

班 级 071271

学 号 07127026

西安电子科技大学

本科毕业设计论文



题 目 科学数据中的序列分类

学 院 数学与统计学院

专 业 统计学

学生姓名 刘精昌

导师姓名 徐林莉

毕业设计（论文）诚信声明书

本人声明：本人所提交的毕业论文《科学数据中的序列分类》是本人在指导教师指导下独立研究、写作成果，论文中所引用他人的无论以何种方式发布的文字、研究成果，均在论文中加以说明；有关教师、同学和其他人员对本文本的写作、修订提出过并为我在论文中加以采纳的意见、建议，均已在我的致谢辞中加以说明并深致谢意。

本文和资料若有不实之处，本人承担一切相关责任。

论文作者：_____（签字） 时间：

指导教师已阅：_____（签字） 时间：

摘要

序列数据在日常生活中特别常见，比如心电图、语音序列等。对序列数据的分类对科学研究有很大的帮助，比如对心电图的分类能得知病人是否患病，语音序列的分类能对语音的来源进行识别。

不同于一般的数据，序列数据一般有很多元素，而且各个元素之间都有一定的依赖关系，因而使用常见的分类方法，比如支持向量机，并不可行。常用的方法是定义序列之间的距离，根据序列之间的远近，把序列归类到离其最近的那个序列。一般序列之间的距离定义为欧式距离，欧式距离在求解时，只会把相同索引的元素对齐。但是比较序列相近与否时，我们更为看重的是序列形态的相似，而相近形态的序列片段大多不在序列的相同位置。DTW 方法就能够有效解决欧式距离的这一缺点。可以定义反映两序列元素距离以及对齐关系的距离矩阵，DTW 方法的关键就是找到从距离矩阵左下角游走到右上角的路径，两序列之间的 DTW 距离通过简单的迭代式就可以得到。此外，本文还具体介绍了一些 DTW 方法的加速技巧。通过人造数据以及官方数据的实验，可以发现，在大多数的情况下，DTW 距离的分类精确度都要优于欧式距离。

DTW 方法考虑的是序列整体，本文还介绍了一种利用序列局部特性进行序列分类的方法——Shapelets。该方法需要找到最能体现序列特性的子序列 Shapelets，再根据待分类序列最具特征的局部片段和 Shapelets 的距离，对待分类序列进行分类。

关键词：序列分类 序列距离 DTW 方法 Shapelets

Abstract

Sequence data is so common, like ECG, speech sequence. Classifying sequence will be helpful for scientific research, we can know suspicious patient sick or not from ECG and recognize speech by classifying speech sequence.

Different from general data, there are so many elements in sequences and these elements are not independent, so some usual methods, such as SVM, will not workable. Generally, we need to define distance between two sequences, and classify the sequence to the nearest one. Commonly, the distance is defined as Euclidean distance, which can only align elements with the same index. But when compare two sequences, what we will pay more close attention to is the similarity between sequence shape, and in most cases, the similar shape has different indexes. DTW is a ingenious methods which can overcome this weakness. We can define the distance matrix which can reflect the distance and Alignment relationship, the key process is to find a warp path which wriggles from one corner to the opposite corner, and we can get the DTW distance by a simple iteration. Furthermore, we introduce several methods to accelerate DTW. By some experiment with artificial and official data, we can know, in most cases, DTW method is more accurate than Euclidean distance.

While DTW considers the whole sequence, there is a method which makes use of local features named Shapelet. By searching the most representative subsequence which is Shapelet, measuring the distance between shapelet and the subsequence of waiting-classify sequence, we can finish classification task.

Key Words: sequence classification distance DTW method Shapelets

Abstract

目录

第一章 绪论	1
1.1 序列以及序列分类	1
1.2 序列分类的主要方法	1
1.2.1 DTW 方法	2
1.2.2 Shapelets 方法	3
第二章 DTW 方法的计算以及加速方法	5
2.1 DTW 方法的计算	5
2.1.1 不带有权重的 DTW 方法	5
2.1.2 带有权重系数的 DTW 方法	7
2.2 DTW 方法的加速技巧	9
2.2.1 warp window	9
2.2.2 序列的分段平均表示法	10
2.2.3 FastDTW 方法	11
2.3 实验结果	13
2.3.1 实验 1	13
2.3.2 实验 2	14
2.3.3 实验 3	14
第三章 Shapelets 方法	17
3.1 Shapelets 简介以及定义	17
3.2 Shapelet 计算以及加速方法	19
3.2.1 Shapelet 的计算	19
3.2.2 Shapelet 计算的加速方法	20
第四章 总结与展望	23
致谢	25

参考文献

27

第一章 绪论

1.1 序列以及序列分类

现实生活中会有各种各样的序列。在语音识别领域，一段语音可以被看作一段序列；在经济领域，股票的走势可以被看作是序列图示；在分子生物学领域，蛋白质、DNA 也可以看成是一段序列；在网络管理中，客户端的登录活动也可以被记录成一段序列。而序列分类在现实生活中也有着广泛的应用。以语音库中的不同样本的语音为训练集，通过对需要测试的语音进行分类，可以识别出语音的来源；在医疗领域，通过对心电图序列的分类，可以得到这个疑似患者的患病信息；在信息检索的研究中，利用序列分类可将文件归类；通过蛋白质、DNA 序列的分类，有助于科学家探索蛋白质、DNA 的新功能^[14]。

一般说来，序列可以看成是事件的集合，而每一个事件可以由符号或者能比较大小的实数表示。比如一段 DNA 序列可以表示成 *ACCCCCGT*，一段简单的时间序列可以表示成 $\langle 0.2, 0.3, 0.5, 0.9, \dots \rangle$ ，两种序列的不同之处在于实数的大小是可以度量的，而一般而言符号之间的可度量性比较差，本文主要关注的是实数序列。对于符号序列，有一类的处理方法是借鉴实数序列，比如对于蛋白质序列和 DNA 序列，有一种基于序列对齐的距离度量方法^[4]。

对于一般的序列分类任务来说，训练集里面的序列都含有描述其类属性的标签。对于语音识别任务来说，语料库里的语音其标签是其发出者；对于文档归类任务来说，文档的标签即是其所属类别。定义 L 为标签集， S 为待分类序列集，序列分类的任务就是通过对含有标签信息的训练集 T 的训练，得到一个可将 S 映射到 L 的分类器 C ，对于进入 C 中的序列，都能够被贴上一个标签。可以表示成 $C: s \rightarrow l, s \in S, l \in L$ 。对于传统的分类任务来说，常用的分类器主要有 k NN，决策树、支持向量机^[12]。

1.2 序列分类的主要方法

对于序列分类问题，因为其灵活性，很多的学者从不同的角度提出了大量的方法^{[1][5][10][11][13][17]}。本文主要介绍其中两种比较典型的方法，一种是基于序列距离度量的 Dynamic Time Warping(DTW) 方法^{[1][2][5][7][8][10][13]}，另一种是基于特征选择的

shapelets 方法^[17]。其中对 DTW 方法做重点介绍。

1.2.1 DTW 方法

对于分类任务而言，一种广泛使用而且特别简单的方法是 k NN，使用 k NN 时，关键是度量测试数据和训练数据的亲近与否，也就是需要定义一种距离。对于一般的实数序列，最常用的距离度量是欧氏距离，对于序列 s 和 s' ，它们的欧式距离是：

$$\text{dist}(s, s') = \sqrt{\sum_{i=1}^L (s[i] - s'[i])^2} \quad (1-1)$$

也就是相同索引的元素的差平方相加再开方。同样地，还可以定义其他诸如街区距离，最大值距离，这里需要重点强调的是，这些距离都是将同索引值元素对齐。

但是对于序列数据而言，这些距离有很大局限性，主要体现在以下两点：

1. 以上的距离度量都要求两序列长度相同，而这一点在实际问题中很难得到满足。比如在基因组分析中，很难保证基因序列的长度相同。
2. 设想这样一种情况。在步态分析中，同一测试者的步速可能不同，或者在某时间段上存在着加速和减速。那么对于其两段步态序列，比较相似的步态之间可能会有一定的时间差，而上面的这些距离测度只会将同一时刻的步态相比较。也就是说，上面的这些距离测度不能反映出序列比较中的错位。

