

基于小波和动态时间弯曲的时间序列相似匹配

曲文龙 张德政 杨炳儒

北京科技大学信息工程学院, 北京 100083

摘要 提出了一种基于小波和动态时间弯曲(DTW)距离的时间序列索引和相似匹配方法. 该方法采用小波变换进行数据降维, 利用 R^{*}-tree 建立多维索引结构, 给出了查询序列的 DTW 距离边界和其在小波空间的查询超矩形的计算方法, 从而将原始空间的基于 DTW 距离的相似匹配转换为小波空间基于欧氏距离的相似匹配. 证明了此匹配方法不会产生漏报, 给出了基于 DTW 距离的范围查询算法和近邻查询算法. 实验结果表明该方法具有较高匹配精度和其较低的计算代价.

关键词 时间序列; 相似匹配; 动态时间弯曲; 小波变换

分类号 TP311.13

时间序列分析与挖掘广泛应用于包括金融、商业、气象、医学、电力、水文、工业等众多领域, 具有重要的研究价值. 相似模式匹配是时间序列分类、聚类、规则获取和预测等挖掘方法的基础, 建立索引是实现时间序列相似匹配的有效方法.

典型的相似性测度多采用欧几里德距离或其改进, 但欧氏距离测度存在局限性, 对数据在时间轴上的形变缺乏辨识能力和对噪声的鲁棒性. 动态时间弯曲(Dynamic Time Warping, DTW)^[1]可以获得很高的识别、匹配精度. 由于 DTW 距离不满足三角不等式, 在低维特征空间中无法保证检索的完整性, 因此基于 DTW 的索引和相似性搜索有待研究. Yi 给出了一个基于 FASTMAP 映射的 DTW 索引方法^[2], 但只是近似索引, 无法保证检索完整性. Park 采用线性分段表示建立 DTW 索引^[3], 但无法保证无漏报且查询精度较低. Kim 采用序列的四个特征建立索引^[4], 保证无漏报, 但未能有效的利用多维索引结构缩减搜索空间, 误报率很高. Keogh 等给出了局部 DTW 的包围边界, 并采用分段近似表示^[5], 保证无漏报并使搜索空间得到一定程度缩减. 由于小波变换可以有效约简特征空间、具有多尺度特性, 因此本文采用小波方法进行时间序列的 DTW 距离索引和相似匹配, 给出了 DTW 的小波变换边界, 进一步缩减搜索空间, 并证明了该方法无漏报, 且给

出了范围查询和近邻查询算法, 最后通过对比试验验证了该方法的优越性.

1 动态时间弯曲距离

1.1 动态时间弯曲概念

设有两个时间序列 $Q = (q_1, \dots, q_i, \dots, q_n)$ 和 $C = (c_1, \dots, c_j, \dots, c_m)$, 长度分别为 n 和 m (如图 1), 为利用 DTW 将两个时间序列对准, 首先定义 DTW 对准矩阵 M .

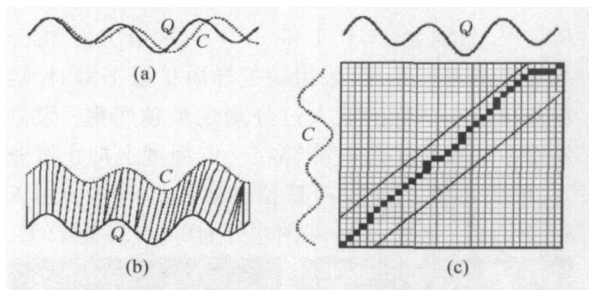


图 1 (a) 时间序列 Q 和 C ; (b) DTW 对准方式; (c) DTW 对准矩阵和最佳弯曲路径

Fig. 1 (a) Time series Q and C ; (b) DTW aligned; (c) DTW aligned matrix and optimum wrapping path

定义 1 n 行 m 列矩阵 M , 矩阵中的元素 (i, j) 为两时间序列数据中对准点 q_i 和 c_j 之间的距离 $d(q_i, c_j)$ ($d(q_i, c_j) = (q_i - c_j)^2$), 定义 M 为时间序列 Q 和 C 的 DTW 对准矩阵.

