

使用位置Burrows-Wheeler变换(PBWT)进行有效的单倍型匹配和存储)

理查德·杜宾

威廉信托桑格研究所, 威廉信托基因组校园, 剑桥CB101SA, 英国

副编辑: 杰弗里·巴雷特

摘要

动机: 近几年来, 基于后缀数组的Burrows-Wheeler变换方法被广泛应用于DNA序列读取匹配和组装。这些算法提供了非常快速的搜索算法, 线性的搜索模式大小, 在一个高度可压缩的表示数据集被搜索。同时, 基因型数据的算法发展主要集中在基于对参考数据的隐马尔可夫模型表示的概率匹配的分阶段和归责统计方法上, 这种方法虽然强大, 但计算效率要低得多。在这里, 提出了一种使用后缀数组思想的单倍型匹配理论, 它应该比目前基因型算法处理的数据集要大得多。

结果: 给定具有N个双等位基因变量位点的M序列, 给出了一种基于位置前缀阵列的数据表示O(NM)算法, 称为位置Burrows-Wheeler变换(PBWT)。在大型数据集上, 这种压缩的运行长度编码比在原始数据上使用gzip小了100多倍。使用这种表示, 给出了一种在O(NM)时间而不是O(NM)时间内找到集合中所有最大单倍型匹配的方法²⁾, 如预期的那样, 从朴素的两两比较, 也是一种快速算法, 经验上独立于M给定索引的足够内存, 以找到新序列和集合之间的最大匹配。讨论中包括一些关于如何将这方法用于估算和分阶段的建议。

可用性: <http://github.com/richarddurbin/pbwt>

联系人: 理查德·杜宾@sanger.ac.uk

于2013年9月14日收到; 2013年12月9日修改;

于2014年1月4日接受

1 引言

给定大量对齐的遗传序列或单倍型集合, 通常感兴趣的是在集合中的序列之间, 或在新的测试序列和集合中的序列之间找到长的匹配。例如, 足够长的相同子串是候选的区域, 它们通过来自共同祖先的下降(IBD)而相同。(我将使用“子串”一词来表示连续子序列, 这在计算机科学文本匹配文献中是标准的。)当使用估算方法推断缺失值时, 您希望识别与被估算位置周围的测试序列尽可能接近的序列, 例如那些是IBD的序列, 或者至少与测试序列共享长匹配。最大限度地利用这种长匹配的数量也可以形成基因型分期的基础。

原始子字符串匹配测试需要 ON^2 每个测试序列的M时间, 其中有N个可变位点和M序列, 因此有 ON^2M 在一组序列中完成所有对比较的时间。通过保持一个运行的匹配分数来找到最大匹配, 就像在BLAST中一样, 可以很容易地将这一点降低到每次测试的O(NM), 因此ONM³Langmead等人, 2009年李和杜宾, 2009年等人, 2009年 在整个集合中, 但对于大型M来说, 这仍然很大。最近, 基于后缀数组的方法在标准序列匹配中被证明是强大的, 例如Bowtie()⁴、BWA()⁵和SOAP2(Li)⁶。在这里, 描述了一种基于后缀数组的方法, 它可以在O(Nm)时间内在一组序列中找到最佳匹配, 在O(Nm)时间内对数据集进行预处理后, 也可以在O(Nm)时间内找到最佳匹配, 并在O(N)时间内找到经验上最佳的单倍型匹配。

这里描述的算法和基于标准后缀数组的序列匹配之间的差异是由于有许多序列是所有的

相同的长度和已经对齐。所以一方面没有

需要考虑测试序列相对于的偏移量

序列在集合中, 但另一方面, 测试序列是长的, 我们正在寻找最大匹配的

测试序列的任意子字符串, 而不是整个测试序列的任意子字符串。

2 方法

当观察来自人类或其他二倍体生物的遗传数据时, 每个人有两个潜在的基因组序列, 一个来自他们的父亲, 一个来自他们的母亲。这些被称为“单倍型”序列。在这里, 我考虑的情况是, 我们分别给出了这两个序列, 而不是非相位二倍体“基因型”序列, 其中两个单倍型序列已经被观察到在一起。

