

基于 Lasso 方法的碳钢土壤腐蚀率预报研究

鲁庆 穆志纯*

(北京科技大学 自动化学院 北京 100083)

摘要 提出了一种基于 Lasso 的 step adaptive Lasso with Pls (SALP) 方法,采用 Bayesian Bootstrap 算法重构样本,通过多模型集成对变量进行预选,以消除数据扰动和离群点对于模型性能的影响。应用偏最小二乘权重系数改善 Lasso 方法在处理小样本数据时的参数估计不准确问题。采用碳钢在土壤中的腐蚀数据为研究对象,建立了腐蚀率模型。实验证明: SALP 方法建立的模型可以准确地拟合和预测土壤中的碳钢腐蚀率变化。该方法适用于自然环境中材料腐蚀率的预测及类似研究领域。

关键词 Lasso 腐蚀率 碳钢 自然环境腐蚀 土壤

中图分类号 TP391.9;

文献标志码 A

材料在自然环境中的腐蚀率预报是腐蚀研究的重点和难点之一。由于自然环境的材料腐蚀实验周期漫长、实验站点分布广,而且自然环境的腐蚀影响因素复杂多变,实验数据往往存在高维度、小样本和数据质量较低等问题。目前腐蚀率预报常用的方法包括多元回归^[1]、神经网络^[2]、灰色系统^[3]等。这些方法虽然各具优势,也存在一些问题,例如多元回归方法需要较大的样本容量,很难处理小样本数据;神经网络方法可以较好的拟合实验数据,但在解释腐蚀影响因素方面却存在局限性;灰色系统处理高维数据时需要复杂的数据预处理过程等。

Lasso 方法最早由 Tishirani^[4]提出,能够在高维变量空间中获取稀疏线性模型,目前广泛应用于图像识别^[5]、基因分析^[6]等领域。它将变量选择和参数估计融合在模型训练过程中,通过 L1 罚函数压缩模型系数,系数被压缩为 0 的变量被视为无关变量,只有显著变量被保留在最终模型中,即在完成模型系数估计的同时也实现了变量选择。

本文将 Lasso 方法引入材料在自然环境中的腐蚀率预报研究中。提出了一种基于 Lasso 的 SALP 方法,以实际项目中的碳钢土壤腐蚀数据检验 SALP 方法的性能,并与采用 Adaptive Lasso、神经网络以及支撑向量回归等方法训练的模型进行比较,实验结果表明, SALP 方法能够较好地处理高维、小样本

数据建模问题,性能优于其他方法。建立的模型可以准确的拟合和预测碳钢在土壤的腐蚀率,并具有较好的可解释性。

1 Lasso 方法

Lasso 方法是一种基于 L1 罚的正则化方法,以下简述其原理。

对于任意样本容量为 n 的数据集 $\{X, Y\}_1^n$, 其中 X 为 p 维预报变量集, $X = (x_{1j}, \dots, x_{nj})^T; j = 1, \dots, p$; Y 为响应变量, $Y = (y_1, \dots, y_n)^T$; 假定线性模型可以合理的描述变量 X 与 Y 的关系,模型的形式如式(1)所示。

$$Y = \sum_{j=1}^p x_j \beta_j + \varepsilon \quad (1)$$

式(1)中 β_j 为模型中待估计的回归系数, ε 是随机误差项,代表了不能用变量 X 解释的因素对于 Y 的影响。

对于回归系数 β 的估计,常用的方法是最小二乘方法:令 $\hat{\beta} = (X^T X)^{-1} X^T Y$, 这时模型残差的平方和 $\sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij} \beta_j)^2$ 为最小值。

但这种方法要求样本容量 n 和变量集维数 p 至少要满足 $n \geq p + 1$ 的关系,对于高维数据和小样本数据,最小二乘方法并不适用。

Lasso 方法的基本原理是对模型回归系数的绝对值之和 $\sum_{j=1}^p |\beta_j|$ 进行惩罚,系数 β 的估值如式(2)所示。

$$\beta_{(\text{Lasso})} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij} \beta_j)^2 \right\} \times \sum_{j=1}^p |\beta_j| \leq t \quad (2)$$

