



(12) 发明专利申请

(10) 申请公布号 CN 105512669 A

(43) 申请公布日 2016. 04. 20

(21) 申请号 201410531233. 4

(22) 申请日 2014. 10. 10

(30) 优先权数据

61/975267 2014. 04. 04 US

(71) 申请人 佰欧迪塞克斯公司

地址 美国科罗拉多州

(72) 发明人 J. 勒德 H. 勒德

(74) 专利代理机构 中国专利代理(香港)有限公司

司 72001

代理人 叶晓勇 徐厚才

(51) Int. Cl.

G06K 9/62(2006. 01)

G06F 19/00(2011. 01)

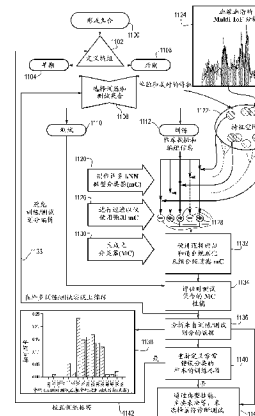
权利要求书5页 说明书53页 附图18页

(54) 发明名称

使用基于血液的样本的质谱的肺癌患者的治疗选择

(57) 摘要

用于预测与化学疗法相比、非小细胞肺部患者是否更可能获益于 EGFR-I 的测试使用对从患者所得到的基于血液的样本的质谱进行操作的计算机实现分类器。分类器利用训练集合,其包括来自作为被预测对 EGFR-I 具有总生存有益效果的一类患者的成员的其它癌症患者、例如在美国专利 7736905 所述的测试下测试 VS 良好的那些患者的基于血液的样本的质谱数据。这个类标记编组进一步细分为两个子集,即,在癌症的治疗中的 EGFR-I 的投药之后呈现疾病的早期(类标签“早期”)和后期(类标签“后期”)进展的那些患者。



1. 一种用于提前预测非小细胞肺癌(NSCLC)患者是否为与化学疗法相比、可能从采取上表皮生长因子受体抑制剂(EGFR-I)的投药的形式所述的NSCLC的治疗得到更大有益效果的一类癌症患者的成员的方法,包括下列步骤:

(a)在计算机可读介质中存储包括从通过基于血液的样本的质谱测定来确定为被预测为从癌症的治疗中的EGFR-I得到总生存有益效果的一类患者的成员的多个癌症患者所得到的类标记质谱数据的训练集合,这类患者还分为两个子类:

1. 在癌症的治疗中投药所述EGFR-I之后呈现疾病的早期进展的那些患者,这类患者的质谱数据具有“早期”或等效物的类标签;以及

2. 在癌症的治疗中投药EGFR-I之后呈现疾病的后期进展的那些患者,这类患者的质谱数据具有类标签“后期”或等效物;

(b)向质谱仪提供来自所述NSCLC患者的基于血液的样本,对所述基于血液的样本进行质谱测定,并且由此生成所述基于血液的样本的质谱;

(c)借助于经编程计算机对于在步骤(b)所得到的所述质谱进行预定义预处理步骤;

(d)在已经执行对步骤(c)所述的所述质谱的所述预处理步骤之后得到在多个预定义m/z范围的所述质谱中的所选特征的累积强度特征值;以及

(e)在所述经编程计算机中运行分类器,其包括用于将在步骤(d)所得到的所述累积强度值与在步骤(a)所存储的所述训练集合进行比较并且作为响应而生成所述基于血液的样本的类标签的分类算法,

其中,如果在步骤(e)所生成的所述类标签对所述基于血液的样本的所述质谱是“后期”或等效物,则将所述患者识别为与癌症的治疗中的化学疗法相比、可能从所述EGFR-I得到更大有益效果。

2. 如权利要求1所述的方法,其中,所述EGFR-I包括吉非替尼、厄洛替尼、第二代EGFR-I、例如达克替尼、阿法替尼或等效物。

3. 如权利要求1或2所述的方法,其中,所述分类器包括退出规则化和逻辑训练之后的经过滤微型分类器的组合(CMC/D分类器)。

4. 如权利要求1-3中的任一项所述的方法,还包括下列步骤:

进行参考样本的质谱测定,并且从所述参考样本的质谱测定来得到参考样本特征值的集合;

检查所述参考样本特征值与特征值的预定义集合的一致性;

从所述参考样本特征值来定义所述样本的所述质谱的特征校正函数;以及

按照所述特征校正函数来校正所述基于血液的样本的所述质谱的所述特征值。

5. 如权利要求1-4中的任一项所述的方法,还包括下列步骤:

a)存储特征相关噪声特性的集合;

b)生成所述基于血液的样本的所述质谱的所述特征值的有噪特征值实现的集合;

c)将所述分类器应用于所述有噪特征值实现,并且核对所述应用步骤的结果;

d)对在步骤c)所核对的所述结果生成统计数据;以及

e)结合在权利要求1的步骤(e)所生成的所述类标签来使用在步骤d)所生成的所述统计数据,以确定所述基于血液的样本的所述质谱的所述类标签。

6. 如权利要求1所述的方法,还包括权利要求4中所述的步骤以及权利要求5中所述的

步骤。

7. 如权利要求1-6中的任一项所述的方法,其中,所述训练集合包括从多个NSCLC患者的基于血液的样本所得到的类标记质谱。

8. 如权利要求1-7中的任一项所述的方法,其中,所述分类算法包括k最近邻分类算法。

9. 如权利要求1-8中的任一项所述的方法,还包括下列预备步骤:确定所述患者是否具有从EGFR-I的预测总生存有益效果的一类患者的成员,并且然后对所述样本进行权利要求1中所述的步骤b)-e)。

10. 如权利要求3所述的方法,其中,所述分类器将所述样本的所述质谱中的至少50个特征的特征值与所述训练集合的相同至少50个特征的特征值进行比较。

11. 如权利要求1-10中的任一项所述的方法,其中,所述特征值包含附录B所列的所述特征。

12. 如权利要求3所述的方法,其中,所述CMC/D分类器采取由于将分类器形成样本集合分为多个训练和测试集合而生成的多个最终分类器的形式。

13. 如权利要求12所述的方法,其中,所述多个最终分类器包括多于100个最终分类器,产生形成样本集合到训练和测试集合的多于100次划分。

14. 如权利要求13所述的方法,其中,所述多个最终分类器包括多于200个最终分类器,产生形成样本集合到训练和测试集合的多于200次划分。

15. 如权利要求3所述的方法,其中,所述CMC/D分类器被选择作为与CMC/D分类器的形成期间所生成的其它主分类器相比具有典型性能的主分类器。

16. 如权利要求1、2或3中的任一项所述的方法,其中,所述训练集合中的所述多个患者由其基于血液的样本的质谱在美国专利7736906所述的测试下测试为VeriStrat良好的那些患者组成。

17. 一种用于处理非小细胞肺癌(NSCLC)患者的基于血液的样本以确定所述患者是否为与所述NSCLC的治疗中的化学疗法相比、可能从采取上表皮生长因子受体抑制剂(EGFR-I)的投药的形式所述NSCLC的治疗得到更大有益效果的一类癌症患者的成员的系统,组合地包括:

(a)质谱仪,生成所述基于血液的样本的质谱;以及

(b)经编程计算机,包括处理单元以及存储来自所述质谱仪的质谱数据的存储器,所述存储器还存储:

1) 采取包括从通过基于血液的样本的质谱测定来确定为被预测为从癌症的治疗中的EGFR-I得到总生存有益效果的一类患者的成员的多于一个癌症患者所得到的类标记质谱数据的训练集合的形式的非暂时数据,这类患者还分为两个子类:

1. 在癌症的治疗中投药所述EGFR-I之后呈现疾病的早期进展的那些患者,这类患者的质谱数据具有“早期”或等效物的类标签;以及

2. 在癌症的治疗中投药EGFR-I之后呈现疾病的后期进展的那些患者(类标签“后期”或等效物);

2) 程序代码,用于实现采取对所述训练集合的退出规则化和逻辑训练的经过滤微型分类器的组合(CMC/D分类器)的形式的分类器;

3) 程序代码,用于对在1)所存储的所述质谱进行预定义预处理步骤,在已经执行对所

述质谱的所述预处理步骤之后得到在多个预定义 m/z 范围的所述质谱中的所选特征的累积强度特征值;以及

4)程序代码,将所述CMC/D分类器应用于在3)得到的所述累积强度值和所述训练集合,并且作为响应而生成所述基于血液的样本的类标签,

其中,如果由程序代码4)所生成的所述类标签对所述基于血液的样本是“后期”或等效物,则将所述患者识别为与癌症的治疗中的化学疗法相比、可能从所述EGFR-I得到更大有益效果。

18.如权利要求17所述的系统,其中,所述EGFR-I包括吉非替尼、厄洛替尼、第二代EGFR-I、例如达克替尼、阿法替尼或等效物。

19.如权利要求17或18所述的系统,其中,所述系统还包括基于血液的参考样本,其中所述质谱仪进行所述参考样本的质谱测定,并且所述存储器还存储用于执行下列步骤的程序:

- a.从所述参考样本的质谱来得到参考样本特征值的集合;
- b.检查所述参考样本特征值与特征值的预定义集合的一致性;
- c.从所述参考样本特征值来定义所述样本的所述质谱的特征校正函数;以及
- d.按照所述特征校正函数来校正所述基于血液的样本的所述质谱的所述特征值。

20.如权利要求17-19中的任一项所述的系统,其中:

- a)所述存储器存储表示特征相关噪声特性的集合的数据;以及
- b)所述存储器存储用于执行下列步骤的程序代码:
 - 1.生成所述基于血液的样本的所述质谱的所述特征值的有噪特征值实现的集合;
 - 2.将所述分类器应用于所述有噪特征值实现,并且核对所述应用步骤的结果;
 - 3.对在步骤2所核对的所述结果生成统计数据;以及
 - 4.结合通过权利要求16的程序代码4)所生成的所述类标签来使用在3所生成的所述统计数据,以确定所述基于血液的样本的所述质谱的所述类标签。

21.如权利要求17-20中的任一项所述的系统,其中,所述训练集合包括从多个NSCLC患者的基于血液的样本所得到的类标记质谱。

22.如权利要求17-21中的任一项所述的系统,其中,所述CMC/D分类器实现对所述测试样本的所述特征值和所述训练集合中的特征值进行操作的 K 最近邻分类算法。

23.如权利要求17-22中的任一项所述的系统,其中,针对包括与没有呈现来自癌症的治疗中的EGFR-I的投药的有益效果的患者关联的质谱数据的数据集(这种数据集具有类标签“不良”或等效物)来测试所述CMC/D分类器。

24.如权利要求17-23中的任一项所述的系统,其中,所述CMC/D分类器将所述样本的所述质谱中的至少50个特征的特征值与所述训练集合的相同至少50个特征的特征值进行比较。

25.如权利要求17-24中的任一项所述的系统,其中,所述特征值包括附录B所列的所述特征。

26.如权利要求17-25中的任一项所述的系统,其中,所述CMC/D分类器采取由于将分类器形成样本集合分为多个训练和测试集合而生成的多个最终分类器的形式。

27.如权利要求26所述的系统,其中,所述多个最终分类器包括多于100个最终分类器,

产生形成样本集合到训练和测试集合的多于100次划分。

28. 如权利要求27所述的系统,其中,所述多个最终分类器包括多于200个最终分类器,产生形成样本集合到训练和测试集合的多于200次划分。

29. 如权利要求26所述的系统,其中,所述CMC/D分类器选择作为与CMC/D分类器的形成期间所生成的其它主分类器相比具有典型性能的主分类器。

30. 一种设备,包括:

计算机存储器,存储采取包括从作为被预测从癌症的治疗中的EGFR-I得到总生存有益效果的一类患者的成员的多个癌症患者所得到的类标记质谱数据的训练集合的形式的非暂时数据,这类患者还分为两个子类:

1. 在癌症的治疗中投药所述EGFR-I之后呈现疾病的早期进展的那些患者,这类患者的质谱数据具有“早期”或等效物的类标签;以及

2. 在癌症的治疗中投药EGFR-I之后呈现疾病的后期进展的那些患者(类标签“后期”或等效物)。

31. 如权利要求30所述的设备,其中:

所述存储器还存储代码,供计算机处理单元执行以实现采取对所述训练集合的退出规则化和逻辑训练的经过滤微型分类器的组合(CMC/D分类器)的形式的分类器。

32. 如权利要求31所述的设备,其中:

所述存储器还存储来自基于血液的样本的质谱数据供所述CMC/D分类器分类。

33. 如权利要求32所述的设备,

其中,所述存储器还存储定义从所述质谱数据所得到的特征值的特征校正函数的例程,所述特征校正函数从自参考样本所得到的质谱数据来得出。

34. 如权利要求33所述的设备,其中:

a) 所述存储器存储表示从所述参考样本所得到的特征相关噪声特性的集合的数据;以及

b) 所述存储器存储用于执行下列步骤的程序代码:

1. 生成所述基于血液的样本的所述质谱的所述特征值的有噪特征值实现的集合;

2. 将所述分类器应用于所述有噪特征值实现,并且核对所述应用步骤的结果;

3. 对在步骤2所核对的所述结果生成统计数据;以及

4. 使用在3所生成的所述统计数据来确定所述基于血液的样本的所述质谱的类标签。

35. 如权利要求30所述的设备,其中,所述存储器存储特征表,其中包括附录B对所述训练集合的各成员所列的所述特征的质谱强度值。

36. 一种编程为CMC/D分类器以及存储采取所述分类器的训练集合的形式的质谱数据的存储器的计算机,所述训练集合包括附录B对所述训练集合的各成员所列的所述特征的每个的质谱数据的特征值。

37. 一种实现用于基于血液的NSCLC患者样本的的分类的分类器的经编程计算机,其中,所述经编程计算机编程为生成样本的下列类标签其中之一:1)后期或等效物,指示预测所述患者与化学疗法相比、从所述NSCLC的治疗中的EGFR-I得到更大有益效果,以及2)中等或等效物,其中预测所述患者对所述NSCLC的治疗中的所述EGFR-I或化学疗法得到相似临床结果。

38. 如权利要求37所述的经编程计算机,其中,所述经编程计算机还编程为生成样本的下列类标签:未知或等效物,在这种情况下,没有进行关于所述患者与化学疗法相比、是否可能从所述NSCLC的治疗中的EGFR-I得到更大有益效果的预测。

39. 一种治疗NSCLC患者的方法,包括下列步骤:

向所述NSCLC患者投药EGFR-I,其中

通过在经编程计算机中运行如本文所述对所述NSCLC患者的基本血液的样本进行操作的分类器,预测所述患者与化学疗法相比、更多地获益于所述EGFR-I。

40. 如权利要求39所述的方法,其中,所述经编程计算机包括权利要求37所述的经编程计算机。

使用基于血液的样本的质谱的肺癌患者的治疗选择

[0001] 相关申请的交叉引用

本申请根据35 U.S.C. § 119来要求2014年4月4日提交的美国临时申请序号61/975267的优先权,通过引用将其结合到本文中。

技术领域

[0002] 本发明涉及生物标志发现和个性化医疗的领域,以及更具体来说,涉及一种用于在治疗之前预测与化学疗法相比、非小细胞肺癌(NSCLC)患者是否可能从上表皮生长因子受体抑制剂(EGFR-I)、例如厄洛替尼或吉非替尼得到更多有益效果的方法。

背景技术

[0003] 非小细胞肺癌是美国的男性和女性死于癌症的主要原因。存在至少四(4)种不同类型的NSCLC,包括腺癌、鳞状细胞、大细胞和支气管肺泡癌(bronchoalveolar carcinoma)。肺部的鳞状细胞(表皮状)癌是与吸烟最频繁相关的显微类型的癌症。肺部的腺癌占美国的所有肺部病例的50%以上。这种癌症在女性中更常见,并且仍然是在非吸烟者中看到的最频繁类型。大细胞癌、特别是具有神经内分泌特征的癌通常与肿瘤向大脑的扩散关联。当NSCLC进入血流时,它能够扩散到不同部位,例如肝、骨、大脑和肺部的其它位置。

[0004] NSCLC的治疗多年来一直较差。化学疗法、即晚期癌症的主要治疗只是轻微有效的,除了局部癌症之外。虽然外科手术是NSCLC的最可能治疗选项,但是根据癌症阶段,它不是始终可能的。

[0005] 用于研制治疗NSCLC患者的抗癌药物的最近方式集中于降低或消除癌细胞生长和分离的能力。这些抗癌药物用来中断送往细胞的信号,以通知它们是生长还是死亡。通常,细胞生长通过细胞接收的信号来严密控制。但是,在癌症中,这个信令出错,并且细胞按照不可控方式继续生长和分离,由此形成肿瘤。当体内称作上表皮生长因子的化学品接合到在体内的许多细胞的表面上发现的受体时,这些信令通路之一开始。称作上表皮生长因子受体(EGFR)的受体经过细胞中发现的、称作酪氨酸激酶(TK)的酶的激发来向细胞发送信号。信号用来通知细胞生长和分离。

[0006] 被研制并且向NSCLC患者开处方的两个EGFR-I抗癌药物称作吉非替尼(商标名“易瑞沙”)和厄洛替尼(商标名“特罗凯”)。这些抗癌药物针对EGFR通路,并且在治疗NSCLC癌症的有效性方面是有希望的。易瑞沙抑制存在于肺癌细胞中的酶酪氨酸激酶以及正常组织中的其它癌症,并且其看来对癌细胞的生长是重要的。易瑞沙一直用作在两种其它类型的化学疗法之后有进展或者无法响应两种其它类型的化学疗法的NSCLC的治疗的单一制剂。在针对使用使用不同化合物的相同EGFR通路的研制和验证中存在其它药物,例如不可逆EGFR-TKI抑制剂阿法替尼(勃林格殷格翰)和达克替尼(辉瑞)。

[0007] 本发明人的受让人开发了称作VeriStrat®的测试,其预测NSCLC患者是否可能获益于EGFR通路靶向药物、包括吉非替尼和厄洛替尼的治疗。在美国专利7736906中描述本文中又称作“VS 1.0”的测试,通过引用将其内容结合到本文中。该测试也在在Taguchi F.等

人(J.Nat.Cancer Institute,2007 v.99(11),838-846)中描述,也通过引用将其内容结合到本文中。在本受让人的其它专利(包括美国专利7858380、7858389和7867774)中描述测试的附加应用,通过引用将其内容结合到本文中。

[0008] 简言之,VeriStrat测试基于癌症患者的血清和/或血浆样本。通过在计算机中实现的MALDI-TOF质谱和数据分析算法的组合,它借助于分类算法将在预定义m/z范围的八个累积峰值强度的集合与来自训练组群的强度进行比较。该分类算法生成患者样本的类标签:VeriStrat“良好”、VeriStrat“不良”或VeriStrat“中等”。在多个临床验证研究中,已经表明,与其样本引起VeriStrat“不良”特征的那些患者相比,其预治疗血清/血浆为VeriStrat“良好”的患者在采用上表皮生长因子受体抑制剂药物来治疗时具有明显更好的结果。在极少病例(少于2%)中,不能进行确定,从而引起VeriStrat“中等”标签。VeriStrat是从本发明的受让人Biodesix,Inc.可购买的,并且用于非小细胞肺癌患者的治疗选择中。

[0009] 从采用吉非替尼治疗的NSCLC患者的多机构研究的分析中开发VeriStrat测试。使用来自遭受长期稳定疾病或者对吉非替尼疗法的早期进展的患者的预治疗血清样本的训练集合,来开发测试。来自这些患者的血清样本的质谱(MS)用来定义12个质谱特征(即,谱峰值),从而区分这两个结果编组。测试基于k最近邻法(KNN)分类方案及其使用来自训练组群的附加谱所优化的参数来利用这些特征中的八个。测试还对采用吉非替尼或厄洛替尼来治疗的患者的两个单独组群的预治疗血清、按照单盲方式(blinded fashion)来证明。这些研究确认,分类为VeriStrat良好(VSG)的患者比分类为VeriStrat不良(VSP)的患者具有更好的结果(在一个组群中死亡风险比[HR] = 0.43 P = 0.004,在另一组群中死亡HR = 0.33 P = 0.0007)。测试表明与沿用表皮状EGFR TKI疗法而不是沿用化学疗法或术后的临床结果相关,因为在接受二线化学治疗之前、在分类为VSG或VSP的患者的总生存者(OS)中没有看到统计上显著差异(在一个组群中HR = 0.74、P = 0.42,以及在另一个组群中HR = 0.81、P = 0.54)。在具有切除早期NSCLC的患者的第三控制组群中,OS的HR为0.90(P=0.79)。

[0010] VeriStrat测试以后在称作PROSE研究的研究中正式、有希望地证明。参见患有不宜手术的非小细胞肺癌的患者中的二线厄洛替尼与化学疗法的随机蛋白质分层阶段III研究(Randomized Proteomic Stratified Phase III Study of Second-Line Erlotinib Versus Chemotherapy in Patients with Inoperable Non-Small Cell Lung Cancer, ClinicalTrials.gov # NCT00989690,向2013 ASCO conference提供的简报,2013年6月)。简言之,PROSE是在一线化学疗法治疗之后有进展的患有晚期NSCLC的285位患者的多中心随机阶段3研究。患者经过1:1随机化接受标准剂量厄洛替尼或化学疗法(以研究人员的判断的多西他赛或培美曲塞),其通过东部肿瘤协作组(ECOG)性能状态、吸烟状态和单盲VeriStrat分类来分层。PROSE结果确认,分类为VSP的患者对化学疗法与厄洛替尼具有更好的生存,并且分类为VSG的患者在采用厄洛替尼或化学疗法来治疗时具有相似OS。研究达到表明治疗结果与VeriStrat分类之间的重要交互的主要目标,其中交互p值为0.031。

[0011] 虽然PROSE结果确认,VeriStrat对厄洛替尼的取消选择是有用的测试(即,测试VSP的那些患者没有从厄洛替尼得到有益效果并且对化学疗法得到更好的生存),数据的进一步审查表明,识别患者可能对厄洛替尼具有优于化学疗法的优良生存的测试具有附加临床值。这种未满足临床需要引起本文档所述的、进行这个识别的新测试的开发。

发明内容

[0012] 在第一方面,一种用于在治疗之前预测非小细胞肺癌(NSCLC)患者是否为与化学疗法、例如多西他赛或培美曲塞相比、可能从采取上表皮生长因子受体抑制剂(EGFR-I)的投药的形式的NSCLC的治疗得到更大有益效果的一类癌症患者的成员的方法。该方法包括步骤(a):在计算机可读介质中存储采取包括从通过基于血液的样本的质谱测定来确定为被预测为从癌症的治疗中的EGFR-I得到总生存有益效果的一类患者的成员的多个癌症患者、例如VS 1.0状态为“良好”的患者所得到的类标记质谱数据的训练集合的形式的非暂时数据,这类患者还分为两个子类:

1.在癌症的治疗中投药EGFR-I之后呈现疾病的早期进展的那些患者,这类患者的质谱数据具有“早期”或等效物的类标签;以及

2.在癌症的治疗中投药EGFR-I之后呈现疾病的后期进展的那些患者(类标签“后期”或等效物)。

[0013] 该方法继续进行步骤(b):从NSCLC患者向质谱仪提供基于血液的样本,对基于血液的样本进行质谱测定,并且由此生成基于血液的样本的质谱。

[0014] 该方法继续进行步骤(c):借助于经编程计算机对于在步骤(b)所得到的质谱进行预定义预处理步骤。

[0015] 该方法继续进行步骤(d):在已经执行对步骤(c)所述的质谱的预处理步骤之后得到在多个预定义m/z范围的所述质谱中的所选特征的累积强度特征值。

[0016] 该方法继续进行步骤(e):在经编程计算机中运行分类器,其包括用于将在步骤(d)所得到的累积强度值与在步骤(a)所存储的训练集合进行比较并且作为响应而生成基于血液的样本的类标签的分类算法。如果在步骤(e)所生成的类标签对基于血液的样本的质谱是“后期”或等效物,则将该患者识别为与癌症的治疗中的化学疗法相比、可能从EGFR-I得到更大有益效果。

[0017] 存储训练集合的步骤(a)优选地在步骤(b)、(c)、(d)和(e)的执行之前执行。例如,训练集合能够使用峰值查找和本文所公开的其它方法从经过质谱测定的样本集合来形成,并且经过适当验证研究,然后存储在计算机系统、便携计算机介质、云存储或其它形式供以后使用。在给定基于血液的样本将要按照步骤(b)-(e)来测试和处理的时候,训练集合被访问并且用于按照步骤(e)的分类。

[0018] 在一个具体实施例中,组合治疗中的EGFR-I是小分子EGFR酪氨酸激酶抑制剂、例如吉非替尼或等效物、如厄洛替尼。在其它可能实施例中,EGFR-I能够采取第二代EGFR-I,例如达克替尼和阿法替尼。

[0019] 在一个实施例中,训练集合采取从多个NSCLC患者所得到的类标记质谱的形式。但是,类标记谱可能从其它类型的固体上皮肿瘤癌症患者、例如直肠癌患者或SCCHN癌症患者来得到。

[0020] 在一个实施例中,分类器采取退出规则化和逻辑训练之后的经过滤微型分类器的组合(CMC/D分类器)的形式。本文描述从样本的形成集合生成这种分类器的方法。

[0021] 在另一实施例中,该方法包括下列步骤:进行参考样本的质谱测定并且从参考样本的质谱来得到参考样本特征值的集合;检查参考样本特征值与特征值的预定义集合的一

致性；从参考样本特征值来定义样本的质谱的特征校正函数；以及按照特征校正函数来校正基于血液的样本的质谱的特征值。

[0022] 在另一个实施例中，该方法包括下列步骤：a) 存储特征相关噪声特性的集合；b) 生成基于血液的样本的质谱的特征值的有噪特征值实现的集合；c) 将分类器应用于有噪特征值实现并且核对应用步骤的结果；d) 对于在步骤c) 所核对的结果生成统计数据；以及e) 结合对样本所生成的类标签来使用在步骤d) 所生成的统计数据，以确定样本的类标签。

