

# 血糖近红外光谱分析的 Savitzky-Golay 平滑模式与偏最小二乘法因子数的联合优选

谢军<sup>1</sup> 潘涛<sup>2</sup> 陈洁梅<sup>3</sup> 陈华舟<sup>1</sup> 任小焕<sup>1</sup>

<sup>1</sup>(广东省高等学校光电信息与传感技术重点实验室(暨南大学), 广州 510632)

<sup>2</sup>(暨南大学生物工程学系, 广州 510632) <sup>3</sup>(上海大学数学系, 上海 200442)

**摘要** 利用偏最小二乘法(PLS)和光谱 Savitzky-Golay(S-G)平滑方法, 建立血清葡萄糖近红外光谱分析的优化模型。基于最优单波数模型的预测效果, 提出划分校正集和验证集的一种新方法。采用 Savitzky-Golay(S-G)平滑和 PLS 的组合波段, 光谱经过 S-G 平滑处理, 利用 PLS 方法建立定标预测模型。将平滑点数扩充为 8, 10, 12, 14, 16, 18, 20 (奇数), 多项式次数扩充为  $n \leq m$ ,  $n \leq 1$  得到包含  $8^m$  个平滑模式的  $m$  个平滑系数表。对所有平滑模式和 PLS 因子数(PLS)分别建立 PLS 模型。按照预测效果进行优选, 得到最优 S-G 平滑模式为 10 阶导数平滑, 16 次多项式类型, S-G 平滑点数为 18, 最优 PLS 因子数为 0, 最优  $R^2$  达到 0.999 6648。所采用的划分校正集和验证集的方法、S-G 平滑模式的扩充、S-G 平滑模式和 PLS 因子数的联合大范围筛选能够有效地应用于近红外光谱分析的模型优化。

**关键词** 血糖; 近红外光谱; 偏最小二乘法; Savitzky-Golay 平滑; 校正集验证集划分

## 1 引言

随着光谱技术和化学计量学的快速发展, 近红外光谱以其分析效率高、速度快、成本低、非破坏性和易于在线分析等特点已广泛应用于农业、食品、烟草、医药等领域<sup>[1-3]</sup>。模型优化对于提高近红外光谱预测能力具有重要意义。偏最小二乘法(PLS)是融合主成分分析和多元线性回归的一种有效的化学计量学方法<sup>[4,5]</sup>, 其中合理选用 PLS 因子数, 对于充分利用光谱信息和消除噪声非常重要。在光谱预处理中, 平滑可以保留光谱轮廓而消除噪声, 求导则可以有效消除基线漂移、倾斜等噪声。S-G 平滑方法是应用十分广泛而有效的平滑和求导预处理方法<sup>[6,7]</sup>。按照导数阶数(平滑看成  $n$  阶求导)和多项式次数和平滑点数的不同, S-G 平滑模式有很多种, 计算公式也各不相同。其中平滑点数的设置非常重要, 点数过少容易产生新误差, 点数过多则容易使包含信息的光谱数据磨光丢失, 都会造成模型精度下降。根据预测效果对 S-G 平滑模式与 PLS 因子数联合筛选是很有必要的, 但由于工作量庞大, 既往的研究很少做到这一步。另一方面, 考虑到有些实际测量体系可能需要更多的平滑点数, 比如测量数据波长间隔小的情形, 相邻波长点的数据过于相似, 点数少的平滑效果往往不够好。为了拓宽适用范围, 有必要按照原始论文的方法<sup>[6]</sup>扩充平滑系数表。

血糖近红外光谱分析及其模型优化是很重要的研究方向<sup>[1,2]</sup>。本实验以血糖近红外光谱分析为例, 研究 S-G 平滑模式与 PLS 因子数的联合优化设计在近红外光谱分析模型优化中的作用。为了改善模型预测能力, 基于最优单波数模型提出了划分校正集和验证集的新方法。

