

# Applications of Gaussian Process Latent Variable Models in Finance \*

Rajbir-Singh Nirwan<sup>†,2</sup> and Nils Bertschinger<sup>‡1,2</sup>

<sup>1</sup>Frankfurt Institute for Advanced Studies, Frankfurt am Main, Germany

<sup>2</sup>Goethe University, Frankfurt am Main, Germany

April 19, 2019

## Abstract

Estimating covariances between financial assets plays an important role in risk management. In practice, when the sample size is small compared to the number of variables, the empirical estimate is known to be very unstable. Here, we propose a novel covariance estimator based on the Gaussian Process Latent Variable Model (GP-LVM). Our estimator can be considered as a non-linear extension of standard factor models with readily interpretable parameters reminiscent of market betas. Furthermore, our Bayesian treatment naturally shrinks the sample covariance matrix towards a more structured matrix given by the prior and thereby systematically reduces estimation errors. Finally, we discuss some financial applications of the GP-LVM.

## 1 Introduction

Many financial problems require the estimation of covariance matrices between given assets. This may be useful to optimize one's portfolio, i.e.: maximize the portfolio returns  $\mathbf{w}^T \mathbf{r}$  and/or minimize the volatility  $\sqrt{\mathbf{w}^T \mathbf{K} \mathbf{w}}$ . Indeed, Markowitz received a Noble Price in economics for his treatment of modern portfolio theory [1]. In practice, estimating historical returns and high-dimensional covariance matrices is challenging and often times equally weighted portfolio outperforms the portfolio constructed from sample estimates [2]. The estimation of covariance matrices is especially hard, when the number of assets is large compared to the number of observations. Sample estimations in those cases are very unstable or can even become singular. To cope with this problem, a wide range of estimators, e.g. factor models such as the single-index model [3] or shrinkage estimators [4], have been developed and employed in portfolio optimization.

With todays machine learning techniques we can even further improve those estimates. Machine learning has already arrived in finance. Nevmyvaka et al. [5] trained an agent via reinforcement learning to optimally execute trades. Gately [6] forecast asset prices with neural networks and Chapados et al. [7] with Gaussian Processes. Recently, Heaton et al. [8] made an ansatz to optimally allocate portfolios using deep autoencoders. Wu et al. [9] used Gaussian Processes to build volatility models and Wilson et al. [10] to estimate time varying covariance matrices. Bayesian

---

\*This is a pre-print of an article accepted at the IntelliSys 2019.

<sup>†</sup>nirwan@fias.uni-frankfurt.de

<sup>‡</sup>bertschinger@fias.uni-frankfurt.de

machine learning methods are used more and more in this domain. The fact, that in a Bayesian framework parameters are not treated as true values, but as random variables, accounts for estimation uncertainties and can even alleviate the unwanted impacts of outliers. Furthermore, one can easily incorporate additional information and/or personal views by selecting suitable priors.

In this paper, we propose a Bayesian covariance estimator based on the Gaussian Process Latent Variable Models (GP-LVMs) [11], which can be considered as a non-linear extension of standard factor models with readily interpretable parameters reminiscent of market betas. Our Bayesian treatment naturally shrinks the sample covariance matrix (which maximizes the likelihood function) towards a more structured matrix given by the prior and thereby systematically reduces estimation errors. We evaluated our model on the stocks of S&P500 and found significant improvements in terms of model fit compared to classical linear models. Furthermore we suggest some financial applications, where Gaussian Processes can be used as well. That includes portfolio allocation, price prediction for less frequently traded stocks and non-linear clustering of stocks into their sub-sectors.

In section 2 we begin with an introduction to the Bayesian non-parametric Gaussian Processes and discuss the associated requirements for learning. Section 3 introduces the financial background needed for portfolio optimization and how to relate it to Gaussian Processes. In section 4 we conduct experiments on covariance matrix estimations and discuss the results. We conclude in section 5.

## 2 Background

In this paper, we utilize a Bayesian non-parametric machine learning approach based on Gaussian Processes (GPs). Combining those with latent variable models leads to Gaussian Process Latent Variable Models (GP-LVMs), that we use to estimate the covariance between different assets. These approaches have been described in detail in [11,12]. We provide a brief review here. Subsequently, we show, how to relate those machine learning approaches to the known models in finance, e.g. the single-index model [3].

### 2.1 Gaussian Processes

A Gaussian Process (GP) is a generalization of the Gaussian distribution. Using a GP, we can define a distribution over functions  $f(\mathbf{x})$ , where  $\mathbf{x} \in \mathbb{R}^Q$  and  $f(\cdot) \in \mathbb{R}$ . Like a Gaussian distribution, the GP is specified by a mean and a covariance. In the GP case, however, the mean is a function of the input variable  $m(\mathbf{x})$  and the covariance is a function of two variables  $k(\mathbf{x}, \mathbf{x}')$ , which contains information about how the GP evaluated at  $\mathbf{x}$  and  $\mathbf{x}'$  covary

$$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})], \quad (1)$$

$$k(\mathbf{x}, \mathbf{x}') = \text{cov}(f(\mathbf{x}), f(\mathbf{x}')). \quad (2)$$

We write  $f \sim \text{GP}(m(\cdot), k(\cdot, \cdot))$ . Any finite collection of function values, at  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , is jointly Gaussian distributed

$$p(f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_N)) = \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}), \quad (3)$$

where  $\boldsymbol{\mu} = (m(\mathbf{x}_1), \dots, m(\mathbf{x}_N))^T$  is the mean vector and  $\mathbf{K} \in \mathbb{R}^{N \times N}$  is the Gram matrix with entries  $\mathbf{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ . We refer to the covariance function as kernel function. The properties of the function  $f$  (i.e. smoothness, periodicity) are determined by the choice of this kernel function. For example, sampled functions from a GP with an exponentiated quadratic covariance function  $k_{\text{SE}}(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp(-0.5\|\mathbf{x} - \mathbf{x}'\|_2^2/l^2)$  smoothly vary with lengthscale  $l$  and are infinitely often differentiable.

