

## A Bayesian Nonlinear Support Vector Machine Error Correction Model

TONY VAN GESTEL,<sup>1,2\*</sup> MARCELO ESPINOZA,<sup>2</sup>  
BART BAESENS,<sup>3</sup> JOHAN A. K. SUYKENS,<sup>2</sup>  
CARINE BRASSEUR<sup>4</sup> AND BART DE MOOR<sup>2</sup>

<sup>1</sup> *Dexia Group, Belgium*

<sup>2</sup> *Katholieke Universiteit Leuven, Belgium*

<sup>3</sup> *School of Management, University of Southampton, UK*

<sup>4</sup> *Fortis Bank Brussels, Belgium*

### ABSTRACT

The use of linear error correction models based on stationarity and cointegration analysis, typically estimated with least squares regression, is a common technique for financial time series prediction. In this paper, the same formulation is extended to a nonlinear error correction model using the idea of a kernel-based implicit nonlinear mapping to a high-dimensional feature space in which linear model formulations are specified. Practical expressions for the nonlinear regression are obtained in terms of the positive definite kernel function by solving a linear system. The nonlinear least squares support vector machine model is designed within the Bayesian evidence framework that allows us to find appropriate trade-offs between model complexity and in-sample model accuracy. From straightforward primal–dual reasoning, the Bayesian framework allows us to derive error bars on the prediction in a similar way as for linear models and to perform hyperparameter and input selection. Starting from the results of the linear modelling analysis, the Bayesian kernel-based prediction is successfully applied to out-of-sample prediction of an aggregated equity price index for the European chemical sector. Copyright © 2006 John Wiley & Sons, Ltd.

**KEY WORDS** financial time series prediction; least squares support vector machines; Bayesian inference; error correction mechanism; kernel-based learning

### INTRODUCTION

Financial time series forecasting is a dynamic field with important contributions coming from many disciplines. The issue of forecasting financial time series has traditionally been seen as a difficult task, since the data generating process is dominated by stochastic rather than deterministic components. Although the efficient markets hypothesis (Bachelier, 1900; Fama, 1965) states that stock returns are unpredictable, recent modelling techniques seem to suggest that stock returns are pre-

\* Correspondence to: Tony Van Gestel, Credit Risk Modelling, Risk Management, Dexia Group, Square Meeus 1, B-1000 Brussels, Belgium. E-mail: tony.vangestel@dexia.com

dictable to some degree (Brock *et al.*, 1992; Campbell *et al.*, 1997; Lo *et al.*, 2000; Sullivan *et al.*, 1999).

Within the linear framework, usually forecasting models are based on stationarity considerations of the series at hand (Campbell *et al.*, 1997; Granger and Newbold, 1986; Hamilton, 1994). In this sense, evidence of possible cointegration can be exploited into an error correction mechanism formulation. This model is typically estimated by ordinary least squares, applying input selection to control the model complexity and avoid overfitting on the training set.

The universal approximation property of multilayer perceptrons (MLPs) motivated their use for nonlinear financial time series prediction (Granger and Terasvirta, 1993; Hutchinson *et al.*, 1994; Refenes and Zapranis, 1999). While powerful design techniques like the Bayesian evidence framework (MacKay, 1995) have been developed, the practical use of neural networks suffers from drawbacks like the nonconvex optimization problem with multiple local minima and the choice of the number of hidden neurons. In support vector machines (SVMs), least squares support vector machines (LS-SVMs) and related kernel-based prediction techniques (Schölkopf and Smola, 2002; Suykens *et al.*, 2002; Vapnik, 1998), the solution follows from a convex optimization problem. Basically these methods map the inputs in a nonlinear way, first into a high kernel-induced feature space, in which ridge regression is applied in the case of LS-SVMs. The solution follows from a linear Karush–Kuhn–Tucker system in the dual space in terms of the positive definite kernel function by applying Mercer’s theorem (Schölkopf and Smola, 2002; Suykens *et al.*, 2002; Vapnik, 1998).

In this paper, a nonlinear error correction model (ECM) formulation is estimated using LS-SVMs to predict an aggregated equity price index for the European chemical sector. First a stationarity and cointegration analysis is performed to define a good linear model formulation. This model is used as a starting point to design the nonlinear kernel-based model within the Bayesian evidence framework. The parameters and inputs of the LS-SVM are estimated<sup>1</sup> in the Bayesian evidence framework (Van Gestel *et al.*, 2001, 2002). The Bayesian framework embodies Occam’s razor to find an optimal trade-off between training set accuracy and model complexity in a similar way as the Akaike and Bayesian information criteria (Akaike, 1974; Schwarz, 1978). The linear and nonlinear predictions are compared in terms of directional accuracy and market timing ability.

This paper is organized as follows. The initial problem definition, stationarity analysis and linear model specification are reviewed and applied in the next section. The design and application of nonlinear kernel-based regression within the evidence framework is presented in the third section. The final results are discussed in the fourth section.

## LINEAR MODELLING

In financial engineering applications, the importance of having good forecasting modelling tools is straightforward. Risk and portfolio management are mainly based on such tools, therefore any improvement over traditional techniques can lead to competitive advantages. Traditional forecasting based on linear models is built upon the concepts of stationarity and cointegration.

### Stationarity

A linear model formulation to predict an output  $y \in \mathbb{R}$  based on  $n$  explanatory input variables  $\mathbf{x} = [x_1; \dots; x_n] = [x_1, \dots, x_n]^T \in \mathbb{R}^n$  can be written as

<sup>1</sup>A Matlab toolbox for the LS-SVM formulation and Bayesian inference is available from <http://www.esat.kuleuven.ac.be/sista/lssvmlab>.

$$y = \mathbf{w}^T \mathbf{x} + b + e \quad (1)$$

with  $\mathbf{w} \in \mathbb{R}^n$  is a coefficient vector and  $b \in \mathbb{R}$  a bias term. Having a set of  $n_D$  observations  $\mathcal{D} = \{(\mathbf{x}_b, y_i)\}_{i=1}^{n_D}$ , the most usual technique to estimate a linear model is by using the ordinary least squares (OLS) estimator

$$\min_{\mathbf{w}, b} \sum_{i=1}^{n_D} (y_i - (\mathbf{w}^T \mathbf{x}_i + b))^2 \quad (2)$$

with error  $e_i = y_i - (\mathbf{w}^T \mathbf{x}_i + b)$ . Defining  $\mathbf{y} = [y_1; y_2; \dots; y_{n_D}] = [y_1, y_2, \dots, y_{n_D}]^T \in \mathbb{R}^{n_D}$ ,  $\mathbf{1} = [1, 1, \dots, 1]^T \in \mathbb{R}^{n_D}$  and  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n_D}]^T \in \mathbb{R}^{n_D \times n}$  the solution to (2) is obtained from the linear set of equations

$$[\hat{\mathbf{w}}; \hat{b}] = ([\mathbf{X}, \mathbf{1}]^T [\mathbf{X}, \mathbf{1}])^{-1} [\mathbf{X}, \mathbf{1}]^T \mathbf{y} \quad (3)$$

Although the general assumptions underlying the application of the OLS are described in the Gauss–Markov conditions, in the particular scope of time series forecasting it is required that the series involved should be stationary. It has been widely recognized that performing linear regression with nonstationary series has the potential to lead to serious inference errors (Granger and Newbold, 1974). Some of the known problems when performing OLS estimations with nonstationary series are, for instance, the identification of spurious relationships between unrelated variables, or the non-convergence of the  $\hat{\mathbf{w}}$  estimates (Maddala and Kim, 1998). Formally, for a process  $y_t$  to be (weakly) stationary, it must satisfy the following set of properties:  $E[y_t] = \mu_y$ ,  $E[(y_t - \mu_y)^2] = \text{var}(y_t) = \sigma_y^2 = \gamma(0)$ ,  $E[(y_t - \mu_y)(y_{t-\tau} - \mu_y)] = \text{cov}(y_t, y_{t-\tau}) = \gamma(\tau)$ , where the mean and variance of  $y_t$  are constant, and the covariances depend only on the time *interval*  $\tau$  and not on the particular moment of time  $t$ .

One of the most common tests for stationarity of a time series  $y_t$  is based on the so-called augmented Dickey–Fuller (ADF) regression (Dickey and Fuller, 1979):  $\Delta y_t = \alpha + \rho y_{t-1} + \sum_{j=1}^q \beta_j \Delta y_{t-j} + e_t$ . Under the null hypothesis of nonstationarity ( $H_0: \rho = 1$ ), the  $t$ -statistic of the estimated coefficient  $\hat{\rho}$  will follow a nonstandard distribution, usually known as the Dickey–Fuller (DF) distribution. If the corresponding  $t$ -statistic for the coefficient  $y_{t-1}$  is above the critical value of the ADF test, then the null hypothesis of nonstationarity cannot be rejected (Rao, 1994).

### Cointegration

If a series  $y_t^{(l)}$  of levels is found to be nonstationary on its original levels, one usual transformation is to take first differences and work with the transformed variable  $y_t^{(d)} = \Delta y_t^{(l)} = y_t^{(l)} - y_{t-1}^{(l)}$ . However, before attempting to transform all nonstationary variables into first differences, it is useful to explore for possible cointegration between the dependent variable and any subset of the explanatory variables. For the case of two nonstationary variables  $y_t^{(l)}$ ,  $x_t^{(l)}$  testing for cointegration involves testing for stationarity of the residuals in the regression

$$y_t^{(l)} = \beta_0 + \beta_1 x_t^{(l)} + e_t \quad (4)$$

Thus, finding stationary residuals from the regression above is equivalent to finding a cointegrating relationship between the variables, where the stationary cointegrating linear combination can be estimated as  $z_t = y_t^{(l)} - \hat{\beta}_0 - \hat{\beta}_1 x_t^{(l)}$ .

If cointegration exists, then it is possible to take advantage of this long-term relationship and use it to model the short-term behaviour of the system. It was proved that any cointegrating system can have an equivalent ECM representation (Engle and Granger, 1987). Following from the example above, if the series  $x_t^{(d)}$  and  $y_t^{(d)}$  do cointegrate, then their corresponding linear ECM equals

$$y_t^{(d)} = \alpha_0 + \alpha_1 z_{t-1} + \sum_{j=1}^p b_j y_{t-j}^{(d)} + \sum_{j=1}^p c_j x_{t-j}^{(d)} + e_t \quad (5)$$

which can be written more generally as

$$y_t^{(d)} = f(z_{t-1}; y_{t-1}^{(d)}, \dots, y_{t-p}^{(d)}; x_{t-1}^{(d)}, \dots, x_{t-p}^{(d)}) + e_t \quad (6)$$

It is possible to also include some additional external variables that can help to improve the model, but the central concept of the ECM is as shown above. The extension<sup>2</sup> using more than two variables is straightforward.

The number of lags  $p$  for an autoregressive AR( $p$ ) model can be heuristically defined based on the partial autocorrelation function  $\text{pacf}(p)$  (Hamilton, 1994). For an autoregressive formulation of the stationary series  $y(t)^{(d)}$  of order  $p$  (like the ECM above), it can be shown that the  $\text{pacf}(p)$  function will drop to zero after  $\rho > p$  (Box and Jenkins, 1970).

## NONLINEAR KERNEL-BASED MODELLING AND PREDICTION

### Least squares support vector machines

A straightforward way to extend the linear models (1), (5) and (19) to a nonlinear model is to pre-process the inputs  $\mathbf{x}$  in a nonlinear way by a mapping

$$\boldsymbol{\varphi}: \mathbb{R}^n \rightarrow \mathbb{R}^{n_D}: \mathbf{x} \mapsto \boldsymbol{\varphi}(\mathbf{x}) \quad (7)$$

where the feature vector  $\boldsymbol{\varphi}(\mathbf{x})$  is typically high (or even infinite)-dimensional. Given the nonlinear mapping, the ECM model (6) is assumed to be of the following form:

$$y = \mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}) + b + e \quad (8)$$

where  $f(\mathbf{x}) \simeq \mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}) + b$ . Given this nonlinear mapping, the coefficient vector  $\mathbf{w}$  is estimated by solving the (regularized) least squares problem in the primal or feature space:

$$\min_{\mathbf{w}, b, e} J_1(\mathbf{w}, b) = \frac{\mu}{2} \mathbf{w}^T \mathbf{w} + \frac{\zeta}{2} \sum_{t=1}^{n_D} e_t^2 \quad (9)$$

$$\text{s.t. } e_t = y_t - (\mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}_t) + b), \quad t = 1, \dots, n_D \quad (10)$$

<sup>2</sup>The model depicted so far is known as the Engle–Granger two-step approach for cointegration. When using vector formulations, the so-called Johansen procedure (Johansen, 1988) is applied.

