if __name__ == "__main__":
31     HOST = "localhost"
32     PORT = 8888
33     callback()

```

Listing 7.8: 音频传输的 server 端代码

```

1 import sys
2 import wave
3 import struct
4 import socketserver
5 class myTCPServer(socketserver.ThreadingMixIn, socketserver.TCPServer):
6     def __init__(self, address, handlerclass):
7         socketserver.TCPServer.__init__(self, address, handlerclass)
8 class myTCPRequestHandler(socketserver.BaseRequestHandler):
9     def handle(self):
10         chunk = self.request.recv(1024) #一次只能接受1024个字节的数据，其他的会
            继续传过来

```



```

11     len_fmt = struct.unpack('>i', chunk[:4])[0]
12     fmt = struct.unpack('>is'%len_fmt, chunk[4:4+len_fmt])[0].decode('
        utf-8')
13     size_fmt = struct.calcsize(fmt)
14     chunk = chunk[(4+len_fmt):]
15     target_length, nchannels, samplewidth, framerate, filename = \
16         struct.unpack(fmt, chunk[:size_fmt])
17     filename = filename.decode('utf-8')
18     cons = chunk[size_fmt:]
19     while len(cons) < target_length:
20         chunk = self.request.recv(1024)
21         cons += chunk
22     filename = self._write_to_file(cons, str(filename), nchannels,
        samplewidth, framerate)
23     self.request.sendall(filename.encode('utf-8'))
24     def _write_to_file(self, data, filename, nchannels, samplewidth,
        framerate):
25         filename = filename.split('.')[0] + '_' + self.client_address[0] + '
            .wav'
26         file = wave.open(filename, 'wb')
27         file.setnchannels(nchannels)
28         file.setframerate(framerate)
29         file.setsampwidth(samplewidth)
30         file.writeframes(data)
31         file.close()
32         return filename
33 if __name__ == "__main__":
34     HOST = "localhost"
35     PORT = 8888
36     server = myTCPServer((HOST, PORT), myTCPRequestHandler)
37     print('-----')
38     print("Server启动")
39     print('-----')
40     server.serve_forever()

```

## 7.4 Python 中的正则表达式

## 7.5 替代 os 的 pathlib

## 第 8 章 C++ 学习笔记

```
1 #include <iostream>
2 #include <string.h>
3 #include <stdlib.h>
4 using namespace std;
5
6 int main()
7 {
8     printf("Hello World");
9     return 0;
10 }
```

## 第 9 章 Docker

### 9.1 安装 docker 和 nvidia-docker

安装 docker:

```
1 wget -O /etc/yum.repos.d/CentOS-Base.repo http://mirrors.aliyun.com/repo/Centos-7.repo
2 yum install epel-release
3 yum install epel-release
4 yum install docker-ce
```

安装 nvidia-docker

```
1
2 distribution=$(. /etc/os-release;echo $ID$VERSION_ID)
3 curl -s -L https://nvidia.github.io/nvidia-docker/$distribution/nvidia-docker.repo | sudo tee /etc/yum.repos.d/nvidia-docker.repo
4
5 sudo yum install -y nvidia-container-toolkit
6 sudo systemctl restart docker
7 # Install nvidia-docker and nvidia-docker-plugin
8 wget -P /tmp https://github.com/NVIDIA/nvidia-docker/releases/download/v1.0.1/nvidia-docker-1.0.1-1.x86_64.rpm
9 sudo rpm -i /tmp/nvidia-docker*.rpm && rm /tmp/nvidia-docker*.rpm
10 sudo systemctl start nvidia-docker
11
12 # Test nvidia-smi
13 nvidia-docker run --rm nvidia/cuda nvidia-smi
```

### 9.2 常用操作

镜像 (image) 和容器 (container) 的关系, 就像是面向对象程序设计中的类和实例一样, 镜像是静态的定义, 容器是镜像运行时的实体。容器可以被创建、启动、停止、删除和暂停等。

容器的实质是进程，但与直接在宿主机执行的进行不同，容器进程运行于属于自己的独立的命名空间。

```
1 #开启镜像
2 sudo nvidia-docker run -it -v $(pwd)/DeepSpeech:/DeepSpeech -v /data1/asr_
   data:/mnt/data -v //data/kaldi/2019_0521_kaldi/kaldi-master:/mnt/kaldi
   paddlepaddle/deep_speech:latest-gpu /bin/bash
3 # 挂起镜像
4 Ctrl+P+Q
5 #运行已挂起的镜像
6 docker attach $CONTAINER_ID
```

## 9.3 实践要求

docker 如果想要挂载上本地的 IP 地址，可以再运行 docker 的时候加上指令 `-net=host`。如果要挂载本地物理磁盘，加上指令 `-v`。如果要将宿主机的端口与容器的端口绑定，可以使用 `-p < 宿主机端口 >:< 容器端口 >`。

```
1 nvidia-docker run -it --net=host -p 50001:22 -v $(pwd)/DeepSpeech:/
   DeepSpeech -v /data1/asr_data:/mnt/data -v /data/kaldi/2019_0521_kaldi/
   kaldi-master:/mnt/kaldi duhu/ds-server /bin/bash
```

1. 容器不应该向其存储层内写入任何数据，容器存储层要保持无状态化。所有的文件写入操作，都应该使用数据卷（Volume）、或者绑定宿主目录，在这些位置的读写会跳过容器存储层，直接对宿主（或网络存储）发生读写，其性能和稳定性更高。
2. 数据卷的生存周期独立于容器，容器消亡，数据卷不会消亡。因此，使用数据卷后，容器删除或者重新运行之后，数据却不会丢失。

## 9.4 docker 安装 TensorFlow

为了避免影响到主机上诸多配置，因此选用 docker 安装 TensorFlow，想怎么造就怎么造。安装 TensorFlow 时，使用以下指令就会启动该镜像安装。