[0023] 在另一方面，一种用于处理非小细胞肺癌NSCLC患者的基于血液的样本以确定患者是否为与NSCLC的治疗中的化学疗法相比、可能从采取上表皮生长因子受体抑制剂(EGFR-I)的投药的形式NSCLC的治疗得到更大有益效果的一类癌症患者的成员的系统。该系统包括：

(a) 质谱仪，生成基于血液的样本的质谱；以及

(b) 经编程计算机，包括处理单元以及存储来自质谱仪的质谱数据的存储器。存储器还存储：

1) 采取训练集合的非暂时数据，包括从作为被预测从癌症的治疗中的EGFR-I得到总生存有益效果的一类患者的成员的多个癌症患者(例如在VS 1.0测试中分类为“良好”的那些患者)所得到的类标记质谱数据，这一类患者还分为两个子类：

1. 在癌症的治疗中投药EGFR-I之后呈现疾病的早期进展的那些患者，这类患者的质谱数据具有“早期”或等效物的类标签；以及

2. 在癌症的治疗中投药EGFR-I之后呈现疾病的后期进展的那些患者(类标签“后期”或等效物)；

2) 程序代码，用于实现采取对训练集合的退出规则化和逻辑训练的经过滤微型分类器的组合(CMC/D分类器)的形式的分类器；

3) 程序代码，用于对在1) 所存储的质谱进行预定义预处理步骤，在已经执行对质谱的预处理步骤之后得到在多个预定义m/z范围的所述质谱中的所选特征的累积强度特征值；以及

4) 程序代码，将CMC/D分类器应用于在3) 得到的累积强度值和训练集合，并且作为响应而生成基于血液的样本的类标签，

其中，如果由程序代码4) 所生成的类标签对基于血液的样本是“后期”或等效物，则将该患者识别为与癌症的治疗中的化学疗法相比、可能从EGFR-I得到更大有益效果。

[0024] 在另一方面，描述一种供分类样本中使用的设备，其包括计算机存储器，计算机存储器存储采取包括从作为被预测从癌症的治疗中的EGFR-I得到总生存有益效果的一类患者的成员的多个癌症患者所得到的类标记质谱数据的训练集合的形式的非暂时数据，这类患者还分为两个子类：

1. 在癌症的治疗中投药EGFR-I之后呈现疾病的早期进展的那些患者，这类患者的质谱数据具有“早期”或等效物的类标签；以及

2. 在癌症的治疗中投药EGFR-I之后呈现疾病的后期进展的那些患者，这类患者具有类标签“后期”或等效物。

[0025] 在又一方面，公开一种治疗NSCLC患者的方法，包括下列步骤：向NSCLC患者投药EGFR-I，其中通过在经编程计算机中运行分类器(其将质谱仪从NSCLC患者的基于血液的样

本所产生的质谱数据与包括从通过基于血液的样本的质谱测定来确定为被预测从癌症的治疗中的EGFR-I得到总生存有益效果的一类患者的成员的多个癌症患者所得到的类标记质谱数据的训练集合进行比较),预测患者与化学疗法相比更多地获益于EGFR-I,这类患者还分为两个子类:

1.在癌症的治疗中投药EGFR-I之后呈现疾病的早期进展的那些患者,这类患者的质谱数据具有“早期”或等效物的类标签;以及

2.在癌症的治疗中投药EGFR-I之后呈现疾病的后期进展的那些患者,这类患者的质谱数据具有类标签“后期”或等效物。

附图说明

[0026] 图1是示出用于生成CMC/D分类器的方法的流程图。

[0027] 图2是示出用于使用按照图1所生成的CMC/D分类器来测试生物样本的测试方法的流程图。

[0028] 图3是类标签的初始指配并且在NSCLC/EGFR-I CMC/D分类器中划分为训练和测试集合的图示。

[0029] 图4A-4F是在CMC/D分类器生成方法(图1中的步骤1134)中生成的PFS和OS的测试集合的早期和后期分类之间的风险比(HR)的分布的图表。图4A-4B用于初始类标签的PFS和OS,而图4C-4F用于频繁错误分类的测试样本的类标签的一个或两个翻转(flip)之后的PFS和OS。

[0030] 图5是形成集合与通过从同一参考样本所得到的等式2的一致性的特征的谱的后续返回之间的特征值比的图表。

[0031] 图6A-6D是示出具有从生成集合谱所指配的标签的NSCLC/EGFR-I CMC/D分类器生成集合的患者的时间-事件结果的Kaplan-Meier曲线。图6A示出吉非替尼治疗患者的OS;图6B示出吉非替尼治疗患者的PFS,图6C示出化学疗法治疗患者的OS,以及图6D示出化学疗法治疗患者的PFS。

[0032] 图7是应用于PROSE样本集合的NSCLC/EGFR-I CMC/D分类器的灵敏度校正的回归曲线的图表。

[0033] 图8A和图8B是采用厄洛替尼(图8A)和化学疗法(图8B)所治疗的患者的编组后期和早期/未知(原始VeriStrat测试中测试VeriStrat良好的那些患者)的总生存的Kaplan-Meier图表。

[0034] 图9A和图9B是采用厄洛替尼(图9A)和化学疗法(图9B)所治疗的患者的编组后期和早期/未知(原始VeriStrat测试中测试VeriStrat良好的那些患者)的无进展生存的Kaplan-Meier图表。

[0035] 图10是通过治疗分类为VeriStrat不良和后期的患者的总生存的Kaplan-Meier图表。

[0036] 图11是通过治疗的VeriStrat良好早期/未知编组中的OS的Kaplan-Meier图表。

[0037] 图12A是通过治疗的后期编组中的OS的Kaplan-Meier图表;图12B是通过治疗的后期编组中的PFS的Kaplan-Meier图表。

[0038] 图13是图2的平均工作流程模块1206的图示。

[0039] 图14是图2的预处理工作流程模块1212的图示。

[0040] 图15是将主分类器应用于校正测试样本特征值和有噪特征值实现的图2的模块1228和1234的图示。

[0041] 图16是包括存储分类器和训练集合的计算机以及得到基于血液的样本的质谱的质谱仪的测试样本处理系统的图示。

具体实施方式

[0042] 以下描述分为四个一般小节：

小节I：描述形成本文称作CMC/D分类器的分类器的新方式(具有退出的微型分类器的组合)。这种新方式在用于创建对进行本文档的测试方法有用的分类器。

[0043] 小节II：描述按照小节I所形成的特定CMC/D分类器(其用于本文档所述的预测测试中)以及示范其预测与化学疗法相比、来自EGFR1-1的NSCLC患者有益效果的结果。

[0044] 小节III：描述用于使用小节II所述的CMC/D分类器对患者样本进行测试的当前优选测试方法。

[0045] 小节IV：描述用于生成分类器并且进行小节I-III所述的测试的实际计算环境。

[0046] 小节I 一般CMC/D分类器形成

本小节描述CMC/D、其生成或形成和优点。一般来说,当受到可用于生成分类器的样本数量限制时,CMC/D分类器是特别适用的。此外,CMC/D分类器在性质上是真正多元的,并且其优点在于它们避免过适于可用样本集合。

[0047] 与在大训练数据集可用时集中于形成分类器的机器学习的标准应用相反,在生物生命科学中,大数据难题、问题设定是不同的。这里的问题是,可用样本的数量通常由于临床研究而受到限制,并且属性的数量通常超过样本的数量。不是从许多实例来得到信息,在这些深层数据问题中,而是尝试从单独实例的深层描述来获得信息。

[0048] CMC/D分类器形成包括第一步骤a):从多个样本得到用于分类的数据,即,测量数据反映样本的某个物理性质或特性。每个样本的数据由多个特征值和类标签组成。这个集合在本文中称作“形成集合”或“形成样本集合”,参见图1的1100。例如,数据可以从使样本经过某种形式的质谱测定、例如MALDI-TOF所得到的、采取特征值(在多个m/z范围/峰值/特征的峰值强度)以及指示样本的某个属性的标签的形式的质谱测定数据。这个标签可能具有诊断或治疗属性,例如诊断标签(癌症/非癌症)、即样本是否来自获益于某种特定药物或者药物的组合(有益/无益)的患者或者指示样本的另外某种性质或特性、例如患者具有疾病的早期或后期复发、具有良好或不良总生存等的标签。类标签能够按照某种自动化方式在先前指配,或者可由人类操作员在分类器的形成时或之前指配。类标签也能够对分类器形成过程的多次迭代重新定义,换言之,类标签结合分类器本身的形成来定义。

[0049] 该方法继续进行步骤b):使用来自样本、一直到预先选择集合大小s(s=整数1...n)的特征值的集合来构成多个单独微型分类器。例如,多个单独微型或原子分类器可能使用单一特征(s=1)或一对特征(s=2)或三个特征(s=3)或者甚至包含3个以上特征的高阶组合来构成。s的值的通常选择将足够小,以允许实现方法的代码在适当时间量中运行,但是在一些情况下或者在较长代码运行时间是可接受的情况下可能较大。

[0050] 该方法继续进行过滤步骤c),即,测试单独微型分类器的每个的性能、例如精度以

分类多个样本的至少一部分,或者根据另外某种量度(例如在通过临床试验的实验和控制臂中的训练集合样本的单独微型分类器的分类所定义的编组之间所得到的风险比(HR)之间的差)来测量单独微型分类器性能,并且仅保留其分类精度或者其它性能量度超过预定义阈值的那些微型分类器,以达到微型分类器的过滤(截取)集合。在这个步骤中,微型分类器的每个将形成样本集合中的微型分类器的特征的(一个或多个)特征值(例如在预定义 m/z 范围中的累积强度值与类标记测量数据的训练集合中的样本的相同特征值进行比较。在这个步骤中,微型分类器对来自基于相对于训练集合中的相同(一个或多个)特征为微型分类器、例如K最近邻分类算法(KNN)所选的特征的(一个或多个)特征值的所选样本集合的样本的数据运行分类算法,并且输出样本的类标签。如果微型分类器过滤的所选性能量度是分类精度,则产生于分类操作的类标签可与提前已知的样本的类标签进行比较。但是,其它性能量度可被使用,并且使用产生于分类操作的类标签来评估。仅保持在分类的所选性能量度下适当地顺利执行的那些微型分类器。可使用备选监控分类算法,例如线性判别、决策树、概率分类方法、基于余量的分类器、例如支持向量机以及训练来自标记训练数据的集合的分类器的任何其它分类方法。

[0051] 为了克服通过某种一元特征选择方法根据子集偏置所偏置的问题,我们将所有可能特征较大比例看作是这个步骤中的微型分类器的候选。然后使用一直到预先选择大小(参数 s)的特征集合来构成所有可能KNN分类器。这给予我们许多“微型分类器”:例如,如果我们开始于各样本的100个特征,则从这些特征对($s=2$)的所有不同组合来获得4950($100 \times 99/2$)个“微型分类器”、使用三个特征($s=3$)的所有可能组合来获得161700个微型分类器,依此类推。探索可能微型分类器以及定义它们的特征的空间的其它方法当然是可能的,并且可用来代替这种分层方式。当然,许多这些“微型分类器”将具有不良性能,并且因此在过滤步骤c)中,仅使用通过预定义性能标准的那些“微型分类器”。这些标准根据特定问题来选择:如果具有二类分类问题,则仅选择其分类精度超过预定义阈值的那些微型分类器。我们选择在某种程度上是预测性的那些分类器,即,其中后期与早期复发编组之间的风险比(HR)在治疗臂中比在控制臂中要小某个预先指定值。甚至通过“微型分类器”的这种过滤,我们也以具有跨越从边界线达到优良性能的整个范围的性能的数千“微型分类器”候选而告终。(在典型示例中,存在数千个这类微型分类器,其通过过滤测试,并且用于具有退出的逻辑训练。)

该方法继续进行步骤d):使用规则或规则化组合方法来组合经过滤微型分类器。在其一个可能示例中,这个步骤涉及对样本的分类标签重复地进行在步骤c)所生成的微型分类器的过滤集合的逻辑训练。这通过因执行来自微型分类器的过滤集合的极端退出而随机选择经过滤微型分类器的一小部分并且对这类所选微型分类器进行逻辑训练来实现。虽然在实质上与标准分类器组合方法(参见例如Tulyakov等人的“Review of Classifier Combination Methods”,*Studies in Computational Intelligence*, Volume 90, 2008, 第361-386页)相似,但是具有某些“微型分类器”可能只通过随机机会是人为完善的并且因此主导组合的特定问题。为了避免这种过适于特定主导“微型分类器”,我们通过仅随机选择这些逻辑训练步骤的每个的“微型分类器”的一小部分,来生成许多逻辑训练步骤。这实质上是如深度学习理论中使用的退出的问题的规则化。在这种情况下,在我们具有许多微型分类器和小训练集合的情况下,我们使用极端退出,其中在各迭代中退出超过99%的预先过

滤微型分类器。

[0052] 可能使用的、用于在步骤(d)执行规则化组合方法的其它方法包括：

(基于Tikhonov regularization, Tikhonov, Andrey Nikolayevich (1943)。“Об устойчивости обратных задач”[关于逆问题的稳定性]。Doklady Akademii Nauk SSSR 39 (5):195-198。)

Lasso方法(Tibshirani, R. (1996), Regression shrinkage and selection via the lasso, J. Royal. Statist. Soc B., Vol. 58, No. 1, 第267-288页)。

[0053] (Nitish Shrivastava, “Improving Neural Networks with Dropout”, Master’s Thesis, Graduate Department of Computer Science, University of Toronto; 在 http://www.cs.toronto.edu/~nitish/msc_thesis.pdf 可得到。)

(Girosi F. 等人, Neural computation, (7), 219 (1995)。通过引用将上述发表物结合到本文中。

[0054] 该方法继续进行步骤e): 在执行规则化组合方法d)之后、例如在逻辑训练和退出迭代之后, 从微型分类器的过滤集合的组合来生成主分类器。在一个实施例中, 这个主分类器是对在步骤d)所述的退出期间所选的过滤微型分类器的那些集合的所有逻辑回归训练的平均。能够针对测试集合划分或者形成集合的子集来评估最终分类器, 该评估还对形成集合到训练和测试集合的多个不同划分来执行, 以及能够通过选择产生于具有“典型”性能的特定训练和测试集合划分其中之一, 或者备选地通过保留来自各训练和测试集合划分的所有主分类器并且使用来自每个主分类器的多数表决将标签指配给被测样本, 来生成最终分类器。这种方式实质上与“退出”规则化、即深度学习社区中用来对神经网络训练添加噪声以避免在目标函数的局部极小中被俘获的方法相似。参见Nitish Shrivastava, “Improving Neural Networks with Dropout”, Master’s Thesis, Graduate Department of Computer Science, University of Toronto; 在 http://www.cs.toronto.edu/~nitish/msc_thesis.pdf 可得到。)我们的方法还能够从整体学习方式来看(参见例如“Ensemble Methods”, Zhi-Hua Zhou, CRC Press, 2012 Boca Raton)。这类方式在避免过适并且增加生成可一般化测试、即能够在单独样本集合中验证的测试的可能性方面有希望。

[0055] CMC/D分类器生成方法和上述所产生分类器具有许多实际优点和用途。在分类形成中、具体来说在健康科学、例如癌症研究或药物研制中, 研究人员常常面临仅具有小样本集合可用的问题, 其在将要沿用对分类器形成的标准方式时引起极小训练和测试集合。例如, 在药物效能研究的样本集合中, 训练集合可由来自治疗臂的也许20个样本组成以及在还将控制臂分为训练和测试集合时由相似大小的训练集合组成。这仅引起通过某个训练标签指配、例如早期或后期所定义的早期和后期复发编组(参见下文)中的大约10个样本。标准方式开始于调查特征(例如质谱测定数据中的峰值), 并且选择表明包含与训练类相干的信息的某种希望的那些特征。这些然后使用k最近邻方法来组合, 以生成多元测试。对于小样本大小, 如同这个示例中一样, 多元测试的构造中包含的特征的选择能够易于由主要因样本到训练和测试集合的特定划分而表明差别能力的某些特征来主导。换言之, 使用一元p值来选择特征对较小样本大小变得较小信息性, 因为p值本身变得较小信息性。人们可能通过检验许多训练/测试集合划分情形来克服这个问题, 但是看来不是避免拣选这些情形的每个的专用特征的实际方式, 这使所形成测试的一般化性能的估计比较困难。在先前工作

中,我们开发了复杂交叉验证技术,其表明给定样本集合允许预测测试的开发的充分希望。但是,这个工作产生许多分类器候选,以及进一步验证的特定分类器的选择仍然比较困难。

[0056] 我们开发了本文所述的方法,其解决两个问题:(a)它不依靠在多元测试中包含的特征的特定选择,以及(b)通信组合许多、甚至数千个可能分类器候选,它提供自动生成一个单一顺利执行分类器(测试)的部件。

[0057] 我们套用术语“具有退出的微型分类器的组合”CMC/D来表示本文档所述的分类器生成方法。下面在节点II和III中说明将CMC/D应用于VS 1.0测试的创建中使用的质谱测定数据集。CMC/D使我们与较小训练集合配合工作,并且因此允许将样本集合分为训练集合和测试集合。这缓解在一些分类器形成问题中会具有的忧虑,即缺乏单独测试集合。CMC/D还允许分类器性能与特定测试/训练划分的相关性的调查,这可能导致小样本集合的偏置。最后,CMC/D对各训练/测试集合划分产生一个主分类器/测试。虽然这个测试可能不是最佳的(其可在给定数据的情况下构成),但是这种测试通过构造将不易受到过适于训练集合数据中的某个人工产物的危险。

[0058] 由CMC/D所生成的分类器由于使用在方法的步骤d)的“微型分类器”的组合中的逻辑回归而实际上是概率的。在给定样本数据的情况下,将CMC/D分类器应用于特定样本测量数据(例如质谱)的结果给出特定类(编组)标签的可能性,在这种情况下为早期或后期。

[0059] 图1是示出本小节中更详细描述的分类器形成过程的流程图。分类器形成过程通常在采取存储数据(例如采取质谱测定数据和实现图中所示模块的可执行代码形式)的分类器形成集合的通用计算机形式的计算系统中实现。

[0060] 如图1所示,该过程开始于数据的分类器形成集合1100、例如从质谱仪(未示出)自人类患者的基于血液的样本所得到的质谱测定数据的集合。图1的流程图所示的过程并不局限于任何特定形式的数据,如先前所述。但是,基于血液的样本的质谱测定的示例适合于本论述,而决不是意在进行限制。在本示例中,从作为被预测从NSCLC癌症的治疗中的EGFR-I得到总生存有益效果的一类患者的成员的患者的、例如VS 1.0状态为“良好”的患者来得到基于血液的样本。这些样本的类标签进一步分为两个子类,即,如下面所述的早期和后期。

[0061] 在步骤1102,定义分类器形成集合1100中的编组(类标签),例如分别为“早期”和“后期”编组1104、1106。在这个示例中,“早期”编组1104由与在抗癌药物的投药之后具有疾病的较早期进展或复发的患者关联的形成集合110中的谱的集合组成。相反,“后期”编组1106由与抗癌药物的注入之后的疾病的较后期复发或进展关联的形成集合1100中的谱的集合组成。下面详细描述在定义早期和后期编组中的其它考虑因素。形成集合110分为早期和后期编组可以分为具有偶数样本的编组。

[0062] 在步骤1108,早期和后期样本编组均分为训练和测试集合。在步骤1108的这种划分不一定分为相等编组。我们可按照2:1或其它比率来划分。如果我们具有极大集合,则可能不希望使用真正较大的训练集合。如果我们具有极受限数量的样本,则可能在训练集合中使用比在测试集合中更多的样本。在1108的这种划分产生两个编组:训练集合1112和测试集合1110(各训练和测试集合包括来自形成集合1100的“早期”和“后期”样本/数据)。

[0063] 如图1所示,训练集合1112然后经过分类器形成步骤1120、1126和1130。在步骤1120,创建多个基于KNN的微型分类器,如先前在以上详细说明。这些微型分类器可以仅使用质谱数据集中的1($s=1$)或者也许2个特征($s=2$)供分类。如气球1122所示,KNN微型分类器

使用比较从整个特征空间所提取的特征(m/Z特征的累积强度值,如框1124所示)的子集。如框1124所示,这些计算特征是谱中的m/Z范围。质谱可采取如我们的先前专利申请序号美国序号13836436(2013年3月15日提交)所述的“深层MALDI”谱的形式,通过引用也结合到本文中。备选地,质谱可采取来自比如2000激光发射的典型“稀释和上样”谱或者具有在谱获取时的谱过滤的实现的若干(例如三个)2000发射谱的各平均的形式。用于微型分类器中的分类的特征是累积强度值,即,在所指定m/Z范围中的预定义峰值位置之下的区域。KNN微型分类器中的分类的累积强度值的生成优选地在已经执行预处理步骤、例如谱的背景减法、归一化和对齐之后执行。这些步骤以及KNN微型分类器的实现通过通用计算机中的计算机代码来执行。

[0064] 在步骤1126,执行在步骤1120所生成的KNN微型分类器的过滤,以便仅保存具有可接受等级的性能的那些微型分类器。这在图1中直观地说明。能够存在作为良好、不良以及各一的特征的重叠集合。特征集合能够重叠,以及一部分将通过过滤而一部分将不会通过。相对于所定义性能量度来评估各微型分类器。在这个步骤中,仅保留具有良好分类性能的那些微型分类器,如在1128的加号所指示。

[0065] 在步骤1130,从在执行规则化组合方法、例如许多逻辑回归和退出规则化迭代之后通过过滤步骤的微型分类器来生成主分类器,如上所述。更详细来说,各微型分类器的结果是两个值其中之一,即“早期”或“后期”。然后能够通过经由标准逻辑回归(参见例如http://en.wikipedia.org/wiki/Logistic_regression)定义得到“早期”的概率,使用逻辑回归来组合实质上的逻辑回归的微型分类器的结果。

[0066] 等式(1)

$$P(\text{"early"} | \text{feature for a spectrum}) = \frac{\exp\left(\sum_{\text{Microclassifiers}} w_{mc} I(\text{mc}(\text{feature values}))\right)}{\text{Normalization}}$$

其中,如果应用于样本的特征值的微型分类器mc返回“早期”,则 $I(\text{mc}(\text{feature values})) = 1$,以及如果微型分类器返回“后期”,则为-1。权重 w_{mc} 是未知的,并且需要对训练集合中的所有样本、对公式的左边、分别对训练集合中的早期标记样本使用+1以及对后期标记样本使用-1、从上式的回归拟合来确定。由于我们具有比样本要多许多的微型分类器并且因此具有要多许多的权重、通常数千个微型分类器以及只有数十个样本,所以这种拟合将始终引起接近完善的分类,并且能够易于由可能随机地很好地拟合特定问题的微型分类器来主导。我们不希望最终测试由仅对这个特定集合顺利执行但不能完全一般化的单一特殊微型分类器来主导。因此,我们设计一种方法来规则化这种行为:不是一个总回归将所有微型分类器的权重同时拟合到训练数据,我们仅将微型分类器的几个用于回归,但是多次重复进行这个过程。例如,我们随机挑选微型分类器的三个,执行其三个权重的回归,挑选三个微型分类器的另一个集合,并且确定其权重,以及多次重复这个过程,从而生成许多随机挑选、即三个微型分类器的实现。定义CMC/D分类器的最终权重则是对这类实现的权重的平均。实现的数量应当足够大,使得各微型分类器很可能在整个过程期间被挑选至少一次。这种方式实质上与“退出”规则化、即深度学习社区中用来对神经网络训练添加噪声以避免在目标函数的局部极小中被俘获的方法相似。这个主分类器可实现为在逻辑回归和退出规则化之后的过滤分类器的组合的平均。形成这个主分类器(MC)的数据集在1132来指示,并且存储在运行图1所示方法的计算机的存储器中。

[0067] 在步骤1134,在步骤1130所生成的主分类器的性能然后通过由主分类器使形成集合数据(1110)的测试集合划分经过分类来测试。(测试集合再次可在主分类器中的分类算法的执行之前经过预处理步骤。)主分类器的性能的结果被存储,并且能够例如表示为风险比分布的直方图,如图1中在1138所示或者在先前描述中所示。

[0068] 步骤1108、1110、11128、1120、1126、1130、1132和1134如循环1136所指示重复进行,其中具有早期和后期样本集合到不同训练和测试集合实现的不同划分。循环1136的目的是避免训练集合/测试集合划分偏置。循环1136的各迭代的结果是不同主分类器。对训练和测试集合划分的各实现的各样本测试集合(1110)来评估主分类器的性能。

[0069] 在步骤1136,分析来自各训练/测试集合划分的分类器性能数据(例如HR直方图)。例如,如图1在1138所示,训练/测试集合划分的各实现产生主分类器,以及能够创建由许多主分类器所产生的分类(早期/后期)的风险比的直方图。风险比的分布能够用来评估分类器性能,如先前所述。将会注意,通过规则化步骤(1132)以及从具有典型性能的主分类器之一选择主分类器或者例如使用来自所有主分类器的多数表决算法对所有主分类器求平均或者将加权应用于所有主分类器,来使最终主分类器过适于训练数据为最小。在分析步骤1136中的最终分类器性能估计的置信度通过具有相似良好性能的许多主分类器的观测来增强。

[0070] 可存在训练集合中的特定样本(通常为少数)常常被主或最终分类器错误分类的情况。在这种情况下,重新定义这类样本的训练标签、例如将标签从“早期”改变成“后期”会有用的。这对其中训练标签例如在治疗有益效果或相对治疗有益效果的测试中难以定义的分类问题是特别相干的。这在步骤1142进行,并且该过程返回到步骤1102,以及按照校正训练标签将形成集合分为“早期”和“后期”编组继续进行。在流程图的步骤1108和后续步骤将这些编组分为训练和测试集合划分的过程继续进行,从而产生在步骤1136和1138的主分类器性能的新评估。步骤1140不一定始终是必需的,例如在极少或没有错误分类的实例的情况下,在这种情况下,在分析步骤1136之后,处理直接进行到步骤1144。

