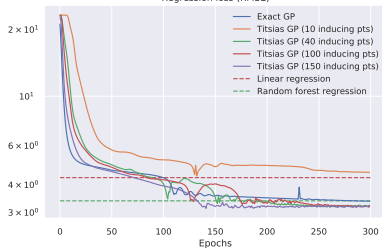
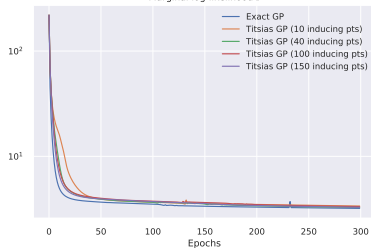


Marginal log-likelihood ℓ

Matérn kernel $\nu = 1.5$

Regression loss (RMSE)



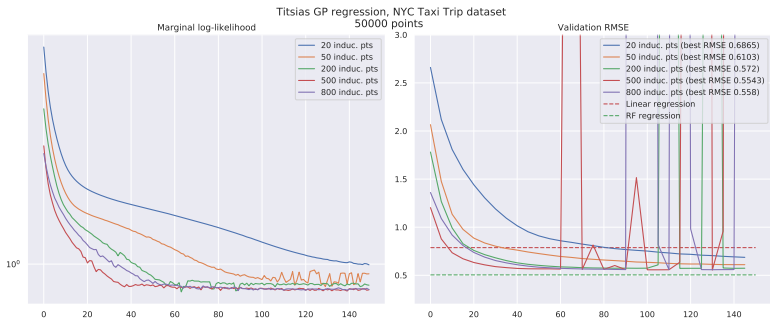


Figure: Comparison on the NYC Taxi dataset.

Deep Neural Network as Gaussian Processes

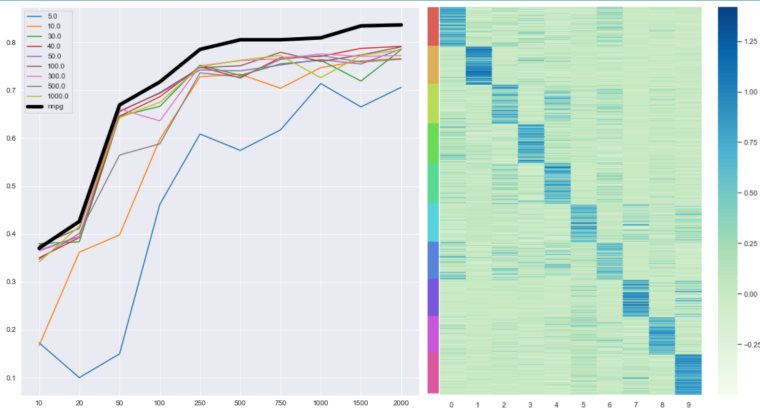
Review of Lee et al. article

Delanoue Pierre, Guillo Clement

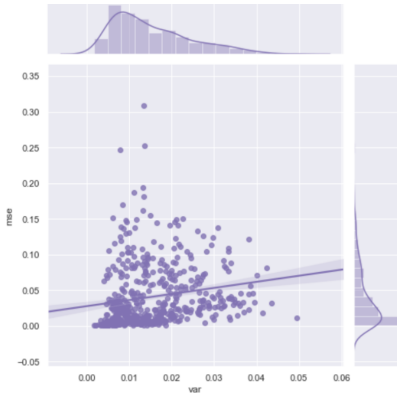
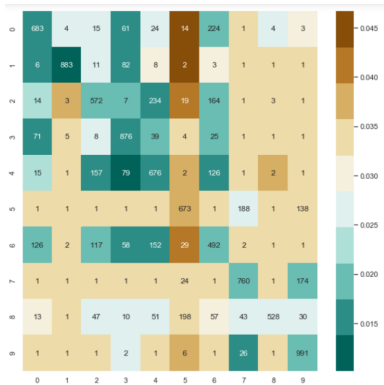
ENS Paris Saclay, MVA Master