Given a dataset  $\mathcal{D}$  of  $N$  input points  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$  and  $N$  corresponding targets  $\mathbf{y} = (y_1, \dots, y_N)^T$ , the predictive distribution for a zero mean GP at  $N_*$  new locations  $\mathbf{X}_*$  reads [12]

$$\mathbf{y}_*|\mathbf{X}_*, \mathbf{y}, \mathbf{X} \sim \mathcal{N}(\mathbf{f}_*, \mathbf{K}_*), \quad (4)$$

where

$$\mathbf{f}_* = \mathbf{K}_{\mathbf{X}_* \mathbf{X}} \mathbf{K}_{\mathbf{X} \mathbf{X}}^{-1} \mathbf{y}, \quad (5)$$

$$\mathbf{K}_* = \mathbf{K}_{\mathbf{X}_* \mathbf{X}_*} - \mathbf{K}_{\mathbf{X}_* \mathbf{X}} \mathbf{K}_{\mathbf{X} \mathbf{X}}^{-1} \mathbf{K}_{\mathbf{X} \mathbf{X}_*}. \quad (6)$$

$\mathbf{K}_{\mathbf{X}_* \mathbf{X}} \in \mathbb{R}^{N_* \times N}$  is the covariance matrix between the GP evaluated at  $\mathbf{X}_*$  and  $\mathbf{X}$ ,  $\mathbf{K}_{\mathbf{X} \mathbf{X}} \in \mathbb{R}^{N \times N}$  is the covariance matrix of the GP evaluated at  $\mathbf{X}$ . As we can see in equations (5) and (6), the kernel function plays a very important role in the GP framework and will be important for our financial model as well.

## 2.2 Gaussian Process Latent Variable Model

Often times we are just given a data matrix  $\mathbf{Y} \in \mathbb{R}^{N \times D}$  and the goal is to find a lower dimensional representation  $\mathbf{X} \in \mathbb{R}^{N \times Q}$ , without losing too much information. Principal component analysis (PCA) is one of the most used techniques for reducing the dimensions of the data, which has also been motivated as the maximum likelihood solution to a particular form of Gaussian Latent Variable Model [13]. PCA embeds  $\mathbf{Y}$  via a linear mapping into the latent space  $\mathbf{X}$ . Lawrence [11] introduced the Gaussian Process Latent Variable Model (GP-LVM) as a non-linear extension of probabilistic PCA. The generative procedure takes the form

$$\mathbf{Y}_{n,:} = \mathbf{f}(\mathbf{X}_{n,:}) + \boldsymbol{\epsilon}_n, \quad (7)$$

where  $\mathbf{f} = (f_1, \dots, f_D)^T$  is a group of  $D$  independent samples from a GP, i.e.  $f_d \sim \text{GP}(0, k(\cdot, \cdot))$ . By this, we assume the rows of  $\mathbf{Y}$  to be jointly Gaussian distributed and the columns to be independent, i.e. each sample  $\mathbf{Y}_{:,d} \sim \mathcal{N}(\mathbf{Y}_{:,d}|\mathbf{0}, \mathbf{K})$  where  $\mathbf{K} = k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}$  and  $\sigma^2$  denotes the variance of the random noise  $\boldsymbol{\epsilon}$ . The marginal likelihood of  $\mathbf{Y}$  becomes [11]

$$p(\mathbf{Y}|\mathbf{X}) = \prod_{d=1}^D \mathcal{N}(\mathbf{Y}_{:,d}|\mathbf{0}, \mathbf{K}) = \frac{1}{(2\pi)^{ND/2} |\mathbf{K}|^{D/2}} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{K}^{-1} \mathbf{Y} \mathbf{Y}^T)\right). \quad (8)$$

The dependency on the latent positions  $\mathbf{X}$  and the kernel hyperparameters is given through the kernel matrix  $\mathbf{K}$ . As suggested by Lawrence [11], we can optimize the log marginal likelihood  $\log p(\mathbf{Y}|\mathbf{X})$  with respect to the latent positions and the hyperparameters.

