

基于粗集理论的气测录井数据归一化处理

刘华莹,孙春雨,张方舟,邱露露,高晓松

(东北石油大学 计算机与信息技术学院,黑龙江 大庆 163318)

摘要:在现代录井工程中,地质情况的不同、钻井工艺的优差等诸多原因都会影响气测录井资料的录取。即使在同一地区、同一层位进行钻井勘探,测量的气测资料结果也会存在很大差异。可想而知如果在不同的地区、不同的层位进行勘探会导致测量气测资料变得更加困难与复杂。因此,快速、有效的规范气测资料及参数选择是现代录井工艺与油气层识别技术中至关重要的步骤。针对 RBF 神经网络算法具有收敛速度慢且不稳定等缺点,无法有效处理气测录井资料,文中提出了一种基于粗集理论的归一化处理方法,利用粗集理论对气测样本数据归一化处理后提高了 RBF 神经网络的训练速度。为了验证方法的可行性,以辽河油田的气测录井数据为背景进行仿真计算,实验结果表明此方法有效地提高了 RBF 神经网络处理气测录井数据速度。

关键词:数据归一化;气测录井;RBF 神经网络;粗集理论

中图分类号:TP39

文献标识码:A

文章编号:1673-629X(2015)07-0189-04

doi:10.3969/j.issn.1673-629X.2015.07.042

Gas Logging Data Normalization Processing Based on Rough Set Theory

LIU Hua-ying, SUN Chun-yu, ZHANG Fang-zhou, QIU Lu-lu, GAO Xiao-song

(College of Computer and Information Technology, Northeast Petroleum University,
Daqing 163318, China)

Abstract: In the modern mud logging engineering, different geological conditions, the good and poor of drilling process, and other reasons will have effects on gas logging data. Drilling exploration, even in the same area, the same horizon, makes the measurement results of gas logging data different highly. It can imagine if in different layers of different area, exploration will result in measuring gas data becomes more difficult and complex. Therefore, fast, effective and standardized gas logging data and parameter selection is essential to modern mud logging technology and oil and gas reservoir recognition technology. In view of the disadvantage of slow convergence and instability of RBF neural networks algorithm, unable to handle the gas effective logging data, present a normalization method based on rough set theory, using rough set theory on gas logging normalized sample data to improve RBF neural network training speed. In order to verify the feasibility of the method, taking the Liaohe oil field gas logging data calculation as the background of the simulation, the experimental results show that this method effectively improves gas logging data speed RBF neural network dealt with.

Key words: data normalization; gas logging; RBF neural network; genetic algorithm; rough set theory

0 引言

在油气勘探与开发工程中,气测录井技术是最基本的方法,在钻井勘探过程中气测录井技术可直接测量出钻井液变化情况与井眼中气体含量与种类情况,对快速与正确地识别油气藏、评价油气藏和合理地优化钻井工艺等方面有重要意义。在国内地区地质情况复杂,油气勘探目标由简单到复杂导致勘探难度日益增大,以及对油气藏控制因素多元性认识的深入^[1-4]

等诸多因素,为了有效勘探复杂的地质情况,气测录井工艺伴随着国内科技的进步与发展,逐渐成为一项至关重要的新一代录井技术。气测录井工艺主要是测量与处理钻井勘探中的相关工程参数,其中包含油、气、水显示信息与钻井液信息等。气测录井是对地层油气层变化情况进行实时监测的技术,它具有连续性、灵敏性等优势。在识别裂缝性油气层、轻质油气层、凝析油气层等方面发挥着重要作用,是发现与评价油、气层的

收稿日期:2014-07-31

修回日期:2014-10-31

网络出版时间:2015-05-06

基金项目:黑龙江省科技攻关项目(F2004-01);黑龙江省教育重大科研项目(10051x0001);黑龙江省教育科学技术研究项目(11551016)

作者简介:刘华莹(1969-),女,教授,硕士生导师,研究方向为智能计算。

网络出版地址:http://www.cnki.net/kcms/detail/61.1450.TP.20150506.1648.043.html

主要手段。但由于在常规钻井过程中会出现下钻、接单根、注样分析、钻井取芯、钻时变化、排量变化等情况,严重影响着气测录井资料的采集,因此利用气测录井资料解释气层的标准和图版很难建立,传统的 RBF 神经网络方法在针对气测录井数据处理时具有样本训练时间较长,收敛速度较慢等缺陷。因此,文中提出了一种基于粗集理论的气测录井数据归一化处理方法,有效提高了 RBF 网络的训练速度。

