



(12) 发明专利

(10) 授权公告号 CN 113035275 B

(45) 授权公告日 2023. 08. 15

(21) 申请号 202110438217.0

G06F 16/26 (2019.01)

(22) 申请日 2021.04.22

G06F 16/28 (2019.01)

(65) 同一申请的已公布的文献号
申请公布号 CN 113035275 A

G06F 17/18 (2006.01)

(43) 申请公布日 2021.06.25

(73) 专利权人 广东技术师范大学
地址 510000 广东省广州市天河区石牌中
山大道293号

(56) 对比文件

CN 105044722 A, 2015.11.11

CN 106980763 A, 2017.07.25

CN 110379460 A, 2019.10.25

US 10052026 B1, 2018.08.21

US 2015197785 A1, 2015.07.16

US 2016042508 A1, 2016.02.11

US 2018060758 A1, 2018.03.01

US 2020297323 A1, 2020.09.24

(72) 发明人 李振彰 罗文 钟祺楠 翁剑波
黄亮雄 陆海威

罗文 等. 基于结合多头注意力机制BiGRU网络的生物医学命名实体识别. 计算机应用与软件. 2020, 151-155.

(74) 专利代理机构 东莞市神州众达专利商标事
务所(普通合伙) 44251
专利代理师 周松强

审查员 杨平勇

(51) Int. Cl.

G16B 20/50 (2019.01)

G06F 16/36 (2019.01)

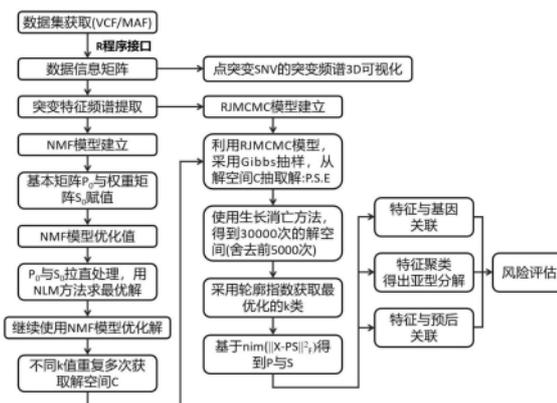
权利要求书5页 说明书11页 附图2页

(54) 发明名称

结合轮廓系数和RJCMC算法的肿瘤基因点突变的特征提取方法

(57) 摘要

本发明提供结合轮廓系数和RJCMC算法的肿瘤基因点突变的特征提取方法, 涉及肿瘤基因特征提取技术领域。该结合轮廓系数和RJCMC算法的肿瘤基因点突变的特征提取方法, 包括以下步骤: S1、数据集获取: 突变数据集突变类型包含 Somatic SNV 和 Somatic INDEL, 对 Somatic SNV 和 Somatic INDEL 进行整体统计使用 MuTect 软件。该结合轮廓系数和 RJCMC 算法的肿瘤基因点突变的特征提取方法, 实现注释文件的输入模式, 方便使用, 节约前期数据处理时间, 提高效率, 突变频谱 3D 可视化展示, 让研究者可以从空间视觉上直观的看到每个类型的突变情况, 增强类型的比较效果展示, 创新性结合轮廓系数, 进行构建了 RJCMC NMF 的模型与算法实现, 以及完成代码软件装置设计, 实现特征图谱与基因关联获取的软件装置。



1. 结合轮廓系数和RJCMC算法的肿瘤基因点突变的特征提取方法,其特征在于,包括以下步骤:

S1、数据集获取:突变数据集突变类型包含Somatic SNV和Somatic INDEL,对Somatic SNV和Somatic INDEL进行整体统计使用MuTect软件,使用MuTect软件来寻找Somatic SNV和Somatic INDEL位点,对Somatic SNV和Somatic INDEL进行注释使用ANNOVAR或者Oncotator软件,利用ANNOVAR或者Oncotator软件将所检测到SNP以及InDel基因组变异与外部数据库进行注释分析,确定与人类疾病高度相关变异的基因组位置、变异频率、蛋白有害性、基因型杂合性以及所在的功能通路信息;

S2、数据信息矩阵获取:采用具有处理器的计算机,可进行并行运算操作,其中处理器配置成一个R脚本程序接口,基于数据集获取中的文件,选取匹配的参考基因组自动生成信息矩阵;

S3、突变频谱3D可视化展示:采用具有处理器的计算机,可进行并行运算操作,其中处理器配置成一个R脚本程序接口,基于数据信息矩阵获取中获取到的信息矩阵文件,生成该数据集的突变频谱可视化3D lego图;

S4、突变特征频谱获取:主要包含两个方面,其一是特征提取算法方法,其二是频谱分析软件装置;

S5、特征图谱与基因关联获取:随着特征图谱的分解出来,根据数据集中注释的基因信息实现基因与特征图谱的关联,其实现的途径是确立每个基因非沉默突变对应到某个样本;

S6、特征图谱亚型聚类与预后关联获取:基于系数矩阵信息,获取每个样本对signature的贡献度,基于这些贡献度,可以使用无监督聚类方法对样本进行分类,得到不同的亚型,然后不同亚型与临床信息关联,做预后生存分析,可以找到与预后相关的图谱特征或者与之关联的预后因子。

2. 根据权利要求1所述的结合轮廓系数和RJCMC算法的肿瘤基因点突变的特征提取方法,其特征在于,根据S1中的操作步骤,获取基于参考基因组GRCH37或者GRCH38的注释结果VCF或MAF格式的文件,注释的文件头应该包含至少五列信息:样本名,染色体编号,变异的位点坐标值,参考基因组的碱基和变异后的碱基。

3. 根据权利要求1所述的结合轮廓系数和RJCMC算法的肿瘤基因点突变的特征提取方法,其特征在于,根据S4中的操作步骤,所述特征提取算法方法,包括以下步骤:

S401、确立分析模型:

$$X_{m \times n} = P_{m \times k} S_{k \times n} + E_{m \times n}$$

约束: $P \geq 0, S \geq 0$

其中 $P \in \mathbb{R}_+^{m \times k}, S \in \mathbb{R}_+^{k \times n}$, n 为样本数目, m 为特征类型, $E \in \mathbb{R}^{m \times n}$;

S402、基于NMF算法的特征解空间的构建: $C_k = \{P, S\}$, 表示分类为 k 的空间集;

S403、可逆跳转蒙特卡罗采样算法模型构建:对于mutational signature分解来讲,其最后得到的类别里边也是96个特征比例图,设最后分解的 k 个signature就是分层,对于每个signature来讲,其特征是固定的,且每个type对应到signature的概率分配是不一样的,但其分配和为1,对于每个样本来讲,其分配到每个signature的贡献度之和为1,对单个样

本而言,设96个特征为: $y = \{y_1, \dots, y_{96}\}$

其中 y_t 为混合数目为k的多元正态混合分布模型 $f(y_t)$ 中抽取的一组随机样本观测值,

则包含未知参数 Θ 的混合模型为: $y_t \sim \sum_{j=1}^k w_j f_j(y_t | \Theta)$

由此可得似然函数模型为:

$$\begin{aligned} p(P, S, E | y, \theta) &\propto p(y, P, S, E, \theta) = p(y | P, S, E) p(P | \theta) p(S | \theta) p(E | \theta) \\ &\propto p(y | P, S, E) \prod_{mk} p(p_{ij}) \prod_{kn} p(s_{ji}) \prod_{t=1}^m p(\varepsilon_t) \end{aligned}$$