## 2 实验部分

### 2.1 实验材料、仪器和测量方法

100 份血清样品由广州市某医院提供, 样品葡萄糖的含量由全自动生化分析仪测定作为光谱分析的参考化学值。全体化学值范围 0.81~1.08 6648 g, 均值、标准偏差分别为 0.93 和 0.08 6648 g。实验仪器为 8000 傅里叶变换型近红外光谱仪(美国 Thermo 公司), 探测器为铟镓砷(InGaAs) 用光程

收稿日期: 2010-08-10; 接受日期: 2010-10-10

本文系国家自然科学基金(40701000)和广东省自然科学基金(07050800)和广东省科技计划项目(200701000000000000)和广州市科技攻关项目(080100000000000000)资助

# 66 的石英比色皿测量光谱, 扫描谱区  $4000 \sim 600 \text{ cm}^{-1}$ , 分辨率  $4 \text{ cm}^{-1}$ , 扫描次数 10。

## 2.2 校正集和验证集的划分方法

基于全体样品最优单波数模型的预测效果给出划分校正集验证集的一种新方法。根据比尔定律, 考虑血清样品吸光度与葡萄糖化学值的单波数线性模型

$$A(v) = k(v)C + \epsilon \quad (1)$$

其中  $A(v)$  为样品在波数  $v$  的吸光度,  $k(v)$  为在波数  $v$  的葡萄糖单位浓度吸光系数,  $C$  为样品的葡萄糖浓度化学值,  $\epsilon$  为其它未知干扰。在每个波数  $v$ , 利用全体样品的吸光度和化学值回归计算  $k(v)$ , 再利用  $k(v)$  和样品吸光度计算样品  $i$  的预测值  $C'_i(v)$  ( $i = 1, 2, \dots, N$ ),  $N$  是全体样品个数。进一步计算预测值与化学值的均方根偏差 (RMSE), 设  $C_i$  为样品  $i$  的化学值, 则

$$\text{RMSE}(v) = \sqrt{\frac{\sum_{i=1}^N (C'_i(v) - C_i)^2}{N - 1}} \quad (2)$$

按 RMSE 值最小选出最优单波数模型和相应波数  $v_{0 \leq v \leq 8}$ 。根据最优单波数模型计算每个样品的浓度预测值与化学值的偏差, 称为单波数预测偏差 ( $\text{K}'_i = |C'_i(v_{0 \leq v \leq 8}) - C_i|$ ,  $i = 1, 2, \dots, N$ )。

$$\text{K}'_i = |C'_i(v_{0 \leq v \leq 8}) - C_i|, \quad i = 1, 2, \dots, N \quad (3)$$

$\text{K}'_i$  是吸光度和化学值的一种关联指标。根据  $\text{K}'_i$  划分校正集检验集。利用计算程序筛选使两个集合的  $\text{K}'_i$  分布一致 (均值和标准偏差相近, 相对误差小于 5%)。将化学值和光谱数据结合起来使校正集验证集具有相似性, 从而具有建模代表性。为了使得校正集浓度范围能够涵盖验证集浓度范围, 将化学值最大和最小的样品放在校正集, 化学值次大次小的样品放在验证集。

