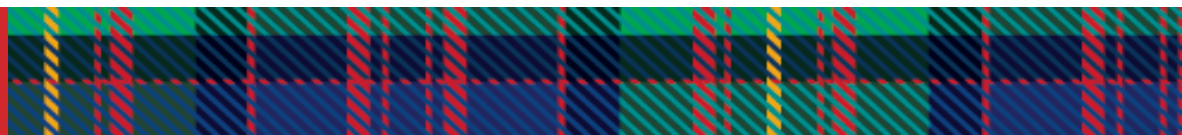


Forecasting and Trading Futures Contract Spreads with Gaussian Process Regression

Xuan Li

April 5, 2018



Outline

- I. Motivation (why GPR)**
- II. Gaussian Process Regression (GPR) in General (Framework)**
- III. Adapting GPR towards Futures Trading Strategy**
 - I.** Kernel Covariance Selection (Augmented Functional Representation)
 - II.** Feature Selection & Engineering & Preprocessing
 - III.** Training (Parameter Learning)
 - IV.** Model Evaluation & Comparison
 - V.** From Prediction to Trading Decision
- IV. Potential Improvements**

I.1 Motivation

1. Definition of futures contracts spread

To profit from the **change in the differential of buying/selling the two related contracts**. Essentially, you **consider the risk in the difference between two prices** (contract) rather than the risk of an immediate futures contract. In other words, **maximizing difference, minimizing risk**.

2. Trading Decision with Risk

Forecast information ratio

$$\widehat{IR} = \frac{\mathbb{E}[\tilde{p}_{t2} - \tilde{p}_{t1} | \mathcal{I}_{t_0}]}{\sqrt{\text{Var}[\tilde{p}_{t2} - \tilde{p}_{t1} | \mathcal{I}_{t_0}]}} \quad \frac{\text{difference}}{\text{risk}}$$

where

$$\begin{aligned} \text{Var}[\tilde{p}_{t2} - \tilde{p}_{t1} | \mathcal{I}_{t_0}] &= \text{Var}[\tilde{p}_{t2} | \mathcal{I}_{t_0}] + \text{Var}[\tilde{p}_{t1} | \mathcal{I}_{t_0}] - 2 \text{Cov}[\tilde{p}_{t1}, \tilde{p}_{t2} | \mathcal{I}_{t_0}] \\ &= \text{Cov}[\tilde{p}_{t2} | \mathcal{I}_{t_0}] + \text{Cov}[\tilde{p}_{t1} | \mathcal{I}_{t_0}] - 2 \text{Cov}[\tilde{p}_{t1}, \tilde{p}_{t2} | \mathcal{I}_{t_0}] \end{aligned}$$

II. Gaussian Processing

Regression (GPR)

1. Property of Gaussian Distribution

1.1.1 Marginal Gaussian Distribution

Suppose $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$ is jointed Gaussian $\sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$$

and

$$\boldsymbol{\Sigma} = \begin{pmatrix} \Sigma_{11}, \Sigma_{12} \\ \Sigma_{21}, \Sigma_{22} \end{pmatrix}$$

Thus, we have

$$\mathbb{P}(x_1|x_2) \sim N(\mu_{1|2}, \Sigma_{1|2})$$

where

$$\mu_{1|2} = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(x_2 - \mu_2)$$

and

$$\Sigma_{1|2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$$

1.1.2 Sum of Gaussian

Given $x_1, x_2 \in \mathbb{R}^d$, $x_1 \perp x_2$, $x_1 \sim N(\mu_1, \Sigma_1)$, $x_2 \sim N(\mu_2, \Sigma_2)$, then we have

$$x_1 + x_2 \sim N(\mu_1 + \mu_2, \Sigma_1 + \Sigma_2)$$

II. Gaussian Processing Regression (GPR)

2. Inference of GPR

Let \mathbf{X} be the matrix of training input, \mathbf{y} the vector of output. Assume

$$y_i = f(\mathbf{x}_i) + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma_n^2)$$

and the functional operator

$$\mathbf{f} = [f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_N)] \sim N(\mathbf{0}, K(\mathbf{X}, \mathbf{X})) \quad K(\mathbf{X}, \mathbf{X})_{ij} = k(\mathbf{X}_i, \mathbf{X}_j)$$

where $k(\mathbf{u}, \mathbf{v})$ is kernel covariance operator that for example

$$k(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|^2}{2\sigma_l^2}\right)$$



Thus, the regression problem is that given \mathbf{X} , \mathbf{y} , and \mathbf{X}_* , how could we estimate $\mathbb{E}[\mathbf{y}_* | \mathbf{y}]$ and $\text{Cov}[\mathbf{y}_* | \mathbf{y}]$?