考虑一组M单倍型序列 x_1, \dots, x_M 的N个变量站点, 由k索引, 编号从0到(N-1)。我们可以把所有的站点都用0或1的值进行双等位, 所以典型的地点 $x_i[k]$ 是0或1。为了任何顺序S, 让我们写作 $s[k_1, k_2]$ 去代表的半开放的子字符串的s开始吧在 k_1 在 k_2 完成。我们说, 从k来的s和t之间有一个“匹配”敬 k_2 如果是 $k_1, k_2 \leq t$, k_2 如果没有扩展也是匹配的, 则此匹配是“本地最大值”。e. 如果 $k_1 \leq s[k_1 - 1]$ 或 $s[k_1 - 1] \leq t[k_1 - 1]$, 和 $k_2 \leq N$ 或 $s[k_2] \leq t[k_2]$ 当将s与序列X的集合进行比较时, 我们说s与x有一个集合最大匹配, 来自 K_1 敬 k_2 如果匹配是局部最大的, 并且不再来自匹配s去任何其他 x_j 那个包括的间隔 $[k_1, k_2]$:

对于一些应用程序，我们将对集合最大值感兴趣在X内匹配，即。每个x的集合最大匹配 g_x 去

$Y \cap x, g$:

我们的方法的基础是考虑序列集合X的子串上的一种特定形式的排序。下面是对这种排序的解释，并给出了为什么它是重要的动机。我们将考虑对0到N之间的每个位置k进行单独的排序。对于给定的k，让我们在X中对序列x进行排序，使它们的反向前缀 x^k ， k 是有序的，我的意思是，从(k-1)返回到0的前缀的反向序列是以自然的方式排序的，如果前缀是相同的，则根据它们在X中的索引i排序。

让我们考虑两个序列之间的集合最大匹配从 k^i 到 k 。如果我们在k处按这个相反的顺序排序，那么最大匹配序列必须是相邻的，因为如果它们之间有另一个前缀排序，那么它

将不得不匹配两者从 k^i 到k，因为排序顺序，它将必须匹配两个位置中的一个 $\delta k^i - 1$ ，因为最大匹配序列必须采取不同的

值在那里，只有两个可能的值。新的因此，序列将形成一个更长的匹配，这与原对之间的匹配被设置为最大值的假设相矛盾。严格地说，这个参数要求原始匹配在两个方向上都是最大的。我们将看到这一点在下面很重要。然而，这只是一个激励段落，所以没有必要考虑它在一个方向上是最大的情况。

那些事先接触后缀数组算法的人会注意到，我在这里谈论的是前缀和反向排序，而不是后缀和词典排序。在本文中，标准理论的方向被逆转，使得算法通过从0到0的序列自然地向前处理(N-1)，而不是向后；这对任何算法或结果都没有实质性影响，但将使我们能够按照它们自然出现的顺序处理非常大的数据集，并使一些表示法更自然。

2.1 前缀数组表示的派生

前一段中的论点表明，在k位置按反向前缀的顺序排序将有助于找到最大匹配。在每个位置k找到所有前缀的排序顺序似乎在计算上是昂贵的，但事实并非如此。如果我们知道k位置的排序顺序，算法1显示了如何导出排序

每个序列的k值。因此，我们可以计算每个在位置(kp1)通过一个简单的过程只看所有k在一次通过所有se的排序集-数量，时间与NM成正比。让 a_k 是序列x的索引m5 m 。从中导出k处反向排序中的第一前缀。数组 a_k 是数字0的排列，... $\delta M-1$ ：因为我们经常想讨论按前缀顺序排序的序列，