[0071] 在步骤1144,定义用于定义待测试样本的最终测试标签的过程。样本的最终测试表标签能够按照若干方式来指定,例如它能够定义为对来自所有训练/测试集合划分的所有最终主分类器的分类标签的多数表决的结果。备选地,它能够定义为由对提供典型性能的给定训练/测试集合划分的所选主分类器或者通过使用主分类器例如使用下列小节所述的过程所产生的分类结果的统计分析所产生的标签。

[0072] 小节II 从患者基于血液的样本的质谱测定来生成CMC/D分类器供EGFR-I药物(VS 2.0)的NSCLC患者选择

在这个小节中将描述对指导NSCLC患者的治疗有用的CMC/D分类器的生成的示例。分类器的生成主要遵循以上小节1所述并且在以上图1的论述中所述的方法。但是,处理测试样本以使用这个示例中的CMC/D分类器进行预测利用参考谱以及对谱的处理的附加调整以考虑对存在的机器证明和谱再现性的限制。被测样本的最终分类标签的生成还利用下面将结合图2更详细描述的特征相关噪声特性和其它技术。然后,本小节将示范从质谱数据生成CMC/D分类器及其用来在治疗之前对NSCLC患者是否可能获益于EGFR-I药物的投药进行预测的另一示例。

[0073] 在我们不知道‘正确’类标签是什么的问题开始(分类器的生成)的意义上,分类问

题是不寻常的。在某种意义上,这个问题更可能是无监控学习问题。我们通过开始于类标签的初始猜测、训练这些猜测的测试并且迭代这个过程以细化类标签来解决这个问题。这个过程输出的是最终类标签,并且该算法从患者的样本检测这些类。

[0074] 如本文档先前所述,先前美国专利7736905(在本文中偶尔称作“VS 1.0”)等中所述的VeriStrat测试在治疗之前进行关于NSCLC患者是否为称作VeriStrat“不良”、不可能获益于NSCLC的治疗中的EGFR-I、例如厄洛替尼和吉非替尼的一类的成员的预测。预测基于来自患者的基于血液的样本的质谱以及计算机中实现的分类器的使用。来自NSCLC的治疗中的最近EGFR-I试验、称作TAILOR和DELTA试验的结果指示厄洛替尼可能是EGFR野生型人口中的较差治疗。因此,特罗凯(厄洛替尼)的使用超出了其肿瘤表明EGFR敏化突变的患者的前线治疗,并且作为高线的挽救治疗。

[0075] ‘905专利中所述的测试没有描述如何进行EGFR-I、例如厄洛替尼是否为优于化学疗法的优良治疗的预测,甚至在VS 1.0测试中的测试VeriStrat“良好”的那些患者中。后续研究、例如PROSE研究[1]没有设计成表明一种治疗优于另一种治疗的优势。此外,虽然PROSE研究中的少量VeriStrat“良好”患者到目前为止过少而无法表明厄洛替尼和化学疗法治疗的等效性,但是从PROSE研究也没有一个治疗优于另一个治疗的证据。

[0076] 本发明人开发我们的新CMC/D分类器形成方法并且将其应用于这个问题。在更深层地探究血清蛋白质组的方式的开发期间,使用所谓的“深层MALDI”,我们还开发了工具和算法,以通过组合来自标准获取、例如VS 1.0测试中使用并且在美国专利7736905中所述的标准“稀释和上样”质谱数据获取的多个技术复制的谱,来增加增强标准质谱获取技术的峰值内容的能力。本小节中描述来自标准“稀释和上样”质谱获取的多个技术复制的谱的这种组合的示例。

[0077] 最近分类工作的目标是开发新测试(在本文中称作VeriStrat 2.0或VS 2.0),其识别具有来自厄洛替尼的比化学疗法要大的有益效果的一组NSCLC患者。在本文档中描述这个新测试以及生成测试中使用的分类器的方法。一测试的一个可能实现中,测试基于标准MALDI-ToF质谱获取,例如2000发射“稀释和上样”谱。作为分类器形成集合(图1,1100),我们具有可用的来自在生成’905专利的VS 1.0测试中使用的原始形成集合和初始验证集合的样本的子集。我们对形成该形成集合的那些样本选择在VS 1.0测试下测试VS良好的那些患者、即被预测从EGFR-I得到总生存有益效果的那些患者。如本文档所述的所产生测试表明在所选子集中优于化学疗法的厄洛替尼的优势,同时保留VS 1.0测试的预测特性。本文档所述的测试说明如何识别NSCLC患者是否为可能从EGFR-I、例如厄洛替尼得到比化学疗法要大的有益效果的患者的这个子集的成员。这个子集在以下论述中与类标签“后期”关联。类标签可被给予另外某种等效名称,以便识别这类患者,例如“EGFR有益效果”、“阳性”、“+”等。因此,类标签的特定名称并不重要。因此,在本公开和权利要求书中,当我们说“后期或等效物”或者“早期或等效物”时,意味着类标签的名称的选择并不重要。

[0078] 本文档所述的测试可能可选地包含分类算法,其中识别为不良等的患者被预测没有获益于NSCLC癌症的治疗中的EGFR-I。第三类标签能够被指配给被测患者样本,在这里称作“中间”,其与被预测按照临床有意义项对化学疗法(多西他赛、培美曲塞)或者EGFR-I、例如吉非替尼或厄洛替尼类似地表现的患者关联。

[0079] 患者人口和可用样本

患者的下列组群具有可用于这个项目的样本：称作“*Italian A*”、“*Italian B*”、“*Italian C*”的样本集合。*Italian A*和*B*是患有晚期NSCLC、采用原始VeriStrat测试的开发和验证中使用的吉非替尼所治疗的患者的组群。一般参见美国专利7736906；Taguchi等人，JNCI 99:838-846(2007)。*Italian C*是在高级线采用多种化学疗法所治疗的患者的组群。

[0080] 初始计划是直接创建预测分类器，以通过使用患者的全部三个组群来识别与化学疗法相比对吉非替尼具有更好结果的患者。但是，由于对其无进展生存(PFS)数据是可用的患者的子集中的*Italian C*组群的总结果比*Italian A*和*B*组群要差，所以这种方法没有正常工作。

[0081] 使用所有样本来创建识别对吉非替尼疗法具有良好结果的患者的分类器的初始工作产生许多分类器，其产生具有与始终VeriStrat分类的极强重叠的分类，即，我们能够产生具有相似性能并且与使用CMC/D方法和不同特征的原始VeriStrat相比产生极相似的样本分类的许多分类器。甚至在从该过程排除与来自VeriStrat的质谱特征重叠的谱的区域中的特征时，这种情况也成立。

[0082] 因此，决定将分类器构造过程限制到由在VS 1.0测试中产生始终“VeriStrat良好”分类的样本组成的分类器形成样本集合，即，设计将VeriStrat良好样本分为对EGFR-I具有更好或更差结果的患者的分类器。最后，由于存在认为具有性能状态(PS)2的患者和疗法的第四线中的患者一般可能接受来自吉非替尼疗法的极少有益效果的原因，所以来自这些患者的样本也没有包含在分类器训练中。来自三个组群的其它样本、包括来自原始形成集合的VeriStrat不良样本、来自*Italian C*组群的样本以及来自具有PS 2并且在第四线疗法的患者的样本仍然在形成过程期间用于分类器评估中。此外，在本小节稍后描述的CMC/D分类器的临床应用中，用于分类的训练集合包括来自具有类标签VeriStrate不良的患者的谱的特征值。

[0083] 在分类器形成期间使用的样本的列表在附录A中给出。

[0084] 新CMC/D分类器的形成在图1所示的简图中示出。以上详细论述该简图。基本上，并且作为第一近似，形成样本集合(附录A)根据与样本关联的患者分为两个编组(“早期”和“后期”)在采用EGFR-I的治疗开始之后遭受疾病的早期或后期进展。参见以下所述的图3。遭受后期进展的那些患者对类标签的初始指配能够被认为是从EGFR-I治疗比备选方案、例如化学疗法获益要多的并且对其标本指配类标签“后期”的那些患者。遭受早期进展的那些患者能够作为初始估计被认为是没有从EGFR-I比化学疗法获益更多的并且对其标本指配类标签“早期”的那些患者。

[0085] 从这两组样本，将编组分离为近似相等大小的训练和测试集合(图1，步骤1108)。训练集合使用其血清样本的MALDI-ToF谱中的特征经过图1右边所示的CMC/D分类器生成步骤1120、1126、1130、1134。测试样本由所产生主分类器(MC)来分类，以及MC性能在步骤1134对样本的测试集合来评估(1110)。该过程对许多训练/测试集合划分实现(在这个示例中为250)循环。经过错误分类的样本被给予重新定义训练标签，以及CMC/D分类和评估步骤重复进行(步骤1140、1142)。这个标签重新定义过程在这个测试的开发中重复进行两次。最终分类器然后从MC来选择，在这种情况下，全部250个分类器的多数表决产生训练/测试划分的每个。最终分类器的备选构造也是可能的，例如提供“典型”性能、250个MC的平均等的一个MC的选择。

[0086] 谱获取和预处理

图1的分类器生成中使用的质谱通过Bruker质谱仪从基于血液的样本来获取。质谱在分类之前经过预处理步骤。在本小节中描述步骤。

[0087] a. 在形成期间使用的质谱的生成

基于血液的样本的谱获取使用用于VeriStrat测试的合格质谱测定机器来执行(对于细节,参见附录H)。机器证明能够使用J. Röder等人的US Patent No. 8467988的专利的方法来执行,通过引用将其内容结合到本文中。

[0088] 谱在2000所获取发射谱的三乘中获取。在这个特定情况下,谱在获取时使用Bruker Flexcontrol设定来过滤,以便仅获取具有预期质量的谱。样本经受的实际发射的次数高于2000,并且逐个样本以及逐个MALDI点改变。对各样本所获取的谱的三乘被对齐并且求平均,以产生每个样本一个6000发射谱。

[0089] b. 背景估计和减法

预处理平均谱中的第一步骤是背景估计和减法。平均谱的背景分量使用单窗口方法和100的乘法器来估计。然后从平均谱中减去估计背景。

[0090] c. 谱对齐

在任何质谱中,存在相对飞行时间数量到 m/Z 值的转换的细微差异。我们识别质谱的绝大多数中存在的峰值的集合,并且重新缩放每个谱的 m/Z 值,使得各单独谱中的公共峰值与参考集合的平方偏差尽可能小。这个过程引起接近(单位为 m/Z)特征的更好分辨率。

[0091] d. 归一化

为了得到区分临床编组的特征,需要测量来自不同样本的峰值的强度,并且将其值进行比较。电离蛋白质的总量在MALDI过程中不是可控的,并且因此我们只能测量相对峰值强度。为了这样做,需要归一化谱。为了避免在归一化期间传播峰值强度从本征可变的或者与患者的临床状态相关的峰值到稳定峰值的可变性,需要注意确定谱的哪些区域能够用于归一化。用于归一化的 m/Z 区域使用部分离子电流归一化工具来选择。本领域已知的部分离子电流归一化和感兴趣读者针对美国专利7736905中的归一化过程的论述。

[0092] e. 特征定义和特征表

为了定义能够区分临床编组的峰值的可能候选(即,KNN分类中使用的 m/Z 特征),我们定位预处理谱中的峰值,并且定义各峰值的最大数周围的单位为 m/Z 的范围。单位为 m/Z 的这些范围定义用于所有其它分析的特征。我们选择76个特征作为用于区分编组的可能候选,并且计算每个谱的这些特征的每个的累积强度。这样,得到每个谱的各特征的特征值。表这些累积强度(特征值)的列表清单(行为谱,列为特征)称作特征表,其存储在实现图1的方法的通用计算机的存储器中。由于缺乏重新检查的充分特征质量(噪声),在 $m/Z=7616$ 和 14392 所定义的特征的两个在CMC/D分类器形成过程期间没有使用。我们观察到,样本的一部分表明充分的氧化水平,从而引起双重峰值结构或者相似峰值的偏移。为了避免遗漏基本多肽的氧化形式,我们使用很广泛的特征定义。CMC/D分类器生成过程中使用的74 m/Z 特征的定义在附录B提供。

[0093] CMC/D分类器形成方法

早期/后期进展编组和测试集合的选择(图1的步骤1102和1108)

从临床数据,不可能肯定地确定哪些患者或多或少地获益于给定疗法。作为定义形成

集合的类标签的第一近似,我们决定在步骤1102(图11)将类标签定义为那些患者是否或多或少地获益于采用EGFR-I的治疗,PFS小于80天的患者定义为“早期”(指示从疗法的可能极少获益的早期进展),以及PFS超过200天的患者定义为“后期”(指示从疗法的可能更大获益的后期进展)。参见图3。这产生“早期”编组中的23个患者以及“后期”编组中的23个患者。这些在附录C中随其所指配类标签来列示。在图1的步骤1108,这些然后分为训练(11“早期”和11“后期”)和测试集合(12“早期”和12“后期”),通过疗法线和性能状态(PS)来分层。有可能的是,一些训练/测试划分能够产生对分类器的创建特别良好或不良的训练集合以及特别容易或难以分类的测试集合。因此,分层训练/测试划分随机进行250次(通过图1中的循环1136所指示)。在图1的步骤1130,各划分提供训练集合1112,引起CMC/D主分类器(MC)的生成,其性能能够对于对应测试集合来评估。(步骤1134) 为了提供表示根据PFS时间的分布的人口的测试集合,PFS在80与200天之间、具有PS 0或1并且在疗法的一线至三线的患者的一半随机选择以包含在测试集合1110中。类标签的初始指配和分为训练和测试集合在图3中示出。

[0094] 微型分类器的创建(步骤1120,图11)

对于给定训练集合,有可能使用74个特征的子集来创建许多单独K最近邻(KNN)分类器。通过训练集合和特征的特定子集中的样本所定义的这些单独KNN分类器定义“微型分类器”(mC)。对于这个项目,KNN算法中的K=5始终是固定的。

[0095] 考虑所有mC,其使用74个特征其中之一($s=1$)或者74个特征的一对($s=2$)。这给出各训练集合的总共2775个mC。

[0096] 微型分类器的过滤(步骤1126,图11)

在步骤1120所生成的微型分类器通过对训练集合执行mC、基于过滤来截取。这使用CMC/D过程的ERRORS方法、以 $J_{min}=0.7$ 和 $J_{max}=0.9$ 进行。这意味着,每个mC应用于其训练集合。计算它被指配“早期”和“后期”标签的精度。如果这个精度在0.7与0.9之间,则mC通过过滤,并且可用来制作主分类器(MC)。如果精度位于这个范围之外,则mC未通过过滤,并且从CMC/D过程中丢弃。通过过滤的MC的数量取决于训练集合、即特定训练/测试集合划分实现,但是通常大约为1000-1500。

[0097] 本质上,ERRORS方法评估mC所给出的分类的精度。在过滤过程中,每个mC应用于训练集合的各成员,并且这给予我们训练集合的各成员的分类。我们知道已经指配给训练集合的各成员的定义(类标签),因此我们只计算各微型分类器的正确分类的比例。我们拣选这个精度(正确分类的比例)必须位于0.7与0.9之间。

[0098] 我们有意不将上限(J_{max})推送到1.0的完善分类。首先,不存在取得这个精度的许多微型分类器,但是其次以及更重要地,我们在生成分类器时设法避免在该过程的各阶段的过适。取得异常高精度的微型分类器可能是‘特殊’而不是‘典型’的,产生于训练集合和特征的某些物质,而不可能顺利一般化。因此,我们选择不将‘过于良好’的微型分类器包含到主分类器中。相当感兴趣的是要注意,当过滤标准设置成过于极端并且组合具有异常良好性能的微型分类器时,所产生的总分类器结果具有更差性能。

[0099] 使用具有退出的逻辑回归来创建主CMC/D分类器(步骤1130)

通过使用具有极端退出的后期和早期训练集合标签作为规则化器以训练逻辑回归,将通过过滤的mC组合为一个主分类器(MC)。执行一万次退出迭代,在其每个中,5个mC使用逻

辑回归随机选择和组合。来自各退出迭代的每个mC的逻辑回归权重(参见上式1)经过求平均,以产生到最终MC的逻辑组合的最终权重。

[0100] CMC/D分类器性能评估(步骤1134,1136,图11)

一旦对给定训练/测试集合实现创建主分类器,则在步骤1134,通过对测试集合(1110)以及对于从Italian C组群的样本所得到的谱运行分类器来对它评估。这个过程对250个训练和测试划分的每个来执行。所评估的量包括测试集合的“早期”与“后期”分类之间的风险比(HR)以及对于Italian C组群的总生存(OS)和PFS以及对测试集合和Italian C组群的“早期”和“后期”分类的中值。所生成的PFS和OS的HR分布在图4A-B中示出。另外,类标记样本的单独分类在处于测试集合中时被检查。许多样本重复地被指配不匹配其PFS定义标签的分类。这些样本在表3中识别和列示。

[0101] 表1. 持续错误分类的样本

样本ID
ICA_11
ICA_12
ICA_18
ICA_20
ICA_21
ICA_22
ICA_36
ICA_38
ICA_39
ICA_45
ICA_51
ICA_68
ICB_22
ICB_3
ICB_38
ICB_49
ICB_61

初始类标签指配的细化(图1的步骤1140)

翻转表1所列示的、对许多训练/测试划分持续错误分类的样本的类标签(“早期”到“后期”以及“后期”到“早期”)。这产生将要再次执行的CMC/D分类器生成过程的训练标签的新集合。

[0102] 使用新标签,“早期”和“后期”样本再次随机化为训练和测试集合250次,如前所述在疗法的线和PS上分层。微型分类器如前所述来创建,并且使用相同标准来过滤。这些过滤mC使用具有退出的逻辑回归来组合,以创建MC,并且MC的性能对新测试集合来评估。所生成的PFS和OS的HR分布在图4C和图4D中示出。在再次翻转之后所生成的PFS和OS的HR的分布在图4D和图4E中示出。

[0103] 识别若干样本,其在测试集合的部分持续错误分类。这些在表2中列出。

[0104] 表2. 样本在类标签的第一集合翻转之后持续错误分类

样本ID
ICA_20
ICA_21
ICA_38
ICA_39
ICA_45
ICB_12
ICB_40

翻转表4所列示的、在CMC/D过程的第二运行之后持续错误分类的样本的类标签(“早期”到“后期”以及“后期”到“早期”)。这产生类标签的新集合,其再次随机化成训练和测试编组250次,通过疗法线和PS来分层。创建mC、过滤、组合成MC并且评估性能的整个过程第三次重复进行。在该过程的第三重复之后,只有两个样本在处于训练集合中时不良地分类,并且判定不要求进一步处理。

[0105] CMC/D过程的第三迭代的250个训练/测试划分的MC性能的分布在图4E-4F中示出。超过90%的训练/测试划分实现产生少于1但是多于实现的一半的测试集合的早期与后期分类之间的HR具有对PFS小于0.76以及对于OS小于0.78的HR。不是选择最终测试/CC/D分类器的这些单独训练/测试划分其中之一,最终分类器而是定义为第三CMC/D迭代的全部250个MC的多数表决。这具有不要求从具有特别有益测试或训练集合的特定训练/测试集合划分中选择主分类器并且还去除进行选择 and 潜在地提供更健壮最终分类器中的人为主观性的任何元素的优点。

[0106] 对于考虑对机器证明和谱再现性的限制的调整

上述最终分类器的实现以生成被测样本的类标签实现质谱数据处理中的某些调整,以便考虑对测试被开发时存在的机器证明和谱再现性的某些限制。在本小节中描述这些调整。稍后还结合图2来描述这个过程。本领域的技术人员将会清楚地知道,这些调整可以不是生成CMC/D分类器或者使用CMC/D分类器来实现预测测试所需的。本小节中描述的调整由用来生成质谱的质谱仪的某些限制以及还由增加测试的稳定的期望而产生。

[0107] A. 质谱仪的m/Z灵敏度的变化的校正

谱使用先前证明以执行原始VeriStrat测试的Bruker质谱仪机器、使用J. Roder等人的美国专利8467988所述的过程来获取。虽然原始VeriStrat 1.0测试仅使用5kDa与13kDa之间的特征,但是除了这个范围中的特征之外,本小节中描述的测试还使用具有更高和更低m/Z位置的特征。对原始VeriStrat测试所证明的谱仪必须具有用于原始测试的质谱特征的足够再现性,但是对这个范围外部的m/Z灵敏度不存在要求。

[0108] 均在先前证明机器上、在与生成用于本测试开发中使用的谱同时地从参考样本所生成的参考谱与在稍后时间从相同参考样本所生成的谱的比较指示,虽然m/Z灵敏度对于5kDa至13kDa特征范围之内的特征是相似的,但是在这个范围之外,m/Z灵敏度表明某些系统差。

[0109] 为了能够比较在不同时间或者在可用于按照这个新测试进行测试的水平的合格设定中的不同机器上生成的谱,需要对m/Z灵敏度的这些差校正特征值。这能够使用从在与

用于本测试开发的谱相同批次以及来自将要使用新VS 2.0测试来分类的患者样本的谱的后续批次中生成的单个参考样本所生成的参考谱进行。在这个示例中(如图2在1202A和1202B所示),参考示例是来自健康人类的血清样本。

[0110] 参考样本的两个准备运行三次,其中谱用于VS 2.0开发。这些三乘使用平均工作流程来求平均,并且使用预处理工作流程来预处理(参见以下图2的论述)。生成特征值,并且在两个准备之间比较特征值。为了避免使用来自一个或另一个准备的异常值特征值,特征减少到其特征值处于两个准备的相互10%之内的那些特征。如果FV1是参考样本的准备1的特定特征的特征值(1202A,图12)并且FV2是参考样本的准备2的同一特征的特征值(1202B,图12),则特征被认为适合于相对m/Z灵敏度的分析,若:

$$|1-(FV1/FV2)| < 0.1 \text{ or } |1-(FV2/FV1)| < 0.1. \text{等式2}$$

这些特征的特征值将与从VS2.0测试的样本的后续批次中的参考样本的准备所生成的相同特征的特征值进行比较。如果两个准备在后一批次中可用,理想地在将要经过VS2.0测试的样本之前和之后运行,则对于能够用于第二批次中的m/Z灵敏度比较的特征也应当满足等式2的阈值。如果参考样本的多于2个准备是可用的,则等式2能够一般化成使用从增加数量的谱可用的信息,使得特征值的标准偏差能够与各特征的平均特征值进行比较,以及能够使用其与平均数的标准偏差的比率低于设置阈值、例如0.1的特征。

[0111] 一旦识别具有适当再现性的特征的子集,则从样本的VS2.0形成批次到样本的任何后续批次的m/Z灵敏度的变化能够在作为m/Z的函数、形成批次(AVO)中的参考谱的平均特征值与后一批次(AVN)中的参考谱的平均特征值的比率的图表中检查。这种图表在图5中示出。

[0112] m/Z灵敏度的系统变化能够在图5中看到,其中与后一批次相比,形成批次在较高m/Z具有较低灵敏度以及在较低m/Z具有较高灵敏度。为了允许m/Z灵敏度的这种系统差的校正,直线在图5中拟合到数据,以及斜率和截距被确定。这给出一个函数,能够用以校正对后一批次中的任何样本所得到的各特征值,以使它与对VS2.0形成批次中的样本所得到的特征值是可比较的。

[0113] B.VS2.0分类与经由VS1.0样本操控和谱获取过程从血清样本获取质谱中固有的噪声的稳定性的分析

VS1.0是高度可再现测试,其中分类的再现性超过95%。获得测试中的再现性的一种方法是谱生成的样本的三元点位的使用以及VS1.0分类的生成之前的三元标签的比较。由于来自样本的三元谱对VS2.0测试来求平均,所以VS1.0的冗余度丢失,并且这种方式无法扩展到VS2.0。但是,已经开发给定测试样本的多个复制的硅中生成的方法,其允许模拟VS1.0样本准备、点位和谱生成的过程中固有的样本和MALDI点相关、非系统不可再现(噪声)的效果。

[0114] 为了表征各特征的噪声,比较对于为VS1.0重新证明的质谱仪执行Italian A、B和C样本集合的两次运行。对于各VS2.0特征,跨再次运行来比较各样本的特征值。这产生每个VS2.0特征的一致性图表。对于各一致性图表,线性回归用来将直线拟合到特征值数据。为了表征围绕这个拟合的噪声,检查线性回归的残差。噪声被指配为主要加性或主要乘性的。对于加性噪声,噪声强度定义为残差的标准偏差。对于乘性噪声,各残差除以对应特征值,以及这个量的标准偏差定义为噪声强度。按照这种方式所估计的VS2.0特征的噪声类型和

噪声强度在附录D中给出。

[0115] 根据其类型和强度 σ 表征了各特征的噪声,具有测量特征值F的各样本的各特征的有噪实现可经由下式来生成:

$$\text{加性噪声: } F_{\text{noisy}} = F + \sigma \varepsilon \quad \text{等式(3)}$$

$$\text{乘性噪声: } F_{\text{noisy}} = F (1 + \sigma \varepsilon) \quad \text{等式(4)}$$

其中 ε 是具有零平均和单位标准偏差的高斯随机数。

[0116] 为了调查特定测试样本的噪声下的VS2.0分类的稳定性,各样本的特征表的160个有噪实现使用等式(3)、等式(4)和附录D中给出的各过滤器的噪声参数来生成。各有噪实现使用在上述CMC/D过程的最终迭代期间生成的250个MC来分类。这产生样本的各有噪实现的“早期”或“后期”的250个分类,即,每个样本40000个“早期”或“后期”分类。设跨250个主分类器的“早期”分类的总数为 N_{Early}^i ,以及跨250个主分类器的“后期”分类的总数为 N_{Late}^i ,其中 $1 \leq i \leq 160$ 。根据定义,对于所有 i , $0 \leq N_{\text{Early}}^i \leq 250$, $0 \leq N_{\text{Late}}^i \leq 250$, and $N_{\text{Early}}^i + N_{\text{Late}}^i = 250$ 。