Mars 13, 2020

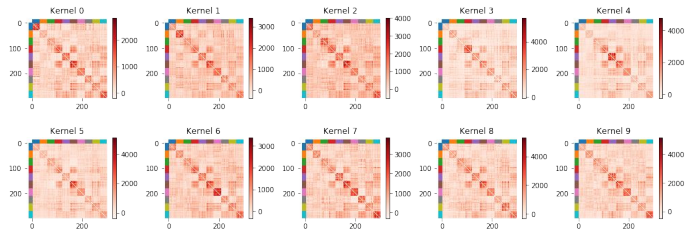
Accuracy and Posterior Mean



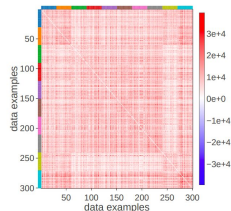
Posterior Variance



Approximate Inference Turns Deep Networks Into Gaussian Processes

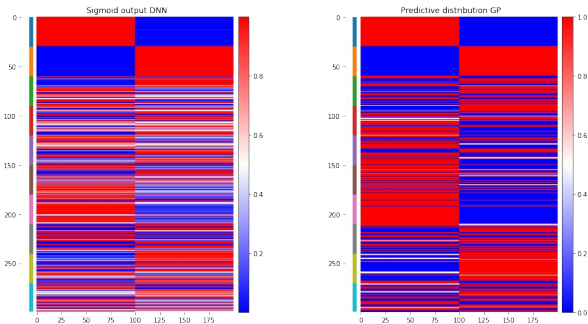


Class-specific kernels trained on MNIST.