1 RBF 神经网络

针对气测录井数据易变化、杂乱无章等特点,采用 RBF 神经网络方法对其进行处理。RBF 神经网络^[5-7]方法相对于其他数据处理方法,具有更快的学习速度与更加简单的结构等特点。

RBF 神经网络具有三层静态前向功能,以及三层拓扑结构^[8-10],其中包含输入层、隐含层、输出层(如图1所示)。

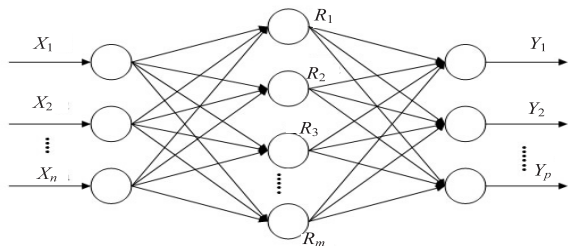


图1 RBF神经网络拓扑结构图

数据经输入层处理后传递到隐含层。网络结构的隐含层中,高斯函数执行信号描述功能。在网络结构的输出层中,线性函数执行信号描绘功能。隐含层节点函数使输入的信号产生连锁反应,当隐含层的中央函数接收到输入信号时,通过计算处理产生更大的输出,然后传递到隐含层节点。因此,RBF神经网络方法具有全局逼近与学习速度快等特点。

$$a_i(x) = \exp\left[-\frac{\|X - c_i\|^2}{2S_i}\right], i = 1, 2, \dots, m$$

其中, $a_i(x)$ 函数为网络结构中隐含层第 i 个节点的输出; X 为样本输入信号; c_i 为和 X 同维数的第 i 层节点的中心; S_i 为标准化长度; m 为隐含层节点的个数。

在径向基函数网络结构中,每个隐含层节点都有自己的一个径向基中心向量 c_i , 它和输入样本 X 拥有同样的维数,即 $c_i = [c_{i1}, c_{i2}, \dots, c_{in}]^T, i = 1, 2, \dots, m$ 。在这个网络结构里具有 n 个中心数,其中输入信号样本 $X = (x_1, x_2, \dots, x_n)^T$ 和 RBF 中心向量 c_i 的欧几里得范数 $\|X - c_i\|^2$ 的距离叫做隐含层节点的“净输入”。在径向基函数网络结构中,输入信号样本 X 到隐含层若干个中心向量的距离程度代表了每个隐含层节点的

输出^[11-12]。相对于径向基函数网络,其他的网络中隐含层的训练就是对权矩阵的调节,但是在径向基学习训练过程当中,函数主要执行的功能是选择合适的中心向量,然后分配给各个隐含层的节点。

在径向基函数网络中,输入层与输出层相连接并且产生连锁反应,分别实现了从 $X \rightarrow \alpha_i(x)$ 的非线性映射和 $\alpha_i(x) \rightarrow y_k$ 的线性映射,即

$$y_k = \sum_{i=1}^m w_{ik} \alpha_i(x), k = 1, 2, \dots, p$$

式中, p 为输出层节点。

2 基于粗集理论的归一化方法

2.1 粗集理论

假设 $K = (U, A)$ 为信息系统,样本的集合为论域 U 。假设样本中输入量的属性集合为 C , 输出量的属性集合为 D , 则样本的属性集 C 与 D 的并集,即表示为 $A = C \cup D$ 。 V_a 为每个属性 $a \in A$ 的属性值^[13-14]。

定义1(不可分辨关系):假设信息系统 $K = (U, A)$, 对于其中的每个子集 $B \subseteq A$ 定义一个等价关系表示为 $IN(B)$, 将其称作不可分辨关系,即 $2IN(B) = \{(x, y) \in U: \forall a \in B(a(x) = a(y))\}$ 。其中, $[x]B$ 是样本 $x \in U$ 中 B 的等价类。

定义2(上、下近似与边界):给定信息系统 $K = (U, A)$, 一个样本的集合表示为 $X \subseteq U$, 属性的集合表示为 $B \subseteq A$ 。则 B -下近似、 B -上近似和 B -边界表示为 $B_-(X) = \{x \in U: [X]B \subseteq X\}; B(X) \subseteq X\}; B_+(X) = \{x \in U: [x]B \cap X \neq \emptyset\}; BNR(X) = B_+(X) - B_-(X)$ 。

