# Application of Partial Least Square Regression in Uncertainty Study Area

Yingying Chen and Karlene A. Hoo‡

Department of Chemical Engineering

Texas Tech University

Lubbock, TX 79409, USA

Yingying.chen@ttu.edu and Karlene.Hoo@ttu.edu

*Abstract*— The aim of this work is to show how partial least squares (PLS) regression when combined with two other techniques Karhunen-Loeve (KL) expansion and Markov chain Monte Carlo (MCMC) can be efficient and effective at addressing parameter uncertainties that affect the predictive ability of a model for critical applications such as monitoring and control. We introduce a combination of PLS regression and KL to develop a reduced-order model (ROM) that captures the uncertain parameters effect on the model outputs, and the combination of PLS regression and MCMC for efficient updates of the uncertain parameter distributions. Two examples, a tubular reactor and an oil producing reservoir are presented to demonstrate these concepts.

## 1. INTRODUCTION

Partial least squares (PLS) regression is a statistical technique that was developed in the 1960's as an econometric technique in the social sciences [1]. The objective of a PLS developed model is to predict a set of dependent variables from a set of independent variables (predictors), which are the factors that have the largest effect on the dependent variables. From the predictors a set of latent variables that has the best predictive power is extracted and used to realize the prediction [2]. The PLS technique performs a simultaneous decomposition of predictors and dependent variables with the constraint that the dominant latent variables explain as much as possible the covariances between predictors and dependent variables.

PLS as a regression technique is used widely today in various and different disciplines (medicine to manufacturing). In this paper, PLS regression is used in the study of uncertainty. Uncertainty always exist in the description of models to represent real process phenomena. For example, the heat transfer process between hot and cold streams is characterized by a heat transfer coefficient. The value of this parameter is not exact, it contains some degree of error in the estimation of its bulk value. The error in this value can affect the final design of the heat transfer device, the overall economics and the estimation of the exiting stream temperatures. This work will demonstrate that PLS when combined with other techniques can predict the relationship between the uncertain parameters and the model outputs. Such knowledge can be useful for applications such as model-based control, online monitoring and fault diagnosis.

The paper is organized as follows. Sections 2 presents an overview of the methods used to address: sampling and sequencing effectiveness to propagate the uncertainties, efficient updating of the uncertain parameter estimates to maintain model accuracy, and improving computational efficiency for online applications by identifying a reduced-order model that captures the relationships between the parameters and the model outputs. Two examples, a chemical tubular reactor and an oil producing reservoir, are introduced in section 3 to demonstrate these concepts. An analysis of the results also is presented. Lastly, section 4 summaries the major contributions of this work.

## 2. METHODOLOGY

### A. PLS Regression Technique

The following was excerpted from [3]. Let $X^{n \times I}$ be a matrix of $I$ predictors collected on $n$ observations that describe $J$ dependent variables, $Y^{n \times J}$. Decompose both $X$ and $Y$ as a product of a set of orthogonal factors ($T$) and a set of loadings ($P$),

$$\begin{aligned} \mathbf{X} &= \mathbf{TP'} + \mathbf{E} \\ \mathbf{Y} &= \mathbf{TBC'} + \mathbf{F} \end{aligned} \tag{1}$$

The columns of $T$ are the latent vectors, $P$ is the coefficient matrix of $X$, the diagonal elements of $B$ are the regression weights, $C$ represents the weights of the dependent variables, and $E$ and $F$ are the matrices of residual errors.

To specify the latent vectors in $T$, two sets of weights $w$ and $c$ are needed to create a linear combination of the columns of $X$ and $Y$ such that their covariance is a maximized. The goal is to obtain a first pair of vectors $t_i = Xw_j$ and $u_i = Yc_j, i = j$ with constraints that $b_i = t_i'u_j, i = j$ is maximal and $w_i'w_j = 1, i = j, t_i't_j = 1, i = j$. It then folllows that $p_i = X't_i$.

Procedurally, let $Q = X$ and $R = Y$. Then column center and normalize $R, Q$.

