

文章编号: 0253-2239(2006)01-0136-5

简化苹果糖度预测模型的近红外光谱预处理方法*

赵杰文¹ 张海东^{1,2} 刘木华^{1,3}

1 江苏大学生物与环境工程学院, 镇江 212013
2 云南农业大学工程技术学院, 昆明 650201
3 江西农业大学工学院, 南昌 330045

摘要: 采用正交信号校正法(OSC)和净分析物预处理法(NAP)分别对苹果的近红外光谱(1300~2100 nm)进行预处理,并结合偏最小二乘法(PLS)建立了糖度预测模型。应用结果显示,随着预处理过程中所用的正交信号校正因子或净分析物预处理因子的逐渐增加,偏最小二乘糖度模型(OSC/PLS模型和NAP/PLS模型)所采纳的最佳因子数也会随之减少,甚至可减至1。当采用10个正交信号校正因子预处理苹果光谱时,OSC/PLS糖度模型达到最佳性能,最佳模型采纳的因子数为2;采用11个净分析物预处理因子预处理光谱时,NAP/PLS糖度模型达到最佳性能,最佳模型采纳的因子数为1。从总体上评价,最佳OSC/PLS糖度模型和最佳NAP/PLS糖度模型的性能都明显优于原始光谱的最佳偏最小二乘模型。这些结果表明,正交信号校正法和净分析物预处理法都能在保证精度的同时有效地简化苹果糖度预测模型。

关键词: 近红外光谱; 模型简化; 正交信号校正; 净分析物预处理

中图分类号: TP242.62; S123; O657 文献标识码: A

Preprocessing Methods of Near-Infrared Spectra for Simplifying Prediction Model of Sugar Content of Apples

Zhao Jiewen¹ Zhang Haidong^{1,2} Liu Muhua^{1,3}

1 School of Biological and Environmental Engineering, Jiangsu University, Zhenjiang 212013
2 Faculty of Engineering and Technology, Yunnan Agricultural University, Kunming 650201
3 Engineering Collage, Jiangxi Agricultural University, Nanchang 330045

Abstract: Orthogonal signal correction (OSC) and net analyte preprocessing (NAP) were respectively used to pretreat the near infrared (NIR) spectra of apples ranging from 1300 nm to 2100 nm, and the models of sugar content were developed from these two pretreated spectra by the partial least square (OSC/PLS models and NAP/PLS models). Results showed that the optimal number of factors used in OSC/PLS models and NAP/PLS models would reduce as the number of OSC factors and NAP factors increased on by one, even to 1 finally. The best OSC/PLS model with 2 PLS factors was obtained when 10 OSC factors were used in pretreatment; and the best NAP/PLS model with 1 PLS factor was obtained when 11 NAP factors were used in pretreatment. Although the best OSC/PLS model and the best NAP/PLS model didn't improve precision to a great extent, they needed fewer factors and became more parsimonious, compared with the model before NAP pretreatment. In a conclusion, OSC and NAP pretreatment could simplify the prediction model of sugar content.

Key words: near-infrared spectrum; model simplification; orthogonal signal correction; net analyte preprocessing

1 引 言

近年来,近红外光谱检测技术因其快速、无损的优点而越来越多地被应用于生物体内部品质的检测

中^[1~3]。但由于近红外区的谱带复杂、重叠多,生物体光谱中对待测品质毫无价值的信息(包括来自生物体内部以及来自外部环境的干扰信息)必定会对

* 国家自然科学基金(30370813)和国家 863 计划(2002AA248051)资助课题。

作者简介: 赵杰文(1945~),男,江苏苏州人,教授,博士,主要从事农产品品质无损检测技术的研究。

E-mail: zhao@ujs.edu.cn

收稿日期: 2005-01-10; 收到修改稿日期: 2005-06-21

待测品质的近红外预测产生影响。为了获得满意的预测精度,通常都要先对原始光谱进行预处理,然后再利用多元统计分析技术建立预测模型。

在利用近红外光谱技术检测苹果糖度的研究中发现,用常用的光谱数据预处理方法如塞维兹-戈莱(Savitzky-Golay)卷积、多元散射校正及二阶求导等对光谱进行预处理后,虽然糖度模型的预测精度较高,但模型依然过于复杂(模型采纳的因子数过多)。为此,尝试采用两种能简化模型的光谱预处理算法——正交信号校正法(Orthogonal signal correction, OSC)和净分析物预处理法(Net analyte preprocessing, NAP)对苹果的近红外原始光谱进行预处理,希望能获得理想的简化模型。