该模型的先验分布为:

$$s_{ki} \sim \text{Expon}(\alpha_{s_{ki}})$$

$$p_{ij} \sim \text{Expon}(\alpha_{p_{ij}})$$

$$\varepsilon_t^{-1} \sim G(\gamma, \delta)$$

$$i \in [1, n], t \in [1, m], j \in [1, k]$$

其中超参数为: $\theta = \{\gamma, \delta, \alpha_{p_{ij}}, \alpha_{s_{ji}}\}$;

$$\text{S404、特征相似性计算方法: } \text{sim}(A, B) = \frac{\sum_{k=1}^k A_k B_k}{\sqrt{\sum_{K=1}^K (A_K)^2} \sqrt{\sum_{K=1}^K (B_K)^2}};$$

S405、轮廓系数计算:将所有k对应的每个特征作为一个类,通过轮廓系数公式进行这k类数据的评估分析,获取轮廓指数;

S406、运行结果可视化展示方式:将基础矩阵进行归一化后,按照百分比把每个特征属性的柱状图刻画出来,采用不同的颜色进行区分。

4. 根据权利要求3所述的结合轮廓系数和RJCMC算法的肿瘤基因点突变的特征提取方法,其特征在于,根据S402中的操作步骤,还包括以下步骤:

S4021、随机选取矩阵P0与S0,并且要求P0与S0均是非负定矩阵,归一化处理信息矩阵V0的列,按照V0矩阵的每一分量概率重新生成新的信息矩阵V;

S4022、定义好目标函数模型,模型如下:

$$\arg \min_{P, S} E(\alpha, \beta) = \frac{1}{2} \|V - PS\|_F^2 + \alpha \sum_{j \neq i} P_j' P_i + \beta \sum \sum S_{ij};$$

S4023、获取最优初始解,将矩阵P0与S0进行按列拉直或者按照行拉直,然后按照P0拉直的向量在前,S0拉直的向量在后组成新的向量,该向量作为第二步模型的初始值输入,然后利用R统计软件中的nlm函数进行求最优解;

S4024、处理第三步的最优解,将小于0的分量替换为R统计软件中默认的双精度型最小的数值,然后根据S4023步骤中的向量拉直规则还原矩阵P与S,得到的P与S作为矩阵分解中的最优初始值;

S4025、获取迭代收敛解,将S4024步骤中得到的P和S,还有S4021步骤中得到的V进行跌

代计算,精度选择为 10^{-10} 次方,迭代次数上限约定为100000,计算公式如下:

$$P_{ik} = P_{ik} \cdot \frac{(VS')_{ik}}{(PSS')_{ik}}$$

$$S_{kj} = S_{kj} \cdot \frac{(P'V)_{kj}}{(P'PS)_{kj}};$$

S4026、选取不同的分解梯度k,重复操作步骤S4021到S4025,针对每个k都重复进行100次试验,记录每次试验的数据结果,结果包括:k,V,P,S,E;

S4027、所有S4026步骤中组成的解空间就是特征提取的解空间。

5.根据权利要求3所述的结合轮廓系数和RJCMC算法的肿瘤基因点突变的特征提取方法,其特征在于,根据S403中的操作步骤,所述模型的Gibbs抽样约定为如下模型:

$$I) : p(s_{ji} | X, P, S_{-ji}, E) \propto N^+(s_{ji} | \mu_{s_{ji}}, \sigma_{s_{ji}}^2)$$

$$\text{其中 } N^+(x | \mu, \sigma^2) = \begin{cases} \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{1}{2}(\frac{x-\mu}{\sigma})^2), 0 \leq x < \infty \\ 0, x < 0 \end{cases}$$

$$X_{ij} = \frac{y_{ij}}{\sum_{z=1}^n y_{iz}}$$

$\mu_{s_{ji}}$ calculate from $\{S \in C_k\}_{ji}$

$$\sigma_{s_{ji}}^2 = \frac{1}{N_c} \sum_{c=1}^{N_c} \left(\frac{\sum_{i=1}^m X_{ij} - \frac{1}{n} \sum_{i=1}^m \varepsilon_i - \sum_{z=1, z \neq j}^k (H_{ij}^c \sum_{i=1}^m P_{iz})}{\sum_{i=1}^m P_{ij}} - \mu_{s_{ji}} \right)^2$$

$$\text{其中: } H_{ji}^c = s_{ji} \cdot \frac{(P_c^T X)_{ji}}{(P_c^T P_c S_c)_{ji}}$$

$$H_{ji}^c = \frac{H_{ji}^c}{\sum_{z=1}^k H_{ji}^c}$$

$$II) : p(p_{ij} | X, P_{-ij}, S, E) \propto N^+(p_{ij} | \mu_{p_{ij}}, \sigma_{p_{ij}}^2)$$

$$\text{其中 } N^+(x | \mu, \sigma^2) = \begin{cases} \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{1}{2}(\frac{x-\mu}{\sigma})^2), 0 \leq x < \infty \\ 0, x < 0 \end{cases}$$

$$X_{ii} = \frac{y_{ii}}{\sum_{z=1}^n y_{iz}}$$

$\mu_{p_{ij}}$ calculate from $\{P \in C_k\}_{ij}$

$$\sigma_{p_{ij}}^2 = \frac{1}{N_c} \sum_{c=1}^{N_c} \left(\frac{\sum_{i=1}^n X_{ii} - \varepsilon_i - \sum_{z=1, z \neq j}^k (A_{iz}^c \sum_{i=1}^n s_{zi})}{\sum_{i=1}^n s_{ij}} - \mu_{p_{ij}} \right)^2$$

$$\text{其中: } A_{ij}^c = p_{ij} \cdot \frac{(XS_c^T)_{ij}}{(P_c S_c S_c^T)_{ij}}$$

$$A_{ij}^c = \frac{A_{ij}^c}{\sum_{z=1}^m A_{ij}^c}$$

$$III) : p(\varepsilon_t^{-1} | y_t, P_t, S) \propto \text{Gamma}(\alpha_t, \beta_t)$$

其中:

$$\beta_t = (\delta^{-1} + \frac{1}{2} \sum_{i=1}^n (y_{ti} - \sum_{j=1}^k p_{ij} s_{ji})^2)^{-1}$$

$$\alpha_t = \gamma + \frac{n}{2}。$$

6. 根据权利要求5所述的结合轮廓系数和RJCMC算法的肿瘤基因点突变的特征提取方法,其特征在于,根据I)-III),所述抽样实现包括以下步骤:

S4031、使用Gibbs抽样,从 $N^+(s_{ji} | \mu_{s_{ji}}, \sigma_{s_{ji}}^2)$ 分布中抽取 s_{ji} ;

S4032、使用Gibbs抽样,从 $N^+(p_{ij} | \mu_{p_{ij}}, \sigma_{p_{ij}}^2)$ 分布中抽取 p_{ij} ;

S4033、使用Gibbs抽样,从Gamma(α_t, β_t)分布中抽取 ε_t^{-1} ;

S4034、更新k,对于k的更新接受规则如下:

设RJCMC NMF的分解过程,分解维度k的变化看做是状态从 C_k 跳跃到 $C_{k'}$ 的过程,则设其跳跃的接受概率为:

$$\alpha(k, k') = \min \left\{ 1, \frac{p(k', \Theta_{k'} | X, \theta) q_k(k, \Theta_k)}{p(k, \Theta_k | X, \theta) q_{k'}(k', \Theta_{k'})} \right\} = \min \left\{ 1, \frac{A(k') B(k)}{A(k) B(k')} \right\}$$