我定义你 k 成为这个排序顺序中的第一个序列 y^k 。当前位置关键观察是，条件是值 y^k ，元素的顺序 k_{b1} 与他们在a中的顺序相同 k 。文中给出了一个例子。图1

算法1构建前缀数组-构建位置前缀数组 a_{k1}
从 a_k

$u0, v0$ ，创建空数组 a^k ， b^k

因为我 $\%0!$ M-1 do

if $y^k \% k$ 然后 $\%0$

$a^k[u]$ $a_k[i]$

其他

$b^k[v]$ $a_k[i]$

a_{k1} 后面接a的级联 b

为了识别最大匹配从哪里开始，我们需要跟踪相邻前缀之间匹配的起始位置。形式上，对于i $\%0$ 定义 d_k 去是的最小值j

在这里和下面的y的k后缀，时间是隐式的k

y^k ，然后设置 d_k ，然后可以显示任何最大匹配的起始以k结

v_i ，是是由马克

有效地扩展算法1来更新 d_k ，如算法2所示。

算法2构建前缀和发散阵列来自 D_k 还有 a_k

$u0, v0, pkp1, qkp1$ 为我创

空数组 a^k ， b^k ， d^k ， e^k

0! M-1是

如果是的话， $a^k[u]$ $a_k[i]$ ， $d^k[u]$ p

其他

$b^k[v]$ $a_k[i]$ ， $e^k[v]$ q

a_{k1} a后面接b的级联

d_{k1} 后面跟着d的级联 e

因为我们在处理双等数据，

和 y_i d_k 我1必须分别为0和

按排序顺序排列。这意味着作为推论，d是不可能的 k

等于 d_k p 只要它们大于0，因为否则 y_i d_k 我1需要同时是1

和0，这是不可能的。]-]

我将数组的集合称为 a_k 对于所有k，X的“位置前缀数组。这些与标准后缀数组有关，但除了作为前缀而不是后缀数组，因为排序是相反的方向，它们的不同是因为它们形成了一组N个数组，每个数组的大小为M，而不是单个大小为NM的数组。医生 k 在标准后缀数组算法中包含相当于“最长公共前缀”值。

2.2 在X长于最小长度L

现在我们可以用 a_k 和 d_k

只计数一次匹配，我们将

每个k处报告匹配

最后在k，我。e。对于哪

要找到与新序列的匹配，如算法5中的匹配，我们还需要存储的数组 u_k 和 v_k 为了评估扩展-

离子函数 $w_k \delta b$ ：这些数组对应于信息存储在所描述的索引中，Ferragina和Manzini（2000年）俗称FM指数。考虑到PBWT，我们可以存储所需的信息，以便以与正常字符串完全类似的方式有效地生成它们。

我们确实需要 a 的值 k 对于报告，但是考虑到PBWT和位置FM索引中的 y 值，我们只能通过存储 a 来有效地做到这一点 k 对于 k 的值，例如每32或64个位置。报告的匹配将比此长，因此通过使用扩展函数 w 扩展到 a 的下一个存储值 $k \delta b$ 相对便宜。

最后，我们需要 d 的紧凑表示 k 数组。目前，提出了对赫夫曼编码的差异-也许只是商店他们在 a 子集的 k 至于 a_k 。这里可能还有进一步改进的余地。

3 结果

在这里，我给出了模拟数据的初始结果。一个10万个单倍型序列的数据集，覆盖20Mb部分基因组序列使用顺序马尔可夫聚结模拟器MACS进行模拟，基本上使用命令MACS1000002e7-t0.001-r0.001（事实上，进行了更大的模拟，它崩溃了一点超过20Mb，其余的材料被削减到这一组）。陈（2009）该数据集中有370264个隔离站点。原始的MACS输出本质上包含以0‘s和1’ s编写的单倍型序列，因此大小约为37GB。用gzip压缩到1.02GB。