2 算法简介

2.1 正交信号校正法

正交信号校正法由 Wold^[4] 首先提出,随后 Sjöblom^[5]、Andersson^[6]、Fearn^[7]、Westerhuis^[8] 及 Feudale^[10] 等先后对这种预处理法作了改进与发展。正交信号校正法的基本思想是:利用数学上空间正交的原理,将原始光谱矩阵中与待测品质不相关的部分信息特别是系统噪声滤除。本研究中采用 Fearn 的正交信号校正法对苹果光谱进行预处理^[7,10]:

设苹果校正集的近红外原始光谱数据矩阵为 \mathbf{X} ($I \times J$, I 为测试集样本数, J 为波长数),该矩阵中的某一元素 x_{ij} ($i=1,2,\dots,I; j=1,2,\dots,J$) 为第 i 个样本在第 j 个波长处的反射率,苹果糖度实测值向量为 \mathbf{y} ($I \times 1$)。

Fearn 算法的关键在于寻求一个 $J \times A$ 阶权重矩阵 \mathbf{W}_{OSC} ,该矩阵的每一列为权重向量 \mathbf{w}_l ($l=1,2,\dots,A$),也称为正交信号校正因子,与之相应的分量 t_l 应最大程度地反映苹果原始光谱 \mathbf{X} 中与糖度实测值向量 \mathbf{y} 正交的那部分变化,即 $\mathbf{y}^T t_l = 0$ (上标“T”表示矩阵或向量的转置)。权重向量 \mathbf{w}_l ($l=1,2,\dots,A$) 的求解过程是:

1) 构造一个与糖度向量 \mathbf{y} 正交的矩阵 $\mathbf{Z} = \mathbf{X}[\mathbf{I} - \mathbf{X}^T \mathbf{y}(\mathbf{y}^T \mathbf{X} \mathbf{X}^T \mathbf{y})^{-1} \mathbf{y}^T \mathbf{X}]$ (\mathbf{I} 为 $J \times J$ 阶单位矩阵);

2) 求出平方矩阵 $[\mathbf{Z}^T \mathbf{Z}]$ 的前 A 个特征向量即为权重向量 \mathbf{w}_l ($l=1,2,\dots,A$)。

求出 \mathbf{w}_l 后,根据 $t_l = \mathbf{X} \mathbf{w}_l$ 及 $\mathbf{p}_l = \mathbf{X}^T t_l (t_l^T t_l)^{-1}$ 得到 $J \times A$ 阶载荷矩阵 \mathbf{P}_{OSC} (该矩阵的每一列为载荷向量 \mathbf{p}_l)。然后将原始光谱矩阵 \mathbf{X} 向权重空间 \mathbf{W}_{OSC} 作正交投影就可得到经 A 个正交信号校正因子处

理后的光谱矩阵 $\mathbf{X}_{\text{O,A}}$,即 $\mathbf{X}_{\text{O,A}} = \mathbf{X}(\mathbf{I} - \mathbf{W}_{\text{OSC}} \mathbf{P}_{\text{OSC}}^T)$,此式中 \mathbf{I} 为 $J \times J$ 阶单位矩阵。

预测集苹果光谱 \mathbf{X}_{UN} 的正交信号校正按式 $\mathbf{X}_{\text{UN,O}} = \mathbf{X}_{\text{UN}}(\mathbf{I} - \mathbf{W}_{\text{OSC}} \mathbf{P}_{\text{OSC}}^T)$ 进行, $\mathbf{X}_{\text{UN,O}}$ 为经正交信号校正法预处理后的预测集苹果光谱。

