

ANALYSIS OF CONTROLLING FACTORS FOR HYDRAULIC FRACTURING PARAMETERS AND ACCUMULATED PRODUCTION USING MACHINE LEARNING

by

**Zhihua ZHU^a, Maoya HSU^b, Chang LI^a, Jiacheng DAI^b, Bobo XIE^a,
Zhengchao MA^b, Tianyu WANG^{b*}, Jie LI^a, and Shouceng TIAN^b**

^a Research Institute of Engineering Technology,
PetroChina Xinjiang Oilfield Company, Karamay, China

^b National Key Laboratory of Petroleum Resources and Prospecting,
China University of Petroleum, Beijing, China

Original scientific paper
<https://doi.org/10.2298/TSCI2404417Z>

This study, based on static data from over a thousand fracturing wells, employs data governance, data mining, and machine learning regression uncover principal controlling factors for production in the fracturing context. Utilizing multiple evaluation methods, the entropy weight method comprehensively scores and ranks the principal controlling factors. A machine learning production prediction model is established for validation. Results show that DBSCAN achieves better accuracy in identifying field anomaly data. For missing data, it is recommended to use tree models or neural networks instead of imputation or constant filling, as incorrect imputation significantly degrades model performance. The entropy weight method effectively integrates various correlation analysis results, providing a better connection with production compared to other approaches. This research utilizes large-scale field data to extract key parameters affecting production, supporting the establishment of high precision prediction models and the optimization of parameters for unconventional reservoir production forecasts

Key words: *hydraulic fracturing design, data-driven model, machine learning, production forecasting, hyperparameter optimization, feature selection*

Introduction

Horizontal well fracturing is crucial for unconventional oil and gas development, necessitating accurate post-fracturing production prediction for optimal parameter adjustments [1-4]. However, mechanistic models are time-consuming, hindering on-site development and rapid optimization, particularly in heterogeneous reservoirs [5]. With 1226 fracturing operations in Xinjiang oilfields, significant data is available for big data research. This study employs big data and artificial intelligence to manage on-site data, analyze production-controlling factors, and establish machine learning models for efficient fracturing design in unconventional oil reservoirs. Researchers have employed correlation algorithms and model evaluation methods to assess the relationship between well parameters and production, identifying principal controlling factors [6, 7]. Machine learning techniques, such as dimensionality reduction and

* Corresponding author, e-mail: wangty@cup.edu.cn

feature synthesis, capture non-linear connections between geological engineering parameters and production [8, 9].

In summary, machine learning-based analysis of principal controlling factors and production forecasting has shown initial effectiveness. However, challenges such as small sample size and limited parameter ranges persist. This paper collects fracturing well data, conducts data governance and mining, comprehends principal controlling factors, and builds a machine learning production forecasting model based on field data from Xinjiang oilfield, contributing to improved predictions and optimization.

Material and method

Data collection

This study compiled a dataset from Xinjiang oilfields, consisting of 1226 hydraulically fractured horizontal wells with 44 feature parameters. After preliminary cleaning, wells with a single well data missing rate exceeding 40% and parameters with a data missing rate exceeding 80% were excluded. The dataset includes geological, engineering, and production data. Parameters such as horizontal section length, porosity, permeability, oil saturation, reservoir length, and reservoir type constitute geological and engineering data. Engineering data include modified section length, number of fractured sections, number of fracture clusters, liquid intensity, proppant intensity, pre-flush fluid ratio, slip water ratio, maximum construction displacement, and average sand ratio. Geological and engineering data serve as input features during model training, while production data mainly involve cumulative oil production over 330 days per unit of modified section length, serving as the output target during model training. The data were normalized to the [0, 1] interval, and kernel functions were employed to describe the data distribution, as shown in fig. 1.