Step 1: Initialize the vector $u$ with random values

Step 2: Estimate weights for $X$, $w \propto Q'u$

Step 3: Estimate $X$ factor scores, $t_1 = Qw_1$

Step 4: Estimate weights for $Y$, $c_1 \propto R't_1$

Step 5: Estimate $Y$ scores, $u_1 = Rc_1$

Step 6: Return to step 2 if $t_1$ has not converged. Otherwise continue

Step 7: Calculate $b = t_1'u_1$

Step 8: Compute the loadings for $X$: $p_1 = Q'T$.

Step 9: Subtract the effect of $t_1$ from both $Q$ and $R$: $Q = Q - t_1 p_1'$ and $R = F - b t_1 c_1'$. $b$ is a diagonal element of $B$.

Step 10: Go to step 1 until the matrix $Q$ becomes null.

The symbol $\propto$ represents a normalization of the result. The above relations show that $w_1$ is the first right singular vector of $X'Y$ and $c_1$ is the first left singular vector of $X'Y$. Similarly, $t_1$ and $u_1$ are the first eigenvectors of $XX'YY'$ and $YY'XX'$, respectively [3].

The prediction of dependent variables is based on a multivariate regression given by,

$$\hat{Y} = TBC' = XB_{PLS} \tag{2}$$

where $B_{PLS} = (P')^{-1}BC'$.

### B. Karhunen-Loève Technique

A system (M) of nonlinear partial differential equations (PDEs) with appropriate initial and boundary conditions can be used to represent a computational model of a physical process. Solutions of these PDEs are infinite series solution that cannot readily be used for real-time model-based control applications. To address this issue, it is not unusual to substitute a reduced-order model (ROM) to overcome this limitation. The Karhunen-Loève (KL) expansion [4] is one such technique that when combined with a suitable solution method gives a satisfactory ROM for the designed operating targets. To use the KL method requires a large collection of data that are generated experimentally or numerically (simulation of a physical model) [5]. As stated above, the parameters in all physical models are never known exactly. This means that the accuracy of the model's predictions may be affected by these uncertainties in some complex and nonlinear fashion. The technique known as Latin hypercube Hammersley sequence sampling (LHHS) has been shown to be computationally efficient when compared to Monte Carlo at sampling multiple uncertain distributions to cover the uncertain regions effectively [6]. In this manner, the LHHS technique provides the conditions for the PLS technique to find the relationships between the uncertain parameters and the model's outputs.

It is recognized that the outputs of M, $y(t,z) \in R^n, t \in R, z \in \Omega_{0,1}$ may be many compared to the number of uncertain parameters, $\Theta_p^m, m \ll n$. Additionally, the set $y$, is not independent. For example, the measured value of the reactor temperature at a spatial location $z(t) = 1/2L$ is not independent of the measurement at $z(t) = 1/3L$ where $L$ is the reactor length. To reduce the size of $y$ and concentrate the information about the relationships between $\Theta_p$ and $y$ a ROM can be generated using the KL technique that is similar to PLS in that it represents the dominant relationships using a small number of empirical eigenfunctions (EEFs) that are dictated by the data [7],

$$M(z,t) = \bar{M}(z,t) + \sum_{k=1}^{K} \sqrt{\lambda_k} \varsigma_k(t) \psi_k(z) \tag{3}$$

where $M(z,t)$ is an element of the data matrix $M$, $\bar{M}$ denotes the mean of $M$, $\psi_j \in \Psi$ and $\lambda_j \in \Lambda$ are the eigenfunctions and eigenvalues of the covariance of $M$, respectively and the coefficients $\varsigma_k$ are the projections of $M$ onto $\psi$, K is the number of $\psi_k$ that represent the dominant characteristics of the data $M$.

It is very convenient to predict the coefficients, $\varsigma$, of $\psi_k$ from knowledge of the $\Theta_p$ by an application of PLS. To apply PLS, let the set of sample points of $\Theta_p$ be $X$ and the set of $\varsigma$ be $Y$ given in Equation (1). Then, given any samples of $\Theta_p$, $\varsigma$ can be calculated from the relationship in Equation (2). It then follows, that with the known $\varsigma$, $y_{ROM}$ can be obtained without resorting to simulations of M. It is remarked that a limitation to the proposed approach is that both PLS and KL are linear techniques.

