

Predictions in Hierarchical Bayesian Models

Georgios Arampatzis

Professorship for Computational Science, ETH-Zurich, CH-8092, Switzerland

September 23, 2018

Let Y_i, Θ_i for $i = 1, \dots, N$ and Ψ be random variables with $Y_i \in \mathbb{R}^{N_{Y_i}}$, $\Theta_i \in \mathbb{R}^{N_{\Theta}}$ and $\Psi \in \mathbb{R}^{N_{\Psi}}$. The dependency of the random variables is described by the Directed Acyclic Graph (DAG) given in Figure 1. The joint probability distribution is given by

$$p(\mathbf{d}, \boldsymbol{\vartheta}, \psi) = \prod_{i=1}^N p(d_i | \vartheta_i) p(\vartheta_i | \psi) p(\psi) = p(\mathbf{d} | \boldsymbol{\vartheta}) p(\boldsymbol{\vartheta} | \psi) p(\psi) = p(\mathbf{d} | \boldsymbol{\vartheta}) p(\boldsymbol{\vartheta}, \psi), \quad (1)$$

with $\mathbf{d} = (d_1, \dots, d_N)$ and $\boldsymbol{\vartheta} = (\vartheta_1, \dots, \vartheta_N)$. In Bayesian inverse problems this graph describes the output of a model, f , that depends on parameters ϑ , under the hypothesis of random errors. The noise is usually modelled and the function $p(y_i | \vartheta_i)$ is fully determined. A usual assumption is the following,

$$d_i = f(\vartheta_i) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \Sigma_i), \quad (2)$$

with $\Sigma_i := \Sigma(\vartheta_i)$ and thus $p(d_i | \vartheta_i) = \mathcal{N}(d_i | f(\vartheta_i), \Sigma(\vartheta_i))$. The probability distributions $p(\vartheta_i | \psi)$ and $p(\psi)$ are modelled and are problem dependent.

The Bayesian inversion problem seeks answers for the following questions: After observing some data for the random variables Y_1, \dots, Y_N , i.e., we know that $\mathbf{Y} = \mathbf{d} = \{d_1, \dots, d_N\}$,

1. What is the probability of the parameters ϑ_i and ψ conditioned on the observations \mathbf{d} , $p(\boldsymbol{\vartheta}, \psi | \mathbf{d})$? What is the probability of ϑ^{new} conditioned on the data \mathbf{d} ?
2. What is the probability distribution of a new random variable y^{new} conditioned to the data \mathbf{d} and prior to observing any data for y^{new} , $p(y^{\text{new}} | \mathbf{d})$?
3. What is the probability distribution of a new random variable y_i^{new} , that depends on ϑ_i , conditioned on the data \mathbf{d} , $p(y_i^{\text{new}} | \mathbf{d})$?

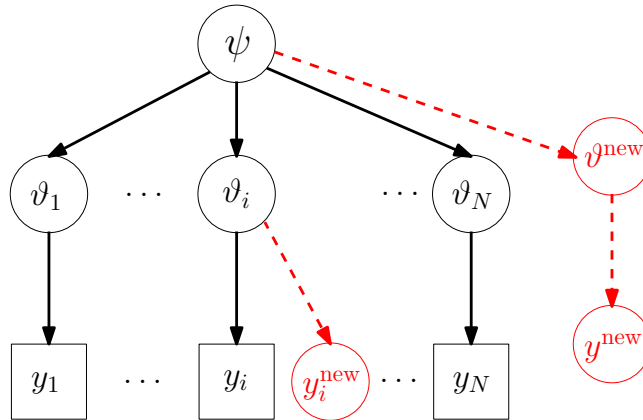


Figure 1: Directed Acyclic Graph (DAG) showing the dependency of the random variables y_i, ϑ_i and ψ . The circles show unobserved and the squares observed random variables. The variables ϑ_i and ϑ^{new} are being inferred after observing the variables y_i . Then, the probability of y_i^{new} and y^{new} conditioned on the observations \mathbf{d} can be quantified.