2.2 净分析物预处理法

净分析物预处理法由 Goicoechea 等^[10] 首先提出,该法基于净分析物信号^[11] (NAS)理论,主要用于提取混合物光谱中某一纯组分的光谱信息^[10,12]。其基本思想是:利用数学上空间正交的原理,将原始光谱矩阵中待测组分的净分析物信号(NAS)提取出来。苹果光谱的净分析物预处理(NAP)算法如下^[10,13]:

同样设苹果校正集的原始光谱矩阵为 \mathbf{X} ($I \times J$),苹果糖度实测值向量为 \mathbf{y} ($I \times 1$)。由于苹果糖度是反映多种物质综合作用的一项指标,因此在运用净分析物预处理法时将苹果的近红外原始光谱矩阵 \mathbf{X} 分为两部分,其中一部分是与糖度相关的信息,而另一部分是与糖度不相关的所有干扰信息(包括来自苹果内部以及来自环境的干扰信息)的综合,即

$$\mathbf{X} = \mathbf{X}_{\text{SC}} + \mathbf{X}_{-\text{SC}}, \quad (1)$$

式中 \mathbf{X}_{SC} 表示苹果光谱中与糖度相关的信息, $\mathbf{X}_{-\text{SC}}$ 则表示光谱中糖度之外的所有其它干扰信息的综合。

寻求一个与 $\mathbf{X}_{-\text{SC}}$ 正交的 $J \times J$ 阶矩阵 \mathbf{F}_{NAP} (即 $\mathbf{X}_{-\text{SC}} \mathbf{F}_{\text{NAP}} = \mathbf{0}$),使(1)式两边同乘以 \mathbf{F}_{NAP} 后有 $\mathbf{X} \mathbf{F}_{\text{NAP}} = \mathbf{X}_{\text{SC}} \mathbf{F}_{\text{NAP}}$ 成立,这一步是该算法的关键步骤。矩阵 \mathbf{F}_{NAP} 的求解过程为:

1) 将原始光谱矩阵 \mathbf{X} 向糖度实测值向量 \mathbf{y} 作正交投影得到 $\mathbf{X}_{-\text{SC}} = [\mathbf{I} - \mathbf{y}(\mathbf{y}^T \mathbf{y})^{-1} \mathbf{y}^T] \mathbf{X}$ (式中 \mathbf{I} 为 $I \times I$ 阶单位矩阵);

2) 求出平方矩阵 $[(\mathbf{X}_{-\text{SC}})^T \mathbf{X}_{-\text{SC}}]$ 的特征向量矩阵 \mathbf{U} (\mathbf{U} 为 $J \times A$ 阶矩阵, \mathbf{U} 中的每一列为一个净分析物预处理因子);

3) 构造矩阵 $\mathbf{F}_{\text{NAP}} = \mathbf{I} - \mathbf{U} \mathbf{U}^T$ (式中 \mathbf{I} 为 $J \times J$ 阶单位矩阵)。

然后即可求出经 A 个净分析物预处理因子处理后的光谱 $\mathbf{X}_{\text{SC}}^* = \mathbf{X} \mathbf{F}_{\text{NAP}} = \mathbf{X}(\mathbf{I} - \mathbf{U} \mathbf{U}^T)$,式中 \mathbf{X}_{SC}^* 为经净分析物预处理法处理后得到的光谱矩阵,也即糖度的净分析物信号矩阵。

预测集苹果光谱 \mathbf{X}_{UN} 的净分析物预处理按式 $\mathbf{X}_{\text{UN,SC}}^* = \mathbf{X}_{\text{UN}}[\mathbf{I} - \mathbf{U} \mathbf{U}^T]$ 进行, $\mathbf{X}_{\text{UN,SC}}^*$ 为预测集苹果光谱中糖度的净分析物信号矩阵。

分析上面的算法步骤可以发现,两种算法都经过了两次正交投影的过程,第一次正交投影分别得到 \mathbf{Z} 和 $\mathbf{X}_{-\text{SC}}$,第二次则得到 $\mathbf{X}_{\text{O,A}}$ 和 \mathbf{X}_{SC}^* 。