### C. Markov Chain Monte Carlo

The prediction characteristic of the PLS regression technique also can be used to update $\Theta_p$ when combined with a technique such as Markov chain Monte Carlo (MCMC). The MCMC technique generates parameter values from a constructed Markov chain that converges to a stationary distribution [8], [9]. The adaptive Metropolis (AM) algorithm that generates the uncertain parameters in a single iteration is used in this work [10].

Start with an arbitrarily chosen initial vector of parameters $\Theta_p = \Theta_p^i$. A candidate set, $\Theta_p^*$, is generated from a proposed density distribution based on the current values $\Theta_p^i$. Next, compute the acceptance probability, $\alpha$ as a function of $\Theta_p^i, \Theta_p^*$ and M. If $\Theta_p^*$ is accepted with acceptance probability

$$\alpha = \min \left\{ 1, \frac{P(X \mid \Theta_p^*)P(\Theta_p^*)}{P(X \mid \Theta_p^i)P(\Theta_p^i)} \right\} \tag{4}$$

it then follows that $\Theta_p^{i+1} = \Theta_p^*$, otherwise $\Theta_p^{i+1} = \Theta_p^i$. Here $P(X \mid \Theta_p)$ is the likelihood function of the observed data $X$ and $P(\Theta_p)$ is the prior distribution of $\Theta_p$.

The likelihood function is a multi-normal joint probability density function of the data,

$$P(X \mid \Theta_p) = (2\pi\sigma_\varepsilon^2)^{-N_X/2} \exp \left\{ -\frac{\sum_{j=1}^{N_X} [X(z_j) - M(z_j; \Theta_p)]^2}{2\sigma_\varepsilon^2} \right\} \tag{5}$$

where $N_X$ is the number of data points, $X(z)$ is an observed datum at location $z$. $M(z; \Theta_p)$ is the result from M(z), $\Theta_p$ is the vector of uncertain parameters to be estimated from the observed data. The error term is given by, $\varepsilon(z) = X(z) - M(z; \Theta_p)$ and $\sigma_\varepsilon^2$ is its variance.

Update of the uncertain parameters distribution in M with MCMC cannot be done without comparing the observed data and $y$. A huge number of sets of $\Theta_p$ must be simulated by M to generate $y$. This step constitutes a large computational burden. To overcome this inefficiency, PLS regression is applied to find the relationship between $\Theta_p$ and $y$. But since the set $\Theta_p$ is smaller than the set $y$, the KL technique is applied first to reduce the size of $y$.

## 3. EXAMPLES

The following examples are cited from [6] and [7].

### A. Tubular Reactor to Produce Benzene

The example is a nonlinear tubular reactor system used for the production of benzene from the hydrodealkylation of toluene [11] (see Figure 1),
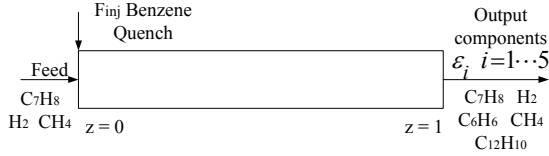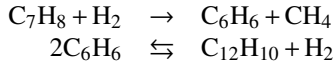


Fig. 1.   Plug Flow Reactor.

$$C_7H_8 + H_2 \quad \rightarrow \quad C_6H_6 + CH_4$$
$$2C_6H_6 \quad \leftrightarrows \quad C_{12}H_{10} + H_2$$

The first reaction is irreversible and the second is an equilibrium reaction. The reaction temperature (deg F) is such that $1150 < T < 1300$ and the reactor pressure is 500 psia. The ratio of $H_2$ and $C_7H_8$ is maintained at 5:1 to prevent coking. The chemical phenomena can be described by a system ($M$) of nonlinear PDEs in dimensionless quantities, [11],

$$\frac{\partial \varepsilon_1}{\partial \tau} = -v[\frac{\partial \varepsilon_1}{\partial \tau_1} + \frac{\varepsilon_1}{\theta}\frac{\partial \theta}{\partial \tau_1}] - \varepsilon_1 \varepsilon_2^{0.5}\theta^{1.5}e^{\gamma_1^{(\theta-1/\theta)}}$$