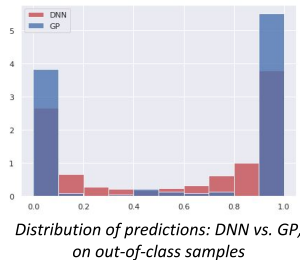


Kernel trained on CIFAR-10

DNN vs. GP in the face of overconfidence

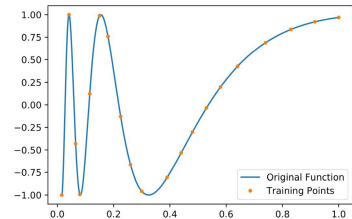


Predictions of DNN vs GP, trained on binary MNIST, on out-of-class samples

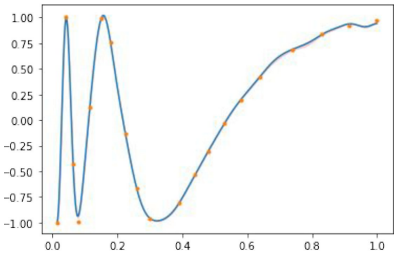


Supervised learning

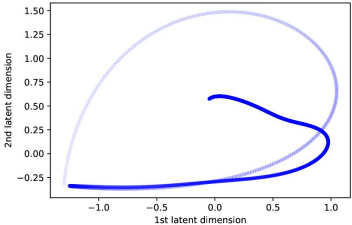
Original function



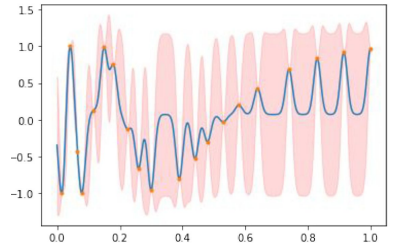
DGP



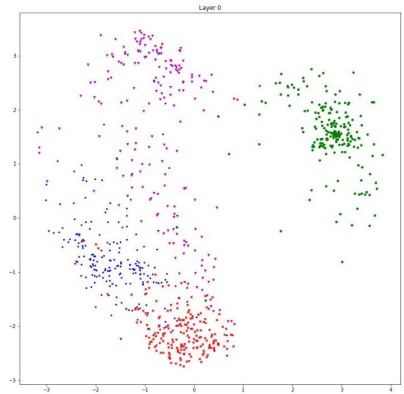
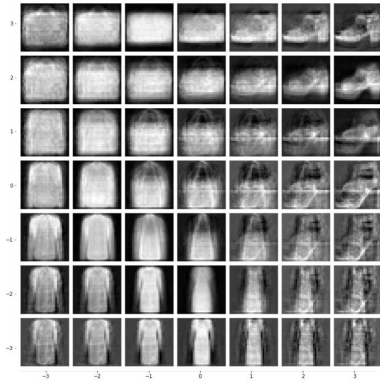
Latent space



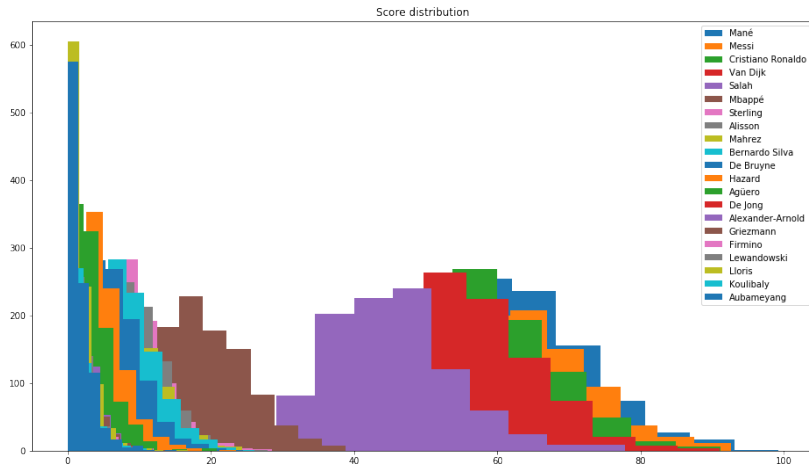
GP



Unsupervised learning



Player score distribution



Film time-dependant scores

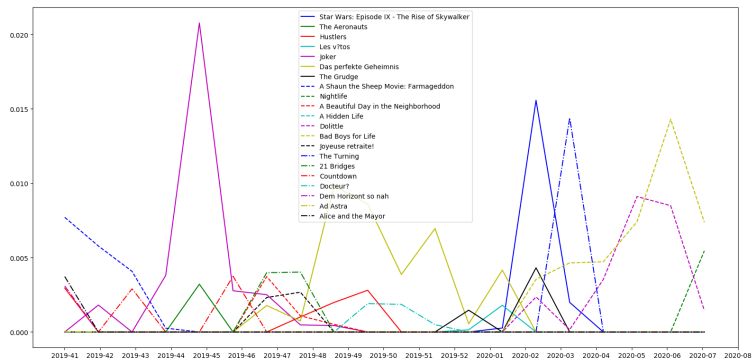


Figure: Evolution of the weights $w_{t,k}$

KPLS: from non-Bayesian to Bayesian

Thibault LAHIRE and Samuel GRUFFAZ

Presentation

Friday, 13 March 2020

KPLS (Kriging using Partial Least Squares)

Reference article:

- Title: Improving Kriging surrogates of high-dimensional design models by Partial Least Squares dimension reduction
- Authors: Mohamed Amine Bouhlel, Nathalie Bartoli, Abdelkader Otsmane, Joseph Morlier
- Date: 2015

→ Main contribution = use of PLS to reduce the number of hyper-parameters

→ Underlying theory is fully **frequentist**

Experiments

Introducing $\Delta = (\mathbf{y} - \mathbf{F}\boldsymbol{\beta})^\top \mathbf{R}^{-1} (\mathbf{y} - \mathbf{F}\boldsymbol{\beta})$, we have:

$$\hat{\sigma}_{MP}^2 = \frac{2\delta + \Delta}{2(\alpha - 1) + n} \quad \text{and} \quad \hat{\sigma}_{MAP}^2 = \frac{2\delta + \Delta}{2(\alpha + 1) + n}$$