## 2.3 Variational Inference

Optimization can easily lead to overfitting. Therefore, a fully Bayesian treatment of the model would be preferable but is intractable. Bishop [14] explains the variational inference framework, which not only handles the problem of overfitting but also allows to automatically select the dimensionality of the latent space. Instead of optimizing equation (8), we want to calculate the posterior using Bayes rule  $p(\mathbf{X}|\mathbf{Y}) = p(\mathbf{Y}|\mathbf{X})p(\mathbf{X})/p(\mathbf{Y})$ , which is intractable. The idea behind variational Bayes is to approximate the true posterior  $p(\mathbf{X}|\mathbf{Y})$  by another distribution  $q(\mathbf{X})$ , selected from a tractable family. The goal is to select the one distribution  $q(\mathbf{X})$ , that is closest to the true posterior  $p(\mathbf{X}|\mathbf{Y})$  in some sense. A natural choice to quantify the closeness is given by the Kullback-Leibler divergence [15]

$$\text{KL}[q(\mathbf{X})||p(\mathbf{X}|\mathbf{Y})] = \int q(\mathbf{X}) \log \frac{q(\mathbf{X})}{p(\mathbf{X}|\mathbf{Y})} d\mathbf{X}. \quad (9)$$

By defining  $\tilde{p}(\mathbf{X}|\mathbf{Y}) = p(\mathbf{Y}|\mathbf{X})p(\mathbf{X})$  as the unnormalized posterior, equation (9) becomes

$$\begin{aligned} \text{KL}[q(\mathbf{X})||p(\mathbf{X}|\mathbf{Y})] &= \int q(\mathbf{X}) \log \frac{q(\mathbf{X})}{\tilde{p}(\mathbf{X}|\mathbf{Y})} d\mathbf{X} + \log p(\mathbf{Y}) \\ &= - \underbrace{\mathbb{E}_{q(\mathbf{X})} \left[ \log \frac{\tilde{p}(\mathbf{X}|\mathbf{Y})}{q(\mathbf{X})} \right]}_{\text{ELBO}} + \log p(\mathbf{Y}) \end{aligned} \quad (10)$$

with the first term on the right hand side being known as the evidence lower bound (ELBO). Equation (10) is the objective function we want to minimize with respect to  $q(\mathbf{X})$  to get a good approximation to the true posterior. Note that on the left hand side only the ELBO is  $q$  dependent. So, in order to minimize (10), we can just as well maximize the ELBO. Because the Kullback-Leibler divergence is non-negative, the ELBO is a lower bound on the evidence  $\log p(\mathbf{Y})$ <sup>1</sup>. Therefore, this procedure not only gives the best approximation to the posterior within the variational family under the KL criterion, but it also bounds the evidence, which serves as a measure of the goodness of our fit. The number of latent dimensions  $Q$  can be chosen to be the one, which maximizes the ELBO.

So, GP-LVM is a model, which reduces the dimensions of our data-matrix  $\mathbf{Y} \in \mathbb{R}^{N \times D}$  from  $D$  to  $Q$  in a non-linear way and at the same time estimates the covariance matrix between the  $N$  points. The estimated covariance matrix can then be used for further analysis.

### 3 Finance

Now we have a procedure to estimate the covariance matrix between different datapoints. This section discusses how we can relate this to financial models.

#### 3.1 CAPM

The Capital Asset Pricing Model (CAPM) describes the relationship between the expected returns of an asset  $\mathbf{r}_n \in \mathbb{R}^D$  for  $D$  days and its risk  $\beta_n$

$$\mathbb{E}[\mathbf{r}_n] = \mathbf{r}_f + \beta_n \mathbb{E}[\mathbf{r}_m - \mathbf{r}_f], \quad (11)$$

where  $\mathbf{r}_f \in \mathbb{R}^D$  is the risk free return on  $D$  different days and  $\mathbf{r}_m$  is the market return on  $D$  different days. The main idea behind CAPM is, that an investor needs to be compensated for the risk of his holdings. For a risk free asset with  $\beta_n = 0$ , the expected return  $\mathbb{E}[\mathbf{r}_n]$  is just the risk free rate  $\mathbf{r}_f$ . If an asset is risky with a risk  $\beta_n \neq 0$ , the expected return  $\mathbb{E}[\mathbf{r}_n]$  is increased by  $\beta_n \mathbb{E}[\tilde{\mathbf{r}}_m]$ , where  $\tilde{\mathbf{r}}_m$  is the excess return of the market  $\tilde{\mathbf{r}}_m = \mathbf{r}_m - \mathbf{r}_f$ .

We can write down equation (11) in terms of the excess return  $\tilde{\mathbf{r}} = \mathbf{r} - \mathbf{r}_f$  and get

$$\mathbb{E}[\tilde{\mathbf{r}}_n] = \beta_n \mathbb{E}[\tilde{\mathbf{r}}_m], \quad (12)$$

where  $\tilde{\mathbf{r}}_n$  is the excess return of a given asset and  $\tilde{\mathbf{r}}_m$  is the excess return of the market (also called a risk factor). Arbitrage pricing theory [16] generalizes the above model by allowing multiple risk factors  $\mathbf{F}$  beside the market  $\tilde{\mathbf{r}}_m$ . In particular, it assumes that asset returns follow a factor structure

$$\mathbf{r}_n = \alpha_n + \mathbf{F}\beta_n + \epsilon_n, \quad (13)$$

<sup>1</sup>The evidence  $\log p(\mathbf{Y})$  is also referred to as log marginal likelihood in the literature. The term marginal likelihood is already used for  $p(\mathbf{Y}|\mathbf{X})$  in this paper. Therefore, we will refer to  $\log p(\mathbf{Y})$  as the evidence.

with  $\epsilon_n$  denoting some independent zero mean noise with variance  $\sigma_n^2$ . Here,  $\mathbf{F} \in \mathbb{R}^{D \times Q}$  is the matrix of  $Q$  factor returns on  $D$  days and  $\beta_n \in \mathbb{R}^Q$  is the loading of stock  $n$  to the  $Q$  factors. Arbitrage pricing theory [17] then shows that the expected excess returns adhere to

$$\mathbb{E}[\tilde{\mathbf{r}}_n] = \mathbb{E}[\tilde{\mathbf{F}}]\beta_n, \quad (14)$$

i.e. the CAPM is derived as special case when assuming a single risk factor (single-index model).

