

Chapter 6 Relationships between GPs and Other Models, Part I

Randy

8/25/2021

1 6.1 Reproducing Kernel Hilbert Spaces

2 6.2 Regularization

3 6.3 Spline Models

Section 1

6.1 Reproducing Kernel Hilbert Spaces

RKHS

Definition 6.1 (Reproducing kernel Hilbert space)

Let \mathcal{H} be a Hilbert space of real functions f defined on an index set \mathcal{X} . Then \mathcal{H} is called a reproducing kernel Hilbert space endowed with an inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ (and norm $\|f\|_{\mathcal{H}} = \sqrt{\langle f, f \rangle_{\mathcal{H}}}$);

If there exists a function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathfrak{R}$ with the following properties:

- $\forall \mathbf{x}, k(\mathbf{x}, \mathbf{x}')$ as a function of \mathbf{x}' belongs to \mathcal{H}
- k has the reproducing property $\langle f(\cdot), k(\cdot, \mathbf{x}) \rangle_{\mathcal{H}} = f(\mathbf{x})$

Note also that as $k(\mathbf{x}, \cdot)$ and $k(\mathbf{x}', \cdot)$ are in \mathcal{H} we have that $\langle k(\mathbf{x}, \cdot), k(\mathbf{x}', \cdot) \rangle_{\mathcal{H}} = k(\mathbf{x}, \mathbf{x}')$. The RKHS uniquely determines k , and vice versa.

Moore-Aronszajn Theorem

Theorem 6.1 (Moore-Aronszajn theorem, Aronszajn [1950])

Let \mathcal{X} be an index set. Then for every positive definite function $k(\langle \cdot, \cdot \rangle)$ on $\mathcal{X} \times \mathcal{X}$ there exists a unique RKHS, and vice versa.

The Hilbert space L_2 (which has the dot product $\langle f, g \rangle_{L_2} = \int f(\mathbf{x})g(\mathbf{x})d\mathbf{x}$) contains many non-smooth functions.

In L_2 (which is not a RKHS) the delta function is the representer of evaluation, i.e. $f(\mathbf{x}) = \int f(\mathbf{x}')\delta(\mathbf{x} - \mathbf{x}')d\mathbf{x}'$. Kernels are the analogues of delta functions within the smoother RKHS.

Note that the delta function is not itself in L_2 ; in contrast for a RKHS the kernel k is the representer of evaluation and is itself in the RKHS.

The key intuition behind the RKHS formalism is that the squared norm $\|f\|_{\mathcal{H}}^2$ can be thought of as a generalization to functions of the n -dimensional quadratic form $\mathbf{f}^\top K^{-1} \mathbf{f}$ we have seen in earlier chapters.

Consider a real positive semidefinite kernel $k(\mathbf{x}, \mathbf{x}')$ with an eigenfunction expansion $k(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^N \lambda_i \phi_i(\mathbf{x}) \phi_i(\mathbf{x}')$ relative to a measure μ .

Recall from Mercer's theorem that the eigenfunctions are orthonormal w.r.t. μ , i.e. we have $\int \phi_i(\mathbf{x}) \phi_j(\mathbf{x}) d\mu(\mathbf{x}) = \delta_{ij}$.

We now consider a Hilbert space comprised of linear combinations of the eigenfunctions, i.e. $f(\mathbf{x}) = \sum_{i=1}^N f_i \phi_i(\mathbf{x})$ with $\sum_{i=1}^N f_i^2 / \lambda_i < \infty$.

We assert that the inner product $\langle f, g \rangle_{\mathcal{H}}$ in the Hilbert space between functions $f(\mathbf{x})$ and $g(\mathbf{x}) = \sum_{i=1}^N g_i \phi_i(\mathbf{x})$ is defined as

$$\langle f, g \rangle_{\mathcal{H}} = \sum_{i=1}^N \frac{f_i g_i}{\lambda_i} \quad (6.1)$$

Thus this Hilbert space is equipped with a norm $\|f\|_{\mathcal{H}}$ where

$$\|f\|_{\mathcal{H}}^2 = \langle f, f \rangle_{\mathcal{H}} = \sum_{i=1}^N f_i^2 / \lambda_i.$$

Note that for $\|f\|_{\mathcal{H}}$ to be finite the sequence of coefficients $\{f_i\}$ must decay quickly; effectively this imposes a smoothness condition on the space.