As the nonlinear mapping  $\boldsymbol{\varphi}$  is high-dimensional, the regularization term  $\frac{\mu}{2}\mathbf{w}^T\mathbf{w}$  is introduced to avoid overfitting the training data. The parameters  $\mu$  and  $\zeta$  determine the trade-off between regularization  $J_w = (1/2)\mathbf{w}^T\mathbf{w}$  and error minimization  $J_e = (1/2)\sum_{i=1}^{n_D} e_i^2$ .

A key element of support vector machines and kernel-based learning methods is that the nonlinear mapping  $\boldsymbol{\varphi}(\mathbf{x})$  is not explicitly known. Instead it is implicitly defined from Mercer's theorem in terms of the positive-definite kernel function

$$K(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\varphi}(\mathbf{x}_i)^T \boldsymbol{\varphi}(\mathbf{x}_j) \tag{11}$$

Some commonly used kernel functions are

1.  $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$  (linear kernel)
  2.  $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j / c)^d$  (polynomial kernel of degree  $d$ )
  3.  $K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{\sigma^2}\right)$  (radial basis function kernel)
- (12)

where  $d \in \mathbb{N}$  and  $c, \sigma \in \mathbb{R}^+$  are tunable parameters.

In order to solve the constrained optimization problem (9) and (10), one constructs the Lagrangian

$$\mathcal{L}(\mathbf{w}, \mathbf{e}, \boldsymbol{\alpha}, b) = \frac{\mu}{2} \mathbf{w}^T \mathbf{w} + \frac{\zeta}{2} \sum_{i=1}^{n_D} e_i^2 + \sum_{i=1}^{n_D} \alpha_i (y_i - (\mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}_i) + b) - e_i) \tag{13}$$

where the scalars  $\alpha_i \in \mathbb{R}$  are the Lagrange multipliers associated with the equality constraints (10) and are called support values. The conditions for optimality are

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0 \leftrightarrow \mu \mathbf{w} = \sum_{i=1}^{n_D} \alpha_i \boldsymbol{\varphi}(\mathbf{x}_i) & \leftrightarrow \mu \mathbf{w} - \boldsymbol{\Phi}^T \boldsymbol{\alpha} = 0 \\ \frac{\partial \mathcal{L}}{\partial b} = 0 \leftrightarrow \sum_{i=1}^{n_D} \alpha_i = 0 & \leftrightarrow \boldsymbol{\alpha}^T \mathbf{1} = 0 \\ \frac{\partial \mathcal{L}}{\partial e_i} = 0 \leftrightarrow \alpha_i = \zeta e_i, \quad (i = 1, \dots, n_D) & \leftrightarrow \boldsymbol{\alpha} - \zeta \mathbf{e} = 0 \\ \frac{\partial \mathcal{L}}{\partial \alpha_i} = 0 \leftrightarrow e_i = y_i - \mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}_i) - b, \quad (i = 1, \dots, n_D) & \leftrightarrow \boldsymbol{\Phi} \mathbf{w} + b \mathbf{1} + \mathbf{e} = \mathbf{y} \end{cases} \tag{14}$$

in which we have defined  $\boldsymbol{\Phi} = [\boldsymbol{\varphi}(\mathbf{x}_1), \dots, \boldsymbol{\varphi}(\mathbf{x}_{n_D})]^T \in \mathbb{R}^{n_D \times n_\varphi}$ ,  $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_{n_D}]^T \in \mathbb{R}^{n_D}$ ,  $\mathbf{e} = [e_1, \dots, e_{n_D}]^T \in \mathbb{R}^{n_D}$ ,  $\mathbf{y} = [y_1, \dots, y_{n_D}]^T \in \mathbb{R}^{n_D}$  and  $\mathbf{1} = [1, \dots, 1]^T \in \mathbb{R}^{n_D}$ . For the linear case (e.g. a linear kernel) one typically has  $n_\varphi = n \ll n_D$  and after elimination of  $\mathbf{e}$  and  $\boldsymbol{\alpha}$ , one solves the  $(n_\varphi + 1) \times (n_\varphi + 1)$  linear system in the primal space

$$\begin{bmatrix} \boldsymbol{\Phi}^T \boldsymbol{\Phi} + \frac{\mu}{\zeta} \mathbf{I}_{n_\varphi} & \boldsymbol{\Phi}^T \mathbf{1} \\ \mathbf{1}^T \boldsymbol{\Phi} & n_D \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} = \begin{bmatrix} \boldsymbol{\Phi}^T \mathbf{y} \\ \mathbf{1}^T \mathbf{y} \end{bmatrix} \tag{15}$$

In nonlinear kernel-based regression, one usually has  $n_\phi \gg n_D$  and, moreover, the feature vector  $\boldsymbol{\phi}(\mathbf{x})$  is only implicitly defined in terms of the kernel function  $K$  from (11). Eliminating  $\mathbf{w}$  and  $\mathbf{e}$  from (14), one obtains the linear Karush–Kuhn–Tucker (KKT) system of dimension  $(n_D + 1) \times (n_D + 1)$  in the dual space (Suykens *et al.*, 2002)

$$\begin{bmatrix} \frac{1}{\mu} \boldsymbol{\Omega} + \frac{1}{\zeta} \mathbf{I}_{n_D} & \mathbf{1} \\ \mathbf{1}^T & 0 \end{bmatrix} \begin{bmatrix} \boldsymbol{\alpha} \\ b \end{bmatrix} = \begin{bmatrix} \mathbf{y} \\ 0 \end{bmatrix} \quad (16)$$

where the Mercer condition (11) is applied in the matrix  $\boldsymbol{\Omega} = \boldsymbol{\Phi}\boldsymbol{\Phi}^T \in \mathbb{R}^{n_D \times n_D}$  with elements  $\Omega_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$ ,  $i, j = 1, \dots, n_D$  and guarantees that  $\boldsymbol{\Omega} \geq 0$ . The primal–dual formulations also allow us to make extensions to nonlinear generalized least squares regression in a straightforward way, as typically used in financial forecasting (Campbell *et al.*, 1997).

Given the support values  $\boldsymbol{\alpha}$  and bias term  $b$ , one obtains the predicted value  $\hat{y}$  corresponding to a new input  $\mathbf{x}$  as a weighted sum of the kernel functions evaluated in the new data point and the training data points:

$$\hat{y} = \hat{\mathbf{w}}^T \boldsymbol{\phi}(\mathbf{x}) + \hat{b} = \frac{1}{\mu} \sum_{i=1}^{n_D} \alpha_i K(\mathbf{x}, \mathbf{x}_i) + \hat{b} \quad (17)$$

### Bayesian inference for model design

Given the primal–dual formulations, it is clear how to estimate the model parameters  $\hat{\mathbf{w}}$ ,  $\hat{b}$  and point prediction  $\hat{y}$ . However, the regularization and kernel function parameters still have to be tuned from the given training data. In this subsection, the design is done within the Bayesian evidence framework (Van Gestel *et al.*, 2001, 2002), depicted in Figure 1.

The model parameters  $\mathbf{w}$ ,  $b$ , the hyperparameters  $\mu$ ,  $\zeta$  and model structure  $\mathcal{H}$  (corresponding, e.g., to the input set and/or tunable kernel parameters) are inferred by applying Bayes' formula on three different levels:

1. On the first level, it is assumed that the hyperparameters  $\mu$ ,  $\zeta$  and model  $\mathcal{H}$  are given. Applying Bayes' formula, the posterior probability of the model parameters  $\mathbf{w}$  and  $b$  is obtained:

$$p(\mathbf{w}, b | \mathcal{D}, \log \mu, \log \zeta, \mathcal{H}) = \frac{p(\mathcal{D} | \mathbf{w}, b, \log \mu, \log \zeta, \mathcal{H}) p(\mathbf{w}, b | \log \mu, \log \zeta, \mathcal{H})}{p(\mathcal{D} | \log \mu, \log \zeta, \mathcal{H})}$$

The evidence shows that  $p(\mathcal{D} | \log \mu, \log \zeta, \mathcal{H})$  does not depend upon the model parameters  $\mathbf{w}$  and  $b$  and is a normalizing constant such that the left-hand side is a probability density function  $\iint \dots \int p(\mathbf{w}, b | \mathcal{D}, \log \mu, \log \zeta, \mathcal{H}) d\mathbf{w}_1 \dots d\mathbf{w}_{n_\phi} db = 1$ . The ridge regression cost function (9) and (10) is obtained by taking the negative logarithm of the posterior  $p(\mathbf{w}, b | \mathcal{D}, \log \mu, \log \zeta, \mathcal{H})$  using proper choices for the prior  $p(\mathbf{w}, b | \log \mu, \log \zeta, \mathcal{H})$  and the likelihood  $p(\mathcal{D} | \mathbf{w}, b, \log \mu, \log \zeta, \mathcal{H})$ . In the dual space, the parameters  $\boldsymbol{\alpha}$  and  $b$  are obtained from the linear KKT-system (16).

2. The hyperparameters  $\mu$  and  $\zeta$  are inferred on the second level:

$$p(\log \mu, \log \zeta | \mathcal{D}, \mathcal{H}) = \frac{p(\mathcal{D} | \log \mu, \log \zeta, \mathcal{H}) p(\log \mu, \log \zeta, \mathcal{H})}{p(\mathcal{D} | \mathcal{H})}$$

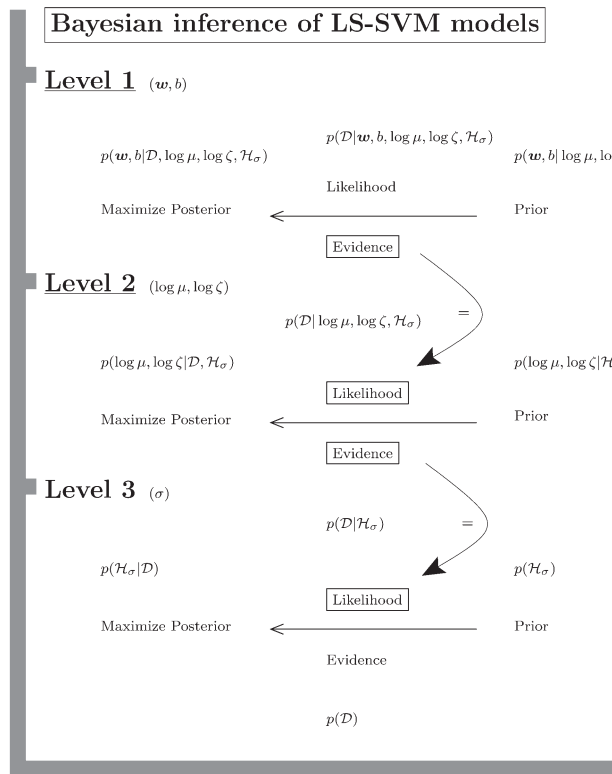


Figure 1. Different levels of Bayesian inference. The posterior probability of the model parameters  $w$  and  $b$  is inferred from the data  $\mathcal{D}$  by applying Bayes' formula on the first level for given hyperparameters  $\mu$  (prior) and  $\zeta$  (likelihood) and the model structure  $\mathcal{H}$ . The model parameters are obtained by maximizing the posterior. The evidence on the first level becomes the likelihood on the second level when applying Bayes' formula to infer  $\mu$  and  $\zeta$  (with  $\gamma = \zeta/\mu$ ) from the given data  $\mathcal{D}$ . The optimal hyperparameters  $\mu_{MP}$  and  $\zeta_{MP}$  are obtained by maximizing the corresponding posterior on level 2. Model comparison is performed on the third level in order to compare different model structures, e.g., with different candidate input sets and/or different kernel parameters

Taking the negative logarithm of the posterior  $p(\log \mu, \log \zeta | \mathcal{D}, \mathcal{H})$ , the cost function (34) is obtained in order to optimize  $\mu$  and  $\zeta$ , which can be further simplified into optimizing  $\gamma = \zeta/\mu$  from (38).