To match the form of the GP-LVM (see equation (7)), we rewrite equation (13) as

$$\mathbf{r}_{n,:} = \mathbf{f}(\mathbf{B}_{n,:}) + \epsilon_n, \quad (15)$$

where  $\mathbf{r} \in \mathbb{R}^{N \times D}$  is the return matrix<sup>2</sup>. Note that assuming latent factors distributed as  $\mathbf{F} \sim \mathcal{N}(0, 1)$  and marginalizing over them, equation (13) is a special case of equation (15) with  $\mathbf{f}$  drawn from a GP mapping  $\beta_n$  to  $\mathbf{r}_n$  with a linear kernel. Interestingly, this provides an exact correspondence with factor models by considering the matrix of couplings  $\mathbf{B} = (\beta_1, \dots, \beta_N)^T$  as the latent space positions<sup>3</sup>. In this perspective the factor model can be seen as a linear dimensionality reduction, where we reduce the  $N \times D$  matrix  $\mathbf{r}$  to a low rank matrix  $\mathbf{B}$  of size  $N \times Q$ . By choosing a non-linear kernel  $k(\cdot, \cdot)$  the GP-LVM formulation readily allows for non-linear dimensionality reductions. Since, it is generally known, that different assets have different volatilities, we further generalize the model. In particular, we assume the noise  $\epsilon$  to be a zero mean Gaussian, but allow for different variances  $\sigma_n^2$  for different stocks. For this reason, we also have to parameterize the kernel (covariance) matrix in a different way than usual. Section 4.1 explains how to deal with that. The model is then approximated using variational inference as described in section 2.3. After inferring  $\mathbf{B}$  and the hyperparameters of the kernel, we can calculate the covariance matrix  $\mathbf{K}$  and use it for further analysis.

### 3.2 Modern Portfolio Theory

Markowitz [1] provided the foundation for modern portfolio theory, for which he received a Nobel Prize in economics. The method analyses how good a given portfolio is, based on the mean and the variance of the returns of the assets contained in the portfolio. This can also be formulated as an optimization problem for selecting an optimal portfolio, given the covariance between the assets and the risk tolerance  $q$  of the investor.

Given the covariance matrix  $\mathbf{K} \in \mathbb{R}^{N \times N}$ , we can calculate the optimal portfolio weights  $\mathbf{w}$  by

$$\mathbf{w}_{\text{opt}} = \min_{\mathbf{w}} (\mathbf{w}^T \mathbf{K} \mathbf{w} - q \bar{\mathbf{r}}^T \mathbf{w}), \quad (16)$$

where  $\bar{\mathbf{r}}$  is the mean return vector. Risk friendly investors have a higher  $q$  than risk averse investors. The model is constrained by  $\sum_{\mathbf{w}} = 1$ . Since  $\bar{\mathbf{r}}$  is very hard to estimate in general and we are primarily interested in the estimation of the covariance matrix  $\mathbf{K}$ , we set  $q$  to zero and get

$$\mathbf{w}_{\text{opt}} = \min_{\mathbf{w}} (\mathbf{w}^T \mathbf{K} \mathbf{w}). \quad (17)$$

This portfolio is called the minimal risk portfolio, i.e. the solution to equation (17) provides the weights for the portfolio, which minimizes the overall risk, assuming the estimated  $\mathbf{K}$  is the true covariance matrix.

<sup>2</sup>To stay consistent with the financial literature, we denote the return matrix with lower case  $\mathbf{r}$ .

<sup>3</sup>Because of the context, from now on we will use  $\beta$  for the latent space instead of  $\mathbf{x}$ .

## 4 Experiments

In this section, we discuss the performance of the GP-LVM on financial data. After describing the data collection and modeling procedure, we evaluate the model on the daily return series of the S&P500 stocks. Subsequently, we discuss further financial applications. In particular, we show how to build a minimal risk portfolio (this can easily be extended to maximizing returns as well), how to fill-in prices for assets which are not traded frequently and how to visualize sector relations between different stocks (latent space embedding).

### 4.1 Data Collection and Modeling

For a given time period, we take all the stocks from the S&P500, whose daily close prices were available for the whole period<sup>4</sup>. The data were downloaded from Yahoo Finance. After having the close prices in a matrix  $\mathbf{p} \in \mathbb{R}^{N \times (D+1)}$ , we calculate the return matrix  $\mathbf{r} \in \mathbb{R}^{N \times D}$ , where  $\mathbf{r}_{nd} = (\mathbf{p}_{n,d} - \mathbf{p}_{n,d-1}) / \mathbf{p}_{n,d-1}$ .  $\mathbf{r}$  builds the basis of our analysis.

