**Faculty of Mathematics**

Mathematical Finance

Alois Pichler

# Regression With Kernels

**Perspective:** *Seminar thesis, Bachelor, Master's thesis.*

## Problem Description

Suppose the function $f\colon \mathbb{R}^d \to \mathbb{R}$ is observed with function values, $(x_i, f_i) \in \mathbb{R}^d \times \mathbb{R}$, $i = 1, \ldots, n$. Employing a kernel function $k\colon \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$, it holds that

$$\hat{f}(x_i) = f_i, \qquad i = 1, \ldots, n, \tag{1}$$

where

$$\hat{f}(x) := \sum_{j=1}^{n} k(x, x_j)\, w_j$$

and provided that the weights $w$ satisfy the linear system of equations

$$\sum_{j=1}^{n} k(x_i, x_j)\, w_j = f_i, \qquad i = 1, \ldots, n.$$

A canonical candidate for the kernel is the Gaussian kernel $k(x, y) := e^{-\|y-x\|^2 / (2\eta^2)}$ (with Euclidean norm $\|\cdot\|$ in $\mathbb{R}^d$ and characteristic length $\eta$).

For observations $f_i$ *with noise*, the *kernel estimator* for $f(x)$ is

$$\hat{f}(x) = \sum_{j=1}^{n} k(x, x_j)\, w_j \tag{2}$$

with

$$\lambda\, w_i + \sum_{j=1}^{n} k(x_i, x_j)\, w_j = f_i, \qquad i = 1, \ldots, n. \tag{3}$$

The regression parameter $\lambda$ is chosen to match the noise: $\lambda = 0$ matches the initial situation so that $\hat{f}(x_i) = f_i$ without noise (interpolation, cf. (1)). A value of $\lambda$ too close to zero leads to *overfitting*, larger values for $\lambda$ lead to *oversmoothing* (underfitting). $\lambda$ acts as a regularization parameter, balancing bias and variance to prevent overfitting while maintaining model flexibility.

# Task: Delayed Embedding of Dimension $d$

We shall employ (2) and (3) to a stationary time series $X = (X_t)$ with $X_t \in \mathbb{R}$, $\mathbb{E}\, X_t = 0$, and constant variance

$$\sigma^2 := \mathrm{var}\big(X_{t+1}|\, X_{t-d+1}, \ldots, X_t\big),$$

where $d \in \mathbb{N}$ is fixed.

Consider the observation (realization)

$$X_{2-d}, \ldots, X_1, X_2, \ldots, X_{n+1} \tag{4}$$

of the time series. In the *sliding window approach*, set

$$x_i := (X_{i-d+1}, X_{i-d+2}, \ldots, X_i) \in \mathbb{R}^d \quad \text{and} \quad f_i := X_{i+1} \in \mathbb{R}, \qquad i = 1, \ldots, n, \tag{5}$$

then $\hat{f}\big((X_{t-d+1}, \ldots, X_t)\big)$ is an *estimate* for the subsequent $X_{t+1}$,

$$\hat{f}\big((X_{t-d+1}, \ldots, X_t)\big) \approx \mathbb{E}\big[X_{t+1}|\, X_{t-d+1}, \ldots, X_t\big].$$

1. Discuss and implement the relations (2), (3) and (5) for some adequate, oscillating time series observations $(X_t)$ as in (4).

2. **Recursive multistep forecasting:** For some chosen starting values $(\hat{X}_1, \ldots, \hat{X}_d)$, extend this time series by setting

$$\hat{X}_{t+1} := \hat{f}\big((\hat{X}_{t-d+1}, \ldots, \hat{X}_t)\big) + \varepsilon_t, \qquad t = d,\, d+1, \ldots$$

   and *plot* the realization $\hat{X}_1, \hat{X}_2, \ldots, \hat{X}_d, \ldots, \hat{X}_n, \ldots, \hat{X}_N$ with some $N \gg d$; the iid random errors $\varepsilon_t$ have $0$-mean, and variance corresponding to the residual variance,

$$\mathrm{var}\, \varepsilon_t \approx \sigma^2 = \mathrm{var}\big(X_{t+1}|\, (X_{t-d+1}, \ldots, X_t)\big).$$

3. Adjust the parameters $\lambda$, $d$, $\sigma$ and the kernel $k$ (or the parameter in the kernel, $\eta$, say) so that the time series $\hat{X}$ *visually* matches the initial observations (4) of $X$.

4. Try the norm $\|x\|_\Sigma := x^\top \Sigma^{-1} x$ (Mahalanobis distance) instead of the Euclidean norm, where $\Sigma$ is the covariance matrix of the random vector $(X_{i-d+1}, \ldots, X_i)$. The covariance matrix can be estimated from the observations (5).

Bonne chance!