[0117] 噪声效果估计器定义为:

$$\begin{aligned} \text{噪声效果估计器} &= N_{\text{Early}}^i / (|\sum_i N_{\text{Early}}^i - \sum_i N_{\text{Late}}^i| / 320) \text{的标准偏差} \\ &= \text{sqrt}(\sum_i (N_{\text{Early}}^i)^2 - (\sum_i N_{\text{Early}}^i)^2) / (|\sum_i N_{\text{Early}}^i - \sum_i N_{\text{Late}}^i| / 320) \\ &= \text{sqrt}(\sum_i (N_{\text{Early}}^i)^2 - (\sum_i N_{\text{Early}}^i)^2) / (|\sum_i N_{\text{Early}}^i - 20000| / 160) \text{等式(5)} \end{aligned}$$

这个“噪声效果估计器”将“早期”主分类器分类的数量的可变性与“早期”和“后期”主分类器分类的总数的差进行比较。如果与实现的“早期”和“后期”主分类之间的典型差相比、噪声实现产生“早期”分类的数量的低可变性,则噪声效果估计器将较小。如果与实现的“早期”和“后期”主分类之间的典型差相比、噪声实现产生“早期”分类的数量的大可变性,则噪声效果估计器将较大。

[0118] 其“早期”和“后期”主分类器分类的数量的差较大的样本能够在产生所返回VS2.0分类的变化之前容许充分可变性,而其这个差较小的样本仅以较小可变性经过所返回总分类的变化。因此,等式5所定义的噪声效果估计器提供关于样本对分类标签变化的敏感程度的量度。

[0119] 将这个过程应用于Italian A、B和C样本集合的两次运行以计算各样本的噪声效果估计器揭示,通过仅对噪声效果估计器低于0.5的阈值的样本返回VS2.0分类器分类,可对样本返回可靠分类。高于这个阈值,在返回被测样本的分类标签中存在充分不定性,并且应当报告中间/未知分类标签。

[0120] 最终分类器应用于形成集合中的样本

VS2.0最终分类器应用于形成集合中的所有样本。注意,这包括分类器的训练中包含的样本。形成集合样本的VS2.0分类在附录E中给出。注意,VS1.0分类为不良的所有样本被指派标签早期。对按照下列项编组的形成集合中的患者来绘制OS和PFS:后期、未知和早期(不包括VS1.0不良)以及图6中的VS 1.0不良。注意,Italian C组群中的若干患者具有OS数据但是没有PFS数据。图6是具有从形成集合谱所指派的标签的形成集合中的患者的时间-事件结果的图表;图6A:吉非替尼治疗患者的OS,图6B:吉非替尼治疗患者的PFS,图6C:化学疗法治疗患者的OS,以及图6D:化学疗法治疗患者的PFS。通过比较图6A和图6C,要注意,其样本测试后期的那些患者从吉非替尼得到比化学疗法要大的益处,如这些患者的总生存曲线

所指示。

[0121] 与图6的图表相关的生成统计在表3和4中提供

表3 与图6关联的中值

端点	编组	n	中值(天数)	95% CI(天数)
OS	后期GEF	32	457	259-680
OS	早期/未知GEF	53	243	144-304
OS	VS1.0不良GEF	44	96.5	60-162
PFS	后期GEF	32	208	90-287
PFS	早期/未知GEF	53	92	69-122
PFS	VS1.0不良GEF	44	61.5	43-83
OS	后期CT	3	80	55-92
OS	早期/未知CT	17	172	132-383
OS	VS1.0不良CT	12	141	60-250
PFS	早期/未知CT	14	78.5	40-113
PFS	VS1.0不良CT	10	82.5	29-93

表4 与图6关联的风险比和p值

端点	比较	对数秩p	CoX HR (95% CI)	CPH p值
OS	GEF:早期/未知与后期	0.025	0.59 (0.37-0.94)	0.027
OS	GEF:不良与后期	<0.001	0.30 (0.18-0.49)	<0.001
OS	GEF:不良与早期/未知	<0.001	0.49 (0.33-0.75)	<0.001
PFS	GEF:早期/未知与后期	0.018	0.58 (0.37-0.91)	0.018
PFS	GEF:不良与后期	<0.001	0.36 (0.22-0.60)	<0.001
PFS	GEF:不良与早期/未知	0.025	0.64 (0.42-0.95)	0.029
OS	CT:不良与早期/未知	0.217	0.61 (0.28-1.35)	0.221
PFS	CT:不良与早期/未知	0.477	0.74 (0.31-1.72)	0.479

来自Italian A、B和C的样本再运行两次。(在最后一次运行中,只有VS1.0良好样本再运行,并且几个样本因缺乏剩余样本容积而被省略。)跨三次运行的结果在附录F中概括。

[0122] 灵敏度校正连同硅内噪声分析一起引起可动作标签的良好再现性。在最后一次运行中运行的93个样本中,16个标记为后期,35个标记为早期,以及42个标记为未知。在第三运行中标记为后期的样本在先前运行中标记为后期或未知。在第三运行中标记为早期的样本在先前运行中标记为早期或未知。在第三运行中标记为早期的35个样本中的24个在全部三个运行中标记为早期。在第三运行中标记为后期的16个样本中的14个在全部三个运行中标记为后期。在第三运行中标记为未知的42个样本中的20个在全部三个运行中标记为未知。虽然未知的大比例不合乎需要,但是看来好像是,如果我们从VS2.0分析中称作标签早期(后期),则这个样本在另一运行中表征为早期(后期)或者称作未知。

[0123] 将最终CMC/D分类器应用于来自PROSE研究的样本

测试过程:单盲

上述最终CMC/D分类器经过对于在验证协议下从PROSE研究的可用样本所得到的质谱的测试。最终CMC/D分类器被认为在这个验证协议之前是固定的。将质谱提供给看不见其临

床数据的分析人员。如上所述来分析谱,以及生成所产生分类(附录G)。然后提供非单盲关键字,并且执行统计分析。

[0124] 测试过程:m/z灵敏度校正计算

分析连同PROSE谱一起生成的血清P2(参考)谱,以提供必要的m/z灵敏度校正。由于PROSE样本跨越5个批次,所以随各批次来收集血清P2的一个准备。对于5个独立准备,使用CV计算方式(以上概述)。PROSE数据的回归曲线在图7中示出。从这个曲线,得到Y轴截距和斜率值,如图7的小图表所示。

[0125] 结果的统计分析

从PROSE试验对样本所得到的VS2.0分类在附录G中列示。对统计分析仅考虑来自PROSE主要分析人口中的患者的样本。对于患者01_044和患者01_080,两个样本是可用的。具有标准标记的样本、而不是标记为‘second_sample’的结果用于统计分析。两个样本也可用于患者06_010,但是均具有VS2.0分类早期。没有样本可用于患者01_050、患者03_006、患者06_004、患者06_021、患者11_043、患者11_048和患者12_014。

[0126] 因此,样本是从PROSE按协议人口的263个患者中的256个可用的:148个分类为早期,39个分类为后期,以及69个分类为未知。分类为后期的所有样本与具有VS1.0良好分类的患者关联。只有在PROSE主要分析中分类为VS1.0不良的患者的两个分类为未知;所有其它分类为早期。在分类为早期的148个患者中,73个具有VS1.0分类VS良好,以及75个具有VS1.0分类VS不良。

[0127] 由VS2.0分类进行的患者特性在表5中示出。

[0128] 表5. 由VS2.0分类在VS1.0良好人口中进行的患者特性

		后期(N=39)	早期/未知(N=140)	p值
组织学	腺	27 (69%)	93 (67%)	0.100
	鳞状	2(5%)	24 (17%)	
	BAC	2 (5%)	1(1%)	
	大	2 (5%)	8 (6%)	
	NOS	2 (5%)	4 (3%)	
	其它/缺失	4 (10%)	10 (7%)	
性别	男性	26 (67%)	94 (67%)	>0.99
	女性	13 (33%)	46 (33%)	
吸烟状况	从不	7 (18%)	23 (16%)	0.968
	以前	23 (59%)	82 (58%)	
	当前	9 (23%)	35 (25%)	
PS	0	24 (62%)	81 (58%)	0.491
	1	15 (38%)	52 (37%)	
	2	0 (0%)	7 (5%)	
EGFR突变	突变	5 (16%)	7 (7%)	0.159
	WT	24 (75%)	84 (86%)	

图8根据治疗示出分类编组后期和早期/未知(VS1.0良好)的OS结果,其中图8A示出厄洛替尼治疗编组的数据,而图8B示出化学疗法治疗编组的数据。图9根据治疗示出分类编组

后期和早期/未知(VS1.0良好)的PFS结果,其中图9A示出厄洛替尼治疗编组的数据,而图9B示出化学疗法治疗编组的数据。

[0129] VS1.0良好人口的多元分析的结果在表6中示出。后期或早期/未知的VS2.0结果在对可能混合因素来调整时保持有意义。

[0130] 表6. VS1.0良好人口的多元分析

端点	协变量	HR (95% CI)	p 值
OS	治疗: CT 与 ERL	1.12 (0.85-1.65)	0.320
	VS2.0: 早期/未知与后期	0.59 (0.39-0.89)	0.012
	性别: 男性与女性	0.83 (0.57-1.20)	0.316
	PS: 0-1 与 2	1.87 (0.86-4.08)	0.114
	吸烟状况: 从不与曾经	1.23 (0.75-2.00)	0.411
PFS	治疗: CT 与 ERL	1.43 (1.05-1.93)	0.023
	VS2.0: 早期/未知与后期	0.57 (0.39-0.83)	0.004
	性别: 男性与女性	1.06 (0.75-1.48)	0.759
	PS: 0-1 与 2	1.30 (0.60-2.81)	0.500
	吸烟状况: 从不与曾经	1.31 (0.85-2.02)	0.230

图10示出根据治疗的编组VS1.0不良和后期的OS的Kaplan-Meier图表连同分类、VS1.0不良和后期以及治疗之间的交互的分析的结果。

[0131] 图11比较VS1.0良好早期/未知编组中的化学疗法与厄洛替尼之间的结果。

[0132] 根据治疗的后期编组中的结果的比较在图12中示出。注意,图12A中,分类为后期并且接受厄洛替尼的那些患者具有17.1个月的中值总生存时间,比接受化学疗法的那些患者要多两个月。

[0133] 在表6中对各治疗臂连同其95%置信度间隔和各编组中的患者的数量来概括各编组的OS和PFS的中值。

[0134] 表6 根据编组和治疗臂的OS和PFS的中值

端点	编组	n	中值(月)	95% CI(月)
OS	后期CT	16	15.1	6.2-24.2
OS	VS1.0不良CT	40	6.4	3.3-7.4
OS	早期/未知(VS1.0良好)CT	69	10.9	7.4-14.1
OS	后期ERL	23	17.1	13.1-27.9
OS	VS1.0不良ERL	37	3.1	2.0-4.0
OS	早期/未知(VS1.0良好)ERL	71	9.6	6.3-11.0
PFS	后期CT	16	6.1	2.6-10.4
PFS	VS1.0不良CT	40	2.8	1.9-4.5
PFS	早期/未知(VS1.0良好)CT	69	4.7	2.5-5.4
PFS	后期ERL	23	3.9	2.4-7.8
PFS	VS1.0不良ERL	37	1.7	1.5-2.2
PFS	早期/未知(VS1.0良好)ERL	71	2.3	2.0-2.8

小节II结论

本小节所述的测试(VS2.0)是利用从基于血液的样本的质谱所得出的74个特征的真实

多元测试,以识别一组对厄洛替尼具有优于化学疗法的优良性能的第二线NSCLC患者。这个测试的开发验证了CMC/D分类器形成方法。VS2.0将我们先前在原始VeriStrat测试编组中识别为“良好”的编组分离为两个小组,即“VS2.0早期”或“早期”以及“VS2.0后期”或“后期”,虽然对于不可识别患者的充分编组,在这里因谱获取的限制而描述为“VS2.0未知”。

[0135] 在其当前实现中,这个测试(VS2.0)依靠在对我们的原始VeriStrat测试所证明的机器上的谱获取。由于VS2.0要求来自VS1.0验证方案之外的m/z范围的特征值,所以需要特别注意通过利用参考样本来校正m/z相关灵敏度的差。标签稳定性使用硅内灵敏度分析来评估,这产生VS2.0未知的充分数量。根据仅指配确切标签的所指配VS2.0标签的再现性通过形成集合的三次运行来评估,并且非常高。对于VS2.0的临床使用,我们分析三个编组:VS2.0后期、VS1.0良好人口中的VS2.0早期和未知以及几乎均匀地分类为VS2.0早期的VS1.0不良。

[0136] VS2.0在PROSE样本的单盲分析中证明(临床验证)。VS2.0后期编组中的样本的可用数量在某些方面限制这个证明的显著性。将VS2.0后期中的总生存与VS1.0良好编组中的VS2.0早期/未知进行比较表明,VS2.0将VS1.0良好编组分为厄洛替尼治疗下的良好和不良执行编组,而存在化学疗法臂中的这种划分的极少证据。然而,样本大小过小而无法取得厄洛替尼优于化学疗法的优良性的统计显著性。VS2.0保留VS1.0的预测能力(根据治疗的VS2.0后期与VS1.0不良),即使样本大小减半。关于PFS的结果与OS中相似。

[0137] VS2.0的成功开发验证测试开发以及一般的CMC/D方法的相关方式。训练标签和识别这类患者的测试的并行迭代开发惊人地顺利工作。CMC/D中固有的、避免过适的量度已经证明是有效的,并且扩展成包括对训练/测试划分MC的多数表决,进一步降低测试/最终分类器选择中的不明确。VS2.0利用我们使用的合计谱中的可观测峰值的大约60%(2000发射谱的3个复制),其中没有明确有利特征。因此,虽然本示例利用附录B中所述的特定特征,但是这些特定特征并不认为是本质或关键的,以及顺利执行测试可基于这些特征的子集或者例如通过从更大数量的发射所得到的谱所发现的可能的附加特征。

[0138] 在商业使用方面,VS2.0提供识别对其能够适当地确信厄洛替尼至少相当于化学疗法并且可能是优良的一组患者的工具。第二线设定中的17个月总生成的中值是壮观的,并且可能引起第2线NSCLC的治疗方案的变化。我们再次能够定义类标签“早期”和“后期”(或者等效体),其实现作为这个过程的一部分的这个预测。

[0139] 小节III 测试环境(图2)中的VS 2.0 CMC/D分类器的使用

在本小节将结合图2来描述如小节II所述的分类来自NSCLC患者的基于血液的样本的CMC/D分类器的应用。如上所述,如果指配给测试样本的类标签为“后期”或等效体,则类标签预测提供样本的NSCLC患者更可能获益于EGFR-I、例如厄洛替尼或吉非替尼,如与化学疗法相比。其测试样本具有与其关联的“中间”标签的患者被预测从化学疗法和EGFR-I得到相似的临床有意义有益效果。

[0140] 在该方法的一个可能实现中,来自样本的质谱首先经过美国专利7736906所述的VS 1.0测试,以及如果向样本指配不良标签,则报告那个测试标签。具有这个标签的患者被预测没有从患者的治疗中的EGFR-I得到有益效果。如果标签为VS良好或等效体,则样本谱经过图2所示的VS 2.0的测试过程,以便确定患者是否具有“后期”标签,其中患者被预测与化学疗法相比、从EGFR-I、厄洛替尼或吉非替尼得到更大有益效果,或者相反,具有“中间”

类标签,其中患者被预测从化学疗法和EGFR-I得到相似的临床有意义有益效果。预期第三类标签、即“未知”或“中间”,其中无法预测与化学疗法相比、患者是否可能从EGFR-I得到有益效果。

[0141] 示出按照图1所生成的CMC/D分类器对测试样本的质谱的使用的工作流程在图2中示出。该过程开始于向质谱仪提供三个基于血液的样本:来自对其执行测试的患者的测试样本1200以及分别示为参考样本1和参考样本2、项1202A和1202B的两个参考样本等分式样。这两个参考样本是来自健康人类患者的参考基于血液的样本的两个等分式样。在这个实施例中,使用参考样本1202A和1202B,以便校正对VS 1.0测试中使用的特定质谱仪的先前证明m/z范围之外的m/z范围的m/z灵敏度变化。有可能的是,通过适当证明的机器,参考样本1和2的使用会是不必要的。

[0142] 在步骤1204,对三个样本1200、1202A和1202B的质谱测定使用MALDI-ToF质谱仪来执行。各样本在仪表中三次经过2000发射“稀释和上样”MALDI-ToF质谱测定,其中具有谱获取过滤(参见先前论述)。三个样本的每个的所产生三个2000发射谱从质谱仪传递给实现图2的工作流程的通用计算机的机器可读存储器。

[0143] 然后调用软件模块平均工作流程1206以执行在步骤1204所得到的三元谱的平均,在步骤1208所示。平均工作流程在图13中示出。基本上,这个模块估计用于对齐的谱中的峰值,执行原始谱的对齐,并且然后从三个样本的每个的三个复制来计算对齐谱的平均值。

[0144] 然后调用预处理工作流程模块1212(图14),以执行平均谱的预处理,并且生成特征值(特征表)供分类中使用,如在步骤1214所示。该步骤包括对预定义m/Z范围的特征值(累积强度值)的背景减法和估计、峰值检测和对齐、部分离子电流归一化和计算。范围在附录B中列示。

[0145] 如在1216所示,将在步骤1214所生成的两个参考样本(1202A和1202B)的特征值提供给模块1218,其检查参考值是否一致。基本上,在模块1218中,执行参考特征值的比较。这涉及下列步骤:

1. 对于在步骤1214所得到的所有特征值F,计算参数 $\delta_F = \min(|1 - (FV_{pre}/FV_{post})|, |1 - (FV_{post}/FV_{pre})|)$ 。这里的思路是在测试样本1200之前(或者在测试样本的批次的开始)运行一个参考样本(1202A),并且从参考样本得到特征值的集合、即 FV_{Pre} ,然后在测试样本1202之后(或者在测试样本的批次的结束)运行参考样本1202B的另一个准备,并且再次从参考样本得到特征值的集合、即 FV_{Post} 。

[0146] 2. 选择那些特征,其中 δ_F is < 0.1 ,将那些特征值添加到特征值的列表(列表L)。

[0147] 3. 将在2所选的特征的列表L与从相同步骤1-2、从随用来生成CMC/D分类器的样本的形成集合运行的参考样本所得到的特征值的列表L'(即,附录B中的特征的列表)进行比较。

[0148] 4. 如果列表L包含在m/Z位置3219和18634的特征,则这些特征值被认为是一致的。

[0149] 如果一致性测试(4)失败,则该过程返回到开始,并且重做测试样本和两个参考样本的谱获取。如果一致性测试(4)成功,则该处理使用特征值1220的标准集合继续进行到定义特征校正函数步骤1222。这些是当生成原始谱时(即,在CMC/D分类器的生成时)随形成集合样本运行的参考样本(1201A和1202B)的两个准备的特征值。它能够是所有特征值的列表,但是一部分没有通过在两个准备之间建立的一致性标准,并且因此这些特征实际上从

未用,并且会从列表中排除。我们寻找在随形成集合谱运行并且对前和后参考谱也是一致的参考样本的两个准备之间是一致的特征。然后,计算原始样本的平均以及这些特征的前和后样本的平均。我们计算出这两者的比率,并且将它作为m/Z的函数来绘制。生成比率的图表的线性回归,并且返回Y轴截距和斜率。参见以上图5的论述。

[0150] 在步骤1224,来自步骤1222的Y轴截距和斜率分别是来自线性回归图表的特征校正函数参数a和b。这些值应用于在步骤1214所生成的测试样本特征值。这个校正能够表达如下:

在步骤1224,这些校正特征值存储在存储器中。校正特征值用于两个独立处理分支中:步骤1228和步骤1232。

[0151] 在步骤1228,表示最终按照图1的过程所生成的CMC/D分类器1226的数据集应用于校正测试样本特征值。在这个示例中,最终CMC/D分类器是在来自分类器生成样本集合1100(图1)的测试和训练样本划分实现的每个中生成并且在图1的步骤1134所创建的250个主分类器的集合。主分类器对校正特征值的这个应用的结果是测试样本分类标签,如在1229所示。

[0152] 如图2在1232所示,在步骤1224所生成的校正特征值也发送给模块1232,其利用预定义特征相关噪声特性1230来生成新特征值实现(“噪声实现”)。基本上,这个模块1232使用从形成样本集合(图1的1100)所得到的噪声参数 σ_i 来生成160个噪声实现:

加性噪声实现:

$$FVN_i = FV_{corrected,i} + \epsilon_i$$

- 乘性噪声实现:

$$FVN_i = FV_{corrected,i} * (1 + \epsilon_i)$$

其中 ϵ_i 是具有零平均和单位标准偏差的高斯随机数(N),其特征在于表达 $N(0, \sigma_i)$,其中 σ_i 是如先前所述从形成集合所确定的噪声参数。

[0153] 在步骤1232所生成的所产生“噪声”特征值采取特征表的形式。所有特征值作为工作流程人工产物来提供。这个过程的结果按照便利形式、例如Excel电子表格来存储。

[0154] 在步骤1234,表示主分类器(1226,以上所述)的数据集应用于在步骤1232所生成的有噪特征值。参见图15。这产生主分类器结果的表(每种类型的类标签的数量)。在这个具体示例中,在主分类器采取产生于250个训练/测试集合划分(如上所述)的250个主分类器的形式的情况下,存在对各噪声实现所生成的250个类标签。噪声实现的主分类器结果如在步骤1236所示来核对,使得对分类结果的统计数据能够如1238所示来得到。在这个步骤1236,我们生成比率R(称作“噪声效果估计器”),其与后期和早期分类的数量之间的差的标准偏差相关。这对特征表的所有有噪实现进行。比率R的这个统计分析和计算的细节如下:

设 N_{Early}^i = 跨对测试样本的各噪声实现i(在这个示例中 $1 \leq i \leq 160$,因为存在160个不同噪声实现)所计算的250个主分类器(MC)的早期分类的数量。对所有i计算总和, $\sum_i N_{Early}^i$ 。

[0155] 设 N_{Late}^i = 跨对测试样本的噪声实现i($1 \leq i \leq 160$)所计算的250个主分类器(MC)的后期分类的数量。对所有i计算总和, $\sum_i N_{Late}^i$ 。

[0156] 因此,对于所有 $i, 0 \leq N_{\text{Early}}^i \leq 250$ 和 $0 \leq N_{\text{Late}}^i \leq 250$ 。

[0157] 以及对于所有噪声实现 $i, N_{\text{Early}}^i + N_{\text{Late}}^i = 250$ 。

[0158] 噪声效果估计器 $R = N_{\text{Early}}^i / (|\sum_i N_{\text{Early}}^i - \sum_i N_{\text{Late}}^i| / 320)$ 的标准偏差
 $= \text{sqrt}(\sum_i (N_{\text{Early}}^i)^2 - (\sum_i N_{\text{Early}}^i)^2) / (|\sum_i N_{\text{Early}}^i - \sum_i N_{\text{Late}}^i| / 320)$
 $= \text{sqrt}(\sum_i (N_{\text{Early}}^i)^2 - (\sum_i N_{\text{Early}}^i)^2) / (|\sum_i N_{\text{Early}}^i - 20000| / 160)$

R 中的分母 $(|\sum_i N_{\text{Early}}^i - \sum_i N_{\text{Late}}^i| / 320)$ 给出我们跨160个噪声实现获得的早期与后期的数量之间的平均差的量度。如果这个数量较小,则多数表决分类接近,以及如果它比较大,则它是单方面表决。本质上,比率 R 将MC标签中的可变性与单方面的方式进行比较,这是重要的,因为我们希望知道在噪声参数 ϵ 中我们测量的可变性是否可能引起不可靠多数表决分类。也就是说,如果我们对所有250个MC对220个早期和30个后期求平均,则我们不在意比如10的可变性,但是如果我们对所有250个MC对130个早期和120后期求平均,则我们在意10的可变性。

[0159] 测试样本(图2的1200)的最终分类标签在步骤1240来生成。在所示实施例中,这个分类将仅对VS1.0分类为良好的样本来执行;即,初步测试使用VS 1.0进行,以及如果患者测试VS不良,则报告那个标签。报告的最终分类标签如下:

1. 如果在步骤1236所确定的比率 $R > 0.5$,则返回标签中间(或等效体)。其样本具有与其关联的中间标签的患者被预测从化学疗法和EGFR-I得到相似的临床有意义有益效果。注意,这与主分类器对校正特征值(1129)所产生的类标签无关。

[0160] 2. 如果在步骤1236所确定的比率 $R \leq 0.5$,则,

A. 如果在1229所生成的测试样本标记为后期,则返回后期标签。

[0161] B. 如果在1229所生成的测试样本标记为早期,则返回早期标签。

[0162] 其测试样本具有2.A中的后期标签的患者被预测与NSCLC癌症的治疗的化学疗法相比、从EGFR-I得到更大有益效果。

[0163] 在一个可能实施例中,中间标签被认为包括那些患者,其中噪声效果估计器 > 0.5 (上述1.)加上早期(≤ 0.5 噪声效果估计器和早期标签)。它们被结合,因为这是临床有用的(如果决定给出后期EGFR-I和作为VS1.0不良化学疗法的那些测试,则它们本质上由剩余的那些患者组成。结果可对化学疗法和TKI相似的结果对这个组合编组结束(噪声效果估计器 > 0.5 (上述1.)加上早期(≤ 0.5 噪声效果估计器和早期标签),而不是单独对任一编组。