3. The posterior probability of the model  $\mathcal{H}$  is obtained on the third level as

$$p(\mathcal{H} | \mathcal{D}) \propto p(\mathcal{D} | \mathcal{H})p(\mathcal{H})$$

As there are infinitely many models, the evidence is omitted here. The candidate models  $\mathcal{H}_i$  (with different kernel parameters  $\sigma_i$  or different sets of explanatory inputs  $I_i$ ) are compared using the expression (44).

Observe that the evidence on level 1 and 2 is equal to likelihoods on level 2 and 3, respectively, which implies that one also needs expressions of the lower levels in order to perform inference on the higher levels. The mathematical details of the Bayesian inference are given in Appendix B.

The LS-SVM is designed in the Bayesian framework using the following steps:

1. Preprocess the data by completing missing values and handling outliers. Standardize the inputs to zero mean and unit variance.
2. Define models  $\mathcal{H}_i$  by choosing a candidate input set  $I_i$ , a kernel function  $K_i$  and a kernel parameter, e.g.,  $\sigma_i$  in the RBF kernel case. For all models  $\mathcal{H}_i$ , with  $i = 1, \dots, n_{\mathcal{H}}$  (with  $n_{\mathcal{H}}$  the number of models to be compared), compute the level 3 posterior:
  - (a) Find the optimal hyperparameters  $\mu_{MP}$  and  $\zeta_{MP}$  by solving the scalar optimization problem (38) in  $\gamma = \zeta\mu$  related to maximizing the level 2 posterior.<sup>3</sup> With the resulting  $\gamma_{MP}$ , compute the effective number of parameters from (35), the hyperparameters  $\mu_{MP}$  and  $\zeta_{MP}$ .
  - (b) Evaluate the level 3 posterior (44) for model comparison.
3. Select the model  $\mathcal{H}_i$  with maximal evidence. If desired, refine the model tuning parameters  $K_i$ ,  $\sigma_i$ ,  $I_i$  to further optimize the classifier and go back to step 2; else, go to step 4.
4. Given the optimal  $\mathcal{H}_i^*$ , calculate  $\alpha$  and  $b$  from (16), with kernel  $K_i$ , parameter  $\sigma_i$  and input set  $I_i$ .

Given the optimized model and its parameters, the prediction for a new observation is obtained by first standardizing the inputs in exactly the same way as the training set and then evaluating (17) for the prediction  $\hat{y}$  and (48) for the variance  $\sigma_y^2$  indicating the uncertainty on the prediction due to the noise and model uncertainty.

## CASE STUDY: STOCK MARKET PREDICTION

In this application, the goal is to predict the performance of an aggregated equity price index for the European chemical sector. The data set consists of 13 variables in weekly values ranging from April 1986 to February 2001, provided by Datastream (787 observations). The dependent variable is labelled CHMCLEM, and the prediction will be made in a one-week-ahead schedule. The set of (candidate) explanatory variables selected by the financial analyst includes macroeconomic indices (industrial production index, gross domestic product, consumer price index), some specific market series (oil price, raw materials price) and some financial variables (bonds, exchange rate dollar/euro, fibor 3-month interest rate). The first 600 observations are selected for initial model estimation (from April 1986 to mid-September 1997). Details on the behaviour of the data (in logarithms) within the estimation sample are described in Table I.

### Performance measures

The model is evaluated in a forward way on the time period  $t = 601, \dots, 787$ . The out-of-sample forecasts for  $\text{CHMCLEM}_t^{(d)}$  are computed as follows. The first forecast ( $t = 601$ ) is computed from the initial model. Then, the first observation is dismissed ( $t = 1$ ) and the new observation is incorporated in the model for re-estimation ( $t = 601$ ). In this way, each forecast is computed from a model estimated with the last 600 observations available. Only the variables found to be relevant are used in this moving window approach, it is assumed that the relevance found in the initial 600 data points will also hold out-of-sample.

<sup>3</sup> Observe that this implies in each iteration step maximizing the level 1 posterior in  $w$  and  $b$ .



Table I. Name and description of the dependent or output variable and of the candidate explanatory or input variables as selected by the financial analyst. The mean, maximum, minimum and standard deviation of the data in the training sample are also reported

Variable	Description	Min	Max	Mean	St.Dev.
CHMCLEM	Index/Chemical Sector	5.6503	6.8340	6.1555	0.2668
CPI	Consumer Price Index	4.2885	4.6201	4.4607	0.1080
FIBOR3M	Interest Rate	1.1309	2.2935	1.7294	0.3750
EUOOCIPDG	Industrial Production Index	4.4224	4.6482	4.5423	0.0586
OILBREN	Oil Price	2.1955	3.7098	2.8871	0.1895
ETYEUSP	Ethylene Price	5.6922	6.7557	6.0901	0.2562
PHREUSN	Specific Polymer Price	-1.0564	-0.0030	-0.3875	0.2582
GDP	European GDP	6.9418	7.2249	7.1061	0.0841
USEURWD	Exchange Rate US/Euro	-0.0739	0.4240	0.2123	0.0857
CHLORED	Chlorine Price	1.7373	2.2330	2.0123	0.1231
BMBD02Y	Bond 2 years (index)	4.5641	4.6774	4.6331	0.0305
BMBD05Y	Bond 5 years (index)	4.4890	4.6659	4.5862	0.0533
BMBD010Y	Bond 10 years (index)	4.4323	4.6955	4.5689	0.0654

The quality of the forecasts is assessed as follows. One usual way to quantify the quality of the forecasts is by using magnitude accuracy measures, such as the mean squared error (MSE) or the mean absolute error (MAE). Additionally, in financial applications, the percentage of correct sign predictions (PCSP) is often used, or the success in forecasting only the *direction* of the change rather than its magnitude (in plain terms, if the stock price rises or falls). The PCSP significance is assessed by using the Pesaran–Timmerman test statistic (PTstat) and the corresponding *p*-value (Pesaran and Timmerman, 1992). This test discriminates if the PCSP is obtained randomly or not. A PTstat above 2 allows us to reject the null hypothesis of no dependency between the predictions and the observations.

Nevertheless, a simple trading exercise is performed using a transaction cost of 0.1% (10bps as in Refenes and Zapranis, 1999) to assess the market timing ability. In investment strategy 1 (IS1), a naive allocation of 100% equities or cash is implemented, based on the sign of the prediction. The corresponding Sharpe ratio (SR, defined as the ratio between the return and the risk of a particular asset), equivalent yearly return (Re) and risk (Ri) on the test set are computed. A more advanced trading rule involves the use of the uncertainty or moderated output on the prediction for doing the actual trade: trading is still based on the sign of the prediction, but only when the ratio  $\hat{y}/\sigma_y$  exceeds a threshold value (investment strategy 2). For comparison, the same indicators for a simple buy&hold strategy (buy the asset today and sell it at the end period) are calculated.

### Stationarity and cointegration analysis

All series were found to be nonstationary by using the ADF test. Nevertheless, evidence of linear cointegration between the variables CHMCLEM<sup>(l)</sup>, FIBOR3M<sup>(l)</sup>, CPI<sup>(l)</sup> and IP<sup>(l)</sup> was found. The residuals  $e_t$  in the linear regression

$$\text{CHMCLEM}_t^{(l)} = a_0 + a_1 \text{CPI}_t^{(l)} + a_2 \text{FIBOR3M}_t^{(l)} + a_3 \text{IP}_t^{(l)} + e_t \quad (18)$$

are found to be stationary, as reported in Table II. It can be seen from the critical values of the ADF test that the null hypothesis of nonstationarity of the residuals can be rejected.

Table II. Estimates of the cointegrating regression, where the variables are used in nonstationary levels. The high  $R^2$  value and the low Durbin–Watson (DW) statistics are clearly a sign of a misleading regression in terms of inference and forecasts. Usually the DW statistic, measuring the serial correlation between the residuals of a regression, is between 0 and 4 with a value of 2.00 showing no serial correlation. The ADF statistic is equal to  $-4.85$ , which is low given the critical values  $-4.64$  (1%),  $-4.10$  (5%) and  $-3.81$  (10%)

Cointegrating regression: $R^2 = 0.95$ , DW = 0.08, ADF = $-4.85$			
Variable	Coeff.	St.Dev.	$t$ -stat
Constant	$-8.4264$	0.2378	$-35.4397$
CPI	0.7897	0.0617	12.8098
FIBOR3M	$-0.3422$	0.0083	$-41.4730$
EUOOCIPDG	2.5658	0.0977	26.2700

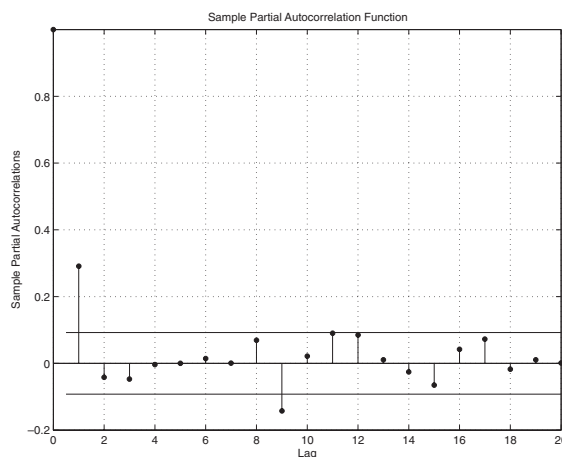


Figure 2. Partial autocorrelation function (pacf) for the output or dependent variable  $\text{CHCLEM}^{(d)}$  in the ECM specification

The evidence of cointegration between the dependent original variable and a subset of the explanatory variables allow us to implement an ECM specification. The series  $\text{CHMCLEM}^{(d)}$ ,  $\text{FIBOR3M}^{(d)}$ ,  $\text{CPI}^{(d)}$  and  $\text{IP}^{(d)}$  are used as first differences. Also, from Figure 2 we see that the partial autocorrelation function of the dependent variable in the ECM ( $\text{CHMCLEM}_t^{(d)}$ ) drops to zero after one lag. Therefore, we will include all variables lagged at  $t - 1$  and  $t - 2$  (one additional lag is used for conservativeness). The remaining explanatory variables (those not included in the cointegration) will also be included in the model up to two lags, in first differences, as exogenous variables.

According to the ECM specification, we have the following model:

$$\begin{aligned} \text{CHMCLEM}_t^{(d)} = & f(z_{t-1}, \text{CHMCLEM}_{t-i}^{(d)}, \text{FIBOR3M}_{t-i}^{(d)}, \text{CPI}_{t-i}^{(d)}, \text{CPI}_{t-i}^{(d)}, \\ & \text{IP}_{t-i}^{(d)}, \text{OILBREN}_{t-i}^{(d)}, \text{ETYEUSP}_{t-i}^{(d)}, \text{PHREUSN}_{t-i}^{(d)}, \text{GDP}_{t-i}^{(d)}, \text{USEURWD}_{t-i}^{(d)}, \\ & \text{CHLORED}_{t-i}^{(d)}, \text{BMBD02Y}_{t-i}^{(d)}, \text{BMBD05Y}_{t-i}^{(d)}, \text{BMBD010Y}_{t-i}^{(d)}) + e_t \end{aligned} \quad (19)$$

Table III. Final estimates of the linear regression based on the ECM specification

Linear ECM regression: $R^2 = 0.10$ , DW = 2.01			
Variable	Coeff.	St.Dev.	$t$ -stat
$z_{t-1}$	-0.0397	0.0100	-3.9718
CHMCLEM $_{t-1}^{(d)}$	0.2213	0.0391	5.6589
FIBOR3M $_{t-1}^{(d)}$	0.1555	0.0513	3.0318
FIBOR3M $_{t-2}^{(d)}$	-0.2216	0.0514	-4.3144
ETYEUSP $_{t-2}^{(d)}$	-0.0659	0.0327	-2.0141

where  $i = 1, 2$  and  $z_{t-1} = \text{CHMCLEM}_{t-1}^{(d)} - (\hat{a}_0 + \hat{a}_1 \text{CPI}_{t-1}^{(d)} + \hat{a}_2 \text{FIBOR3M}_{t-1}^{(d)} + \hat{a}_3 \text{IP}_{t-1}^{(d)})$ , where the coefficients are estimated from (18).