```
1 docker pull tensorflow/tensorflow
```

## 第 10 章 端到端语音识别汇总

### 10.1 CTC

CTC 说实话，就是比较麻烦……我已经看过很多遍了，不过看的原理居多，代码这块，前后向的实现以及梯度的求取这些都没看过……先总结下 CTC 的基本原理吧还是。

列出本节参考的一些文献和博客：

1. CTC 的开山之作<sup>(2)</sup> 《[Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks](#)》，不建议看这篇……因为这篇在讲到求输出序列的时候使用的前后向算法中前向和后向算法的时候对当前时刻  $t$  的输出概率算了两遍，后面还得再去掉，感觉没啥意思。其次在介绍前后向算法的时候不够简洁，比他的博士论文复杂不少；
2. Graves 大佬的博士毕业论文<sup>(1)</sup> 《[Supervised Sequence Labelling with Recurrent Neural Networks](#)》里面的推导过程写的就相对来说更清楚一些，过程非常简洁，虽然还有一点小错误，但是不影响整体的理解，所以推荐看这个来搞清楚 CTC 的前世今生；
3. 关于 CTC 的原理图形化描述请参考[sequence modeling with ctc](#)<sup>(9)</sup>，用很多小例子来展现 CTC 的原理和解码等，有助于去理解结果；
4. 关于 CTC 求导部分<sup>(17)</sup> 的得到最后一步结果可以参考[教程：Connectionist Temporal Classification 详解补充](#)；
5. 关于 CTC 前后向 python 代码实现请参考[CTC 原理及实现](#)。

以上资料看完，我觉得就完全可以理解 CTC 的原理了，当然手推公式是少不了的，接下来进入正题。

#### 10.1.1 白话 CTC

首先明确：语音识别是一个序列分类的任务，那么其输入是逐帧的，如果不知道每一帧对应的标签，那么我们想要用 connectionist network 去干这件事，就得有目标函数。整个神经网络学习的准则就来自于目标函数，所以什么都可以没有，不能没有目标函数；话又说回来了，语音识别在训练的时候，逐帧输入，则必然对应着逐帧的输出，如果不要对齐，想要端到端，那么就必须想办法把这些逐帧的输出映射成序列输出。真正难的地方就在这儿，怎么去映射，映射了之后怎么定义目标函数：

CTC 就是来干这个事情的。

我们先回忆一下使用传统的 HMM-DNN 模型来搭建语音识别框架，其中包含很多个子任务。首先我们需要使用 EM 算法训练 HMM-GMM 模型以拿到对齐的数据，即一帧对应一个音

素。其次以这些数据来训练 DNN 模型，训练好了 DNN 模型之后，通过 HMM 和 vocabulary 映射到更高的建模单元，再结合词典、语言模型来用 WFST 进行解码。

好的，这个过程很复杂，而且每一块信息量都很足，即要花很多的精力去学习和设计；

不用担心，现在端到端已经越来越火热，效果看上去也是越来越好了。那么咱们说一说 CTC 是怎么干的。

CTC 呢，有一些很重要的设定：（1）每一帧的输出是相互独立的；（2）只要最终映射出来的序列是对的，在任何时刻输出某个标签都可以。有个图很好玩，来自 PilgrimHui 的博客<sup>(14)</sup>，如图10.1。我觉得这个图很好的说明了 CTC 的内涵和特点。

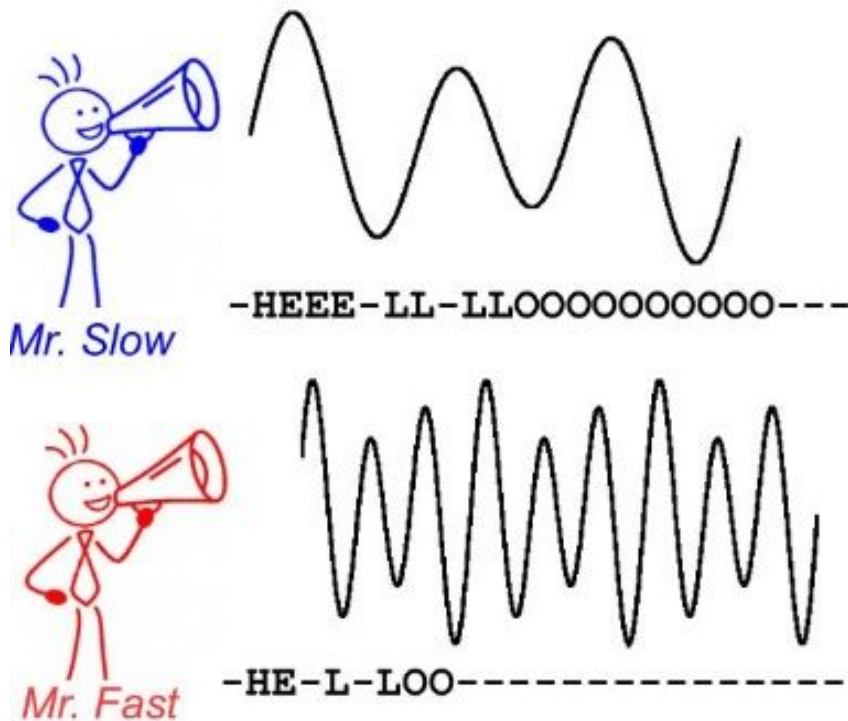


图 10.1: 同一单词不一样的输出音频的图解

第一个设定是为了后面使用最大似然估计去定义 loss 函数的，第二个设定呢就是彻底扔掉对齐数据的限制，按照这个设定，是没有固定的对齐结果的。在讲怎么从连续帧的输出映射到一整个序列之前，我们还要再讲一讲一些规则。为了避免相邻输出相同和更好的对建模单元的边界进行建模，CTC 引入了一个非常重要的输出标签：blank。你可以认为这个东西就是没有标签的意思，因为它的出现对句中无意义的片段比如静音段提供了建模单元。

假设字母表为  $A$ ，则  $A' = A \cup \{blank\}$ ，激活函数的输出  $y_k^t$  为  $t$  时刻输出  $A'$  中元素  $k$  的概率值，给定输入长度为  $T$  的音频序列  $x$ ；基于输出标签集合  $A'$ ，定义  $A'^T$  为长度为  $T$  的序列集合。

CTC 有如下假设：每一个时间步输出的标签概率值与其他时间步之间相互独立，或者说在

$x$  的条件下, 相互独立。

那么对于  $\pi \in A'^T$ , 其条件概率如公式10.1。

$$P(\pi|x) = \prod_{t=1}^T y_{\pi_t}^t \quad (10.1)$$

为了和域  $A$  内的标签序列  $l$  区分开来, 我们称  $\pi$  为域  $A'$  内的一条路径。由此定义一个 many-to-one 的函数:  $\mathcal{F}: A'^T \mapsto A^{\leq T}$ , 因为最终系统要输出的是整个正常的序列, 而不是穿插着一堆 blank 还有一堆重复 label 的奇奇怪怪的东西, 我们得到的  $\pi$  是沿着时间线严格输出的标签值, 而  $l$  是实际上某句话的内容, 所以需要  $\mathcal{F}$  来映射成我们需要的结果, 这样才能叫做端到端呀。举个例子, 比如某一段音频有 15 帧, 这 15 帧说的是英文单词 "beef", 这个就是  $l$ , 那实际上有 15 帧就会输出 15 帧的输出结果, 可能是这样的 "\_ \_ b b \_ e \_ \_ e e \_ f f \_ \_", 这个就是  $\pi$ ,  $\mathcal{F}$  干的事就是把  $\pi$  映射成  $l$ , 其规则是:

1. 合并重复的标签: "\_ \_ b b \_ e \_ \_ e e \_ f f \_ \_"  $\longrightarrow$  "\_ b \_ e \_ e \_ f \_ \_";
2. 移除 blank: "\_ b \_ e \_ e \_ f \_ \_"  $\longrightarrow$  "beef".

这样的路径有很多很多条, 而每一条路径都是排他的, 因此我们就可以通过公式10.2计算出域  $A$  内某个标签序列的概率值了, 即从一段音频中, 输出一句话的概率值, CTC 牛逼。

$$P(l|x) = \sum_{\pi \in \mathcal{F}^{-1}(l)} P(\pi|x) \quad (10.2)$$

这种整合的意义在于我不在乎你在什么时间点输出什么, 我只在乎最后映射的结果是我想要的。只要最终结果是一样的, 中间有多少个 blank, 同样的标签重复了多少次, 我不在乎。这就使得 CTC 可以使用没有对齐的数据来进行语音识别。

好的, 现在还剩下一些问题, 公式10.2这玩意怎么计算, 穷举吗? 计算机表示: 无能为力。序列越长, 标签越多, 尤其是碰到汉语这种用字建模的, 基本上想都不要想。回想起 HMM 中讲到的前后向算法, CTC 也可以这么干啊, 毕竟都是路径概率求解问题。

### 10.1.2 CTC 中的前后向算法

对于公式10.2, 假设标签序列的长度为  $U$ , 音频的长度为  $T$ , 共有  $2^{T-U^2+U(T-3)}3^{(U-1)(T-U)-2}$  条路径。

.....

说实话, 我不知道这个是怎么算出来的, 但是还是挺恐怖的, 还是按照上面举的例子,  $U = 4$ ,  $T = 15$ , 算一下就是  $2^{47}3^{31}$ 。

.....

再见!



前后向算法的核心是标签序列  $l$  对应的所有路径概率和可以转换成标签前缀的迭代求和。

为了能让 **blank** 出现在输出路径中，我们将  $l$  修改下，在  $l$  的前后都加上 **blank**，同时在每个标签之间也都加上 **blank**，记作  $l'$  即 "b e e f"  $\rightarrow$  "\_ b \_ e \_ e \_ f \_"。若  $l$  的长度为  $U$ ，则  $l'$  的长度为  $U' = 2U + 1$ 。为计算  $l'$  的前缀概率，我们允许 **blank** 和 **non-blank**、非重复的 **non-blank** 之间直接转换，多解释下这句话的意思："beef" 这个词，"b" 在下一个时刻直接跳转到 "b"、"e" 或者 "\_" 都是可以的，但是第一个 "e" 只能跳转到自己本身 "e" 和 "\_"，不可以跳转到下一个 "e"，不然的话，在进行合并操作的时候，第二个 "e" 就没了，这个需要好好理解，因为后面的前后向算法里面有用到这个特点。

对于一个标签序列  $l$ ，前向变量  $\alpha(t, u)$  表示所有长度为  $t$  的路径的概率和，这些路径通过  $\mathcal{F}$  映射到  $l$  中长度为  $u/2$  的前缀， $u/2$  向下取整。对于序列  $s$ ， $s_{p:q}$  为  $s$  的子序列  $s_p, s_{p+1}, \dots, s_{q-1}, s_q$ ，定义集合  $V(t, u) = \{\pi \in A^n : \mathcal{F}(\pi) = l_{1:u/2}, \pi_t = l'_u\}$ ，则  $\alpha(t, u)$  的计算如公式10.3。

$$\alpha(t, u) = \sum_{\pi \in V(t, u)} \prod_{i=1}^t y_{\pi_i}^i \quad (10.3)$$

我们可以通过  $t-1$  时刻的前缀来计算  $\alpha(t, u)$ 。

给定上面的公式，我们就可以求得  $l$  的概率，其等于  $T$  时刻标签为 **blank** 和 **non-blank** 的前向概率和，如公式10.4。

$$p(l|x) = \alpha(T, U') + \alpha(T, U' - 1) \quad (10.4)$$