[0164] 小节IV 用于生成CMC/D的有形系统的实际示例 分类器和进行预测测试

分类器生成系统和样本测试系统

小节I和II所述的CMC/D分类器开发方法能够实现为采取质谱仪(或者其它测量仪表)形式的有形分类器开发系统,其用来从多个样本(例如样本的形成集合)来得到质谱(或其它)数据,以及实现为具有运行实现CMC/D分类方法的代码的处理单元的通用计算机。具体来说,计算机包括存储测量数据的机器可读存储器(例如硬盘)。计算机还存储可执行代码,其执行测量数据的预处理,例如背景减法、谱对齐和归一化,如上所述,以及存储在用于分类的特定特征的累积强度值,例如附录B中所列特征的累积强度值。

[0165] 计算机还存储可执行代码,其用于使用来自一直到预先选择特征集合大小(s , 整数)的样本的特征的集合来构成多个单独微型分类器。在一个实施例中,该代码包括KNN分类算法(本领域已知),其应用于质谱测定数据中的特定或者多个特征,并且将特征值与样

本的形成集合的子集(例如类标记质谱数据的训练集合)进行比较。KNN算法基于特征空间中的最近邻来生成类标签。

[0166] 该代码然后测试单独微型分类器的每个的分类精度或者某个备选性能量度,以分类样本的给定集合(例如训练集合)中的生物样本,并且保留其性能超过预定义阈值或者处于预定义极限之内以达到微型分类器的过滤集合的那些分类器。

[0167] 代码然后通过随机选择过滤微型分类器的一小部分并且对这类所选微型分类器进行逻辑训练,使用极端退出、对样本的分类标签(使用等式1)重复进行微型分类器的过滤集合的逻辑训练。

[0168] 该代码然后继续生成最终分类器,例如作为对退出迭代的逻辑回归训练的平均。在一个示例中,最终分类器在计算机存储器中表示为使用分类的单一特征($s=1$)的微型分类器和使用分类的两个特征($s=2$)(其通过过滤标准)的加权组合。

[0169] 能够针对测试集合划分或者形成集合的子集来评估最终分类器,该评估还对形成集合到训练和测试集合的多个不同划分来执行,以及能够通过选择产生于特定训练和测试集合划分其中之一,或者备选地通过保留来自各训练和测试集合划分的所有主分类器并且使用来自每个主分类器的多数表决将标签指配给被测样本,来生成最终分类器。

[0170] 这个最终分类器然后用于测试样本、例如NSCLC癌症患者的基于血液的样本的分类,以在治疗之前预测NSCLC患者是否可能获益于EGFR-I。如果指配给样本的质谱的类标签为后期,则那表示患者可能获益。

[0171] 上述分类系统能够在商业上测试样本并且为诊所、医院、肿瘤学家和其它健康护理提供者提供服务的实验室测试中心来实现,其中具有关于从癌症靶向药物的患者获益的测试结果。当然,分类器开发方法能够用于其它目的、例如诊断目的。

[0172] 测试系统

图16是用于使用按照图1所生成的分类器来处理测试样本的有形系统的另一个示例,其中包括质谱仪2606以及通用计算机2610,其实现编码为机器可读指令的CMC/D分类器2620以及形成存储器2614中存储的类标记质谱测定数据2622的训练集合的特征表2622。将会理解,图16的测量仪表2606和计算机2610可用来按照图1生成CMC/D分类器。

[0173] 在小节III的具体实施例中,质谱仪和计算机2610实现图2所示并且以上详细描述的工作流程。

[0174] 现在将描述一备选实施例。图16的系统得到多个样本2600、例如来自癌症患者的基于血液的样本(血清或血浆)。2600用来进行关于患者是否可能获益于特定药物或者药物的组合的预测。样本可作为血清卡等得到,其中基于血液的样本涂抹到纤维素或其它类型的卡上。得到样本的三个等分式样。在一个可能实施例(如小节III所述)中,也可使用参考样本2604。

[0175] 样本的三个等分式样放置到MALDI-ToF样本“板”2602上,以及板插入测量仪表,在这种情况下为MALDI-ToF质谱仪2606。质谱仪2606从样本的三个等分式样的每个来获取质谱2608。质谱采取数字形式来表示,并且提供给编程通用计算机2610。计算机2610包括运行编程指令的中央处理器2612。存储器2614存储表示质谱2608的数据。

[0176] 存储器2614还存储主或最终CMC/D分类器2620,其包括:a) 采取N个类标记谱的特征表的形式训练集合2622,其中N是某个整数,在这个示例中,来自在临床试验中如先前

所述登记的患者的类标记谱,并且各样本被指配类标签、例如“早期”、“后期”、“+”、“-”、“良好”、“不良”等;b)表示KNN分类算法的代码;c)用于对患者的质谱运行按照图1所生成的最终分类器的程序代码;以及d)用于存储分类结果和测试样本的最终类标签的数据结构2628。存储器2614还存储用于实现在2650所示的处理的程序代码2630,包括:用于在步骤2652从质谱仪来获取质谱数据的代码(未示出);用于实现背景减法、归一化和对齐步骤2654的预处理例程2632;用于在背景所减去、归一化和对齐谱中的预定义m/Z位置得到累积强度值(步骤2654)的模块(未示出);以及用于对在步骤2656所得到的值使用训练集合2622来实现分类器2620的代码例程2638。过程2658在步骤2660产生类标签。程序代码2642包括进行检查(步骤2662)以确定样本的所有三个等分式样是否产生相同类标签的代码。如果不是,则报告类标签“未定义”或等效体。如果对患者样本2600的三个等分式样产生相同类标签,则模块2640如在2666所示报告类标签(即,“早期”、“后期”、“+”、“-”、“良好”、“不良”或等效体)。

[0177] 程序代码2630能够包括附加和可选模块,例如特征校正函数代码2632(图2所述)、用于处理来自参考样本2604的谱以定义特征校正函数的例程的集合、存储特征相关噪声特征和所生成有噪特征值实现(参见图2)并且分类这类有噪特征值实现的模块以及存储用于得到关于分类器对有噪特征值实现的性能的统计数据的统计算法的模块。可包含又一些可选软件模块,如本领域的技术人员将会清楚地知道。

[0178] 图16的系统能够实现为实验室测试处理中心,其从肿瘤学家、患者、诊所等得到多个患者样本,并且作为收费服务生成患者样本的类标签。质谱仪2606无需物理上位于实验室测试中心,计算机2610而是可通过计算机网络来得到表示测试样本的质谱的数据。

[0179] NSCLC患者的治疗方法

还将会理解,我们描述了治疗NSCLC患者的方法。治疗采取向NSCLC患者注射EGFR-I的形式,其中通过在编程计算机中运行分类器(其将质谱仪从NSCLC患者的基于血液的样本所产生的质谱数据与包括从通过基于血液的样本的质谱测定来确定为被预测从癌症的治疗中的EGFR-I得到总生存有益效果的一类患者的成员的多个癌症患者所得到的类标记质谱数据的训练集合进行比较),预测患者与化学疗法相比更多地获益于EGFR-I。这类患者还分为两个子类:

1.在癌症的治疗中注射所述EGFR-I之后呈现疾病的早期进展的那些患者,这类患者的质谱数据具有“早期”或等效体的类标签;以及

2.在癌症的治疗中注射EGFR-I之后呈现疾病的后期进展的那些患者,这类患者的质谱数据具有类标签“后期”或等效体。此外,编程计算机能够采取实现如本文档的先前小节详细描述的分类算法的分类器的形式。例如,编程计算机实现采取退出规则化和逻辑训练之后的过滤微型分类器的组合(CMC/D分类器)的形式分类器。EGFR-I可采取具有按照所建立协议的剂量的吉非替尼、厄洛替尼、第二代EGFR-I、例如达克替尼、阿法替尼或等效物的形式。

[0180] 所附权利要求书作为所公开发明的进一步描述来提供。

[0181] 附录

附录A:分类器形成中使用的样本

样本ID

ICA_1
ICA_10
ICA_11
ICA_12
ICA_13
ICA_14
ICA_15
ICA_17
ICA_18
ICA_19
ICA_2
ICA_20
ICA_21
ICA_22
ICA_23
ICA_24
ICA_25
ICA_26
ICA_27
ICA_28
ICA_29
ICA_3
ICA_30
ICA_31
ICA_32
ICA_34
ICA_35
ICA_36
ICA_38
ICA_39
ICA_4
ICA_40
ICA_41
ICA_42
ICA_43
ICA_44
ICA_45
ICA_46
ICA_47

ICA_48
ICA_49
ICA_5
ICA_50
ICA_51
ICA_52
ICA_54
ICA_55
ICA_56
ICA_57
ICA_58
ICA_59
ICA_6
ICA_60
ICA_61
ICA_63
ICA_64
ICA_65
ICA_67
ICA_68
ICA_69
ICA_7
ICA_70
ICA_8
ICB_1
ICB_10
ICB_11
ICB_12
ICB_13
ICB_14
ICB_15
ICB_16
ICB_17
ICB_18
ICB_19
ICB_2
ICB_20
ICB_21
ICB_22

ICB_23
ICB_24
ICB_25
ICB_26
ICB_27
ICB_28
ICB_29
ICB_3
ICB_30
ICB_31
ICB_32
ICB_33
ICB_34
ICB_35
ICB_36
ICB_37
ICB_38
ICB_39
ICB_4
ICB_40
ICB_41
ICB_42
ICB_43
ICB_44
ICB_45
ICB_46
ICB_47
ICB_48
ICB_49
ICB_5
ICB_50
ICB_51
ICB_52
ICB_53
ICB_54
ICB_55
ICB_56
ICB_57
ICB_58

ICB_59
ICB_6
ICB_60
ICB_61
ICB_62
ICB_63
ICB_64
ICB_65
ICB_66
ICB_67
ICB_8
ICB_9
ICC_1
ICC_10
ICC_11
ICC_12
ICC_13
ICC_14
ICC_15
ICC_16
ICC_17
ICC_18
ICC_19
ICC_2
ICC_20
ICC_21
ICC_22
ICC_23
ICC_24
ICC_25
ICC_26
ICC_27
ICC_28
ICC_29
ICC_3
ICC_30
ICC_31
ICC_32
ICC_4

ICC_5
ICC_6
ICC_7
ICC_8
ICC_9

附录B:CMC/D分类器中使用的特征

中心	左	右
3218.7386	3206.9871	3230.49
3315.4528	3302.6206	3328.285
4409.1599	4400.38	4417.94
4466.5671	4453.3297	4479.805
4715.9166	4700.9233	4730.91
4790.6135	4764.6789	4816.548
4862.7438	4846.8049	4878.683
5740.33	5689.9468	5790.713
5851.6323	5796.3864	5906.878
5945.9151	5914.4425	5977.388
6291.0333	6276.175	6305.892
6436.5097	6410.7103	6462.309
6531.4679	6517.0148	6545.921
6647.2276	6606.9751	6687.48
6835.523	6823.2312	6847.815
6859.0262	6849.9761	6868.076
6887.3988	6871.2103	6903.587
6942.638	6907.3833	6977.893
7044.8902	7019.7662	7070.014
7195.2294	7176.9942	7213.465
7388.9278	7374.8799	7402.976
7567.903	7548.4521	7587.354
7663.6716	7641.9244	7685.419
7765.1134	7750.9304	7779.296
7940.7116	7914.2368	7967.187
8019.8659	7975.8313	8063.901
8222.2092	8194.6538	8249.765
8582.8611	8556.6564	8609.066
8633.3793	8615.0091	8651.75
8696.8649	8673.0916	8720.638
8771.1565	8751.5705	8790.742
8819.6486	8800.1977	8839.1

8874.8945	8858.5504	8891.239
8934.0576	8900.4238	8967.692
9023.3426	9004.2969	9042.388
9147.2069	9108.5753	9185.839
9296.8707	9269.4504	9324.291
9359.8159	9331.8553	9387.777
9440.8613	9401.8245	9479.898
9584.3116	9553.2442	9615.379
9654.0106	9619.7014	9688.32
9731.9492	9696.4243	9767.474
9939.5604	9899.9833	9979.138
10641.5484	10617.64	10665.46
10828.7631	10808.2317	10849.29
11395.5404	11375.4141	11415.67
11440.1153	11427.013	11453.22
11512.9211	11464.564	11561.28
11699.0553	11597.2083	11800.9
11884.9193	11831.2943	11938.54
12112.5217	12062.4086	12162.63
12449.5353	12424.2762	12474.79
12577.8361	12557.5686	12598.1
12615.0568	12600.6529	12629.46
12727.1157	12712.9328	12741.3
12864.8928	12838.1478	12891.64
13125.0484	13107.6237	13142.47
13312.3983	13293.3526	13331.44
13577.2816	13556.615	13597.95
13749.638	13693.4466	13805.83
13883.9032	13816.0952	13951.71
13982.3733	13959.5455	14005.2
14048.2902	14021.0049	14075.58
14096.9174	14079.0874	14114.75
14156.3507	14130.146	14182.56
14484.7195	14462.432	14507.01
14777.5634	14759.4632	14795.66
17268.0853	17235.6355	17300.54
17401.8418	17364.907	17438.78
17607.8848	17577.5456	17638.22
18634.4067	18591.1403	18677.67

21071.3078	21030.6796	21111.94
22316.6349	22129.9002	22503.37
23220.6291	22951.4507	23489.81

附录C:分类器形成的第一阶段的初始类标签

样本ID	类标签
36HSR	早期
38HSR	早期
39HSR	早期
40HSR	早期
45HSR	早期
51HSR	早期
56HSR	早期
63HSR	早期
68HSR	早期
ICB_03	早期
ICB_06	早期
ICB_10	早期
ICB_12	早期
ICB_13	早期
ICB_22	早期
ICB_26	早期
ICB_34	早期
ICB_38	早期
ICB_40	早期
ICB_43	早期
ICB_45	早期
ICB_60	早期
ICB_63	早期
10HSR	后期
11HSR	后期
12HSR	后期
13HSR	后期
14HSR	后期
17HSR	后期
18HSR	后期
19HSR	后期
1HSR	后期
20HSR	后期
21HSR	后期

22HSR	后期
2HSR	后期
4HSR	后期
7HSR	后期
8HSR	后期
ICB_05	后期
ICB_28	后期
ICB_31	后期
ICB_41	后期
ICB_57	后期
ICB_61	后期
ICB_64	后期

附录D:VS2/0特征的噪声类型和噪声强度

特征的m/Z中心	噪声类型	噪声强度
3218.7386	加性	0.449589
3315.4528	加性	0.705299
4409.1599	加性	0.372679
4466.5671	加性	0.558918
4715.9166	乘性	0.215793
4790.6135	加性	0.871467
4862.7438	乘性	0.224417
5740.33	乘性	0.219152
5851.6323	乘性	0.250464
5945.9151	乘性	0.671156
6291.0333	加性	0.204162
6436.5097	加性	1.674129
6531.4679	加性	0.19534
6647.2276	加性	3.511696
6835.523	加性	0.369546
6859.0262	加性	0.216011
6887.3988	加性	0.449448
6942.638	加性	1.17939
7044.8902	加性	0.435487
7195.2294	加性	0.222608
7388.9278	加性	0.163982
7567.903	乘性	0.156163
7663.6716	乘性	0.195681
7765.1134	加性	0.319943
7940.7116	加性	0.419978

8019.8659	加性	0.356489
8222.2092	加性	0.431253
8582.8611	加性	0.347085
8633.3793	加性	0.268113
8696.8649	乘性	0.274013
8771.1565	加性	0.692564
8819.6486	乘性	0.38203
8874.8945	加性	0.514021
8934.0576	乘性	0.29018
9023.3426	加性	0.416469
9147.2069	乘性	0.233822
9296.8707	乘性	2.007367
9359.8159	乘性	0.15884
9440.8613	乘性	0.155807
9584.3116	乘性	0.280165
9654.0106	乘性	0.200748
9731.9492	乘性	0.200652
9939.5604	乘性	0.240092
10641.5484	加性	0.246795
10828.7631	加性	0.374312
11395.5404	加性	0.511211
11440.1153	乘性	0.240577
11512.9211	乘性	0.316491
11699.0553	乘性	0.402835
11884.9193	乘性	0.190473
12112.5217	乘性	1.367853
12449.5353	乘性	2.019671
12577.8361	乘性	0.163202
12615.0568	乘性	0.50929
12727.1157	乘性	0.212812
12864.8928	乘性	0.116047
13125.0484	加性	0.143445
13312.3983	加性	0.144914
13577.2816	加性	0.136992
13749.638	加性	1.208693
13883.9032	加性	2.503822
13982.3733	加性	0.517253
14048.2902	加性	1.393395
14096.9174	加性	0.595363

14156.3507	加性	0.837603
14484.7195	加性	0.22863
14777.5634	加性	0.091024
17268.0853	加性	0.353217
17401.8418	加性	0.574893
17607.8848	加性	0.142937
18634.4067	加性	0.133441
21071.3078	加性	0.133543
22316.6349	加性	1.392056
23220.6291	加性	0.776561

附录E:形成集合样本的VS2.0分类

样本ID	总分类	VS1.0分类
ICA_1	后期	良好
ICA_10	后期	良好
ICA_11	早期	良好
ICA_12	早期	良好
ICA_13	后期	良好
ICA_14	后期	良好
ICA_15	后期	良好
ICA_17	后期	良好
ICA_18	早期	良好
ICA_19	后期	良好
ICA_2	后期	良好
ICA_20	后期	良好
ICA_21	后期	良好
ICA_22	早期	良好
ICA_23	早期	良好
ICA_24	早期	不良
ICA_25	早期	良好
ICA_26	早期	良好
ICA_27	后期	良好
ICA_28	早期	良好
ICA_29	早期	良好
ICA_3	早期	不良
ICA_30	早期	不良
ICA_31	早期	良好
ICA_32	早期	良好
ICA_34	后期	良好
ICA_35	早期	良好

ICA_36	后期	良好
ICA_38	早期	良好
ICA_39	早期	良好
ICA_4	后期	良好
ICA_40	早期	良好
ICA_41	后期	良好
ICA_42	早期	良好
ICA_43	早期	不良
ICA_44	后期	良好
ICA_45	早期	良好
ICA_46	早期	良好
ICA_47	早期	不良
ICA_48	后期	良好
ICA_49	早期	不良
ICA_5	后期	良好
ICA_50	后期	良好
ICA_51	后期	良好
ICA_52	早期	不良
ICA_54	早期	不良
ICA_55	后期	良好
ICA_56	早期	良好
ICA_57	早期	不良
ICA_58	早期	不良
ICA_59	早期	不良
ICA_6	早期	不良
ICA_60	早期	不良
ICA_61	早期	不良
ICA_63	早期	良好
ICA_64	早期	不良
ICA_65	早期	不良
ICA_67	早期	良好
ICA_68	后期	良好
ICA_69	早期	不良
ICA_7	后期	良好
ICA_70	早期	良好
ICA_8	后期	良好
ICB_1	早期	不良
ICB_10	早期	良好
ICB-11	早期	不良

ICB_12	后期	良好
ICB_13	早期	良好
ICB_14	早期	良好
ICB_15	早期	良好
ICB_16	后期	良好
ICB_17	后期	良好
ICB_18	早期	不良
ICB_19	早期	不良
ICB_2	后期	良好
ICB_20	早期	不良
ICB_21	早期	良好
ICB_22	后期	良好
ICB_23	早期	不良
ICB_24	早期	不良
ICB_25	早期	不良
ICB_26	早期	良好
ICB_27	早期	不良
ICB_28	后期	良好
ICB_29	早期	不良
ICB_3	后期	良好
ICB_30	早期	不良
ICB_31	后期	良好
ICB_32	早期	不良
ICB_33	早期	不良
ICB_34	早期	良好
ICB_35	早期	不良
ICB_36	后期	良好
ICB_37	早期	不良
ICB_38	后期	良好
ICB_39	早期	良好
ICB_4	早期	不良
ICB_40	后期	良好
ICB_41	后期	良好
ICB_42	早期	不良
ICB_43	早期	良好
ICB_44	早期	不良
ICB_45	早期	良好
ICB_46	早期	不良
ICB_47	后期	良好

ICB_48	早期	良好
ICB_49	后期	良好
ICB_5	后期	良好
ICB_50	后期	良好
ICB_51	早期	不良
ICB_52	后期	良好
ICB_53	早期	不良
ICB_54	早期	良好
ICB_55	早期	不良
ICB_56	早期	不良
ICB_57	后期	良好
ICB_58	早期	不良
ICB_59	早期	不良
ICB_6	早期	良好
ICB_60	早期	良好
ICB_61	早期	良好
ICB_62	早期	良好
ICB_63	早期	良好
ICB_64	后期	良好
ICB_65	早期	良好
ICB_66	早期	不良
ICB_67	后期	良好
ICB_8	早期	不良
ICB_9	后期	良好
ICC_1	早期	不良
ICC_10	早期	良好
ICC_11	后期	良好
ICC_12	早期	不良
ICC_13	早期	不良
ICC_14	早期	良好
ICC_15	早期	不良
ICC_16	早期	不良
ICC_17	后期	良好
ICC_18	早期	不良
ICC_19	早期	良好
ICC_2	早期	不良
ICC_20	早期	不良
ICC_21	后期	良好
ICC_22	早期	良好

ICC_23	后期	良好
ICC_24	后期	良好
ICC_25	早期	良好
ICC_26	早期	良好
ICC_27	后期	良好
ICC_28	后期	良好
ICC_29	后期	良好
ICC_3	早期	不良
ICC_30	早期	良好
ICC_31	早期	良好
ICC_32	早期	不良
ICC_4	早期	良好
ICC_5	早期	良好
ICC_6	早期	不良
ICC_7	早期	良好
ICC_8	早期	不良
ICC_9	早期	良好

附录F:跨三次运行的形成集合样本的VS2.0分类

样本ID	形成运行分类	形成运行噪声效应估计器	Feb 3分类	Feb 3噪声效应估计器	Feb 25分类	Feb 25噪声计量器
ICA 1	后期	0.2508903	后期	0.466734822	未知	1.25354
ICA 10	后期	0.3138037	未知	1.964538835	未知	3.23176
ICA 11	早期	0.080601	早期	0.31109509	早期	0.18127
ICA 12	早期	0.0355124	早期	0.00909397	早期	0.1501
ICA 13	后期	0.0047174	后期	0.030926878	后期	0.08849
ICA 14	未知	2.7555361	未知	6.009376135	未知	0.57061
ICA 15	后期	0.0149085	后期	0.187318654	后期	0.08098
ICA 17	后期	0.0451973	后期	0.130183945	后期	0.10486
ICA 18	早期	0.3983651	早期	0.134071541	早期	0.2023
ICA 19	后期	0.0826776	后期	0.027922277	后期	0.03699
ICA 2	后期	0.0115269	后期	0.014803894	后期	0.01478
ICA 20	后期	0.2883118	后期	0.468349356	未知	1.55056
ICA 21	后期	0.3249368	后期	0.197541409	后期	0.42881
ICA 22	早期	0.4547106	未知	408.6471898	未知	10.2749
ICA 23	早期	0.0748141	未知	1.064878786		
ICA 24	未知	0.5213397	早期	0.273862348		
ICA 25	未知	0.5367448	未知	0.576202188	未知	2.14736
ICA 26	未知	1.4825573	未知	1.176456598	未知	1.14433
ICA 27	后期	0.4851147	未知	0.823851604	未知	0.54047
ICA 28	早期	0.024537	早期	0.041470212	早期	0.04415
ICA 29	早期	0.0684268	早期	0.199645029	早期	0.23878
ICA 3	早期	0.0449748	早期	0		
ICA 30	早期	0.1134967	早期	0		
ICA 31	未知	1.1973862	未知	2.017268589	未知	7.40837
ICA 32	未知	0.9744799	未知	3.705512439	未知	1.88644
ICA 34	后期	0.0513075	后期	0.075731492	后期	0.15651
ICA 35	早期	0.2933299	早期	0.191894212	早期	0.0942
ICA 36	后期	0.0405301	后期	0.207008265		

ICA 38	未知	0.6299707	早期	0.286152473	未知	1.39855
ICA 39	未知	0.6493858	未知	2.07717748	未知	1.02573
ICA 4	后期	0	后期	0.038223058	后期	0.06442
ICA 40	早期	0.1460363	未知	2.460497465	早期	0.11424
ICA 41	后期	0.359934	后期	0.401264716	未知	0.757
ICA 42	未知	2.2944611	早期	0.123948659	早期	0.27961
ICA 43	早期	0.0967663	早期	0.000632487		
ICA 44	未知	1.6734598	早期	0.169833656	早期	0.40807
ICA 45	未知	1.0538265	未知	0.584840142	早期	0.21289
ICA 46	早期	0.4287061	未知	2.926113519	未知	0.6906
ICA 47	早期	0.0535227	早期	0		
ICA 48	后期	0.4357615	未知	2.0349327	未知	2.07714
ICA 49	早期	0	早期	0		
ICA 5	未知	0.9192309	未知	0.653490123	后期	0.21708
ICA 50	未知	2.6894001	早期	0.158682214	未知	0.51338
ICA 51	后期	0.1653643	后期	0.31185332	未知	0.9165
ICA 52	早期	0.0045497	早期	0		
ICA 54	早期	0.0918534	早期	0		
ICA 55	后期	0.009786	未知	0.556007152	未知	1.96082
ICA 56	早期	0.0022435	早期	0.050034194	早期	0.0091
ICA 57	早期	0.0050177	早期	0.000632487		
ICA 58	早期	0	早期	0		
ICA 59	早期	0.0020317	早期	0.001887201		
ICA 6	早期	0.0010887	早期	0		
ICA 60	早期	0	早期	0		
ICA 61	早期	0	早期	0		
ICA 63	早期	0.0304895	早期	0.046816893	早期	0.14536
ICA 64	早期	0	早期	0		
ICA 65	早期	0	早期	0		
ICA 67	未知	0.7938756	未知	0.826523764	未知	0.60441
ICA 68	后期	0.2370179	未知	2.282512088	未知	2.00963
ICA 69	早期	0.0061302	早期	0.014126042		
ICA 7	后期	0.2874263	后期	0.092535875	后期	0.17229
ICA 70	未知	0.8459228	未知	0.592744714	早期	0.19042
ICA 8	后期	0.3185725	未知	0.524389074	未知	1.06012
ICB 1	早期	0.001642	早期	0		
ICB 10	早期	0.1244703	早期	0.071776831	早期	0.04976
ICB 11	早期	0	早期	0		
ICB 12	后期	0.4010251	未知	3.819985778	未知	2.46467
ICB 13	早期	0.0335419	早期	0.239284331	早期	0.20115
ICB 14	未知	0.7794731	未知	1.064463653	早期	0.20933
ICB 15	未知	1.402295	早期	0.005996916	早期	0.05784
ICB 16	后期	0.49193	未知	3.18288305		
ICB 17	未知	15.495518	未知	2.770598757	未知	0.75083
ICB 18	早期	0.0104891	早期	0		
ICB 19	早期	0.0044287	早期	0		
ICB 2	未知	1.8319861	未知	0.574145865	未知	1.11314
ICB 20	早期	0.1010281	早期	0.001265038		
ICB 21	早期	0.3837118	早期	0.047678494	早期	0.42108
ICB 22	后期	0.24719	未知	1.296687602	未知	2.0375
ICB 23	早期	0.0080037	早期	0		
ICB 24	早期	0	早期	0		
ICB 25	早期	0.4691525	早期	0.374906318		