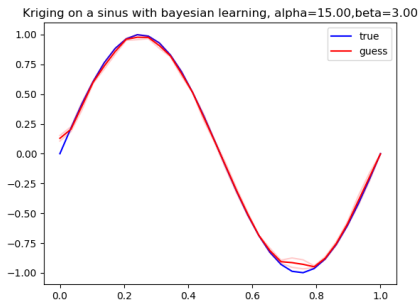
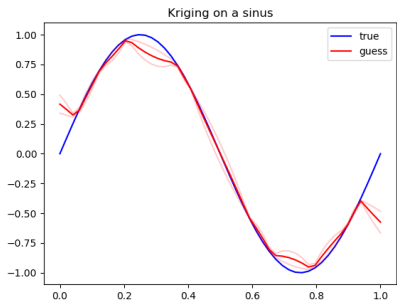


Figure: Example where $n = 13$ points

Data and Training

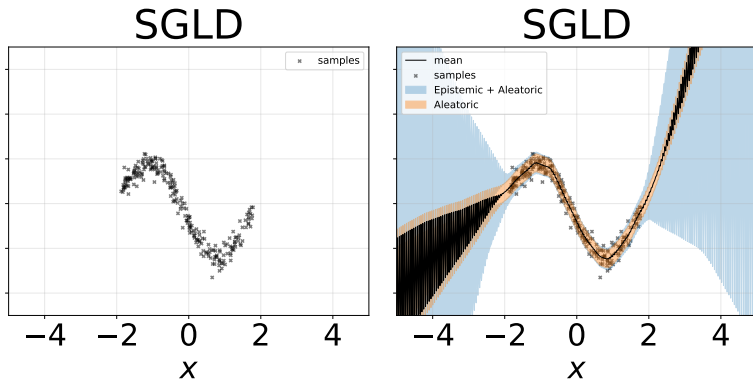


Figure: On the right: The data to be learned / On the left: The learned distribution of the data

Visualizing ϵ_c

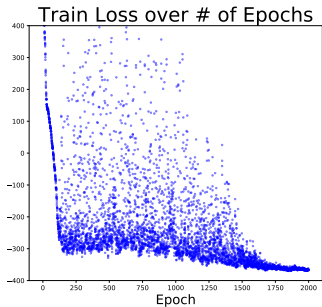
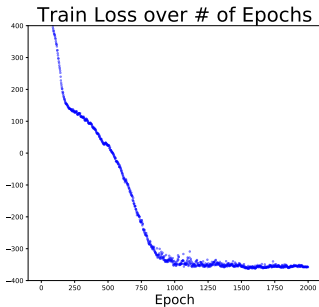


Figure: Comparing the training with a singleton batch or five batches

Latent Dirichlet Allocation graphical model

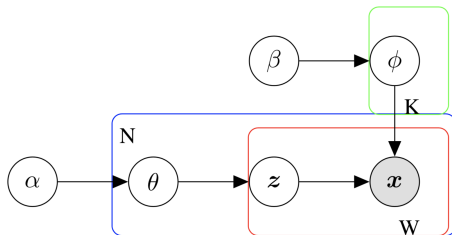


Figure: LDA graphical model

$$\begin{aligned}\theta_j &\sim \text{Dirichlet}(\alpha) \\ \theta_j &\sim \text{Dirichlet}(\alpha) \\ z_{ij} &\sim \text{Multinomial}(\theta_j) \\ x_{ij} &\sim \text{Multinomial}(\phi_{z_{ij}})\end{aligned}$$

$$\begin{aligned}1 &\leq j \leq N \\ 1 &\leq k \leq K \\ 1 &\leq i \leq W\end{aligned}$$