3 实验方法与数据

实验选用市售陕西白水水晶富士 39 个,从中随机地选取 28 个作为校正集,余下的 11 个作为预测集,分别编号后置于 4 °C 冰柜中贮藏。光谱检测实验在环境温度可控的实验室(实验环境温度控制为 26 °C)内进行。实验前,将冰柜中取出的苹果置于试验室中 12 h,以使苹果整体温度达到与环境温度的一致;由近红外光谱仪(Nexus 670 FT-IR,美国 Nicolet 公司)在每个苹果的最大横径上等距离地选

取四个点进行光谱扫描,扫描波长范围为 1300~2100 nm,波长间隔为 0.5 nm(即波长数 $J=1600$);扫描时光纤探头与苹果果皮直接接触,并尽量避开表面缺陷处;取这四个点的平均光谱作为整个苹果的原始光谱。然后将该苹果削皮,取可食用部分榨汁,并用手持式糖度计(WYT0-32 型,泉州韦达计量仪器厂)测定其糖度值。图 1 为校正集 28 个苹果的原始光谱图。表 1 列出了所有 39 个被测苹果糖度实测值的变化范围、平均值及标准偏差。

表 1 苹果糖度实测值的变化范围、平均值和标准偏差

Table 1 Summary of the variation ranges, means and standard deviations of sugar content

Sample set	Number of samples	Variation range (° Brix)	Mean (° Brix)	Standard deviation
Calibration set	28	9.35~16.65	12.98	1.511
Prediction set	11	11.25~15.75	12.80	1.377

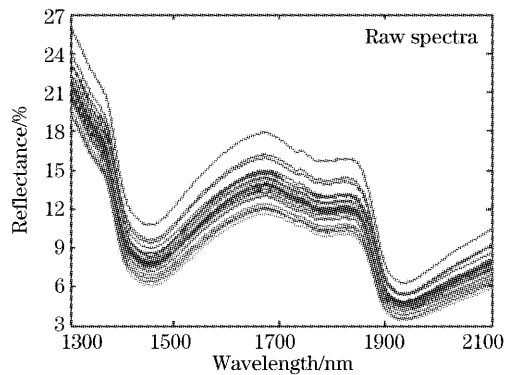


图 1 校正集苹果的近红外光谱

Fig. 1 Raw NIR spectra of apples in the calibration set

4 实验结果与讨论

4.1 正交信号校正法和净分析物预处理法预处理后的苹果光谱比较

分别采用正交信号校正法和净分析物预处理法对苹果的原始光谱进行预处理,图 2 分别为采用前 3

个正交信号校正因子和 3 个净分析物质处理因子处理后的校正集苹果光谱。由图可以看出,经正交信号校正法预处理后的光谱[图 2(a)]形状与原始光谱(见图 1)没有太大的区别,只是在排列上更为紧密。而经净分析物预处理法预处理后的光谱[图 2(b)]形状则发生了巨大的变化,光谱曲线变得更加粗糙。同时,在约 1400 nm、1740 nm 和 1880 nm 处出现了三个较为明显的峰(苹果的原始光谱是平缓、光滑的曲线)。这是因为正交信号校正法主要用于滤除原始光谱中的部分系统噪声(如光的散射及光程差异^[4]等),该法在去噪的过程中仍保留了光谱中的主要信息,因此处理后的光谱形状与原始光谱形状相比没有太大的区别。而净分析物预处理法则主要用于提取光谱中糖度的净分析物信号。在信号提取过程中,与糖度不相关的所有信息(包括来自苹果内部的其它成分的信息以及来自外部环境的干扰信息)被最大程度地从原始光谱中剔除,因此,光谱中只含有糖度信息和少量干扰信息,这就导致光谱形