ICB 26	早期	0.2842823	未知	18.84274386	未知	1.65263
ICB 27	早期	0.1090687	早期	0.026120232		
ICB 28	后期	0.0106621	后期	0.174473568	后期	0.11698
ICB 29	早期	0.0235619	早期	0.009862237		
ICB 3 rerun	后期	0.0304724	后期	0.067773006		
ICB 30	早期	0.0210381	早期	0.007672574		
ICB 31	后期	0.1671391	未知	1.269484668	未知	2.60353
ICB 32	早期	0.0504194	早期	0.006513994		
ICB 33	早期	0.0022743	早期	0		
ICB 34	未知	0.7717411	早期	0.235015835	早期	0.23868
ICB 35	早期	0.1187116	未知	0.684071314		
ICB 36	未知	0.6113689	早期	0.495122448		
ICB 37	早期	0	早期	0.000632487		
ICB 38	未知	0.7252647	后期	0.327507909	未知	7.41886
ICB 39	早期	0.0873692	未知	0.538723703	未知	0.69525
ICB 4	早期	0.0583902	早期	0		
ICB 40	未知	1.5221366	未知	1.376172237	未知	4.11934
ICB 41	后期	0.2281209	未知	2.942393151		
ICB 42	早期	0.016582	早期	0.001265038		
ICB 43	早期	0.008667	早期	0.014663441	早期	0.00617
ICB 44	早期	0.026458	早期	0.001253172		
ICB 45	早期	0.3637465	早期	0.19639466	早期	0.17223
ICB 46	早期	0	早期	0		
ICB 47	后期	0.3112708	后期	0.37180672	未知	0.53511
ICB 48	未知	0.6104345	未知	0.695956133	未知	1.19754
ICB 49	未知	0.8091827	未知	1.921287211		
ICB 5	未知	0.5610236	未知	1.791500069	未知	19.3159
ICB 50	未知	1.5210721	早期	0.322646083		
ICB 51	早期	0.2798399	早期	0.411311501		
ICB 52	后期	0.0913128	未知	0.995984435	后期	0.0946
ICB 53	早期	0.0177726	早期	0		
ICB 54	未知	3.9796933	未知	0.729611954		
ICB 55	早期	0.2673627	早期	0.016808751		
ICB 56	早期	0.016083	早期	0.001660149		
ICB 57	后期	0.0495004	后期	0.454621578	未知	5.38489
ICB 58	早期	0	早期	0		
ICB 59	早期	0.099419	早期	0		
ICB 6	早期	0.0926929	早期	0.010137147	早期	0.01514
ICB 60	早期	0.024118	早期	0.045176626	早期	0.22779
ICB 61	早期	0.0207761	早期	0.098978496	早期	0.05717
ICB 62	早期	0.1123475	早期	0.038795663		
ICB 63	早期	0.3143604	未知	0.5577347	早期	0.17666
ICB 64	后期	0.2135021	未知	0.981560369		
ICB 65	早期	0.4912493	未知	0.975042177	早期	0.48021
ICB 66	早期	0.0471047	早期	0.046567508		
ICB 67	未知	0.5234719	早期	0.322026183		
ICB 8	早期	0.0052102	早期	0		
ICB 9	后期	0.1080207	后期	0.042361028	后期	0.04029
ICC 1	早期	0.2070783	早期	0.085396794		
ICC 10	早期	0.1236901	早期	0.004740175	早期	0.01399
ICC 11	未知	1.1814412	未知	2.209011682	未知	1.34544
ICC 12	早期	0.0054516	早期	0		
ICC 13	早期	0	早期	0		
ICC 14	未知	0.9532531	早期	0.208090801	早期	0.40234

ICC 15	早期	0.0046228	早期	0.000632487		
ICC 16	早期	0.0006325	早期	0		
ICC 17	未知	1.060111	未知	0.503778812	后期	0.33919
ICC 18	早期	0.001265	早期	0.010079649		
ICC 19	早期	0.0946116	早期	0.034253636	早期	0.21303
ICC 2	早期	0	早期	0		
ICC 20	早期	0.0392832	早期	0.101833857		
ICC 21	后期	0.1985239	后期	0.269895491	未知	1.26594
ICC 22	早期	0.1766128	未知	1.01724785	未知	2.29042
ICC 23	未知	2.3518283	未知	4.747822355	未知	36.0979
ICC 24	后期	0.4498147	未知	1.641647487	后期	0.23851
ICC 25	早期	0.2547183	早期	0.026712614	早期	0.20825
ICC 26	早期	0.0183961	早期	0.177587583	早期	0.06516
ICC 27	未知	2.6560691	未知	0.894522603	未知	4.03214
ICC 28	未知	5.162227	未知	1.585391499	未知	1.17993
ICC 29	后期	0.0907799	后期	0.134559673	后期	0.30603
ICC 3	早期	0.0006325	早期	0		
ICC 30	早期	0.0374486	早期	0.025356686	早期	0.03447
ICC 31	早期	0.2820449	早期	0.145453279	早期	0.23148
ICC 32	早期	0.0045497	早期	0		
ICC 4	未知	2.6580968	未知	0.635164246	未知	5.92408
ICC 5	早期	0.1713111	未知	0.519211365	未知	0.51357
ICC 6	早期	0.0193609	早期	0		
ICC 7	早期	0.0008917	早期	0	早期	0.0272
ICC 8	早期	0.0873546	早期	0		
ICC 9	早期	0.0085559	早期	0.002577784	早期	0.00956

附录G:对PROSE样本所返回的VS2.0分类

单盲ID	VS2.0分类	PROSE样本号
3001	未知	01_024_1
3009	早期	11_046_1
3023	未知	01_055_1
3038	未知	16_005_1
3053	早期	04_001_1
3058	未知	10_002_1
3065	早期	16_013_1
3098	早期	11_055_1可能重复
3099	未知	06_014_1
3116	早期	01_059_1
3170	未知	01_013_1
3194	后期	10_005_1
3200	后期	01_074_1
3204	早期	01_010_1
3214	样本对MS生成不可用	11_043_1
3246	早期	16_012_1
3262	早期	01_039_1

3306	早期	01_044_1
3336	后期	16_017_1
3344	后期	06_012_1
3382	早期	01_075_1
3402	早期	06_043_1
3410	早期	06_002_1
3412	未知	11_050_1
3413	早期	01_008_1
3421	早期	06_010_1
3423	早期	01_066_1
3435	未知	11_044_1
3437	早期	11_003_1
3438	未知	08_001_1
3444	早期	11_047_1
3470	后期	01_021_1
3481	未知	01_025_1
3508	早期	01_001_1
3521	早期	16_006_1
3526	早期	01_034_1
3535	早期	01_062_1
3553	未知	01_082_1
3563	早期	06_040_1
3592	早期	11_005_1
3600	未知	14_001_1
3609	早期	14_012_1
3646	早期	11_030_1
3655	早期	07_012_1
3670	未知	06_030_1
3678	早期	01_052_1
3686	未知	01_080_1
3698	早期	01_029_1
3701	早期	01_060_1
3704	未知	01_049_1
3727	早期	12_007_1
3739	早期	11_008_1
3763	未知	01_061_1
3764	早期	06_020_1
3767	未知	12_013_1
3780	早期	12_009_1

3792	早期	12_003_1
3798	未知	01_089_1
3801	早期	07_011_1
3806	未知	04_013_1
3821	早期	16_016_1
3850	早期	11_056_1
3854	早期	14_013_1
3874	早期	01_093_1
3882	未知	12_006_1
3903	早期	07_007_1
3920	早期	11_026_1
3943	早期	11_012_1
3945	早期	11_033_1
3953	早期	11_042_1
3955	未知	04_005_1
3962	未知	12_013_1第二样本
3969	未知	14_006_1
3973	早期	13_005_1
3978	未知	03_001_1
3993	未知	02_005_1
4001	早期	06_016_1
4009	未知	16_009_1
4014	后期	04_003_1
4034	早期	12_008_1
4042	早期	06_013_1
4049	未知	06_009_1
4053	早期	01_007_1
4055	早期	11_039_1
4062	未知	12_001_1
4076	后期	01_035_1
4083	早期	11_015_1
4120	早期	11_053_1
4136	后期	07_008_1
4161	未知	16_011_1
4200	未知	06_022_1
4202	未知	07_006_1
4227	未知	01_030_1
4308	后期	01_067_1
4331	样本对MS生成不可用	01_040_1重复(原始样本号在pdf文档中未列示)

4345	后期	11_024_1
4349	未知	13_004_1
4353	后期	11_051_1
4364	早期	11_029_1
4381	早期	01_015_1
4385	早期	01_083_1
4419	未知	11_001_1
4426	早期	01_069_1
4431	未知	01_019_1
4445	早期	11_041_1
4446	未知	01_032_1
4455	早期	11_028_1
4462	早期	01_090_1
4499	早期	02_002_1
4504	早期	01_073_1
4505	未知	16_015_1
4509	早期	11_016_1
4510	后期	01_033_1
4515	早期	12_002_1
4540	早期	11_034_1
4562	早期	01_014_1
4564	早期	04_002_1
4607	未知	01_047_1
4618	早期	06_042_1
4634	早期	01_053_1
4667	未知	13_003_1
4683	早期	14_010_1
4694	后期	06_024_1
4697	早期	06_038_1
4699	早期	11_037_1
4713	后期	01_016_1
4730	早期	01_028_1
4753	早期	06_015_1
4770	早期	06_034_1
4780	后期	06_018_1
4783	后期	01_027_1
4786	未知	04_010_1
4803	早期	01_026_1
4826	早期	01_006_1

4851	早期	01_086_1
4873	未知	12_012_1
4876	早期	11_022_1
4880	早期	01_077_1
4900	早期	01_020_1
4910	早期	06_031_1
4936	早期	01_088_1
4961	后期	01_072_1
4976	早期	01_037_1
4986	后期	15_002_1
5007	未知	01_079_1
5072	未知	11_035_1
5079	早期	03_004_1
5090	早期	11_049_1
5091	早期	01_087_1
5101	未知	01_063_1
5134	早期	12_010_1
5158	后期	07_014_1
5195	早期	01_080_1第二样本
5196	早期	16_014_1
5214	未知	14_009_1
5228	未知	11_036_1
5239	早期	04_009_1
5250	后期	11_021_1
5254	早期	06_026_1
5292	早期	11_004_1
5295	早期	07_005_1
5307	早期	06_025_1
5330	后期	11_045_1
5336	未知	10_003_1
5351	早期	06_033_1
5352	后期	16_010_1
5358	未知	13_001_1
5362	后期	04_004_1
5374	未知	02_003_1
5391	早期	01_064_1
5395	早期	06_032_1
5401	后期	01_092_1
5411	早期	13_002_1

5424	后期	01_043_1
5431	未知	02_004_1
5440	早期	06_029_1
5443	未知	12_011_1
5444	早期	11_006_1
5447	未知	01_003_1
5448	未知	04_006_1
5456	早期	14_011_1
5466	早期	14_004_1
5497	未知	16_003_1
5505	早期	01_002_1
5507	早期	12_005_1
5512	后期	01_070_1
5567	未知	02_001_1
5573	早期	01_022_1
5583	早期	04_012_1
5587	早期	12_004_1
5594	早期	06_041_1
5638	早期	11_023_1
5658	早期	01_011_1
5663	早期	01_094_1
5671	早期	11_031_1
5672	早期	01_056_1
5673	早期	01_004_1
5680	后期	14_003_1
5713	早期	01_009_1
5714	后期	06_005_1
5721	未知	01_071_1
5724	早期	08_002_1
5725	未知	06_019_1
5747	早期	01_065_1
5755	早期	01_042_1
5767	未知	07_004_1
5791	早期	06_037_1
5801	后期	11_018_1
5813	早期	11_027_1
5820	后期	01_018_1
5842	后期	03_005_1
5847	未知	11_054_1

5869	早期	14_005_1
5874	早期	15_001_1
5910	未知	01_091_1
5911	早期	06_035_1
5913	早期	03_002_1
5935	早期	16_018_1
5963	早期	06_039_1
5970	后期	01_054_1
5975	早期	01_046_1
5976	早期	01_085_1
5997	未知	14_002_1
6048	早期	01_017_1
6056	未知	16_007_1
6082	早期	11_014_1
6093	早期	07_001_1
6098	后期	11_017_1
6105	未知	16_002_1
6122	早期	06_010_1第二样本
6130	早期	14_007_1
6140	未知	07_003_1
6156	后期	11_011_1
6161	早期	01_068_1
6182	早期	11_020_1
6193	未知	16_008_1
6203	早期	11_013_1
6235	未知	11_010_1
6260	早期	01_045_1
6270	早期	11_052_1
6278	早期	06_008_1
6281	早期	04_008_1
6282	未知	06_022_1
6295	早期	11_009_1
6296	早期	01_041_1
6297	未知	01_081_1
6299	早期	14_014_1
6321	早期	11_057_1
6336	后期	01_023_1
6349	后期	10_001_1
6361	未知	03_003_1

6390	早期	01_078_1
6398	未知	06_001_1
6419	后期	01_044_1第二样本
6424	早期	06_023_1
6438	未知	16_001_1
6439	早期	01_036_1
6442	早期	10_004_1
6476	早期	01_084_1
6487	样本对MS生成不可用	11_048_1
6492	后期	01_057_1
6572	未知	13_006_1
6585	早期	01_076_1
6604	早期	11_002_1
6622	早期	01_031_1
6625	早期	06_011_1
6626	早期	06_003_1
6667	未知	11_025_1
6712	早期	01_038_1
6718	早期	07_013_1
6729	早期	06_036_1
6737	早期	06_006_1
6741	早期	16_004_1
6752	早期	11_019_1
6761	后期	06_027_1
6770	早期	11_007_1
6795	未知	11_038_1
6797	早期	01_058_1
6824	未知	04_007_1
6827	早期	06_007_1
6847	早期	04_011_1
6854	早期	07_002_1
6886	未知	01_012_1
6887	后期	01_051_1
6932	早期	01_005_1
6939	后期	14_008_1
6947	早期	11_032_1
6977	早期	07_009_1
6981	未知	06_028_1
6982	早期	13_007_1

6992	后期	11_040_1
6998	未知	06_017_1

附录H:用于谱获取的仪表的细节

运行	日期	序列号	合格日期
140131_ItalianABC	2/3/2014-2/4/2014	260	1/30/2014 NRS 1/27/2014 RuO
140225_ItalianABC	2/25/2014	260	2/25/2014 NRS
140130_Furb_PROSE* ²	1/30/2014-1/31/2014	260	1/30/2014 NRS
140115_PROSE	1/15/2014-1/17/2014	258	12/11/2013 *
131118_ItalianABC	11/18/2013-11/19/2013	258	11/12/2013 RuO

*这是两个样本具有无法获取的点的快速一致性校验,但是如果丢弃这两个样本则是一致的。

[0182] *² 这个运行对于与从仪表258的140115_PROSE运行相同的板进行。

[0183] [1] 参见V. Gregorc等人的 Randomized Proteomic Stratified Phase III Study of Second-Line Erlotinib Versus Chemotherapy in Patients with Inoperable Non-Small Cell Lung Cancer,在ASCO annual meeting提交的简报(2013年6月)。

