

Additive Kernel GPPVAE

Santeri Mentu

Aalto University CSB group

Summer 2019

Contents

- Motivation for the work
- Mathematical description
- Software implementation
- Experimental results

Motivation for the work

- Variational Autoencoders (VAE) are a powerful method for unsupervised learning
- However the i.i.d. assumption for latent representations is too strong
- We would like to include label covariates, such as time, in the model
- Multiple extensions and variations on the VAE have been published
- ...

- Gaussian Process Variational Autoencoder (GPPVAE) is one such extension
- Gaussian Processes (GPs) act as prior for the latent space
- The GPs are indexed with label information
- Once the model is trained, we can perform label inference on new data

Bayesian inference

- In a generative model the probability of observations is defined by a likelihood $p(x|z)$ and prior on the latent space $p(z)$
- Ideally we could simply maximize the marginal probability of the posterior

$$p(z \mid x) = \frac{p(x \mid z)p(z)}{p(x)}$$

- However, the evidence term in the denominator is intractable
- Therefore we use a variational approximation where the latent distribution and the likelihoods are parametrized by NNs

Structure of a VAE

- We optimize the variational approximation of the latent distribution by minimizing the KL divergence between it, and the prior
- This is equivalent to maximizing the evidence lower bound, or ELBO

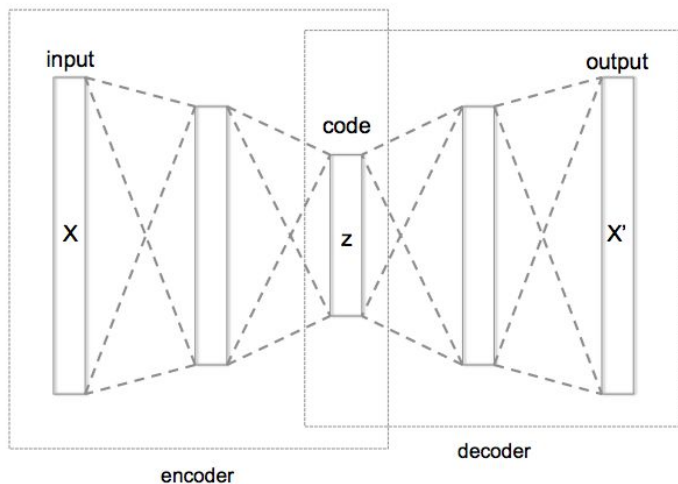
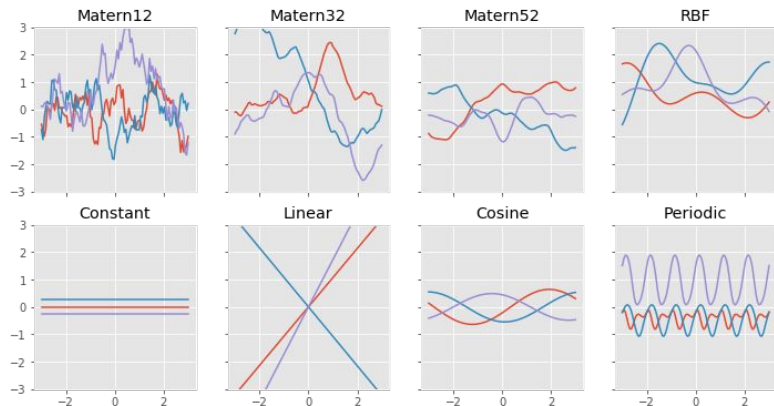


Figure from Wikipedia

(Very) short intro to Gaussian Processes

- A Gaussian process is a stochastic process where each finite collection of variables has a multivariate normal distribution
- A zero mean GP is completely defined by its covariance function, or kernel
- The choice of kernel defines the function space in GP regression