2014年7月30日收到

国家科技基础性工作专项基金项目
(2012FY113000) 资助

第一作者简介:鲁庆(1975—),男,博士研究生。研究方向:数据挖掘、系统建模。E-mail: luqing_qhd@hotmail.com。

* 通信作者简介:穆志纯,教授,博士生导师。研究方向:复杂系统的建模与智能控制、信息自动化。E-mail: mu@ies.ustb.edu.cn。

数。Adaptive Lasso 方法的变量系数估计值如式 (3)。

$$\beta_{(\text{Adaptive-Lasso})} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N (y_i - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p \omega_j |\beta_j| \right\} \quad (3)$$

式(3)中权重系数 $\omega_j = |\beta_{\text{ols}}|^{-\gamma}$, $\gamma > 0$ 。 β_{ols} 为最小二乘回归系数, 当变量存在多重共线性问题, Zou 建议采用岭回归^[14]系数替代 β_{ols} 。

将 Adaptive Lasso 应用到 SALP 方法中, 但对于小样本数据, 使用最小二乘回归或岭回归权重系数的模型预报偏差较大, 为寻求改善, 采用偏最小二乘算法 (partial least squares, PLS) 系数 β_{pls} 替代 β_{ols} 。偏最小二乘算法由 S. Wold 等人提出^[15], 它可以在样本个数少于变量个数的条件下进行回归建模。在处理样本容量小、自变量多、变量间存在严重相关性问题的方面有独特的优势。

SALP 方法的基本流程描述如下:

Input: 数据集 $\{X, Y\}_n^1$, 重构样本数量 k ;
 Output: 模型中待估计的回归系数 β ;
 Step1: 按照 Lasso 方法的要求进行数据预处理, 令处理后的数据集满足: $\sum_{i=1}^n y_i = 0$, $\sum_{i=1}^n x_{ij} = 0$, $\sum_{i=1}^n x_{ij}^2 = 1$;
 Step2: 使用 Bayesian Bootstrap 算法重构数据集, 并生成训练集 $X(l)$ $l=1, \dots, k$;
 Step3: 定义变量 $\text{Vote}(j) = 0$, 用于记录变量 x_j 被模型选择的次数;
 For $l = 1$ to k
 以 $X(l)$ 为训练集, 使用 Adaptive Lasso 方法训练模型 $M(l)$;
 For $j = 1$ to p //统计变量入选模型次数
 If 变量 x_j 被模型 $M(l)$ 选中
 $\text{Vote}(j) = \text{Vote}(j) + 1$
 End If
 End For
 Step4: 根据变量入选模型次数统计结果对预报变量进行预选, 生成新的数据集 $\{X^*, Y\}_n^1$, X^* 是 X 的子集;
 Step5: 在 $\{X^*, Y\}_n^1$ 中抽取训练集, 使用基于 PLS 权重系数的 Adaptive Lasso 方法训练模型, 输出结果。

3 实验结果与讨论

3.1 数据说明

以碳钢在土壤中的腐蚀数据为研究对象, 数据的样本容量为 70 个, 包含了 104 个变量。其中变量 x_1 至 x_{13} 为碳钢埋件实验测量记录, 变量 y 是碳钢在土壤的腐蚀率记录。 x_{14} 至 x_{104} 为检测预报变量与响应变量的非线性关系以及预报变量的交互效应是否存在而设置的虚拟变量。其中变量 x_{14} 至 x_{26} 是原变

量的二次项集合, 即 $\{x_1^2, \dots, x_i^2, \dots, x_{13}^2\}$; 变量 x_{27} 至 x_{104} 是原变量的交互项集合, 即 $\{x_1 \times x_2, \dots, x_{(i-1)} \times x_i, \dots, x_{12} \times x_{13}\}$ 。研究的目的是训练模型来拟合、预报腐蚀率 y , 并考察其影响因素。

表 1 报告了实验数据的变量名称、变量定义和变量的均值、标准差。

表 1 碳钢的土壤腐蚀实验数据
Table 1 Experimental data of carbon steel corrosion in soil

变量名称	定义	单位	均值	标准差
x_1	碳钢埋试时间	a	7.37	10.20
x_2	土壤电阻值	欧姆	141.39	199.28
x_3	含水量	%	20.54	7.05
x_4	土壤 PH 值	—	7.60	1.41
x_5	有机质含量	%	0.67	0.53
x_6	氮含量	%	0.04	0.02
x_7	HCO_3^- 离子含量	%	0.02	0.02
x_8	Cl^- 离子含量	%	0.07	0.20
x_9	SO_4^{2-} 离子含量	%	0.15	0.33
x_{10}	Ca^{2+} 离子含量	%	0.05	0.09
x_{11}	Mg^{+} 离子含量	%	0.01	0.01
x_{12}	K^{+} 离子含量	%	0.00	0.01
x_{13}	Na^{+} 离子含量	%	0.06	0.14
$x_{14} \sim x_{26}$	$x_1 \sim x_{13}$ 的二次项	—	—	—
$x_{27} \sim x_{104}$	$x_1 \sim x_{13}$ 的交互项	—	—	—
y	腐蚀率	mm/a	3.20	1.79