## 2.3 SG 平滑方法

S-G 平滑的参数包括导数阶数  $s$ 、多项式次数  $n$  和平滑点数  $m$ 。S-G 平滑把光谱区间的若干个连续点作为一个窗口, 窗口内每点用多项式 (以点的编号  $r$ ,  $r = 1, 2, \dots, m$  为变量) 来做实测数据的最小二乘拟合。拟合后, 多项式在编号为  $s$  (中心点) 的值就是 S-G 平滑值, 多项式对编号求导后在编号为  $s$  (中心点) 的值就是 S-G 导数值。按上述程序, 窗口中心点的平滑值和各阶导数值都可以表示为窗口内各点实测数据的线性组合。线性组合的系数 (即平滑系数) 由平滑点数 (即窗口内的点数)  $m$ 、多项式次数和导数阶数唯一确定。通过窗口移动, 得到每个窗口中心点的平滑值和各阶导数值, 从而得到原谱的 S-G 平滑谱和 S-G 导数谱。为了拓宽应用范围, 本研究将平滑点数从原有的  $8 \sim 16$  之间奇数<sup>[1]</sup>扩充为  $8 \sim 20$  之间的奇数, 多项式次数扩充为  $n = m - 1$  (原为  $n = m - 1$ ,  $C \leq 8$ ), 按照原方法<sup>[1]</sup>编写程序计算, 得到 100 个涵盖原有平滑系数的平滑系数表, 共有 100 个平滑模式 (原有 10 个), 是适用范围更宽的 S-G 平滑预处理群。

## 2.4 模型的评价指标

模型评价指标主要包括预测均方根偏差 (RMSE) 和预测相关系数 ( $R_0$ )

$$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^M (C'_{ip} - C_{ip})^2}{M - 1}}, \quad R_0 = \frac{\sum_{i=1}^M (C_{ip} - C_{mp})(C'_{ip} - C'_{mp})}{\sqrt{\sum_{i=1}^M (C_{ip} - C_{mp})^2 \sum_{i=1}^M (C'_{ip} - C'_{mp})^2}} \quad (4)$$

其中  $C'_{ip}$  和  $C_{ip}$  分别为验证集中第  $i$  个样品的预测值和化学值,  $C'_{mp}$  和  $C_{mp}$  分别为验证集样品的预测值均值和化学值均值,  $M$  为验证集的样品个数。 $R_0$  与 RMSE 是有一定关联的, RMSE 值低,  $R_0$  一般也较高。本研究以 RMSE 为优化目标来进行参数设计和模型优选。

# 3 结果与讨论

## 3.1 样品光谱、校正集和验证集的划分

100 个血清样品的近红外光谱如图 1 所示。光谱在  $1450 \sim 850 \text{ cm}^{-1}$  附近有水分子的强烈吸收, 除了水的吸收峰外没有其它显著的吸收峰, 光谱重叠严重, 吸收较弱。考虑到在  $850 \sim 1450 \text{ cm}^{-1}$  附近

吸收强烈,光谱能量低,信息含量差,噪音大,故把这两段(吸光度高于 0.2 的波段)光谱数据扣除后用于建模。用于建模的光谱波段是 4000-5300  $\text{cm}^{-1}$  和 5300-6087.5  $\text{cm}^{-1}$  两段的组合。

按照 2.1 节的方法,建立每个波数点的吸光度和化学值的单波数模型,按照  $\alpha_1:2$  最小找到最优波数  $\nu_{0 \leq \alpha_1 \leq 6}$  为 0-6。根据 0-6 对应的最优单波数模型计算每个样品的  $K_9'$ ,全体样品的  $K_9'$  和化学值分布如图 3 所示。由图 3 可见,全体样品的化学值和  $K_9'$  分布均匀,无显著的异常样品。因此,全体样品都用于建模。按照大约 1% 的比例,校正集 1% 个样品,验证集 1% 个样品,按照 2.1 节方法划分校正集验证集,得到的校正集验证集的  $K_9'$ 、化学值的均值和标准偏差如表 3 所示。表 3 和图 3 都表明,校正集验证集的化学值和  $K_9'$  分布都非常一致。