1 Conditional probability of $\boldsymbol{\vartheta}$ and ψ

The joint distribution approach. We ask to draw samples from the distribution $p(\boldsymbol{\vartheta}, \psi | \mathbf{d})$. Using Bayes theorem,

$$p(\boldsymbol{\vartheta}, \psi | \mathbf{d}) = \frac{p(\mathbf{d} | \boldsymbol{\vartheta}, \psi) p(\boldsymbol{\vartheta}, \psi)}{p(\mathbf{d})} = \frac{p(\mathbf{d} | \boldsymbol{\vartheta}) p(\boldsymbol{\vartheta} | \psi) p(\psi)}{p(\mathbf{d})}, \quad (3)$$

where by (1) the likelihood function is simplified as $p(\mathbf{d} | \boldsymbol{\vartheta}) = \prod_{i=1}^N p(d_i | \vartheta_i)$. The problem with the sampling if this distribution is the dimension of the state space of the target probability which is equal to $NN_{\boldsymbol{\vartheta}} + N_{\psi}$. Unless this number is small the sampling of (3) is highly inefficient. This approach was considered in [1] in low dimensional problems. The conditional probability of ϑ_i is given by marginalizing (3),

$$p(\vartheta_i | \mathbf{d}) = \int p(\boldsymbol{\vartheta}, \psi | \mathbf{d}) d\psi d\boldsymbol{\vartheta}_{\setminus i}, \quad (4)$$

where $\boldsymbol{\vartheta}_{\setminus i} = (\vartheta_1, \dots, \vartheta_{i-1}, \vartheta_{i+1}, \dots, \vartheta_{N_{\boldsymbol{\vartheta}}})$ and the conditional distribution of ϑ^{new} ,

$$\begin{aligned} p(\vartheta^{\text{new}} | \mathbf{d}) &= \int p(\vartheta^{\text{new}}, \psi | \mathbf{d}) d\psi \\ &= \int p(\vartheta^{\text{new}} | \psi, \mathbf{d}) p(\psi | \mathbf{d}) d\psi = \int p(\vartheta^{\text{new}} | \psi) p(\psi | \mathbf{d}) d\psi \\ &\approx \frac{1}{N_s} \sum_{k=1}^{N_s} p(\vartheta^{\text{new}} | \psi^{(k)}), \quad \psi^{(k)} \sim p(\psi^{(k)} | \mathbf{d}), \end{aligned} \quad (5)$$

and $p(\psi | \mathbf{d})$ is given by marginalizing (3) over $\boldsymbol{\vartheta}$.

The two step approach. In this approach the conditional distribution of ψ is sampled first,

$$p(\psi | \mathbf{d}) = \frac{p(\mathbf{d} | \psi) p(\psi)}{p(\mathbf{d})}, \quad (6)$$

where

$$\begin{aligned} p(\mathbf{d} | \psi) &= \prod_{i=1}^N p(d_i | \psi) \\ &= \prod_{i=1}^N \int p(d_i | \vartheta_i) p(\vartheta_i | \psi) d\vartheta_i \\ &\approx \prod_{i=1}^N \frac{1}{N_s} \sum_{k=1}^{N_s} p(d_i | \vartheta_i^{(k)}), \quad \vartheta_i^{(k)} \sim p(\vartheta_i^{(k)} | \psi). \end{aligned} \quad (7)$$

Notice that one evaluation of $p(\mathbf{d} | \psi)$ involves NN_s evaluations of the likelihood $p(y | \vartheta)$ and thus NN_s evaluations of the model. In order to alleviate the huge computational cost an importance sampling scheme was proposed in [Wu2016] with $p(\vartheta_i | d_i, \mathcal{M}_i)$ as instrumental density in each of the integrals in (7). Here, $p(\vartheta_i | d_i, \mathcal{M}_i)$ is the posterior distribution of ϑ_i after observing the data d_i and by considering d_i and ϑ_i independent of ψ ,

$$p(\vartheta_i | d_i, \mathcal{M}_i) = \frac{p(d_i | \vartheta_i, \mathcal{M}_i) p(\vartheta_i | \mathcal{M}_i)}{p(d_i | \mathcal{M}_i)}, \quad (8)$$

where the likelihood function $p(d_i | \vartheta_i, \mathcal{M}_i)$ is chosen equal to the likelihood function $p(d_i | \vartheta_i)$. Finally, the likelihood (7) is written as