We can feed  $\mathbf{r}$  into the GP-LVM. The GP-LVM procedure, as described in section 2, assumes the likelihood to be Gaussian with the covariance given by the kernel function for each day and assumes independency over different days. We use the following kernel functions

$$\begin{aligned} k_{\text{noise}}(\beta_i, \beta_j) &= \sigma_{\text{noise},i}^2 \delta_{i,j}, \\ k_{\text{linear}}(\beta_i, \beta_j) &= \sigma^2 \beta_i^T \beta_j, \end{aligned} \quad (18)$$

and the stationary kernels

$$\begin{aligned} k_{\text{se}}(\beta_i, \beta_j) &= k_{\text{se}}(d_{ij}) = \exp\left(-\frac{1}{2l^2} d_{ij}^2\right), \\ k_{\text{exp}}(\beta_i, \beta_j) &= k_{\text{exp}}(d_{ij}) = \exp\left(-\frac{1}{2l} d_{ij}\right), \\ k_{\text{m32}}(\beta_i, \beta_j) &= k_{\text{m32}}(d_{ij}) = \left(1 + \frac{\sqrt{3}d_{ij}}{l}\right) \exp\left(-\frac{\sqrt{3}}{2l} d_{ij}\right), \end{aligned} \quad (19)$$

where  $d_{ij} = \|\beta_i - \beta_j\|_2$  is the Euclidean distance between  $\beta_i$  and  $\beta_j$ .  $\sigma^2$  is the kernel variance and  $l$  kernel lengthscale. Note that since the diagonal elements of stationary kernels are the same, they are not well suited for an estimation of a covariance matrix between different financial assets. Therefore, in the case of stationary kernel we decompose our covariance matrix  $\mathbf{K}_{\text{cov}}$  into a vector of coefficient scales  $\boldsymbol{\sigma}$  and a correlation matrix  $\mathbf{K}_{\text{corr}}$ , such that  $\mathbf{K}_{\text{cov}} = \boldsymbol{\Sigma} \mathbf{K}_{\text{corr}} \boldsymbol{\Sigma}$ , where  $\boldsymbol{\Sigma}$  is a diagonal matrix with  $\boldsymbol{\sigma}$  on the diagonal. The full kernel function  $k(\cdot, \cdot)$  at the end is the sum of the noise kernel  $k_{\text{noise}}$  and one of the other kernels. In matrix form we get

$$\begin{aligned} \mathbf{K}_{\text{linear}} &= k_{\text{linear}}(\mathbf{B}, \mathbf{B}) + k_{\text{noise}}(\mathbf{B}, \mathbf{B}), \\ \mathbf{K}_{\text{se}} &= \boldsymbol{\Sigma} k_{\text{se}}(\mathbf{B}, \mathbf{B}) \boldsymbol{\Sigma} + k_{\text{noise}}(\mathbf{B}, \mathbf{B}), \\ \mathbf{K}_{\text{exp}} &= \boldsymbol{\Sigma} k_{\text{exp}}(\mathbf{B}, \mathbf{B}) \boldsymbol{\Sigma} + k_{\text{noise}}(\mathbf{B}, \mathbf{B}), \\ \mathbf{K}_{\text{m32}} &= \boldsymbol{\Sigma} k_{\text{m32}}(\mathbf{B}, \mathbf{B}) \boldsymbol{\Sigma} + k_{\text{noise}}(\mathbf{B}, \mathbf{B}), \end{aligned} \quad (20)$$

where  $\mathbf{B} = (\beta_1, \dots, \beta_N)^T$ . We chose the following priors

$$\mathbf{B} \sim \mathcal{N}(0, 1) \quad l, \sigma \sim \text{InvGamma}(3, 1) \quad \boldsymbol{\sigma}, \sigma_{\text{noise}} \sim \mathcal{N}(0, 0.5). \quad (21)$$

<sup>4</sup>We are aware that this introduces survivorship bias. Thus, in some applications our results might be overly optimistic. Nevertheless, we expect relative comparisons between different models to be meaningful.

The prior on  $\mathbf{B}$  determines how much space is allocated to the points in the latent space. The volume of this space expands with higher latent space dimension, which make the model prone to overfitting. To cope with that, we assign an inverse gamma prior to the lengthscale  $l$  and  $\sigma$  ( $\sigma$  in the linear kernel has a similar functionality as  $l$  in the stationary kernels). It allows larger values for higher latent space dimension, thereby shrinking the effective latent space volume and exponentially suppresses very small values, which deters overfitting as well. The parameters of the priors are chosen such that they allow for enough volume in the latent space for roughly 100-150 datapoints, which we use in our analysis. If the number of points is drastically different, one should adjust the parameters accordingly. The kernel standard deviations  $\sigma_{\text{noise}}$  and  $\sigma$  are assigned a half Gaussian prior with variance 0.5, which is essentially a flat prior, since the returns are rarely above 0.1 for a day.

Model inference under these specifications (GP-LVM likelihood for the data and prior for  $\mathbf{B}$  and all kernel hyperparameters  $\sigma, \sigma_{\text{noise}}, l$  and  $\sigma$ , which we denote by  $\theta$ ) was carried out by variational inference as described in section 2.3. To this end, we implemented all models in the probabilistic programming language Stan [18], which supports variational inference out of the box. We approximate the posterior by independent Gaussians. The source code is available on Github<sup>5</sup>. We tested different initializations for the parameter (random, PCA solution and Isomap solution), but there were no significant differences in the ELBO. So, we started the inference 50 times with random initializations and took the result with highest ELBO for each kernel function and  $Q$ .