II. Gaussian Processing

Regression (GPR)

2. Inference of GPR, cont

First, use '*sum of Gaussian*' property, suppose $\mathbf{y} \in \mathbb{R}^N, \mathbf{y}_* \in \mathbb{R}^M$, we know that

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{y}_* \end{pmatrix} = \begin{pmatrix} f(\mathbf{X}) \\ f(\mathbf{X}_*) \end{pmatrix} + \begin{pmatrix} \sigma_n^2 \mathbf{1}_N \\ \sigma_n^2 \mathbf{1}_M \end{pmatrix} \sim N\left(\mathbf{0}, \begin{pmatrix} K(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I}_N & K(\mathbf{X}, \mathbf{X}_*) \\ K(\mathbf{X}_*, \mathbf{X}) & K(\mathbf{X}_*, \mathbf{X}_*) + \sigma_n^2 \mathbf{I}_M \end{pmatrix}\right)$$

Second, use '*Marginal Gaussian Distribution*' property, we have

$$\mathbb{E}[\mathbf{y}_* | \mathbf{y}] = K(\mathbf{X}_*, \mathbf{X})(K(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{y}$$

and

$$\text{Cov}[\mathbf{y}_* | \mathbf{y}] = K(\mathbf{X}_*, \mathbf{X}_*) - K(\mathbf{X}_*, \mathbf{X})(K(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I}_N)^{-1} K(\mathbf{X}, \mathbf{X}_*)$$

III. Adapting GPR towards Futures Trading Strategy

1. Kernel covariance selection (Augmented functional representation)

The key idea of GPR lies in the choice of kernel covariance operator (function), which models the relationship between corresponding observations (input, x). In this paper, the authors choose:

$$k(u, v | \ell, \alpha, \sigma_f, \sigma_{TS}) = \sigma_f^2 \left(1 + \frac{1}{2\alpha} \sum_{k=1}^d \frac{(u_k - v_k)^2}{\ell_k^2} \right)^{-\alpha} + \sigma_{TS}^2 \delta(u_i, v_i)$$

where $\delta(u_i, v_i)$ (or δ_{u_i, v_i}) is the Kronecker delta that

$$\delta(u_i, v_i) = \begin{cases} 1, & u_i = v_i \\ 0, & u_i \neq v_i \end{cases}$$

The authors considered the covariance based on contribution of each dimension of the input features separately, and the effect of a dominant feature using Kronecker delta function.

III. Adapting GPR towards Futures Trading Strategy

2. Feature Selection & Engineering & Preprocessing

In all, the feature used in the paper include:

- current time-series index or year: the Kronecker delta term is applied on this dimension
- operation time: the time at which the forecast is made
- target time: the time at which the forecast is being made
- current spread price
- the price of the three nearest future contracts
- stock-to-use ratio:
- year-over-year difference in total ending stocks

All of these features are standardized with zero mean and unit deviation. For the price target y , the trajectory is standardized to start at zero at the start of every year, by subtracting the first price value.