3.2 建立模型

由于碳钢的土壤腐蚀数据采集周期长, 过程复杂, 其中可能混杂离群点、错误记录等。为消除数据扰动的影响, 首先进行 SALP 方法的第一步, 对变量进行预选。

按照算法要求, 将实验数据中预报变量集 X 的 104 个变量标准化处理, 响应变量 Y 对中处理。然后采用 Bayesian Bootstrap 算法对预处理后的数据进行抽样, 重构 500 个数据集。对于每个重构的数据集, 采用基于 PLS 权重系数的 Adaptive Lasso 方法训练模型。

由于变量入选模型的次数分布不属于正态分布、卡方分布等常见分布, 考察其经验分布。次数分布的最小值、10% 分位数、25% 分位数、中位数和最大值分别为 0、11、34、75、432。从分位数可以看出, 有 10% 的变量被模型选中的次数少于 11 次, 而 25% 的变量被模型选中的次数少于 34 次。

为确定变量预选门限, 将被模型选择次数少于 11 次、34 次、75 次的变量分别移除, 形成 3 个新的数据集。每个新数据集的预报变量是原预报变量集合的子集, 样本容量为 70 个。新数据集记为数据集 1、数据集 2 和数据集 3。对于每个新数据集, 随机

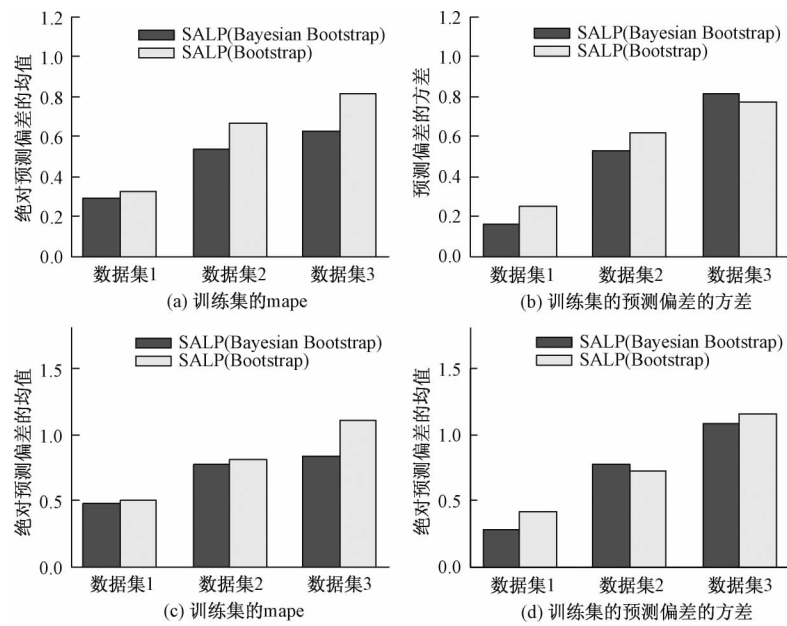


图 2 SALP(with Bayesian Bootstrap) 和 SALP(with Bootstrap) 方法在不同数据集的性能对比
Fig. 2 The performance comparison of SALP(with Bayesian Bootstrap) and SALP
(with Bootstrap) methods on different training set

抽取 60 个样本生成训练集 ,其余 10 个样本作为检验集。

应用基于 Pls 权重系数的 Adaptive Lasso 方法在 3 个新训练集上分别训练模型以确定变量预选门限。作为对比 ,本文还采用 Bootstrap 算法重构样本并重复了上述过程。实验结果如图 2 所示 ,对应的数据记录在表 2 中。

本文采用二个指标检测模型的性能 ,一是绝对预测偏差的均值 *mape*(mean of absolute prediction error) ,定义为 $mape = \frac{1}{m} \sum_i |y_i - y_i^*|$,其中 m 为数据集的样本容量 $y_i - y_i^*$ 是实际值和模型预报值的偏差;二是预测偏差的方差。