上述的第一个缺点容易克服，两不同长度的时间序列，可以通过插值的方法使得序列等长，从而应用传统的距离度量方法。比如，对于两序列较短的那个，相邻点的之间可插入一定数目的点，插入点的数值可以用这两点的均值代替，以此小技巧，使得短序列扩展到与长序列长度。同样地，分段平均处理能够减小长序列至短序列长度^[5]。但是对于对于大多数的序列而言，在描述序列之间的相似性时，更多的是考虑序列形态、走势方面的相似与否，而不会要求相似形态子序列的索引必须一致。传统的距离测度则只会比较同索比较。

DTW 方法能够很好地克服上述的两个缺点。

图1.1是 DTW 示意图，图中的两序列是某一长序列的两段，两子序列之间的索引有一定错位，示意图反映了有索引错位的两序列之间的对齐情况^[3]。

图中的虚线表示了两序列点的对齐情况，这里对齐的两点是指在距离求解时用来比较距离的两点，在传统的距离度量中，对齐的两点索引值相同。而从图中可以看出，对于 DTW 方法，点的对齐还参考了点邻近曲线的形态，而不仅仅局限于相同索引点，此

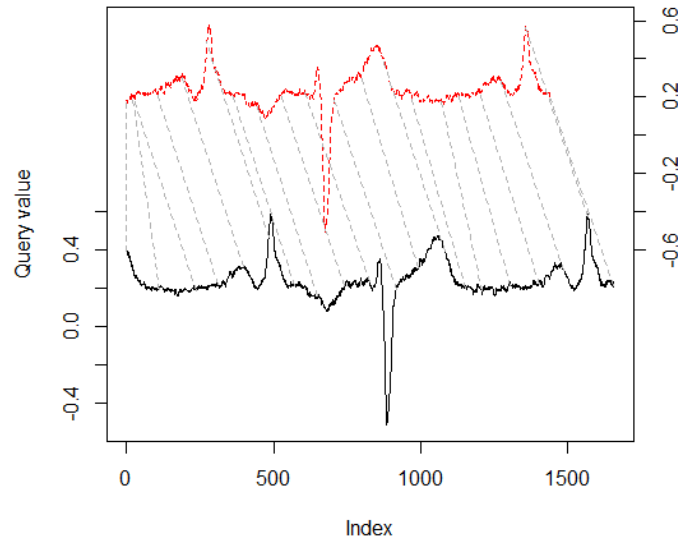


图 1.1: DTW 示意图

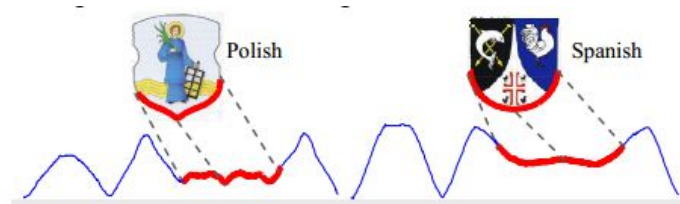


图 1.2: Shapelets 示意图

外序列上的某点可以和另一个序列中的多个点对齐。因而 DTW 方法还能够直接处理两序列长度不同的情况。DTW 的具体求解将在下文中介绍。

1.2.2 Shapelets 方法

Shapelets 方法最早是由 (Ye *et al* 2009) 提出, 作者经过大量实验, 验证了该方法的优秀性^[16]。

简单来说, Shapelet 是一小段序列, 选择待分类序列中最能代表该序列的那一部分, 比较其与 Shapelet 的相似度, 就可对其分类。最具特征子序列如图1.2所示。

图中的样本是两个徽章, 徽章的边缘可以转换成序列。对于这两个徽章而言, 徽章的下部的轮廓最能够反映该徽章的特征, 因而可以选择徽章下部轮廓对应的序列, 比较其与 Shapelet 的相似度。可以看出, 两徽章的下部明显不同, 因而可以由下部序列之间的差异, 利用分类器, 对徽章进行分类。在这里, 一般最常用的分类器是 k NN 分类器, 它简单而有效, 通常取 $k = 1$ ^[13]。如图1.3所示。

Shapelets 方法有着广泛的应用, 比如在植物学中, 利用叶柄与叶片的角度, 对植



图 1.3: 利用 Shapelets 进行 1NN 分类示意图

物进行分类；在考古学领域，通过对出土器件的特征部位的识别，利用决策树进行分类，能得到考古学发现^[6]。如图1.4所示。

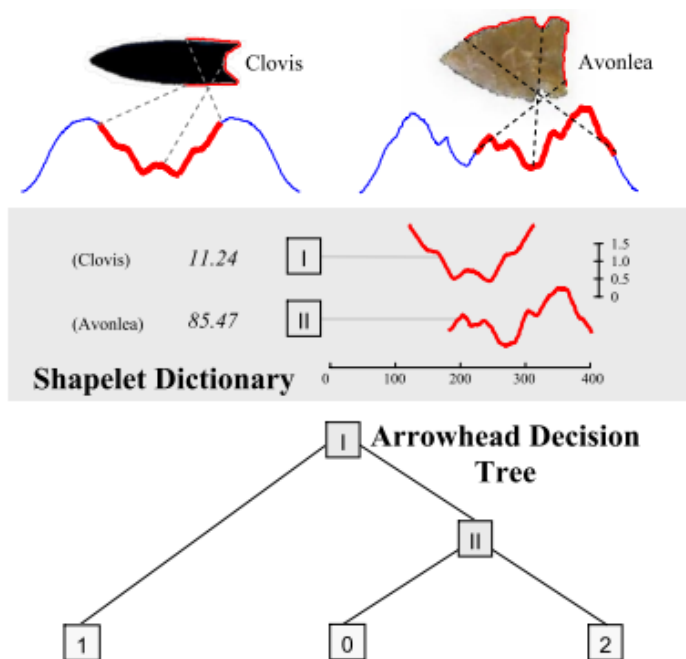


图 1.4: 基于 shapelets 的决策树方法对考古器件进行分类

第二章 DTW 方法的计算以及加速方法

2.1 DTW 方法的计算

2.1.1 不带有权重的 DTW 方法

首先定义两个序列 $Q = q_1, q_2, \dots, q_i, \dots, q_m$, $C = c_1, c_2, \dots, c_j, \dots, c_n$, 并定义 Q 序列的第 i 个元素与 C 序列的第 j 个元素之间的距离为它们的欧式距离, 即:

$$D(i^{th}, j^{th}) = d(q_i, c_j) = (q_i - c_j)^2 \quad (2-1)$$

DTW 方法中很重要的一点就是找到一条从两序列起点, 到两序列终点的, 能反映两序列对齐关系的 warping path。下面, 首先给出 DTW 方法的 warping path 的示意图, 如图2.1所示。示意图的左方和上方表示两个序列 Q 和 C , 序列之间的方格是能够反映序

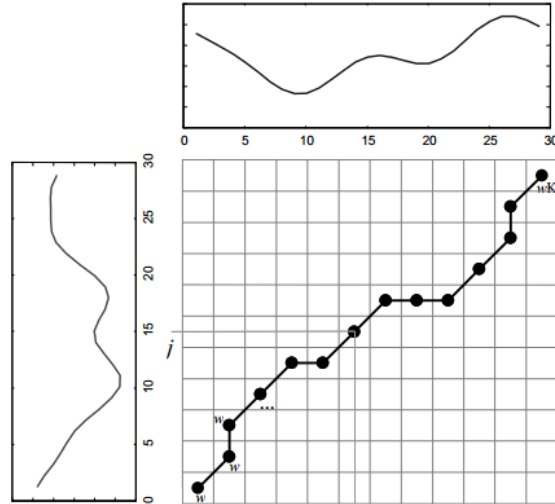


图 2.1: warping path 示意图

列对齐关系以及两序列对齐点之间距离信息的度量矩阵 M , M 的大小为 $m \times n$ 。 M 中的黑色实线即是 DTW 方法中的 warping path。 warping path 记为

$$W = w_1, w_2, \dots, w_k, \dots, w_K \quad \max(m, n) \leq K \leq m + n - 1 \quad (2-2)$$

w_k 表示路径上的第 k 个点, $w_k = (i, j)_k$, 表示在路径的第 k 步, q_i 点和 c_j 点进行了对齐。反映在距离矩阵上就是路径的第 k 步经过了坐标为 (i, j) 的方格。warping path 主要要满足下面三个条件:

1. **Boundary conditions:** $w_1 = (1, 1)$ 以及 $w_K = (m, n)$ 。这条约束的含义是两序列最开始是起点对齐, 最后是两序列的终点对齐。也就是说, warping path 是从距离矩阵的左下角蜿蜒到右上角。
2. **Continuity:** 给定路径的第 k 点 $w_k = (a, b)$ 以及路径的第 $k-1$ 点 $w_{k-1} = (a', b')$, 那么 $a - a' \leq 1$ 以及 $b - b' \leq 1$ 。也就是说, 路径游走的下一步点只能是当前点周围的那个点而不能跳跃。
3. **Monotonicity:** 给定路径的第 k 点 $w_k = (a, b)$ 以及路径的第 $k-1$ 点 $w_{k-1} = (a', b')$, 那么 $a - a' \geq 0$ 以及 $b - b' \geq 0$ 。表示路径只能向前游走。

若 $w_k = (i, j)$, 那么根据上面的三条约束, 路径下一步游走的点只能是 $w_{k+1} = (i+1, j+1)$ 或 $w_{k+1} = (i+1, j)$ 或 $w_{k+1} = (i, j+1)$ 。在这三个约束下, DTW 方法需要的优化目标是:

$$DTW(W) = \min_W \left\{ \sum_{k=1}^K d(w_{ki}, w_{kj}) \right\} \quad (2-3)$$

(w_{ki}, w_{kj}) 表示路径 W 上的第 k 个点, 说明第 k 步 q_i 和 c_j 点进行了对齐, $d(w_{ki}, w_{kj})$ 表示 q_i 和 c_j 点的距离。也就是说, 路径上的每点都对应着两序列对齐的点之间的距离, 现在 DTW 的目标就是使路径上所有的这样的距离之和最小。

求解这个最优化问题很简单, 只需要用贪心的方法, 使得每一次的迭代得到的距离和最小即可。定义这样的距离和为累计距离 $\gamma(i, j)$ 。直观意义是路径从距离矩阵的起始位置(左下方)蜿蜒到 (i, j) 方格的最小距离。路径蜿蜒到 (i, j) 方格的上一步是 $(i-1, j)$, $(i, j-1)$ 或 $(i-1, j-1)$, 因而迭代公式为:

$$\gamma(i, j) = d(q_i, c_j) + \min(\gamma(i-1, j), \gamma(i, j-1), \gamma(i-1, j-1)) \quad (2-4)$$

根据这个迭代公式, DTW 方法的伪代码如 Algorithm1所示。

伪代码中 DTW 矩阵是累计距离矩阵, 每个元素记录了 Q, C 两序列中对应元素之前的累计距离, 矩阵的最后一个元素即为最终 Q 和 C 的 DTW 距离。

根据上面的迭代式, 最终得到的仅仅是两序列最终的 DTW 距离, 除此之外, 在很多情况下, 还需要得到两序列的 DTW 对齐情况, 也就是 DTW 路径, 为了得到 DTW 路径, 采用回溯的方法。如图2.2。

Algorithm 1 Calculate DTW**Require:** $Q : \text{array}[1..m], C : \text{array}[1..n]$ **Ensure:** $DTW[m, n]$

```

1.  $DTW := [0..m, 0..n]$ 
2. for  $i := 0$  to  $m$  do
3.    $DTW[i, 0] := \text{inf}$ 
4. end for
5. for  $j := 0$  to  $n$  do
6.    $DTW[0, j] := \text{inf}$ 
7. end for
8.  $DTW[0, 0] := 0$ 
9.
10. for  $i := 1$  to  $m$  do
11.   for  $j := 1$  to  $n$  do
12.      $\text{cost} := d(Q[i], C[j])$ 
13.      $DTW[i, j] := \text{cost} + \min(DTW[i-1, j], DTW[i, j-1], DTW[i-1, j-1])$ 
14.   end for
15. end for

```

回溯从累计距离矩阵的右上角开始，回溯的上一步是当前方格左、下、斜下三个元素的最小值所在的方格，这样便能找到当前路径的上一步，通过不断这样地回溯，最终将到达路径的起始点，即得到 DTW 路径。

2.1.2 带有权重系数的 DTW 方法

对于 DTW 方法的路径走势，显然，路径总体斜向上游走，而不是呈阶梯形状游走更加合理。所以可以根据当前路径的走势，对路径的不同蜿蜒方向增加一个‘权重’，以增加路径斜向上蜿蜒的可能性。此时，2-4中的最优化条件调整为：

$$DTW(W) = \min_W \left\{ \frac{\sum_{k=1}^K d(w_{ki}, w_{kj}) \cdot \alpha_k}{\sum_{k=1}^K \alpha_k} \right\} \quad (2-5)$$

$\sum_{k=1}^K \alpha_k$ 是正则化项，可以被设定为常数 C ，因而2-5可以写成：

$$DTW(W) = \frac{1}{C} \min_W \left\{ \sum_{k=1}^K d(w_{ki}, w_{kj}) \cdot \alpha_k \right\} \quad (2-6)$$

α_k 表示第 k 步的权重系数，一般 α_k 有以下几种取值：

1. 对称形式： $\alpha_k = (i_k - i_{k-1}) + (j_k - j_{k-1})$ ，也就是说若路径下一步是沿着斜向上方向游走，那么权重系数 $\alpha_k = 2$ ，路径直向右或直向上游走， $\alpha = 1$ ，于是增加了路

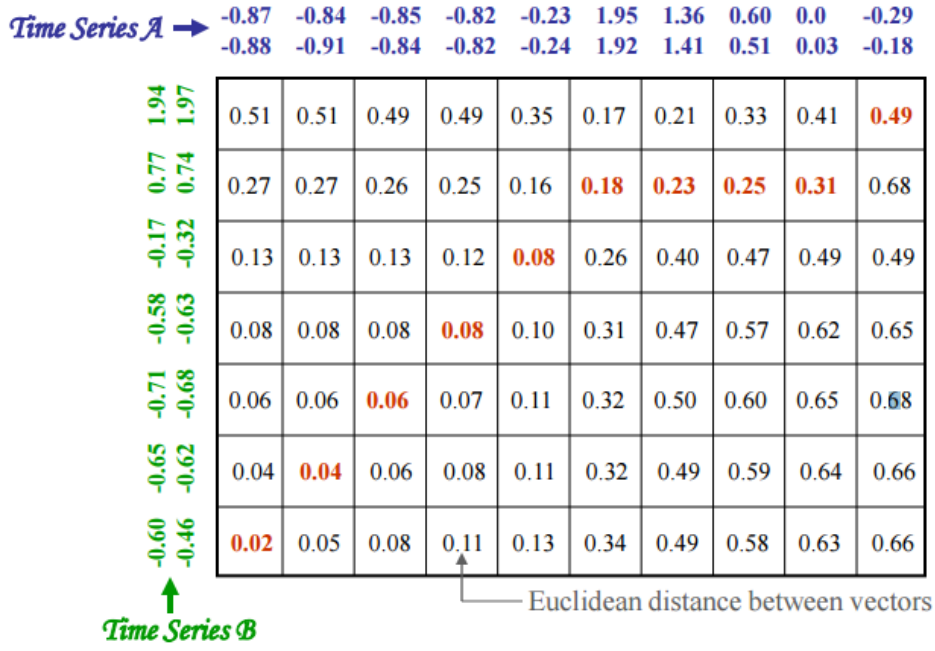


图 2.2: 回溯法得到 DTW 路径示意图

径斜向上游走的可能性。这种形式下， $C = m + n$ 。

2. 非对称形式:

- $\alpha = (i_k - i_{k-1})$, 这时 $C = m$ 。
- $\alpha = (j_k - j_{k-1})$, 这时 $C = n$ 。

也就是说，路径朝上或朝右以及斜向游走，权重系数为 2。

对于带有权重系数的 DTW 方法，迭代公式进行很小地调整即可。对于对称情形，迭代公式为：

$$\gamma(i, j) = d(q_i, c_j) + \min(\gamma(i-1, j), \gamma(i, j-1), \gamma(i-1, j-1) + d(q_i, c_j)) \quad (2-7)$$

两种非对称情形的迭代公式分别为：

$$\gamma(i, j) = d(q_i, c_j) + \min(\gamma(i-1, j) + d(q_i, c_j), \gamma(i, j-1), \gamma(i-1, j-1) + d(q_i, c_j)) \quad (2-8)$$

以及

$$\gamma(i, j) = d(q_i, c_j) + \min(\gamma(i-1, j), \gamma(i, j-1) + d(q_i, c_j), \gamma(i-1, j-1) + d(q_i, c_j)) \quad (2-9)$$

在实际的 DTW 方法的应用中,可以根据实际情况,对路径的每一步赋予特定的权重,而这只需要在迭代公式上做很小的修改。

2.2 DTW 方法的加速技巧

从上文中可以看到 DTW 方法的时间和空间复杂度都为 $O(n^2)$,该节主要关注 DTW 方法的加速技巧。首先提出两种加速小技巧,最后引入能够将 DTW 时间和空间复杂度降低一个数量级的方法——FastDTW^[10]。

2.2.1 warp window

虽然 DTW 方法能够对两序列中不同索引的元素进行对齐,但是直观上看,将两索引值相差过大的元素对齐在一起显然不合理。反映在距离矩阵上,也就是,路径不能偏离距离矩阵的对角线过大。基于此,可以定义 warp window,使得路径蜿蜒在 warp window 所限定的区域之中。warp window 如图2.3 所示。

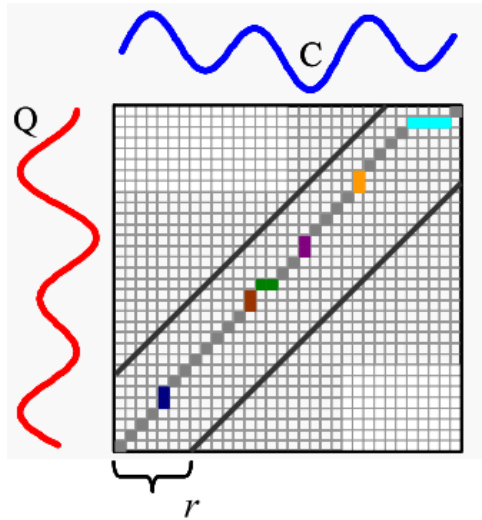


图 2.3: warp window 示意图

在不使用 warp window 情况下,需要计算整个整个距离矩阵,也就是说,需要填满整个距离矩阵的方格。引入 warp window,只需要填充 warp window 所限定的、对角线周围的方格即可。所以引入 warp window,在一定的程度上,能够降低计算量。此外,因为引入 warp window 后,路径的蜿蜒更符合直观实际,因而理论上说,引入 warp window 能够提高 DTW 方法的精确度。引入 warp window 后的 DTW 方法的伪代码如 Algorithm2 所示。

引入 warp window,与原始 DTW 方法的主要不同在第 11 行,通过控制 j 的取值范围来控制路径在距离矩阵中的蜿蜒范围。

Algorithm 2 Calculate DTW with warp window**Require:** $Q : \text{array}[1..m], C : \text{array}[1..n], w : \text{warp window}$ **Ensure:** $DTW[m, n]$

```

1.  $DTW := \text{array}[0..m, 0..n]$ 
2.  $w := \max(w, |m - n|)$ 
3. for  $i := 0$  to  $m$  do
4.   for  $j := 0$  to  $n$  do
5.      $DTW[i, j] := \text{inf}$ 
6.   end for
7. end for
8.  $DTW[0, 0] := 0$ 
9.
10. for  $i := 1$  to  $n$  do
11.   for  $j := \max(1, i - w)$  to  $\min(m, i + w)$  do
12.      $\text{cost} := d(Q[i], C[j])$ 
13.      $DTW[i, j] := \text{cost} + \min(DTW[i - 1, j], DTW[i, j - 1], DTW[i - 1, j - 1])$ 
14.   end for
15. end for

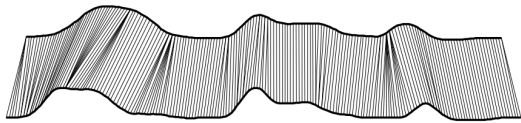
```

2.2.2 序列的分段平均表示法

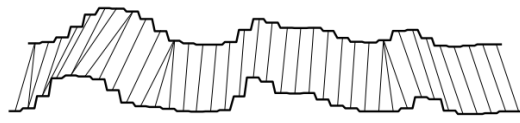
本节中的加快 DTW 计算的方法是从减小序列长度的角度考虑的，序列的分段平均表示法^[5]是一种常用的，能够有效减小序列长度的方法。如图2.4所示。

给定序列 $X = x_1 x_2 \cdots x_n$ ，它的分段平均表示下的序列记为 $\bar{X} = \bar{x}_1 \bar{x}_2 \cdots \bar{x}_N$ ，这里假定 n 能被 N 整除。其中 \bar{X} 的第 i 个元素是按照下面的计算式得到的：

$$\bar{x}_i = \frac{N}{n} \sum_{j=\frac{n}{N}(i-1)+1}^{\frac{n}{N}i} x_j \quad (2-10)$$



(a) 原始序列



(b) 分段平均处理后的新序列

图 2.4: 序列分段平均处理示意图

也就是说 \bar{X} 序列中的元素是对 X 序列中对应范围的元素求均值得到的。经过序列分段平均处理后，再对新的较短的序列应用 DTW 方法。

在某些情况下，原始序列中的某些元素是序列中的噪声，经过序列分段平均处理，能从某种程度上消除序列噪声元素的影响，提取出序列特性。从这个角度看，对于某些序列，应用较为合适的序列分段平均处理，能够很大地减小序列长度，从而极大的加快 DTW 的计算，并且可以增加 DTW 方法的精度。如图2.5所示，原始序列是长度为 9900000 的时间序列，序列的扰动很大，分段平均处理后的序列长度为 495，可以看出，

分段平均处理能很好地提取出原始序列中的有价值信息。

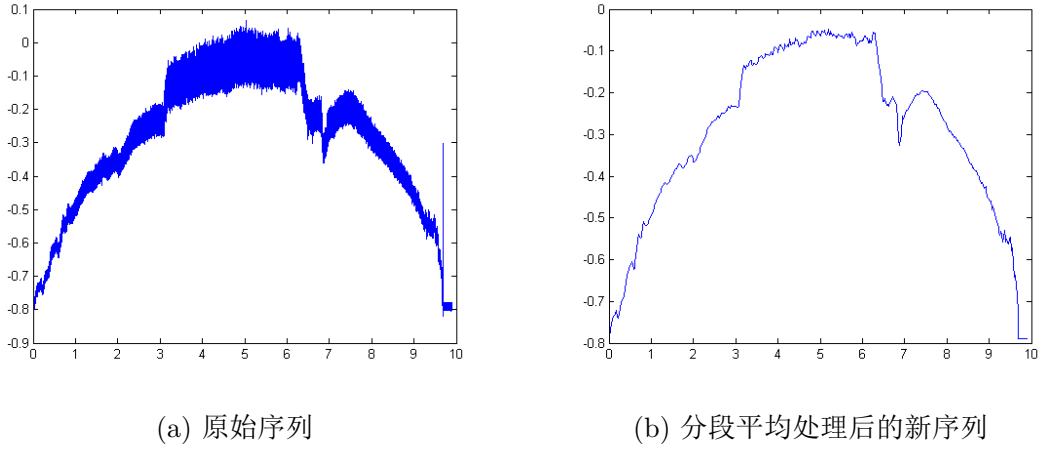


图 2.5: 序列分段平均处理能提取序列的信息

2.2.3 FastDTW 方法

FastDTW 方法简介

上面介绍的两种方法只能在 DTW 的实际应用中，经验地加快 DTW 的计算，本节引入一种新的能够从算法的角度有效地把 DTW 方法的时间、空间复杂度从 $O(n^2)$ 降低到 $O(n)$ 的方法——FastDTW^[10]。FastDTW 的思想是首先降低距离矩阵的分辨率，在较低分辨率下，得到 DTW 路径，然后投影较低分辨率的路径所经过的方格，得到在较高分辨率下路径所能够走的方格的范围。算法主要包括以下三个关键步骤：