$$p(\mathbf{d} | \psi) = \prod_{i=1}^N p(d_i | \psi), \quad (9)$$

with

$$\begin{aligned} p(d_i | \psi) &= \int p(d_i | \vartheta_i) p(\vartheta_i | \psi) d\vartheta_i \\ &= \int \frac{p(d_i | \vartheta_i) p(\vartheta_i | \psi)}{p(\vartheta_i | d_i, \mathcal{M}_i)} p(\vartheta_i | d_i, \mathcal{M}_i) d\vartheta_i \\ &= \int \frac{p(d_i | \vartheta_i) p(\vartheta_i | \psi) p(d_i | \mathcal{M}_i)}{p(d_i | \vartheta_i, \mathcal{M}_i) p(\vartheta_i | \mathcal{M}_i)} p(\vartheta_i | d_i, \mathcal{M}_i) d\vartheta_i \\ &= \int \frac{p(\vartheta_i | \psi) p(d_i | \mathcal{M}_i)}{p(\vartheta_i | \mathcal{M}_i)} p(\vartheta_i | d_i, \mathcal{M}_i) d\vartheta_i, \end{aligned} \quad (10)$$

and approximated by

$$p(\mathbf{d}|\psi) \approx \prod_{i=1}^N \frac{p(d_i|\mathcal{M}_i)}{N_s} \sum_{k=1}^{N_s} \frac{p(\vartheta_i^{(k)}|\psi)}{p(\vartheta_i^{(k)}|\mathcal{M}_i)}, \quad \vartheta_i^{(k)} \sim p(\vartheta_i^{(k)}|d_i, \mathcal{M}_i). \quad (11)$$

With this scheme the model function is evaluated NN_s times and the $\vartheta_i^{(k)}$ samples are stored. Then, the evaluation of $p(\mathbf{d}|\psi)$ involves no evaluations of the likelihood function and thus the computationally expensive model function.

The posterior samples for a new parameter ϑ^{new} can be obtained by,

$$\begin{aligned} p(\vartheta^{\text{new}}|\mathbf{d}) &= \int p(\vartheta^{\text{new}}, \psi|\mathbf{d}) d\psi \\ &= \int p(\vartheta^{\text{new}}|\mathbf{d}, \psi) p(\psi|\mathbf{d}) d\psi \\ &= \int p(\vartheta^{\text{new}}|\psi) p(\psi|\mathbf{d}) d\psi \\ &\approx \frac{1}{N_s^\psi} \sum_{k=1}^{N_s^\psi} p(\vartheta^{\text{new}}|\psi^{(k)}), \quad \psi^{(k)} \sim p(\psi|\mathbf{d}). \end{aligned} \quad (12)$$

Notice that (12) is easy to be sampled once $p(\psi|\mathbf{d})$ is sampled, since $p(\vartheta^{\text{new}}|\psi^{(k)})$ is a known function.

The posterior samples for the parameter ϑ_i can be obtained by,

$$\begin{aligned} p(\vartheta_i|\mathbf{d}) &= \int p(\vartheta_i, \psi|\mathbf{d}) d\psi \\ &= \int p(\vartheta_i|\mathbf{d}, \psi) p(\psi|\mathbf{d}) d\psi \\ &= \int p(\vartheta_i|d_i, \psi) p(\psi|\mathbf{d}) d\psi \\ &= p(d_i|\vartheta_i) \int \frac{p(\vartheta_i|\psi)}{p(d_i|\psi)} p(\psi|\mathbf{d}) d\psi \\ &\approx \frac{p(d_i|\vartheta_i)}{N_s^\psi} \sum_{k=1}^{N_s^\psi} \frac{p(\vartheta_i|\psi^{(k)})}{p(d_i|\psi^{(k)})}, \quad \psi^{(k)} \sim p(\psi|\mathbf{d}). \end{aligned} \quad (13)$$

Notice that the term $p(d_i|\psi^{(k)})$ can be precomputed using (10) since it is independent of ϑ_i ,

$$\begin{aligned} p(d_i|\psi^{(k)}) &\approx p(d_i|\mathcal{M}_i) \frac{1}{N_s^\vartheta} \sum_{\ell=1}^{N_s^\vartheta} \frac{p(\vartheta_i^{(\ell)}|\psi^{(k)})}{p(\vartheta_i^{(\ell)}|\mathcal{M}_i)} \\ &:= p(d_i|\mathcal{M}_i) b_{i,k}. \end{aligned} \quad (14)$$

Finally, the log-posterior is given by,

$$\log p(\vartheta_i|\mathbf{d}) \approx \log p(d_i|\vartheta_i) - \log p(d_i|\mathcal{M}_i) + \log \frac{1}{N_s^\psi} \sum_{k=1}^{N_s^\psi} \frac{p(\vartheta_i|\psi^{(k)})}{b_{i,k}}. \quad (15)$$