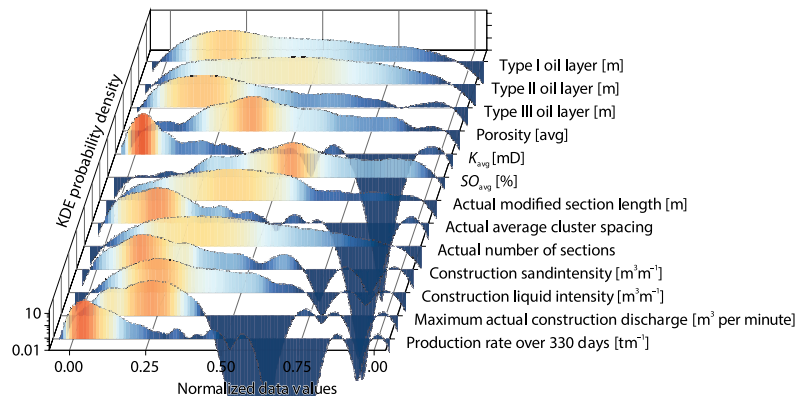


Figure 1. Waterfall plot depicting the kernel density distribution of selected data points

Upon observation of fig. 1, it is evident that the overall data distribution is uneven, with some data exhibiting a left-skewed pattern. Anomalies are particularly noticeable in parameters such as porosity, proppant intensity, and construction displacement, where significant outliers are present.

Data governance

Data governance encompasses outlier identification, missing value imputation, and categorical data encoding, as illustrated in fig. 2.

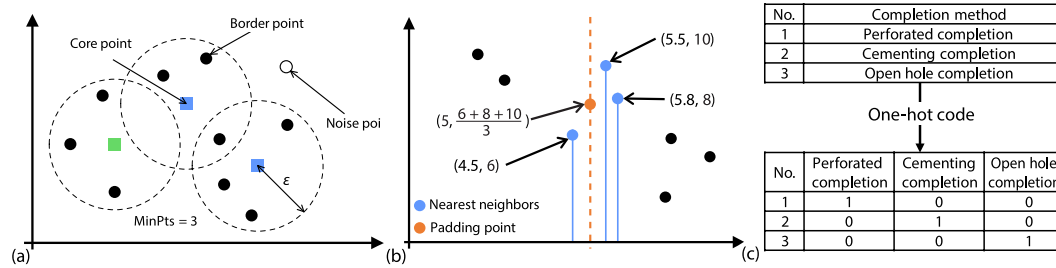


Figure 2. Data governance methods; (a) DBSCAN outlier identification, (b) the KNN missing value imputation, and (c) one-hot encoding of label data

To handle outlier data, we used the DBSCAN unsupervised clustering method. DBSCAN, a density-based clustering algorithm, categorizes and filters outliers by assessing the density of neighboring data points. For missing data, various methods were employed, including zero filling, mean imputation, and KNN modelling imputation. The KNN imputation references known data from neighboring samples, improving imputation quality with different reference sample sizes.

Feature importance evaluation

This study employs grey correlation and maximum mutual information for correlation analysis, along with embedded feature evaluation and SHAP post-interpretation assessment for importance evaluation. Utilizing results from various principal control analysis methods, it establishes an evaluation matrix for a comprehensive assessment of parameters through weighted calculations.

The entropy weight method is an objective weighting technique that determines the weights of indicators based on the information entropy of each indicator. The calculation formula:

$$n_{ab} = \frac{e_{ab} - \min(E_b)}{\max(E_b) - \min(E_b)}, P_{ab} = \frac{n_{ab}}{\sum_{a=1}^n n_{ab}}, \sum P_{ab} = 1, r_b = \sum_{a=1}^n P_{ab} \ln P_{ab} \quad (1)$$

where E_b is the data vector of evaluation indicator b in the evaluation matrix, n_{ab} – the normalized score value of evaluation parameter a for the positive evaluation indicator b , n – the number of evaluation parameters, m – the number of evaluation indicators, and Rb – the information entropy value of evaluation indicator b . The entropy weight of evaluation indicator b reads:

$$\omega_b = \frac{1 - r_b}{\sum_{b=1}^m (1 - r_b)} \quad (2)$$

Finally, the evaluation matrix E is multiplied by the indicator weight vector ω to obtain the comprehensive evaluation vector s , given:

$$S = E \times \omega = \begin{bmatrix} \zeta_1 & M_1 & J_1 & \varphi_1 \\ \dots & \dots & \dots & \dots \\ \zeta_n & M_n & J_n & \varphi_n \end{bmatrix} \begin{bmatrix} \omega_\zeta \\ \omega_M \\ \omega_J \\ \omega_\varphi \end{bmatrix} = \begin{bmatrix} \omega_\zeta \zeta_1 + \omega_M M_1 + \omega_J J_1 + \omega_\varphi \varphi_1 \\ \dots \\ \omega_\zeta \zeta_n + \omega_M M_n + \omega_J J_n + \omega_\varphi \varphi_n \end{bmatrix} = \begin{bmatrix} s_1 \\ \dots \\ s_n \end{bmatrix} \quad (3)$$

where ζ is the evaluation vector obtained from the grey relational analysis, M – the evaluation vector obtained from the maximum mutual information, J – the feature importance vector obtained from embedded feature evaluation, and φ – the evaluation vector obtained from the SHAP method.