图 2(a) 表示基于不同抽样算法的 SALP 方法在 3 个训练集上的 *mape* 对比;图 2(b) 表示了对应的预测偏差的方差对比;图 2(c) ~ (d) 则分别表示了对应指标在测试集上的表现。图 2 显示 ,采用 Bayesian Bootstrap 算法抽样的模型性能略优于采用 Bootstrap 算法的对应模型。这可能是由于 Bootstrap 的应用在取决于经验分布的选取和样本数的大小^[16]。当样本量小不足以提供总体分布信息时 ,其结果并不可靠。

从图 2 和表 2 的记录可以发现 ,如果在变量预选过程中移除的变量过多 ,如数据集 2 和数据集 3 ,模型的预报性能会变得较差。而适当的变量预选 ,如数据集 1 ,却可以明显的提高模型性能。因此可以认为在变量预选过程中 ,应该以减少无关变量入

选模型可能性为目标 ,如果移除变量过多可能造成最终模型漏掉重要变量或过于稀疏 ,进而影响其预报性能。

表 2 不同变量选择门限的模型性能对比
Table 2 Model performance Comparison with
different variable selection strategy

数据		指标	算法	
			SALP(with Bayesian Bootstrap)	SALP(with Bootstrap)
数据集 1	训练集	<i>mape</i>	0. 288	0. 328
		偏差方差	0. 160	0. 249
	检验集	<i>mape</i>	0. 483	0. 501
		偏差方差	0. 288	0. 414
数据集 2	训练集	<i>mape</i>	0. 534	0. 655
		偏差方差	0. 531	0. 617
	检验集	<i>mape</i>	0. 771	0. 810
		偏差方差	0. 776	0. 723
数据集 3	训练集	<i>mape</i>	0. 629	0. 817
		偏差方差	0. 813	0. 776
	检验集	<i>mape</i>	0. 934	1. 018
		偏差方差	1. 076	1. 151

根据实验结果 ,将变量被模型选择次数的门限定为 11 次。

在完成变量的预选后 ,进行 SLAP 方法的第二步:以数据集 1 作为变量预选结果 ,随机抽取 60 个样本组成训练集 ,其余 10 个样本作为检验集。采用基于 Pls 权重系数的 Adaptive Lasso 方法训练模型。

3.3 模型性能检验与讨论

作为对比,在原数据集上选择与 SALP 方法第二步相同的样本组成训练集和检验集,训练了基于岭回归权重系数的 Adaptive Lasso 模型、支持向量回归(SVR)模型和神经网络模型。这里的数据集包含全部的变量。其中 Adaptive Lasso 方法的调控参数 λ 和 γ 通过 10 折交叉认证选取;SVR 模型使用 R 软件包 E1071 计算;神经网络模型使用 R 软件包 NNET 实现。后两者的参数主要采用了默认值。

在检验集上测试各模型的结果如图 3 所示。

从图 3 可以看出 SALP 方法所建立的模型可以较好的预测碳钢在自然环境土壤中的腐蚀率,尽管模型中包含较多的变量,但没有出现过拟合情况。在处理与实现对象数据类似的高维小样本数据时, SALP 方法的性能优于未经特殊优化的 Adaptive Lasso、SVR 和神经网络等方法,而且由于最终输出为线性模型,具有较好的可解释性。

在验证算法的有效性后,为了充分利用数据中的信息,采用全部 70 个样本重新训练模型。最终模型形式如式(1),共包含 35 个变量。回归系数 β 的估值记录在表 3 中。

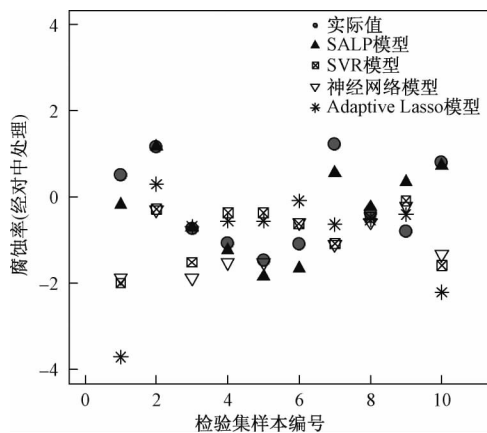


图3 SALP 模型与其他模型在测试集的预测结果对比

Fig. 3 The prediction comparison of SALP model with other model based on test set