其中

$$A(k) = \ln p(k, \Theta_k | X, \theta) \propto \ln p(X | k, \theta) + \ln p(P, S, E | k, \theta) + \ln p(k)$$

$$\propto -\frac{mn}{2} \ln(2\pi) - \frac{1}{2} n \sum_{i=1}^m \ln(\varepsilon_i) - \frac{1}{2} \sum_{i=1}^n (X[i, i] - PS[i, i])^T E^{-1} (X[i, i] - PS[i, i]) +$$

$$mk \ln(\alpha_p) - \alpha_p \sum_{i=1}^m \sum_{j=1}^k p_{ij} + nk \ln(\alpha_s) - \alpha_s \sum_{i=1}^n \sum_{j=1}^k s_{ji} + (\gamma - 1) \sum_{i=1}^m \ln(\varepsilon_i^{-1}) - \frac{1}{\delta} \sum_{i=1}^m \varepsilon_i^{-1} + \ln\left(\frac{1}{k_{\max} - k_{\min}}\right)$$

$$B(k) = \ln q_k(k, \Theta_k)$$

$$\propto -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^k (\ln(2\pi\sigma_{p_{ij}}^2) + \left(\frac{p_{ij} - \mu_{p_{ij}}}{\sigma_{p_{ij}}}\right)^2)$$

$$-\frac{1}{2} \sum_{j=1}^k \sum_{i=1}^n (\ln(2\pi\sigma_{s_{ji}}^2) + \left(\frac{s_{ji} - \mu_{s_{ji}}}{\sigma_{s_{ji}}}\right)^2) +$$

$$\sum_{i=1}^m ((\alpha_t - 1) \ln(\varepsilon_i^{-1}) - (\beta_t \varepsilon_i)^{-1} - \alpha_t \ln(\beta_t) - \ln \Gamma(\alpha_t))。$$

7. 根据权利要求6所述的结合轮廓系数和RJCMC算法的肿瘤基因点突变的特征提取方法,其特征在于,在S4034中的操作步骤中,所述RJCMC NMF实现包括步骤如下:

1)、设定初始值 k_0 ;

2)、计算收敛的初始 S_0, P_0 ;

3)、通过公式抽样 P, S, E ;

4)、使用生长消亡方法, $u \sim U(0, 1)$, 如果 $u \leq b_k$, 则进行生长步骤, 如果 $b_k < u \leq b_k + d_k$, 则进行消亡步骤;

5)、重复以上步骤至设定的迭代步骤(step=10000, 其中前1000次为燃烧期)。

8. 根据权利要求7所述的结合轮廓系数和RJCMC算法的肿瘤基因点突变的特征提取方法,其特征在于,在4)中的操作步骤中,所述生长步骤包括以下步骤:

4011)、 $k = k_0 + 1$;

4012)、执行2), 收敛则继续以下步骤;

4013)、从 C_k 中抽取 q_k , 即执行3);

4014)、计算 $\alpha(k_0, k)$;

4015)、计算特征之间的相似性;

4016)、 $u \sim U(0, 1)$;

4017)、如果 $u \leq \alpha(k_0, k)$, 且两两相似性均小于0.3, 则接受 k , 否则不接受。

9. 根据权利要求7所述的结合轮廓系数和RJCMC算法的肿瘤基因点突变的特征提取方法, 其特征在于, 在4) 中的操作步骤中, 所述消亡步骤包括以下步骤:

4021)、 $k = k_0 - 1$;

4022)、执行2), 收敛则继续以下步骤;

4023)、从 C_k 中抽取 q_k , 即执行3);

4024)、计算 $\alpha(k_0, k)$;

4025)、计算特征之间的相似性;

4026)、 $u \sim U(0, 1)$;

4027)、如果 $u \leq \alpha(k_0, k)$, 且两两相似性均小于0.3, 则接受 k , 否则不接受。

结合轮廓系数和RJCMC算法的肿瘤基因点突变的特征提取方法

技术领域

[0001] 本发明涉及肿瘤基因特征提取技术领域,具体为结合轮廓系数和RJCMC 算法的肿瘤基因点突变的特征提取方法。

背景技术

[0002] 癌症是基因疾病,是由生物体细胞突变引起的。随着基因检测技术例如下一代测序(NGS)的发展,人们发现这些突变是由特定突变特征的组合引起的,这些突变特征通常具有已知的基础过程,它可以更好地提供癌症机制信息,也有助于癌症的预防和治疗。人类基因组是由染色体组成,每个染色体由四种不同的核苷酸组成——A/C/G/T。四个核苷酸实际上形成两对A-T、C-G,当A位于一个链上时,T位于另一个链上,当G位于一个链上时,C必须在同一位置组成。当癌症基因组发生突变时,其中一个核苷酸被另一个核苷酸交换,例如,T被A取代。除了替换(如插入和删除)之外,还有其他突变。突变可能是有缺陷的DNA修复或不同的突变过程的结果,如突变暴露(辐射,吸烟),DNA的酶修饰等。实际上大多数突变都是无害的。按照突变的类型可以分为六大类,分别为C>A(表示有C变异成A),C>G,C>T,T>A,T>C和T>G,按照三碱基核算则可以分为96种不同的突变类型。突变性特征是由不同的突变过程引起的突变类型的某种组合,然后除以该签名引起的突变总数,以便最终考虑每种突变类型的比例贡献。研究表明,某些突变类型在特定癌症中发生更为频繁。例如对肺和皮肤肿瘤中突变的肿瘤基因的分析表明,发现的突变类型与烟草致癌物和紫外光的实验结果相匹配,这主要是已知的外源性致癌物质影响着这些突变类型。值得注意的是,C:G>A:T突变在吸烟相关的肺癌中占主导地位,而C:G>T:A主要发生在dipyrimidines 和CC:GG>TT:AA双核苷酸替代是常见的紫外线光相关皮肤癌的变化特征。因此,从基因组突变数据中寻找这些特征对于发现癌症的基本机制,做好预防和治疗非常重要。

[0003] 目前,NMF即非负矩阵分解法是很多研究者关注的重点。NMF的基本原理是将信号矩阵分解为基本矩阵和相应的系数矩阵,根据代价函数来计算各个信号成分所对应的基本矩阵和系数矩阵,从而实现信号的分离。当下,研究工作者合理地认为在细胞中发生的生化过程通常是独立作用的,因此,可以假设基因组中的突变是细胞中所有突变过程活动的总和,其数据是所有检测样本的不同突变类型的突变计数和,即为观测到的信号矩阵Y。给定模型, $Y=WX$,其中W为系数矩阵,也就是不同签名的集合,可以理解为Mutational Signature,X为基本矩阵,也就是决定其活动的强度,代表的是每个样本在每个Mutational Signature的贡献度。

[0004] NMF的优点是稳定性功能,它很好地确定了正确的签名数,由其衍生出一些生物学方法,专门应用于肿瘤特征图谱的提取的,比如NMF、BayeNMF、SigProfiler以及SignatureAnalyzer等。但在大多数人类癌症类型中,DNA 损伤和修复过程所印迹的突变特征受到非常有限的表征,而且这些方法存在一定的局限性,功能相对单一,且对于一些数据集的分析结果差强人意,尤其是小样本数据或者低深度数据的结果,误差比较大。

发明内容

[0005] 本发明提供的发明目的在于提供结合轮廓系数和RJCMC算法的肿瘤基因点突变的特征提取方法,该提供结合轮廓系数和RJCMC算法的肿瘤基因点突变的特征提取方法具有更全面的分析内容,适用于大小样本数据集,稳定性高,操作方便,包括从基因突变数据注释结果(MAF或VCF格式,参考基因组可以是GRCH37或GRCH38)生成信号矩阵和突变集合的三维可视化图,基于轮廓系数与RJCMC方法的图谱特征提取,和基因与突变特征图谱的关联研究功能。

[0006] 为了实现上述的效果,本发明提供如下技术方案:结合轮廓系数和RJCMC 算法的肿瘤基因点突变的特征提取方法,包括以下步骤:

[0007] S1、数据集获取:突变数据集突变类型包含Somatic SNV和Somatic INDEL,对Somatic SNV和Somatic INDEL进行整体统计使用MuTect软件,使用MuTect 软件来寻找Somatic SNV和Somatic INDEL位点,对Somatic SNV和Somatic INDEL进行注释使用ANNOVAR或者Oncotator软件,利用ANNOVAR或者 Oncotator软件将所检测到SNP以及InDel基因组变异与外部数据库进行注释分析,确定与人类疾病高度相关变异的基因组位置、变异频率、蛋白有害性、基因型杂合性以及所在的功能通路信息;