所有正确的路径一定是以 **blank** 或者  $l$  中的第一个元素开头的，因此我们可以得到  $\alpha(t, u)$  的初始值，如公式10.5。

$$\begin{aligned} \alpha(1, 1) &= y_b^1 \\ \alpha(1, 2) &= y_{l_1}^1 \\ \alpha(1, u) &= 0, \forall u > 2 \end{aligned} \quad (10.5)$$

初始化的值得到之后，我们可以通过对前缀进行迭代求和的方式得到当前时刻标签索引为  $u$  的路径概率值，如公式10.6。

$$\alpha(t, u) = y_{l'_u}^t \sum_{i=f(u)}^u \alpha(t-1, i) \quad (10.6)$$

其中  $f(u)$  如公式10.7。

$$f(u) = \begin{cases} u-1 & l'_u = l'_{u-2} \\ u-2 & otherwise \end{cases} \quad (10.7)$$

我觉得有必要解释下公式10.7。

公式10.6左边表示的是在  $t$  时刻标签索引为  $u$  的输出概率，右边是在已经确定了  $t$  时刻输出索引为  $u$  的情况下， $t-1$  时刻的前向概率。因为已经确定了  $t$  时刻的输出，因此我们需要知道在  $t-1$  时刻，哪些标签能够在下一个时刻转移到标签  $u$ 。对这些可能的标签前向概率求和，再乘以当前时刻的概率，就得到了我们想要的  $\alpha(t, u)$ 。

那么  $t-1$  时刻究竟可能是哪些标签呢？

我们还是拿"beef" 举例子。

$l = \text{"b e e f"} , l' = \text{"_ b _ e _ e _ f _"} ,$  为了说的更清楚，我们来给  $l'$  里面的字母都来标个号， $l' = \text{"_(1) b(2) _(3) e(4) _(5) e(6) _(7) f(8) _(9)"} ,$  先来看第一种情况，如果  $t$  时刻的标签  $u$  有  $l'_u = l'_{u-2}$ ，那么  $l'_u = \text{"e(6)"} ,$  或者  $l'_u = \text{"_(i)"} , i = \{3, 5, 7, 9\} .$

1.  $l'_u = \text{"e(6)"} :$  那么  $t-1$  时刻的标签等于  $\text{"e(6)"} ,$  是完全有可能的；等于  $_(5)$  也是完全可以的；那么能等于  $\text{e(4)}$  可以吗？假设可以，连起来看  $t-1$  和  $t$  时刻的子序列，就是  $\text{"e(4)e(6)"} ,$  乍一看好像没问题啊，但是要知道我们在训练的时候  $\text{"e"}$  是标签，是没有  $\text{"e(4)"} ,$  和  $\text{"e(6)"} ,$  这种形式的标签的，他们都是  $\text{"e"} ,$  根据上面说的 CTC 的合并规则，这俩就会合并成一个，无法表征两个相同的连续元素。也就是说，CTC 是不允许两个相同的连续元素直接跳转过去的， $\text{"e(4)"} ,$  是不可以直接跳转到  $\text{"e(6)"} ,$  所以啊， $f(u) = u-1$  表示前一个时刻的标签索引只可能和当前时刻的索引相同或者是当前时刻的索引的前一个；
2.  $l'_u = \text{"_(i)"} , i = \{3, 5, 7, 9\} :$  这个解释起来和上面是一样一样的。同样我们的输出标签里面就只有  $\text{"_"} ,$  这个东西，所以按照合并规则，但凡重复了的  $\text{"_"} ,$  都看作是同一个  $\text{"_"} ,$  因此是没有办法表征序列  $l'$  中的上一个  $\text{"_"} ,$  因此  $f(u) = u-1 .$

当  $l'_u \neq l'_{u-2} ,$  这个就比较无所谓了，咱们举个例子，假设  $l'_u = \text{"f(8)"} ,$  同样  $t-1$  时刻的输出为  $\text{"f(8)"} ,$  没得问题， $\text{"_(7)"} ,$  也没得问题， $\text{"e(6)"} ,$  当然也没得问题，反正他们不会乱合并。所以有  $f(u) = u-2 .$

对应  $t$  时刻的  $u$  不能太小，太大无所谓大，大不了剩下的  $T-t$  帧都重复输出某一个标签或者 blank 就可以了。如果太小，剩下的  $T-t$  帧的输出都不够剩下还没有输出的标签了，所以考虑最极端的一种情况就是剩下的每一帧都输出一个  $l$  中的标签，且这些标签不会自身重复，也就是说每一帧输出的标签都只会输出一次，这样刚好够完整的输出  $l$ 。这就是临界条件，因此我们可以得到  $u$  的取值范围：

$$T-t \geq \frac{U' - u - 1}{2} \quad (10.8)$$

解得：

$$u \geq U' - 2(T-t) - 1 \quad (10.9)$$

反过来说，也就是意味着：

$$\alpha(t, u) = 0, \forall u < U' - 2(T - t) - 1 \quad (10.10)$$

前向算法讲完了，后向算法就很类似了。同样我们定义一个后向变量  $\beta(t, u)$ ，如果说想要得到输出序列  $l$ ，假定  $t$  时刻的输出是  $l'_u$ ， $\beta(t, u)$  可以看成是这一些路径的后部分，我们定义  $\beta(t, u)$  为  $t + 1$  时刻配合上  $\alpha(t, u)$  恰好构成完整的  $l$  的所有的后缀路径的概率总和。比照着  $V(t, u)$ ，定义集合  $W(t, u) = \{\pi \in A^{T-t} : \mathcal{F}(\hat{\pi} + \pi) = l, \forall \hat{\pi} \in V(t, u)\}$ 。那么直接计算  $\beta(t, u)$  的公式如10.11。

$$\beta(t, u) = \sum_{\pi \in W(t, u)} \prod_{i=1}^{T-t} y_{\pi_i}^{t+i} \quad (10.11)$$

后向算法的初始化如公式10.12。

$$\begin{aligned} \beta(T, U') &= 1 \\ \beta(T, U' - 1) &= 1 \\ \beta(T, u) &= 0, \forall u < U' - 1 \\ \beta(T, U' + 1) &= 0 \end{aligned} \quad (10.12)$$

公式10.13描述了  $\beta(t, u)$  的计算方式。

$$\beta(t, u) = \sum_{i=u}^{g(u)} \beta(t + 1, i) y_{l'_i}^{t+1} \quad (10.13)$$

其中：

$$g(u) = \begin{cases} u + 1 & l'_u = l'_{u+2} \\ u + 2 & otherwise \end{cases} \quad (10.14)$$

公式10.14表明了对  $t + 1$  时刻能达到的标签的限制，这部分和前向很类似，在此就不再赘述了。

同样后向算法中的  $u$  不能太大，因为  $u$  过大的意思就是前  $t$  个时刻得输出  $u$  个标签，如果说  $t$  稍微小了点，那就不够输出这么多标签了。同样，考虑临界条件：前  $t$  个时刻恰好能够无重复的输出 non-blank 标签。我们就可以得到如下不等式：

$$t \geq \frac{u}{2} \quad (10.15)$$

解得：

$$u \leq 2t \quad (10.16)$$

即：

$$\beta(t, u) = 0, \forall u > 2t \quad (10.17)$$

上述迭代过程，中间会涉及到概率的求和，如果直接就在计算机中按照上述方式实现的话，很容易就会导致 **underflows**，所以需要将概率变成 **log** 概率值，对于两个概率的和，我们一般使用公式10.18来实现，等到算完之后，再对其进行指数操作，就可以还原回来了。此外大部分的编程语言都会有一个稳定的函数去计算  $\log(1 + x)$ ，当  $x$  接近于 0 的时候。

$$\log(e^a + e^b) = \max\{a, b\} + \log(1 + e^{-|a-b|}) \quad (10.18)$$

到这儿，CTC 中的前后向算法就讲完了，下一步是定义 **loss** 函数和求梯度。

### 10.1.3 CTC 中的 **loss** 函数和梯度

对于数据集  $S$ ，CTC 的损失函数  $\mathcal{L}(S)$  定义为数据集中所有输出序列正确的负 **log** 概率，如公式10.19。

$$\begin{aligned} \mathcal{L}(S) &= -\ln \prod_{(x,z) \in S} p(z|x) \\ &= -\sum_{(x,z) \in S} \ln p(z|x) \end{aligned} \quad (10.19)$$

因为损失函数是可微分的，所以我们可以通过 **BPTT**（以后这块也得写个总结）去求权重对的梯度。所以任何一个基于梯度的非线性优化函数都可用来训练 **CTC** 网络。我们先看下单个样本的 **loss**，如公式10.20。

$$\mathcal{L}(x, z) = -\ln p(z|x) \quad (10.20)$$

则有：

$$\mathcal{L}(S) = \sum_{(x,z) \in S} \mathcal{L}(x, z) \quad (10.21)$$



从而有：