III. Adapting GPR towards Futures Trading Strategy

3. Training (parameter learning)

Given the kernel covariance function, the total hyper-parameter needed to be learned are

- $\alpha \in \mathbb{R}$
- $\sigma_f \in \mathbb{R}$
- $\sigma_{TS} \in \mathbb{R}$
- $\sigma_n \in \mathbb{R}$
- $\ell \in \mathbb{R}^d$

$$k(u, v | \ell, \alpha, \sigma_f, \sigma_{TS}) = \sigma_f^2 \left(1 + \frac{1}{2\alpha} \sum_{k=1}^d \frac{(u_k - v_k)^2}{\ell_k^2} \right)^{-\alpha} + \sigma_{TS}^2 \delta(u_i, v_i)$$

As the authors have claimed in the paper, one efficient approach to find the hyper-parameter for Gaussian Process is to maximize the **marginal likelihood** of the observed data (given the hyper-parameter). To be specific, define θ as the hyper-parameter, $\mathbf{y} \in \mathbb{R}^N$ are the observation, and $\mathbf{X} \in \mathbb{R}^{N \times d}$ are the input features, then

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \log \mathbb{P}(\mathbf{y} | \mathbf{X}; \theta)$$

Since from Part one we know that $\mathbb{P}(\mathbf{y} | \mathbf{X}; \theta)$ follows Multivariate Gaussian (normal) Distribution $N(\mathbf{0}, \Sigma)$ ($\Sigma = K(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I}_N$) given the kernel function $k(u, v; \theta)$, which is

III. Adapting GPR towards Futures Trading Strategy

3. Training (parameter learning) cont.

$$\mathbb{P}(\mathbf{y}|\mathbf{X};\boldsymbol{\theta}) = \frac{1}{\sqrt{(2\pi)^N|\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}\mathbf{y}^T\boldsymbol{\Sigma}^{-1}\mathbf{y}\right)$$

where $|\boldsymbol{\Sigma}| = \det(\boldsymbol{\Sigma})$. Thus, to put it back to equation above, we have that

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \quad -\frac{1}{2}\mathbf{y}^T(K(\mathbf{X}, \mathbf{X};\boldsymbol{\theta}) + \sigma_n^2\mathbf{I}_N)^{-1}\mathbf{y} - \frac{1}{2}\log|K(\mathbf{X}, \mathbf{X};\boldsymbol{\theta}) + \sigma_n^2\mathbf{I}_N| - \frac{N}{2}\log(2\pi)$$

which is equivalent to

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \quad -\frac{1}{2}\mathbf{y}^T(K(\mathbf{X}, \mathbf{X};\boldsymbol{\theta}) + \sigma_n^2\mathbf{I}_N)^{-1}\mathbf{y} - \frac{1}{2}\log|K(\mathbf{X}, \mathbf{X};\boldsymbol{\theta}) + \sigma_n^2\mathbf{I}_N|$$

III. Adapting GPR towards Futures Trading Strategy

3. Training (parameter learning) cont.

Gradient Descent

Given the expression of objective and form of variables, the problem is a non convex optimization problem and could be potentially solved by gradient descent. First, we can get the gradient of the objective function as

$$\frac{\partial \log \mathbb{P}(y|X; \theta)}{\partial \theta} = \frac{1}{2} y^T \Sigma^{-1} \frac{\partial \Sigma}{\partial \theta} \Sigma^{-1} y - \frac{1}{2} \text{tr}(\Sigma^{-1} \cdot \frac{\partial \Sigma}{\partial \theta})$$

Thus, the key component to get the gradient is to calculate the matrix $\frac{\partial \Sigma}{\partial \theta}$. Given the definition in Section 2.2.2, and

$$\Sigma = K(X, X) + \sigma_n^2 I_N$$

for each of the parameter in θ , we have

$$\left(\frac{\partial \Sigma}{\partial \alpha}\right)_{i,j} = \sigma_f^2 \left(1 + \frac{c}{2\alpha}\right)^{-\alpha} \cdot \left(-\ln\left(1 + \frac{c}{2\alpha}\right) + \frac{1}{1 + 2\alpha/c}\right), \quad c = \sum_{k=1}^d \frac{(X_k^i - X_k^j)^2}{\ell_k^2}$$

$$\left(\frac{\partial \Sigma}{\partial \sigma_f}\right)_{i,j} = 2\sigma_f \left(1 + \frac{1}{2\alpha} \sum_{k=1}^d \frac{(X_k^i - X_k^j)^2}{\ell_k^2}\right)^{-\alpha}$$

$$\left(\frac{\partial \Sigma}{\partial \sigma_{TS}}\right)_{i,j} = 2\sigma_{TS} \delta(X_d^i, X_d^j)$$