对关键算法进行了初步实现。这对PBWT使用单字节运行长度编码，顶部位编码值，接下来的两位选择长度是以1、64还是2048为单位，其余5位给出单位数。为了跑464次但是这通常需要2个字节，对于运行42048但是564k这通常需要3个字节。所有实验都进行了表1在苹果MacAir笔记本电脑与2.13GHz英特尔核心2Duo处理器使用一个单一的核心。对上述10万个序列的数据集进行编码需要1070秒（用户加系统），生成一个PBWT表示，其大小为7.7MB，比原始数据的gzip压缩小130倍以上。进一步的结果包括对数据子集的应用，表明相对增益随序列数的增加而增加，清楚地表明了该算法的非线性优点。这一点可以通过观察清楚地看到，对于序列数中10倍的每增加一倍，PBWT用于将单倍型值存储在站点上的平均字节数只有大约两倍。作为对真实数据的测试，对来自1000个基因组项目第一阶段数据发布1000)的染色体1数据应用了类似的措施，其中包括3007196个位点的2184个单倍型。基因组计划(2012年 该数据的gzip文件花费了303MB，而PBWT使用了51.1MB，几乎是6个较小的因子，与基于模拟数据的预期因子不远。

接下来的算法4被实现，以找到所有的集合最大匹配在模拟数据集。正如预期的时间

表1. pbwt在不断增大的数据集上的压缩性能

序列数目	1000	10 000	100 000
序列. gz大小(KB)	10 515	105 559	1 024 614
PBWT大小(KB)	1686	3372	7698
比率. gz/PBWT	6.2	31.3	133.1
PBWT字节/站点	4.6	9.1	20.8

表2在不断增大的数据集上，pbwt

序列数目	1000	10 000	100 000
集-最大时间)	12.1	120.3	1213.7
集-最大平均长度(Mb)	0.27	1.48	3.98

在序列数（2）中，取的序列是线性的，只需20分钟就能找到10万个序列中的所有最大共享子串。表

最后，评估了三种不同的方法将新序列匹配到预索引参考面板的比较性能，找到了每个新序列与参考集的所有集合最大匹配。对于这一评估，我对模拟序列数据集进行了亚采样，以近似基因型阵列实验中的典型数据，只保留了等位基因频率为45%的位点的一小部分(10%。减少了站点数为5940个，约一个前3.4kb，cor-响应于人类基因组中的850K，与标准基因分型阵列的含量相当。表3给出了从模拟中匹配1000个序列的结果，并将它们与大小在1000到50000之间的非重叠子集进行了比较。

首先实现了一种“朴素”算法，将每个序列与面板中的所有序列进行一次比较，保持覆盖测试序列中每个碱基的最佳匹配段。正如预期的那样，这种方法在面板的大小上需要线性的时间，在面板中每个序列~0.05s。第二种算法5被实现，称为“索引。表3 这需要0.9个用户秒的恒定用户时间来匹配1000个序列到10000个大小的参考面板，但对于5万个参考序列，存储的 u 、 v 和 d 数组（在本实现中没有压缩）的大小比可用内存大，导致系统时间从0.2s增加到15s，用户时间减少到1.7s。因此，我得出结论，基于PBWT的方法可以比直接搜索方法快数百倍，并在与参考面板大小无关的时间内找到匹配，正如上面推测的那样，只要相关的索引数组适合内存。

对于索引不适合内存的情况，我们仍然可以使用PBWT数据结构来提供第三个“批处理”匹配过程，这仍然比天真的要快得多

表3。以秒为单位匹配1000个新序列的时间，分为用户 (U) 和系统 (S) 对索引和批处理方法的贡献

序列数目	1000	5000	10 000	50
	000			
Nat 有	52.1	258.9	519.2	2582.6
索引	0.9u b 0.1s	0.9u b 0.1s	0.9u b 0.2s	<u>1.7ub15s</u>
批次	2.3u b 0.1	3.5u b 0.1	4.8u b 0.1	12.1u
<u>b0.1</u>				