Reproducing property

$$\langle f(\cdot), k(\cdot, \mathbf{x}) \rangle_{\mathcal{H}} = \sum_{i=1}^N \frac{f_i \lambda_i \phi_i(\mathbf{x})}{\lambda_i} = f(\mathbf{x}) \quad (6.2)$$

$$\langle k(\mathbf{x}, \cdot), k(\mathbf{x}', \cdot) \rangle_{\mathcal{H}} = \sum_{i=1}^N \frac{\lambda_i \phi_i(\mathbf{x}) \lambda_i \phi_i(\mathbf{x}')}{\lambda_i} = k(\mathbf{x}, \mathbf{x}') \quad (6.3)$$

Notice also that $k(\mathbf{x}, \cdot)$ is in the RKHS as it has norm

$$\sum_{i=1}^N (\lambda_i \phi_i(\mathbf{x}))^2 / \lambda_i = k(\mathbf{x}, \mathbf{x}) < \infty.$$

The Hilbert space comprised of linear combinations of the eigenfunctions with the restriction $\sum_{i=1}^N f_i^2 / \lambda_i < \infty$ fulfills the two conditions given in Definition 6.1. As there is a unique RKHS associated with $k(\cdot, \cdot)$, this Hilbert space must be that RKHS.

The advantage of the abstract formulation of the RKHS

The eigenbasis will change with different measures μ in Mercer's theorem. However, the RKHS norm is in fact solely a property of the kernel and is invariant under this change of measure.

Notice the analogy between the RKHS norm

$\|f\|_{\mathcal{H}}^2 = \langle f, f \rangle_{\mathcal{H}} = \sum_{i=1}^N f_i^2 / \lambda_i$ and the quadratic form $\mathbf{f}^\top K^{-1} \mathbf{f}$ if we express K and f in terms of the eigenvectors of K we obtain exactly the same form (but the sum has only n terms if f has length n).

If we sample the coefficients f_i in the eigenexpansion

$f(\mathbf{x}) = \sum_{i=1}^N f_i \phi_i(\mathbf{x})$ from $\mathcal{N}(0, \lambda_i)$ then

$$E[\|f\|_{\mathcal{H}}^2] = \sum_{i=1}^N E[f_i^2] \lambda_i = \sum_{i=1}^N 1 \quad (6.4)$$

Thus if N is infinite the sample functions are not in \mathcal{H} (with probability 1) as the expected value of the RKHS norm is infinite

However, note that although sample functions of this Gaussian process are not in \mathcal{H} , the posterior mean after observing some data will lie in the RKHS, due to the smoothing properties of averaging.

Another view of the RKHS can be obtained from the reproducing kernel map construction. We consider the space of functions f defined as

$$\{f(\mathbf{x}) = \sum_{i=1}^N \alpha_i k(\mathbf{x}, \mathbf{x}_i) : n \in N, \mathbf{x}_i \in \mathcal{X}, \alpha_i \in \mathfrak{R}\} \quad (6.5)$$

Now let $g(\mathbf{x}) = \sum_{j=1}^{n'} \alpha'_j k(\mathbf{x}, \mathbf{x}'_j)$. Then we define the inner product

$$\langle f, g \rangle_{\mathcal{H}} = \sum_{i=1}^n \sum_{j=1}^{n'} \alpha_i \alpha'_j k(\mathbf{x}_i, \mathbf{x}'_j) \quad (6.6)$$

Clearly condition 1 of Definition 6.1 is fulfilled under the reproducing kernel map construction. We can also demonstrate the reproducing property, as $\langle k(\cdot, \mathbf{x}), f(\cdot) \rangle_{\mathcal{H}} = \sum_{i=1}^N \alpha_i k(\mathbf{x}, \mathbf{x}_i) = f(\mathbf{x}) \quad (6.7)$

Section 2

6.2 Regularization

Inferring from a finite dataset without any assumption is clearly “ill posed”. For example, in the noise-free case, any function that passes through the given data points is acceptable.

Under a Bayesian approach our assumptions are characterized by a prior over functions, and given some data, we obtain a posterior over functions.

The problem of bringing prior assumptions to bear has also been addressed under the regularization viewpoint, where these assumptions are encoded in terms of the smoothness of f : $J[f] = \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2 + Q(\mathbf{y}, \mathbf{f})$ (6.8)

- The first term is regularizer and represents smoothness assumptions on f as encoded by a suitable RKHS
- The second term is a data-fit term assessing the quality of the prediction $f(x_i)$ for the observed datum y_i , e.g. the negative log likelihood.