1. **粗化**：首先利用上面所提到的序列分段平均表示法，得到较短的序列。假设较短序列长度为原序列长度的 $\frac{N}{n}$ ($n \geq N$)，那么这些短序列构成的距离矩阵的宽度为原距离矩阵宽度的 $\frac{N}{n}$ ，也就是说，现在得到 $\frac{1}{n}$ 分辨率的距离矩阵。
2. **投影**：即找到低分辨率下的 DTW 路径。
3. **细化**：由低分辨率下 DTW 路径所经过的距离矩阵所在的方格，找到较高分辨率下这些方格的对应，以得到高分辨率下 DTW 路径的所能蜿蜒的范围。

FastDTW 方法的示意图如2.6所示。

以图中的 $\frac{1}{2}$ 分辨率到 $\frac{1}{1}$ 分辨率为例。 $\frac{1}{2}$ 低分辨率是 $\frac{1}{1}$ 高分辨率下，通过对原始序列进行分段平均处理后得到的。在 $\frac{1}{2}$ 分辨率下，利用原始 DTW 方法可以得到如图所示的路径，路径周围黑色的方块是该分辨率下路径所走过的区域。在将分辨率调高到 $\frac{1}{1}$ 分辨率后， $\frac{1}{2}$ 分辨率下得到的黑色方块区域也相应地投影到 $\frac{1}{1}$ 分辨率对应区域，即 $\frac{1}{1}$ 分

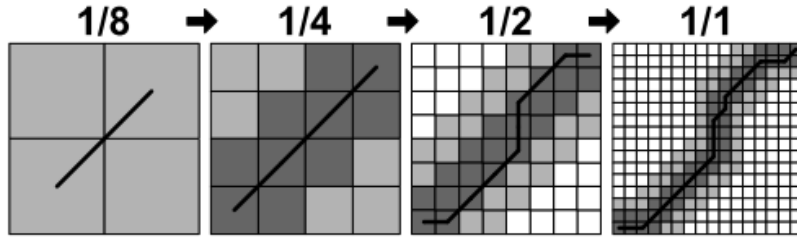


图 2.6: FastDTW 示意图

分辨率中的黑色区域中。但是在实际的操作中,为了使得得到的路径更为精确,这时会人为地设定一个参数:半径 r ,它是一个倍率,表示将黑色区域扩大了多少倍。 $\frac{1}{1}$ 分辨率下的黑色区域外围的灰色区域即是加入半径参数后对黑色区域的扩大所得,在该图所示的情况下, $r = 1$,在 $\frac{1}{1}$ 分辨率的距离矩阵中,路径蜿蜒区域被限定在黑色和灰色区域中。

在实际的 FastDTW 操作中,首先会把分辨率调低到一个合理的水平,然后逐步调高分辨率,求出在每种分辨率下的路径,最终将得到 $\frac{1}{1}$ 分辨率下的路径。FastDTW 对于每一种分辨率取最优路径,但是最终得到的不一定是和原始 DTW 方法相同的最优路径,因为从低分辨率路径得到的高分辨率情形下的路径限制方格,只是高分辨率中,路径最有可能经过的方格,真实的最优路径可能会‘突破’这些限制方格,所以 FastDTW 方法最终所得到的路径可能会和原始 DTW 路径有所不同。增大半径 r ,会降低路径约束,使得 FastDTW 路径和原始 DTW 路径更为接近,但增加半径会使得需要填充的距离矩阵的方格更多,从而增加了计算量。

FastDTW 方法的时间、空间复杂度分析

首先计算 FastDTW 方法需要计算距离矩阵的元素的总计数目,直观上看,也就是需要填充多少距离矩阵的方格。首先,考虑从 $\frac{1}{2}$ 分辨率到 $\frac{1}{1}$ 分辨率需要填充方格数目最多的情况,如图2.7所示。

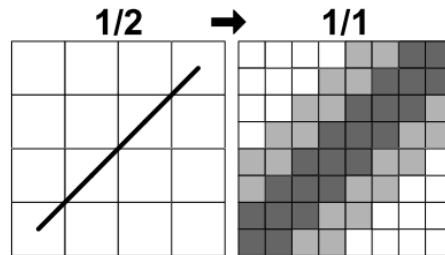


图 2.7: 最大填充方格示意图

$\frac{1}{1}$ 分辨率的图示中,黑色和灰色的区域都需要填充,因而在 $\frac{1}{1}$ 分辨率情形下,需要填充的最大方格数目上界为: $3N + 2 \cdot (2Nr) = N(4r + 3)$ 。因为分辨率会有 $\frac{1}{1}, \frac{1}{2}, \frac{1}{2^2}, \frac{1}{2^3}, \dots$

这些，因而在所有分辨率下总共需要填充的方格数目为：

$$N(4r+3) + \frac{N}{2}N(4r+3) + \frac{N}{2^2}N(4r+3) + \frac{N}{2^3}N(4r+3) \cdots = 2N(4r+3) \quad (2-11)$$

时间复杂度由下面三部分组成：

1. 需要计算的方格数目。刚刚已经计算得到，是 $2N(4r+3)$ 。
2. 计算得到不同的分辨率。对于一段序列（长度为 N ）来说，不同分辨率下需要的时间复杂度之和为： $N + \frac{N}{2} + \frac{N}{2^2} + \frac{N}{2^3} + \cdots = 2N$ ，因为有两段序列，因而这部分的时间复杂度合计为 $4N$ 。
3. 回溯得到不同分辨率下的路径。 $\frac{1}{r}$ 分辨率下，路径长度最多为 $2N$ ，也就是说最多需要回溯比较 $2N$ 次。所以所有分辨率下，回溯以得到路径的时间复杂度是： $2N + \frac{2N}{2} + \frac{2N}{2^2} + \frac{2N}{2^3} + \cdots = 4N$ 。

将这三项加起来，可得 FastDTW 方法的时间复杂度为： $N(8r+14)$ 。

可以看到，FastDTW 是线性时间复杂度。

再来估计空间复杂度。空间复杂度由下面三部分组成：

1. 存储不同分辨率下的序列。因为不同分辨率下的序列长度分别为 $N, \frac{N}{2}, \frac{N}{2^2}, \frac{N}{2^3}, \cdots$ ，而进行距离比较的序列有两条，再考虑到所有分辨率，空间复杂度总计为 $4N$ 。
2. 存储不同分辨率下的距离矩阵。前面已经估计到，总共需要填充的距离矩阵方格数目为 $N(4r+3)$ 。
3. 存储路径。前面也已经说明，存储路径的空间复杂度的上限为 $4N$ 。

因而 FastDTW 方法总的空间复杂度为： $N(4r+11)$ 。可以看出，同样地，FastDTW 方法也为线性空间复杂度。当序列的长度比较小，也就是 N 比较小时， $N(8r+14)$ 以及 $N(4r+11)$ 的复杂度与 N^2 相比并没有优势， N 很大时， $N(8r+14)$ 以及 $N(4r+11)$ 将远远小于 N^2 ，因而，对于长序列，FastDTW 的时间、空间复杂度将远远优于原始 DTW 方法，当然，这牺牲了一定的精确度。