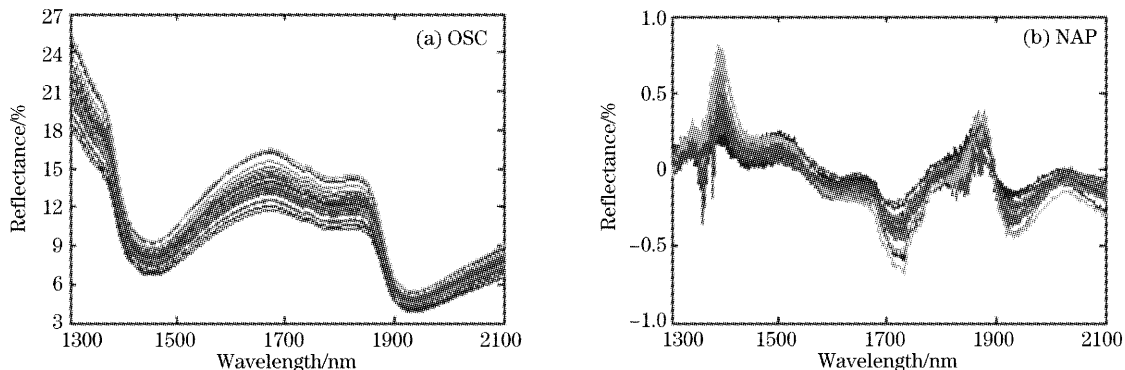


图 2 经正交信号校正法和净分析物预处理法预处理后的校正集苹果光谱

Fig. 2 Apple's spectra in the calibration set being pretreated by OSC (a) and NAP (b) respectively

状发生了变化,并出现了较为明显的峰。

4.2 偏最小二乘校正(PLS)结果比较

在采用偏最小二乘法(PLS)建立模型前,先按校正集和预测集分别将经过正交信号校正法和净分析物预处理法预处理的苹果光谱中心化。校正模型的最佳因子数(主成分数)由交互验证法(Cross-

Validation)确定,即由最小的预测残差平方和(PRESS,物理量用 s 表示)确定。表 2 和表 3 分别列出了经正交信号校正法和净分析物预处理法预处理前后,实验数据的偏最小二乘校正结果(表 2、表 3 中分别列出了采用 0~12 个两种因子预处理苹果光谱后的校正结果)。

表 2 正交信号校正法预处理前后的偏最小二乘法建模结果

Table 2 PLS results before and after apple spectra being pretreated by orthogonal signal correction (OSC)

OSC factors	PLS factors	Calibration set		Prediction set	
		Correlation coefficient (r^2)	Standard error of calibration (SEC)	Correlation coefficient (r^2)	Standard error of prediction (SEP)
0	11	0.92190	0.41473	0.86572	0.50473
1	10	0.92190	0.41473	0.86572	0.50473
2	9	0.92190	0.41473	0.86571	0.50473
3	8	0.92191	0.41471	0.86571	0.50475
4	7	0.92194	0.41464	0.86566	0.50482
5	6	0.92198	0.41453	0.86560	0.50495
6	5	0.92203	0.41438	0.86549	0.50515
7	4	0.92212	0.41414	0.86531	0.50550
8	3	0.92241	0.41337	0.86414	0.50767
9	3	0.92175	0.41513	0.86757	0.50122
10	2	0.92644	0.40250	0.86701	0.50229
11	2	0.92137	0.41615	0.85937	0.51653
12	1	0.92156	0.41562	0.86950	0.49757

表 3 净分析物预处理法预处理前后的偏最小二乘法建模结果

Table 3 PLS results before and after apple spectra being pretreated by net analyte preprocessing (NAP)

OSC factors	PLS factors	Calibration set		Prediction set	
		Correlation coefficient (r^2)	Standard error of calibration (SEC)	Correlation coefficient (r^2)	Standard error of prediction (SEP)
0	11	0.92190	0.41473	0.86572	0.50473
1	10	0.92192	0.41469	0.86564	0.50487
2	9	0.92195	0.41461	0.86576	0.50465
3	8	0.92203	0.41440	0.86597	0.50425
4	7	0.92206	0.41430	0.86608	0.50405
5	6	0.92209	0.41422	0.86602	0.50416
6	5	0.92216	0.41405	0.86589	0.50440
7	4	0.92229	0.41370	0.86553	0.50507
8	3	0.92306	0.41164	0.86531	0.50549
9	3	0.92160	0.41553	0.86803	0.50036
10	2	0.92605	0.40356	0.86818	0.50008
11	1	0.93089	0.39014	0.87260	0.49161
12	1	0.92820	0.39765	0.86255	0.51063