$$\frac{\partial \varepsilon_2}{\partial \tau} = -v[\frac{\partial \varepsilon_2}{\partial \tau_1} + \frac{\varepsilon_2}{\theta}\frac{\partial \theta}{\partial \tau_1}] - \varepsilon_1 \varepsilon_2^{0.5}\theta^{1.5}e^{\gamma_1^{(\theta-1/\theta)}}$$
$$+ k_2(\varepsilon_3\theta)^2 e^{\gamma_2^{(\theta-1/\theta)}} - k3\varepsilon_2\varepsilon_5\theta^2 e^{\gamma_3^{(\theta-1/\theta)}}$$

$$\frac{\partial \varepsilon_3}{\partial \tau} = -v[\frac{\partial \varepsilon_3}{\partial \tau_1} + \frac{\varepsilon_3}{\theta}\frac{\partial \theta}{\partial \tau_1}] + \varepsilon_1 \varepsilon_2^{0.5}\theta^{1.5}e^{\gamma_1^{(\theta-1/\theta)}} -$$
$$2k_2(\varepsilon_3\theta)^2 e^{\gamma_2^{(\theta-1/\theta)}} + 2k3\varepsilon_2\varepsilon_5\theta^2 e^{\gamma_3^{(\theta-1/\theta)}} + F_{B_m}$$

$$\frac{\partial \varepsilon_4}{\partial \tau} = -v[\frac{\partial \varepsilon_4}{\partial \tau_1} + \frac{\varepsilon_4}{\theta}\frac{\partial \theta}{\partial \tau_1}] + \varepsilon_1 \varepsilon_2^{0.5}\theta^{1.5}e^{\gamma_1^{(\theta-1/\theta)}}$$

$$\frac{\partial \varepsilon_5}{\partial \tau} = -v[\frac{\partial \varepsilon_5}{\partial \tau_1} + \frac{\varepsilon_5}{\theta}\frac{\partial \theta}{\partial \tau_1}] + k_2(\varepsilon_3\theta)^2 e^{\gamma_2^{(\theta-1/\theta)}}$$
$$- k3\varepsilon_2\varepsilon_5\theta^2 e^{\gamma_3^{(\theta-1/\theta)}}$$

$$\frac{\partial \theta}{\partial \tau} = \frac{1}{\zeta}(H_{r1}\frac{\partial \varepsilon_1}{\partial \tau} - H_{r2}\frac{\partial \varepsilon_5}{\partial \tau} + Q(\theta_F - \theta)$$
$$- v(\zeta\frac{\partial \theta}{\partial \tau_1} - H_{r1}\frac{\partial \varepsilon_1}{\partial \tau_1} + H_{r2}\frac{\partial \varepsilon_5}{\partial \tau_1})) - F_{B_m}\zeta_B$$

The outputs of $M$, $y$, are benzene concentration and reactor temperature. The set $\Theta_p$ that affects $y$ includes the reaction rate and heat of reaction of the first reaction, and the fresh benzene injection rate used for quench purposes. A means of validating the ROM is to compare the ROM results, $y_{ROM}$ to that of $y$. Figure 3 shows $y_{ROM}$ ($\circ$) and $y$ ($\triangle$) when a +5% bias in the mean values of $\Theta_p$ is introduced. Three $\psi$ are used, because 99.49% of the output data characteristics can be explained. The results show that $y_{ROM}$ tracks $y$ satisfactorily. The maximum relative errors between $M$ and the ROM are listed in Table I.

TABLE I
MAXIMUM RELATIVE ERRORS BETWEEN $y_M$ AND $y_{ROM}$.

| ROM (+5%) | Benz | Temp |
|---|---|---|
| With uncertainty | +3.47% | +0.22% |
| Without uncertainty | -12.21% | -0.74% |

### B. Oil Producing Reservoir

The example is a five-spot pattern reservoir for oil production as shown in Figure 2 [12]. The water wells are located at the four corners of the reservoir and oil production is located in the middle. The reservoir covers an area of $630\times630$ ft$^2$ and has a thickness of 30 ft which is modeled by a $9\times9\times1$ horizontal two dimensional (2D) grid. The set $\Theta_p$ consists of porosity, $\phi$, and permeability, $K$ [12], [13].