2.2 数据归一化方法

数据归一化方法中条件属性为样本的输入量信息,决策属性为输出量信息,最后应用离散方法对各属性执行离散功能,然后获得该信息系统的决策表。再通过约简信息的属性,得到简化的决策表。

方法1:设信息系统 $L = (M, N)$, 决策属性集定义为 $D = \{d\}$, $Z_d = 1, 2, \dots, r$, 定义 Z_d 对应于不同的类。

方法2:设信息系统 $L = (M, N)$, 条件属性集定义为 $C = \{c_1, c_2, \dots, c_n\}$, 规定 k 类样本的平均能量为

$$W_k = \frac{1}{m} \sum_{j=1}^m \sum_{i=1}^n V_{kf_i}^2, k = 1, 2, \dots, r$$

其中, V_{kf_i} 表示 k 类的第 j 个样本的 c_i 输入量对应的属性值。

方法3:对于信息系统 $L = (M, A)$, 条件属性集定义为 $C = \{c_1, c_2, \dots, c_n\}$, 规定 p 类第 j 个样本与 q 类第 k 个样本间的距离为

$$d_{p/q_k} = \left(\sum_{i=1}^n (V_{p_i c_i} - V_{q_i c_i})^2 \right)^{1/2}, p \neq q, p, q = 1, 2, \dots,$$

$r;j=1,2,\cdots,m;k=1,2,\cdots,t$

方法 4: p 类的第 j 个样本到 q 类中每一个样本间的最小距离定义为

$$d_{p,q_i} = \min d_{p,q_i}; p \neq q, p, q = 1, 2, \cdots, r;$$
$$j = 1, 2, \cdots, m; k = 1, 2, \cdots, t$$

方法 5:原输入样本按下式进行伸缩预处理:

$$y_{p,i} = x_{p,i} * (1 + \frac{S_1}{d_{p,q}}), W_p > W_q$$
$$y_{p,i} = x_{p,i} * (1 - \frac{S_1}{d_{p,q}}), W_p < W_q$$

其中, $p \neq q, p, q = 1, 2, \cdots, r; j = 1, 2, \cdots, m; k = 1, 2, \cdots, t$ 。

3 气测录井数据归一化处理

3.1 气测数据影响因素分析

气测仪器采集的气体浓度数据主要受外界环境与采集气体装置的影响。目前,常见的采集气体装置大多数都是电动脱气仪器,它可以逐步分离钻井液中的气体,但是它的脱气速度相对较低,而且钻井液与外界环境的变化对其采集的气体浓度也影响较大。

主要影响采集数据的原因在以下几个方面:

- (1)设备功能缺陷,导致气体采集量无法达标,造成错误的组分;
- (2)气候与温度环境随时变化,影响采集数据;
- (3)当脱气仪器中的气体采集满后,继续采集气体导致气测值超过 100%,气体平面过高;
- (4)设备吸入大量空气,因为空气中含有大量的 CO₂ 与 N₂,影响仪器中非烃类气体(CO₂和 N₂)的测量;
- (5)脱气仪器与采集终端处于较远距离时造成脱气仪器无法灵活的反应。

综合以上影响因素,以辽河油田为实验地区,随机选取气测仪进行试验操作。表 1 与表 2 是试验的数据和与结果。用 4 000 ml,1.15% 的标准气样进行实验,表 1 中的三台仪器不在同一地点进行实验,但是发现在时间上有不同,由于它们使用的是不同的样品泵与管线,所以造成时间上的差异。

表 1 使用 4 000 ml、1.15% 的球胆标准气样试样时间

录井队	仪器型号	样品泵	脱气器	管路延时
712 队	SK-IV	1'38"	3'23"	6'
813 队	SK-IV	55"	1'12"	49"
914 队	SK2000c	1'26"	5'33"	2'67"

在钻井液连续脱气仪器工作过程中,气泵首先抽汲钻井液中的气体后经由传输管道传送到脱气仪器内,然后仪器进行鉴定与分析。因为脱气时仪器需要对气体进行分析与数据的采集,所以样品泵的抽气速

度经常快于仪器脱出气体的速度。为了克服这一困难,工作人员会从仪器的空气入口补充一定量的空气,以满足样品泵的抽气量。经空气稀释后的气体再给全烃仪和色谱仪进行处理。此时仪器中气体所含量是脱出的气体与吸入空气含量的总和。因此,造成各台仪器所测量结果有差别,如表 2 所示。