$$\frac{\partial \mathcal{L}(\mathcal{S})}{\partial w} = \sum_{(x,z) \in \mathcal{S}} \frac{\partial \mathcal{L}(x,z)}{\partial w} \quad (10.22)$$

设  $l = z$ ，定义集合  $X(t,u) = \{\pi \in A^T : \mathcal{F}(\pi) = z, \pi_t = z'_u\}$ ，这个  $X(t,u)$  表示的是所有最终整合之后的序列为  $z$ ，且在  $t$  时刻输出索引为  $u$  的路径集合。根据公式10.3和公式10.11，我们可以得到在集合  $X(t,u)$  中的所有路径元素的概率和，如公式10.23。

$$\alpha(t,u)\beta(t,u) = \sum_{\pi \in X(t,u)} \prod_{t=1}^T y_{\pi_t}^t \quad (10.23)$$

公式10.23中说明了当  $t$  时刻，输出的标签索引为  $u$  的所有路径的可能性，那么我们希望获得的是输出  $z$  的概率值，那么我们只需要知道  $t$  时刻共有多少个可能的标签，然后对这些标签按照公式10.23去求解，最后将这些可能性都加起来，不就得到了输出序列  $z$  的概率值了吗？

所以我们得到了输出序列  $z$  的概率值如公式10.24。

$$p(z|x) = \sum_{u=1}^{|z'|} \alpha(t,u)\beta(t,u) \quad (10.24)$$

那么单个样本的损失函数我们就可以算出来了，如公式10.25。

$$\mathcal{L}(x,z) = -\ln \sum_{u=1}^{|z'|} \alpha(t,u)\beta(t,u) \quad (10.25)$$

ok，接下来，我们就求一下  $t$  时刻单个样本的损失函数对 softmax 层神经元输出的梯度，只要这个求出来了，其他的就比较简单了。其梯度如公式10.26。

$$\begin{aligned} \frac{\partial \mathcal{L}(x,z)}{\partial y_k^t} &= -\frac{\partial \ln p(z|x)}{\partial y_k^t} \\ &= -\frac{1}{p(z|x)} \frac{\partial p(z|x)}{\partial y_k^t} \\ &= -\frac{1}{p(z|x)} \frac{\partial}{\partial y_k^t} \sum_{u=1}^{|z'|} \alpha(t,u)\beta(t,u) \\ &= -\frac{1}{p(z|x)} \sum_{u=1}^{|z'|} \frac{\partial \alpha(t,u)\beta(t,u)}{\partial y_k^t} \end{aligned} \quad (10.26)$$

因为  $p(z|x)$  是已知的，我们把重点放在求解  $\frac{\partial p(z|x)}{\partial y_k^t}$  上。由于我们是对标签  $k$  所对应的神经

元输出求导，所以我们只需要去理会那些在  $t$  时刻输出为  $k$  的路径即可，那么我们定义个集合  $B(z, k) = \{u : z'_u = k\}$ ，结合公式10.23我们可以得到结论10.27。

$$\frac{\partial \alpha(t, u) \beta(t, u)}{\partial y_k^t} = \begin{cases} \frac{\alpha(t, u) \beta(t, u)}{y_k^t} & \text{if } k \text{ occurs in } z' \\ 0 & \text{otherwise} \end{cases} \quad (10.27)$$

这里我又得嘴碎多说一句，为啥这个导数是这个样子的，因为从公式10.23中，我们可以看出来那些求和单元每一个路径中都包含了一个  $y_k^t$ ，这求导以后就把这个给弄没了，那么为了后续的运算，我们分子分母同时乘以  $y_k^t$ ，这样就可以得到这一块。

首先我们回忆下 softmax 的输出  $y_{k'}^t$  对  $a_k^t$  的导数，见公式10.28，这块的求导细节不赘述了。

$$\frac{\partial y_{k'}^t}{\partial a_k^t} = y_{k'}^t (\delta_{k'k} - y_k^t) \quad (10.28)$$

接下来我们结合上面的这些信息来对 softmax 层前的那一层神经元的输出  $a_k^t$  求导，如公式10.29。

$$\begin{aligned} \frac{\partial \mathcal{L}(x, z)}{\partial a_k^t} &= \sum_{k'} \frac{\partial \mathcal{L}(x, z)}{\partial y_{k'}^t} \frac{\partial y_{k'}^t}{\partial a_k^t} \\ &= -\frac{1}{p(z|x) y_{k'}^t} \sum_{k'} \sum_{u \in B(z, k)} \alpha(t, u) \beta(t, u) y_{k'}^t (\delta_{k'k} - y_k^t) \\ &= -\frac{1}{p(z|x)} \sum_{k'} \sum_{u \in B(z, k)} \alpha(t, u) \beta(t, u) (\delta_{k'k} - y_k^t) \\ &= -\frac{1}{p(z|x)} \sum_{u \in B(z, k)} \alpha(t, u) \beta(t, u) + \frac{1}{p(z|x)} \sum_{k'} \sum_{u \in B(z, k)} \alpha(t, u) \beta(t, u) y_k^t \\ &= -\frac{1}{p(z|x)} \sum_{u \in B(z, k)} \alpha(t, u) \beta(t, u) + y_k^t \frac{1}{p(z|x)} \sum_{k'} \sum_{u \in B(z, k)} \alpha(t, u) \beta(t, u) \\ &= -\frac{1}{p(z|x)} \sum_{u \in B(z, k)} \alpha(t, u) \beta(t, u) + y_k^t \end{aligned} \quad (10.29)$$

解释下倒数第二步是怎么转换到倒数第一步的。 $\sum_{u \in B(z, k)} \alpha(t, u) \beta(t, u)$  表示的是  $t$  时刻经过标签索引为  $u$  的所有路径之和，加上限定  $u \in B(z, k)$  之后，这个式子指的就是当第  $u$  个标签为  $k$  的所有路径的概率和。 $\sum_{k'}$  中的  $k'$  表示的就是标签，这一求和就意味着算了一下在  $t$  时刻

所有可能的标签的所有路径的概率和，咱们翻过头看看公式10.24，我们可以得到公式10.30。

$$\begin{aligned} p(z|x) &= \sum_{u=1}^{|z'|} \alpha(t, u) \beta(t, u) \\ &= \sum_{k'} \sum_{u \in B(z, k)} \alpha(t, u) \beta(t, u) \end{aligned} \quad (10.30)$$

所以……以上就是梯度的求导过程。另外原论文中在链式求导的那块有个公式错误如图10.2，多了个负号。

$$\frac{\partial \mathcal{L}(\mathbf{x}, \mathbf{z})}{\partial y_k^t} = -\frac{1}{p(\mathbf{z}|\mathbf{x})y_k^t} \sum_{u \in B(\mathbf{z}, k)} \alpha(t, u) \beta(t, u). \quad (7.31)$$

Finally, to backpropagate the gradient through the output layer, we need the loss function derivatives with respect to the outputs  $a_k^t$  before the activation function is applied:

$$\frac{\partial \mathcal{L}(\mathbf{x}, \mathbf{z})}{\partial a_k^t} = - \sum_{k'} \frac{\partial \mathcal{L}(\mathbf{x}, \mathbf{z})}{\partial y_{k'}^t} \frac{\partial y_{k'}^t}{\partial a_k^t} \quad (7.32)$$

where  $k'$  ranges over all the output units. Recalling that for softmax outputs

$$y_k^t = \frac{e^{a_k^t}}{\sum_{k'} e^{a_{k'}^t}}$$

图 10.2: 原 Alex 博士论文中链式求导部分的错误

## 10.1.4 CTC 的解码

CTC 的解码常用的有两种方式，greedy search 和 prefix beam search。greedy 解码速度很快，但是很容易出错，但是 prefix beam search 解码速度慢，准确率较高。接下来挨个介绍两种解码方式的算法和流程，以及对应的代码解释。

### 10.1.4.1 greedy search

greedy search，也就是 best path search，就是说找个每个时间点输出的最大概率值对应的标签，然后去重，再去掉 blank。这种方式很快，因为不需要考虑整个路径输出，只考虑当前时刻的输出即可。。但是这种方式效果不太好，输出可能比较奇怪。代码如下，来自<sup>(6)</sup>，这个库中还包含了一些关于 CTC 解码的论文。