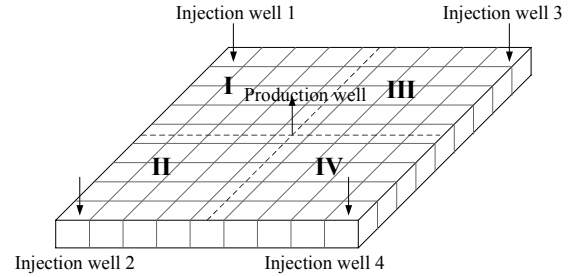Fig. 2.   Schematic of a 2D reservoir. ↓: water injection; ↑: oil production.

Assume the grids on the left (I and II) have the same porosity $\phi_l$ and permeability $K_l$ and the same assumption holds for the grids on the right (III and IV), $\phi_r$ and $K_r$. Additionally, $K$ in the x and y directions are assumed to have the same values. Clearly the value of the properties are uncertain. However, prior knowledge of the reservoir's geologic structure speculates that $\phi_l$ and $K_l$ are about 0.1 and 10 mDarcy, respectively and $\phi_r$ and $K_r$ are about 0.3 and 500 mDarcy.

The LHHS technique can be used to sample the assumed distributions of $\Theta_p = \{\phi, K\}$. The set, $\Theta_p$ are inputs to $M$ to generate $y$ = {oil production, water production and the bottom hole pressures of wells}. Here, $M$ is an ECLIPSE (Schlumberger, Houston, TX) model of the reservoir. PLS can be applied to determine the relationship between $\Theta_p$ and $y$ supplanting the huge number of model executions needed to compare with the observed data. The execution time to generate the PLS results is more than 10 times less than one execution of $M$. The error in the PLS predictions is within 1% of the true values.

Figure 4 describes the prior and posterior probability density of $\phi_l$. The posterior distribution shows that $\bar\phi_l$ and $\sigma(\phi_l)$ are 0.115 and 0.051, respectively. Figure 5 describes the prior and posterior probability density of $K_l$. The posterior distribution's $\bar K_l$ and $\sigma(K_l)$ are 8.8 mDarcy and 0.6, respectively. Figures 6 and 7 give the prior and posterior distributions of $\phi_r$ and $K_r$. The posterior's $\bar K_r, \bar\phi_r$ are 0.32

and 585 mDarcy, respectively; and $\sigma(\phi_r), \sigma(K_r)$ are 0.02 and 40, respectively. Table II lists the root mean square errors (RMS) of the prior and posterior values of $\Theta_p$ that are calculated by Equation (6) ,

$$\text{RMS} = \left( \frac{\sum_{j=1}^{M}(x_j - x_j^t)^2}{p} \right)^{0.5} \tag{6}$$

where $p$ is the number of parameters; $x_j$ is an estimate of the parameters; $x_j^t$ is the true value of the parameters.

TABLE II

ROOT MEAN SQUARE ERRORS

| RMS | $\phi$ | $K$ |
|---|---|---|
| Prior | 0.004 | 11.11 |
| Posterior | 0.0012 | 1.67 |

## 4. SUMMARY

This work desribed new applications for the use of PLS regression combined with other techniques to address model parameter uncertainty. For a high-dimension set of outputs, it was shown that the combination of PLS and Karhunen-Loeve expansion (also known as proper orthogonal decomposition) can identify the relationships between the uncertain parameters and the coefficients of the KL model so that an estimate of the outputs can be determined quuickly thereby avoiding the potentially large computational overhead associated with the simulation of a fundamental model. It also was shown that PLS when combined with Markov Chain Monte Carlo technique can provide fast updates of the parameters to maintain model prediction accuracy.

### REFERENCES

[1] H. Wold, *Estimation of principle components and related models by iterative least squares*, in p.r. krishnaiaah ed. New York: Academic Press, 1966, multivariate Analysis. (pp.391-420).

[2] A. J. Burnham, J. F. MacGregor, and R. Viveros, "A statistical framwwork for multivariate latent variable regression methods based on maximum likelihood," *Journal of Chemometrics*, vol. 13, pp. 49–65, 1999.

[3] H. Abdi, *Partial Least Squares (PLS) Regression*. Thousand Oaks (CA): Sage, 2003, pp. 792–795.