*Figure from GPflow by
Matthews et al.*

Gaussian prior

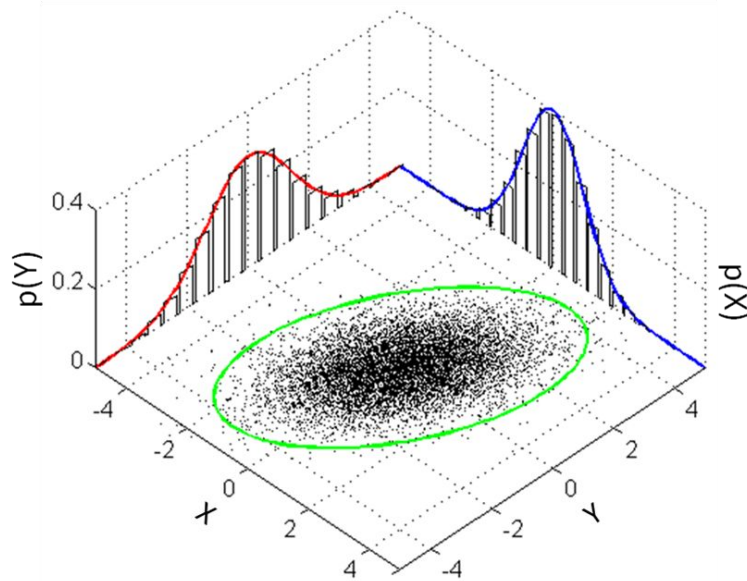


Figure from Wikipedia

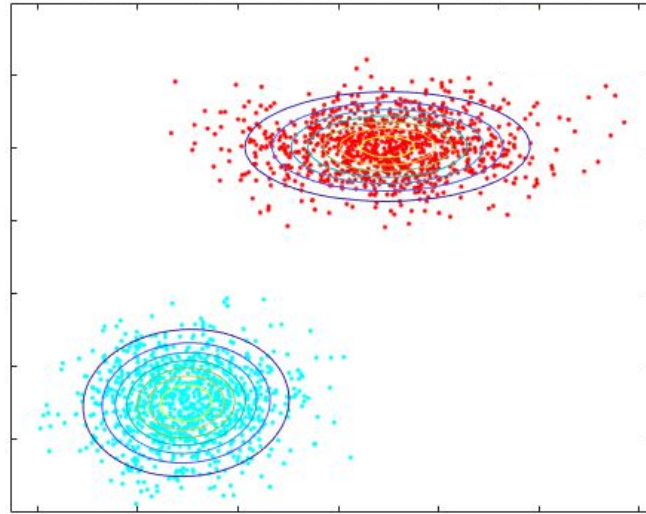


Figure by MathWorks

ELBO for GPPVAE

- ELBO derived by Casale et al.

$$\begin{aligned} \log p(\mathbf{Y} \mid \mathbf{X}, \mathbf{W}, \phi, \sigma_y^2, \boldsymbol{\theta}) \geq & \mathbb{E}_{\mathbf{Z} \sim q_{\boldsymbol{\psi}}} \left[\sum_n \log \mathcal{N}(\mathbf{y}_n \mid g_{\phi}(\mathbf{z}_n), \sigma_y^2 \mathbf{I}_K) + \log p(\mathbf{Z} \mid \mathbf{X}, \mathbf{W}, \boldsymbol{\theta}, \alpha) \right] + \\ & + \frac{1}{2} \sum_{nl} \log(\sigma_{\boldsymbol{\psi}}^{z^2}(\mathbf{y}_n)_l) + \text{const.} \end{aligned}$$

- Which gives the loss for SGD

$$\begin{aligned} l(\phi, \boldsymbol{\psi}, \boldsymbol{\theta}, \alpha, \sigma_y^2) = \\ = NK \log \sigma_y^2 + \underbrace{\sum_n \frac{\|\mathbf{y}_n - g_{\phi}(\mathbf{z}_{\boldsymbol{\psi}_n})\|^2}{2\sigma_y^2}}_{\text{reconstruction term}} - \underbrace{\log p(\mathbf{Z}_{\boldsymbol{\psi}} \mid \mathbf{X}, \mathbf{W}, \boldsymbol{\theta}, \alpha)}_{\text{latent-space GP term}} + \underbrace{\frac{1}{2} \sum_{nl} \log(\sigma_{\boldsymbol{\psi}}^{z^2}(\mathbf{y}_n)_l)}_{\text{regularization term}}, \end{aligned}$$

GP likelihoods

The marginal likelihood for observations

$$\begin{aligned} p(\mathbf{y}, f_*) &= \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}, f_*)d\mathbf{f} \\ &= \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_{ff} + \sigma^2 \mathbf{I} & \mathbf{K}_{f_*f} \\ \mathbf{K}_{f_*f} & K_{f_*f_*} \end{bmatrix}\right) \end{aligned}$$

The marginal likelihood of kernel parameters

$$\begin{aligned} p(\mathbf{y}|\boldsymbol{\theta}) &= \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\boldsymbol{\theta})d\mathbf{f} \\ &= \int \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}) d\mathbf{f} \\ &= \mathcal{N}(\mathbf{y}|\mathbf{0}, \sigma^2 \mathbf{I} + \mathbf{K}) \end{aligned}$$

Additive and multiplicative kernels

- It would be most beneficial if kernels could be configured in a modular and flexible manner
- My implementations allows for kernel addition and multiplication thanks to the GPyTorch library