```
1 def ctcBestPath(mat, classes):
2     "implements best path decoding as shown by Graves (Dissertation, p63)"
3
4     # dim0=t, dim1=c
5     maxT, maxC = mat.shape #blank的索引是最后一个
```

```

6  label = '' #初始化输出标签序列
7  blankIdx = len(classes)
8  lastMaxIdx = maxC # init with invalid label
9
10 for t in range(maxT):
11     maxIdx = np.argmax(mat[t, :]) #找到当前帧最大标签的索引
12
13     if maxIdx != lastMaxIdx and maxIdx != blankIdx: #去重去blank
14         label += classes[maxIdx]
15
16     lastMaxIdx = maxIdx #重置上一个标签的索引
17
18 return label

```

#### 10.1.4.2 prefix beam search

First-Pass Large Vocabulary Continuous Speech Recognition using Bi-Directional Recurrent DNNs<sup>(10)</sup> 中针对 CTC 提出了一种 prefix beam search 的算法, 这种算法能够避免全局搜索的复杂度过高无法实现的问题, 同时比 greedy search 的结果要好很多。

首先介绍下一些公式。公式10.31是最终的解码公式, 配合语言模型和网络输出(声学模型), 得到最有可能的  $W$  词序列。

$$W = \arg \max_W p_{net}(W; X) p_{lm}^\alpha(W) |W|^\beta \quad (10.31)$$

其中  $p_{net}$  是网络的输出,  $p_{lm}$  是语言模型的输出,  $\alpha$  和  $\beta$  分别为语言模型的权重和补偿系数。一般情况下, 我们会降低语言模型对整体输出的影响, 所以  $\alpha$  一般取 0.2 ~ 0.7。

接下来我们讲一下 prefix beam search 的 demo<sup>(11)</sup>。

总体是有三个循环, 第一个是时间维度上的, 时间  $t$  从 1 到  $T$ , 也就是逐帧解码。第二个是对应候选输出序列的, 这个是 beam search, 那么一定得设置一个 beam size, 那么会考虑所有的候选序列跟当前输出结合起来的概率值, 那当前输出的话被剪枝之后还有很多个标签, 再挨个的把每一个候选和每一个当前帧的标签进行结合计算。然后再利用这个结合的概率值进行重新排序得到新的候选。

```

1  pruned_alphabet = [alphabet[i] for i in np.where(ctc[t] > prune)[0]]

```



这一步是为了减少计算，也就是说先设定一个阈值，当前  $t$  时刻的时候，会做一个判断，只有当前时刻各个标签概率值大于 `prune` 的时候才会去做后续的操作，这就意味着当前时刻所有概率低于 `prune` 的标签都会被抛弃，不再参与到当前时刻的解码过程中。因为这些标签概率值太小，不太可能是正确的输出标签，留着只会增加计算量，还不如删掉，省时省心省力！

```

1 if len(l) > 0 and l[-1] == '>':
2     Pb[t][l] = Pb[t - 1][l]
3     Pnb[t][l] = Pnb[t - 1][l]
4     continue

```

这一步就是判断下是不是到结尾了，结尾的表示符号是">". 如果到了结尾呢，说明这个时候输出的标签序列和  $t - 1$  时刻是一毛一样滴。 $t$  时刻输出序列  $l$  以 `blank` 结尾的概率跟  $t - 1$  时刻以 `blank` 结尾的概率是一样的， $t$  时刻输出序列  $l$  以 `non-blank` 结尾的概率跟  $t - 1$  时刻以 `non-blank` 结尾的概率是一样的。然后就跳出当前时刻，因为当前时刻表示这个序列已经到了结尾了，没必要再折腾下去了。

```

1 if c == '%':
2     Pb[t][l] += ctc[t][-1] * (Pb[t - 1][l] + Pnb[t - 1][l])

```

我们假设 `%` 代表 `blank` 这个标签。剪枝之后，会对还剩下的那些标签走一遍遍历，每一个标签都会尝试着和之前存起来的候选序列进行结合，算出来一个概率值。那么既然是遍历，当然会轮到牛逼的 `blank`。所以首先看看当前的这个标签是不是 `blank`。如果是 `blank` 的话，我们就没必要对当前的这个候选序列做啥子改动了，也就是当前时刻的输出标签序列还是  $l$ ，因为最终输出的时候，`blank` 也不会出现。既然当前这个标签是 `blank`，那么  $t$  时刻的标签序列  $l$  的概率需要和当前时刻输出为 `blank` 的概率结合一下，变成当前时刻的  $Pb[t][l]$ 。其计算公式从代码里就可以看到。

```

1 l_plus = l + c
2 if len(l) > 0 and c == l[-1]:
3     Pnb[t][l_plus] += ctc[t][c_ix] * Pb[t - 1][l]
4     Pnb[t][l] += ctc[t][c_ix] * Pnb[t - 1][l]

```

如果说当前时刻的标签不是 `blank`，而是上个时刻的这个候选序列的最后一个输出标签，也就是说当前时刻的输出和上一个时刻的候选序列尾部产生了重复，这个时候其实是有两种情况的，第一种情况是上一个时刻的输出标签正好是 `blank`，因为从上面那一步代码中我们可以看出

来, 候选序列中是不会出现 blank 的, 那如果是这样的情况, 说明实际的输出序列就是有两个一样的字母, 输出就是  $l\_plus$ , 其尾部有两个一样的字母, 这个时候候选序列的概率就等于当前时刻的标签概率乘以上一个时刻输出为 blank 的序列概率, 也就是  $Pb[t-1][l]$ ; 第二种情况是上一个时刻的输出确实也是这个标签, 那么说明这个时候的候选序列不需要做啥变动, 但是概率值需要变一下, 当前时刻标签概率乘以上一个时刻输出是 non-blank 的序列概率值。

