# Machine Learning Techniques for Quantifying Characteristic Geological Feature Difference

Xiaojin Tan, Wenyue Sun

Stanford University

xtan1@stanford.com wenyue@stanford.com

**Abstract**

*This paper presents a novel methodology of quantifying the model mismatch part in the Bayesian framework for reservoir model inversion. During model inversion, the updated models are required to obey the prior characteristic geological information, which poses the issue of violation quantification. Recent approaches for this penalty normally require covariance matrix of model parameters, which are often impractical for three reasons: 1) they require gaussian assumption and 2) the computing of the inverse of the covariance is very time consuming or the approximation of the covariance skews the estimation of the uncertainty. 3) quite large amounts of models are required to achieve a stationary covariance estimation. In this paper, we propose an alternative approach that while circumventing the computing or approximation of the covariance, can provide a reliable estimation on the model misfit penalty as well as relax gaussian assumption. In order to render the model misfit quantitative, it is argued that this paper relies on a statistical learning for two totally different sets of models: prior models and reference models. After defining within-model distance and between-model distance, misfit between any two models can be trained through machine learning technique. SVM algorithm is used in this paper and turns out to be very effective in quantifying the similarities with respect to both the number of training models required and relaxation of the gaussian assumption of the model.*

## I. Introduction

Oil companies often need to make predictions on oil/gas productions for further marketing decision making, and this will involve three steps: reservoir characterization, modeling, flow simulation, and future model updating. For the modeling part, normally we are able to generate lots of reservoir models that follow the prior information (usually represented through some characteristic geological features) gotten from the characterization stage, but during the model updating step, however, there is no guarantee that the updated posterior model will follow the prior information. Thus one of the most critical problems for model updating is that we need a reliable method to quantify the violation of updated models to the prior characteristic geological feature.

In a general sense, model updating is a typical inverse problem, which is ill-posed if only historical production data (which is the output from the simulation over our generated reservoir models) is considered. Additional information, usually prior knowledge about the reservoir, is required to constraint the solutions. Tikhonov regularization is one popular method in which the problem is reformulated not only to match the production data, but also our prior knowledge about the reservoir (which we'll refer as model mismatch part later on).

There have been some works from literatures in approaching this problem on the model mismatch part, most of them, however, normally require gaussian assumption on the model, which is quite impractical for most of the cases. At the same time, they require inverse of the covariance matrix, which is used to describe the two-point statistical behavior of the reservoir. For real cases often involving

models of over millions of grid blocks (which can be considered as the number of freedoms), the requirement for inverse of the covariance is usually infeasible given nowadays computational power.

In order to circumvent the computation of the inverse of covariance matrix and, most importantly, to relax the requirement for gaussian assumption, in this paper, a machine learning based methodology is proposed that allows quantifying model mismatch part regardless of the nature of the reservoir models, namely, whether they are gaussian or nongaussian, discrete or continuous, have small or large grid dimensions. Our method starts with a definition of prior and reference models: the prior models are those we generated during the reservoir modeling part to represent our prior understanding of the reservoir, which is the one that we want to match with; the reference models are those we generated purposely with different geological features from the prior. Then a machine learning technique – SVM is used in this paper to give an estimation about the distance between updated models and prior models. And this distance turns out to be able to act as the model mismatch term naturally.

While one can debate the particular subjective choices about the reference models are made in the presented methodology, the variety of examples illustrate that the results obtained from this methodology are consistent with expectations. It should be understood that we are not seeking the best reference models that should be created, but an reference model that's able to match our demands here, which is to estimate the relative model mismatch between different updated models and the prior models. However, there is no doubt that better reference models may improve the overall reliability of this method, thus future work may be interested in this part.

## II. Methods

Support Vector Machine learning algorithm is one of the best "off-the-shelf" supervised learn-