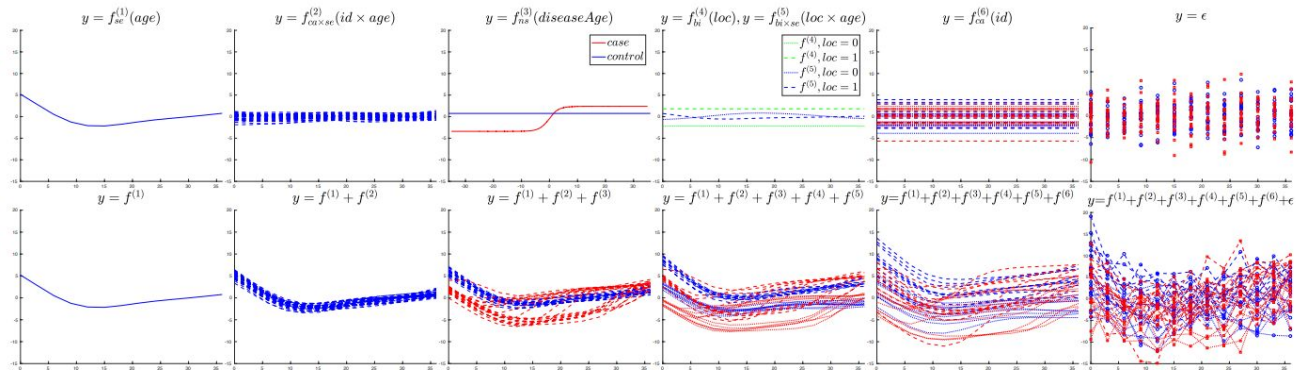


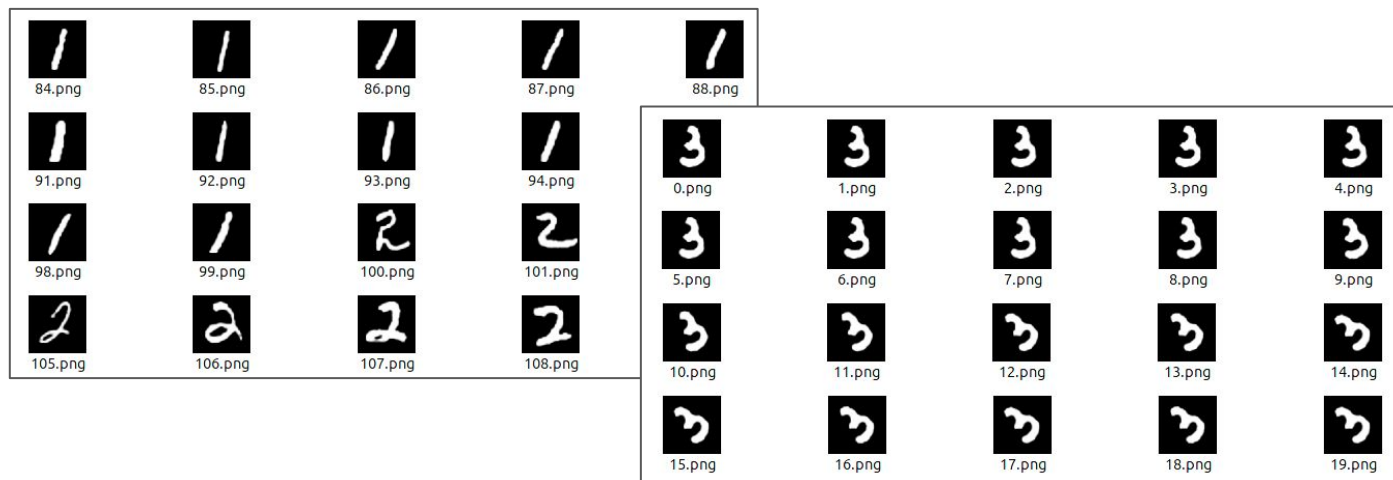
Figure by Cheng et al.