```

 $\ell^+ \leftarrow \text{concatenate } \ell \text{ and } c$ 
if  $c = \ell_{\text{end}}$  then
     $p_{nb}(\ell^+; x_{1:t}) \leftarrow p(c; x_t) p_b(\ell; x_{1:t-1})$ 
     $p_{nb}(\ell; x_{1:t}) \leftarrow p(c; x_t) (p_b(\ell; x_{1:t-1})$ 
else if  $c = \text{space}$  then
     $p_{nb}(\ell^+; x_{1:t}) \leftarrow p(W(\ell^+) | W(\ell))^\alpha p(c; x_t) (p_b(\ell; x_{1:t-1}) + p_{nb}(\ell; x_{1:t-1}))$ 
else
     $p_{nb}(\ell^+; x_{1:t}) \leftarrow p(c; x_t) (p_b(\ell; x_{1:t-1}) + p_{nb}(\ell; x_{1:t-1}))$ 
end if

```

图 10.3: Prefix Beam Search 原论文中算法错误地方

另外原论文中关于这一块的计算步骤写错了, 也就是算法中的这一步, 如图10.3。简直坑死个人。

```

1 elif len(l.replace(' ', '')) > 0 and c in (' ', '>'):
2     lm_prob = lm(l_plus.strip(' >')) ** alpha
3     Pnb[t][l_plus] += lm_prob * ctc[t][c_ix] * (Pb[t - 1][l] + Pnb[t - 1][l]
        ])

```

那还有可能当前的输出是' '(space), 也就是说输出是空格或者是结尾, 这个时候说明一个完整的词出现了, 我们就可以利用语言模型 (LM) 来对输出进行修正和约束, 避免出现毫无意义的结果。那么这个词代入到语言模型中会出现一个概率值, 当前候选序列的概率值就通过公式10.31来计算, 从代码中也可以看出来。

```

1 Pnb[t][l_plus] += ctc[t][c_ix] * (Pb[t - 1][l] + Pnb[t - 1][l])

```

还有最后一种情况, 就是既不是 blank, 又不是 space, 当前输出标签和候选标签序列的最后一个字母也不一样, 这个时候, 就直接算出候选标签序列和当前标签的概率乘积就行, 候选标签序列也有两种情况: 上一个时刻以 blank 或者以 non-blank 结尾。综上集中基本的情况都已经讲完了。

```

1 if l_plus not in A_prev:
2     Pb[t][l_plus] += ctc[t][-1] * (Pb[t - 1][l_plus] + Pnb[t - 1][l_plus])
3     Pnb[t][l_plus] += ctc[t][c_ix] * Pnb[t - 1][l_plus]

```

按照原始论文中，还有上面这个公式。也就是说算出来的  $l_{plus}$  不在候选标签序列里面，就会去之前时刻的候选序列里面去找，再利用之前的后续序列概率计算当前的概率值。百度的 Deep Speech2 代码里面说：这个部分不知道在干啥，还没啥用，所以就给去掉了。

我觉得……也不太好理解……

讲到这儿核心的代码部分已经讲完了，剩下的就是把当前时刻的标签，不管是以 **blank** 结尾的还是非 **blank** 结尾的综合起来，然后进行排序，根据 **beam size** 的大小得到新的候选序列，如此循环往复，直到这个序列输出到了尽头，就可以得到最终的结果啦~

完整代码如下：

```
1 from collections import defaultdict, Counter
2 from string import ascii_lowercase
3 import re
4 import numpy as np
5
6 def prefix_beam_search(ctc, lm=None, k=25, alpha=0.30, beta=5, prune
   =0.001):
7     """
8     Performs prefix beam search on the output of a CTC network.
9
10    Args:
11        ctc (np.ndarray): The CTC output. Should be a 2D array (timesteps x
12            alphabet_size)
13        lm (func): Language model function. Should take as input a string and
14            output a probability.
15        k (int): The beam width. Will keep the 'k' most likely candidates at
16            each timestep.
17        alpha (float): The language model weight. Should usually be between 0
18            and 1.
19        beta (float): The language model compensation term. The higher the '
20            alpha', the higher the 'beta'.
21        prune (float): Only extend prefixes with chars with an emission
22            probability higher than 'prune'.
23
24    Returns:
```

```

19     string: The decoded CTC output.
20     """
21
22     lm = (lambda l: 1) if lm is None else lm # if no LM is provided, just set
        to function returning 1
23     W = lambda l: re.findall(r'\w+[\s|>]', l)
24     alphabet = list(ascii_lowercase) + ['<math>\epsilon</math>', '>', '%']
25     F = ctc.shape[1]
26     ctc = np.vstack((np.zeros(F), ctc)) # just add an imaginative zero'th
        step (will make indexing more intuitive)
27     T = ctc.shape[0]
28
29     # STEP 1: Initiliazation
30     O = ''
31     Pb, Pnb = defaultdict(Counter), defaultdict(Counter)
32     Pb[0][0] = 1
33     Pnb[0][0] = 0
34     A_prev = [0]
35     # END: STEP 1
36
37     # STEP 2: Iterations and pruning
38     for t in range(1, T):
39         pruned_alphabet = [alphabet[i] for i in np.where(ctc[t] > prune)[0]]
40         for l in A_prev:
41
42             if len(l) > 0 and l[-1] == '>':
43                 Pb[t][1] = Pb[t - 1][1]
44                 Pnb[t][1] = Pnb[t - 1][1]
45                 continue
46
47         for c in pruned_alphabet:
48             c_ix = alphabet.index(c)
49             # END: STEP 2
50

```

```

51     # STEP 3: "Extending" with a blank
52     if c == '%':
53         Pb[t][1] += ctc[t][-1] * (Pb[t - 1][1] + Pnb[t - 1][1])
54     # END: STEP 3
55
56     # STEP 4: Extending with the end character
57     else:
58         l_plus = 1 + c
59         if len(l) > 0 and c == l[-1]:
60             Pnb[t][l_plus] += ctc[t][c_ix] * Pb[t - 1][1]
61             Pnb[t][1] += ctc[t][c_ix] * Pnb[t - 1][1]
62         # END: STEP 4
63
64         # STEP 5: Extending with any other non-blank character and LM
65         # constraints
66         elif len(l.replace(' ', '')) > 0 and c in (' ', '>'):
67             lm_prob = lm(l_plus.strip(' >')) ** alpha
68             Pnb[t][l_plus] += lm_prob * ctc[t][c_ix] * (Pb[t - 1][1] + Pnb[t
69             - 1][1])
70         else:
71             Pnb[t][l_plus] += ctc[t][c_ix] * (Pb[t - 1][1] + Pnb[t - 1][1])
72         # END: STEP 5
73
74         # STEP 6: Make use of discarded prefixes
75         if l_plus not in A_prev:
76             Pb[t][l_plus] += ctc[t][-1] * (Pb[t - 1][l_plus] + Pnb[t - 1][
77             l_plus])
78             Pnb[t][l_plus] += ctc[t][c_ix] * Pnb[t - 1][l_plus]
79         # END: STEP 6
80
81     # STEP 7: Select most probable prefixes
82     A_next = Pb[t] + Pnb[t]
83     sorter = lambda l: A_next[l] * (len(W(l)) + 1) ** beta
84     A_prev = sorted(A_next, key=sorter, reverse=True)[:k]

```