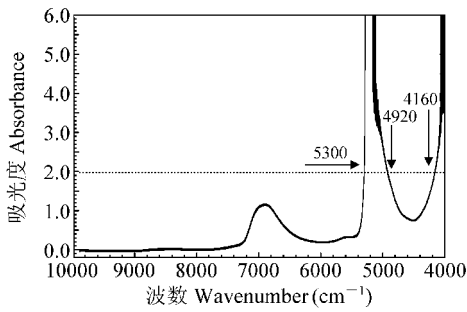


图 2 100 个血清样品的近红外光谱  
—(0% . 7;L);FbL;L7R 507M=L; 4b %D% 57LH6 5; 60B75

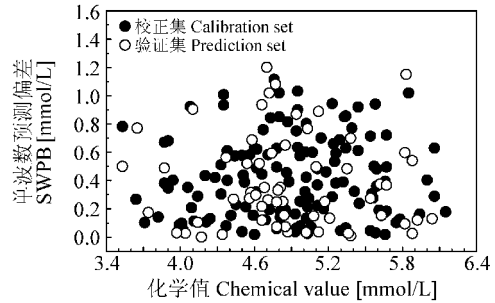


图 3  $K_9'$  与化学值的分布  
—(0# -+5-L+JH=4F 4b 5+F (B7 P;<7FH6J7L 0L7R+=4F J+;5  
( $K_9'$ );FR MN76+=;B <;BH75

### 3.2 SG 平滑模式与 PLS 因子数的联合优选

为了比较,在  $\alpha$  平滑前直接  $\alpha$  方法建模。采用 4000-5300  $\text{cm}^{-1}$  和 5300-6087.5  $\text{cm}^{-1}$  组合波段,  $\alpha$  因子数设置从 0 到 5,按照  $\alpha_1:29$  最小遴选最优因子数为 0,最优  $\alpha_1:29$  值为 0.0016648。此结果优于既往的血清葡萄糖近红外光谱分析效果<sup>[1,2]</sup>。由此说明,所采用的组合波段(4000-5300  $\text{cm}^{-1}$  和 5300-6087.5  $\text{cm}^{-1}$ )和校正集验证集的划分方法具有良好建模代表性和预测效果。

建立计算机算法平台,把全部 8 种  $\alpha$  平滑模式和不同  $\alpha$  因子数(0-5)组合分别建立  $\alpha$  模型,按照预测效果优选  $\alpha$  平滑模式和  $\alpha$  因子数。各阶导数平滑、各平滑点数的最优模型的  $\alpha_1:29$  值(从不同多项式模型、不同  $\alpha$  因子数中优选)如图 4 所示。各阶导数平滑(分开对应不同多项式)的最优模型的  $\alpha_1:29$  值、最优平滑点数、最优  $\alpha$  因子数如表 4 所示。未做  $\alpha$  平滑直接  $\alpha$  方法建模的结果也列在表 4 中。全局最优的  $\alpha$  平滑模式为 0 阶导数平滑、1 次多项式类型、0 平滑点数对应的最优因子数为 0,最优  $\alpha_1:29$  为 0.0016648,预测相关系数  $R_p$  为 0.999,预测效果明显优于未做  $\alpha$  平滑处理的结果。表 4 和图 4 表明,不同的导数平滑和不同的多项式次数类型对应的最优平滑点数、最优  $\alpha$  因子数一般是不相同的。不同的导数平滑、不同的多项式次数类型和采用不同的平滑点数,对应的最优  $\alpha_1:29$  值也是差别比较大的。如果根据既往文献或者其它研究对象所采用的平滑

表 3 校正集验证集  $K_9'$ 、化学值的均值和标准偏差  
\*:JB7 % 1 7;F ;FR 5=;FR;LR R7<+;=4F 4b : $K_9'$  ;FR MN76+=;B <;BH7 +F M;B+JL;=4F 57= ;FR 0L7R+=4F 57=

	化学值 TN76+=;B <;BH7		单波数预测偏差 : $K_9'$	
	均值 1 7;F	标准偏差 :=;FR;LR R7<+;=4F	均值 1 7;F	标准偏差 :=;FR;LR R7<+;=4F
校正集 T;B+JL;=4F 57=	0.0016648	0.080	0.0016648	0.080
验证集 9L7R+=4F 57=	0.0016648	0.080	0.0016648	0.080