Software Implementation

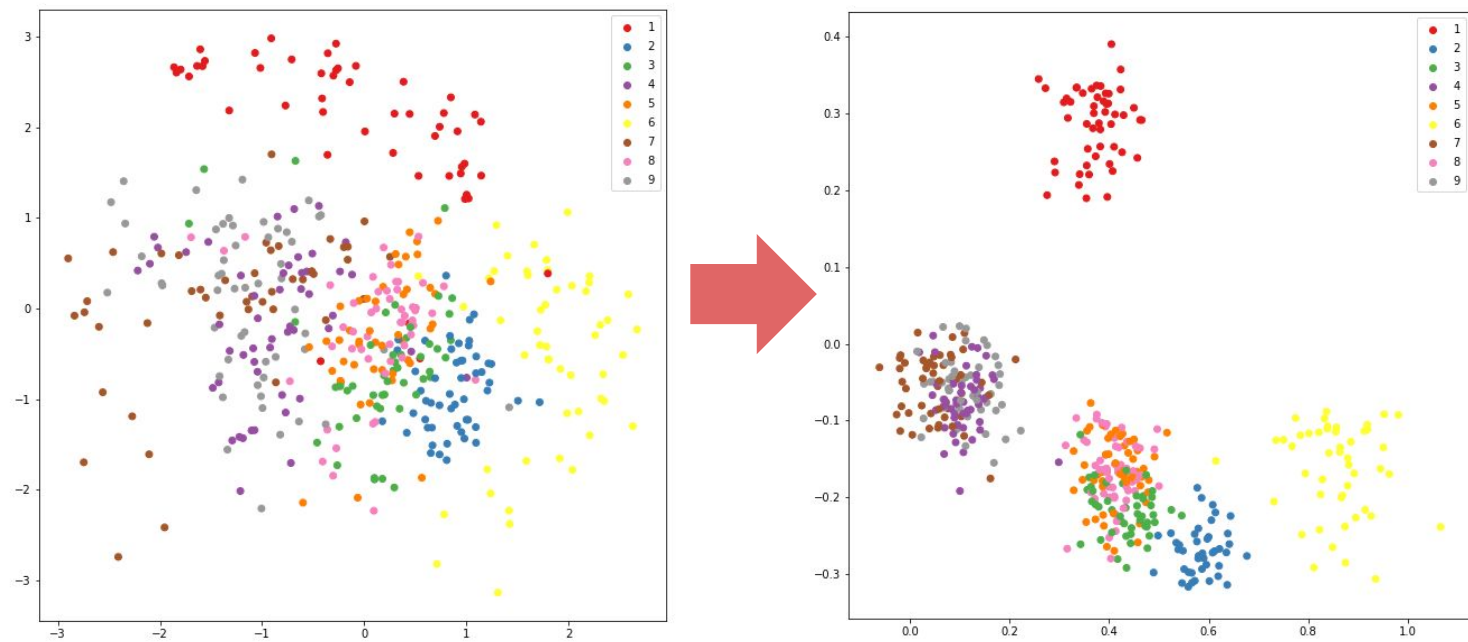
- The implementation was built using PyTorch and the GPyTorch GP library
- GPyTorch provides a highly modular and flexible framework for building GP models, which can be modified to suit specific requirements
- The implementation is available at Aalto Version

Test datasets

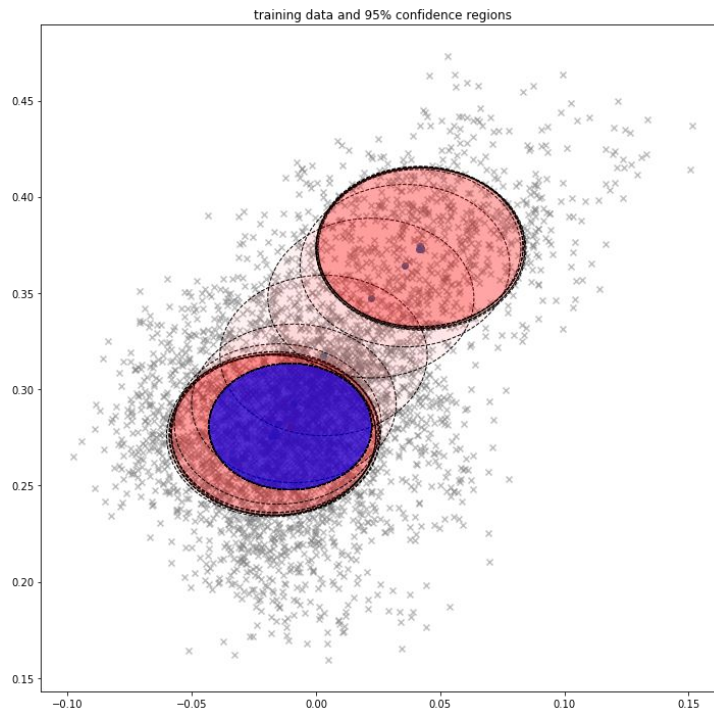
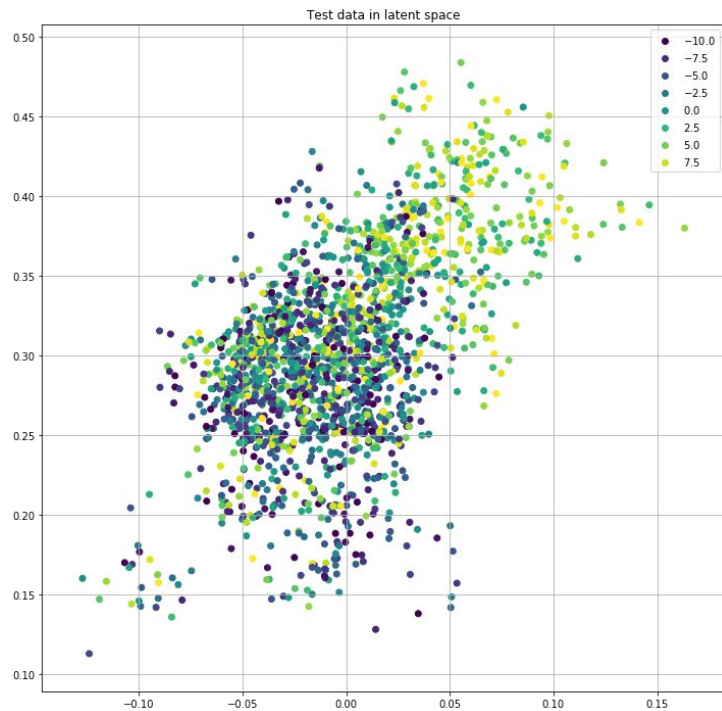
- The implementation was tested using two toy datasets