[0008] S2、数据信息矩阵获取:采用具有处理器的计算机,可进行并行运算操作,其中处理器配置成一个R脚本程序接口,基于数据集获取中的文件,选取匹配的参考基因组自动生成信息矩阵;

[0009] S3、突变频谱3D可视化展示:采用具有处理器的计算机,可进行并行运算操作,其中处理器配置成一个R脚本程序接口,基于数据信息矩阵获取中获取到的信息矩阵文件,生成该数据集的突变频谱可视化3D lego图;

[0010] S4、突变特征频谱获取:主要包含两个方面,其一是特征提取算法方法,其二是频谱分析软件装置;

[0011] S5、特征图谱与基因关联获取:随着特征图谱的分解出来,根据数据集中注释的基因信息实现基因与特征图谱的关联,其实现的途径是确立每个基因非沉默突变对应到某个样本;

[0012] S6、特征图谱亚型聚类与预后关联获取:基于系数矩阵信息,获取每个样本对signature的贡献度,基于这些贡献度,可以使用无监督聚类方法对样本进行分类,得到不同的亚型,然后不同亚型与临床信息关联,做预后生存分析,可以找到与预后相关的图谱特征或者与之关联的预后因子。

[0013] 进一步的,根据S1中的操作步骤,获取基于参考基因组GRCH37或者GRCH38的注释结果VCF或MAF格式的文件,注释的文件头应该包含至少五列信息:样本名,染色体编号,变异的位点坐标值,参考基因组的碱基和变异后的碱基。

[0014] 进一步的,根据S4中的操作步骤,所述特征提取算法方法,包括以下步骤:

[0015] S401、确立分析模型:

$$[0016] X_{m \times n} = P_{m \times k} S_{k \times n} + E_{m \times n}$$

[0017] 约束: $P \geq 0, S \geq 0$

[0018] 其中 $P \in \mathfrak{R}_+^{m \times k}, S \in \mathfrak{R}_+^{k \times n}$, n 为样本数目, m 为特征类型, $E \in \mathfrak{R}^{m \times n}$;

[0019] S402、基于NMF算法的特征解空间的构建: $C_k = \{P, S\}$, 表示分类为 k 的空间集;

[0020] S403、可逆跳转蒙特卡罗采样算法模型构建:对于mutational signature 分解来讲,其最后得到的类别里边也是96个特征比例图,设最后分解的k个 signature就是分层,对于每个signature来讲,其特征是固定的,且每个type对应到signature的概率分配是不一样的,但其分配和为1,对于每个样本来讲,其分配到每个signature的贡献度之和为1,对单个样本而言,设 96个特征为: $y = \{y_1, \dots, y_{96}\}$

[0021] 其中 y_t 为混合数目为k的多元正态混合分布模型 $f(y_t)$ 中抽取的一组随机样本观测值,则包含未知参数 Θ 的混合模型为: $y_t \sim \sum_{j=1}^k w_j f_j(y_t | \Theta)$

[0022] 由此可得似然函数模型为:

$$p(P, S, E | y, \theta) \propto p(y, P, S, E, \theta) = p(y | P, S, E) p(P | \theta) p(S | \theta) p(E | \theta)$$

$$[0023] \propto p(y | P, S, E) \prod_{mk} p(p_{ij}) \prod_{kn} p(s_{ji}) \prod_{i=1}^m p(\varepsilon_i) ,$$

[0024] 该模型的先验分布为:

$$[0025] s_{ki} \sim \text{Expon}(\alpha_{s_{ki}})$$

$$[0026] p_{ij} \sim \text{Expon}(\alpha_{p_{ij}})$$

$$[0027] \varepsilon_i^{-1} \sim G(\gamma, \delta)$$

$$[0028] i \in [1, n], t \in [1, m], j \in [1, k]$$

[0029] 其中超参数为: $\theta = \{\gamma, \delta, \alpha_{p_{ij}}, \alpha_{s_{ji}}\}$;

$$[0030] \text{S404、特征相似性计算方法: } \text{sim}(A, B) = \frac{\sum_{k=1}^K A_k B_K}{\sqrt{\sum_{K=1}^K (A_K)^2} \sqrt{\sum_{K=1}^K (B_K)^2}} ;$$

[0031] S405、轮廓系数计算:将所有k对应的每个特征作为一个类,通过轮廓系数公式进行这k类数据的评估分析,获取轮廓指数;

[0032] S406、运行结果可视化展示方式:将基础矩阵进行归一化后,按照百分比把每个特征属性的柱状图刻画出来,采用不同的颜色进行区分。

[0033] 进一步的,S4021、随机选取矩阵P0与S0,并且要求P0与S0均是非负定矩阵,归一化处理信息矩阵V0的列,按照V0矩阵的每一分量概率重新生成新的信息矩阵V;

[0034] S4022、定义好目标函数模型,模型如下:

$$[0035] \arg \min_{P, S} E(\alpha, \beta) = \frac{1}{2} \|V - PS\|_F^2 + \alpha \sum_{j \neq i} P_j^i P_i + \beta \sum \sum S_{ij} ;$$

[0036] S4023、获取最优初始解,将矩阵P0与S0进行按列拉直或者按照行拉直,然后按照P0拉直的向量在前,S0拉直的向量在后组成新的向量,该向量作为第二步模型的初始值输入,然后利用R统计软件中的nlm函数进行求最优解;

[0037] S4024、处理第三步的最优解,将小于0的分量替换为R统计软件中默认的双精度型最小的数值,然后根据S4023步骤中的向量拉直规则还原矩阵P 与S,得到的P与S作为矩阵分解中的最优初始值;

[0038] S4025、获取迭代收敛解,将S4024步骤中得到的P和S,还有S4021步骤中得到的V进行迭代计算,精度选择为 10^{-10} 次方,迭代次数上限约定为100000,计算公式如下:

$$[0039] \quad P_{ik} = P_{ik} \cdot \frac{(VS')_{ik}}{(PSS')_{ik}}$$

$$[0040] \quad S_{kj} = S_{kj} \cdot \frac{(P'V)_{kj}}{(P'PS)_{kj}};$$

[0041] S4026、选取不同的分解梯度k,重复操作步骤S4021到S4025,针对每个k都重复进行100次试验,记录每次试验的数据结果,结果包括:k,V,P, S,E;

[0042] S4027、所有S4026步骤中组成的解空间就是特征提取的解空间。

[0043] 进一步的,根据S403中的操作步骤,所述模型的Gibbs抽样约定为如下模型:

$$[0044] \quad \text{I)} : p(s_{ji} | X, P, S_{-ji}, E) \propto N^+(s_{ji} | \mu_{s_{ji}}, \sigma_{s_{ji}}^2)$$

$$[0045] \quad \text{其中 } N^+(x | \mu, \sigma^2) = \begin{cases} \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{1}{2}(\frac{x-\mu}{\sigma})^2), & 0 \leq x < \infty \\ 0, & x < 0 \end{cases}$$

$$[0046] \quad X_{ji} = \frac{y_{ji}}{\sum_{z=1}^n y_{iz}}$$

[0047] $\mu_{s_{ji}}$ calculate from $\{S \in C_k\}_{ji}$

$$[0048] \quad \sigma_{s_{ji}}^2 = \frac{1}{N_c} \sum_{c=1}^{N_c} \left(\frac{\sum_{i=1}^m X_{ji} - \frac{1}{n} \sum_{i=1}^m \varepsilon_i - \sum_{z=1, z \neq j}^k (H_{zi}^c \sum_{i=1}^m p_{iz})}{\sum_{i=1}^m P_{ij}} - \mu_{s_{ji}} \right)^2$$