## 4.2 Model Comparison

The GP-LVM can be evaluated in many different ways. Since, it projects the data from a  $D$ -dimensional space to a  $Q$ -dimensional latent space, we can look at the reconstruction error. A suitable measure of the reconstruction error is the R-squared ( $R^2$ ) score, which is equal to one if there is no reconstruction error and decreases if the reconstruction error increases. It is defined by

$$R^2 = 1 - \frac{\sum_i (y_i - f_i)^2}{\sum_i (y_i - \bar{y})^2}, \quad (22)$$

where  $y = (y_1, \dots, y_N)^T$  are the true values,  $f = (f_1, \dots, f_N)^T$  are the predicted values and  $\bar{y}$  is the mean of  $y$ . In the following, we look at the  $R^2$  as a function of the latent dimension  $Q$  for different kernels. Figure 1 (left plot) shows the results for three non-linear kernels and the linear kernel. Only a single dimension in the non-linear case can already capture almost 50% of the structure, whereas the linear kernel is at 15%. As one would expect, the higher  $Q$ , the more structure can be learned. But at some point the model will also start learning the noise and overfit.

As introduced in section 2.3, the ELBO is a good measure to evaluate different models and already incorporates models complexity (overfitting). Figure 1 (right plot) shows the ELBO as a function of the latent dimension  $Q$ . Here we see, that the model selects just a few latent dimensions. Depending on the used kernel, latent dimensions from three to five are already enough to capture the structure. If we increase the dimensions further, the ELBO starts dropping, which is a sign of overfitting. As can be seen from Figure 1, we do not need to go to higher dimensions.  $Q$  between two and five is already good enough and captures the structure that can be captured by the model.

## 4.3 Applications

The GP-LVM provides us the covariance matrix  $\mathbf{K}$  and the latent space representation  $\mathbf{B}$  of the data. We can build a lot of nice applications upon that, some of which are discussed in this section.

<sup>5</sup><https://github.com/RSNirwan/GPLVMsInFinance>

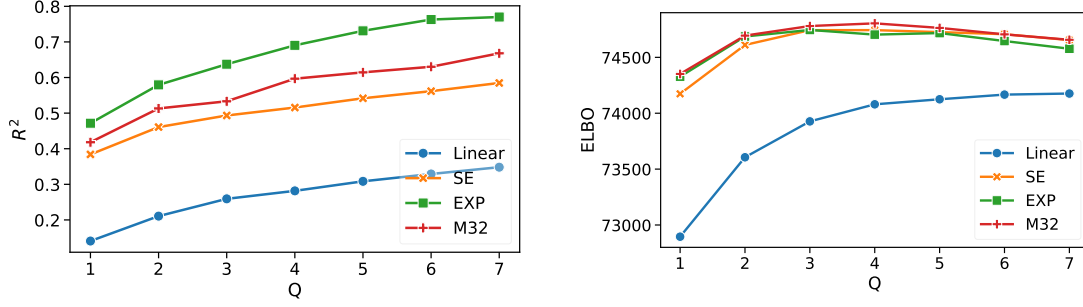


Figure 1: Left:  $R^2$ -score as a function of the latent dimension  $Q$  for different kernel functions. Right: ELBO as a function of the latent dimension  $Q$ . We randomly chose 120 stocks from the S&P500 and made the analysis on returns from Jan 2017 to Dec 2017.

#### 4.3.1 Portfolio Allocation

After inferring  $\mathbf{B}$  and  $\boldsymbol{\theta}$ , we can reconstruct the covariance matrix  $\mathbf{K}$  using equation (20). Thereafter, we only need to minimize equation (17), which provides the weights  $\mathbf{w}$  for our portfolio in the future. Minimization of (17) is done under the constraints:  $\sum_n \mathbf{w}_n = 1$  and  $0 < \mathbf{w}_n < 0.1, \forall n$ . These constraints are commonly employed in practice and ensure that we are fully invested, take on long positions only and prohibit too much weight for a single asset.

For our tests, we proceed as follows: First, we get data for 60 randomly selected stocks from the S&P500 from Jan 2008 to Jan 2018. Then, we learn  $\mathbf{w}$  from the past year and buy accordingly for the next six months. Starting from Jan 2008, this procedure is repeated every six months. We calculate the average return, standard deviation and the Sharpe ratio. [19] suggested the Sharpe ratio as a measure for a portfolio's performance, which is the average return earned per unit of volatility and can be calculated by dividing the mean return of a series by its standard deviation. Table 1 shows the performance of the portfolio for different kernels for  $Q = 3$ . For the GP-LVM we chose the linear, SE, EXP and M32 kernels and included the performance given by the sample covariance matrix, i.e.  $\mathbf{K} = \frac{1}{D}(\mathbf{r} - \hat{\boldsymbol{\mu}})(\mathbf{r} - \hat{\boldsymbol{\mu}})^T$ , where  $\hat{\boldsymbol{\mu}}_n = \frac{1}{D} \sum_{d=1}^D \mathbf{r}_{nd}$ , the shrunk Ledoit-Wolf covariance matrix<sup>6</sup> [4] and the equally weighted portfolio, where  $\mathbf{w} = (1, 1, \dots, 1)/N$ .

Table 1: Mean returns, standard deviation and the Sharpe ratio of different models on a yearly basis.