Experiments: CVB-LDA script on Sequoia dataset

COLLAPSED VARIATIONAL INFERENCE LDA

SKLEARN LDA COMPARISON

Collapsed variational inference LDA exec time: 0.38117408 Sklearn LDA exec time: 6.565928936s

Topics found via collapsed variational inference LDA:

Topics found via Sklearn LDA:

Topic 1: 10 most important words, with $p(w|z)$:

affaire, 0.8127481%
paris, 0.4366119%
president, 0.4365667%
deux, 0.3874088%
etait, 0.3427752%
commission, 0.3388877%
ans, 0.3265192%
dont, 0.2801829%
juge, 0.2785372%
ancien, 0.2730116%

Topic 2: 10 most important words, with $p(w|z)$:

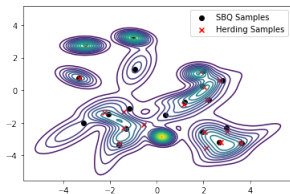
patients, 1.5233556%
aclasta, 1.3945385%
angiox, 0.71278%
bivalirudine, 0.6824946%
perfusion, 0.6370641%
traitement, 0.637059%
mg, 0.5916037%
doit, 0.4835385%
effets, 0.4552107%
voir, 0.4395397%

Topic 1: 10 most important words, with $p(w|z)$:

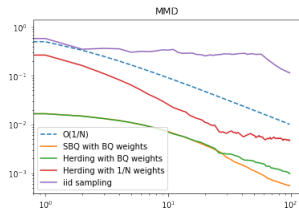
rrb, 1.9544326%
lrb, 1.9532899%
patients, 0.9570906%
aclasta, 0.8762633%
angiox, 0.4479074%
bivalirudine, 0.4297532%
traitement, 0.4013508%
perfusion, 0.4008852%
mg, 0.3725444%
doit, 0.3497594%

Topic 2: 10 most important words, with $p(w|z)$:

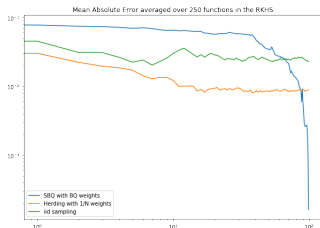
affaire, 0.6752488%
commission, 0.4087772%
president, 0.3725456%
paris, 0.3635402%
jean, 0.3146196%
conseil, 0.2580898%
2006, 0.2577587%
solution, 0.2411031%
etait, 0.2382889%
juge, 0.2326179%



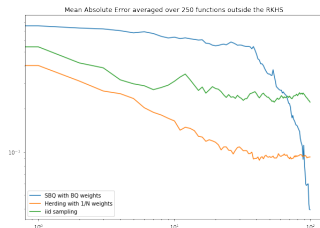
(a) 20 first samples of SBQ and Herding



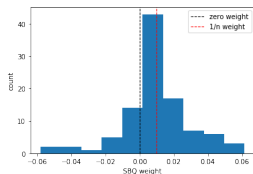
(b) MMD



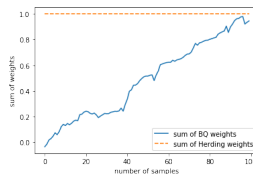
(c) Mean Absolute Error for 250 random functions within the RKHS



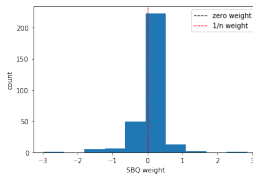
(d) Mean Absolute Error for 250 random functions outside of the RKHS



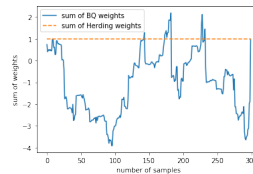
(e) Distribution of weights for 100 samples