### Linear ECM model

With this initial definition of input variables, the function  $f$  will be estimated by the linear OLS regression. The data set contains 787 observations. The model is first estimated using the initial 600 data points, in order to obtain the relevant variables. In the linear regression, the selection of relevant inputs is based on the asymptotic  $t$ -tests of individual significance. With this methodology, the relevant variables are found to be the following:  $z_{t-1}$ , CHMCLEM $_{t-1}^{(d)}$ , FIBOR3M $_{t-1}^{(d)}$ , FIBOR3M $_{t-2}^{(d)}$  and ETYEUSP $_{t-2}^{(d)}$ . The detailed results for the linear regression are reported in Table III.

The linear forecasts yield the following performance indicators: MSE =  $6.63 \times 10^{-4}$ , MAE = 0.021, PCSP = 53.2%, PTstat = 0.73,  $p$ -value = 0.462. In this case, the low PTstat for the 53% of prediction accuracy shows that it is not significantly different from a random case. The results for investment strategy 1 are Sharpe ratio = 0.596, return = 8.94 and risk = 15.00, while investment strategy 2 yields Sharpe ratio = 0.487, return = 7.19 and risk = 14.76, respectively. Compared to the buy&hold strategy (Sharpe ratio = 0.208, return = 3.98, risk = 19.14), we can see that the linear model defined in this section allows us to increase the possible profits compared to a simple buy&hold strategy.

### Nonlinear ECM model

Nonlinear modelling was performed using the ECM specification with the same candidate input set as in (19). The model is designed on the same training data set and is evaluated on the remaining test set using the same moving window approach.

Backward input selection is applied, removing in each step one input until the model probability  $p(\mathcal{H}|\mathcal{D})$  stops increasing. The evolution of the level 3 and level 2 cost function as a function of the number of input pruning steps is depicted in Figure 3, together with the evolution of  $d_{\text{eff}}$ ,  $\gamma$  and of the directional accuracy measures PTstat and PCSP. From the initial 27 input variables, 21 inputs have been removed. The optimal input set for the nonlinear model is reported in Table IV, from which it is observed that the variables found to relevant for the nonlinear models are almost the same as for the linear models. The cointegrating vector  $z_t$  is kept as an important input, its removal would lead to a significantly worse predictive performance of the nonlinear model. The optimal regularization and kernel function parameters that are inferred from the training data  $\mathcal{D}$  are:  $\mu_{MP} = 258.29$ ,  $\zeta_{MP} = 2.93 \times 10^3$ ,  $\gamma_{MP} = 11.34$  and  $\sigma_{MP} = 0.25$ .

With the nonlinear forecasts, we obtain MSE =  $6.87 \times 10^{-4}$ , MAE = 0.021, while the performance measures for directional accuracy are PCSP = 60.2%, PTstat = 2.66,  $p$ -value = 0.8%. The high PTstat

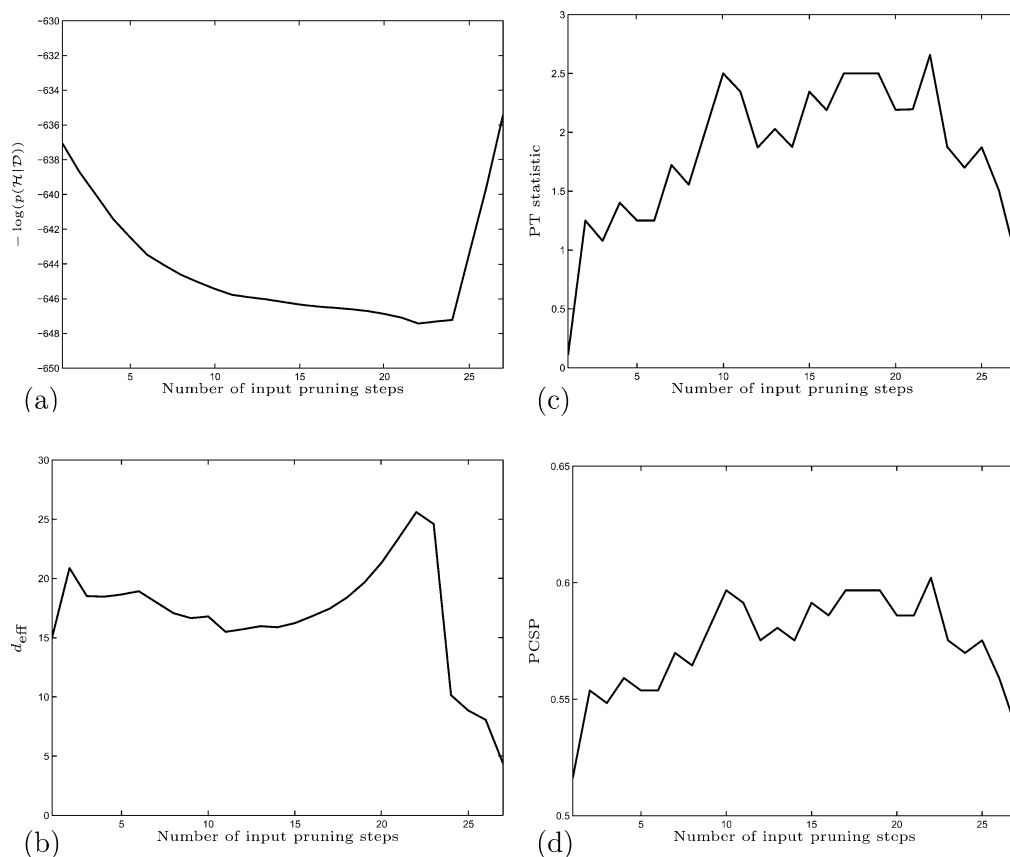


Figure 3. Backward input selection for nonlinear kernel-based regression. The evolution of the level 3 cost function (a) as a function of the number of input pruning steps yields an optimum at step 22. The corresponding values for the effective number of parameters  $d_{\text{eff}}$ , the out-of-sample Pesaran–Timmerman test statistic PTstat and the percentage of correct sign predictions (PCSP) are reported in panels (b)–(d). Notice that the PT statistic and PCSP become maximal at the minimum of the level 3 cost function

for the 60.2% of predictional accuracy shows that it is significantly different from random sign predictions. This is also observed in better performances with both trading strategies. With investment strategy 1 we obtain Sharpe ratio = 0.826, return = 12.75 and risk = 15.44, while investment strategy 2 yields Sharpe ratio = 0.841, return = 12.78 and risk = 15.21. The results for the buy&hold strategy, the linear and nonlinear ECM models are summarized in Table V, while the cumulative profits are depicted in Figure 4. Using almost the same input variables, the nonlinear model achieves clearly better out-of-sample performance.

## CONCLUSIONS

In financial time series modelling and prediction, it is important to have reliable forecasting and modelling techniques. Based on stationarity and cointegration analysis, a linear ECM model is spec-

Table IV. Optimal input sets for the linear and nonlinear models

Linear		LS-SVM	
Variables	Lags	Variables	Lags
$z$	$t - 1$	$z$	$t - 1$
CHMCLEM	$t - 1$	CHMCLEM	$t - 1$
FIBOR3M	$t - 1, t - 2$	FIBOR3M	$t - 1, t - 2$
ETYEUSP	$t - 2$	ETYEUSP	$t - 2$
		BMBD010Y	$t - 2$

Table V. Test set performances of the LS-SVM model obtained on the one-week-ahead prediction of the aggregated chemical index. The LS-SVM time series model with RBF-kernel is compared with linear ECM and a buy&hold strategy. The RBF-LS-SVM clearly achieves a better directional accuracy, better return (Re), risk (Ri) and resulting Sharpe ratio (SR) in combination with investment strategies IS1 and IS2

	Residuals		Directional accuracy			IS1			IS2		
	MSE	MAE	PCSP	PT	$p$ -value	SR <sub>1</sub>	Re <sub>1</sub>	Ri <sub>1</sub>	SR <sub>2</sub>	Re <sub>2</sub>	Ri <sub>2</sub>
LS-SVM	$6.87 \times 10^{-4}$	0.021	60.2	2.66	0.008	0.826	12.75	15.44	0.841	12.78	15.21
Linear	$6.63 \times 10^{-4}$	0.021	53.2	0.73	0.462	0.596	8.94	15.00	0.487	7.19	14.76
B&H	—	—	—	—	—	0.208	3.98	19.14	0.208	3.98	19.14

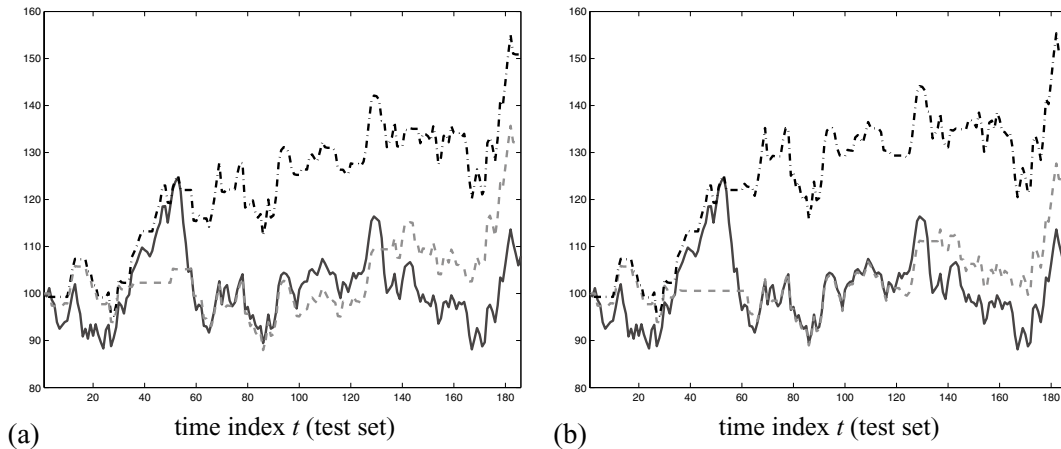


Figure 4. Cumulative returns using the sign predictions (transaction cost 0.1%) on the out-of-sample test set obtained with: (1) LS-SVM regressor with nonlinear RBF-kernel (dash-dotted line); (2) linear model (dashed line); and (3) buy&hold strategy (full line). The LS-SVM regressor yields the highest annualized return and corresponding Sharpe ratio as reported in Table V. Panels (a) and (b) depict the results of investment strategies 1 and 2, respectively

ified and estimated in order to produce out-of-sample linear stock market forecasts. The specified input variables of the linear ECM formulation are used in this paper as an initial candidate input set for nonlinear kernel-based regression. The nonlinear model was designed within the Bayesian evidence framework, getting an appropriate trade-off between model complexity and in-sample model accuracy. The regularization, kernel function parameters and relevant inputs are obtained by applying Bayes' formula on different levels of inference. For the prediction of an aggregated index for the European chemical sector, it was found that the optimal linear and nonlinear forecasts are based on almost the same set of relevant variables including the cointegrating vector. Comparing both techniques, the nonlinear model achieves significantly nonrandom out-of-sample sign predictions and also yields better a Sharpe ratio when implemented in a simple trading strategy.