方法。这使用了一个修改版本的内集算法4，它联合通过面板和组合的新序列集在一起，只是考虑新序列和旧序列之间的匹配。如3所示，这种批处理方法所花费的时间，其内存需求较低，与站点数量无关，随着参考面板大小的增加而增加，但仍然比简单方法中的直接搜索效率高许多倍。表在渐近上，当前实现所花费的时间将线性地依赖于M，但在不需要PBWT压缩的情况下，可以通过小心避免PBWT压缩来减少这一点。

4 讨论

在这里，我提出了一系列算法，从单倍型序列生成位置前缀数组数据结构，并将它们用于非常强的单倍型数据压缩，以及时间和空间高效的单倍型匹配。特别是，匹配算法从直接成对比较方法所花费的搜索时间中去除M的一个因子，即被匹配的单倍型序列集的大小。就这样

可能在数万以内找到所有最好的匹配序列以分钟为单位，并产生实用的潜力扩展到数百万序列的软件。尽管提出了二进制数据的算法，它们可以扩展对多个等位基因的数据有点小心。

这些算法与基于后缀数组的类似算法共享其设计的各个方面，用于一般的字符串匹配，但它们是按字符串的位置构造的，从而产生了很大的差异。一个结果是，与suf-固定阵列方法不同的是，线性时间排序算法是非平凡的，使用算法1在线性时间中构建排序位置前缀数组是简单的。这里使用的方法让人联想到)从非常大的短字符串集合生成字符串BWT。Bauer等人。(2011

布伦斯和惠勒(1994)关于高效表示，有趣的是，原始BWT是由字符串数据压缩而不是搜索引入的，它实际上构成了bzip压缩算法的基础。Valimak. 等人(2007年) 以前曾探索过使用BWT压缩自指数方法来有效地压缩和搜索来自许多个体的遗传序列数据，但这不需要像这里介绍的工作中那样对可变位点进行固定的对齐，而且本质上是不同的。

这里描述的所有算法都需要精确匹配，没有错误或丢失数据。至于序列匹配，如果更多

敏感搜索需要允许错误，仍然可以使用精确匹配算法来查找种子匹配，然后通过直接测试加入或扩展这些内容。这通常是由生产软件采取的方法，但具有强大的功能识别种子的方法是性能的关键。

顺序使用后缀/前缀数组方法的替代方法匹配是建立一个哈希表来识别精确的种子匹配。类似于这里描述的位置前缀数组的创建，可以为单倍型序列中的每个位置构建一组位置哈希表。当良好调整时，基于哈希的方法可以比基于后缀数组的方法更快，因为基本操作更简单，但它们通常需要更大的内存，特别是在后缀表示可以压缩的情况下。在标准序列匹配中不存在基因型数据的一个问题是，位置的信息含量变化很大，大多数罕见的位置信息很少，这意味着哈希词的长度需要根据序列中的位置而变化。另一种选择是基于等位基因频率大于某些值(如10%)或在某些频率范围内的位点子集构建散列，但这将丢失导致错误种子匹配的信息。

大多数关于分析大量单倍型或基因型数据的算法的研究都集中在统计方法上，这些方法对推断很有帮助，但只扩展到几个单位-

砂位和序列；例如。参见最近开发的加速方法可以处理多达数万个序列的数据，例如。马奇尼和豪伊(2010年)。Williams等人。(2012) Delaneau等人。(2012 或)。然而，这些方法提供了统计匹配方法的近似，并且仍然比这里提出的算法重得多。已有100多万人进行了基因分型，虽然在汇集该比例的数据集方面存在后勤问题，但4100000人的基因型数据正在提供Hoffman et al., 2011()。一种提高效率的方法