定义 2 对于两个时间序列的 DTW 对准矩阵, 定义矩阵中一组连续的矩阵元素的集合 $W = \{w_1, \dots, w_k, \dots, w_K\}$, ($w_k = d(q_i, c_j)$), 称满足如下条件的 W 为时间序列 Q 和 C 的弯曲路

收稿日期: 2005-02-21 修回日期: 2005-04-24

基金项目: 北京市自然科学基金资助项目(No. 4022008); 国家科技攻关项目(No. 2004BA616A-11)

作者简介: 曲文龙(1970-), 男, 博士研究生

径:

(1) 有界性: $\max(m, n) \leq K \leq m + n + 1$;

(2) 边界性: $w_1 = d(q_1, c_1), w_K = d(q_n, c_m)$;

(3) 连续性: 给定 $w_k = d(q_a, c_b), w_{k-1} = d(q_a^{\circ}, c_b^{\circ})$, 必有 $a - a^{\circ} \leq 1$ 且 $b - b^{\circ} \leq 1$;

(4) 单调性: 给定 $w_k = d(q_a, c_b), w_{k-1} = d(q_a^{\circ}, c_b^{\circ})$, 必有 $a - a^{\circ} \geq 0$ 且 $b - b^{\circ} \geq 0$.

除上述限制条件外, 在实际应用中还需要限制弯曲路径的宽度, 防止病态弯曲和提高 DTW 算法的速度, 通常对弯曲窗口加以限制, 即对于路径中任一点 ($w_k = d(q_i, c_j)$), 要求 $|i - j| \leq r$. 弯曲路径存在多解, DTW 距离取弯曲路径总长

度的最小值, 即 $d_{DTW}(Q, C) = \min \left[\sum_{k=1}^K w_k \right]$, 最佳路径可以由时间起始点 (1, 1) 到终点 (m, n) 之间的局部最优解通过递归获得, 公式如下:

$$\begin{cases} S(1, 1) = d(q_1, c_1) \\ S(i, j) = d(q_i, c_j) + \min(S(i-1, j), S(i, j-1), S(i-1, j-1)) \end{cases} \quad (1)$$

式中, $S(i, j)$ 为累积距离, 由当前对准点的距离和相邻点的累积 DTW 距离计算得到, 则 $d_{DTW}(Q, C) = \sqrt{S(n, m)}$. 欧几里德距离可视为 DTW 距离的特例, 即第 k 路径为 $w_k = d(q_k, c_k)$, 欧氏距离要求两个序列长度相等.

1.2 动态时间弯曲距离低限边界

在基于 DTW 的相似匹配中, 为防止路径病态弯曲和提高计算速度, 通常需对弯曲路径进行限制.

定义 3 设 $w_k = d(i, j)$ 为 DTW 的第 k 个路径, 设弯曲路径限制为 $i - r \leq j \leq i + r$, r 称为弯曲范围 (r 为常数或为 i 和 j 的函数).

定义 4 $Q = (q_1, \dots, q_i, \dots, q_n)$ 为待匹配的时间序列, r 为序列中每个点允许的最大弯曲范围, 构造如下两个新序列 $U (u_i = \max(q_{i-r}, \dots, q_{i+r}))$ 和 $L (l_i = \min(q_{i-r}, \dots, q_{i+r}))$ (图 2), 称 U 和 L 分别为序列 Q 的 DTW 上边界序列和 DTW 下边界序列.

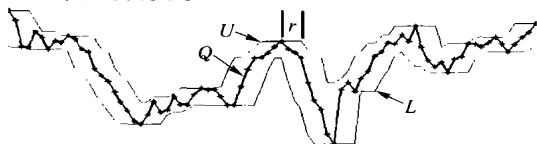


图 2 Q 的 r 弯曲上边界序列和下边界序列

Fig. 2 r -warping upper and lower bound sequences of Q

定义 5 对任意两个时间序列 Q 和 C, U 和 L 为 Q 的 DTW 上边界序列和 DTW 下边界序列, 定义 Q 和 C 的低限距离为:

$$d_{LB}(Q, C) = \sqrt{\sum_{i=1}^n \begin{cases} (c_i - u_i)^2, & c_i \geq u_i \\ (l_i - c_i)^2, & c_i \leq l_i \\ 0, & l_i < c_i < u_i \end{cases}} \quad (2)$$

定理 1 对任意两个时间序列 Q 和 C, r 为允许弯曲范围, 则其 d_{LB} 距离是两序列 DTW 距离的低限距离, 即 $d_{LB}(Q, C) \leq d_{DTW}^r(Q, C)$ (定理 1 证明参见文献[5]).

定理 2 任给时间序列 Q 和 C, U 和 L 为 Q 的 DTW 上边界序列和 DTW 下边界序列, 必存在 $L \leq T \leq U$ (即 $l_i \leq t_i \leq u_i$) 使得 $d_{LB}(Q, C) = d(T, C)$ ($d(*, *)$ 表示两时间序列的欧氏距离, 下同).