## APPENDIX A: BAYESIAN INFERENCE FOR LS-SVM REGRESSION

### Inference of model parameters (level 1)

#### Bayes' formula

Applying Bayes' formula on level 1, one obtains the posterior probability of the model parameters  $\mathbf{w}$  and  $b$ :

$$\begin{aligned} p(\mathbf{w}, b | \mathcal{D}, \log \mu, \log \zeta, \mathcal{H}) &= \frac{p(\mathcal{D} | \mathbf{w}, b, \log \mu, \log \zeta, \mathcal{H}) p(\mathbf{w}, b | \log \mu, \log \zeta, \mathcal{H})}{p(\mathcal{D} | \log \mu, \log \zeta, \mathcal{H})} \\ &\propto p(\mathcal{D} | \mathbf{w}, b, \log \mu, \log \zeta, \mathcal{H}) p(\mathbf{w}, b | \log \mu, \log \zeta, \mathcal{H}) \end{aligned} \quad (20)$$

where the last step is obtained since the evidence  $p(\mathcal{D} | \log \mu, \log \zeta, \mathcal{H})$  is a normalizing constant that does not depend upon  $\mathbf{w}$  and  $b$ .

For the *prior*, no correlation between  $\mathbf{w}$  and  $b$  is assumed:  $p(\mathbf{w}, b | \log \mu, \mathcal{H}) = p(\mathbf{w} | \log \mu, \mathcal{H}) p(b | \mathcal{H}) \propto p(\mathbf{w} | \log \mu, \mathcal{H})$ , with a multivariate Gaussian prior on  $\mathbf{w}$  with zero mean and covariance matrix  $\mu^{-1} \mathbf{I}_{n_\phi}$  and an uninformative, flat prior on  $b$ :

$$\begin{aligned} p(\mathbf{w} | \log \mu, \mathcal{H}) &= \left( \frac{\mu}{2\pi} \right)^{\frac{n_\phi}{2}} \exp\left( -\frac{\mu}{2} \mathbf{w}^T \mathbf{w} \right) \\ p(b | \mathcal{H}) &= \text{constant} \end{aligned} \quad (21)$$

The uniform prior distribution on  $b$  can be approximated by a Gaussian distribution with standard deviation  $\sigma_b \rightarrow \infty$ . The negative logarithm of (21) corresponds to the regularization term  $\mu \mathbf{J}_w = \mu/2 \mathbf{w}^T \mathbf{w}$ . The prior states a belief that without any learning from data, the coefficients are zero with an uncertainty denoted by the variance  $1/\mu$ ; *a priori* we do not expect a functional relation between the feature vector  $\phi$  and the observation  $y$ . Before the data are available, the most likely model has zero weights  $w_k = 0$  ( $k = 1, \dots, n_\phi$ ), corresponding to the efficient market hypothesis (Bachelier, 1900; Campbell *et al.*, 1997; Fama, 1965).

It is assumed that the errors  $e_t = y_t - (\mathbf{w}^T \phi(\mathbf{x}_t) + b)$  are independently identically normally distributed with zero mean and variance  $1/\zeta$  for expressing the *likelihood*

$$p(\mathcal{D}|\mathbf{w}, b, \log \zeta, \mathcal{H}) \propto \prod_{i=1}^{n_D} p(e_i|\mathbf{x}_i, \mathbf{w}, b, \log \zeta, \mathcal{H}) \tag{22}$$

with

$$p(e_i|\mathbf{w}, b, \log \zeta, \mathcal{H}) = \sqrt{\frac{\zeta}{2\pi}} \exp\left(-\frac{\zeta}{2}(y_i - \mathbf{w}^T \varphi(\mathbf{x}_i) - b)^2\right) \tag{23}$$

The negative logarithm of the likelihood (22) corresponds to the sum squared error term  $\zeta J_e = \frac{\zeta}{2} \sum_{i=1}^{n_D} e_i^2$ .

Substituting (21) and (22) into (20), neglecting all constants and taking the negative logarithm, Bayes' rule at the first level of inference corresponds to the constrained minimization problem (9) and (10) that can be solved for  $\mathbf{w}$  and  $b$  in the primal space from (15) in the linear case when  $n \leq n_D$  or  $\boldsymbol{\alpha}$  and  $b$  in the dual space from (16) in the nonlinear kernel-based regression case and in the linear case when  $n \geq n_D$ . In the remainder of this paper, the maximum *a posteriori* parameter estimates are denoted by the subscript 'MP', e.g.,  $\mathbf{w}_{MP}$  and  $b_{MP}$ .

Given that the prior (21) and likelihood (22) are multivariate distributions, the *posterior* (20) is a multivariate normal distribution<sup>4</sup> in  $[\mathbf{w}; b]$  with mean  $[\mathbf{w}_{MP}; b_{MP}] \in \mathbb{R}^{n_{\varphi}+1}$  and covariance matrix  $\mathbf{Q} \in \mathbb{R}^{(n_{\varphi}+1) \times (n_{\varphi}+1)}$ . An alternative expression for the posterior is obtained by substituting (21) and (22) into (20). These approaches yield

$$p(\mathbf{w}, b|\mathcal{D}, \log \mu, \log \zeta, \mathcal{H}) = \sqrt{\frac{\det(\mathbf{Q}^{-1})}{(2\pi)^{n_{\varphi}+1}}} \exp\left(-\frac{1}{2}[\mathbf{w} - \mathbf{w}_{MP}; b - b_{MP}]\mathbf{Q}^{-1}[\mathbf{w} - \mathbf{w}_{MP}; b - b_{MP}]\right) \tag{24}$$

$$\propto \left(\frac{\mu}{2\pi}\right)^{\frac{n_f}{2}} \exp\left(-\frac{\mu}{2}\mathbf{w}^T\mathbf{w}\right) \left(\frac{\zeta}{2\pi}\right)^{\frac{n_D}{2}} \exp\left(-\frac{\zeta}{2}\sum_{i=1}^{n_D} e_i^2\right) \tag{25}$$

respectively.

The *evidence* is a normalizing constant in (20) independent of  $\mathbf{w}$  and  $b$  such that  $\int \dots \int p(\mathbf{w}, b|\mathcal{D}, \log \mu, \log \zeta, \mathcal{H}) d\mathbf{w}_1 \dots d\mathbf{w}_{n_{\varphi}} db = 1$ . Substituting the expressions for the prior (21), likelihood (22) and posterior (25) into (20), one obtains

$$p(\mathcal{D}|\log \mu, \log \zeta, \mathcal{H}) = \frac{p(\mathbf{w}_{MP}|\log \mu, \mathcal{H})p(\mathcal{D}|\mathbf{w}_{MP}, b_{MP}, \log \zeta, \mathcal{H})}{p(\mathbf{w}_{MP}, b_{MP}|\mathcal{D}, \log \mu, \log \zeta, \mathcal{H})} \tag{26}$$

*Computation and interpretation*

The model parameters with maximum posterior probability are obtained by minimizing the negative logarithm of (24) and (25):

<sup>4</sup>The notation  $[x; y] = [x, y]^T$  is used here.

$$\begin{aligned}
(\mathbf{w}_{MP}, b_{MP}) &= \arg \min_{\mathbf{w}, b} J_1(\mathbf{w}, b) \\
&= J_1(\mathbf{w}_{MP}, b_{MP}) + \frac{1}{2}([\mathbf{w} - \mathbf{w}_{MP}; b - b_{MP}]^T \mathbf{Q}^{-1}[\mathbf{w} - \mathbf{w}_{MP}; b - b_{MP}])
\end{aligned} \tag{27}$$

$$= \frac{\mu}{2} \mathbf{w}^T \mathbf{w} + \frac{\zeta}{2} \sum_{i=1}^{n_D} e_i^2 \tag{28}$$

where constants are neglected in the optimization problem. Both expressions yield the same optimization problem and the covariance matrix  $\mathbf{Q}$  is equal to the inverse of the Hessian  $\mathbf{H}$  of  $J_1$ . The Hessian is expressed in terms of the matrix  $\Phi = [\varphi(\mathbf{x}_1), \dots, \varphi(\mathbf{x}_{n_D})]^T$  with regressors, as derived in Appendix B.

The optimal  $\mathbf{w}_{MP}$  and  $b_{MP}$  are computed in the dual space from the linear KKT-system (16), while the prediction  $\hat{y} = \mathbf{w}_{MP}^T \varphi(\mathbf{x}) + b_{MP}$  is expressed in terms of the dual parameters  $\alpha$  and bias term  $b_{MP}$  via (17).

Substituting (21), (22) and (25) into (26), one obtains

$$p(\mathcal{D} | \log \mu, \log \zeta, \mathcal{H}) \propto \left( \frac{\mu^{n_\phi} \zeta^{n_D}}{\det \mathbf{H}} \right)^{\frac{1}{2}} \exp(-J_1(\mathbf{w}_{MP}, b_{MP})) \tag{29}$$

As  $J_1(\mathbf{w}, b) = \mu J_w(\mathbf{w}) + \zeta J_e(\mathbf{w}, b)$ , the evidence can be rewritten as

$$\underbrace{p(\mathcal{D} | \log \mu, \log \zeta, \mathcal{H})}_{\text{evidence}} \propto \underbrace{p(\mathcal{D} | \mathbf{w}_{MP}, b_{MP}, \log \zeta, \mathcal{H})}_{\text{likelihood} | \mathbf{w}_{MP}, b_{MP}} \underbrace{p(\mathbf{w}_{MP} | \log \mu, \mathcal{H})}_{\text{Occam factor}} (\det \mathbf{H})^{-1/2}$$

The model evidence consists of the likelihood of the data and an Occam factor that penalizes for too complex models. The Occam factor consists of the regularization term  $1/2 \mathbf{w}_{MP}^T \mathbf{w}_{MP}$  and the ratio  $(\mu^{n_\phi} / \det \mathbf{H})^{1/2}$ , which is a measure for the volume of the posterior probability divided by the volume of the prior probability. Strong contractions of the posterior versus prior space indicates too many free parameters and, hence, overfitting on the training data. The evidence will be maximized on level 2, where also dual space expressions are derived.

## Inference of hyperparameters (level 2)

### Bayes' formula

The optimal regularization parameters  $\mu$  and  $\zeta$  are inferred from the given data  $\mathcal{D}$  by applying Bayes' formula on the second level (Van Gestel *et al.*, 2001, 2002)

$$p(\log \mu, \log \zeta | \mathcal{D}, \mathcal{H}) = \frac{p(\mathcal{D} | \log \mu, \log \zeta, \mathcal{H}) p(\log \mu, \log \zeta)}{p(\mathcal{D} | \mathcal{H})} \tag{30}$$

The prior  $p(\log \mu, \log \zeta) = p(\log \mu | \mathcal{H}) p(\log \zeta | \mathcal{H}) = \text{constant}$  is taken to be a flat uninformative prior ( $\sigma_{\log \mu}, \sigma_{\log \zeta} \rightarrow \infty$ ). The level 2 likelihood  $p(\mathcal{D} | \log \mu, \log \zeta, \mathcal{H})$  is equal to the level 1 evidence (29). In this way, Bayesian inference implicitly embodies Occam's razor: on level 2 the evidence of level 1 is optimized so as to find a trade-off between the model fit and a complexity term to avoid overfitting (MacKay, 1995). The level 2 evidence is obtained in a similar way as on level 1, as the like-



likelihood for the maximum *a posteriori* times the ratio of the volume of the posterior probability and the volume of the prior probability:

$$p(\mathcal{D}|\mathcal{H}) \simeq p(\mathcal{D}|\log \mu_{MP}, \log \zeta_{MP}, \mathcal{H}) \frac{\sigma_{\log \mu|\mathcal{D}} \sigma_{\log \zeta|\mathcal{D}}}{\sigma_{\log \mu} \sigma_{\log \zeta}} \quad (31)$$

where one typically approximates the posterior probability by a multivariate normal probability function with diagonal covariance matrix  $\text{diag}([\sigma_{\log \mu|\mathcal{D}}^2, \sigma_{\log \zeta|\mathcal{D}}^2]) \in \mathbb{R}^{2 \times 2}$ .