$$[0049] \quad \text{其中: } H_{ji}^c = s_{ji} \cdot \frac{(P_c^T X)_{ji}}{(P_c^T P_c S_c)_{ji}}$$

$$[0050] \quad H_{ji}^c = \frac{H_{ji}^c}{\sum_{z=1}^k H_{ji}^c}$$

$$[0051] \quad \text{II)} : p(p_{ij} | X, P_{-ij}, S, E) \propto N^+(p_{ij} | \mu_{p_{ij}}, \sigma_{p_{ij}}^2)$$

$$[0052] \quad \text{其中 } N^+(x | \mu, \sigma^2) = \begin{cases} \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{1}{2}(\frac{x-\mu}{\sigma})^2), & 0 \leq x < \infty \\ 0, & x < 0 \end{cases}$$

$$[0053] \quad X_{ji} = \frac{y_{ji}}{\sum_{z=1}^n y_{iz}}$$

[0054] $\mu_{p_{ij}}$ calculate from $\{P \in C_k\}_{ij}$

$$[0055] \quad \sigma_{p_{ij}}^2 = \frac{1}{N_c} \sum_{c=1}^{N_c} \left(\frac{\sum_{i=1}^n X_{ji} - \varepsilon_i - \sum_{z=1, z \neq j}^k (A_{iz}^c \sum_{i=1}^n s_{zi})}{\sum_{i=1}^n S_{ji}} - \mu_{p_{ij}} \right)^2$$

$$[0056] \quad \text{其中: } A_{ij}^c = p_{ij} \cdot \frac{(XS_c^T)_{ij}}{(P_c S_c S_c^T)_{ij}}$$

$$[0057] \quad A_{ij}^c = \frac{A_{ij}^c}{\sum_{z=1}^m A_{ij}^c}$$

$$[0058] \quad \text{III)} : p(\varepsilon_t^{-1} | y_t, P_t, S) \propto \text{Gamma}(\alpha_t, \beta_t)$$

[0059] 其中:

$$[0060] \quad \beta_t = (\delta^{-1} + \frac{1}{2} \sum_{i=1}^n (y_{ii} - \sum_{j=1}^k p_{ij} s_{ji})^2)^{-1}$$

[0061] $\alpha_t = \gamma + \frac{n}{2}$ 。

[0062] 进一步的,根据I) - III),所述抽样实现包括以下步骤:

[0063] S4031、使用Gibbs抽样,从 $N^+(s_{ji} | \mu_{s_{ji}}, \sigma_{s_{ji}}^2)$ 分布中抽取 s_{ji} ;

[0064] S4032、使用Gibbs抽样,从 $N^+(p_{tj} | \mu_{p_{tj}}, \sigma_{p_{tj}}^2)$ 分布中抽取 p_{tj} ;

[0065] S4033、使用Gibbs抽样,从Gamma(α_t, β_t)分布中抽取 ε_t^{-1} ;

[0066] S4034、更新k,对于k的更新接受规则如下:

[0067] 设RJCMCNMF的分解过程,分解维度k的变化看做是状态从 C_k 跳跃到 $C_{k'}$ 的过程,则设其跳跃的接受概率为:

[0068] $\alpha(k, k') = \min\{1, \frac{p(k', \Theta_{k'} | X, \theta)q_k(k, \Theta_k)}{p(k, \Theta_k | X, \theta)q_{k'}(k', \Theta_{k'})}\} = \min\{1, \frac{A(k')B(k)}{A(k)B(k')}\}$

[0069] 其中

[0070] $A(k) = \ln p(k, \Theta_k | X, \theta) \propto \ln p(X | k, \theta) + \ln p(P, S, E | k, \theta) + \ln p(k)$

[0071] $\propto -\frac{mn}{2} \ln(2\pi) - \frac{1}{2}n \sum_{i=1}^m \ln(\varepsilon_i) - \frac{1}{2} \sum_{i=1}^n (X[i, i] - PS[i, i])^T E^{-1} (X[i, i] - PS[i, i]) +$
 $mk \ln(\alpha_p) - \alpha_p \sum_{i=1}^m \sum_{j=1}^k p_{ij} + nk \ln(\alpha_s) - \alpha_s \sum_{i=1}^n \sum_{j=1}^k s_{ij} + (\gamma - 1) \sum_{i=1}^m \ln(\varepsilon_i^{-1}) - \frac{1}{\delta} \sum_{i=1}^m \varepsilon_i^{-1} + \ln(\frac{1}{k_{\max} - k_{\min}})$

[0072] $B(k) = \ln q_k(k, \Theta_k)$
 $\propto -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^k (\ln(2\pi\sigma_{p_{ij}}^2) + (\frac{p_{ij} - \mu_{p_{ij}}}{\sigma_{p_{ij}}})^2)$
 $-\frac{1}{2} \sum_{j=1}^k \sum_{i=1}^n (\ln(2\pi\sigma_{s_{ij}}^2) + (\frac{s_{ij} - \mu_{s_{ij}}}{\sigma_{s_{ij}}})^2) +$
 $\sum_{i=1}^m ((\alpha_i - 1) \ln(\varepsilon_i^{-1}) - (\beta_i \varepsilon_i)^{-1} - \alpha_i \ln(\beta_i) - \ln \Gamma(\alpha_i))$ 。

[0073] 进一步的,在S4034中的操作步骤中,所述RJCMCNMF实现包括步骤如下:

[0074] 1)、设定初始值 k_0 ;

[0075] 2)、计算收敛的初始 S_0, P_0 ;

[0076] 3)、通过公式抽样 P, S, E ;

[0077] 4)、使用生长消亡方法, $u \sim U(0, 1)$, 如果 $u \leq b_k$, 则进行生长步骤, 如果 $b_k < u \leq b_k + d_k$, 则进行消亡步骤;

[0078] 5)、重复以上步骤至设定的迭代步骤(step=10000, 其中前1000次为燃烧期)。

[0079] 进一步的,在4)中的操作步骤中,所述生长步骤包括以下步骤:

[0080] 4011)、 $k = k_0 + 1$;

[0081] 4012)、执行2),收敛则继续以下步骤;

[0082] 4013)、从 C_k 中抽取 q_k ,即执行3);

[0083] 4014)、计算 $\alpha(k_0, k)$;

[0084] 4015)、计算特征之间的相似性;

[0085] 4016)、 $u \sim U(0, 1)$;

[0086] 4017)、如果 $u \leq \alpha(k_0, k)$,且两两相似性均小于0.3,则接受k,否则不接受。

[0087] 进一步的,在4)中的操作步骤中,所述消亡步骤包括以下步骤:

[0088] 4021)、 $k = k_0 - 1$;

- [0089] 4022)、执行2),收敛则继续以下步骤;
- [0090] 4023)、从 C_k 中抽取 q_k ,即执行3);
- [0091] 4024)、计算 $\alpha(k_0, k)$;
- [0092] 4025)、计算特征之间的相似性;
- [0093] 4026)、 $u \sim U(0, 1)$;
- [0094] 4027)、如果 $u \leq \alpha(k_0, k)$,且两两相似性均小于0.3,则接受 k ,否则不接受。
- [0095] 本发明提供了结合轮廓系数和RJCMC算法的肿瘤基因点突变的特征提取方法,具备以下有益效果:

[0096] 该结合轮廓系数和RJCMC算法的肿瘤基因点突变的特征提取方法,实现注释文件的输入模式,方便使用,节约前期数据处理时间,提高效率,突变频谱3D可视化展示,让研究者可以从空间视觉上直观的看到每个类型的突变情况,增强类型的比较效果展示,创新性结合轮廓系数,进行构建了RJCMCNMF 的模型与算法实现,以及完成代码软件装置设计,实现特征图谱与基因关联获取的软件装置,实现特征图谱亚型聚类与预后关联获取的软件装置。