III. Adapting GPR towards Futures Trading Strategy

3. Training (parameter learning) cont.

Gradient Descent

$$\left(\frac{\partial \Sigma}{\partial \sigma_n}\right)_{i,j} = 2\sigma_n \delta(i,j)$$

$$\left(\frac{\partial \Sigma}{\partial \ell_d^2}\right)_{i,j} = \sigma_f^2 \left(1 + \frac{1}{2\alpha} \sum_{k=1}^d \frac{(X_k^i - X_k^j)^2}{\ell_k^2}\right)^{-\alpha-1} \cdot \frac{(X_d^i - X_d^j)^2}{\ell_d^3}$$

Algorithm 1: Gradient Descent (Gaussian Process)

Data: X, y

Parameter: θ

initialization: $\text{epoch}^{\max}, \text{epoch} = 0, \theta, \omega, \epsilon$;

while $\text{epoch} < \text{epoch}^{\max}$ do

$\theta_i = \theta_i + \omega \cdot \frac{\partial \log \mathbb{P}(y|X;\theta)}{\partial \theta_i}$;

 if $|\log \mathbb{P}(y|X;\theta) - \log \mathbb{P}(y|X;\theta')| > \epsilon$ then

 epoch ++;

 else

 break;

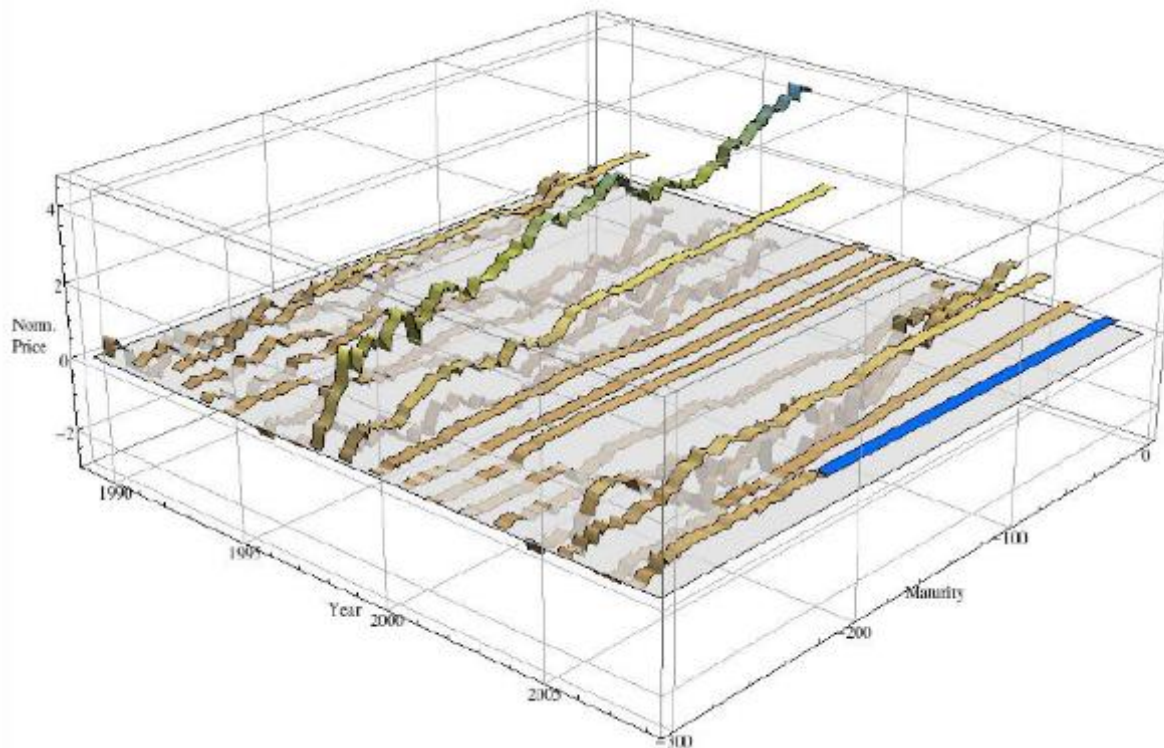
 end

end

III. Adapting GPR towards Futures Trading Strategy

3. Training (parameter learning) cont.

multi-stage training + sub-sampling



III. Adapting GPR towards Futures Trading Strategy

4. Model Evaluation & Comparison between models

4.1 Significance Test

cross-covariance-corrected Diebold-Mariano test for sequential correlated sequence (error trend), with the assumption of covariance stationarity.