Model	Linear	SE	EXP	M32	Sample Cov	Ledoit Wolf	Eq. Weighted
Mean	0.142	0.151	0.155	0.158	0.149	0.148	0.182
Std	0.158	0.156	0.154	0.153	0.159	0.159	0.232
Sharpe ratio	0.901	0.969	1.008	1.029	0.934	0.931	0.786

Non-linear kernels have the minimal variance and at the same time the best Sharpe ratio values. Note that we are building a minimal variance portfolio and therefore not maximizing the mean returns as explained in section 3.2. For finite  $q$  in equation (16) one can also build portfolios, which not only minimize risk but also maximize the returns. Another requirement for that would be to have a good estimator for the expected return as well.

<sup>6</sup>Here, we have used the implementation in the Python toolbox scikit-learn [20].



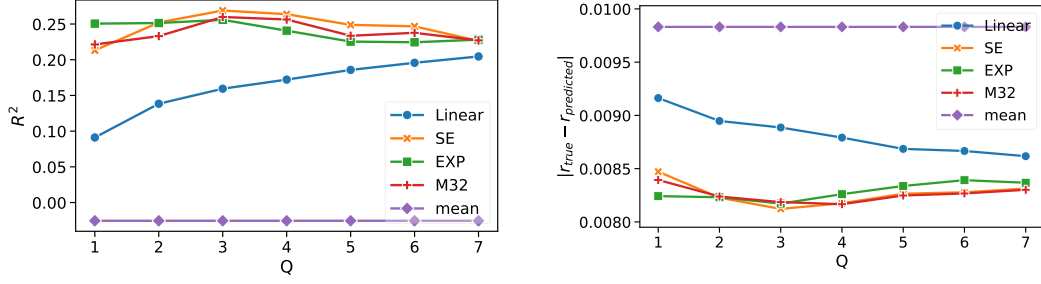


Figure 2: Left:  $R^2$ -score of the predicted values as a function of the latent dimension  $Q$ . Right: The average absolute deviation of suggested return to the real return evaluated by Leaving-one-out cross-validation. Historical mean is indicated by “mean” and is  $Q$ -independent.

#### 4.3.2 Fill in Missing Values

Regulation requires fair value assessment of all assets [21], including illiquid and infrequently traded ones. Here, we demonstrate how the GP-LVM could be used to fill-in missing prices, e.g. if no trading took place. For illustration purposes, we continue working with our daily close prices of stocks from the S&P500, but the procedure can be applied to any other asset class and time resolution.

First, we split our data  $\mathbf{r}$  into training and test set. The test set contains days where the returns of assets are missing and our goal is to accurately predict the return. In our case, we use stock data from Jan 2017 to Oct 2017 for training and Nov 2017 to Dec 2017 to test. The latent space  $\mathbf{B}$  and the hyperparameters  $\theta$  are learned from the training set. Given  $\mathbf{B}$  and  $\theta$ , we can use the standard GP equations (eq. (5) and (6)) to get the posterior return distribution. A suggestion for the value of the return at a particular day can be made by the mean of this distribution. Given  $N$  stocks for a particular day  $d$ , we fit a GP to  $N - 1$  stocks and predict the value of the remaining stock. We repeat the process  $N$  times, each time leaving out a different stock (Leave-one-out cross-validation). Figure 2 shows the  $R^2$ -score (eq. (22)) and the average absolute deviation  $\frac{1}{ND} \sum_{nd} |r_{nd} - r_{nd}^{\text{pred}}|$  of the suggested return to the real return. The average is build over all stocks and all days in the test set.

The predictions with just the historical mean have a negative  $R^2$ -score. The linear model is better than that. But if we switch to non-linear kernels, we can even further increase the prediction score. For  $Q$  between 2 and 4 we obtain the best results. Note that to make a decent suggestion for the return of an asset, there must be some correlated assets in the data as well. Otherwise, the model has no information at all about the asset, we want to predict the returns for.

#### 4.3.3 Interpretation of the Latent Space

The 2-dimensional latent space can be visualized as a scatter plot. For a stationary kernel function like the SE, the distance between the stocks is directly related to their correlation. In this case, the latent positions are even easier to interpret than market betas. As an example, Figure 3 shows the 2-D latent space from 60 randomly selected stocks from the S&P500 from Jan 2017 to Dec 2017. Visually stocks from the same sector tend to cluster together and we consider our method as an alternative to other methods for detecting structure in financial data [22].

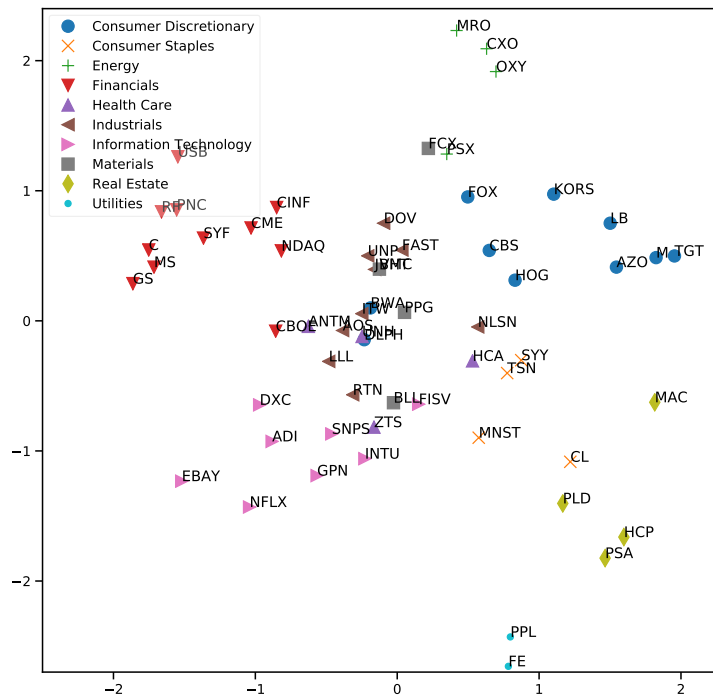


Figure 3: Stocks visualised in the 2-D latent space for the SE kernel.