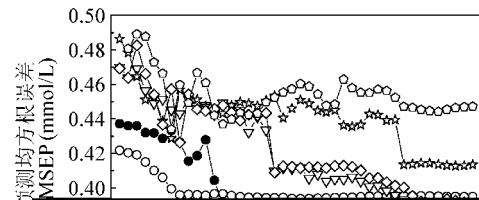


图 4 各阶导数的平滑点数对应的最优  $\alpha_1:29$   
—(0! , 0=6;B L4= 67;F 5cH;L7 7LL4L 4b 0L7R+=4F  
( $\alpha_1:29$ ) M4LL7504FR+F (= 4 5644-N 04+F=5 FH6J7L b4L 7;MN  
4LR7L R7L<;=4F  
(○) 零阶(0, LR7L); (●) 一阶(1, LR7L); (▽) 二阶(2, LR7L),  
(☆) 三阶(3, LR7L); (◇) 四阶(4, LR7L); (○) 五阶(5, LR7L)!

模式的经验, 不经过大范围比较筛选, 很难得到最优的 A 平滑模式和 9G: 因子数。另一方面, 从表 # 和图 ! 还可以看出, 最优平滑点数一般都不在 #8 以内, 如果采用 #8 以内的平滑点数, 就达不到现在的最优预测效果( #8 点以内最优 al : 29 为 \$0!D1 664B&G) , 这说明 : A 平滑模式的扩充是非常有必要的。

表 # 各阶导数平滑最优模型的预测效果

\*: JB7 # 9L7R:M=4F 7bb7M= 4b 40=6; B 64R7B M4LL7504FR+F( =4 7; MN 4LR7L R7L<; =4F

	多项式次数 94B@F46+F; B R7(L77	平滑点数 : 644=N+F( 04+F-5	9G: 因子数 9G: b; M=4L	预测均方根偏差 al : 29
未平滑 . 4 5644=N+F(	V	V	"	\$0C#!
零阶 \$ , LR7L	# , ! C , 8 1	!D CO OD	D D %0	\$0!DC \$0!DC \$0!D"
一阶 % , LR7L	# ! , C 8 , 1	!! !B! "%	0 0 "	\$0!OD \$0!O1 \$0!OD
二阶 # , LR7L	# , ! C , 8 1	"" "0 00	0 0 %\$	\$0! "C \$0C\$D \$0!D%
三阶 ! , LR7L	! , C 8 , 1	"0 "! "!	0 C C	\$0C%8 \$0C%! \$0C!0
四阶 C , LR7L	C , 8 1	%8 "0	1 D	\$0C!0 \$0!D\$
五阶 8 , LR7L	8 , 1	%0	0	\$0C!C

为了观察 9G: 因子数对模型效果的影响, 对不做 : A 平滑的直接 9G: 模型和最优 : A 平滑 9G: 模型 (% 阶导数平滑 !、C 次多项式类型 B! 平滑点数) 分别给出 9G: 因子数对应的 al : 29 , 如图 C 所示。直接 9G: 模型的最优因子数为 " , al : 29 为 \$0C#! 664B&G。最优 : A 平滑 9G: 模型的最优因子数为 0 , al : 29 为 \$0!O1 664B&G。最优因子数不相同, 平滑后预测效果提升较大, 在 : A 平滑 9G: 模型中因子数的影响更为显著。

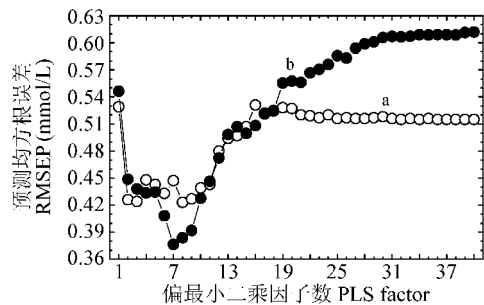


图 C 平滑前后的 9G: 因子数对应的 al : 29  
\_+(0C al : 29 M4LL7504FR+F( =4 9G: b; M=4L J7b4L7 ; FR ; b=7L 5644=N+F( ; 0 9G: 模型( 9G: 64R7B) ; J0 最优 : A 平滑 9G: 模型( , 0=6; B : A 5644=N+F( 9G: 64R7B) !