ing algorithm. In this paper, SVM is used to quantifying the violation level of a specific model to prior models. Fig 1 shows a reservoir model (each block contains a value indicating the permeability at that location) that follows the prior information we specified, for which a 45 degree long correlation trend is observed. Fig 2 shows one reference model, which has the same histogram with the prior model but a different long correlation trend . In this paper, we'll use SVMs to quantifying the difference between a set up updated models and the prior model. For details about how SVM algorithms works, we refer to [2, 3]. Here we label all prior models as 1, and reference models -1. The training data is a single column composed of permeability at all the blocks, this can be easily extended to other types of data, like porosity or z-transmissibility. Since we can generate multiple prior and reference models easily using research software GSLIB [1], training data sets can be easily achieved. We use the same amount of training data for prior models and reference models respectively, and for prediction stage, we used three testing models which we know to be clearly different from the prior models. For each test model, we generate multiple realizations that share the same geological feature. Then these testing models are used to test and the reliability of SVM algorithm for quantifying geological feature difference. In this report, we only use linear kernel instead of higher order kernels, the reason is that even linear kernel turns out to be quite effective, and using higher order kernel doesn't improve the performance significantly, however, for larger and more complex models, higher order kernels might be interested. Also, KKT violation is allowed to alleviate the influence of casual outliers.
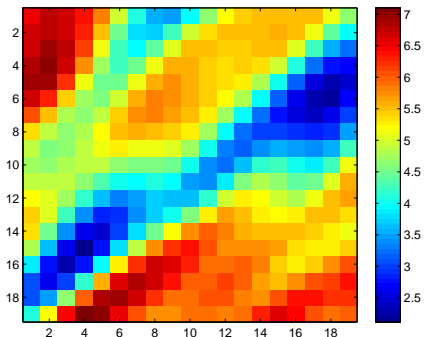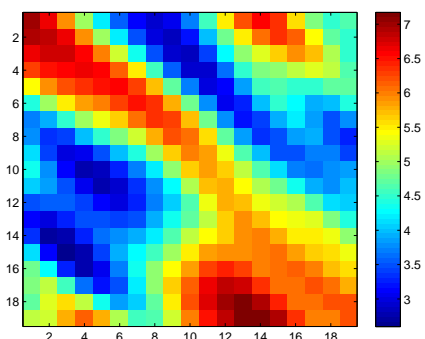
**Figure 1:** *A Prior Reservoir Model*



**Figure 2:** *A Reference Reservoir Model*

## III.  Results

Fig 3-5, gives three test cases we generated, from which we can see that test model 3 is closer to prior model, whereas test model 4 and 5 are closer to the reference model. Table 1 shows the results of classifying and quantifying model difference using SVMs. The training data contains 30 prior and reference models respectively. The testing data set contains 100 realizations for each model respectively, since we know that these 100 realizations follows the same geological feature, thus taking the mean value of them is reasonable. The prediction results are the average of those for each type of model as shown in table 1. Model 1 and 2 in the table are simply the prior and reference models we used to train the algorithm, and the reason to test them is to check the self-consistency of this algorithm, which partially used the idea of "cross-validation".

The prediction data in table 1 is between 0

and 1, where 0 means it's classified as closer to prior model, and 1 closer to reference model. The function margin for each model here is simply treated as the average of all of the realizations, which intrinsically indicates how confident we are in making that prediction, and this margin can act as a measurement of the model mismatch term.
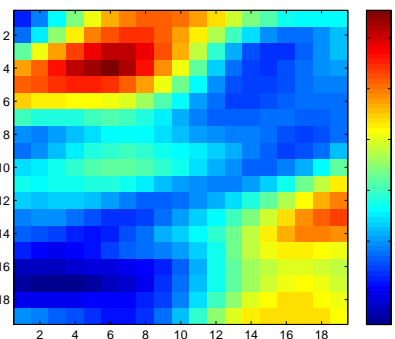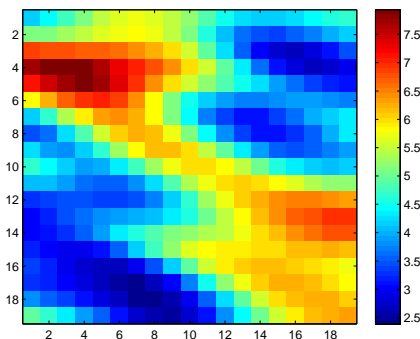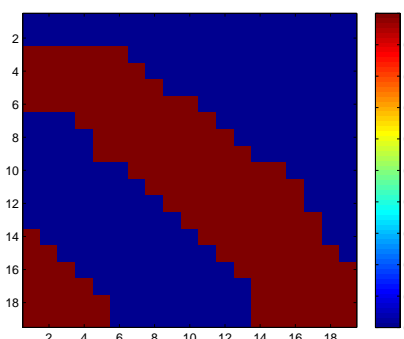


**Figure 3:** *Test Model 3*



**Figure 4:** *Test Model 4*



**Figure 5:** *Test Model 5*

**Table 1:** *Prediction Using SVM*

| Model | Prediction | Margin |
|-------|-----------|--------|
| 1 | 0.03 | -4.09 |
| 2 | 0.99 | 2.04 |
| 3 | 0.09 | -3.41 |
| 4 | 0.75 | 1.29 |
| 5 | 0.99 | 0.87 |

The last important issue left is to test whether the functional margin can provide a reliable (stable) estimation about the mismatch term. Given the observation that model 3 is closer to model 1 compared with model 4, we can use the following algorithm to generate a set of mixed gaussian model as testing models.