表 2 使用 4 000 ml、1.15% 的球胆标准气样气测值

录井队	仪器型号	全烃值/mv		组分值/mv		
		样品泵	脱气器	样品泵	脱气器	
712 队	SK-IV	55	35	36	30	25
813 队	SK-IV	45	20	66.6	26	17
9 队	SK2000c	130	100	27.6	120	85

通过以上分析可以得到以下结论:

- (1)大气环境、工作人员的操作与仪器工作效率都影响气测数据的采集。
- (2)不同型号仪器,采集的数据往往不同。

3.2 气测录井数据归一化方法

文中设计的气测数据归一化方法如下:

- (1)选取合适的样本输入量与输出量,条件属性作为输入量,决策属性作为输出量,通过离散的方法处理输入量与输出量,得到信息系统决策表;
- (2)通过约简属性与约简规则,得到简化的决策表;
- (3)应用方法 1~5 对原始气测数据样本进行归一化处理;
- (4)利用归一化后的数据对 RBF 神经网络方法进行训练;
- (5)利用归一化后的数据样本对 RBF 神经网络方法进行验证。

流程如图 2 所示。

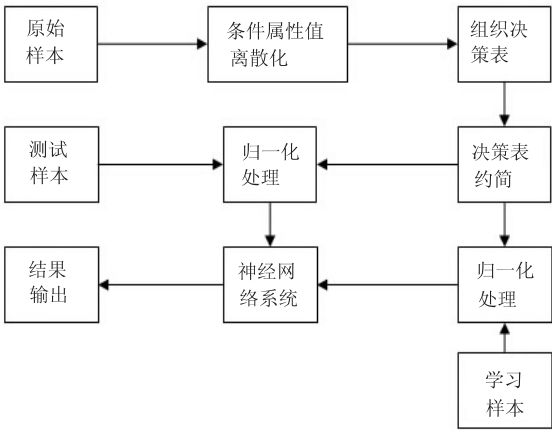


图 2 数据归一化方法

利用文中所设计的归一化方法处理实验气测数据,处理的内容包括:全烃及各组分(C₁、C₂、C₃)参数。表 3、表 4 是辽河油田气测资料全烃值进行校正

前后数据表。可以看出在未归一化处理时杂乱无章,通过归一化处理后的全烃数据可以有效地划分出油气水层,效果非常明显。

表 3 校正前气测资料全烃值

气测值		气水同层		含气层		含气水层	
孔隙度	数值	孔隙度	数值	孔隙度	数值	孔隙度	数值
11.31	0.24	1.57	3.54	11.31	0.24	1.57	3.54
7.91	0.35	2.43	0.26	7.91	0.35	2.43	0.26
5.85	0.47	3.55	4.31	5.85	0.47	3.55	4.31
5.89	0.61	8.93	5.24	5.89	0.61	8.93	5.24
5.97	0.71	7.65	0.98	5.97	0.71	7.65	0.98
7.77	5.61	5.44	3.67	7.77	5.61	5.44	3.67
4.53	7.42	6.53	0.84	4.53	7.42	6.53	0.84
11.13	11.54	9.83	5.29	11.13	11.54	9.83	5.29
8.53	9.78	4.32	0.78	8.53	9.78	4.32	0.78
6.01	10.12	5.96	1.23	6.01	10.12	5.96	1.23
8.15	9.78	6.12	5.54	8.15	9.78	6.12	5.54
7.54	9.34	11.48	8.94	7.54	9.34	11.48	8.94

表 4 校正后气测资料全烃值

气测值		气水同层		含气层		含气水层	
孔隙度	数值	孔隙度	数值	孔隙度	数值	孔隙度	数值
3.84	4.78	2.56	3.48	3.84	4.78	2.56	3.48
4.34	8.41	3.15	4.24	4.34	8.41	3.15	4.24
4.97	8.68	3.55	4.31	4.97	8.68	3.55	4.31
6.12	9.34	4.61	5.24	6.12	9.34	4.61	5.24
5.97	9.85	4.87	5.14	5.97	9.85	4.87	5.14
6.76	11.54	5.10	4.98	6.76	11.54	5.10	4.98
7.43	11.42	5.65	5.39	7.43	11.42	5.65	5.39
7.95	12.31	6.23	5.89	7.95	12.31	6.23	5.89
8.53	13.56	7.12	6.54	8.53	13.56	7.12	6.54
9.56	14.58	7.84	6.32	9.56	14.58	7.84	6.32
11.25	15.34	8.01	7.55	11.25	15.34	8.01	7.55
12.36	15.88	9.23	8.94	12.36	15.88	9.23	8.94