证明: 设 $U (u_i = \max(q_{i-r}, \dots, q_{i+r}))$ 和 $L (l_i = \min(q_{i-r}, \dots, q_{i+r}))$ 为 Q 的 DTW 上、下边界序列, 构造序列 T :

$$t_i = \begin{cases} u_i, & c_i > u_i \\ l_i, & c_i < l_i \\ c_i, & l_i \leq c_i \leq u_i \end{cases}$$

显然 $l_i \leq t_i \leq u_i$ 成立, 容易验证 $d_{LB}(Q, C) = d(T, C)$.

时间序列 DTW 相似匹配的基本思想, 是使用低限距离函数剪除那些不可能匹配的序列以缩减搜索空间, 并保证不产生漏报, 然后对候选序列再采用 DTW 距离进一步精确匹配. 采用的低限距离要比使用原距离计算复杂性低, 其次要使用紧低限距离 (定义的低限距离尽可能接近原距离), 以减少使用低限距离检索结果中的候选数据的数量.

小波变换具有保矩性, 采用小波变换将时间序列变换到时频域, 利用低频段的小波系数建索引结构, 可以提高检索效率.

2 离散小波变换

2.1 基本概念

小波变换是一种非平稳信号分析方法^[4], 它通过一个满足条件 $\int_R \varphi(x) dx = 0$ 的基本小波函数 $\varphi(x)$ 的平移和伸缩构成一族小波函数系去表示或逼近一个函数. 二进制小波是由伸缩因子和 2 的幂次因子满足一定条件的一组函数:

$$\varphi_{j,k}(x) = 2^{-j/2} \varphi(2^{-j}x - k), j, k \in \mathbf{Z} \quad (3)$$

对任意平方可积函数 $f(x)$ 来说, 其离散小波变换为:

$$W_{\varphi}f(j, k) = \int_{-\infty}^{+\infty} f(x) \overline{\varphi_{j,k}(x)} dx = \langle f, \varphi_{j,k} \rangle, f(x) \quad (4)$$

其中, $j, k \in \mathbf{Z}$, j 为尺度因子, k 为平移因子. 信号可以用各尺度的小波系数重构:

$$f(x) = \sum_{j,k} W_{\varphi}f(j, k) \varphi_{j,k}(x) \quad (5)$$

给定时间序列 $X = \{x(t_i)\} (i = 1, \dots, n)$ (设 $n = 2^J$, 不足位补 0), 利用 Haar 小波变换对 X 进行 J 尺度分解, 分解系数为 $A_J, D_J, \dots, D_j, \dots, D_1$ (J 为最大分解尺度). A_J 为原始序列在尺度 J 下的逼近信号, D_j 为原始序列在尺度 $j (1 \leq j \leq J)$ 下的细节信号. j 尺度信号的长度为 2^{J-j} (即 $k = 0, \dots, 2^{J-j} - 1$), 小波分解系数的总长度为 $1 + \sum_{j=1}^J 2^{J-j} = 2^J = n$, 与原序列长度相等.

小波分解系数的排列次序为: J 尺度近似信号、按尺度 j 从大到小的各尺度细节信号 (对于同一尺度的小波系数按平移因子 k 从小到大排列). 由于在时间序列相似匹配中一般首先对序列进行归一化, J 尺度近似值 A_J 均为 0, 因此只取各尺度细节信号即可.

定义 6 对长为 $n (n = 2^J)$ 的序列 X 进行 J 层小波分解, 小波系数按尺度 j 值从大到小排列 (同一尺度的小波系数按平移因子 k 从小到大排列), 取前 $m (1 \leq m \leq n - 1)$ 个系数可得到序列 $X^w = \{x_i^w\} (i = 1, \dots, m)$, 定义为 X 的小波变换序列.

利用小波变换序列近似表示原序列, 构建索引结构, 使索引维数得到约简, 对得到的匹配候选序列, 可通过后处理 (计算 DTW 距离) 来滤除匹配误报. 研究中采用简单易用的 Haar 小波, 也可以采用其他正交小波, 如 Daubechies 紧支集小波等, Haar 小波是正交小波, 其小波函数 $\varphi(x)$ 为:

$$\varphi(x) = \begin{cases} 1, & 0 < x \leq 0.5 \\ -1, & 0.5 < x \leq 1 \\ 0, & x \leq 0 \text{ 或 } x > 1 \end{cases} \quad (6)$$