相位和归责可能是使用计算效率高的方法，如位置前缀数组方法来播种布朗宁和布朗宁(2007年) 匹配统计基因型算法，或在其他计算瓶颈。例如，在他们的BEAGLE软件中，从局部单倍型序列的可变长度马尔可夫模型建立了一个概率隐藏马尔可夫模型，该模型本质上是从位置前缀数组的动态截断导出的。虽然他们使用这个概率模型的算法将他们带到了一个不同的方向，但我建议这里描述的方法可以大大加快BEAGLE的模型构建阶段。

或者，也可以采取更直接的办法。大多数分阶段和归责算法从整个数据集建立一个模型，然后依次线程每个序列，以提供一个新的分阶段有效地基于一系列匹配。相反，位置前缀算法沿着所有序列共同前进。如果我们从数据的两端开始，那么在某个位置k，我们有基于当前相位的双向匹配的信息，并且可以在一个步骤中为k处的所有序列在增量之前提出等位基因的分配基于这一思想的k.方法可能是快速的，并与目前的方法相辅相成。

确认

我感谢 HengLi 和 Jared Simpson 给了我对基于 Burrows-Wheeler 变换和后缀数组的数据结构和算法的见解，以及 Tomislav Illicic 为算法 5 的早期实现的开发提供了见解。这项工作是从 2010 年对 NCBI 的一次休假访问中发现的想法发展而来的。

资金来源：惠康信托（赠款 098051）；ORAU。

利益冲突：无申报。

参考资料

鲍尔, M. j. 等人。 (2011 年) 轻型 BWT 结构, 用于非常大的字符串集合。 In: 詹卡洛, R. 和曼齐尼, G. 组合模式匹配。 斯普林格, 柏林, pp. 219 - 231.

布朗宁, S. R. 和布朗宁, B. L. (2007) 利用局部单倍型聚类, 快速、准确地进行全基因组关联研究的单倍型分期和缺失数据推断。 *Am. J. 哼。 Genet.*, 81, 1084 - 1097.

布伦斯, M. 和惠勒, D. j. (1994 年) 块排序无损数据压缩 a1- 算法。 技术报告 124, 数字设备公司, 帕洛阿尔托, 加利福尼亚州。

陈, G. k. 等人。 (2009 年) 快速灵活地模拟 DNA 序列数据。 *基因组研究.*, 19, 136 - 42.

德莱诺, 奥. 等人。 (2012 年) 数千个基因组的线性复杂性分阶段方法。 *纳特. 方法*, 9, 179-181.

费拉吉纳, P. 和曼齐尼, G. (2000 年) 机会数据结构与应用- 问题。 In: 第 41 届计算机科学基础研讨会论文集 (FOCS2000), 加利福尼亚州雷东多海滩。 *IEEE 计算机学会*, pp. 390 - 398.

1000 基因工程。 (2012 年) 1092 个人类基因组的遗传变异综合图谱。 自然, 491, 56-65.

霍夫曼, T. j. 等人。 (2011 年) 设计和覆盖高通量基因分型阵列优化的个人东亚, 非裔美国人和拉丁美洲种族/族裔使用估算和一种新的混合 SNP 选择算法。 *基因组学*, 98, 422-430.

朗米德, B. 等人。 (2009 年) 短 DNA 的超快和记忆效率比对人类基因组序列。 *基因组生物.*, 10, r25.

李, H 和杜宾, R. (2009) 快速和准确的短读对齐与 Burrows-Wheeler 变换。 *生物信息学*, 25, 1754-1760.

李, R. 等人。 (2009) SOAP2: 一种改进的用于短读对齐的超快工具。 *生物信息学*, 1966-1967 年, 25.

马奇尼, J. 和豪伊, B. (2010 年) 全基因组关联研究的基因型估算。 *纳特. Rev. Genet.*, 11, 499 - 511.

威廉姆斯, A. l. 等人。 (2012 年) 逐步收集数千个基因分型样本。 *Am. J. 哼。 Genet.*, 91, 238 - 251.

瓦利马基, N. 等人。 (2007) 压缩后缀树-基因组规模序列分析的基础。 *生物信息学*, 23, 629-630.