4 针对气测录井数据设计 RBF 神经网络

4.1 RBF 神经网络设计

(1)网络参数的选择。

RBF 神经网络结构具体划分为输入层、隐含层、输出层。确定这三层结构的单元数目与初始权值的选择十分重要。根据文中研究的实际需要,设计 RBF 神经网络的节点数与输入层、输出层的节点数相同,隐含层节点数与输入层的节点数相同。由于为非线性系统,初始值对学习是否达到局部最小与是否能够收敛以及训练时间长短的关系有很大影响。太大的初始权值,

易导致加权后的输入落入激活函数的饱和区,进而导致网络调节过程陷入停顿。因此,希望每个神经元经过初始加权后的输出值都接近于零。这样方便使所有神经元的权值都能够在它们的激活函数变化最大之处进行调节。因此,设所有权值初始值取 $[-1,1]$ 之间的随机数。

(2)传递函数选取。

用高斯核函数做隐含层的传递函数,用线性函数做输出层的传递函数。

(3)输入输出向量设计。

在得到输入与输出变量后,采用文中所设计的归一化方法进行处理,数据经过处理后明显提高了 RBF 神经网络的学习效率。

4.2 RBF 神经网络训练与仿真

采用 RBF 神经网络方法对归一化处理后的气测录井数据进行模拟训练。采用 newrbe () 函数作为 MATLAB 神经网络模拟学习的创建函数。在对网络创建过程中同时进行数据训练过程,使网络的误差趋近于 0。在训练过程中,径向基函数中的分布常数 Spread 至关重要。Spread 数值大小影响网络的预测性能平滑。但是 Spread 数值并不能过大,过大可能导致计算上出现问题。这里先将 Spread 数值设定为 1 并且以 0.2 的间隔递增。不同的 Spread 数值导致网络的训练误差不同。

5 结束语

文中在 RBF 神经网络方法处理气测录井数据的基础上,运用粗集理论提出了一种新的气测录井数据归一化处理方法。首先,对原始数据样本利用粗集理论进行约简;其次,选择合适的输入端,计算所选择样本的输入量属性值与其他样本输入量属性值间的最小距离。此最小距离则为样本需要伸缩的比例,用文中所提到的归一化方法对伸缩后的样本进行处理,然后用 RBF 神经网络方法对其进行训练。通过观察仿真结果,明显发现该方法缩短了训练时间,有效提高了 RBF 神经网络处理气测录井数据的效率,表明该方法是可行的。

参考文献:

[1] 裴亦楠. 石油开发地质方法论(二)[J]. 石油勘探与开发, 1996,23(3):48-51.
[2] 罗英俊. 油田开发生产中的保护油层技术[M]. 北京:石油工业出版社,1996.
[3] Appel M. Nuclear magnetic resonance and formation porosity [J]. Petrophysics,2004,45(3):296-307.

则数量决定了它们都只能发现部分缺陷。为了保证和提高软件测试的质量,笔者认为软件代码规则检查工具在软件静态分析时是行之有效的,但同时也要意识到测试工具的不足之处。

软件代码规则检查工具优点包括:

- (1)开发早期发现软件编码规则错误,易于修改,降低开发成本;
- (2)代码分析使用与开发的任何阶段,不需要牵扯太多其他因素;
- (3)不需要设计测试用例,不需要代码插装,节约时间;
- (4)发现问题时直接指引到问题所在行,易于修改;
- (5)开发早期发现错误,有助于开发人员发现个人编码风格的缺点。

软件代码规则检查工具缺点包括:

- (1)有可能会测试不全,造成测试质量无法 100% 保证;
- (2)需要辅以人工判断;
- (3)各家算法不尽相同,存在分析结果上的差异。

3 结束语

软件代码规则检查工具可以在软件开发的任何阶段使用,有助于降低软件成本和开发时间。随着此类工具分析技术及集成规则数量的不断发展,测试质量不断提高,渐渐受到软件开发人员、测试人员的钟爱。但无论如何,工具自身存在误报、漏报及算法片面性的因素,只是通过代码规则检查工具来进行检查是远远不够的,还需要在软件开发阶段实施动态测试。动态

测试与静态分析相结合,黑盒测试与白盒测试相结合,进而提升测试质量。不同测试工具都有自身不同的特点,必须结合自身需求,挑选趁手的兵器。

参考文献:

[1] 王雅文,宫云战,杨朝红. 软件测试工具综述[J]. 北京化工大学学报:自然科学版,2007,34(A01):1-4.