## 5 Conclusion

We applied the Gaussian Process Latent Variable Model (GP-LVM) to estimate the covariance matrix between different assets, given their time series. We then showed how the GP-LVM can be seen as a non-linear extension to the CAPM with latent factors. Based on the  $R^2$ -score and the ELBO, we concluded, that for fixed latent space dimension  $Q$ , every non-linear kernel can capture more structure than the linear one.

The estimated covariance matrix helps us to build a minimal risk portfolio according to Markowitz Portfolio theory. We evaluated the performance of different models on the S&P500 from year 2008 to 2018. Again, non-linear kernels had lower risk in the suggested portfolio and higher Sharpe ratios than the linear kernel and the baseline measures. Furthermore, we showed how to use the GP-LVM to fill in missing prices of less frequently traded assets and we discussed the role of the latent positions of the assets. In the future, one could also put a Gaussian Process on the latent positions and allow them to vary in time, which would lead to a time-dependent covariance matrix.

## Acknowledgments

The authors thank Dr. h.c. Maucher for funding their positions.

## References

- [1] Harry Markowitz. Portfolio selection. *The Journal of Finance*, 7(1):77–91, 1952.

- [2] J. D. Jobson and Robert M Korkie. Putting markowitz theory to work. *The Journal of Portfolio Management*, 7(4):70–74, 1981.
- [3] William F. Sharpe. Capital asset prices with and without negative holdings. *The Journal of Finance*, 46(2):489–509, 1991.
- [4] Olivier Ledoit and Michael Wolf. Honey, I shrunk the sample covariance matrix. *The Journal of Portfolio Management*, 30(4):110–119, 2004.
- [5] Yuriy Nevmyvaka, Yi Feng, and Michael Kearns. Reinforcement learning for optimized trade execution. In *Proceedings of the 23rd International Conference on Machine Learning, ICML '06*, pages 673–680, New York, NY, USA, 2006. ACM.
- [6] Edward Gately. *Neural Networks for Financial Forecasting*. John Wiley & Sons, Inc., New York, NY, USA, 1995.
- [7] Nicolas Chapados and Yoshua Bengio. Augmented functional time series representation and forecasting with gaussian processes. In J. C. Platt, D. Koller, Y. Singer, and S. T. Roweis, editors, *Advances in Neural Information Processing Systems 20*, pages 265–272. Curran Associates, Inc., 2008.
- [8] J. B. Heaton, N. G. Polson, and J. H. Witte. Deep Portfolio Theory. Papers 1605.07230, arXiv.org, May 2016.
- [9] Yue Wu, José Miguel Hernández-Lobato, and Zoubin Ghahramani. Gaussian process volatility model. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 1044–1052. Curran Associates, Inc., 2014.
- [10] Andrew Gordon Wilson and Zoubin Ghahramani. Generalised wishart processes. In *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence, UAI'11*, pages 736–744, Arlington, Virginia, United States, 2011. AUAI Press.
- [11] Neil Lawrence. Probabilistic non-linear principal component analysis with gaussian process latent variable models. *J. Mach. Learn. Res.*, 6:1783–1816, December 2005.
- [12] Carl Edward Rasmussen. *Gaussian processes for machine learning*. MIT Press, 2006.
- [13] Michael E. Tipping and Chris M. Bishop. Probabilistic principal component analysis. *Journal of the Royal Statistical Society, Series B*, 61:611–622, 1999.
- [14] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, Berlin, Heidelberg, 2006.
- [15] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Wiley-Interscience, New York, NY, USA, 1991.
- [16] Richard Roll and Stephen Ross. The arbitrage pricing theory approach to strategic portfolio planning. *Financial Analysts Journal*, 51:122–131, Jan 1995.
- [17] Stephen A Ross. The arbitrage theory of capital asset pricing. *Journal of Economic Theory*, 13(3):341 – 360, 1976.
- [18] Bob Carpenter, Andrew Gelman, Matthew Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt, Marcus Brubaker, Jiqiang Guo, Peter Li, and Allen Riddell. Stan: A probabilistic programming language. *Journal of Statistical Software, Articles*, 76(1):1–32, 2017.

- [19] William F. Sharpe. Mutual fund performance. *The Journal of Business*, 39(1):119–138, 1966.
- [20] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [21] Financial Accounting Standards Board. Statement of financial accounting standards no. 157 – fair value measurements, 2006. Norwalk, CT.
- [22] Michele Tumminello, Fabrizio Lillo, and Rosario N. Mantegna. Correlation, hierarchies, and networks in financial markets. *Journal of Economic Behavior & Organization*, 75:40–58, 2010.