附图说明

- [0097] 图1为整体流程图;
- [0098] 图2为突变频谱3D可视化展示图;
- [0099] 图3为运行结果可视化展示图。

具体实施方式

[0100] 本发明提供一种技术方案:结合轮廓系数和RJCMC算法的肿瘤基因点突变的特征提取方法,包括以下步骤:

[0101] 步骤一:数据集获取:突变数据集突变类型包含Somatic SNV和Somatic INDEL,对Somatic SNV/InDel进行整体统计使用MuTect软件,使用MuTect 软件来寻找Somatic SNV和InDel位点;对Somatic SNV/InDel进行注释使用ANNOVAR或者Oncotator软件,利用ANNOVAR/Oncotator软件将所检测到 SNP以及InDel等基因组变异与外部数据库进行注释分析,以确定与人类疾病高度相关变异的基因组位置、变异频率、蛋白有害性、基因型杂合性以及所在的功能通路等信息,获取基于参考基因组GRCH37或者GRCH38的注释结果VCF或MAF格式的文件,注释的文件头应该包含至少五列信息:样本名,染色体编号,变异的位点坐标值,参考基因组的碱基,变异后的碱基。

[0102] 步骤二:数据信息矩阵获取:采用具有处理器的计算机,可进行并行运算操作,其中处理器配置成一个R脚本程序接口,步骤一)中的文件,选取匹配的参考基因组就可以自动生成信息矩阵,信息矩阵包含三部分:a)突变信息矩阵,其中行代表属性,比如以6种碱基突变类型为中心,各取5' 和3' 各一个碱基形成多种组合,该组合有96种类型,以这96种突变类型为基础,确定肿瘤基因组的突变特征信息矩阵,矩阵的列代表每一个样本;b)样本列表文件,与a)中的列一致;c)行属性名称列表,与a)中的行一致。

[0103] 步骤三:突变频谱3D可视化展示:采用具有处理器的计算机,可进行并行运算操作,其中处理器配置成一个R脚本程序接口,步骤二)获取到的信息矩阵文件,生成该数据集

合的突变频谱可视化3D lego图。该部分主要功能是展示该样本数据集中,突变类型在每Mp基因组发生的突变频率,主要计算公式如下:该突变类型的每Mp基因组的突变频率=该突变类型的突变数据集总数/基因组长度(Mp);空间转化勾勒函数主要采取勾股定理,根据按比例进行缩放的每Mp基因组发生的突变频率,进行空间描点,从而实现3D的方柱,代表该突变类型的每Mp基因组的突变频率,结果展示图如附图2所示。

[0104] 步骤四:突变特征频谱获取:该部分主要包含两个方面,其一是特征提取算法方法,其二是频谱分析软件装置。

[0105] 关于特征提取算法方法,具体技术方案如下:

[0106] 确立分析模型:

$$[0107] \quad X_{m \times n} = P_{m \times k} S_{k \times n} + E_{m \times n}$$

[0108] 约束: $P \geq 0, S \geq 0$

[0109] 其中 $P \in \mathfrak{R}_+^{m \times k}, S \in \mathfrak{R}_+^{k \times n}$, n为样本数目, m为特征类型, $E \in \mathfrak{R}^{m \times n}$ 。

[0110] 基于NMF算法的特征解空间的构建:

[0111] $C_k = \{P, S\}$, 表示分类为k的空间集;

[0112] 其中解空间的定义求解如下:

[0113] 第一步:随机选取矩阵P0与S0,并且要求P0与S0均是非负定矩阵,归一化处理信息矩阵V0的列,按照V0矩阵的每一分量概率重新生成新的信息矩阵V;

[0114] 第二步:定义好目标函数模型,模型如下:

$$[0115] \quad \arg \min_{P, S} E(\alpha, \beta) = \frac{1}{2} \|V - PS\|_F^2 + \alpha \sum_{j \neq i} P_j' P_i + \beta \sum \sum S_{ij};$$

[0116] 第三步:获取最优初始解,将矩阵P0与S0进行按列拉直或者按照行拉直,然后按照P0拉直的向量在前,S0拉直的向量在后组成新的向量,该向量作为第二步模型的初始值输入,然后利用R统计软件中的nlm函数进行求最优解;

[0117] 第四步:适当处理第三步的最优解,将小于0的分量替换为R统计软件中默认的双精度型最小的数值,然后根据第三步的向量拉直规则还原矩阵P与S,这步得到的P与S作为矩阵分解中的最优初始值;

[0118] 第五步:获取迭代收敛解,将第四步得到的P,S还有第一步得到的V进行迭代计算,精度选择为 10^{-10} 次方,迭代次数上限约定为100000,计算公式如下:

$$[0119] \quad P_{ik} = P_{ik} \cdot \frac{(VS')_{ik}}{(PSS')_{ik}}$$

$$[0120] \quad S_{kj} = S_{kj} \cdot \frac{(P'V)_{kj}}{(P'PS)_{kj}};$$

[0121] 第六步:选取不同的分解梯度k(范围应该固定在1到30),重复操作步骤第一到第五步,针对每个k都重复进行100次试验,记录每次试验的数据结果,结果包括:k,V,P,S,E;

[0122] 第七步:所有第六步组成的解空间就是特征提取的解空间。

[0123] 可逆跳转蒙特卡罗采样(RJMCMC)算法模型构建:

[0124] 对于mutational signature分解来讲,其最后得到的类别里边也是96个特征比

例图,因此这里假设最后分解的k个signature就是分层。理想状态下,对于每个signature来讲,其特征是固定的,且每个type对应到 signature的概率分配是不一样的,但其分配和为1,对于每个样本来讲,其分配到每个signature的贡献度之和为1。对单个样本而言,假设96个特征为:

$$[0125] \quad y = \{y_1, \dots, y_{96}\}$$

[0126] 其中 y_t 为混合数目为k的多元正态混合分布模型 $f(y_t)$ 中抽取的一组随机样本观测值,则包含未知参数 Θ 的混合模型为:

$$[0127] \quad y_t \sim \sum_{j=1}^k w_j f_j(y_t | \Theta);$$

[0128] 由此可得似然函数模型为:

$$[0129] \quad \begin{aligned} p(P, S, E | y, \theta) &\propto p(y, P, S, E, \theta) = p(y | P, S, E) p(P | \theta) p(S | \theta) p(E | \theta) \\ &\propto p(y | P, S, E) \prod_{mk} p(p_{ij}) \prod_{kn} p(s_{ji}) \prod_{t=1}^m p(\varepsilon_t) \end{aligned};$$

[0130] 该模型的先验分布为:

$$[0131] \quad s_{ki} \sim \text{Expon}(\alpha_{s_{ki}})$$

$$[0132] \quad p_{ij} \sim \text{Expon}(\alpha_{p_{ij}})$$

$$[0133] \quad \varepsilon_t^{-1} \sim G(\gamma, \delta)$$

$$[0134] \quad i \in [1, n], t \in [1, m], j \in [1, k];$$

[0135] 其中超参数为: $\theta = \{\gamma, \delta, \alpha_{p_{ij}}, \alpha_{s_{ji}}\}$ 。

[0136] 该模型的Gibbs抽样约定为如下模型:

$$[0137] \quad \text{I)} : p(s_{ji} | X, P, S_{-ji}, E) \propto N^+(s_{ji} | \mu_{s_{ji}}, \sigma_{s_{ji}}^2)$$

$$[0138] \quad \text{其中 } N^+(x | \mu, \sigma^2) = \begin{cases} \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{1}{2}(\frac{x-\mu}{\sigma})^2), & 0 \leq x < \infty \\ 0, & x < 0 \end{cases}$$