```
82     # END: STEP 7
83
84     return A_prev[0].strip('>')
```

## 10.2 RNN-Tranducer

RNN-Tranducer, 简称 RNN-T, 同样是 Alex Graves 大神的作品<sup>(8;7)</sup>。

## 10.3 Attention

## 10.4 Transformer

## 10.5 CNNs

## 10.6 Mixed Models

### 10.6.1 Self-Attention Transducers for End-to-End Speech Recognition

这篇论文的作者是田正坤, 来自中国科学院自动化所。本论文的主要贡献有:

1. 用 self-attention 模块替代了原来 RNN-T 中的 RNN 部分, 可以用于并行计算;
2. 利用 path-aware regularization 帮助 SA-T 学习对齐;
3. 使用了 chunk-flow 机制来进行解码。

#### 10.6.1.1 SA-T 的基本结构

#### 10.6.1.2 path-aware regularization

#### 10.6.1.3 chunk flow mechanism

## 10.7 如何计算 WER?

## 第 11 章 论文阅读笔记

### 11.1 Light Gated Recurrent Units for Speech Recognition

**Li-GRU**是 Mirco Ravanelli 于 2018 年发表的论文，他也是**pytorch-kaldi**的作者。这篇论文主要针对的是对**GRU**(Gated Recurrent Units) 的改进，而且 Li-GRU 是专门为语音识别去设计的。本论文主要的工作有两方面：

(1) 去掉了重置门 (reset gate)，去掉重置门对于模型的效果没什么影响，而且原始的 GRU 的重置门和更新门之间有冗余，因此去掉重置门的模型结构更合理；

(2) 将原始 GRU 中的激活函数 Tanh 换成了 Relu，由于 Relu 本身函数具备的特性，其效果比 Relu 要好多。之所以以前的 RNN（包括 GRU 和 LSTM）不用 Relu，是因为 Relu 的值可以任意大，RNN 不停的迭代中，Relu 的值无法控制，容易导致数值不稳定（numerical instability）。本文作者采用了批量正则（Batch Normalization）的方式来避免数值不稳定的情况。这么做既可以避免梯度消失，又可以加速网络收敛，减少网络的时间。

本节的论文笔记分为三个部分：（1）GRU 的介绍；（2）Li-GRU 的学习；（3）实验配置和结果；（4）个人心得体会。

#### 11.1.1 GRU 的介绍

语音识别是一个序列任务，那么上下文的信息对当前时刻信息的影响很大，RNN 的结构表明其可以动态的决定对于当前时间步使用多少上下文的信息。但是 RNN 存在的梯度消失和爆炸问题使得其学习长期依赖变得困难。所以一般我们都会使用一种门控 RNN(Gated RNN) 来解决这个问题。门控 RNN 的核心思想是引入一种门机制来控制不同时间步之间的信息流动。

常用的门控 RNN 有两种：LSTM 和 GRU。LSTM 的结构复杂且运算效率比较低，LSTM 有三个门，而 GRU 只有两个，所以运算起来快很多。而本论文也是基于 GRU 做的改进，所以不讨论 LSTM。

GRU 的计算公式如下：

$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z) \quad (11.1)$$

$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r) \quad (11.2)$$

$$\tilde{h}_t = \tanh(W_h x_t + U_h (h_{t-1} \odot r_t) + b_h) \quad (11.3)$$

$$h_t = z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t \quad (11.4)$$

其中

- $x_t$ :  $t$  时刻输入特征向量
- $r_t$ :  $t$  时刻重置门向量
- $u_t$ :  $t$  时刻更新门向量
- $h_t$ :  $t$  时刻状态向量
- $\tilde{h}_t$ :  $t$  时刻候选状态向量
- $W, U, b$ : 参数矩阵和向量

由公式11.4可知, 当前状态向量  $h_t$  是前一刻状态向量  $h_{t-1}$  和当前时刻的状态候选向量  $\tilde{h}_t$  之间的一个线性插值。两者之间的权重由更新门  $z_t$  决定, 权重的值表示了更新信息的多少。这个线性插值就是 GRU 学习长期依赖的核心。如果  $z_t$  接近于 1, 那么先前状态的信息就得以保留, 以此学习到间隔长的时间步之间的信息关联。如果  $z_t$  接近于 0, 那么网络更倾向于候选状态  $\tilde{h}_t$ , 而候选状态更依赖于当前输入和临近时间步的状态。同时候选状态还依赖于重置门  $r_t$ , 其使得模型通过忘记之前计算的状态来清除过去的记忆。

GRU 的模型结构图如11.1所示。

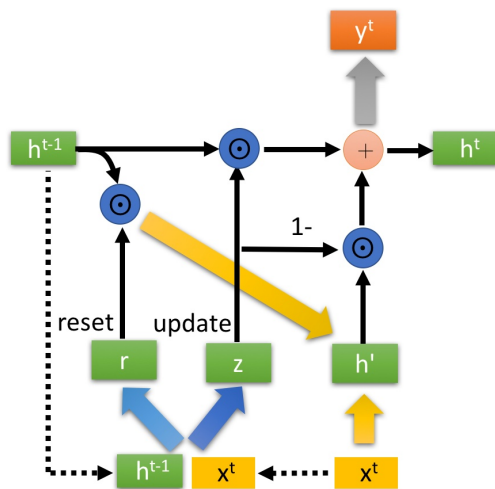


图 11.1: GRU 模型结构图

### 11.1.2 Li-GRU 的学习

本论文对 GRU 的改进主要涉及到三个部分: 重置门、ReLU 激活函数和 BN。

#### 1. 移除重置门:

对于序列中可能出现的 **significant discontinuity**, 重置门可以起到一个清除过去信息的作用。比如说语言建模, 当输入从一个句子跳转到另一个语义无关的句子的时候, 重置门就可以起到很好的作用: 避免过去无关信息对当前状态的干扰。但是对于语音识别来说, 重置门的作用可能就不明显了, 语音识别中的输入变化都比较小 (一般的偏移量才 10ms), 这表明过去的信息



还是挺有用的。即便是元音（Vowel）和擦音（Fricative）间的边界有很强不连续现象，完全去掉过去信息也可能是有害的。另外基于一些音素的转移更相似，存住 phonotactic features 还是很有用的。

与此同时，重置门和更新门之间存在着某种冗余。也就是说  $z_t$  和  $r_t$  的变化比较同步，当前输入信息比较重要的时候， $z_t$  和  $r_t$  都比较小；过去时刻信息比较重要的时候， $z_t$  和  $r_t$  都比较大。拿 TIMIT 中的一段音频来看，更新门和重置门的平均激活值有着时域上的关联性，如图 11.2。

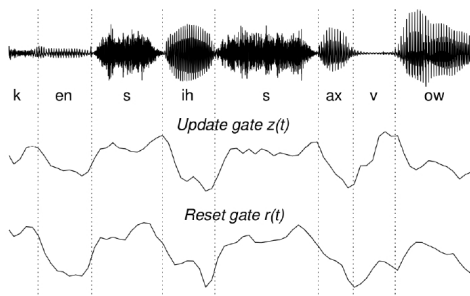


Fig. 1: Average activations of the update and reset gates for a GRU trained on TIMIT in a chunk of the utterance “sx403” of speaker “faks0”.

图 11.2: TIMIT 中更新门与重置门音频上的时域关联

两个门之间的冗余程度可以量化描述，其公式 cross-correlation  $C(z, r)$  如下：

$$C(z, r) = \bar{z}_t \star \bar{r}_t \quad (11.5)$$

其中：

- $\bar{z}_t$ : 更新门神经元的平均激活值
- $\bar{r}_t$ : 重置门神经元的平均激活值
- $\star$ : cross-correlation 的算子

图 11.3 显示了重置门和更新门的 cross-correlation。门激活值计算的是所有的输入帧，所有的隐藏神经元的激活值计算了个平均。从图中我们可以看到重置门和更新门之间相似度还挺高，因此冗余程度也很高。

因此我们决定去掉重置门，那么公式 11.3 就变成了公式 11.6。这么处理之后，运算效率就提高了，就剩下一个门，所以 GRU 的结构变得更紧凑了。

$$\tilde{h}_t = \tanh(W_h x_t + U_h h_{t-1} + b_h) \quad (11.6)$$

## 2. ReLU 激活函数

我们将  $\tanh$  替换成  $\text{ReLU}$ 。 $\tanh$  其实在前馈神经网络中很少使用，因为当神经网络的层数变深时，它就容易陷入左右的边界值  $-1$  和  $1$ ，此时的梯度接近于  $0$ ，网络参数更新缓慢，收敛

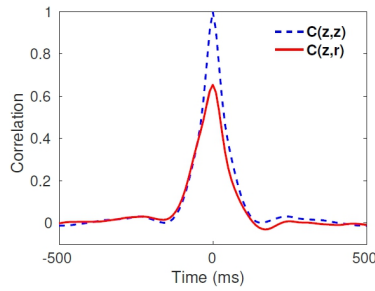


Fig. 2: Auto-correlation  $C(z, z)$  and cross-correlation  $C(z, r)$  between the average activations of the update ( $z$ ) and reset ( $r$ ) gates. Correlations are normalized by the maximum of  $C(z, z)$  for graphical convenience.

图 11.3: Auto-correlation  $C(z, z)$  和 cross-correlation  $C(z, r)$

的也就很慢。ReLU 就不存在这样的问题，但是 ReLU 的值域没有边界，因此容易出现一些数值问题，为了解决这个问题，我们将 ReLU 和 BN 一起使用，这样就很不存在数值问题了。将 tanh 改为 ReLU 之后，GRU 的公式如下：

$$\tilde{h}_t = \text{ReLU}(W_h x_t + U_h h_{t-1} + b_h) \quad (11.7)$$