[2] 黄茂生. 软件自动化测试工具的评估与选择[J]. 电子质量,2007(12):22-25.

[3] 杨 宇,张 健. 程序静态分析技术与工具[J]. 计算机科学,2004,31(2):171-174.

[4] 邓青华. 软件自动化测试工具研究[J]. 软件导刊,2011,10(1):57-59.

[5] 周 涛. 航天型号软件测试[M]. 北京:宇航出版社,1999.

[6] 鞠秀娟,赵 明. 软件自动化测试概述及应用工具分析[J]. 计算机应用,2007,27(B06):317-318.

[7] 周伟明. 软件测试实践[M]. 北京:电子工业出版社,2008.

[8] 韩 柯,杜旭涛. 软件测试[M]. 北京:机械工业出版社,2003.

[9] 王 倩. 软件自动化测试工具的分类与选择[J]. 玻璃,2008,35(8):51-53.

[10] 亢 勇,陈自力,李 鹏,等. 面向对象的软件测试[J]. 测试技术学报,1999,13(2):80-88.

[11] 罗 娜,林和平,袁福宇. 面向对象软件测试的方法研究[J]. 东北师大学报:自然科学版,2004,36(1):39-45.

[12] 高艳霞. 软件测试过程要因分析[J]. 中原工学院学报,2004,15(4):73-75.

[13] 王 云. 基于软件测试的软件质量分析研究[J]. 警察技术,2013(2):40-41.

[14] 赵振宇. 嵌入式软件测试技术研究及应用[D]. 北京:北京邮电大学,2011.

(上接第 192 页)

[4] Morgan W B,Prison S J. Theeffect of fractional wettability on the archie saturation exponent[C]//Proc of fifth annual logging symposium. Midland,TX:[s. n.],1964.

[5] 张海传,刘钟阳,许东卫,等. 基于 RBF 神经网络模型的臭氧浓度软测量研究[J]. 大连理工大学学报,2010,50(6):1020-1023.

[6] 刘军霞,阳春华,王雅琳. 螺旋分级过程数学模型研究及应用[J]. 计算机工程与应用,2010,46(4):230-232.

[7] 周 勇,胡中功. RBF 神经网络理论及其在控制中的应用[J]. 武汉科技学院学报,2007,20(5):40-42.

[8] 董长虹. Matlab 神经网络与应用[M]. 北京:国防工业出版社,2005.

[9] Deng Julong. Spectrum mapping in grey theory[J]. The Journal of Grey System,2000(2):116-124.

[10] Deng J L. Grey forecasting model[M]//Grey system. Beijing: China Ocean Press,1988:54-69.

[11] Deng J L. Properties of the grey forecasting model GM(1,1)[M]//Grey system. Beijing: China Ocean Press,1988:70-78.

[12] 赵振勇,王 力,王保华,等. 遗传算法改进策略的研究[J]. 计算机应用,2006(S2):189-191.

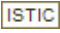
[13] 王 玮,蔡莲红. 基于粗集理论的神经网络[J]. 计算机工程,2001,27(5):65-67.

[14] 郝丽娜,徐心和. 粗糙集神经网络系统在故障诊断中的应用[J]. 控制理论与应用,2001,18(5):681-685.

基于粗集理论的气测录井数据归一化处理

作者：[刘华莹](#)，[孙春雨](#)，[张方舟](#)，[邱露露](#)，[高晓松](#)，[LIU Hua-ying](#)，[SUN Chun-yu](#)，[ZHANG Fang-zhou](#)，[QIU Lu-lu](#)，[GAO Xiao-song](#)

作者单位：[东北石油大学 计算机与信息技术学院, 黑龙江 大庆, 163318](#)

刊名：[计算机技术与发展](#)

英文刊名：[Computer Technology and Development](#)

年，卷(期)：2015(7)

引用本文格式：[刘华莹](#). [孙春雨](#). [张方舟](#). [邱露露](#). [高晓松](#). [LIU Hua-ying](#). [SUN Chun-yu](#). [ZHANG Fang-zhou](#). [QIU Lu-lu](#). [GAO Xiao-song](#) [基于粗集理论的气测录井数据归一化处理](#) [期刊论文] - [计算机技术与发展](#) 2015(7)