4.3 分析与讨论

对原始光谱进行预处理前,偏最小二乘(PLS)模型已能较好地预测苹果的糖度。其校正时的相关系数 r^2 和标准偏差分别为 0.92190 和 0.41473,预测时的相关系数 r^2 和标准偏差分别为 0.86572 和 0.50473。但该模型采纳的最佳因子数为 11,这使得模型显得过于复杂。对原始光谱分别进行预处理

后建立的偏最小二乘法糖度模型(相应称为 OSC/PLS 模型和 NAP/PLS 模型)所采纳的最佳因子数会随着预处理过程中所用正交信号校正因子和净分析物预处理因子数的逐个增加而逐渐减小,甚至可减少至 1。最初时,预处理过程中所用的因子每增加 1 个,模型所采纳的最佳因子就减少 1 个,采用 8 个因子分别对苹果原始光谱进行预处理后,两种模

型的最佳因子数都减到 3。采用 12 个正交信号校正因子和 11 个净分析物预处理因子时,两种模型的最佳因子数均减到了 1。此后即使采用更多的因子预处理光谱,模型的最佳因子个数都保持为 1(表中未列出)。之所以都能使糖度模型的最佳因子数减少,是因为这两种方法在预处理苹果原始光谱的过程中,随着正交信号校正因子或净分析物预处理因子的逐渐增加,与糖度不相关的干扰信息也逐渐减少的缘故。

从表 2 和表 3 中还可以看出,与原始光谱的偏最小二乘法糖度模型相比较,OSC/PLS 糖度模型和 NAP/PLS 糖度模型的精度都只有很小的波动。OSC/PLS 糖度模型在采用 10 个正交信号校正因子预处理光谱时达到最佳性能,最佳模型采纳的偏最小二乘因子数为 2,校正时的相关系数 r^2 和标准偏差分别为 0.92644 和 0.40250,用于预测时的相关系数 r^2 和标准偏差分别为 0.86701 和 0.50229。而 NAP/PLS 糖度模型在采用 11 个净分析物预处理因子预处理光谱时达到最佳性能,最佳模型采纳的偏最小二乘因子数为 1,校正时的相关系数和标准偏差分别为 0.93089 和 0.39014,用于预测时的相关系数 r^2 和标准偏差分别为 0.87260 和 0.49161。可以看出,与原始光谱的糖度模型相比,最佳 OSC/PLS 糖度模型和最佳 NAP/PLS 糖度模型都显得更加简洁,因此从总体上评价,性能都明显优于原始光谱的糖度模型。而前两者中,最佳 NAP/PLS 模型又略优于最佳 OSC/PLS 模型。图 3 为以最小预测残差平方和值确定模型最佳因子数的示意图[为更易看清,图中采用 $\lg(s)$ 值],图中下部自左至右的“☆”、“○”和“□”处的凹点分别对应着最佳 NAP/PLS 模型、最佳 OSC/PLS 模型和原始光谱偏最小二乘法模型的因子

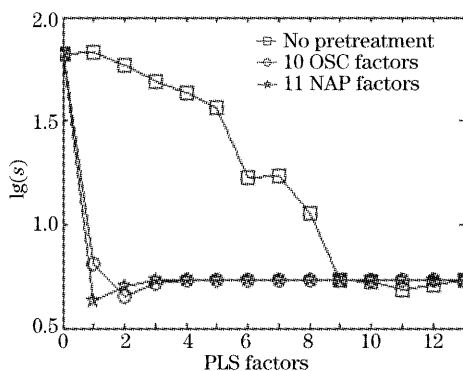


图 3 $\lg(s)$ 值确定最佳 PLS 因子数示意图

Fig. 3 Diagram of $\lg(s)$ value as a function determining the optimal number of PLS model