2 Propagation of Uncertainty in y

The conditional distribution of the output of the model on the observations \mathbf{d} is given by,

$$\begin{aligned} p(y^{\text{new}}|\mathbf{d}) &= \int p(y^{\text{new}}|\vartheta^{\text{new}}, \mathbf{d}) d\vartheta^{\text{new}} \\ &= \int p(y^{\text{new}}|\vartheta^{\text{new}}) p(\vartheta^{\text{new}}|\mathbf{d}) d\vartheta^{\text{new}} \\ &\approx \frac{1}{N_s} \sum_{k=1}^{N_s} p(y^{\text{new}}|\vartheta^{\text{new}(k)}), \quad \vartheta^{\text{new}(k)} \sim p(\vartheta^{\text{new}(k)}|\mathbf{d}), \end{aligned} \quad (16)$$

and $p(\vartheta^{\text{new}(k)}|\mathbf{d})$ can be sampled by (12). Under the assumption that $p(y|\vartheta) = \mathcal{N}(y|f(\vartheta), \Sigma(\vartheta))$ the distribution $p(y^{\text{new}}|\mathbf{d})$ is a Gaussian mixture given by,

$$p(y^{\text{new}}|\mathbf{d}) = \sum_{k=1}^{N_s} \frac{1}{N_s} \mathcal{N}(y^{\text{new}}|f(\vartheta^{(k)}), \Sigma(\vartheta^{(k)})) . \quad (17)$$

The posterior distribution of a new observation on the i -th branch of the DAG in Figure 1, conditioned on the observations \mathbf{d} , is given by

$$\begin{aligned} p(y_i^{\text{new}}|\mathbf{d}) &= \int p(y_i^{\text{new}}, \vartheta_i|\mathbf{d}) d\vartheta_i \\ &= \int p(y_i^{\text{new}}|\vartheta_i) p(\vartheta_i|\mathbf{d}) d\vartheta_i \\ &\approx \frac{1}{N_s} \sum_{k=1}^{N_s} p(y_i^{\text{new}}|\vartheta_i^{(k)}), \quad \vartheta_i^{(k)} \sim p(\vartheta_i^{(k)}|\mathbf{d}) , \end{aligned} \quad (18)$$

where ϑ_i can be sampled by (15). In the special case that $p(y|\vartheta) = \mathcal{N}(y|f(\vartheta), \Sigma(\vartheta))$ the distribution $p(y^{\text{new}}|\mathbf{d})$ is a Gaussian mixture similar to (17).

3 Optimization in Bayesian Networks

In Mixed Effects Models (MEM) the likelihood of the data conditioned on the hyper-parameters (6) is optimized on ψ ,

$$\psi^* = \arg \max_{\psi} p(\mathbf{d}|\psi) . \quad (19)$$

Then, the probability distributions of the parameters ϑ_i conditioned on the data are given by (15) for $p(\psi|\mathbf{d}) = \delta(\psi - \psi^*)$,

$$p(\vartheta_i|\mathbf{d}) = \frac{p(d_i|\vartheta_i) p(\vartheta_i|\psi^*)}{p(d_i|\psi^*)} . \quad (20)$$

Finally, the posterior distributions for the predictions of y^{new} and y_i^{new} are given by (16) and (18), respectively.

3.1 Algorithms for Optimization in Bayesian Networks

Expectation Maximization. It can be shown that the following iterative procedure maximizes the likelihood function (19) starting from $\psi^{(0)}$:

- Expectation step: Evaluate the function defined by,

$$\begin{aligned} Q(\psi|\psi^{(k)}) &= \mathbb{E}_{\boldsymbol{\vartheta}|\mathbf{d}, \psi^{(k)}} [\log p(\mathbf{d}, \boldsymbol{\vartheta}|\psi)] \\ &= \int \log(p(\mathbf{d}, \boldsymbol{\vartheta}|\psi)) p(\boldsymbol{\vartheta}|\mathbf{d}, \psi^{(k)}) d\boldsymbol{\vartheta} \\ &= \int \log(p(\mathbf{d}|\boldsymbol{\vartheta}) p(\boldsymbol{\vartheta}|\psi)) p(\boldsymbol{\vartheta}|\mathbf{d}, \psi^{(k)}) d\boldsymbol{\vartheta} \\ &= \prod_{i=1}^N \int \log(p(d_i|\vartheta_i) p(\vartheta_i|\psi)) p(\vartheta_i|d_i, \psi^{(k)}) d\vartheta_i \\ &\approx \prod_{i=1}^N \frac{1}{N_s} \sum_{\ell=1}^{N_s} \log p(d_i|\vartheta_i^{(\ell)}) + \log p(\vartheta_i^{(\ell)}|\psi), \quad \vartheta_i^{(\ell)} \sim p(\vartheta_i^{(\ell)}|d_i, \psi^{(k)}) , \end{aligned} \quad (21)$$

where $p(\vartheta_i^{(\ell)}|d_i, \psi^{(k)})$ can be sampled as in (20).