尺度函数为:

$$\psi(x) = \begin{cases} 1, & 0 \leq x \leq 1 \\ 0, & x \leq 0 \text{ 或 } x > 1 \end{cases} \quad (7)$$

2.2 小波变换域的动态弯曲低限距离

利用时间序列的前 m 个小波变换系数建立

索引并进行相似匹配, 需给出两时间序列 Q 和 C 的 DTW 距离在小波域的低限距离定义 $d_{LB}^w(Q, C)$, 并证明 $d_{LB}^w(Q, C)$ 为其 DTW 距离的低限距离, 从而保证相似匹配无漏报. Haar 小波二进制平移系数为:

$$h_{j,k}(t) = \begin{cases} 2^{-j/2}, & 2^j k \leq t < (2k+1)2^{j-1} \\ -2^{-j/2}, & (2k+1)2^{j-1} \leq t < (k+1)2^j \\ 0, & t < 2^j k \text{ 或 } t \geq (k+1)2^j \end{cases} \quad (8)$$

其中, $j, k \in \mathbf{Z}$, j 为尺度因子, k 为平移因子. 对于离散时间序列 $Q = (q_1, \dots, q_i, \dots, q_n)$, 其 Haar 的小波系数为:

$$W_Q(j, k) = \sum_{i=1+2^j k}^{(2k+1)2^{j-1}} q(i) \left[2^{-\frac{j}{2}} \right] - \sum_{i=1+(2k+1)2^{j-1}}^{(k+1)2^j} q(i) \left[2^{-\frac{j}{2}} \right] \quad (9)$$

定义 7 设时间序列 Q 的 DTW 上下包围边界序列为 U 和 L , 按如下方法构造两个序列 U^w 和 L^w :

$$U(j, k) = \sum_{i=1+2^j k}^{(2k+1)2^{j-1}} u(i) \left[2^{-\frac{j}{2}} \right] - \sum_{i=1+(2k+1)2^{j-1}}^{(k+1)2^j} l(i) \left[2^{-\frac{j}{2}} \right] \quad (10)$$

$$L(j, k) = \sum_{i=1+2^j k}^{(2k+1)2^{j-1}} l(i) \left[2^{-\frac{j}{2}} \right] - \sum_{i=1+(2k+1)2^{j-1}}^{(k+1)2^j} u(i) \left[2^{-\frac{j}{2}} \right] \quad (11)$$

式中, $j = 1, \dots, \log_2 n, k = 0, \dots, j - 1$.

将 $U(j, k)$ 按 j 从大到小 (j 相同时按 k 从小到大) 排列, 得到序列 $U^w = \{u_i^w\}$; 将 $L(j, k)$ 按 j 从大到小 (j 相同时按 k 从小到大) 排列, 取前 m 个系数得到序列 $L^w = \{l_i^w\} (i = 1, \dots, m)$. 则称 U^w 和 L^w 为序列 Q 在小波变换域的 DTW 上下边界序列.

定理 3 对任意两个时间序列 Q 和 X , U 和 L 为 Q 的 DTW 上边界序列和下边界序列, U^w 和 L^w 为序列 Q 在小波变换域的 DTW 上、下边界序列, X^w 为 X 的小波变换序列. 如果 $L \leq X \leq U (l_i \leq x_i \leq u_i)$, 必有 $L^w \leq X^w \leq U^w$.

证明: 设 X 的小波系数为 $W_X(j, k)$, 由于 $l_i \leq x_i \leq u_i$, 所以有

$$W_X(j, k) = \sum_{i=1+2^j k}^{(2k+1)2^{j-1}} x(i) \left[2^{-\frac{j}{2}} \right] - \sum_{i=1+(2k+1)2^{j-1}}^{(k+1)2^j} x(i) \left[2^{-\frac{j}{2}} \right] \leq \sum_{i=1+2^j k}^{(2k+1)2^{j-1}} u(i) \left[2^{-\frac{j}{2}} \right] - \sum_{i=1+(2k+1)2^{j-1}}^{(k+1)2^j} l(i) \left[2^{-\frac{j}{2}} \right] = U(j, k)$$

和

$$W_X(j, k) = \sum_{i=1+2^j k}^{(2k+1)2^{j-1}} x(i) \begin{pmatrix} -j \\ 2 \end{pmatrix} - \sum_{i=1+(2k+1)2^{j-1}}^{(k+1)2^j} x(i) \begin{pmatrix} -j \\ 2 \end{pmatrix} \geq$$

$$\sum_{i=1+2^j k}^{(2k+1)2^{j-1}} l(i) \begin{pmatrix} -j \\ 2 \end{pmatrix} - \sum_{i=1+(2k+1)2^{j-1}}^{(k+1)2^j} u(i) \begin{pmatrix} -j \\ 2 \end{pmatrix} = L(j, k),$$

所以 $L(j, k) \leq W_X(j, k) \leq U(j, k)$,

即 $l_i^{wave} \leq x_i^{wave} \leq u_i^{wave}$,

可证 $L^w \leq X^w \leq U^w$.