Ridge Regression

Ridge regression can be seen as a particular case of regularization:

$\|f\|_{\mathcal{H}}^2 = \sum_{i=1}^N f_i^2 / \lambda_i$ where f_i is the coefficient of eigenfunction $\phi_i(\mathbf{x})$, we see that we are penalizing the weighted squared coefficients.

This is taking place in feature space, rather than simply in input space, as per the standard formulation of ridge regression, so it corresponds to **kernel ridge regression**.

Hastie and Tibshirani eq(3.41) Ridge regression

The ridge coefficients minimized the penalized residual sum of squares:

$$\hat{\beta}^{ridge} = \operatorname{argmin}_{\beta} \left\{ \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum \beta_j^2 \right\}$$

The representer theorem shows that each minimizer $f \in \mathcal{H}$ of $J[f]$ has the form $f(\mathbf{x}) = \sum_{i=1}^N \alpha_i k(\mathbf{x}, \mathbf{x}_i)$. The representer theorem was first stated by **Kimeldorf and Wahba [1971]** for the case of squared error. **O'Sullivan et al. [1986]** showed that the representer theorem could be extended to likelihood functions arising from generalized linear models.

If the data-fit term is convex (see section A.9) then there will be a unique minimizer \hat{f} of $J[f]$. (we can regard the $J[f]$ as the penalized likelihood, in ESL [Hastie and Tibshirani] the $J[f]$ is the regularization term)

For Gaussian process prediction with likelihoods that involve the values of f at the n training points only (so that $Q(y, f)$ is the negative log likelihood up to some terms not involving f), the analogue of the representer theorem is obvious.

This is because the predictive distribution of $f(x_*) \triangleq f_*$ at test point x_* given the data \mathbf{y} is $p(f_*|\mathbf{y}) = \int p(f_*|\mathbf{f})p(\mathbf{f}|\mathbf{y})d\mathbf{f}$. As derived in eq. (3.22) we have $E[f_*|\mathbf{y}] = \mathbf{k}(\mathbf{x}_*)^\top K^{-1}E[\mathbf{f}|\mathbf{y}]$ (6.9) due to the formulae for the conditional distribution of a multivariate Gaussian.

Thus $E[f_*|\mathbf{y}] = \sum_{i=1}^N \alpha_i k(\mathbf{x}_*, \mathbf{x}_i)$, where $\boldsymbol{\alpha} = K^{-1}E[\mathbf{f}|\mathbf{y}]$.

6.2.1 Regularization Defined by Differential Operators

For $\mathbf{x} \in \mathfrak{R}^D$ define $\|O^m f\|^2 = \int \sum_{j_1 + \dots + j_D = m} \left(\frac{\partial^m f(\mathbf{x})}{\partial x_1^{j_1} \dots \partial x_D^{j_D}} \right)^2 dx$ (6.10)

For example for $m = 2$ and $D = 2$:

$$\|O^2 f\|^2 = \int \left[\left(\frac{\partial^2 f}{\partial x_1^2} \right)^2 + 2 \left(\frac{\partial^2 f}{\partial x_1 \partial x_2} \right)^2 + \left(\frac{\partial^2 f}{\partial x_2^2} \right)^2 \right] dx_1 dx_2 \quad (6.11)$$

Now set $\|Pf\|^2 = \sum_{m=0}^M a_m \|O^m f\|^2$ with non-negative coefficients a_m .
Notice that $\|Pf\|^2$ is translation and rotation invariant.

$\|O^m f\|$ is the penalty term; $\|Pf\|$ is the a projection in null space

In this section we assume that $a_0 > 0$; if this is not the case and a_k is the first non-zero coefficient, then there is a null space of functions that are unpenalized. For example if $k = 2$ then constant and linear functions are in the null space. This case is dealt with in section 6.3.