### 3. BN

神经网络训练时，计算每一个 mini-batch 每层激活前输出的均值和方差，再利用均值和方差对激活前输出进行归一化能够解决所谓的 internal covariate shift 问题，这个就是传说中的 Batch Normalization。BN 既可以加快网络训练速度，又可以提高模型的效果。本文中只对前馈部分进行 BN 操作，因为其只对前馈神经网络进行操作，完全可以实现并行计算。其公式如下：

$$z_t = \sigma(\text{BN}(W_z x_t) + U_z h_{t-1}) \quad (11.8)$$

$$\tilde{h}_t = \text{ReLU}(\text{BN}(W_h x_t) + U_h h_{t-1}) \quad (11.9)$$

$$h_t = z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t \quad (11.10)$$

其中  $\text{BN}(\cdot)$  如公式 11.11 所示， $\mu_b$  和  $\sigma_b$  分别为当前 mini-batch 的均值和方差。

$$\text{BN}(a) = \gamma \odot \frac{a - \mu_b}{\sqrt{(\sigma_b)^2 + \epsilon}} + \beta \quad (11.11)$$

因为 BN 中已经包含了  $\beta$ ，因此之前公式中的偏置  $b_z$  和  $b_h$  就不再需要了。Li-GRU 将 ReLU 和 BN 结合起来，既利用了 ReLU 和 BN 两者的优点，同时还避免了 ReLU 的数值不稳定问题。

11.1.3 个人心得体会

综上所述，Li-GRU 通过减少了一个重置门来达到轻量级的效果，与此同时根据语音识别任务的特殊性，其认为语音序列任务中，对过往记忆清零对模型效果是有伤害的，而且重置门和更新门之间关联比较深，两个门显得冗余，因此其去掉了重置门。为了让网络更新更快，效果更好，以 Relu 函数代替了 tanh 作为候选状态的激活函数，同时以 BN 来解决 Relu 无边界的数值问题。BN 还可以帮助快速训练和提升效果。

附上一些音素的分类，如表11.1

表 11.1: 音素分类及示例

Phonetic Cat.	音素类别	Phone Lists
Vowels	元音	{iy, ih, eh, ae, ..., oy, aw, ow, er}
Liquids	流音	{l, r, y, w, el}
Nasals	鼻音	{en, m, n, ng}
Fricatives	擦音	{ch, jh, dh, z, v, f, th, s, sh, hh, ,zh}
Stops	塞音	{b, d, g, p, t, k, dx, cl, vcl, epi}



## 第 12 章 深度学习框架笔记

### 12.1 paddlepaddle

本处记录一些用 paddlepaddle 跑百度的 DeepSpeech2 的笔记和问题：

1. 基于一些原因，服务器 A 上的 GPU 暂时不可用，因此在服务器 B 上重新部署了百度基于 paddlepaddle 的DS2。A 和 B 有共享目录，因此这个 repository 是放在一个共享目录下的，但是在运行 `./run.sh logA checkpointAB checkpointAB vocab.txt vocab.txt vocab.txt vocab.txt`
2. DS2 提取的音频特征是 80 维的 fbank。

### 12.2 pytorch

### 12.3 tensorflow

## 参考文献

- [1] Graves. Alex. Supervised sequence labelling with recurrent neural networks. *Studies in Computational Intelligence*, 385:52–81, 2008.
- [2] Graves. Alex, Santiago Fernández, Faustino J. Gomez, and Jürgen Schmidhuber. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In *International Conference on Machine Learning*, 2006.
- [3] colah. Understanding lstm networks. <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>.
- [4] Jason Eisner. An interactive spreadsheet for teaching the forward-backward algorithm. In *Proceedings of the ACL-02 Workshop on Effective Tools and Methodologies for Teaching Natural Language Processing and Computational Linguistics - Volume 1*, ETMTNLP '02, pages 10–18, Stroudsburg, PA, USA, 2002. Association for Computational Linguistics.
- [5] L. Gillick and S. J. Cox. Some statistical issues in the comparison of speech recognition algorithms. In *in Proceedings ICASSP*, pages 532–535, 1989.
- [6] githubharald. demo for greedy-search. <https://github.com/githubharald/CTCDecoder>.
- [7] Alex Graves. Sequence transduction with recurrent neural networks. *CoRR*, abs/1211.3711, 2012.
- [8] Alex Graves, Abdel-rahman Mohamed, and Geoffrey E. Hinton. Speech recognition with deep recurrent neural networks. *CoRR*, abs/1303.5778, 2013.
- [9] Awni Hannun. Sequence modeling with ctc. *Distill*, 2017. <https://distill.pub/2017/ctc>.
- [10] Andrew L. Maas, Awni Y. Hannun, Daniel Jurafsky, and Andrew Y. Ng. First-pass large vocabulary continuous speech recognition using bi-directional recurrent dnns. *CoRR*, abs/1408.2873, 2014.
- [11] Timothy I Murphy. demo for prefix-beam-search. <https://github.com/corticph/prefix-beam-search>.
- [12] Andrew Ng. 统计学习方法. <http://cs229.stanford.edu/notes/cs229-notes8.pdf>.
- [13] D. S. Pallet, W. M. Fisher, and J. G. Fiscus. Tools for the analysis of benchmark speech recognition tests. In *International Conference on Acoustics*, 1990.
- [14] PilgrimHui. Ctc 自由输出的图解. <https://www.cnblogs.com/liaohuiqiang/p/9953978.html>.
- [15] Daniel Ramage. Hidden markov models fundamentals. <http://cs229.stanford.edu/section/cs229-hmm.pdf>.
- [16] RossYoung. 语音识别中的决策苏. <https://blog.csdn.net/luashui2009120/article/details/79545164>.
- [17] 农民小飞侠. 教程：connectionist temporal classification 详解补充. <https://blog.csdn.net/w5688414/article/details/77867786>.
- [18] 明无梦. 编辑距离 (edit distance) . <https://www.dreamxu.com/books/dsa/dp/edit-distance.html>.
- [19] 李航. 统计学习方法. 2012. <http://www.dgt-factory.com/uploads/2018/07/0725/%E7%BB%9F%E8%AE%A1%E5%AD%A6%E4%B9%A0%E6%96%B9%E6%B3%95.pdf>.
- [20] 菜鸟. awk 命令详解. <https://www.runoob.com/linux/linux-comm-awk.html>.
- [21] 菜鸟. sed 命令详解. <https://www.runoob.com/linux/linux-comm-sed.html>.