定义 8 对任意时间序列 Q 和 C , C^w 为 C 的小波变换序列, U^w 和 L^w 为序列 Q 在小波域的 DTW 上、下边界序列. 如下定义 Q 和 C 的 d_{LB}^w 距离:

$$d_{LB}^w(Q, C) = \sqrt{\sum_{i=1}^m \begin{cases} (c_i^w - u_i^w)^2, & c_i^w \geq u_i^w \\ (l_i^w - c_i^w)^2, & c_i^w \leq l_i^w \\ 0, & l_i^w < c_i^w < u_i^w \end{cases}} \quad (12)$$

定理 4 任给时间序列 Q, C 和 X, U 和 L 为 Q 的 DTW 上边界序列和下边界序列, C^w 为 C 的小波变换序列, X^w 为 X 的小波变换序列. 如果 $L \leq X \leq U$ ($l_i \leq x_i \leq u_i$), 则:

$$d_{LB}^w(Q, C) \leq d(X^w, C^w).$$

证明:

$$d(C^w, X^w) = \sqrt{\sum_{i=1}^m (c_i^w - x_i^w)^2} =$$

$$\sqrt{\sum_{i=1}^m \begin{cases} (c_i^w - x_i^w)^2, & c_i^w \geq u_i^w \\ (c_i^w - x_i^w)^2, & c_i^w \leq l_i^w \\ (c_i^w - x_i^w)^2, & l_i^w < c_i^w < u_i^w \end{cases}}$$

因为 $L \leq X \leq U$, 由定理 3 可得 $L^w \leq X^w \leq U^w$, 即 $l_i^w \leq x_i^w \leq u_i^w$, 所以当 $c_i^w \geq u_i^w$ 时, 有 $(c_i^w - x_i^w)^2 \geq (c_i^w - u_i^w)^2$. 当 $c_i^w \leq u_i^w$ 时, 有 $(c_i^w - x_i^w)^2 \geq (c_i^w - l_i^w)^2$. 当 $l_i^w < c_i^w < u_i^w$ 时, 显然有 $(c_i^w - x_i^w)^2 \geq 0$. 对照式 (12), 可得 $d_{LB}^w(Q, C) \leq d(X^w, C^w)$.

定理 5 任给时间序列 Q 和 C , 必有 $d_{LB}^w(Q, C) \leq D_{DTW}^r(Q, C)$.

证明: 设 U 和 L 为 Q 的 DTW 上边界序列和下边界序列, 由定理 2, 必存在 $L \leq T \leq U$ (即 $l_i \leq t_i \leq u_i$), 使得 $d_{LB}(Q, C) = d(T, C)$.

设 C^w 为 C 的小波变换序列, T^w 为 T 的小波变换序列, 由定理 4 可得:

$$d_{LB}^w(Q, C) \leq d(T^w, C^w).$$

由于 Haar 为正交小波满足 parseval 定理, 可得

$$d(T^w, C^w) \leq d(T, C).$$

所以,

$$d_{LB}^w(Q, C) \leq d_{LB}^w(Q, C).$$

又由定理 1, 可得

$$d_{LB}^r(Q, C) \leq D_{DTW}^r(Q, C),$$

所以,

$$d_{LB}^w(Q, C) \leq D_{DTW}^r(Q, C).$$

定理 6 采用小波域的 d_{LB}^w 距离进行 DTW 距离相似匹配不会产生漏报.

证明: 设相似阈值为 ϵ , Q 为待查询序列, 任给时间序列 C , 其弯曲范围为 r 的 DTW 距离为 $D_{DTW}^r(Q, C)$.

如果 $D_{DTW}^r(Q, C) < \epsilon$, 由定理 5, 必有 $d_{LB}^w(Q, C) < \epsilon$, 即使用 d_{LB}^w 距离在小波域进行相似匹配得到的候选序列必包含所有 D_{DTW}^r 距离的匹配序列, 不会产生漏报(即在小波域没有检索到却满足查询条件的序列).