Let d_t be the sequence of error difference between two models to be compared. The tested statistic

$$DM = \bar{d} / \sqrt{\hat{v}_{\text{CCC-DM}}}$$

is asymptotically distributed as $N(0, 1)$, where

$$\bar{d} = \frac{1}{M} \sum_t d_t$$

and

$$\hat{v}_{\text{CCC-DM}} = \frac{1}{M^2} \left(\sum_i M_i \sum_{k=-K}^K \hat{\gamma}_k^i + \sum_i \sum_{j \neq i} M_{i \cap j} \sum_{k=-K'}^{K'} \hat{\gamma}_k^{i,j} \right)$$

where M_i is the number of samples in test set i , M is the total number of samples, $\hat{\gamma}_k^i$ is the estimated lag- k autocovariance, and $\hat{\gamma}_k^{i,j}$ is the estimated lag- k cross-covariance.

III. Adapting GPR towards Futures Trading Strategy

4. Model Evaluation & Comparison, cont.

4.2 Criterion

The criterion used are square error (SE) and negative log-likelihood (NLL). SE is normalized by dividing it with the standard deviation of the test target. NLL is normalized by subtracting the likelihood of a univariate Gaussian distribution estimated on the test target.

4.3 Baseline Models

- AugRQ/all-inp: proposed model
- AugRQ/less-inp: do not include the economic variables
- AugRQ/no-inp: only with time relevant variables
- stdRQ/all-inp: single length parameter
- stdRQ/no-inp: single time variable Δ
- Linear/all-inp: Bayesian linear regression

III. Adapting GPR towards Futures Trading Strategy

5. From Prediction to Trading Decision

5.1 Prediction

holding the time series index constant to N , the operation time constant to the time M_N of the last observation, the other input variable constant to their last-observed values $\mathbf{x}_{M_N}^N$ (slow-moving variables that represent a “level”), and varying the target time over the forecasting period Δ .

5.2 Trading Decision

Forecast information ratio

$$\widehat{IR} = \frac{\mathbb{E}[\tilde{p}_{t2} - \tilde{p}_{t1} | \mathcal{I}_{t_0}]}{\sqrt{\text{Var}[\tilde{p}_{t2} - \tilde{p}_{t1} | \mathcal{I}_{t_0}]}}$$

where

$$\text{Var}[\tilde{p}_{t2} - \tilde{p}_{t1} | \mathcal{I}_{t_0}] = \text{Var}[\tilde{p}_{t2} | \mathcal{I}_{t_0}] + \text{Var}[\tilde{p}_{t1} | \mathcal{I}_{t_0}] - 2 \text{Cov}[\tilde{p}_{t1}, \tilde{p}_{t2} | \mathcal{I}_{t_0}]$$

IV. Potential Improvements

- **Scalability** with large dataset, and efficient training algorithm. Implemented in the paper: random sampling subset + retraining, might introduce bias of prediction accuracy towards the selected input. Potential refinement:

(a) approximate Σ with lower rank matrix [2,3]. For example,

$$\Sigma \approx \hat{\Sigma} = \Sigma[:, \mathcal{I}] \Sigma[\mathcal{I}, \mathcal{I}]^{-1} \Sigma[\mathcal{I}, :]$$

- **Assumption** in prediction: the slow-moving assumption of features. “for slow-moving variables that represent a level, one can conceivably keep their value constant to the last known realization across the forecasting period”. Unused features, could be calibrated in the kernel function as well.
- **Feature selection**: any better, or more (temporal) features that are not ‘slow-moving’? any better modelling of kernel?
- **extro-covariance** with other type of futures contract.