2.3 实验结果

2.3.1 实验 1

首先简单地使用人造数据比较欧几里得距离和 DTW 距离的分类精确度。实验有四类，每一类中训练集、测试集均都有四个序列，序列的长度均为 100，相邻类的两序

列都有间隔参数 $interval$ ，序列是在相邻类间隔为 $interval$ 的常数序列中增加随机噪声而产生的，现通过 1NN 的方法对测试集中的序列分类，即将测试集中的序列归类到距离最近的那个序列，然后和其真实的 $label$ 相比较，计算得到测试集的分类精确度。最终的分类精确度是对上面的数据集分类 100 次得到的平均精确度。

将间隔参数 $interval$ 从 0 变化到 2，可以得到相应的欧几里得距离和 DTW 距离的分类精确度曲线，如图 2.8 所示。实验所得到的精确度曲线中，大部分情况下，DTW

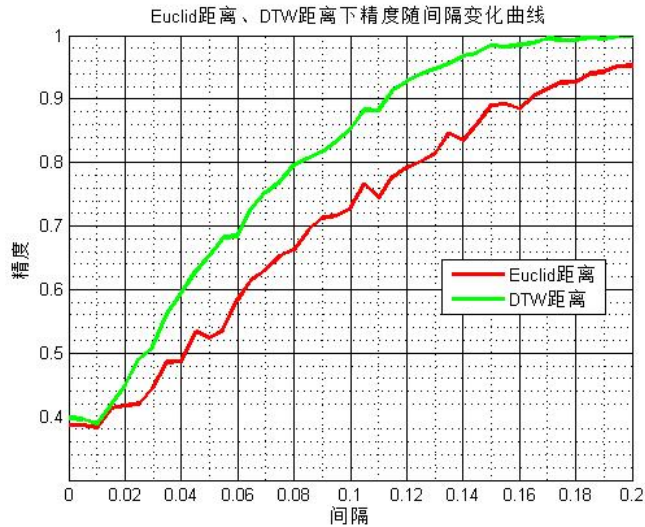


图 2.8: Euclid、DTW 距离下分类精确度随间隔变化曲线

距离的分类精度曲线都在欧几里得距离的分类精度曲线上，这说明了 DTW 距离相对于传统欧几里得距离的优越性。

2.3.2 实验 2

实验 1 是使用人造数据进行实验，实验 2 把数据换为权威的官方数据。UCR Time Series Classification Archive 是加州大学河边分校相关研究人员收集的时间序列数据集^[15]，现采用其中的部分数据集，使用 DTW 以及欧几里得距离进行分类，得到结果如表所示。

从实验结果中的两个表格可以看出，在大多数情况下，DTW 方法的分类精确度都要高于传统的欧几里得距离。

2.3.3 实验 3

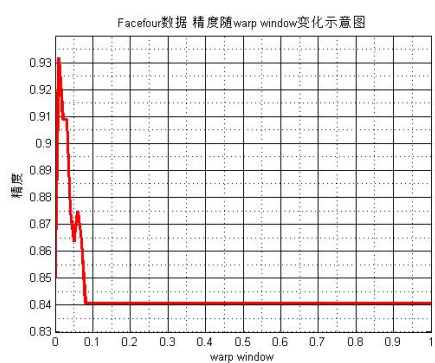
上文已经说过，通过定义 warp window，不但能够减小计算量，而且能够在某种程度上增大分类精度。现在对于部分 UCR 数据集以及 $interval = 0.875$ 时的人造数据集，从 0 到 1 变化 warp window 的比例，得到各种比例下的分类精度，得到分类精度曲线，

表 2.1 UCR 数据集实验结果

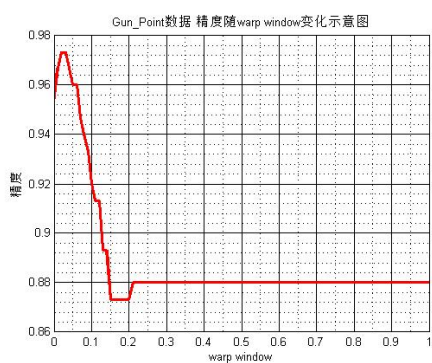
name	Gun_Point	Plane	StrawBerry
lasses	2	7	2
training set size	50	105	370
test set size	150	105	613
sequence length	150	144	235
error rate (Euclid)	0.086667	0.038095	0.06199
error rate (DTW)	0.12	0	0.066884

name	Computers	Trace	FaceFour	WordsSynonyms
lasses	2	4	4	25
training set size	250	100	24	267
test set size	250	100	88	638
sequence length	720	275	350	270
error rate (Euclid)	0.424	0.24	0.21591	0.38245
error rate (DTW)	0.332	0.01	0.15909	0.32445

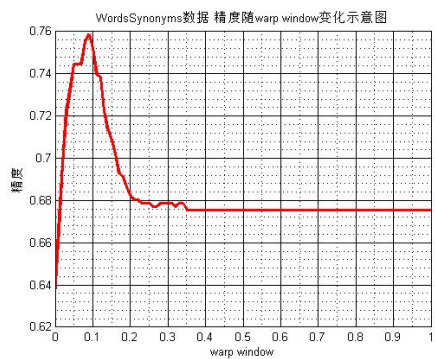
如图2.9 所示。从分类精度随 warp window 变化曲线可以看出, 较小的 warp window (大部分情况下小于 10%) 会提升分类精度。事实上, 很多文献^{[15][9]} 也说明了这点。



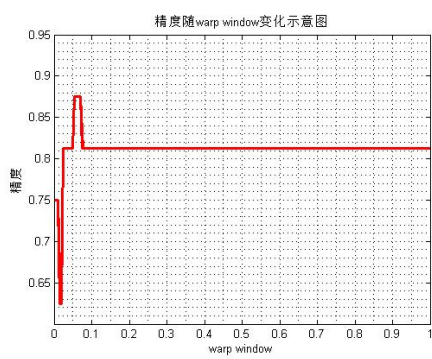
(a) facefour 数据



(b) gun_point 数据



(c) WordsSunonyms 数据



(d) interval 为 0.875 时的人造数据

图 2.9: 不同数据集分类精度随 warp window 变化曲线

第三章 Shapelets 方法

上面介绍的 DTW 方法是从构造衡量两序列距离的角度分析，而本章通过序列最具特征子序列来完成分类。

3.1 Shapelets 简介以及定义

在日常的生活中，人们在对一个物体进行分类的时候，会根据这一类物体所具有的共同特征，比如我们想要区分一支笔是铅笔还是钢笔，我们只需根据笔尖的种类即可，笔尖带有圆珠，那么这支笔就是圆珠笔，笔尖中心有缝隙，那么这支笔就是钢笔。也就是说，对笔分类时，我们仅仅需要利用最能代表笔特征的部位——笔头即可。同样地，对于某些序列数据来说，我们也可以利用这种共同特征来对序列数据分类，这种体现共同特征的，序列集中的某序列的子序列就是 Shapelets^[17]，通过衡量待分类序列的最具特征子序列和 Shapelets 的远近，即可对待分类序列进行分类。

一般而言，Shapelets 方法有以下三点好处：

- 得到可解释的结果：应用 Shapelets 方法的时候，结果中将包含一类序列的共同特征，这些共同特征将有助于对分类结果进行解释。
- 对于一些数据而言分类的结果更为精确。Shapelets 方法和大多数方法的不同之处在于，Shapelets 方法考虑的是序列的局部特征，而其他方法多考虑序列整体，这些方法比较容易受到序列中的无用部分以及噪声、失真等影响。例如把不具区分度的笔杆也考虑进分类中，那么分类结果必不如只考虑笔尖精确。
- Shapelets 方法的时间复杂度比较低。在已知 Shapelets 的情况下，分类的时间复杂度为 $O(ml)$ ， m 是待分类序列的长度， l 是 Shapelets 的长度。因为 Shapelets 是序列片段，故在一般情况下， l 都比较小。

一般利用信息增益来定义 Shapelets，给定顺序序列 $T = t_1, t_2, \dots, t_m$ ，下面来逐步地给出 Shapelets 的定义。

Definition 1: 子序列. 序列 T 的长度为 l 的子序列 S 是由 T 中 l 个连续元素所组成的， $S = t_p, t_{p+1}, \dots, t_{p+l-1}$ $p + l - 1 \leq m$ 。

Definition 2: 长度为 l 的子序列集合. 定义 $S_p^l = t_p, t_{p+1}, \dots, t_{p+l-1}$ 是 T 的起始点为 t_p , 长度为 l 的子序列. 那么 T 的所有长度为 l 的子序列集合是 $\mathbf{S}_T^{[l]} = \{S_p^l, 1 \leq p \leq m - l + 1\}$.