$$[0139] \quad X_{ii} = \frac{y_{ii}}{\sum_{z=1}^n y_{iz}}$$

$$[0140] \quad \mu_{s_{ji}} \text{ calculate from } \{S \in C_k\}_{ji}$$

$$[0141] \quad \sigma_{s_{ji}}^2 = \frac{1}{N_C} \sum_{t=1}^{N_C} \left(\frac{\sum_{i=1}^m X_{ii} - \frac{1}{n} \sum_{t=1}^m \varepsilon_t - \sum_{t=1, t \neq j}^k (H_{ii}^c \sum_{i=1}^m p_{iz})}{\sum_{i=1}^m p_{ij}} - \mu_{s_{ji}} \right)^2$$

$$[0142] \quad \text{其中: } H_{ji}^c = s_{ji} \cdot \frac{(P_c^T X)_{ji}}{(P_c^T P_c S_c)_{ji}}$$

$$[0143] \quad H_{ji}^c = \frac{H_{ji}^c}{\sum_{z=1}^k H_{ji}^c}$$

$$[0144] \quad \text{II)} : p(p_{ij} | X, P_{-ij}, S, E) \propto N^+(p_{ij} | \mu_{p_{ij}}, \sigma_{p_{ij}}^2)$$

$$[0145] \quad \text{其中 } N^+(x | \mu, \sigma^2) = \begin{cases} \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{1}{2}(\frac{x-\mu}{\sigma})^2), & 0 \leq x < \infty \\ 0, & x < 0 \end{cases}$$

$$[0146] \quad X_{ii} = \frac{y_{ii}}{\sum_{z=1}^n y_{iz}}$$

[0147] $\mu_{p_{ij}}$ calculate from $\{P \in C_k\}_g$

$$[0148] \quad \sigma_{p_{ij}}^2 = \frac{1}{N_c} \sum_{c=1}^{N_c} \left(\frac{\sum_{i=1}^n X_{ii} - \varepsilon_t - \sum_{z=1, z \neq j}^k (A_{ij}^c \sum_{z=1}^n s_{zi})}{\sum_{i=1}^n s_{ji}} - \mu_{p_{ij}} \right)^2$$

$$[0149] \quad \text{其中: } A_{ij}^c = p_{ij} \cdot \frac{(\lambda S_c^T)_{ij}}{(P_c S_c S_c^T)_{ij}}$$

$$[0150] \quad A_{ij}^c = \frac{A_{ij}^c}{\sum_{c=1}^m A_{ij}^c}$$

[0151] III) : $p(\varepsilon_t^{-1} | y_t, P_t, S) \propto \text{Gamma}(\alpha_t, \beta_t)$

[0152] 其中:

$$[0153] \quad \beta_t = \left(\delta^{-1} + \frac{1}{2} \sum_{i=1}^n (y_{ii} - \sum_{j=1}^k p_{ij} s_{ji})^2 \right)^{-1}$$

$$[0154] \quad \alpha_t = \gamma + \frac{n}{2}。$$

[0155] 以上的I, II, III的具体抽样实现步骤如下:

[0156] 1)、使用Gibbs抽样,从 $N^+(s_{ji} | \mu_{s_{ji}}, \sigma_{s_{ji}}^2)$ 分布中抽取 s_{ji} ;

[0157] 2)、使用Gibbs抽样,从 $N^+(p_{ij} | \mu_{p_{ij}}, \sigma_{p_{ij}}^2)$ 分布中抽取 p_{ij} ;

[0158] 3)、使用Gibbs抽样,从 $\text{Gamma}(\alpha_t, \beta_t)$ 分布中抽取 ε_t^{-1} ;

[0159] 4)、更新k,

[0160] 注意:对每一个 $k \in [k_{\min}, k_{\max}]$,存在一个与之匹配的参数 $\Theta_k = \{P, S, E\}$,那么对于同一个k值,则存在此k值的参数集 $C_k = \{\Theta_k\}$,那么对于所有的k,则有参数集为 $C = \bigcup_{k_{\min}}^{k_{\max}} C_k$ 。

[0161] 对于以上的k的更新接受规则如下:

[0162] 假设RJCMCMNF的分解过程,分解维度k的变化看做是状态从 C_k 跳跃到 $C_{k'}$ 的过程,则设其跳跃的接受概率为:

$$[0163] \quad \alpha(k, k') = \min \left\{ 1, \frac{p(k', \Theta_{k'} | X, \theta) q_k(k, \Theta_k)}{p(k, \Theta_k | X, \theta) q_{k'}(k', \Theta_{k'})} \right\} = \min \left\{ 1, \frac{A(k') B(k)}{A(k) B(k')} \right\}$$

[0164] 其中

$$[0165] \quad \begin{aligned} A(k) &= \ln p(k, \Theta_k | X, \theta) \propto \ln p(X | k, \theta) + \ln p(P, S, E | k, \theta) + \ln p(k) \\ &\propto -\frac{mn}{2} \ln(2\pi) - \frac{1}{2} n \sum_{i=1}^m \ln(\varepsilon_i) - \frac{1}{2} \sum_{i=1}^n (X[i, i] - PS[i, i])^T E^{-1} (X[i, i] - PS[i, i]) + \\ &mk \ln(\alpha_p) - \alpha_p \sum_{i=1}^m \sum_{j=1}^k p_{ij} + nk \ln(\alpha_s) - \alpha_s \sum_{i=1}^n \sum_{j=1}^k s_{ji} + (\gamma - 1) \sum_{i=1}^m \ln(\varepsilon_i^{-1}) - \frac{1}{\delta} \sum_{i=1}^m \varepsilon_i^{-1} + \ln \left(\frac{1}{k_{\max} - k_{\min}} \right) \end{aligned}$$

$$[0166] \quad \begin{aligned} B(k) &= \ln q_k(k, \Theta_k) \\ &\propto -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^k (\ln(2\pi\sigma_{p_{ij}}^2) + \left(\frac{p_{ij} - \mu_{p_{ij}}}{\sigma_{p_{ij}}} \right)^2) \\ &-\frac{1}{2} \sum_{j=1}^k \sum_{i=1}^n (\ln(2\pi\sigma_{s_{ji}}^2) + \left(\frac{s_{ji} - \mu_{s_{ji}}}{\sigma_{s_{ji}}} \right)^2) + \\ &\sum_{i=1}^m ((\alpha_i - 1) \ln(\varepsilon_i^{-1}) - (\beta_i \varepsilon_i)^{-1} - \alpha_i \ln(\beta_i) - \ln \Gamma(\alpha_i))。 \end{aligned}$$

[0167] 具体RJCMCMNF实现步骤如下:

[0168] 1)、设定初始值 k_0 ;

[0169] 2)、计算收敛的初始 S_0, P_0 ;

[0170] 3)、通过公式抽样 P, S, E ;

[0171] 4)、使用生长消亡方法, $u \sim U(0, 1)$, 如果 $u \leq b_k$, 则进行生长步骤, 如果 $b_k < u \leq b_k + d_k$, 则进行消亡步骤;

[0172] 5)、重复以上步骤至设定的迭代步骤 (step=10000, 其中前1000次为燃烧期。

[0173] 生长步骤包括以下步骤:

[0174] a)、 $k = k_0 + 1$;

[0175] b)、执行2), 收敛则继续以下步骤;

[0176] c)、从 C_k 中抽取 q_k , 即执行3);

[0177] d)、计算 $\alpha(k_0, k)$;

[0178] e)、计算特征之间的相似性;

[0179] f)、 $u \sim U(0, 1)$;

[0180] g)、如果 $u \leq \alpha(k_0, k)$, 且两两相似性均小于0.3, 则接受 k , 否则不接受消亡步骤包括以下步骤:

[0181] a)、 $k = k_0 - 1$;

[0182] b)、执行2), 收敛则继续以下步骤;

[0183] c)、从 C_k 中抽取 q_k , 即执行3);

[0184] d)、计算 $\alpha(k_0, k)$;

[0185] e)、计算特征之间的相似性;

[0186] f)、 $u \sim U(0, 1)$;