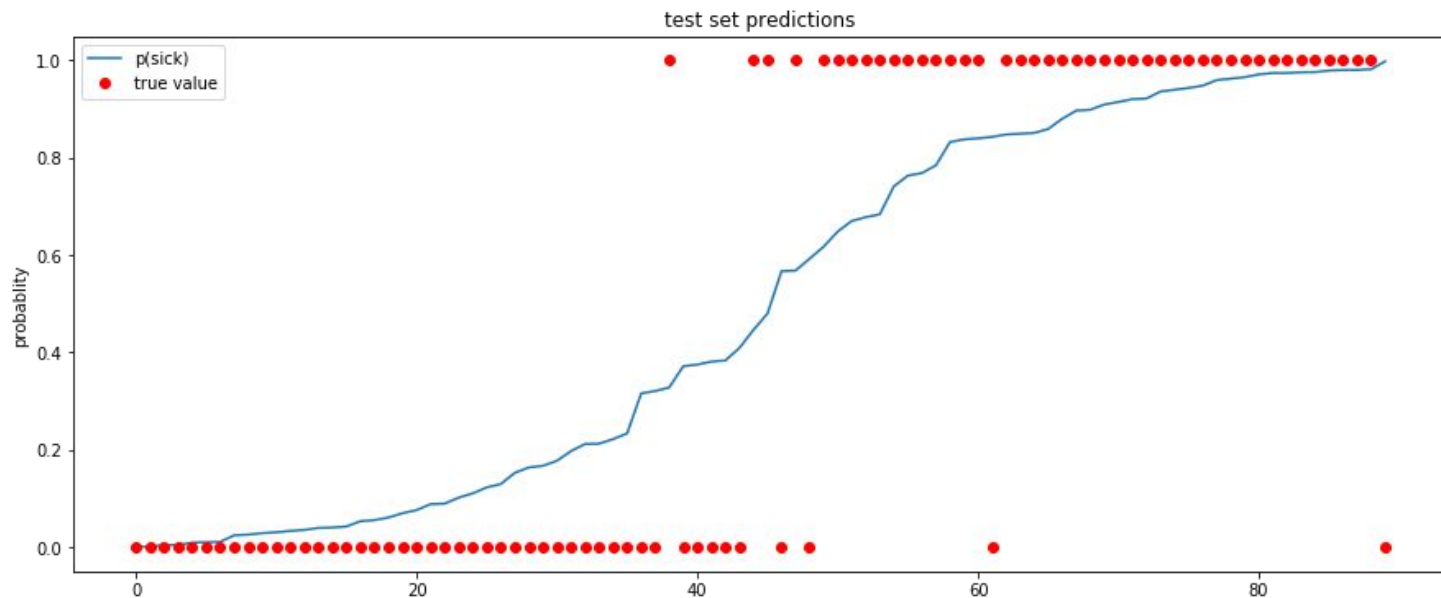
MNIST with categorical kernel



Rotating MNIST with binary + warping kernel



Rotating MNIST label inference



All code and notebooks available at:

https://version.aalto.fi/gitlab/mentus1/HIT_VAE

Special thanks to Gleb Tikhonov