[4] L. Sirovich, *New Perspectives in Turbulence*, 1st ed. New York, NY: Springer-Verlag, 1991.

[5] A. J. Newman, "Model reduction via the karhunen-loève expansion Part I: An exposition technical report t.r. 96-32," University of Maryland, College Park,MD, Tech. Rep., 1996.

[6] Y. Chen and K. Hoo, "Uncertainty propagation for effective reduced-order model generation," *Computers and Chemical Engineering(2010)*, vol. 34, pp. 1597–1605, 2010.

[7] ——, "Uncertainty propagation for efficient model-based control solutions," in *American Control Conference*, Baltimore, Maryland, June & July 2010, pp. 3112–3117, thA22.1.

[8] D. R. Brouwer and J. D. Jansen, "Dynamic optimization of water flooding with smart wells using optimal control theory," in *Proc. 13th European Petroleum Conference*. Society of Petoleum Engineers, 2002, aberdeen, Scotland, U.K., SPE 78278.

[9] H. Haario, E. Saksman, and J. Tamminen, "An adaptive metropolis algorithm," *Bernoulli*, vol. 7, pp. 223–242, 2001.

[10] L. Marshall, D. Nott, and A. Sharma, "A comparative study of Markov chain Mote Carlo methods for conceptual rainfall-runoff modeling," *Water Resources Research*, vol. 40, 2004, w02501, doi:10.1029/2003WR002378.

[11] D. Zheng and K. A. Hoo, "System identification and model-based for distributed parameter systems," *Computers and Chemical Engineering*, vol. 28, pp. 1361–1375, 2004.

[12] Y. Chen and K. Hoo, "Optimal model-based reservoir management with model parameter uncertainty updates." Hangzhou, China: Subbmitted to The 4th International Symposium on Advanced Control of Industrial Processes, May 2011.

[13] ——, "Optimization of reservoirs with model parameter uncertainty updates." Salt Lake City, UT: AIChE, Nov 2010.
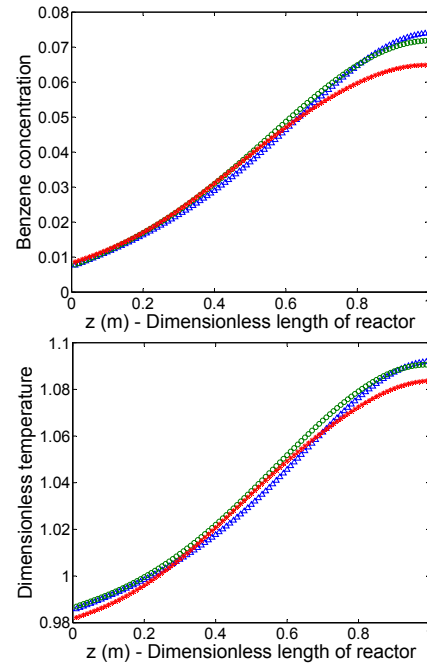
Fig. 3. Output of the ROM and the physical model in the presence of a 5% bias in the mean value of the uncertain parameters [6]. Top: Benzene concentration. Bottom: Reactor temperature. △: Physical model. ○: ROM with uncertainty. +: ROM without uncertainty.
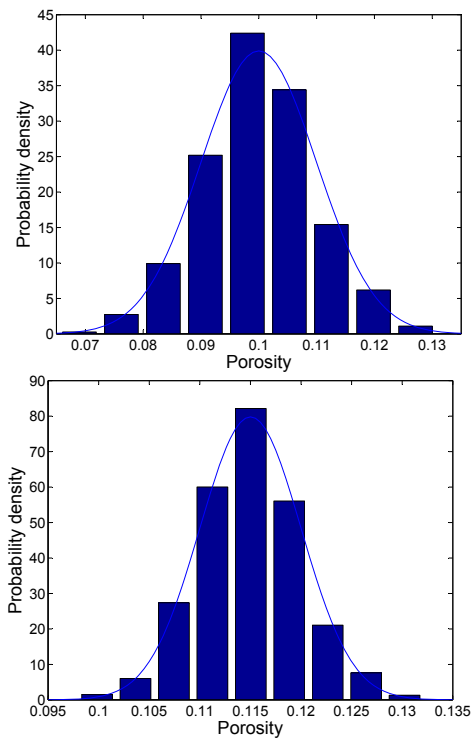
Fig. 4. Distribution of porosity in the left part of the reservoir. Top: prior probability density. Bottom: posterior probability density.
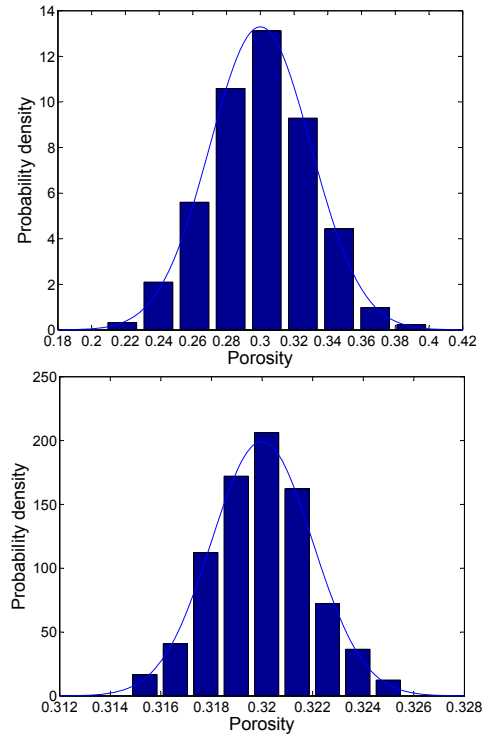