$$m_{mixed}^{(i)} = \frac{m_3^{(1)} + \gamma^{(i)} m_4^{(1)}}{1 + \gamma^{(i)}}$$

where $\gamma^{(i)} = 0.01 \times i, \quad i = 0, 1, ..., 50.$

The superscript for $m_3, m_4$ means the realization we choose from model 3 and model 4, $\gamma$ is a set of coefficients ranges from 0 to 0.5, which stands for a linear combination coefficient for mixing model 4 into model 3, as expected, the larger $\gamma$ is, the more dissimilar the mixed model will be with the prior models. Thus the margin value should be able to change smoothly with respect to $\gamma$ within a certain range, which is required for stable estimation.
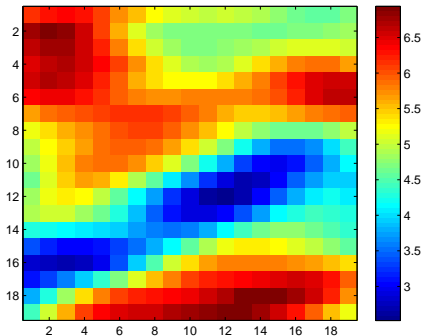


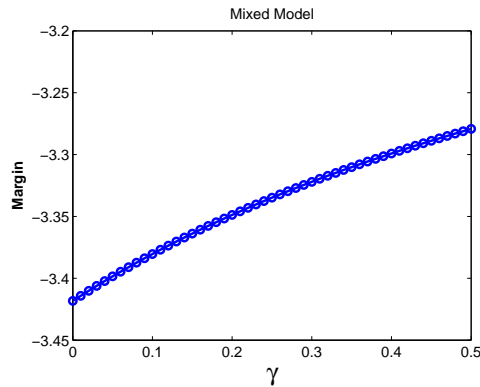**Figure 6:** *Mixed Model at $\gamma = 0.5$*



**Figure 7:** *Margin Value for Mixed Models*

Figure 6 shows a mixed gaussian model when $\gamma = 0.5$, from which we can clearly see that not only a obvious 45 degree trend observed, but also a -45 degree trend observed at the upper left corner, which is the effect carried by mixing model 4. Figure 7 gives the margin value as a function of $\gamma$.

## IV. Discussion

In table 1, the first two is accomplished by performing cross-validation, using selected prior models and references models as test set, the prediction for them are 0.03 and 0.99 respectively, which means that for prior models as test model, the probability of getting correct prediction is about 97%, and for reference model as test model, the probability is about 99%. Both of them gives very good predictions. Also, we can view this from the averaged functional margin value. Support vectors have margin value to be 1 all the time, for prior models regarded as test models, we have averaged margin value about 4 (absolute value), which is considerably larger that 1 (this judgement is quite obvious if we assume that the feature vector follows Gaussian distribution). Model 3, which is generally closer to the prior models as observed from figure 1 and figure 2, is predicted to be prior model 91% the time, and for those predicted to be 0, the function margin value is about 3.4, which, as illustrated, shows how closer model 3 is to model 1. Also, the algorithm gives bigger margin val-

ue of model 1 compare with model 3, which is reasonable in the sense that prior models are always the closest to themselves. Similarly, model 4 is predicted to be reference model 75% of the times, which again is as expected once comparing figure 2 with figure 4.

The most interesting results are for test model 5, which is clearly a non-gaussian model. Comparing figure 5 with figure 1 and 2, we find that model 5 is definitely closer to reference models since they both have a -45 degree oriented correlation feature, but on the other side, unlike model 3 and model 4, it's not largely closer to reference model than it is to the prior model. The reason is that if we take a walk at a 45 degree direction, we will also see a long correlation of the color (here color represents permeability of that block), thus what we should expect is that SVM may have low prediction error, but also with a relatively lower margin value, which is exactly what we observed in table 1 for the value corresponding to model 5. SVM predicts correctly 99% of the time, but on the other size, the margin value is only about 0.87, which is relatively lower than that for model 4.

From figure 7, we observed that the margin value is a smooth function of the $\gamma$ value, and the absolute value of the margin tends to decrease as $\gamma$ increases. This indicates that as we add more -45 degree trend features, the mixed model violated more with respect to the prior 45 degree trend feature, which is as expected. Thus the margin value predicted can potentially act as a measurement for the model mismatch term.

## V. Conclusions

In the report, we initiated the usage of machine learning technique in quantifying the ge-

ological feature difference, which gives relatively good results. The following conclusions can be made from this report:

1. Even though the reservoir model can have large dimension, not too many training models may be required for SVM algorithms. In our particular case, the model dimensions are over 300, but 30 number of training sets can already give great results based on the high prediction accuracy of model 1 and model 2 in table 1.

2. For models of different geological features compared with the prior models and reference models, SVM method can also give quite good prediction results.

3. Function margin varies smoothly with changes of test models, and captures the change of mismatch correctly, which proves its potential to act as an effective tool for quantifying the model mismatch term in the inverse modeling formulation.

## References

[1] Deutsch, C. V. and Journel, A. G. (1998). *GSLIB: Geostatistical Software Library and User's Guide, 2nd edition.* New York: Oxford University Press

[2] Vapnik, V. (1995) *The Nature of Statistical Learning Theory.* Springer-Verlag: New York

[3] Vapnik, V. (1998) *Statistical Learning Theory.* John Wiley. New York