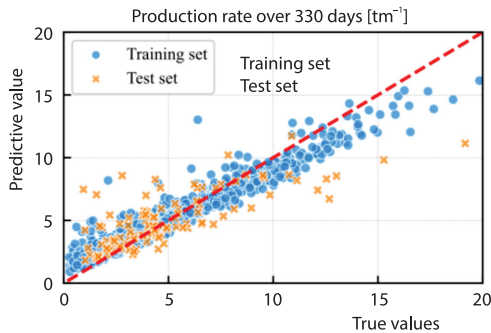


Figure 3. Comparative chart of AutoGluon prediction results

Results and discussion

Model performance

The training data underwent DBSCAN outlier removal, with no imputation for missing values. A test set, comprising 20% of the total dataset, was separated and excluded from training. Utilizing the AutoGluon framework, the training process involved 5-fold cross-validation bagging and a 2-layer stacking approach, using RMSE as the loss function. The training results, as depicted in fig. 3., show an average error of 0 and an R2 of 1 for the training set. For the test set, the RMSE is 2.630, MSE is 6.919, MAE is 1.588, and R2 is 0.768.

Comparison of data governance methods

The dataset was managed using DBSCAN, three standard deviations, and IQR methods to identify outliers, along with KNN nearest neighbor prediction, constant, and mean value imputation for handling missing values or not handling them, respectively. The model training results are compared in fig. 4.

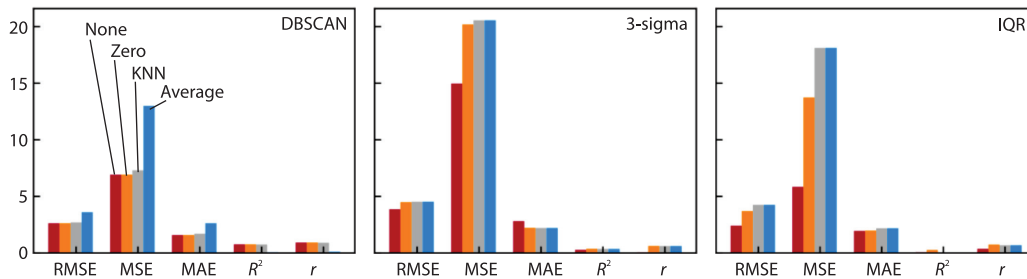


Figure 4. Test metrics for different data governance combinations

Analysis of controlling factors

The results of grey relational analysis, maximum mutual information, embedded feature evaluation, SHAP analysis, and entropy weight analysis are depicted in fig. 5.

In the experiment with varying input feature quantities, the model error is minimized when the feature quantity is set to the top 4 parameters, as illustrated in fig. 6. As the feature quantity further increases, the model error gradually increases. For different feature evaluation methods, subsets of data were created by taking the top four parameters and inputting them into the model for training, T .