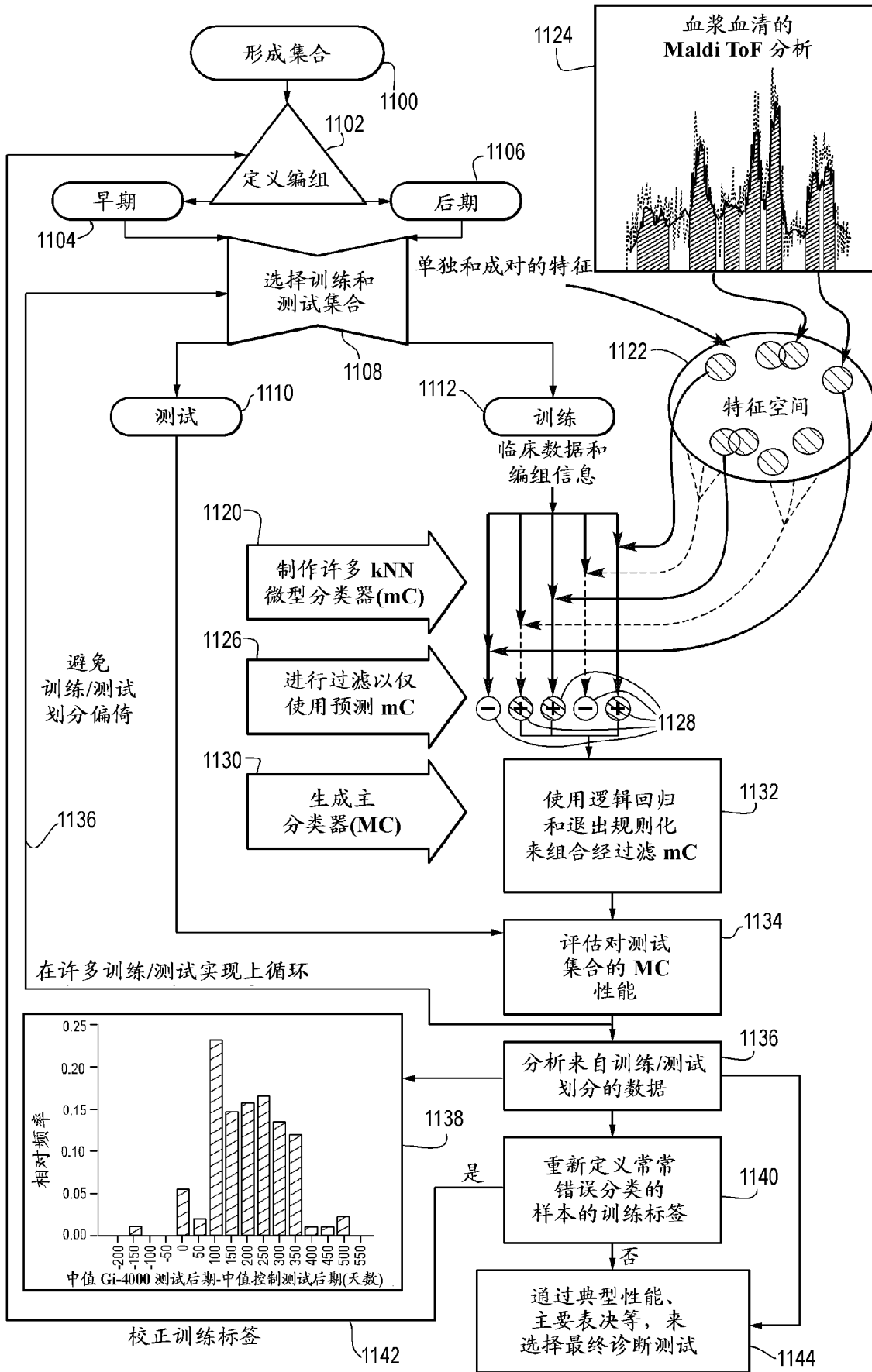


图 1

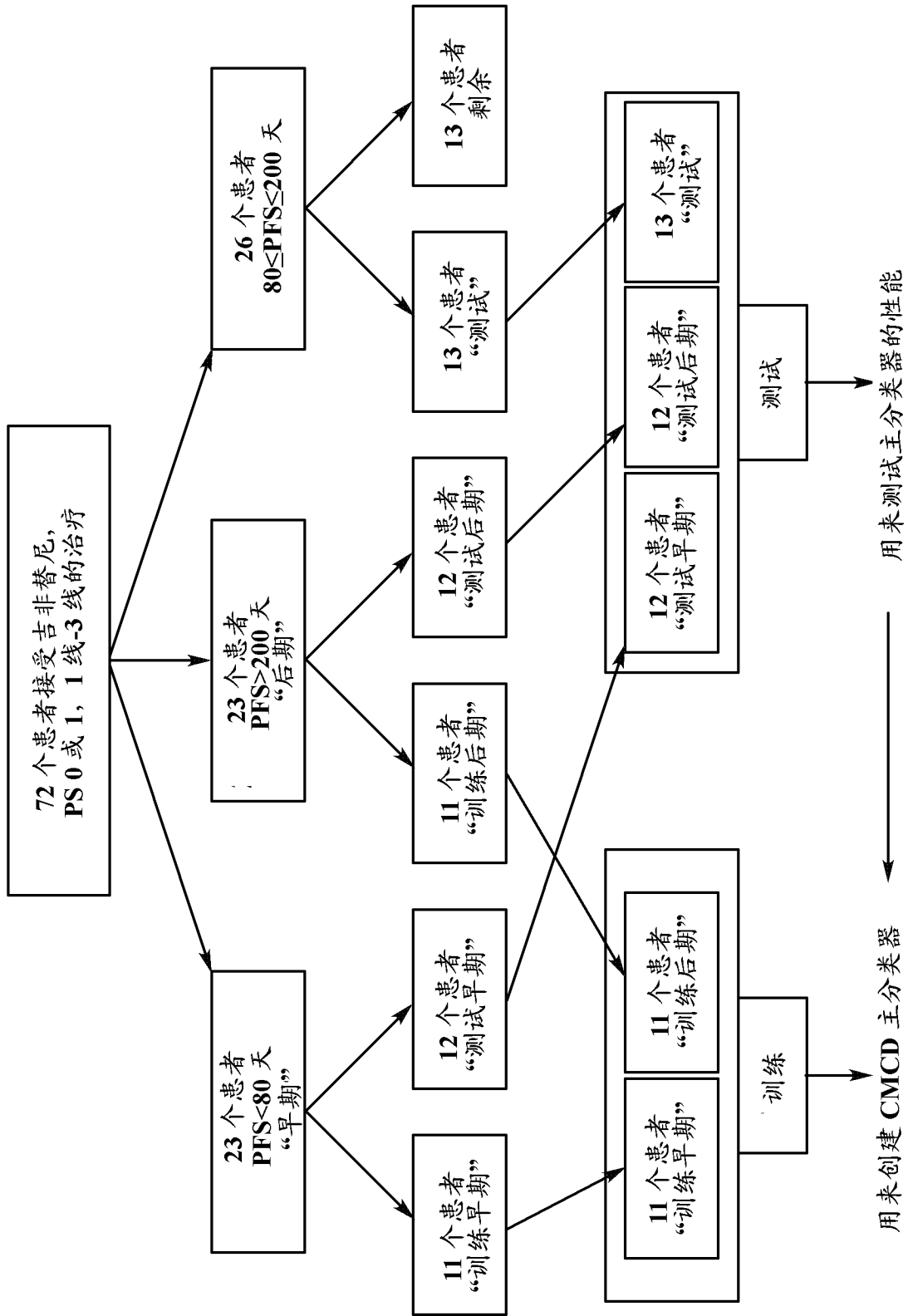


图 3

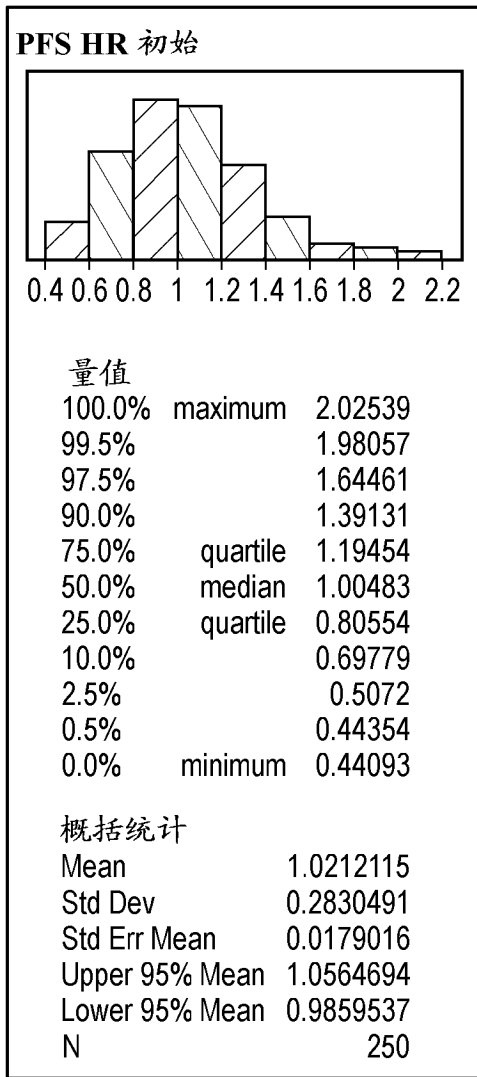


图 4A

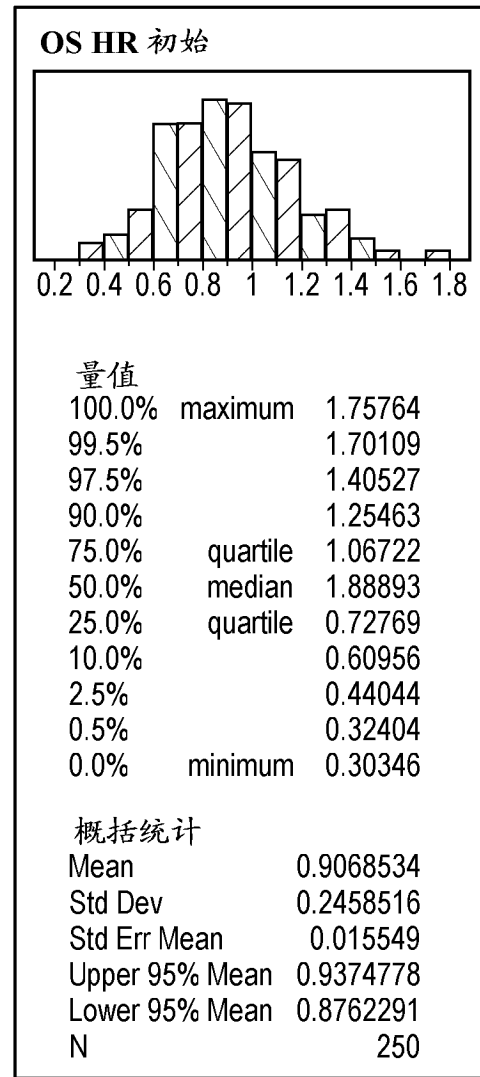


图 4B

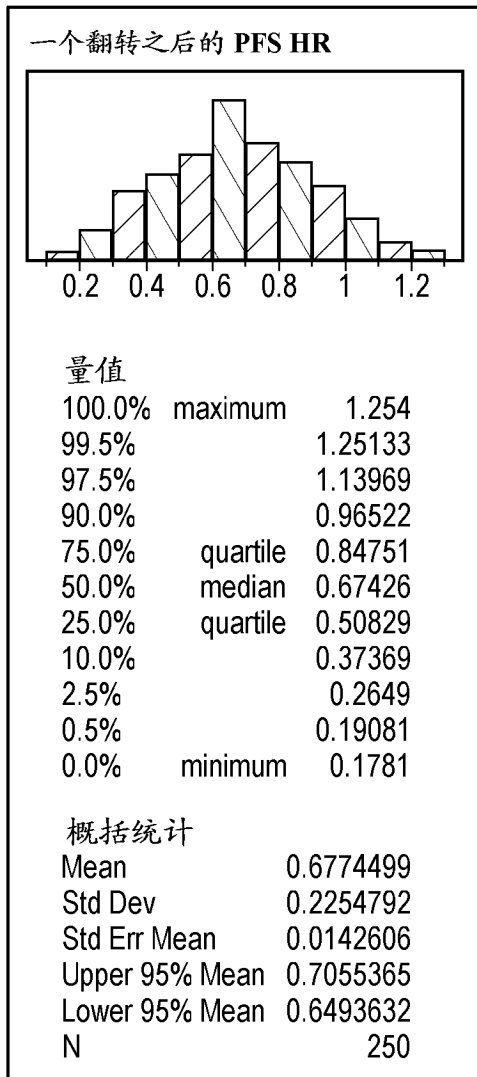


图 4C

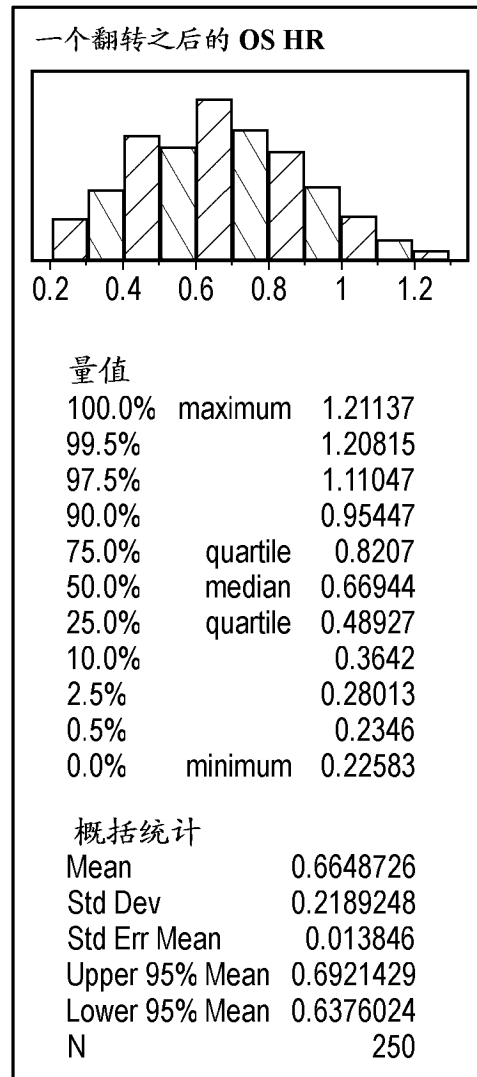


图 4D

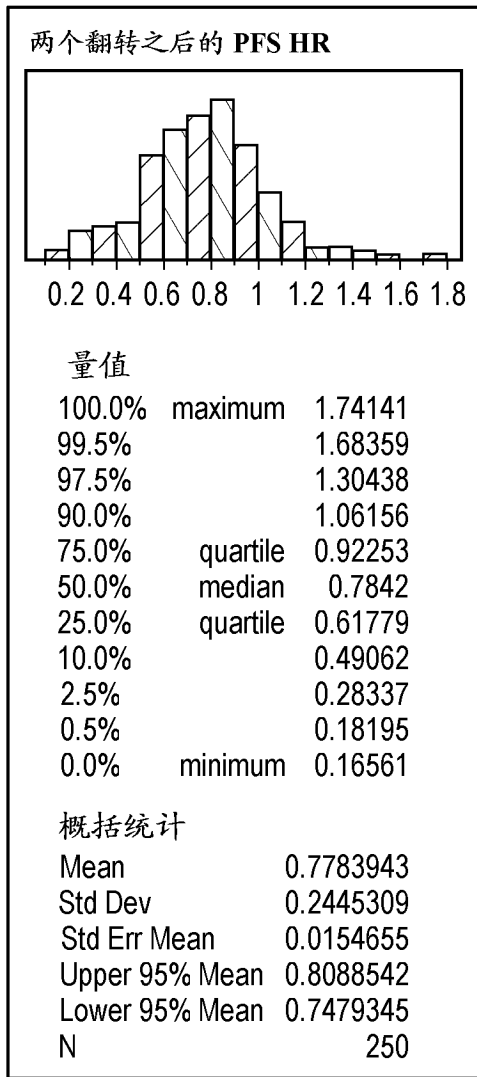


图 4E

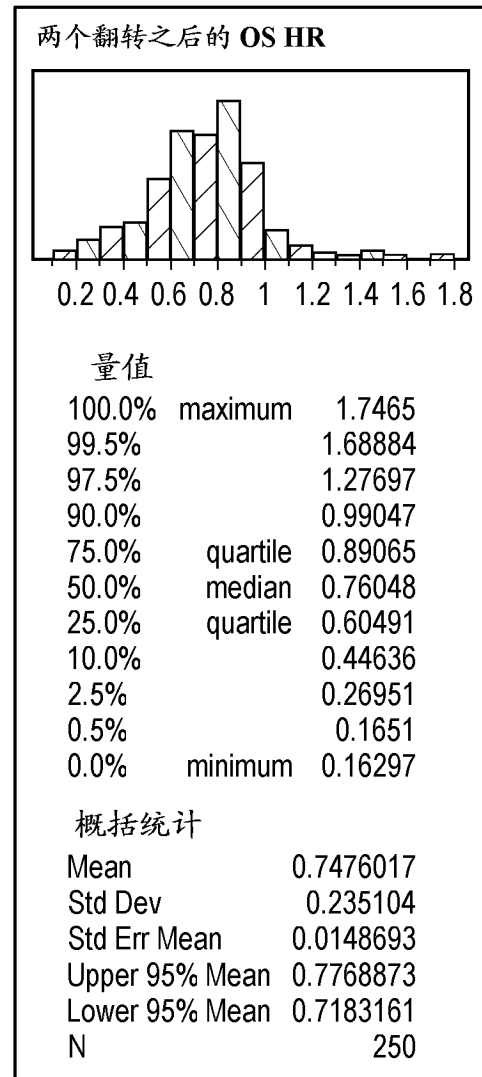


图 4F

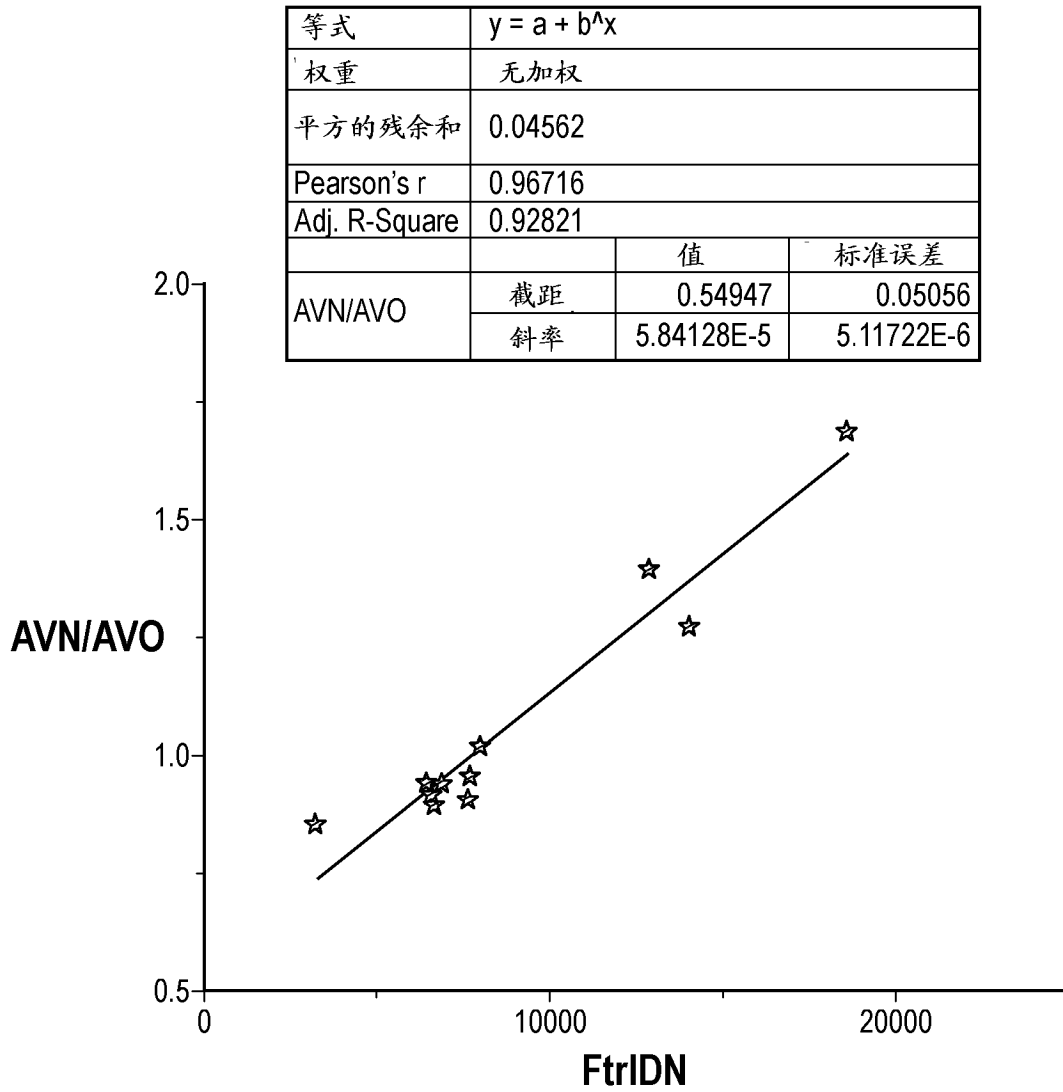


图 5

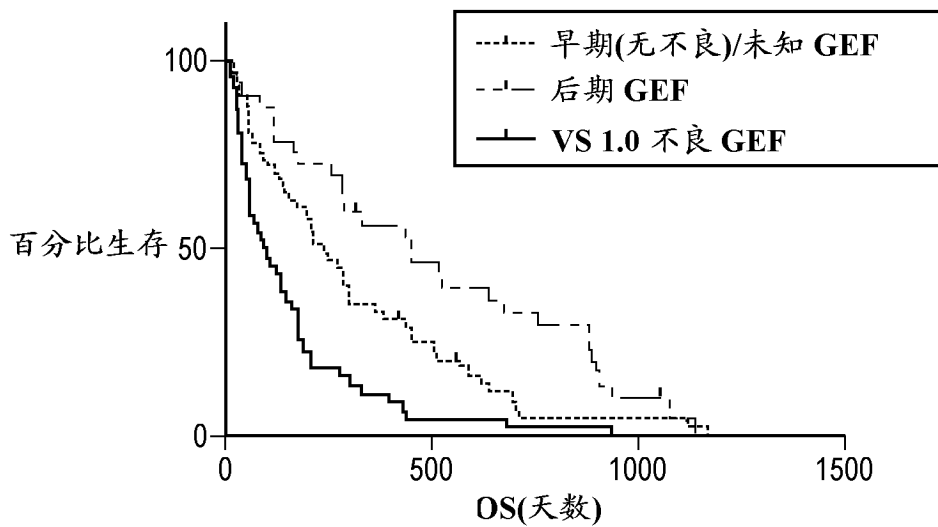


图 6A

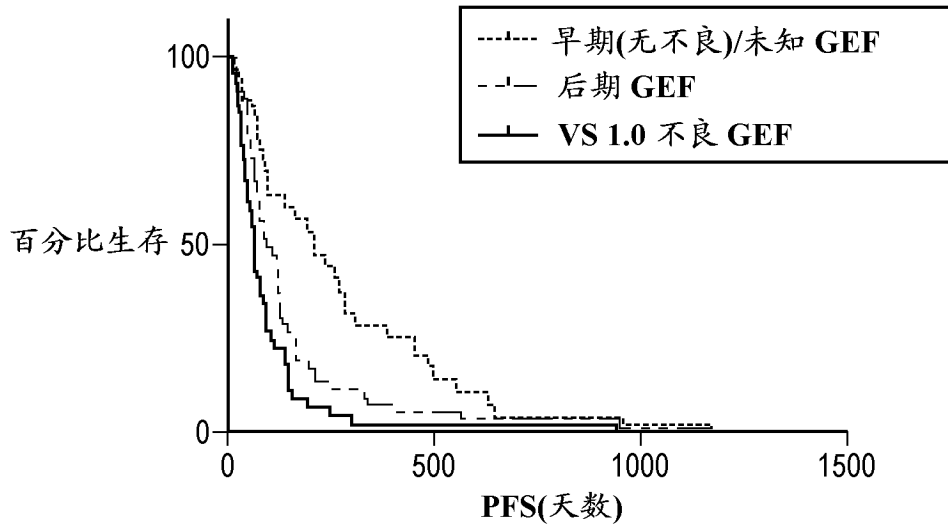


图 6B

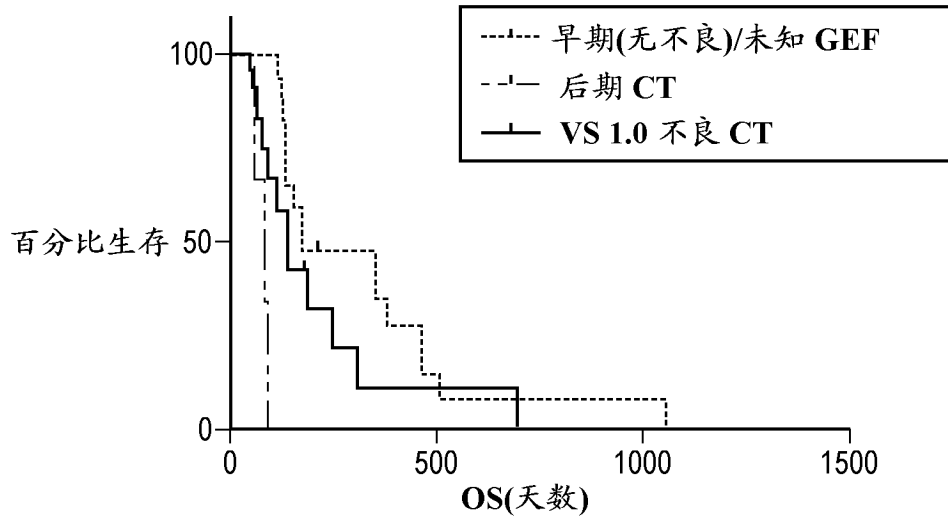


图 6C

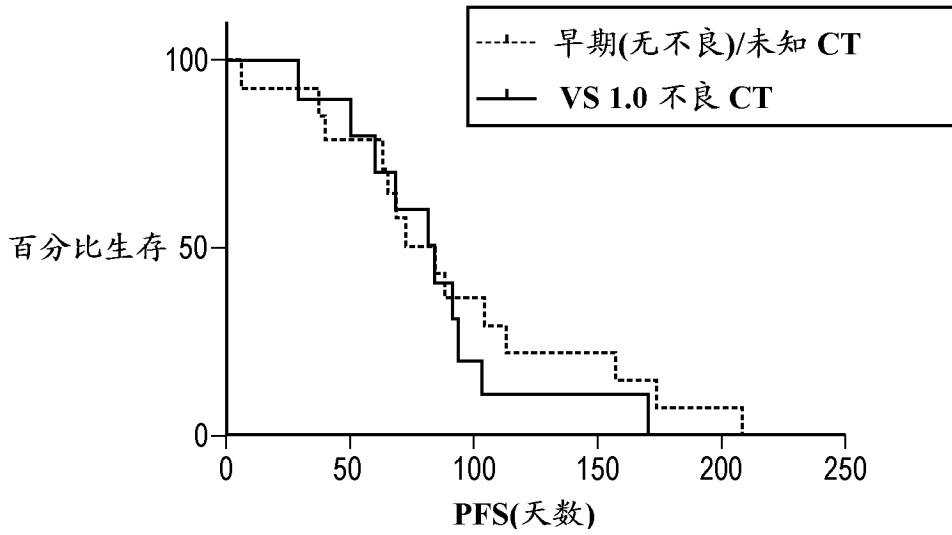


图 6D