- Maximization step: Maximize the function in (21),

$$\psi^{(k+1)} = \arg \max_{\psi} Q(\psi|\psi^{(k)}) . \quad (22)$$

Importance Sampling. The likelihood function in (19) can be approximated by the importance sampling scheme described in Section 1 and summarized in equation (11).

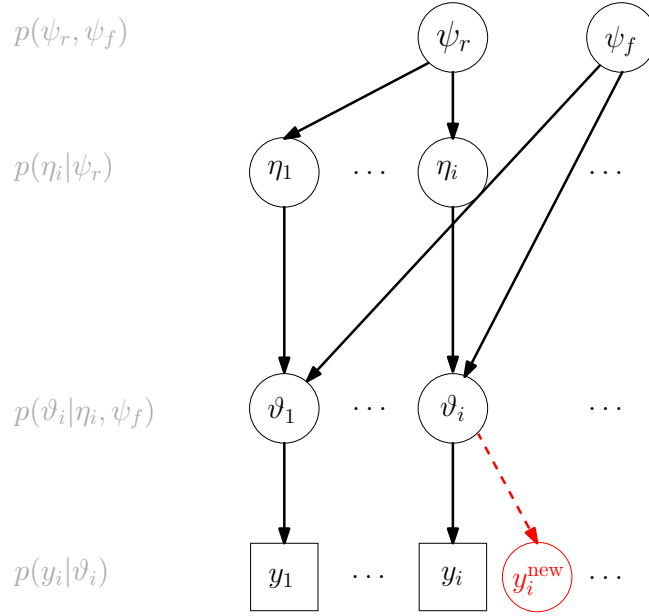


Figure 2: Directed Acyclic Graph (DAG) showing the dependency of the random variables in a two layer hierarchical model.

4 A two layer hierarchical model

Working in the same way as in (9) and (10)

$$p(\mathbf{d}|\psi_r, \psi_f) = \prod_{i=1}^N \int p(d_i|\vartheta_i) p(\vartheta_i|\psi_r, \psi_f) d\vartheta_i, \quad (23)$$

where

$$\begin{aligned} p(\vartheta_i|\psi_r, \psi_f) &= \int p(\vartheta_i|\eta_i, \psi_f) p(\eta_i|\psi_r) d\eta_i \\ &\approx \frac{1}{N_s} \sum_{k=1}^{N_s} p(\vartheta_i|\eta_i^{(k)}, \psi_f), \quad \eta_i^{(k)} \sim p(\eta_i|\psi_r). \end{aligned} \quad (24)$$

Let $\psi = (\psi_r, \psi_f)$.

Similarly with eq. (15), the posterior distribution of parameters ϑ_i is given by,

$$\begin{aligned} p(\vartheta_i|\mathbf{d}) &= \int p(\vartheta_i|\psi, \mathbf{d}) p(\psi|\mathbf{d}) d\psi \\ &= p(d_i|\vartheta_i) \int \frac{p(\vartheta_i|\psi)}{p(d_i|\psi)} p(\psi|\mathbf{d}) d\psi \\ &\approx p(d_i|\vartheta_i) \frac{1}{N_s} \sum_{k=1}^{N_s} \frac{p(\vartheta_i|\psi^{(k)})}{p(d_i|\psi^{(k)})}, \quad \psi^{(k)} \sim p(\psi|\mathbf{d}) \\ &\approx p(d_i|\vartheta_i) \frac{1}{N_s} \sum_{k=1}^{N_s} \frac{1}{p(d_i|\psi^{(k)})} \frac{1}{N_s} \sum_{\ell=1}^{N_s} p(\vartheta_i|\eta_i^{(\ell)} \psi_f^{(k)}), \quad \eta_i^{(\ell)} \sim p(\eta_i|\psi_r^{(k)}), \quad \psi^{(k)} \sim p(\psi|\mathbf{d}). \end{aligned} \quad (25)$$