(f) Sum of the weights for 100 samples



(g) Distribution of weights for 300 samples



(h) Sum of the weights for 300 samples

Plan of the Presentation

- Introduction
- Origin of the model
 - Stochastic Gradient Descent (SGD)
 - Stochastic Weight Averaging (SWA)
 - Assumptions
- Presentation of SWAG
 - Theoretical presentation of SWAG
 - SWAG Diagonal
 - SWAG Low Rank + Diagonal Structure
- Experimental Results
- Conclusion

A Simple Baseline for Bayesian Uncertainty

Algorithm 1 Continuous learning of SWAG model

θ_{pre} : pretrained weights ; η : learning rate ; T : number of steps ; c : moment update frequency ; K : required rank ; S : number of samples

$\bar{\theta} \leftarrow \theta_0, \quad \bar{\theta}^2 \leftarrow \theta_0^2$ {Initialize moments}

for $i \leftarrow 1, 2, \dots, T$ **do**

$\theta_i \leftarrow \text{SGD}(\theta_{i-1})$ {Perform SGD update}

if `update_time` = `True` **then**

$n \leftarrow i/c$ {Number of models}

$\bar{\theta} \leftarrow \frac{n\bar{\theta} + \theta_i}{n+1}, \quad \bar{\theta}^2 \leftarrow \frac{n\bar{\theta}^2 + \theta_i^2}{n+1}$ {Update moments}

if `nbr_stored_param` = K **then**

`forget_first_param`

`store_new_param`($\theta_i - \bar{\theta}$) {Store deviation}

if `estimate_time` = `True` **then**

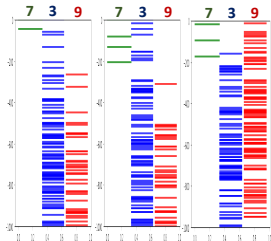
for $i \leftarrow 1, 2, \dots, S$ **do**

Draw $\tilde{\theta}_i \sim \mathcal{N}\left(\theta_{\text{SWA}}, \frac{1}{2}\Sigma_{\text{diag}} + \frac{\hat{D}\hat{D}^\top}{2(K-1)}\right)$

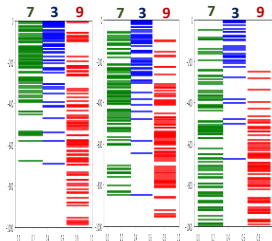
$p(y^*|\text{Data}) += \frac{1}{S}p(y^*|\tilde{\theta}_i)$

return $p(y^*|\text{Data})$

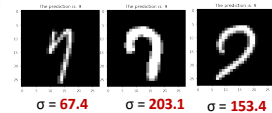
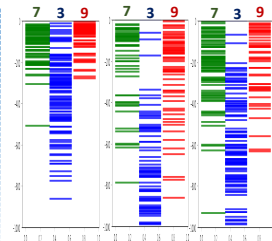
Results obtained on the classification task



Prediction: 7
CORRECT



Prediction: 3
INCORRECT



Prediction: 9
INCORRECT

Results obtained on the regression task

RMSE				
Dataset	Dropout (our imp.)	Dropout (paper)	VI	PBP
BostonHousing	2.85	2.97	4.32	3.01
energy	2.56	1.66	2.65	1.80
naval-propulsion	0.01	0.01	0.01	0.01
power-plant	4.40	4.02	4.33	4.12
wine-quality-rd	0.63	0.62	0.65	0.64

Predictive log-likelihood				
Dataset	Dropout (our imp.)	Dropout (paper)	VI	PBP
BostonHousing	-2.49	-2.46	-2.90	-2.57
energy	-2.39	-1.99	-2.39	-2.39
naval-propulsion	-1.72	3.80	3.73	3.73
power-plant	-3.63	-2.80	-2.89	-2.84
wine-quality-rd	-1.76	-0.93	-0.98	-0.97

Table 1: Test performance for our implementation of the dropout NN, the implementation of the paper, the variational inference method (VI, Graves [2011]), the Probabilistic back-propagation method (PBP, Hernandez-Lobato and P. [2015]).