数 1、2、11 及 $\lg(s)$ 值 0.62961、0.65669、0.68270,因此从图 3 中也可以得出上述结论。

5 结 论

采用正交信号校正法和净分析物预处理法对苹果的近红外进行了预处理,并利用偏最小二乘法(PLS)分别建立了苹果糖度的预测模型,结果表明,两种预处理方法都能在保证精度的前提下有效地简化苹果糖度的预测模型。

参 考 文 献

- Dong Xiaopeng, H. O. Edwards, J. P. Dakin. Detection of methane gas with fibre correlation spectroscopy, using semiconductor laser diode pumped Tm-doped fibre source[J]. *Chin. J. Laser*, 1994, **A21**(10): 789~794 (in Chinese)
- 董小鹏, H. O. Edwards, J. P. Dakin. 采用 Tm 光纤光源的甲烷气体相关光谱检测[J]. *中国激光*, 1994, **A21**(10): 789~794
- Li Hui, Xie Shusen, Lu Zukang *et al.*. A new model of the light scattering in biological tissue for visible and near infrared region[J]. *Acta Optica Sinica*, 1999, **19**(12): 1661~1666 (in Chinese)
- 李 晖, 谢树森, 陆祖康 等. 生物组织的可见光与近红外光散射模型[J]. *光学学报*, 1999, **19**(12): 1661~1666
- Xiao Lifeng, Hu Hongzhang, Zhang Mei *et al.*. A near-infrared spectrometer based on an integrated optical AOTF[J]. *Chin. J. Lasers*, 2004, **31**(3): 269~272 (in Chinese)
- 肖立峰, 胡鸿璋, 张 梅 等. 一种基于集成光学声光可调谐滤波器的近红外光谱仪[J]. *中国激光*, 2004, **31**(3): 269~272
- Svante Wold, Henrik Antti, Fredrik Lindgren *et al.*. Orthogonal signal correction of near-infrared spectra[J]. *Chemometrics and Intelligent Laboratory Systems*, 1998, **44**: 175~185
- Jonas Sjöblom, Olof Svensson, Mats Josefson *et al.*. An evaluation of orthogonal signal correction applied to calibration transfer of near infrared spectra [J]. *Chemometrics and Intelligent Laboratory Systems*, 1998, **44**: 229~244
- Claus A. Andersson. Direct orthogonalization[J]. *Chemometrics and Intelligent Laboratory Systems*, 1999, **47**: 51~63
- Tom Fearn. On orthogonal signal correction[J]. *Chemometrics and Intelligent Laboratory Systems*, 2000, **50**: 47~52
- Johan A. Westerhuis, Sijmen de Jong, Age K. Smilde. Direct orthogonal signal correction[J]. *Chemometrics and Intelligent Laboratory Systems*, 2001, **56**: 13~25
- Robert N. Feudale, Huwei Tan, Steven D. Brown. Piecewise orthogonal signal correction[J]. *Chemometrics and Intelligent Laboratory Systems*, 2002, **63**: 129~138
- Héctor C. Goicoechea, Alejandro C. Olivieri. A comparison of orthogonal signal correction and net analyte preprocessing methods. Theoretical and experimental study[J]. *Chemometrics and Intelligent Laboratory Systems*, 2001, **56**: 73~81
- Avraham Lorber. Error propagation and figures of merit for quantification by solving matrix equations [J]. *Analytical Chemistry*, 1986, **58**(6): 1167~1172
- A. Muñoz de la Peña, A. Espinosa-Mansilla, M. I. Acedo Valenzuela *et al.*. Comparative study of net analyte signal-based methods and partial least squares for the simultaneous determination of amoxicillin and clavulanic acid by stopped-flow kinetic analysis[J]. *Analytica Chimica Acta*, 2002, **463**(1): 75~88
- Avraham Lorber, Klaas Faber, Brace R. Kowalski. Net analyte signal calculation in multivariate calibration [J]. *Analytical Chemistry*, 1997, **69**(8): 1620~1626