### 4 结 论

实验结果表明 : A 平滑模式扩充以及 : A 平滑模式和 9G: 因子数的联合全局筛选都是非常必要的。所提出的划分校正集验证集的方法 : A 平滑模式扩充 : A 平滑模式和 9G: 因子数的联合全局筛选能够有效地应用于近红外光谱分析的模型优化。

### References

% ' HLF5 - ' , T+HLM>; ? 2K0 Handbook of Near-infrared Analysis , #nd ed , . 7P Z4L?: I ; LM7B R7??7L +FM , 2001: 1!! e1C0  
# TX Y ]+; 4)G+( 褚小立) ]Y ZH)97F(( 许育鹏) ,GY K; F)WN7F( 陆婉珍) Q Chinese J. Anal. Chem. (分析化学) , 2008 , !1( 8) : 0\$# e0\$D  
! \; 5765H6L; F : , -H Z 9 , I ; LH4 \ , , >; ?+ Z0 Chemometrics Intell. Lab. Syst. , 2006 , "#( % )# : D0 e%\$!  
C TX2. XH4T);+( 陈华才) , Z' . A WN4F()AH4( 杨仲国) , TX2. ]+F()-; F( 陈星旦) Q Chinese J. Anal. Lab. (分析实验室) , 2005 , C( 0) : %0 e#\$  
8 ]Y XH+)a4F(( 徐惠荣) , K' . A XH+): N7F(( 汪辉君) , XY' . A \; F(( 黄 康) , Z[. A Z+) ' +F( 应义斌) , Z' . A TN7F( ( 杨 诚) , ^[' . X; 4( 钱 豪) , XY EHF( 胡 俊) Q Spectroscopy and Spectral Analysis( 光谱学与光谱分析) , 2008 , #"( % ) : #8#! e#8#1

1 TX2. ]H7)Z+F(( 陈雪英) ,G[ Z7)aH+( 李页瑞) ,TX2. Z4F(( 陈勇) ,K' . A G4F()XH( 王龙虎) ,--F( G+F(( 丁玲) 0  
*Chinese J. Anal. Chem.* (分析化学) ,**2009** ,!0(%\$): %C8% e%C81

0 ZY Z;F)' 4( 于燕波) ,W' . A 97F(( 臧鹏) ,\_Y ZH;F)XH;( 付元华) ,WX' . A GH)-;( 张录达) ,Z' . Z;F)GH( 严衍  
 禄) ,TX2. ' +F( 陈斌) 0 *Spectroscopy and Spectral Analysis* (光谱学与光谱分析) ,**2008** ,#" (0): %88C e%88"  
 " :;<+>?@ ' ,A4B;@ I E 20 *Anal. Chem.* ,**1964** ,!1("): %1#0 e%1!0

D TXY ]+;4)G+( 褚小立) ,ZY' . X4F()\_H( 袁洪福) ,GY K;F)WN7F( 陆婉珍) 0 *Prog. Chem.* (化学进展) ,**2004** ,%1(C):  
 8#" e8C#

\$\$ TX2. E+7)I 7+( 陈洁梅) ,9' . \*;4( 潘涛) ,TX2. ]+F()-;F( 陈星旦) 0 *Optics Preci. Eng.* (光学精密工程) ,**2006** ,  
 %C(%): %e0

%% T' , 9H( 曹璞) ,9' . \*;4( 潘涛) ,TX2. ]+F()-;F( 陈星旦) 0 *Optics Preci. Eng.* (光学精密工程) ,**2007** ,%8(%#):  
 %D8# e%D8"

## Joint Optimization of Savitzky-Golay Smoothing Models and Partial Least Squares Factors for Near-infrared Spectroscopic Analysis of Serum Glucos