Neglecting all constants, Bayes' formula (30) becomes

$$p(\log \mu, \log \zeta|\mathcal{D}, \mathcal{H}) \propto p(\mathcal{D}|\log \mu, \log \zeta, \mathcal{H}) \quad (32)$$

where the expressions for the level 1 evidence are given by (26) and (29).

#### Computation and interpretation

In the primal space, the hyperparameters are obtained by minimizing the negative logarithm of (29) and (32)

$$\begin{aligned} (\mu_{MP}, \zeta_{MP}) = \arg \min_{\mu, \zeta} J_2(\mu, \zeta) &= \mu J_w(\mathbf{w}_{MP}) + \zeta J_e(\mathbf{w}_{MP}, b_{MP}) \\ &+ \frac{1}{2} \log \det \mathbf{H} - \frac{n_\varphi}{2} \log \mu - \frac{n_{\mathcal{D}}}{2} \log \zeta \end{aligned} \quad (33)$$

Observe that in order to evaluate (33) one needs also to calculate  $\mathbf{w}_{MP}$  and  $b_{MP}$  for the given  $\mu$  and  $\zeta$  and evaluate the level 1 cost function. The determinant of  $\mathbf{H}$  is equal to (see Appendix B for details)

$$\det(\mathbf{H}) = (\zeta^{n_{\mathcal{D}}}) \det(\mu \mathbf{I}_{n_\varphi} + \zeta \Phi^T \mathbf{N}_c \Phi)$$

with the idempotent centring matrix  $\mathbf{N}_c = \mathbf{I}_{n_{\mathcal{D}}} - 1/n_{\mathcal{D}} \mathbf{1}\mathbf{1}^T = \mathbf{N}_c^2 \in \mathbb{R}^{n_{\mathcal{D}} \times n_{\mathcal{D}}}$ . The determinant is also equal to the product of the eigenvalues. The  $n_e$  nonzero eigenvalues  $\lambda_1, \dots, \lambda_{n_e}$  of  $\Phi^T \mathbf{N}_c \Phi$  are equal to the  $n_e$  nonzero eigenvalues of  $\mathbf{N}_c \Phi \Phi^T \mathbf{N}_c = \mathbf{N}_c \mathbf{\Omega} \mathbf{N}_c \in \mathbb{R}^{n_{\mathcal{D}} \times n_{\mathcal{D}}}$ , which can be calculated in the dual space. Substituting the determinant  $\det(\mathbf{H}) = \zeta^{n_{\mathcal{D}}} \mu^{n_\varphi - n_e} \prod_{i=1}^{n_e} (\mu + \zeta \lambda_i)$  into (33), one obtains the optimization problem in the dual space

$$J_2(\mu, \zeta) = \mu J_w(\mathbf{w}_{MP}) + \zeta J_e(\mathbf{w}_{MP}, b_{MP}) + \sum_{i=1}^{n_e} \frac{\log(\mu + \zeta \lambda_i)}{2} - \frac{n_e}{2} \log \mu - \frac{n_e - 1}{2} \log \zeta \quad (34)$$

where it can be shown by matrix algebra (see Appendix B) that  $\mu J_w(\mathbf{w}_{MP}) + \zeta J_e(\mathbf{w}_{MP}, b_{MP}) = \frac{1}{2} \mathbf{y}^T \mathbf{N}_c \left( \frac{1}{\mu} \mathbf{N}_c \mathbf{\Omega} \mathbf{N}_c + \frac{1}{\zeta} \mathbf{I}_{n_{\mathcal{D}}} \right)^{-1} \mathbf{N}_c \mathbf{y}$ .

An important concept in neural networks and Bayesian learning in general is the *effective number of parameters*. Although there are  $n_\varphi + 1$  free parameters  $w_1, \dots, w_{n_\varphi}, b$  in the primal space, the use of these parameters (28) is restricted by the use of the regularization term  $1/2 \mathbf{w}^T \mathbf{w}$ . The effective number of parameters  $d_{\text{eff}}$  is equal to  $d_{\text{eff}} = \sum_i \lambda_{i,u} / \lambda_{i,r}$ , where  $\lambda_{i,u}$ ,  $\lambda_{i,r}$  denote the eigenvalues of the Hessian of the unregularized cost function  $J_{1,u} = \zeta E_{\mathcal{D}}$  and the regularized cost function  $J_{1,r} = \mu E_w$

+  $\zeta E_D$  (Bishop, 1995; MacKay, 1995). For LS-SVMs, the effective number of parameters is equal to

$$d_{\text{eff}} = 1 + \sum_{i=1}^{n_e} \frac{\zeta \lambda_i}{\mu + \zeta \lambda_i} = 1 + \sum_{i=1}^{n_e} \frac{\gamma \lambda_i}{1 + \gamma \lambda_i} \quad (35)$$

with  $\gamma = \zeta/\mu \in \mathbb{R}^+$ . The term +1 appears because no regularization is applied on the bias term  $b$ . As shown, one has that  $n_e \leq n_D - 1$  and, hence, also that  $d_{\text{eff}} \leq n_D$ , even in the case of high-dimensional feature spaces.

The conditions for optimality for (34) are obtained by putting  $\partial J_2/\partial \mu = \partial J_2/\partial \zeta = 0$ . One obtains<sup>5</sup>

$$\partial J_2/\partial \mu = 0 \rightarrow 2\mu_{MP} J_w(\mathbf{w}_{MP}; \mu_{MP}, \zeta_{MP}) = d_{\text{eff}}(\mu_{MP}, \zeta_{MP}) - 1 \quad (36)$$

$$\partial J_2/\partial \zeta = 0 \rightarrow 2\zeta_{MP} J_e(\mathbf{w}_{MP}; b_{MP}; \mu_{MP}, \zeta_{MP}) = n_D - d_{\text{eff}} \quad (37)$$

where the latter equation corresponds to the unbiased estimate of the noise variance  $1/\zeta_{MP} = \frac{1}{2} \sum_{i=1}^{n_D} e_i^2/(n_D - d_{\text{eff}})$ .

Instead of solving the optimization problem in  $\mu$  and  $\zeta$ , one may also reformulate (34) using (36) and (37) in terms of  $\gamma = \zeta/\mu$  and solve the following scalar optimization problem (Van Gestel *et al.*, 2002):

$$\min_{\gamma} \sum_{i=1}^{n_D-1} \log\left(\lambda_i + \frac{1}{\gamma}\right) + (n_D - 1) \log(J_w(\mathbf{w}_{MP}) + \gamma J_e(\mathbf{w}_{MP}, b_{MP})) \quad (38)$$

with

$$J_e(\mathbf{w}_{MP}, b_{MP}) = \frac{1}{2\gamma^2} \mathbf{y}^T \mathbf{N}_c \mathbf{V} (\Lambda + \mathbf{I}_N/\gamma)^{-2} \mathbf{V}^T \mathbf{N}_c \mathbf{y} \quad (39)$$

$$J_w(\mathbf{w}_{MP}) = \frac{1}{2} \mathbf{y}^T \mathbf{N}_c \mathbf{V} \Lambda (\Lambda + \mathbf{I}/\gamma)^{-2} \mathbf{V}^T \mathbf{N}_c \mathbf{y} \quad (40)$$

$$J_w(\mathbf{w}_{MP}) + \gamma J_e(\mathbf{w}_{MP}, b_{MP}) = \frac{1}{2} \mathbf{y}^T \mathbf{N}_c \mathbf{V} (\Lambda + \mathbf{I}_N/\gamma)^{-1} \mathbf{V}^T \mathbf{N}_c \mathbf{y} \quad (41)$$

and with the eigenvalue decomposition  $\mathbf{N}_c \mathbf{\Omega} \mathbf{N}_c = \mathbf{V}^T \mathbf{\Lambda} \mathbf{V}$ . Given the optimal  $\gamma_{MP}$  from (38) one finds the effective number of parameters  $d_{\text{eff}}$  from  $d_{\text{eff}} = 1 + \sum_{i=1}^{n_e} \gamma \lambda_i/(1 + \gamma \lambda_i)$ . The optimal  $\mu_{MP}$  and  $\zeta_{MP}$  are obtained from  $\mu_{MP} = (d_{\text{eff}} - 1)/(2J_w(\mathbf{w}_{MP}))$  and  $\zeta_{MP} = (n_D - d_{\text{eff}})/(2J_e(\mathbf{w}_{MP}, b_{MP}))$ .

<sup>5</sup> In this derivation, one uses that (MacKay, 1995; Suykens *et al.*, 2002; Van Gestel *et al.*, 2002)  $\partial(J_1(\mathbf{w}_{MP}, b_{MP}))/\partial \mu = \delta(J_1(\mathbf{w}_{MP}, b_{MP}))/\delta \mu + \delta(J_1(\mathbf{w}_{MP}, b_{MP}))/\delta[\mathbf{w}; b]_{[\mathbf{w}_{MP}; b_{MP}]} \times \delta([\mathbf{w}_{MP}; b_{MP}])/ \delta \mu = J_w(\mathbf{w}_{MP})$ , since  $\delta J_1(\mathbf{w}_{MP}, b_{MP})/\delta[\mathbf{w}; b]_{[\mathbf{w}_{MP}; b_{MP}]} = 0$ .

**Model comparison (level 3)**

*Bayes' formula*

The model structure  $\mathcal{H}$  determines the remaining parameters of the kernel-based model: the selected kernel function (linear, RBF, . . .), the kernel parameter (RBF kernel parameter  $\sigma$ ) and selected explanatory inputs. The model structure is inferred on level 3.

Consider, for example, the inference of the RBF-kernel parameter  $\sigma$ , where the model structure is denoted by  $\mathcal{H}_\sigma$ . Bayes' formula for the inference of  $\mathcal{H}_\sigma$  is equal to

$$p(\mathcal{H}_\sigma|\mathcal{D}) \propto p(\mathcal{D}|\mathcal{H}_\sigma)p(\mathcal{H}_\sigma) \tag{42}$$

where no evidence  $p(\mathcal{D})$  is used in the expression on level 3 as it is in practice impossible to integrate over all model structures. The prior probability  $p(\mathcal{H}_\sigma)$  is assumed to be constant. The likelihood is equal to the level 2 evidence (31).

*Computation and interpretation*

Substituting the evidence (31) into (42) and taking in the constant prior, the Bayes' rule (31) becomes

$$p(\mathcal{H}|\mathcal{D}) \simeq p(\mathcal{D}|\log \mu_{MP}, \log \zeta_{MP}, \mathcal{H}) \frac{\sigma_{\log \mu|\mathcal{D}} \sigma_{\log \zeta|\mathcal{D}}}{\sigma_{\log \mu} \sigma_{\log \zeta}} \tag{43}$$

As uninformative priors are used on level 2, the standard deviations  $\sigma_{\log \mu}$  and  $\sigma_{\log \zeta}$  of the prior distribution both tend to infinity and are omitted in the comparisons of different models in (43). The posterior error bars can be approximated analytically as  $\sigma_{\log \mu|\mathcal{D}}^2 \simeq 2/(d_{\text{eff}} - 1)$  and  $\sigma_{\log \zeta|\mathcal{D}}^2 \simeq 2/(n_{\mathcal{D}} - d_{\text{eff}})$ , respectively (MacKay, 1995). The level 3 posterior becomes

$$p(\mathcal{H}_\sigma|\mathcal{D}) \simeq p(\mathcal{D}|\log \mu_{MP}, \log \zeta_{MP}, \mathcal{H}_\sigma) \frac{\sigma_{\log \mu|\mathcal{D}} \sigma_{\log \zeta|\mathcal{D}}}{\sigma_{\log \mu} \sigma_{\log \zeta}} \tag{44}$$

$$\propto \sqrt{\frac{\mu_{MP}^{n_e} \zeta_{MP}^{n_{\mathcal{D}}-1}}{(d_{\text{eff}} - 1)(N - d_{\text{eff}}) \prod_{i=1}^{n_e} (\mu_{MP} + \zeta_{MP} \lambda_i)}}$$

where all expressions can be calculated in the dual space. A practical way to infer the kernel parameter  $\sigma$  is to calculate (44) for a grid of possible kernel parameters  $\sigma_1, \dots, \sigma_m$  and to compare the corresponding posterior model parameters  $p(\mathcal{H}_{\sigma_1}|\mathcal{D}), \dots, p(\mathcal{H}_{\sigma_m}|\mathcal{D})$ .