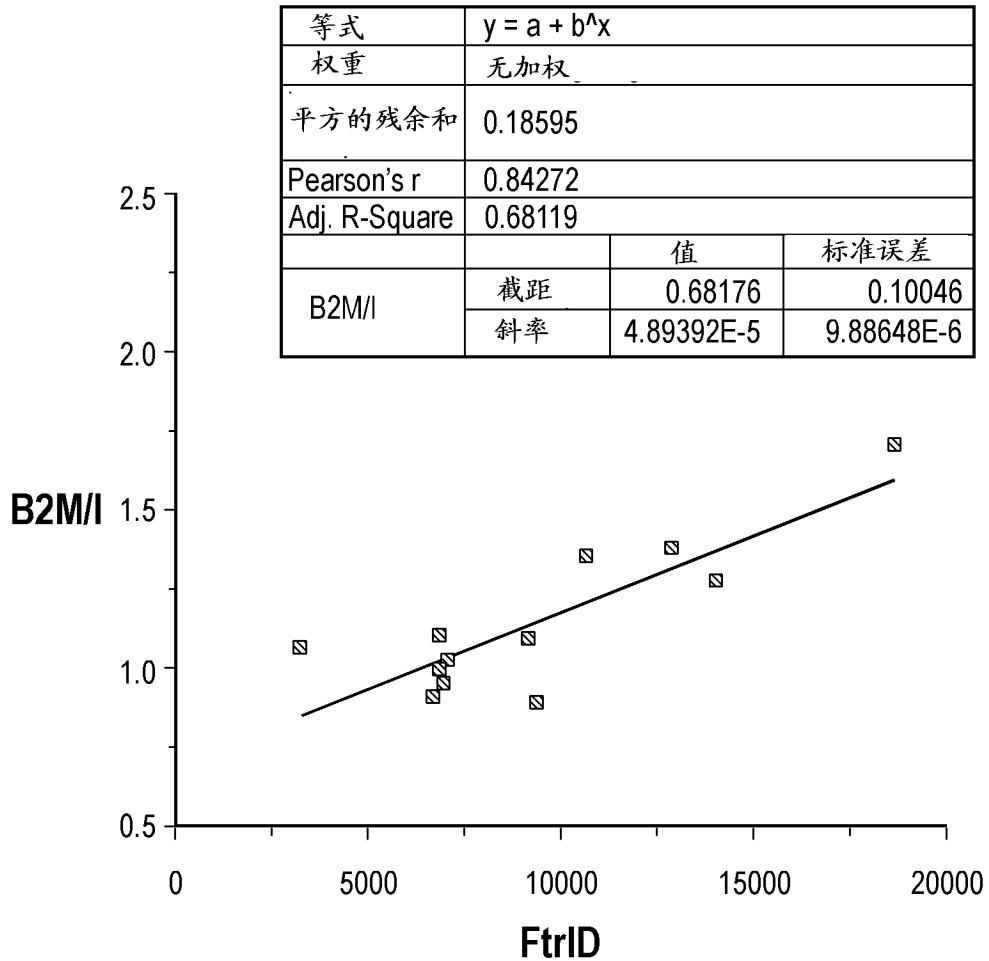


图 7

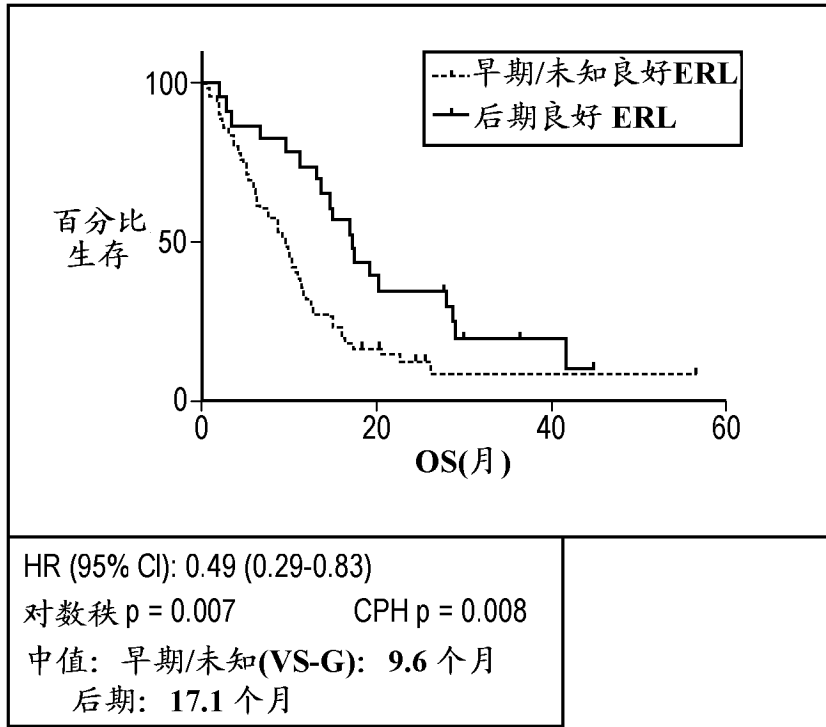


图 8A

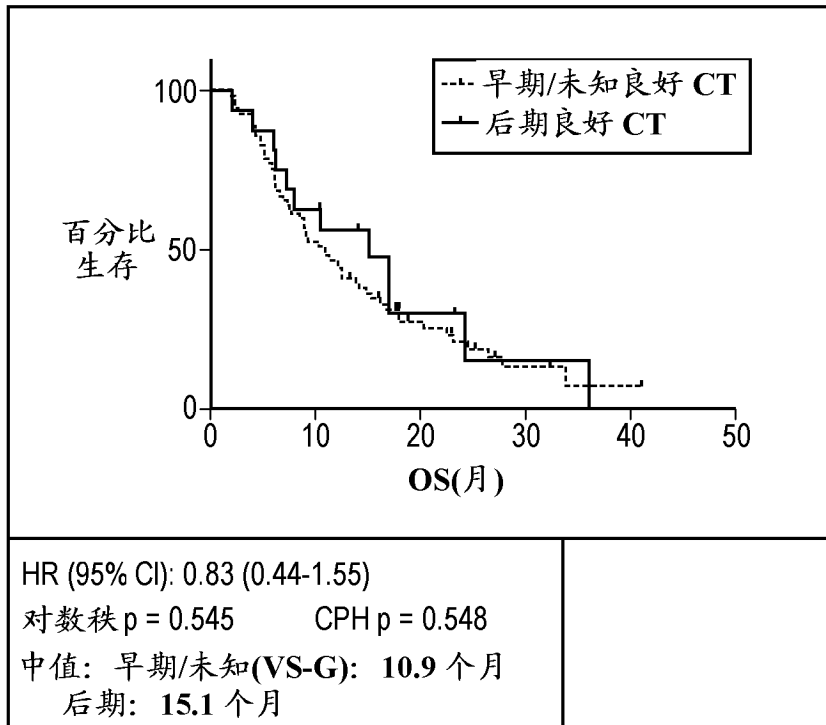


图 8B

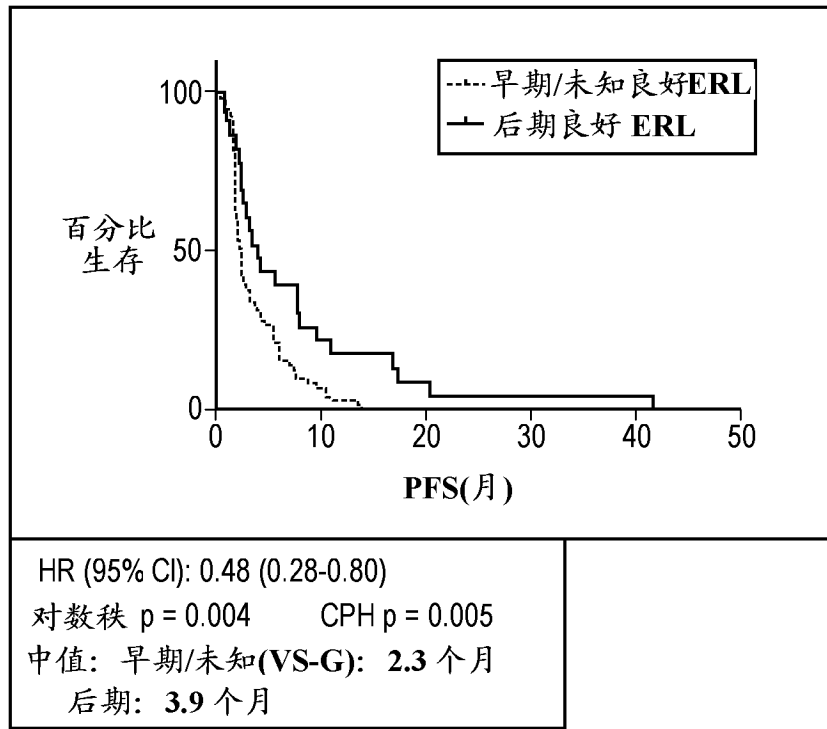


图 9A

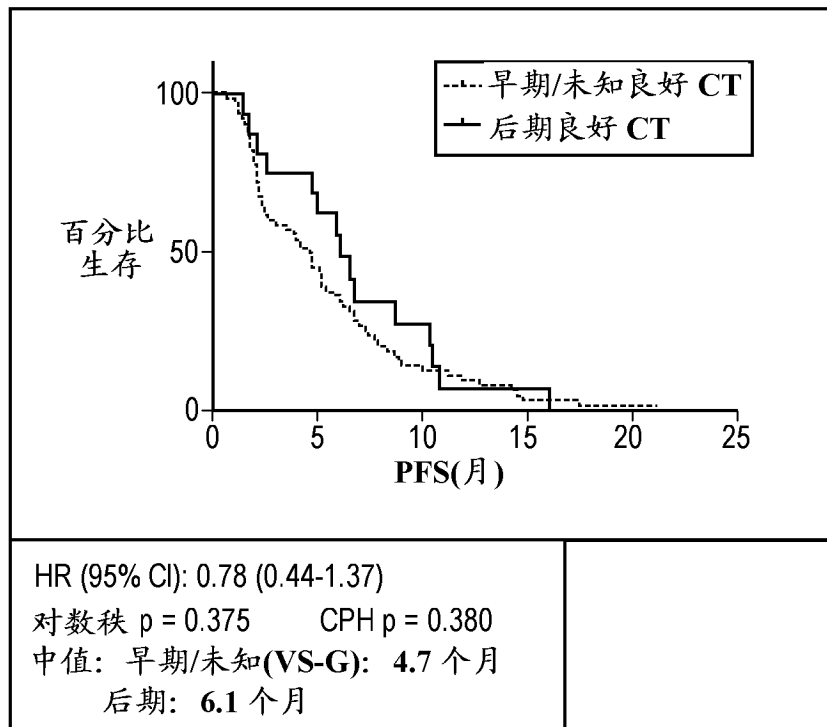


图 9B

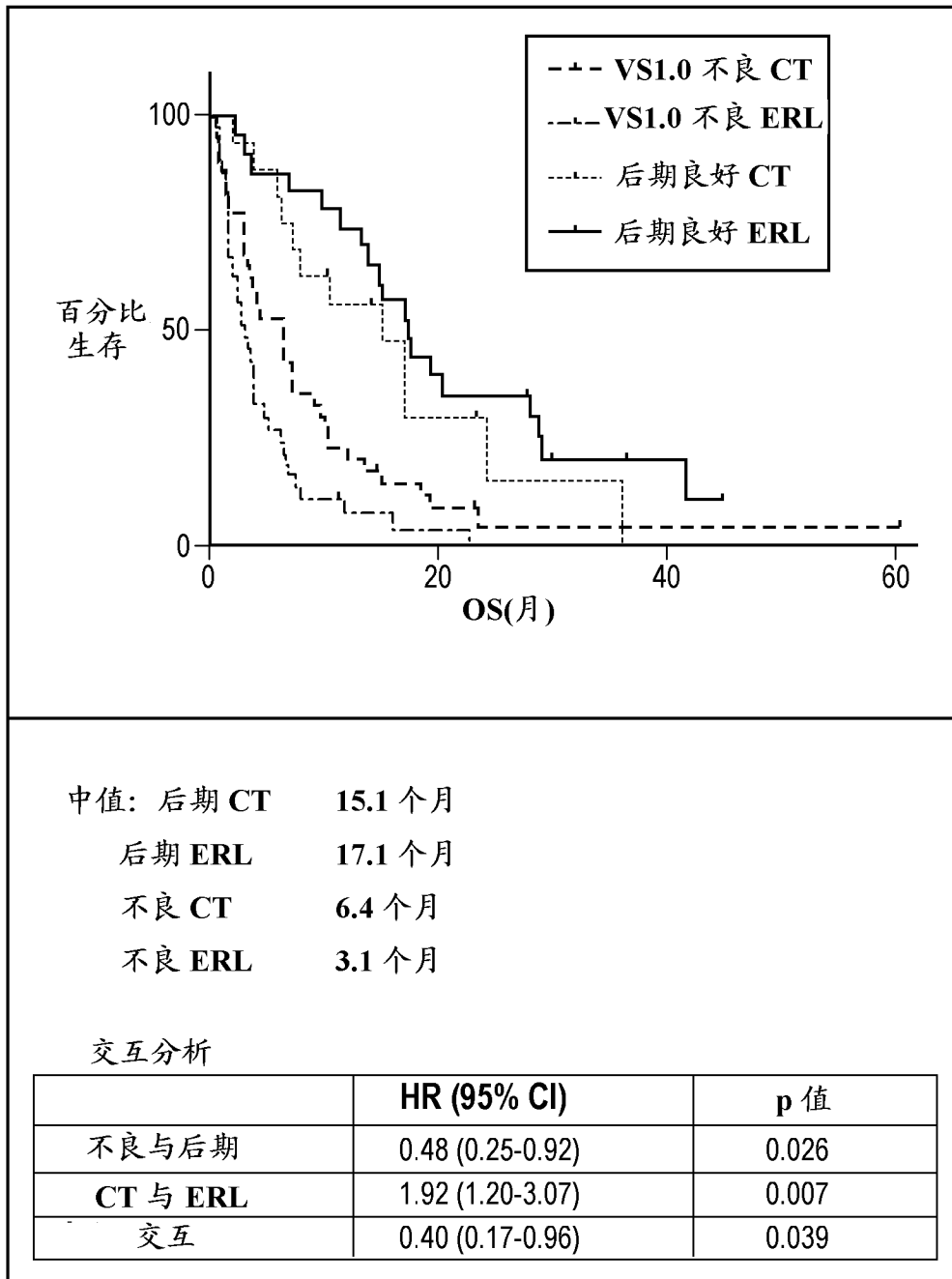


图 10

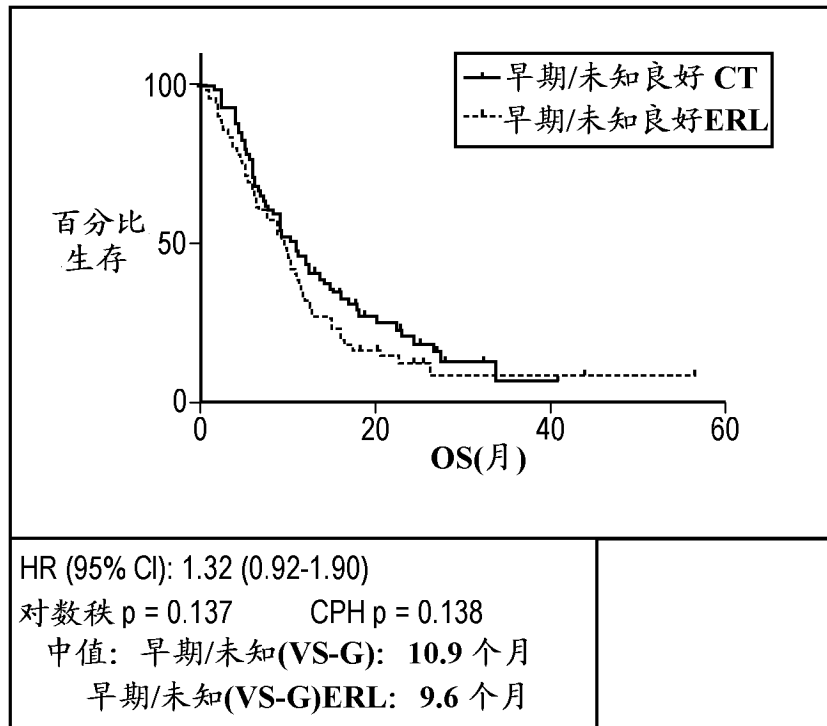


图 11

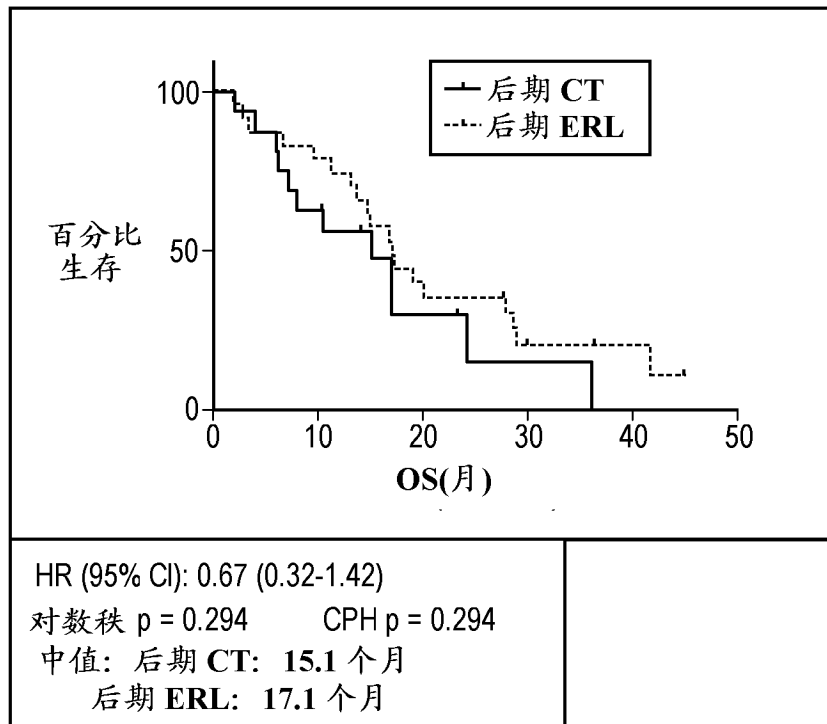


图 12A

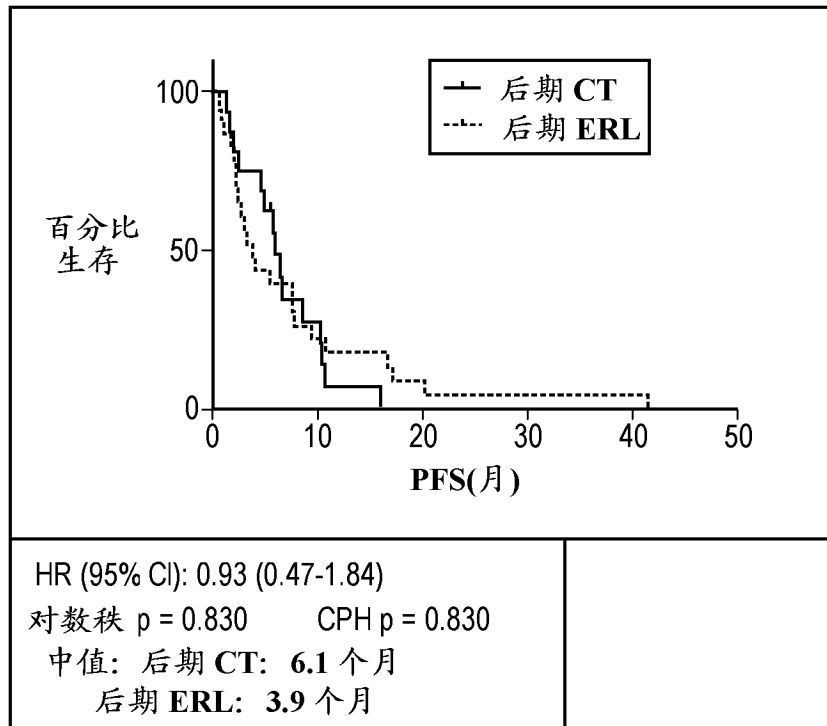


图 12B

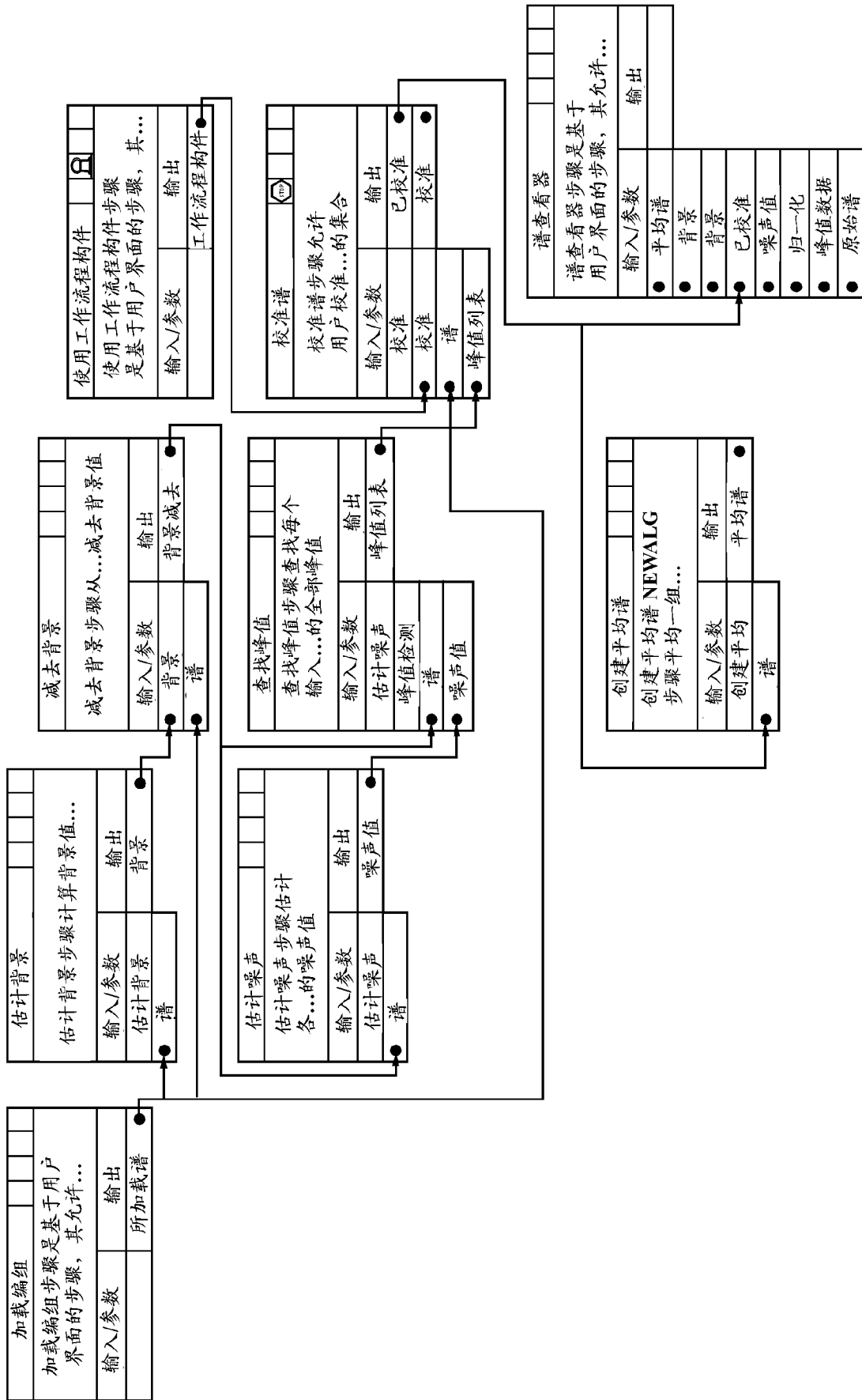


图 13

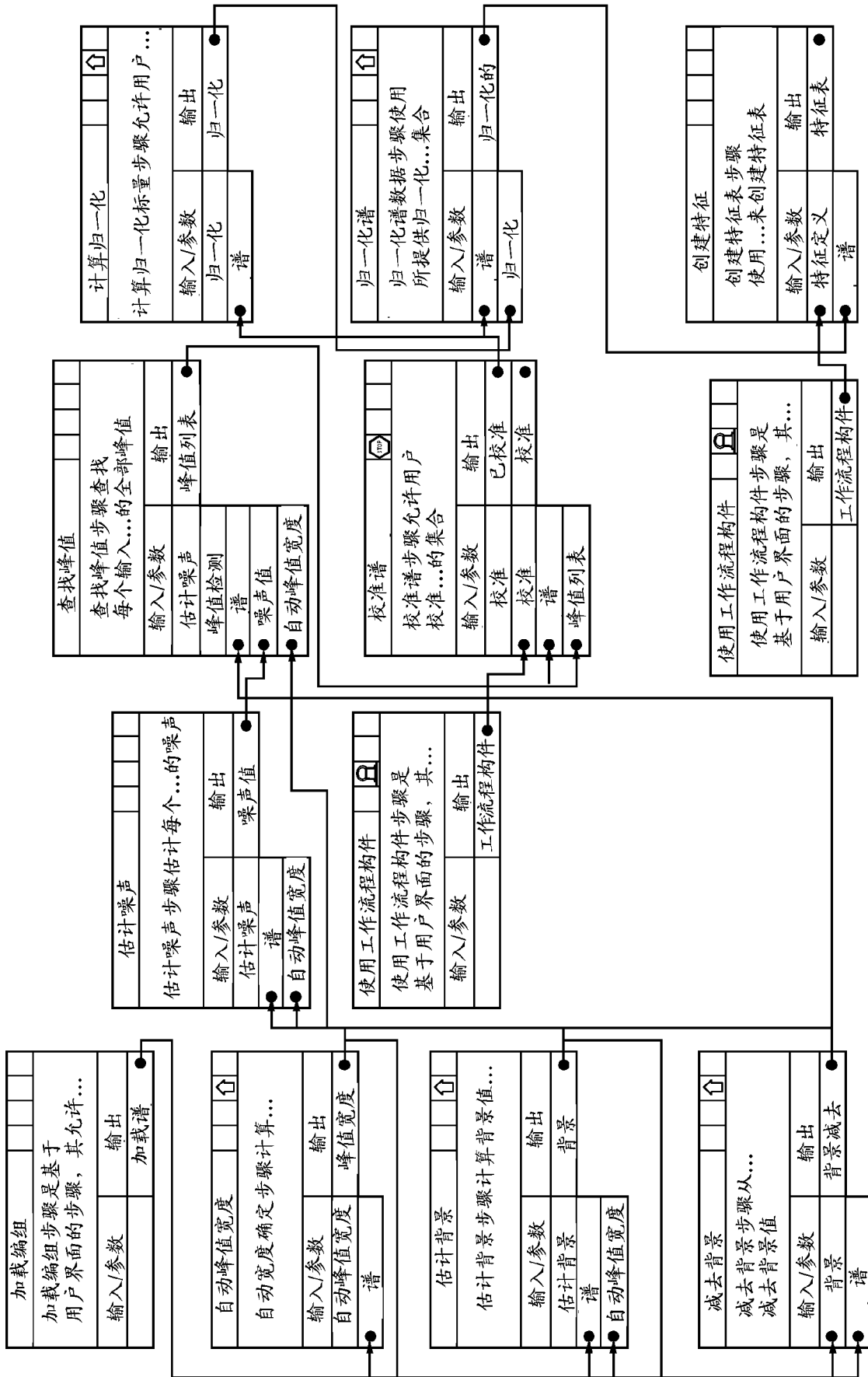


图 14

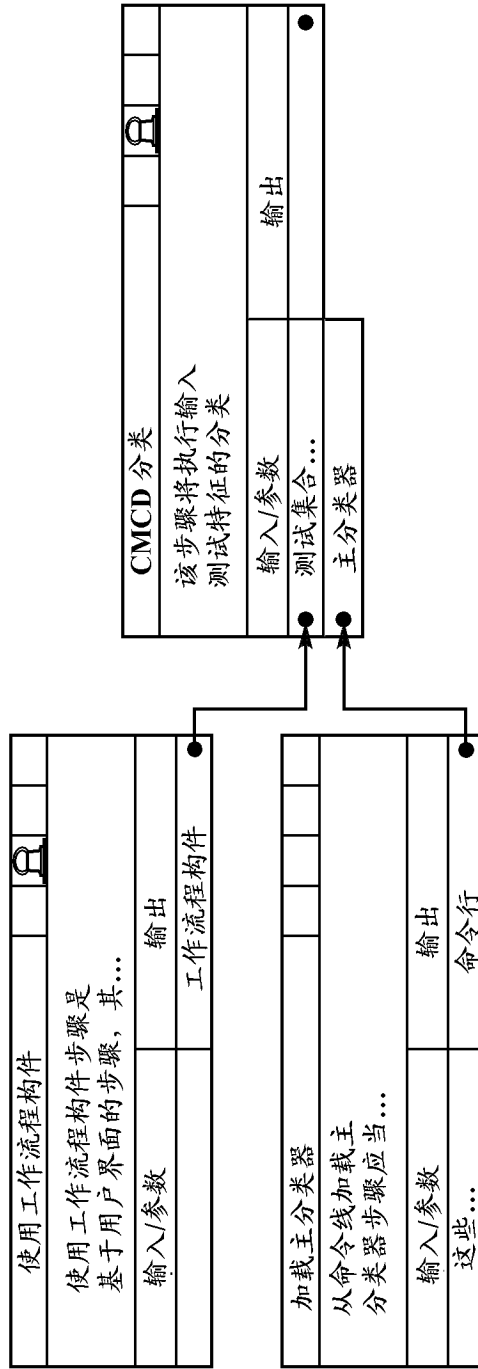


图 15

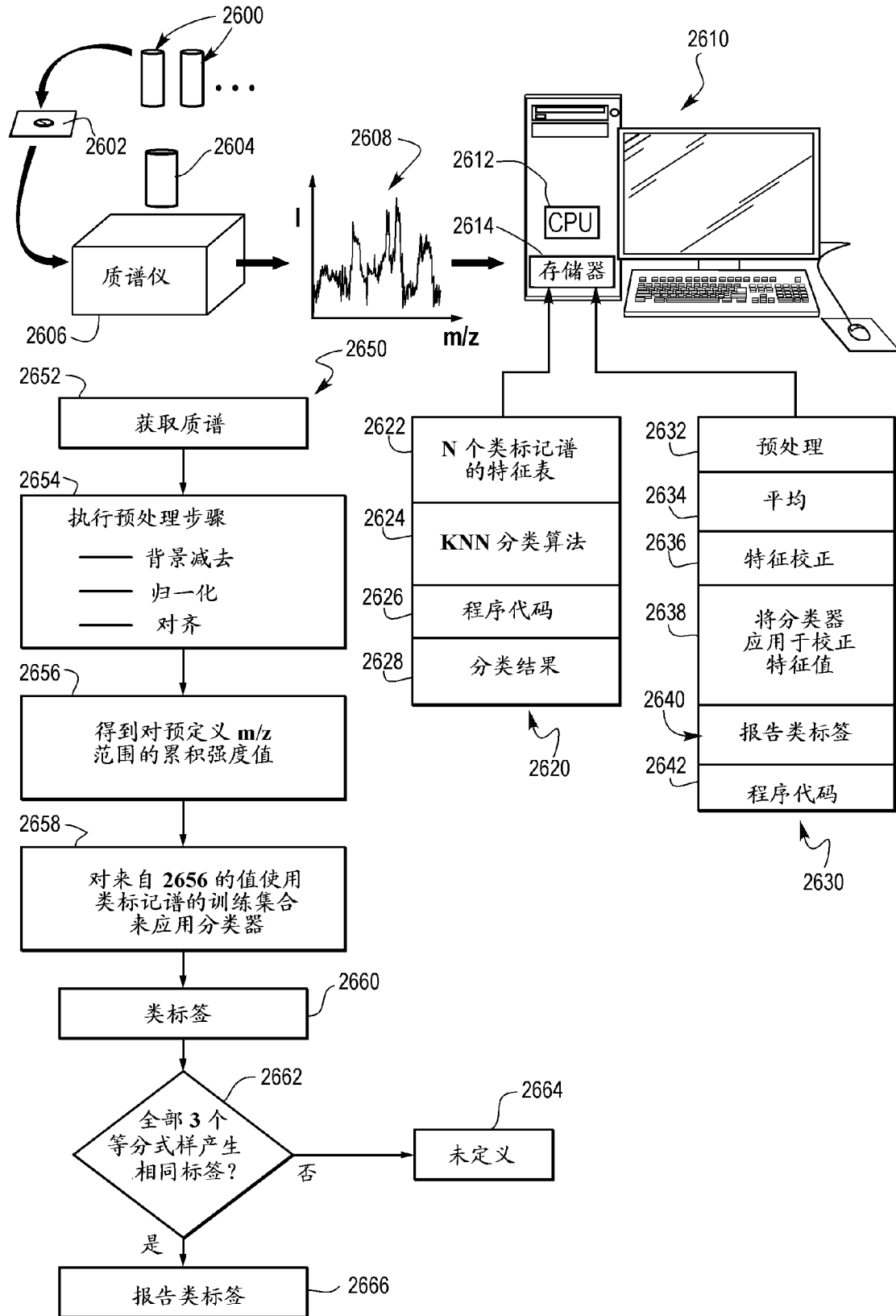


图 16