The joint posterior distribution of parameter (η_i, ϑ_i) is given by,

$$\begin{aligned} p(\eta_i, \vartheta_i|\mathbf{d}) &= \int p(\eta_i, \vartheta_i|\psi, \mathbf{d}) p(\psi|\mathbf{d}) d\psi \\ &= \int p(\eta_i, \vartheta_i|\psi, d_i) p(\psi|\mathbf{d}) d\psi \\ &= \int \frac{p(d_i|\eta_i, \vartheta_i, \psi) p(\eta_i, \vartheta_i|\psi)}{p(d_i|\psi)} p(\psi|\mathbf{d}) d\psi \\ &= \int \frac{p(d_i|\vartheta_i) p(\eta_i, \vartheta_i|\psi)}{p(d_i|\psi)} p(\psi|\mathbf{d}) d\psi. \end{aligned} \quad (26)$$

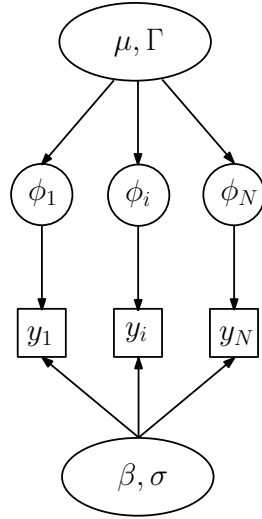


Figure 3: Directed Acyclic Graph (DAG) showing the dependency of the random variables in a two layer hierarchical model.

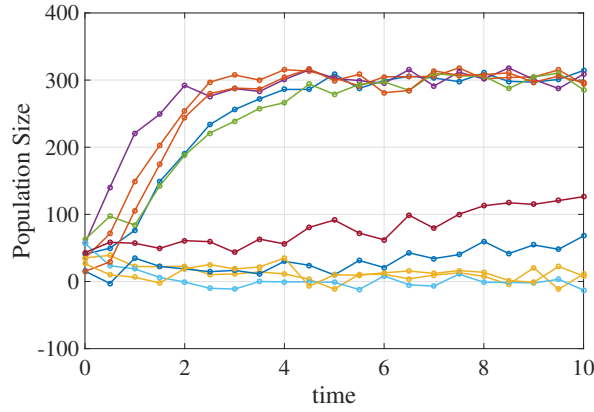


Figure 4: Independent synthetic data sets from equation (29) using the parameters in (31) .

By noticing that

$$p(\eta_i, \vartheta_i | \psi) = p(\vartheta_i | \eta_i, \psi_f) p(\eta_i | \psi_r) , \quad (27)$$

the posterior distribution (26) can be written as,

$$p(\eta_i, \vartheta_i | \mathbf{d}) = p(d_i | \vartheta_i) \int \frac{p(\vartheta_i | \eta_i, \psi_f) p(\eta_i | \psi_r)}{p(d_i | \psi)} p(\psi | \mathbf{d}) d\psi . \quad (28)$$

5 Mixed Effects Model

6 Numerical example

We consider the following function that models the growth of population

$$f(t; \vartheta) = K P_0 \frac{e^{\lambda t}}{K + P_0(e^{\lambda t} - 1)} , \quad (29)$$

where $\vartheta = (K, P_0, \lambda)$ and K is the capacity, P_0 is the initial population and λ the growth rate. We create synthetic data at times $t_k = 0.5k$, for $k = 1, \dots, 20$,

$$d_k = f(t_k; \vartheta) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_d^2) , \quad (30)$$

using the following nominal values for the parameters,

$$K = 300, \quad P_0 \sim \mathcal{N}(40, 10^2), \quad \lambda \sim \mathcal{N}(0.5, 1), \quad \sigma_d = 10 , \quad (31)$$

where $\mathcal{N}(\mu, \sigma^2)$ is the normal distribution with mean μ and standard deviation σ . See Figure 4 for an example of 10 data sets. From now on we consider the extended parameter vector $\vartheta = (K, P_0, \lambda, \sigma_d)$ that contains the standard deviation of the noise.

References

- [1] J. B. Nagel and B. Sudret. Hamiltonian Monte Carlo and borrowing strength in hierarchical inverse problems. *ASCE-ASME Journal of Risk and Uncertainty in Engineering Systems, Part A: Civil Engineering*, 0(0):B4015008, 2015.
- [2] J. Stirnemann, A. Samson, and J.-C. Thalabard. Individual predictions based on nonlinear mixed modeling: application to prenatal twin growth. *Statistics in Medicine*, 31(18):1986–1999, 2012.