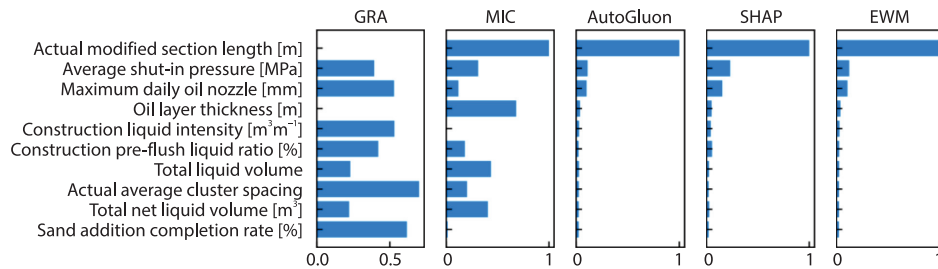


Figure 5. Evaluation of parameter importance

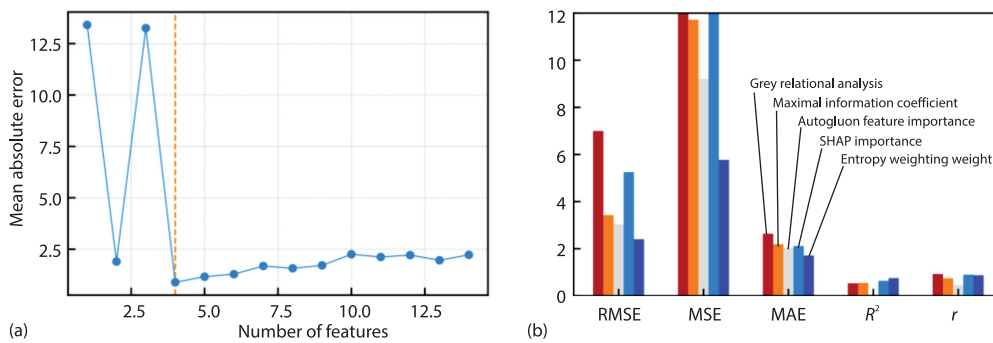


Figure 6. Impact of different principal evaluation methods on model accuracy

Conclusion

The results demonstrate that DBSCAN can achieve better identification accuracy for field anomaly data, whereas the IQR method and 3-sigma method exhibit a broader range of misjudgments, significantly compromising data integrity. It is not advisable to impute or use constant filling for missing data. Instead, tree models or neural networks should be employed to handle missing values automatically, as incorrect data imputation can lead to a substantial decline in model performance. The entropy weight method effectively integrates the results of various correlation analysis methods. The combination of principal controlling factors selected using this method exhibits a stronger correlation with production compared to other approaches. These principal controlling factors can be combined into a feature subset for input into the model, and the reasonableness of the principal evaluation can be judged based on the changes in model accuracy. This research effectively utilizes large scale field data to extract key parameters affecting production, providing technical support for the establishment of high precision prediction models and the optimization of parameters for unconventional reservoir production forecasts.

Acknowledgment

The authors would like to thank the financial project supported by the National Natural Science Foundation of China (Grant/Award No. 5231001009).

References

[1] Liu, H., *et al.*, Application Status and Prospects of Artificial Intelligence in The Refinement of Waterflooding Development Program, *Acta Petrolei Sinica*, 44 (2023), 9, pp. 1574-1586
 [2] Li, Y., *et al.*, Application Status and Prospect of Big Data and Artificial Intelligence in Oil and Gas Field Development, *Journal of China University of Petroleum (Edition of Natural Science)*, 44 (2020), 4, pp. 1-11

- [3] Kuang, L., *et al.*, Application and Development Trend of Artificial Intelligence in Petroleum Exploration and Development, *Petroleum Exploration and Development*, 48 (2021), 1, pp. 1-11
- [4] Sheng, M., *et al.*, Research Status and Prospect of Artificial Intelligence in Reservoir Fracturing Stimulation, *Drilling and Production Technology*, 45 (2022), 4, pp. 1-8
- [5] Anton, D. M., *et al.*, Data-Driven Model for Hydraulic Fracturing Design Optimization: Focus on Building Digital Database and Production Forecast, *Journal of Petroleum Science and Engineering*, 194 (2020), 107504
- [6] Yue, X., *et al.*, Research on Main Control Factors Influencing Fracturing Effect of Jiaoshiba Area Based on Grey Relational Analysis, *Proceedings, ARMA-CUPB Geothermal International Conference*, ARMA, Beijing, China, 2019
- [7] Pu, X., *et al.*, The Main Controlling Factor Analysis and Comprehensive Evaluation on Mid-deep Clastic Reservoir in Qikou Sag, Bohai Bay Basin, North China, SEG Technical Program Expanded Abstracts, Society of Exploration Geophysicists, Houston, Tex., USA
- [8] Shelley, R., *et al.*, Machine Learning and Artificial Intelligence Provides Wolfcamp Completion Design Insight, *Proceedings, 9th Unconventional Resources Technology Conference*, American Association of Petroleum Geologists, Houston, Tex., USA, 2021
- [9] Song, X., *et al.*, Productivity Forecast Based on Support Vector Machine Optimized by Grey Wolf Optimizer, *Lithologic Reservoirs*, 32 (2020), 2, pp. 134-140