SALP 方法对预报变量进行了标准化处理,因此表 3 的模型系数 β 为标准化系数,可以通过系数的绝对值排序来判断预报变量对于响应变量的重要性^[17]。对本文的使用数据而言,碳钢在土壤中的埋试时间、土壤中 Ca^{2+} 离子含量和 Mg^{+} 离子含量的交互项、 K^{+} 离子含量、 Mg^{+} 离子含量的二次项、埋试时间和土壤 pH 值的交互项等是碳钢在土壤中腐蚀率的主要影响因素。

由于模型为线性模型,通过控制其他预报变量,可以考察某个变量对于腐蚀率的净效应。例如根据表 3 中的记录,碳钢在土壤中的埋试时间每增加 1

个标准差,腐蚀率就会下降 309.1 个标准差;土壤含水量和 CL^{-} 离子含量的交互项系数为 214.8,表现出了变量间可能存在较强的相互促进作用;与之相反,土壤氮含量和 Ca^{2+} 离子含量的交互项系数为 -143.08,说明二个变量可能对腐蚀率的作用存在相互削弱的关系等等。当然,任何数学模型都不可能精确地描述现实世界和自然^[18],SALP 模型的可解释性能够为腐蚀因素分析提供线索,但模型产生的结果是否符合现实规律,仍然需要通过实践进行验证。

表 3 碳钢的腐蚀率模型系数估值

Table 3 Coefficients of carbon steel corrosion rate model

变量描述	系数 β	变量描述	系数 β
埋试时间	-309.31	电阻 $\times \text{CL}^{-}$	46.96
含水量	-24.47	电阻 $\times \text{Ca}^{2+}$	-156.17
HCO_3^{-}	70.24	电阻 $\times \text{K}^{+}$	69.16
K^{+}	-266.38	含水 $\times \text{PH}$	41.34
(埋试时间) ²	95.88	含水量 \times 有机质	150.00
(电阻) ²	-29.08	含水量 \times 氮含量	-124.67
(pH) ²	-22.40	含水量 $\times \text{CL}^{-}$	214.48
(HCO_3^{-}) ²	-116.16	PH \times 有机质	-124.15
(SO_4^{2-}) ²	-70.73	PH \times 氮含量	135.24
(Mg^{+}) ²	-279.11	有机质 \times 氮含量	-29.92
埋试时间 \times 电阻	24.00	有机质 $\times \text{HCO}_3^{-}$	22.80
埋试时间 $\times \text{PH}$	217.45	有机质 $\times \text{SO}_4^{2-}$	198.11
埋试时间 $\times \text{HCO}_3^{-}$	-70.61	有机质 $\times \text{K}^{+}$	-49.03
埋试时间 $\times \text{Mg}^{+}$	94.11	氮含量 $\times \text{Ca}^{2+}$	-143.08
电阻 $\times \text{PH}$	37.48	氮含量 $\times \text{Mg}^{+}$	-113.91
电阻 \times 氮含量	21.78	氮含量 $\times \text{K}^{+}$	153.83
$\text{HCO}_3^{-} \times \text{K}^{+}$	108.44	氮含量 $\times \text{Na}^{+}$	-50.00
$\text{Ca}^{2+} \times \text{Mg}^{+}$	299.98		

4 结束语

提出了一种基于 Lasso 的 SALP 方法,改善了 Lasso 方法求解过程中对数据扰动敏感的问题和处理小样本数据时的参数估计问题。采用碳钢在自然环境土壤中的腐蚀数据进行方法检验,实验结果表明,SALP 方法建立的模型能够准确的预测碳钢在土壤中的腐蚀率,具有较好的解释性。该方法适用于自然环境中材料腐蚀率的预测及类似研究领域。

参 考 文 献

- 邓永先,吕国志,张有宏,等. 对结构腐蚀增长模型的探讨. 科学技术与工程, 2007; 7(19): 4956—4960
Deng Yongxian, Lü Guozhi, Zhang Youhong, et al. Discussion of structure corrosion growth. Science Technology and Engineering, 2007; 7(19): 4956—4960
- Sadowski L. Non-destructive investigation of corrosion current density in steel reinforced concrete by artificial neural networks. Archives of Civil and Mechanical Engineering 2013; 13(1): 104—111