Fig. 6. Distribution of porosity in the right part of the reservoir. Top: prior probability density. Bottom: posterior probability density.
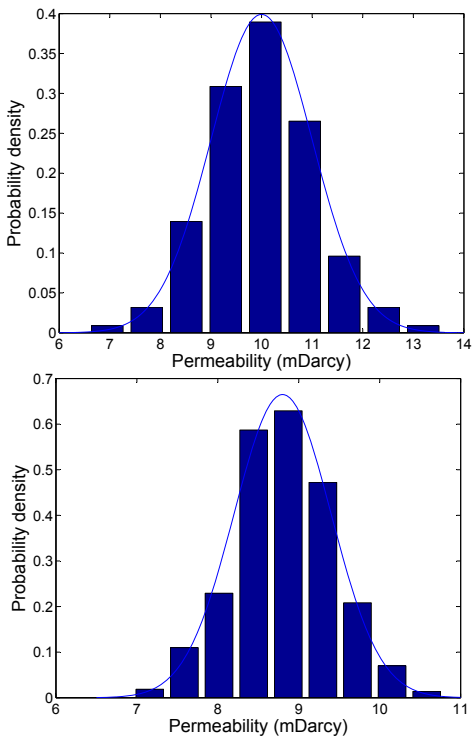


Fig. 5. Distribution of permeability in the left part of the reservoir. Top: prior probability density. Bottom: posterior probability density.
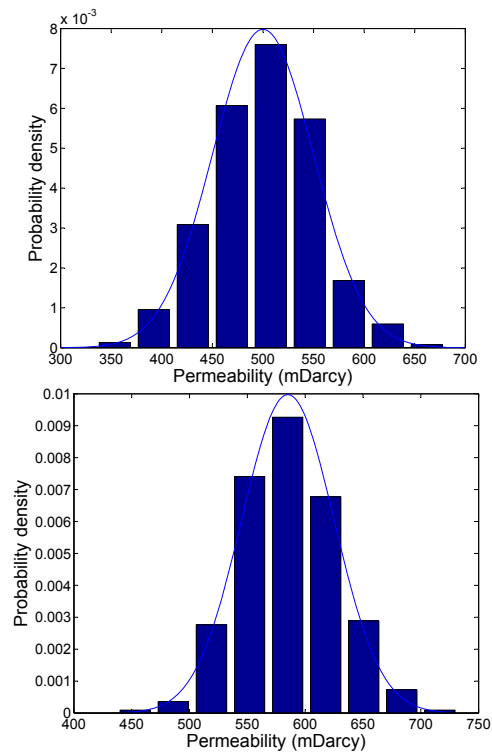


Fig. 7. Distribution of permeability in the right part of the reservoir. Top: prior probability density. Bottom: posterior probability density.