定理 6 保证对于采用 DTW 距离的相似匹配, 可以使用小波变换域定义的 d_{LB}^w 距离检索, 不会产生漏报, 但可能有误报(即在小波域检索到却不满足查询条件的序列), 误报可以通后处理来过滤以得到正确的检索结果. 该方法定义的低限距离比先前方法相比更紧, 可最大限度的减小搜索空间. 采用小波系数方法建立索引, 利用其降维能力和多尺度特性提高搜索效率.

3 基于小波和 DTW 距离的相似匹配算法

3.1 建立索引

首先对长度为 n (设 $n=2^j$, 不足位补 0) 时间序列 X 实施归一化, 以校正序列在 Y 轴的平移和幅度伸缩. 即 $x_i = \frac{x_i - \bar{x}}{x_{\max} - x_{\min}}$, 其中 $\bar{x} =$

$$\frac{1}{n} \sum_{i=0}^{n-1} x_i, x_{\max} \text{ 和 } x_{\min} \text{ 为序列的最大值和最小值.}$$

然后对标准化后的序列实施离散小波变换, 取前 m 个小波系数组成小波系数序列, 视为 m 维空间的一个点, 并将降维之后的小波系数序列组织为多维索引结构 $R^* \text{-tree}$ 索引.

$R^* \text{-tree}$ 是一种多级平衡树^[9], 树中的每个非叶节点对应一个多维超矩形, 该超矩形为其子节点代表超矩形的最小包围超矩形(MBR). 非叶

节点由多个(SRECT, P) 结构组成, 其中 P 为子节点指针, SRECT 为与子节点 P 相关的 MBR. 叶节点含有多个(SRECT, O) 结构, 其中 O 为空间对象的标识号. R^* -tree 中每个节点最多包含的结构个数称为 R^* -tree 的度. 索引建立过程如下:

(1) 对每一序列实施归一化.

(2) 每一序列实施离散 Haar 小波变换, 取前 m 个小波系数, 得到降维之后的小波系数序列, 视为多维空间中的一个点.

(3) 利用小波系数序列构建 R^* -tree 多维索引结构.

3.2 范围查询算法

范围查询定义为对一个查询点 $Q=(q_1, q_2, \dots, q_n)$, 检索 Q 与的 DTW 距离不超过 ϵ 的点集, 算法如下.

Step 1 对长为 n 的查询序列(n 维点) Q , 按定义 7 的公式计算其小波域的 DTW 上下边界序列 U^w 和 L^w .

Step 2 给定范围查询阈值 ϵ , 则对应查询超矩形为:

$$SRect(Q) = ([l_1^w - \epsilon, u_1^w + \epsilon], \dots, [l_m^w - \epsilon, u_m^w + \epsilon]) \quad (12)$$

Step 3 在 R^* -tree 索引结构内检索与查询超矩形 $SRect(Q)$ 相交的 m 维小波域点的集合, 设共有 L 个点, 并且它们对应的原始空间的(n 维)点集 $\{C_1, C_2, \dots, C_L\}$ 作为查询候选集.

Step 4 对查询候选集每一个点 $C_i (1 \leq i \leq L)$, 计算其与 Q 的 DTW 距离判断是否 $D_{DTW}^*(C_i, Q) < \epsilon$, 如果为假则认为 C_i 是误报, 为真则将 C_i 加入到查询结果集.

3.3 最近邻查询算法

k 最近邻查询定义为对一个查询点 $Q=(q_1, q_2, \dots, q_n)$, 检索与 Q 的 DTW 距离最近的 k 个点, 算法如下.

Step 1 对长度为 n 的查询序列 Q 进行小波变换, 取前 m 个小波系数组成 Q 的小波变换序列 Q^w .

Step 2 在 R^* -tree 索引中调用快速最近邻算法^[7], 得到小波域与 Q 的欧氏距离最近的 k 个点, 对应到原始空间, 按照它们与 Q 的距离从小到大的顺序排列为 $(C_1^*, C_2^*, \dots, C_k^*)$.

Step 3 取 $\epsilon = d(Q, C_k^*)$, 调用范围查询算法进行点 Q 的 DTW 距离的 ϵ 范围查询, 得到原始空间中的 L 个点, 分别计算其与 Q 的 DTW 距

离并按从小到大的顺序排列为 (C_1, C_2, \dots, C_L) .

Step 4 取 (C_1, C_2, \dots, C_k) 作为 k 最近邻查询结果.