- 3 Liang Ping , Du Cuwei , Li Xiaogang. Grey relational space analysis of effect of environmental factors on corrosion resistance of x70 pipeline steel in ying tan soil simulated solution. *Corrosion & Protection* , 2009; (4) : 23—27
- 4 Tibshirani R. The lasso method for variable selection in the Cox model. *Statistics in Medicine* , 1997 , 16(4) : 385—395
- 5 杨 关, 张向东, 冯国灿 等. 图模型在彩色纹理分类中的应用. *计算机科学* , 2011; 38(10) : 273—277
Yang Guan , Hang Xiangdong , Feng Guocan , *et al.* Applications of graphical models in color texture classification. *Journal of Frontiers of Computer Science & Technology* , 2011; 38(10) : 273—277
- 6 Tibshirani R , Saunders M , Rosset S , *et al.* Sparsity and smoothness via the fused lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* , 2005; 67(1) : 91—108
- 7 Hastie T , Tibshirani R , Friedman J H. 统计学习基础—数据挖掘、推理与预测. 范 明, 柴玉梅, 译. 北京: 电子工业出版社, 2004: 337—384
Hastie T , Tibshirani R , Friedman JH. The elements of statistical learning: data mining , inference , and prediction. Tr. by Fan Ming , Cgao Yumei. Beijing: Electronic Industry Press , 2004: 337—384
- 8 Osborne M R , Presnell B , Turlach B A. On the Lasso and its dual. *Journal of Computational and Graphical Statistics* , 2000; 9 (2) : 319—337
- 9 Efron B , Hastie T , Johnstone I , *et al.* Least angle regression. *The Annals of Statistics* , 2004; 32(2) : 407—499
- 10 Khan J A , van Aelst S , Zamar R H. Robust linear model selection based on least angle regression. *Journal of the American Statistical Association* , 2007; 102(480) : 1289—1299
- 11 Rubin D B. The bayesian bootstrap. *The Annals of Statistics* , 1981; 9(1) : 130—134
- 12 Clyde M A , Lee H K H. Bagging and the bayesian bootstrap. *Artificial Intelligence and Statistics. Key West* : 2001: 169—174
- 13 Zou H. The adaptive lasso and its oracle properties. *Journal of the American Statistical Association* , 2006; 101(476) : 1418—429
- 14 Hoerl A E , Kennard R W. Ridge regression: applications to non-orthogonal problems. *Technometrics* , 1970; 12(1) : 69—82
- 15 Wold S , Ruhe A , Wold H. The collinearity problem in linear regression , The partial least square (PLS) approach to generalized inverses. *SIAM J Science Statistics Computer* , 1984; 5(2) : 735—743
- 16 Roff D A. Introduction to computer-intensive methods of data analysis in biology. New York: Cambridge University Press , 2006: 126—127
- 17 王海燕 杨方廷 刘 鲁. 标准化系数与偏相关系数的比较与应用. *数量经济技术经济研究* , 2006; (9) : 150—155
Wang Hai-yan , Yang Fang-ting , Loi Lu. Comparison and application of standardized regressive coefficient & partial correlation coefficient. *The Journal of Quantitative & Technical Economics* , 2006; (9) : 150—155
- 18 吴喜之, 马景义, 吕晓玲 等. 数据挖掘前沿问题. 北京: 中国统计出版社 2009: 7—13
Wu Xizhi , Ma Jingyi , Lv Xiaoling , *et al.* The frontier issues of data mining. Beijing: China Statistics Press , 2009: 7—13.

Corrosion Rate Prediction of Carbon Steel in Soil based on Lasso Method

LU Qing , MU Zhi-chun*

(School of Information Engineering , University of Science and Technology Beijing , Beijing 100083 , P. R. China)

[Abstract] A Lasso-based method was proposed named as Step Adaptive Lasso with Pls (SALP) . To eliminate the data perturbation or outlier influence on model training , SALP method uses Bayesian Bootstrap algorithm to reconstruct data sets and establishes models based on them. Then the predictors can be prescreened through integrating the result of those models. In order to deal with small sample size problem of data set , the partial least squares weighting factor is applied. By using SALP method , this study established a corrosion rate model based on the corrosion data of carbon steel in soil. The result shows that , the model established by SALP could precisely describe and predict the corrosion rate of carbon steel in soil. The method can be applied to prediction of material corrosion rate in natural environment or similar research field.

[Key words] Lasso corrosion rate carbon steel corrosion of the natural environment soil