Model comparison is also used to infer the set of most relevant inputs (Van Gestel *et al.*, 2001) out of the given set of candidate explanatory variables by making pairwise comparisons of models with different input sets. In a backward input selection procedure, one starts from the full candidate input set and removes in each input pruning step that input that yields the best model improvement (or smallest decrease) in terms of the model probability (44). The procedure is stopped when no significant decrease of the model probability is observed. In the case of equal prior model probabilities  $p(\mathcal{H}_i) = p(\mathcal{H}_j) (\forall i, j)$  the models  $\mathcal{H}_i$  and  $\mathcal{H}_j$  are compared according to their Bayes factor

$$\mathcal{B}_{ij} = \frac{p(\mathcal{D}|\mathcal{H}_i)}{p(\mathcal{D}|\mathcal{H}_j)} = \frac{p(\mathcal{D}|\log \mu_i, \log \zeta_i, \mathcal{H}_i)}{p(\mathcal{D}|\log \mu_j, \log \zeta_j, \mathcal{H}_j)} \frac{\sigma_{\log \mu_i|\mathcal{D}} \sigma_{\log \zeta_i|\mathcal{D}}}{\sigma_{\log \mu_j|\mathcal{D}} \sigma_{\log \zeta_j|\mathcal{D}}} \tag{45}$$

According to Jeffreys (1961), a value  $2\ln \mathcal{B}_{ij}$  corresponding to 0–2, 2–5, 5–10 and >10 indicates a very weak, positive, strong and decisive evidence against the null hypothesis of no difference in model performance between the models  $\mathcal{H}_i$  and  $\mathcal{H}_j$ .

**Moderated output**

The uncertainty on the estimated model parameters results in an additional uncertainty for the one-step-ahead prediction  $y_{MP} = \mathbf{w}_{MP}^T \boldsymbol{\varphi}(\mathbf{x}) + b_{MP} = \sum_{t=1}^{n_D} \frac{\alpha_t}{\mu} K(\mathbf{x}, \mathbf{x}_t) + b_{MP}$ . Given the normal distribution (27) of the model parameters  $[\mathbf{w}; b]$  with mean  $[\mathbf{w}_{MP}; b_{MP}]$  and covariance matrix  $\mathbf{Q}$ , and the additive noise with mean zero and noise variance  $\zeta$ , it is well known that the predicted output is normally distributed with mean

$$\hat{y}_{MP} = \mathbf{w}_{MP}^T \boldsymbol{\varphi}(\mathbf{x}) + b_{MP} \tag{46}$$

and variance

$$\sigma_y^2 = \frac{1}{\zeta} + [\boldsymbol{\varphi}(\mathbf{x}); 1]^T \mathbf{Q} [\boldsymbol{\varphi}(\mathbf{x}); 1] \tag{47}$$

where the first term is due to the additive noise  $e_t$  and the second term is due to the posterior uncertainty (27) on the model parameters  $\mathbf{w}$  and  $b$ .

The dual space expression for  $y_{MP}$  is given in (17). The expression for the variance  $\sigma_y^2$  involves the inversion of the Hessian  $\mathbf{H} = \mathbf{Q}^{-1}$  in the feature space. Given the expressions (51) and (53), the following practical expression is obtained in the dual space by applying linear matrix algebra:

$$\begin{aligned} \sigma_y^2 = & \frac{1}{\zeta} + \frac{1}{n_D \zeta} + \frac{1}{\mu} \left( K(\mathbf{x}, \mathbf{x}) + \frac{1}{n_D^2} \mathbf{1}^T \boldsymbol{\Omega} \mathbf{1} - \frac{2}{n_D} \mathbf{k}(\mathbf{x})^T \mathbf{1} \right) \\ & - \frac{\zeta}{\mu} \left( \mathbf{k}(\mathbf{x}) - \frac{1}{n_D} \boldsymbol{\Omega} \mathbf{1} \right)^T \mathbf{N}_c(\mu \mathbf{I}_{n_D} + \zeta \boldsymbol{\Omega})^{-1} \mathbf{N}_c \left( \mathbf{k}(\mathbf{x}) - \frac{1}{n_D} \boldsymbol{\Omega} \mathbf{1} \right) \end{aligned} \tag{48}$$

with the vector  $\mathbf{k}(\mathbf{x}) = \boldsymbol{\Phi} \boldsymbol{\varphi}(\mathbf{x}) = [K(\mathbf{x}, \mathbf{x}_1), \dots, K(\mathbf{x}, \mathbf{x}_{n_D})]^T \in \mathbb{R}^{n_D}$ . This dual space expression allows us to compute the variance on the point prediction  $y_{MP}$  when the nonlinear mapping  $\boldsymbol{\varphi}$  is implicitly defined by the (nonlinear) kernel function  $K$  or when  $n \geq n_D$  in the linear case.

APPENDIX B: MATHEMATICS

**Expression for the Hessian and covariance matrix**

The level 1 posterior probability  $p([\mathbf{w}; b]|\mathcal{D}, \mu, \zeta, \mathcal{H})$  is a multivariate normal distribution in  $\mathbb{R}^{n_\theta}$  with mean  $[\mathbf{w}_{MP}; b_{MP}]$  and covariance matrix  $\mathbf{Q} = \mathbf{H}^{-1}$ , where  $\mathbf{H}$  is the Hessian of the least squares cost function (9). Defining the matrix of regressors  $\boldsymbol{\Phi}^T = [\boldsymbol{\varphi}(\mathbf{x}_1), \dots, \boldsymbol{\varphi}(\mathbf{x}_{n_\theta})]$ , the identity matrix  $\mathbf{I}$  and the vector with all ones  $\mathbf{1}$  of appropriate dimension; the Hessian is equal to

$$\mathbf{H} = \begin{bmatrix} \mathbf{H}_{11} & \mathbf{h}_{12} \\ \mathbf{h}_{21} & h_{22} \end{bmatrix} = \begin{bmatrix} \mu \mathbf{I}_{n_\phi} + \zeta \Phi^T \Phi & \zeta \Phi^T \mathbf{1} \\ \zeta \mathbf{1}^T \Phi & \zeta n_{\mathcal{D}} \end{bmatrix} \quad (49)$$

with corresponding block matrices  $\mathbf{H}_{11} = \mu \mathbf{I}_{n_\phi} + \zeta \Phi^T \Phi$ ,  $\mathbf{h}_{12} = \mathbf{h}_{21}^T = \Phi^T \mathbf{1}$  and  $h_{22} = n_{\mathcal{D}}$ . The inverse Hessian  $\mathbf{H}^{-1}$  is then obtained via a Schur complement type argument:

$$\begin{aligned} \mathbf{H}^{-1} &= \left( \begin{bmatrix} \mathbf{I}_{n_\phi} & X \\ \mathbf{0}^T & 1 \end{bmatrix} \begin{bmatrix} \mathbf{I}_{n_\phi} & -X \\ \mathbf{0}^T & 1 \end{bmatrix} \begin{bmatrix} \mathbf{H}_{11} & \mathbf{h}_{12} \\ \mathbf{h}_{12}^T & h_{22} \end{bmatrix} \begin{bmatrix} \mathbf{I}_{n_\phi} & 0 \\ -X^T & 1 \end{bmatrix} \begin{bmatrix} \mathbf{I}_{n_\phi} & 0 \\ -X^T & 1 \end{bmatrix} \right)^{-1} \\ &= \left( \begin{bmatrix} \mathbf{I}_{n_\phi} & X \\ \mathbf{0}^T & 1 \end{bmatrix} \begin{bmatrix} \mathbf{H}_{11} - \mathbf{h}_{12} h_{22}^{-1} \mathbf{h}_{12}^T & 0 \\ \mathbf{0}^T & h_{22} \end{bmatrix} \begin{bmatrix} \mathbf{I}_{n_\phi} & 0 \\ X^T & 1 \end{bmatrix} \right)^{-1} \end{aligned} \quad (50)$$

$$= \begin{bmatrix} (\mathbf{H}_{11} - \mathbf{h}_{12} h_{22}^{-1} \mathbf{h}_{12}^T)^{-1} & -\mathbf{F}_{11}^{-1} \mathbf{h}_{12} h_{22}^{-1} \\ -h_{22}^{-1} \mathbf{h}_{12}^T \mathbf{F}_{11}^{-1} & h_{22}^{-1} + h_{22}^{-1} \mathbf{h}_{12}^T \mathbf{F}_{11}^{-1} \mathbf{h}_{12} h_{22}^{-1} \end{bmatrix} \quad (51)$$

with  $X = \mathbf{h}_{12} h_{22}^{-1}$  and  $\mathbf{F}_{11} = \mathbf{H}_{11} - \mathbf{h}_{12} h_{22}^{-1} \mathbf{h}_{12}^T$ . In matrix expressions, it is useful to express  $\Phi^T \Phi - \frac{1}{n_{\mathcal{D}}}$

$\Phi^T \mathbf{1} \mathbf{1}^T \Phi$  as  $\Phi^T \mathbf{N}_c \Phi$  with the idempotent centring matrix  $\mathbf{N}_c = \mathbf{I}_{n_{\mathcal{D}}} - \frac{1}{n_{\mathcal{D}}} \mathbf{1} \mathbf{1}^T \in \mathbb{R}^{n_{\mathcal{D}} \times n_{\mathcal{D}}}$  having  $\mathbf{N}_c = \mathbf{N}_c^2$ .

Given that  $\mathbf{F}_{11}^{-1} = (\mu \mathbf{I}_{n_\phi} + \zeta \Phi^T \mathbf{N}_c \Phi)^{-1}$ , the inverse Hessian  $\mathbf{H}^{-1} = \mathbf{Q}$  is equal to

$$\mathbf{Q} = \begin{bmatrix} (\mu \mathbf{I}_{n_\phi} + \zeta \Phi^T \mathbf{N}_c \Phi)^{-1} & -\frac{1}{n_{\mathcal{D}}} (\mu \mathbf{I}_{n_\phi} + \zeta \Phi^T \mathbf{N}_c \Phi)^{-1} \Phi^T \mathbf{1} \\ -\frac{1}{n_{\mathcal{D}}} \mathbf{1}^T \Phi (\mu \mathbf{I}_{n_\phi} + \zeta \Phi^T \mathbf{N}_c \Phi)^{-1} & \frac{1}{\zeta n_{\mathcal{D}}} + \frac{1}{n_{\mathcal{D}}} \mathbf{1}^T \Phi (\mu \mathbf{I}_{n_\phi} + \zeta \Phi^T \mathbf{N}_c \Phi)^{-1} \Phi^T \mathbf{1} \end{bmatrix}$$

### Expression for the determinant

The determinant of  $\mathbf{H}$  is obtained from (50) using the fact that the determinant of a product is equal to the product of the determinants and is thus equal to

$$\begin{aligned} \det(\mathbf{H}) &= \det(\mathbf{H}_{11} - \mathbf{h}_{12} h_{22}^{-1} \mathbf{h}_{12}^T) \times \det(h_{22}) \\ &= \det(\mu \mathbf{I}_{n_\phi} + \zeta \Phi^T \mathbf{N}_c \Phi) \times (\zeta n_{\mathcal{D}}) \end{aligned} \quad (52)$$

which is obtained as the product of  $\zeta n_{\mathcal{D}}$  and the eigenvalues  $\lambda_i$  ( $i = 1, \dots, n_\phi$ ) of  $\mu \mathbf{I}_{n_\phi} + \zeta \Phi^T \mathbf{N}_c \Phi$ , denoted as  $\lambda_i(\mu \mathbf{I}_{n_\phi} + \zeta \Phi^T \mathbf{N}_c \Phi)$ . Because the matrix  $\Phi^T \mathbf{N}_c \Phi \in \mathbb{R}^{n_\phi \times n_\phi}$  is rank deficient with rank  $n_e \leq n_{\mathcal{D}} - 1$ ,  $n_\phi - n_e$  eigenvalues are equal to  $\mu$ .