一个点在小波域中基于欧氏距离的 k 最近邻点未必是基于 DTW 距离的 k 最近邻点(两序列的欧氏距离即弯曲度 $r=0$ 的 DTW 距离, 不小于两序列弯曲度 $r>0$ 的 DTW 距离). 因此要使用范围查询进行修正. Step 3 中得到小波域的最大欧氏距离 $\epsilon = d(Q, C_k^*)$, 为进一步的 DTW 范围查询提供了足够小的上界, 可以避免过多的剪枝并保证不出现漏报.

4 实验结果

实验采用上交所股票数据, 选取 2003 年上证 180 指数股和其它随机选取的 70 种股票 2003-01-02 到 2004-12-31 共 241 d 的交易数据. 选取每日收盘价得到 250 个维数为 256 的时间序列, 进行归一化、小波降维并构建 R^* -tree 索引. 由于早期的基于 DTW 的索引和相似匹配算法效率很低, 不列入比较范围, 只将本文方法和 PAA 方法^[3] 进行比较. 采用 Haar 小波, 尺度取 4, 得到小波系数为 16 维. 使用同样的 R^* -tree 维数比较两方法的低限距离函数、范围查询精度、所需 DTW 距离计算比率.

4.1 低限距离比较

两序列的低限距离值小于其 DTW 距离, 低限距离越接近 DTW 距离, 则该低限距离的剪枝能力越强. 随机选取 50 对序列, 计算每两序列的 Haar 低限距离、PAA 低限距离和 DTW 距离, 取 50 次计算的平均值. 图 3 是弯曲路径 ρ 取不同值时的结果, 可以看出低限距离的紧度随弯曲路径增大而减小, 但 Haar 低限距离值高于 PAA 低限距离, 表明 Haar 低限距离可以剪除更多得误报.

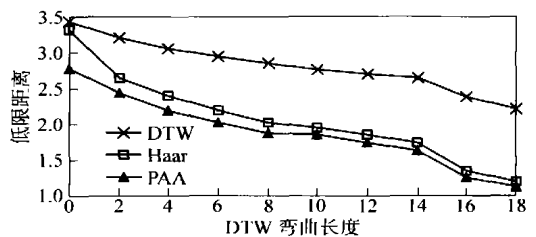


图 3 Haar 和 PAA 低限距离

Fig. 3 Haar and PAA lower bound distances

4.2 范围查询精度比较

范围查询精度定义为基于 DTW 的 ϵ 查询的

正确结果集 $A(Q)$ 与小波域内的 ϵ 查询候选集 $E(Q)$ 数目之比:

$$\text{Precision}(Q) = |A(Q)| / |E(Q)|.$$

该参数反应了低限距离的剪枝能力和需要进行 DTW 距离计算的代价, 查询精度越高则剪枝能力越强, 所需 DTW 计算的数量越小. 任选其中 50 个序列进行 DTW 的范围检索, ϵ 取为 0.7, 计算 50 个查询序列的距离小于 ϵ 的记录总数. 给出了采用 DTW 距离进行 ϵ 范围查询计算得到的真实相似记录数和用 PAA, Haar 低限距离得到的候选记录数之比. 由图 4 可知 Haar 方法候选记录数少于 PAA 方法, 因此查询中需实际计算的 DTW 距离的数目更少, 检索效率更高.

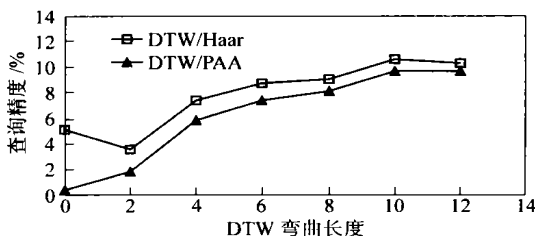


图 4 Haar 和 PAA 范围查询精度

Fig. 4 Range query accuracies of Haar and PAA

4.3 DTW 距离计算比率比较

DTW 距离计算比 = DTW 计算次数 / 记录总数, 即查询算法得到的候选集数目 (需计算其与查询序列的 DTW 距离) 与总记录数之比. 由图 5 可见 Haar 小波方法的所需的 DTW 距离计算次数小于 PAA 方法.

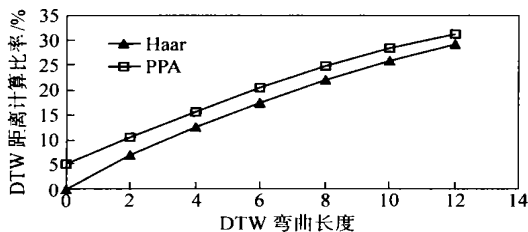


图 5 Haar 和 PAA 方法 DWT 距离计算比

Fig. 5 DTW computation ratios of Haar and PAA

以上实验表明基于小波的 DTW 相似匹配方法在低限距离、范围查询精度、所需 DTW 距离计算比率都优于 PAA 方法. 分析其原因一是该方法采用了更紧的查询边界, 二是该方法充分利

用了小波变换的多尺度特性.