Definition 3: 子序列 S 和长序列 T 之间的距离 $SubSeqDist(S, T)$. 定义 $SubSeqDist(T, S) = \min(Dist(S, S'), S' \in \mathbf{S}_T^{[|S|]})$, 即 S 和 T 之间的距离是 S 和 T 中所有等长子序列之间距离的最小值, 如图3.1所示。

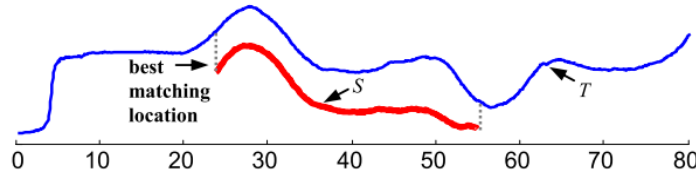


图 3.1: 子序列和长序列距离示意图

之所以把 S 和 T 之间的距离定义成 S 和 T 的子序列之间的距离, 是因为 Shapelets 是对序列子序列比较的方法。

下文中, 本文将利用信息增益来定义、衡量 Shapelets 方法分类的好坏。

Definition 4: 熵. 假定序列集合 \mathbf{D} 由两类 A, B 组成, 比例分别为 $p(A)$ 和 $p(B)$, 那么 \mathbf{D} 的熵的定义是 $I(\mathbf{D}) = -p(A)\log(p(A)) - p(B)\log(p(B))$ 。

求得熵后, 那么 Shapelets 分类前后的信息增益即为熵的减少量。

Definition 5: 信息增益. 特定的划分策略 sp 将数据集 \mathbf{D} 划分成两个子序列集 \mathbf{D}_1 、 \mathbf{D}_2 , 那么信息增益定义为:

$$Gain(sp) = I(\mathbf{D}) - \hat{I}(\mathbf{D}) \quad (3-1)$$

也就是:

$$Gain(sp) = I(\mathbf{D}) - f(\mathbf{D}_1)I(\mathbf{D}_1) + f(\mathbf{D}_2)I(\mathbf{D}_2) \quad (3-2)$$

如图3.2所示, 序列集 \mathbf{D} 中有 10 个序列, 包括两类, 分别用红色和蓝色表示, 每类的序列数目分别为 4 和 6, 计算 Shapelets 候选序列和 \mathbf{D} 中序列距离, 按照大小关系排列在坐标轴上, 在如图所示的最优分割点的情况下, 信息增益为:

$$\left[-\frac{6}{10} \log\left(\frac{6}{10}\right) - \frac{4}{10} \log\left(\frac{4}{10}\right) \right] - \left[\frac{5}{10} \left[-\frac{5}{5} \log\left(\frac{5}{5}\right) \right] + \frac{5}{10} \left[-\frac{4}{5} \log\left(\frac{4}{5}\right) - \frac{1}{5} \log\left(\frac{1}{5}\right) \right] \right] = 0.482 \quad (3-3)$$

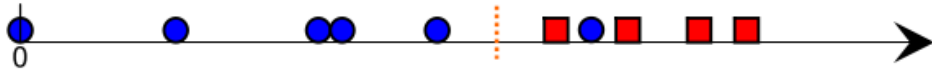


图 3.2: 信息增益计算的示意图

需要说明的是 Shapelet 是序列集 D 中所有子序列中的某个, 那么 Shapelets 候选集就是 D 的所有子序列。

Definition 6: 最优分隔点 (OSP). 数据集 D 中包含两类序列 A 和 B 。对于 Shapelets 候选序列 S , 根据上面的距离的定义可以得到 S 与 D 中所有序列之间的距离, 选择距离阈值 d_{th} , 可以将 D 中的序列分为 D_1 和 D_2 两类, D_1 中的序列 $T_{1,i}$ 满足 $SubSeqDist(S, T_{1,i}) \leq d_{th}$, D_2 中的序列 $T_{2,i}$ 满足 $SubSeqDist(S, T_{2,i}) \geq d_{th}$ 。最优分隔点是所有分割点 d'_{th} 中, 使得信息增益最大的那个分割点:

$$Gain(S, d_{OPS(D,S)}) \geq Gain(S, d'_{th}) \quad (3-4)$$

有了以上这些定义, 最后将给出 Shapelets 的定义。

Definition 7: Shapelets. 给定含有两类序列的序列集合 D , $Shapelet(D)$ 是 D 的子序列, 满足对于 D 的所有子序列 S :

$$Gain(Shapelet(D), d_{OPS(D, Shapelet(D))}) \geq Gain(S, d_{OPS(D,S)}) \quad (3-5)$$

也就是说, 在最优分割的情况下, 利用 Shapelet 对序列集合进行分割, 信息增益最大。

3.2 Shapelet 计算以及加速方法

3.2.1 Shapelet 的计算

利用 Shapelets 进行序列分类, 首先需要找到序列集合的 Shapelets, 算法如下。

伪代码的第一行是构造 Shapelets 候选集, 该步中, Shapelets 候选集是长度介于 $MINLEN$, $MAXLEN$ 之间的所有子序列。然后遍历候选集中的所有候选序列, 能够在最优分割情况下使得信息增益最大的候选序列即是序列集合的 Shapelet。

Algorithm 3 Brute force algorithm for finding shapelet

Require: dataset $D, MAXLEN, MINLEN$ **Ensure:** bsf_shapelet

1. candidates := GenerateCandidates($D, MAXLEN, MINLEN$)
 2. bsf_gain := 0
 3. **for** S in candidates **do**
 4. gain := CheckCandidate(D, S)
 5. **if** gain > bsf_gain **then**
 6. bsf_gain := gain
 7. bsf_shapelet := S
 8. **end if**
 9. **end for**
-

3.2.2 Shapelet 计算的加速方法

在寻找 Shapelet 的过程中，需要遍历序列集所有的子序列，那么候选序列的总数目将特别地大，而且在计算候选序列和原长序列的距离时，需要遍历原长序列中和候选序列等长的子序列，寻找最优分隔点也需要很大的计算量，所以寻找 Shapelet 的计算量特别大，下面介绍能够减小计算量的两种方法。

SubSeqDist 计算的提早舍弃

对于 Shapelets 候选序列 S 和原始序列 T ，假设当前是取出 T 的子序列 S' 与 S 求距离。计算 S 与 S' 之间的距离，是按索引从小到大逐步计算的，我们可以把之间计算得到的最小距离给记录下来，当索引值推进到 S 和 S' 的某一步时，比较当前的距离和之前已记录的最小距离之间的大小，如果当前的距离比记录的最小距离大，那么就没有必要再拿 S' 和 S 比较下去，那么选择 T 的下一个子序列和 S 求距离。如图3.3所示。

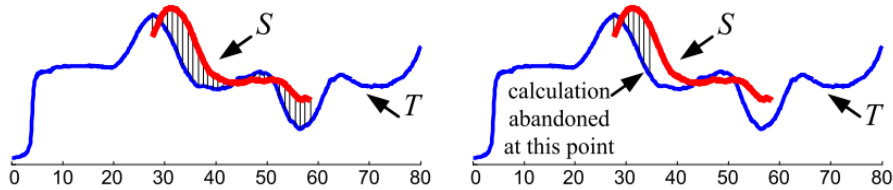


图 3.3: 计算候选序列和原始序列距离时的提早舍弃

计算熵的时候的提早舍弃

还是接着图3.2进行说明。对于另一个 Shapelets 候选序列 \bar{S} ，假设当前已经计算了 \bar{S} 和 D 中的 5 个序列（2 个红色序列，3 个蓝色序列）之间的距离，并按照大小顺序排列好，如图3.4所示。