The dual space expressions can be obtained in terms of the singular value decomposition

$$\Phi^T \mathbf{N}_c = \mathbf{U} \mathbf{S} \mathbf{V}^T = [\mathbf{U}_1 \quad \mathbf{U}_2] \begin{bmatrix} \mathbf{S}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} [\mathbf{V}_1 \quad \mathbf{V}_2] \quad (53)$$

with  $\mathbf{U} \in \mathbb{R}^{n_\phi \times n_\phi}$ ,  $\mathbf{S} \in \mathbb{R}^{n_\phi \times n_{\mathcal{D}}}$ ,  $\mathbf{V} \in \mathbb{R}^{n_{\mathcal{D}} \times n_{\mathcal{D}}}$  and with the block matrices  $\mathbf{U}_1 \in \mathbb{R}^{n_\phi \times n_e}$ ,  $\mathbf{U}_2 \in \mathbb{R}^{n_\phi \times (n_\phi - n_e)}$ ,  $\mathbf{S}_1 = \text{diag}([s_1, s_2, \dots, s_{n_e}]) \in \mathbb{R}^{n_e \times n_e}$ ,  $\mathbf{V}_1 \in \mathbb{R}^{n_{\mathcal{D}} \times n_e}$  and  $\mathbf{V}_2 \in \mathbb{R}^{n_{\mathcal{D}} \times (n_{\mathcal{D}} - n_e)}$ , with  $0 \leq n_e \leq n_{\mathcal{D}} - 1$ . Due to the

orthonormality property we have  $UU^T = U_1U_1^T + U_2U_2^T = \mathbf{I}_{n_\phi}$  and  $VV^T = V_1V_1^T + V_2V_2^T = \mathbf{I}_{n_D}$ . Hence, one obtains the primal and dual eigenvalue decompositions

$$\Phi^T \mathbf{N}_c \Phi = U_1 S_1^2 U_1^T \quad (54)$$

$$\mathbf{N}_c \Phi \Phi^T \mathbf{N}_c = \mathbf{N}_c \Omega \mathbf{N}_c = V_1 S_1^2 V_1^T \quad (55)$$

The  $n_\phi$  eigenvalues of  $\mu \mathbf{I}_{n_\phi} + \zeta \Phi^T \mathbf{N}_c \Phi$  are equal to  $\lambda_1 = \mu + \zeta s_1^2, \dots, \lambda_{n_e} = \mu + \zeta s_{n_e}^2, \lambda_{n_e+1} = \mu, \dots, \lambda_{n_\phi} = \mu$ , where the nonzero eigenvalues  $s_i^2$  ( $i = 1, \dots, n_e$ ) are obtained from the eigenvalue decomposition of  $\mathbf{N}_c \Phi \Phi^T \mathbf{N}_c$  from (55). The expression for the determinant is equal to  $N \zeta \mu^{n_D - n_e} \prod_{i=1}^{n_e} (\mu + \zeta \lambda_i(\mathbf{N}_c \Omega \mathbf{N}_c))$ , with  $\mathbf{N}_c \Omega \mathbf{N}_c = V_1 \text{diag}([\lambda_1, \dots, \lambda_{n_e}]) V_1^T$  and  $\lambda_i = s_i^2, i = 1, \dots, n_e$ .

### Expression for the level 1 cost function

The dual space expression for  $J_1(\mathbf{w}_{MP}, b_{MP})$  is obtained by substituting  $[\mathbf{w}_{MP}; b_{MP}] = \mathbf{H}^{-1}[\Phi^T \mathbf{y}; \mathbf{1}^T \mathbf{y}]$  in (9). Applying a similar reasoning and algebra as for the calculation of the determinant, one obtains the dual space expression:

$$J_1(\mathbf{w}, b) = \mu J_w(\mathbf{w}_{mp}) + \zeta J_e(\mathbf{w}_{MP}, b_{MP}) = \frac{1}{2} \mathbf{y}^T \mathbf{N}_c (\mu^{-1} \mathbf{N}_c \Omega \mathbf{N}_c + \zeta^{-1} \mathbf{I}_{n_D})^{-1} \mathbf{N}_c \mathbf{y} \quad (56)$$

Given that  $\mathbf{N}_c \Omega \mathbf{N}_c = \mathbf{V} \Lambda \mathbf{V}^T$ , with  $\Lambda = \text{diag}([s_1^2, \dots, s_{n_e}^2, 0, \dots, 0])$ , one obtains (41). In a similar way, one obtains (39) and (40).

## ACKNOWLEDGEMENTS

The authors would like to thank Peter Van Dijke (Dexia Bank), Joao Garcia, Luc Leonard, Eric Hermann (Dexia Group) and Dirk Baestaens (Fortis Bank) for many helpful comments. This work was partially supported by grants and projects from the K.U.Leuven (GOA-Mefisto 666, GOA-Ambiorics), the Flemish Government (FWO Projects G.0407.02, G.0499.04, G.0211.05., ICCoS, ANMMM, IWT, GBOU), the Belgian Federal Government (IUAP V-22, PODO-II) and the EU (Ernsi, Eureka 2063, 2419). Scientific responsibility is assumed by its authors.

## REFERENCES

- Akaike H. 1974. A new look at statistical model identification. *IEEE Transactions on Automatic Control* **19**: 716–723.
- Bachelier L. 1900. *Théorie de la spéculation*. Gauthier-Villars: Paris.
- Bishop CM. 1995. *Neural Networks for Pattern Recognition*. Oxford University Press: Oxford.
- Box GEP, Jenkins GM. 1970. *Time Series Analysis, Forecasting and Control*. Holden-Day: San Francisco.
- Brock W, Lakonishok J, Le Baron B. 1992. Simple technical trading rules and the stochastic properties of econometrics. *Journal of Finance* **47**: 1731–1764.
- Campbell JY, Lo AW, MacKinlay AC. 1997. *The Econometrics of Financial Markets*. Princeton University Press: Princeton, NJ.
- Dickey DA, Fuller WA. 1979. Distribution of the estimators for autoregressive time series with a unit root. *Journal of the American Statistical Association* **74**: 427–431.

- Engle RF, Granger CWJ. 1987. Cointegration and error correction: representations, estimation and testing. *Econometrica* **55**: 252–276.
- Fama EF. 1965. The behaviour of stock market prices. *Journal of Business* **38**: 34–105.
- Granger CWJ, Newbold P. 1986. *Forecasting Economic Time Series*. Academic Press: New York.
- Granger CWJ, Newbold P. 1974. Spurious regression in econometrics. *Journal of Econometrics* **2**: 111–120.
- Granger CWJ, Terasvirta T. 1993. *Modelling Nonlinear Economic Relationships*. Oxford University Press: Oxford.
- Hamilton J. 1994. *Time Series Analysis*. Princeton University Press: Princeton, NJ.
- Hutchinson JM, Lo AW, Poggio T. 1994. A nonparametric approach to pricing and hedging derivative securities via learning networks. *Journal of Finance* **49**: 851–889.
- Jeffreys H. 1961. *Theory of Probability*. Oxford University Press: Oxford.
- Johansen S. 1988. Statistical analysis of cointegration vectors. *Journal of Economics Dynamics and Control* **12**: 231–254.
- Lo A, Mamaysky H, Wang J. 2000. Foundations of technical analysis: computational algorithms, statistical inference, and empirical implementation. *Journal of Finance* **55**: 1705–1765.
- MacKay DJC. 1995. Probable networks and plausible predictions—a review of practical Bayesian methods for supervised neural networks. *Network: Computation in Neural Systems* **6**: 469–505.
- Maddala GS, Kim IM. 1998. *Cointegration, Unit Roots and Structural Change*. Cambridge University Press: Cambridge.
- Pesaran MH, Timmerman A. 1992. A simple nonparametric test of predictive performance. *Journal of Business and Economic Statistics* **10**: 461–465.
- Rao B. 1994. *Cointegration for the Applied Economist*. MacMillan: London.
- Refenes AP, Zapranis AD. 1999. Neural model identification, variable selection and model adequacy. *Journal of Forecasting* **18**: 299–332.
- Schölkopf B, Smola A. 2002. *Learning with Kernels*. MIT Press: Cambridge, MA.
- Schwarz G. 1978. Estimating the dimension of a model. *Annals of Statistics* **6**: 461–464.
- Sullivan R, Timmerman A, White H. 1999. Data-snooping, technical trading rule performance, and the bootstrap. *Journal of Finance* **54**: 1647–1691.
- Suykens JAK, Van Gestel T, De Brabanter J, De Moor B, Vandewalle J. 2002. *Least Squares Support Vector Machines*. World Scientific: Singapore.
- Van Gestel T, Suykens JAK, Baestaens DE, Lambrechts A, Lanckriet G, Vandaele B, De Moor B, Vandewalle J. 2001. Predicting financial time series using least squares support vector machines within the evidence framework. *IEEE Transactions on Neural Networks (Special Issue on Financial Engineering)* **12**: 809–821.
- Van Gestel T, Suykens JAK, Lanckriet G, Lambrechts A, De Moor B, Vandewalle J. 2002. A Bayesian framework for least squares support vector machine classifiers, Gaussian processes and kernel Fisher discriminant analysis. *Neural Computation* **14**: 1115–1147.
- Vapnik V. 1998. *Statistical Learning Theory*. John Wiley & Sons: New York.

#### Authors' biographies:

**Tony Van Gestel** obtained an electromechanical engineering degree and a PhD in applied sciences (subject: mathematical modelling for financial engineering) in 1997 and 2002 at the Katholieke Universiteit Leuven. He currently works as a senior quantitative analyst at Dexia Group and is a free postdoctoral researcher at the Katholieke Universiteit Leuven.

**Marcelo Espinoza** obtained a degree in civil engineering and an MSc in applied economics in 1998 at the University of Chile, and the degree of Master in Artificial Intelligence in 2002 at the Katholieke Universiteit Leuven. After 4 years' experience in international commodity trading, he is currently pursuing a PhD at the Katholieke Universiteit Leuven.

**Bart Baesens** obtained the degree of Master in Management Informatics and a PhD in applied economic sciences at the K.U.Leuven (Belgium) in 1998 and 2003, respectively. He currently works as a lecturer (assistant professor) at the School of Management, University of Southampton (UK), and the Department of Applied Economics, K.U.Leuven (Belgium).

**Johan Suykens** obtained a degree in electro-mechanical engineering and a PhD (subject: artificial neural networks) in applied sciences from the Katholieke Universiteit Leuven in 1989 and 1995, respectively. In 1996

he was a visiting postdoctoral researcher at the University of California, Berkeley. He is currently an associated professor at the Department of Electrical Engineering, Katholieke Universiteit Leuven.

**Bart De Moor** received his doctoral degree in applied sciences in 1988 at the Katholieke Universiteit Leuven, Belgium. He was a visiting research associate (1988–1989) at the Department of Computer Science and Electrical Engineering of Stanford University, California. Bart De Moor is a full professor at the Katholieke Universiteit Leuven.

**Carine Brasseur** received her doctoral degree in economic sciences in 2000 at the Université catholique de Louvain, Belgium. She currently works as a senior strategist in the Research Department at Global Markets, Fortis Bank Belgium.

*Authors' addresses:*

**Tony Van Gestel**, Credit Risk Modelling, Risk Management, Dexia Group, Square Meeus 1, B-1000 Brussels, Belgium.

**Tony Van Gestel, Marcelo Espinoza, Johan A. K. Suykens and Bart De Moor**, Katholieke Universiteit Leuven, Department of Electrical Engineering ESAT-SISTA, Kasteelpark Arenberg 10, B-3001 Leuven, Belgium.

**Bart Baesens**, School of Management, University of Southampton, Southampton SO17 1BJ, UK.

**Carine Brasseur**, Financial Markets, Fortis Bank Brussels, Warandeberg 3, B-1000 Brussels, Belgium.