5 结论

本文提出了一种基于小波变换和动态时间弯曲距离的时间序列索引和相似匹配方法, 它首先对时间序列进行归一化和 Haar 小波变换, 取前 m 个变换系数, 将小波变换序列采用多维索引结构 R^* -tree 存储. 对查询序列给出其小波变换域的 DTW 上、下边界序列以计算查询超矩形, 从而将原空间基于 DTW 距离的相似匹配转换为在小波空间的索引查询, 证明了此方法查询不会产生漏报, 给出了基于小波变换和 R^* -tree 索引的 DTW 距离范围查询和近邻查询算法. 实验结果表明该方法进一步提高了基于 DTW 的相似检索精度, 减少了计算代价, 优于已有方法. 方法可应用于多个领域, 进一步工作包括探索使用其他正交小波, 证明其可行性并与本文采用的基于 Haar 小波的 DTW 相似匹配方法进行比较.

参考文献

- [1] Rabiner L, Rosenberg A, Levinson S. Considerations in dynamic time warping algorithms for discrete word recognition. *IEEE Trans Acoust Speech Signal Process*, 1978, 26(12): 575
- [2] Yi B, Jagadeish K, Faloutsos H. Efficient retrieval of similar time sequences under time warping // *Proceedings of the 14th International Conference on Data Engineering*, 1998: 23
- [3] Park S, Lee D, Chu W. Fast retrieval of similar subsequences in long sequence databases // *3rd IEEE Knowledge and Data Engineering Exchange Workshop*, Chicago, USA, 1999: 60
- [4] Kim S, Park S, Chu W. An index-based approach for similarity search supporting time warping in large sequence databases // *Proceedings of the 17th International Conference on Data Engineering (ICDE)*, Heidelberg, Germany, 2001: 607
- [5] Keogh E. Exact indexing of dynamic time warping // *Proceedings of 28th International Conference on Very Large Data Bases*, Hong Kong, 2002: 406
- [6] Beckman N, Kriegel H, Schneider R. The R^* -tree: an efficient and robust access method for points and rectangles // *Proceedings of ACM SIGMOD International Conference on Management of Data*, New Jersey, 1990: 322
- [7] Seidl T, Kriegel H. Optimal multi-step k-nearest neighbor search // *Proceedings of ACM SIGMOD International Conference on Management of Data*, Washington, 1998: 154

Time series similar pattern matching based on wavelet and dynamic time warping

QU Wenlong, ZHANG Dezheng, YANG Bingru

Information Engineering School, University of Science and Technology Beijing, Beijing 100083, China

ABSTRACT The paper proposed a dynamic time warping (DTW) indexing and similar matching method of time series based on discrete wavelet transform, which reduced the dimensionality of time series by discrete wavelet transform and constructed multi-dimensional index structure by R^* -tree. The DTW lower bound and its discrete wavelet transform of query sequence were computed to form a query super-rectangle, thus the similar matching in original space based on DTW was converted to that in wavelet transform space based on Euclidian distance. It was proved that the method guaranteed no false dismissals and proposed the range query algorithm and nearest neighbor query algorithm. The result showed that it was a higher query precision and lower computing cost.

KEY WORDS time series; pattern matching; dynamic time warping; wavelet transform

(上接第 395 页)

Infrared image watershed segmentation based on the preprocess by CNN-PDE bias-anisotropic diffusion filter

JU Lei^{1, 2)}, ZHENG Deling¹⁾, ZHANG Lei³⁾

1) Information Engineering School, University of Science and Technology Beijing, Beijing 100083, China

2) Department of Communications Engineering, Beijing Electronic Science Technology Institute, Beijing 150027, China

3) Shandong Shengli Vocational College, Dongying 257062, China

ABSTRACT The watershed algorithm leads to over-segment when it was used to segment an infrared image. An adjustable bias-anisotropic diffusion filter based on CNN-PDE for smoothing an infrared image was studied. In order to remove the residual noise which can not be smoothed by filter, the smoothing coefficient and constrain coefficient were used to threshold the gradient image. The result of watershed segmentation of a practical infrared image shows the presented method can restrain over-segmentation effectively.

KEY WORDS image segmentation; infrared image processing; cellular neural networks; anisotropic diffusion; filter