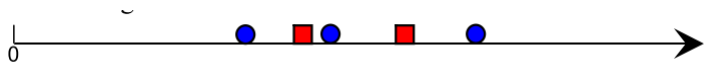


图 3.4: 候选序列与序列集中 5 个序列距离大小排列情况

对于下面的 5 个序列，当 3 个蓝色序列在最左边、2 个红色序列在最右边或者 3 个蓝色序列在最右边、2 个红色序列在最左边时，能获得的信息增益最大，如图3.5 所示，可以计算这种最大信息增益与之前记录的最大信息增益的大小关系，如果小于之前记录的最大信息增益，那么说明当前的 Shapelets 候选序列肯定不是最终的 Shapelet，所以应该选择 Shapelets 候选集中的下一个候选序列进行计算。

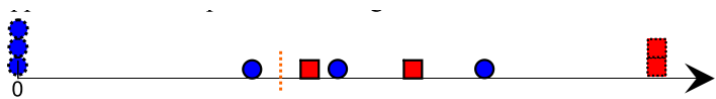


图 3.5: 使得熵最大的极端情况示意图

在得到序列集的 Shapelet 后，比较待分类序列和 Shapelet 之间的距离，根据距离远近，对待分类序列进行分类。

第四章 总结与展望

序列分类中, DTW 是一种经典的方法, 它能够很好地克服两序列中相似部分的索引值不同的问题, 能有效地比较两序列的相似性, 其分类精确度也得到广泛地认可。DTW 方法最大的一个缺点是其较高的时间、空间复杂度, 一些小技巧, 比如对原始序列的分段平均处理^[5], 选择较小的 warp window 参数等能有效地降低 DTW 方法的时间、空间复杂度。这些方法还能够提高 DTW 方法的分类精确度。原始序列的分段平均处理能够有效地提取出原始长序列的有效信息, 反映出原始序列的形态, 克服原始序列中噪声的干扰。较小的 warp window 更符合序列分类的实际, 也就是序列相似的形态之间不会间隔太大的位置, 因而这些方法能够带来分类精确度的提升。[Salvador et.al, 2007] 从算法的角度提出了一种能够有效地将 DTW 的时间、空间复杂度从 $O(n^2)$ 降低到 $O(n)$ 的改进方法——FastDTW, FastDTW 最终得到的结果是原始 DTW 方法结果的较好近似, 可以通过调整 FastDTW 方法的参数, 以平衡 FastDTW 近似结果好坏和时空复杂度。

近些年来, [Ye et.al, 2009] 从另外一个角度提出了序列分类的新方法——Shapelet。DTW 方法比较的是序列整体的相似性, 而 Shapelet 方法只选择序列中最具代表性的子序列, 比较它们之间的相似性。对于一些依赖局部特征进行区分的序列数据, 比如树叶边缘序列, Shapelet 方法将得到较好的分类结果。

在很多实际场景, 比如根据用户登陆情况判别用户异常行为, 越早判别出序列标签越好, 也就是根据尽可能短的前缀序列进行判别, 这是当前序列分类的一个研究方向——早期分类^[14]。实际得到的序列可能是有标签数据, 也可能是无标签数据, 半监督序列分类就是研究怎么结合这两类数据进行分类。实际的序列数据可能含有多个标签, 对这种序列的分类属于多标签序列分类的研究范畴。

致谢

这次的毕业论文设计总结是在我的指导老师徐林莉老师亲切关怀和悉心指导下完成的。从毕业设计选题到设计完成，徐老师给予了我耐心指导与细心关怀，在此表示感谢。徐老师有严肃的科学态度，严谨的治学精神和精益求精的工作作风，这些都是我所需要学习的！感谢研究组里提供的优良的科研资源，这让我享受到了难忘的学习时光，感谢师兄们的指导。另外，也特别感谢崔元顺同学提供的 latex 模板。

参考文献

- [1] Ghazi Al-Naymat, Sanjay Chawla, and Javid Taheri. SparseDTW: A novel approach to speed up dynamic time warping. *Conferences in Research and Practice in Information Technology Series*, 2009, 101(December 2003): 117–127.
- [2] Gustavo E. a. P. a. Batista, Xiaoyue Wang, and Eamonn J Keogh. A Complexity-Invariant Distance Measure for Time Series. *SIAM International Conference on Data Mining*, 2011, pages 699–710.
- [3] Toni Giorgino. Computing and Visualizing Dynamic Time Warping Alignments in R: The dtw Package. *Journal Of Statistical Software*, 2009, 31(7): 1–24.
- [4] László Kaján, Attila Kertész-Farkas, Dino Franklin, Neli Ivanova, András Kocsor, and Sándor Pongor. Application of a simple likelihood ratio approximant to protein sequence classification. *Bioinformatics*, 2006, 22(23): 2865–2869.
- [5] E. J Keogh and M. J Pazzani. Scaling up dynamic time warping for datamining applications. *Knowledge discovery and data mining*, 2000, In 6th ACM: 285–289.
- [6] Eamonn Keogh. Ucr computational anthropology site. <http://www.cs.ucr.edu/~eamonn/anthropology/>.
- [7] CHRISTINA LESLIE, ELEAZAR ESKIN, and WILLIAM STAFFORD NOBLE. the Spectrum Kernel: a String Kernel for Svm Protein Classification. *Biocomputing 2002*, 2001, pages 564–575.
- [8] Jessica Lin, Eamonn Keogh, Li Wei, and Stefano Lonardi. Experiencing SAX: A novel symbolic representation of time series. *Data Mining and Knowledge Discovery*, 2007, 15(2): 107–144.
- [9] Chotirat Ann Ratanamahatana and Eamonn Keogh. Three Myths about Dynamic Time Warping Data Mining. *Proceedings of the 2005 SIAM International Conference on Data Mining*, 2005, page 5.
- [10] Stan Salvador and Philip Chan. FastDTW : Toward Accurate Dynamic Time Warping in Linear Time and Space. *Intelligent Data Analysis*, 2007, 11: 561–580.
- [11] Sebastian Thrun, T O M Mitchell, and Kamal Nigam. Text Classification from Labeled and Unlabeled Documents using EM. 2000, 34: 103–134.

- [12] Xindong Wu, Vipin Kumar, Quinlan J. Ross, Joydeep Ghosh, Qiang Yang, Hiroshi Motoda, Geoffrey J. McLachlan, Angus Ng, Bing Liu, Philip S. Yu, Zhi Hua Zhou, Michael Steinbach, David J. Hand, and Dan Steinberg. *Top 10 algorithms in data mining*. volume 14. 2008.
- [13] Xiaopeng Xi, Eamonn Keogh, Christian Shelton, Li Wei, and Chotirat Ann Ratanamahatana. Fast time series classification using numerosity reduction. *Proceedings of the 23rd international conference on Machine learning (ICML)*, 2006, pages 1033—1040.
- [14] Zhengzheng Xing, Jian Pei, and Eamonn Keogh. A brief survey on sequence classification. *ACM SIGKDD Explorations Newsletter*, 2010, nov, 12(1): 40.
- [15] Bing Hu. Yanping Chen, Eamonn Keogh. The ucr time series classification archive. http://www.cs.ucr.edu/~eamonn/time_series_data/.
- [16] L Ye. The time series shapelets webpage. <http://alumni.cs.ucr.edu/~lexiangy/shapelet.html>.
- [17] Lexiang Ye and Eamonn Keogh. Time series shapelets: a new primitive for data mining. *Proceedings of the ACM International Conference on Knowledge Discovery and Data Mining*, 2009, pages 947–956.