[0187] g)、如果 $u \leq \alpha(k_0, k)$, 且两两相似性均小于0.3, 则接受 k , 否则不接受

[0188] 特征相似性计算方法:

$$[0189] \quad sim(A, B) = \frac{\sum_{k=1}^k A_k B_k}{\sqrt{\sum_{K=1}^K (A_K)^2} \sqrt{\sum_{K=1}^K (B_K)^2}}。$$

[0190] 轮廓系数计算: 将所有 k 对应的每个特征作为一个类, 通过轮廓系数公式进行这 k 类数据的评估分析, 获取轮廓指数;

[0191] 运行结果可视化展示方式: 将基础矩阵进行归一化后, 按照百分比把每个特征属性的柱状图刻画出来, 采用不同的颜色进行区分, 如附图3所示:

[0192] 步骤五: 特征图谱与基因关联获取方法: 这部分主要随着特征图谱的分解出来, 根据数据集中注释的基因信息实现基因与特征图谱的关联, 其实现的途径是确立每个基因非沉默突变对应到某个样本, 该样本对各个signature的贡献可算, 选定贡献度大于20%作为阈值, 定义样本出现signature特征的条件, 从而统计检验 (Fisher检验) 基因与signature出现的可能性。结合基因的功能特征, 研究基因在癌症发生发展中的作用, 从而间接研究特征图谱在癌症中的作用, 甚至可以知道个体用药。比如COSMIC数据库的signature 3这个图谱特征, 跟基因BRCA1/2关联, 其与铂类化疗敏感相关。基于分解的特征图谱, 计算每个非沉默突变对signature的累积贡献概率, 从而寻找一些经典的癌基因热点突变与突变特征图谱潜在的因果关系, 有助于研究癌症发生发展的机制与变化的过程。同

时可以研究该图谱特征关联密切的热点突变富集在通路 (pathway) 的情况,有助于寻找潜在的治疗靶点与方法。

[0193] 步骤六:特征图谱亚型聚类与预后关联获取方法:基于系数矩阵信息,获取每个样本对signature的贡献度,基于这些贡献度,可以使用无监督聚类方法对样本进行分类,得到不同的亚型,然后不同亚型与临床信息关联,做预后生存分析,可以找到与预后相关的图谱特征或者与之关联的预后因子(内因或者外因)。

[0194] 尽管已经示出和描述了本发明的实施例,对于本领域的普通技术人员而言,可以理解在不脱离本发明的原理和精神的情况下可以对这些实施例进行多种变化、修改、替换和变型,本发明的范围由所附权利要求及其等同物限定。

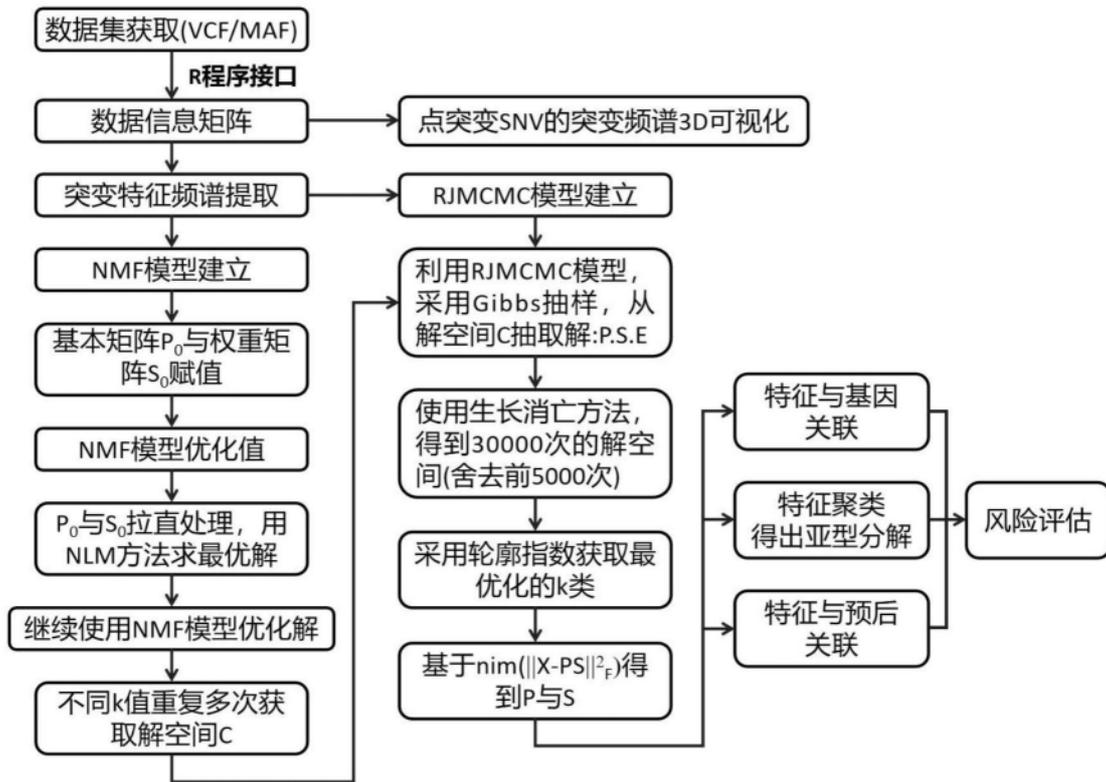


图1

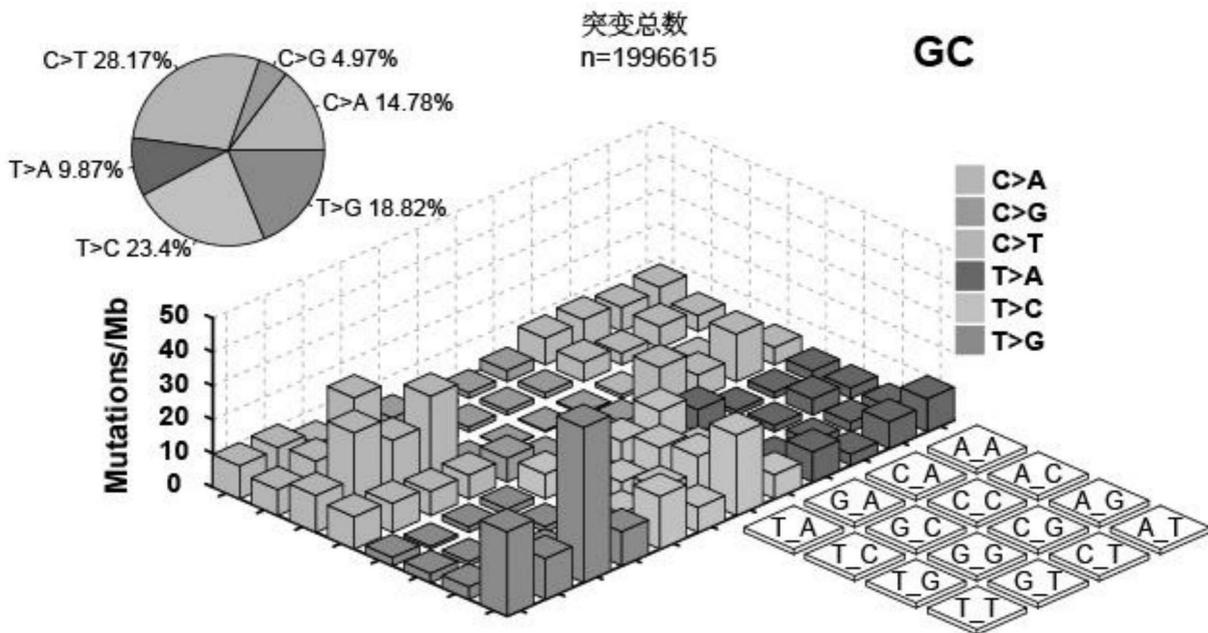


图2