$\|Pf\|^2$ penalizes f in terms of the variability of its function values and derivatives up to order M . The key is to recognize that the complex exponentials $\exp(2\pi i \mathbf{s} \cdot \mathbf{x})$ are eigenfunctions of the differential operator if $\mathcal{X} = \mathfrak{R}^D$.

$$\|Pf\|^2 = \int \sum_{m=0}^M a_m (4\pi^2 \mathbf{s} \cdot \mathbf{s})^m |\tilde{f}(\mathbf{s})|^2 d\mathbf{s} \quad (6.12)$$

where $\tilde{f}(\mathbf{s})$ is the Fourier transform of $f(\mathbf{x})$. Comparing eq. (6.12) with eq. (6.1) we see that the kernel has the power spectrum

$$S(s) = \frac{1}{\sum_{m=0}^M a_m (4\pi^2 \mathbf{s} \cdot \mathbf{s})^m} \quad (6.13) \text{ and thus by Fourier inversion we obtain}$$

$$\text{the stationary kernel } k(\mathbf{x}) = \int \frac{e^{2\pi \mathbf{s} \cdot \mathbf{x}}}{\sum_{m=0}^M a_m (4\pi^2 \mathbf{s} \cdot \mathbf{s})^m} d\mathbf{s} \quad (6.14), \text{ by Bochner}$$

Theorem and Wiener-Khintchin Theorem (explained why penalize higher frequency functions)

A slightly different approach to obtaining the kernel is to use calculus of variations to minimize $J[f]$ with respect to f . The Euler-Lagrange equation leads to: $f(\mathbf{x}) = \sum_{i=1}^n \alpha_i G(\mathbf{x} - \mathbf{x}_i)$ (6.15) and

$\sum_{m=0}^M (-1)^m a_m \nabla^{2m} G = \delta(\mathbf{x} - \mathbf{x}')$ (6.16), where $G(\mathbf{x}, \mathbf{x}')$ is known as a Green's function. Notice that the Green's function also depends on the boundary conditions.

For the case of $\mathcal{X} = \Re^D$ by Fourier transforming eq. (6.16) we recognize that G is in fact the kernel k .

The differential operator $\sum_{m=0}^M (-1)^m a_m \nabla^{2m}$ and the integral operator $k(\langle \cdot, \cdot \rangle)$ are in fact inverses, as shown by eq. (6.16). Arfken [1985] provides an introduction to calculus of variations and Green's functions. RKHSs for regularizers defined by differential operators are **Sobolev spaces**.

Sidenotes for Green's function

If L is the linear differential operator, then:

- the Green's function G is the solution of the equation $LG = \delta$, where δ is Dirac's delta function;
- the solution of the initial-value problem $Ly = f$ is the convolution $(G * f)$, where G is the Green's function.

Green's function

$G(x, s)$ of a linear differential operator $L = L(x)$ acting on distributions over a subset of the Euclidean space \mathfrak{R}^n , at a point s , is any solution of $L G(x, s) = \delta(s - x)$

where δ is the Dirac delta function.

This property of a Green's function can be exploited to solve differential equations of the form $L u(x) = f(x)$

If the operator is translation invariant, that is, when L has constant coefficients with respect to x , then the Green's function can be taken to be a convolution kernel, $G(x, s) = G(x - s)$.

Motivation

A function G can be found for the operator L , then multiply the Green's function by $f(s)$, and then integrate with respect to s , we obtain

$$\int L G(x, s) f(s) ds = \int \delta(x - s) f(s) ds = f(x).$$

Because the operator $L = L(x)$ is linear and acts only on the variable x (and not on the variable of integration s), one may take the operator L outside of the integration, yielding $L \int G(x, s) f(s) ds = f(x)$

This means that $u(x) = \int G(x, s) f(s) ds$ is a solution to the equation $L u(x) = f(x)$

Example 1

Set $a_0 = \alpha^2$, $a_1 = 1$ and $a_m = 0$ for $m \geq 2$ in $D = 1$. Using the Fourier pair $e^{-\alpha|x|} \Leftrightarrow 2\alpha/(\alpha^2 + 4\pi^2 s^2)$ we obtain $k(x - x') = \frac{1}{2\alpha} e^{-\alpha|x-x'|}$.

Note that this is the covariance function of the Ornstein-Uhlenbeck process, see section 4.2.1.

Example 2

By setting $a_m = \frac{\sigma^{2m}}{m!2^m}$ and using the power series $e^y = \sum_{k=0}^{\infty} y^k/k!$ we obtain

$$k(\mathbf{x} - \mathbf{x}') = \int \exp(2\pi i \mathbf{s} \cdot (\mathbf{x} - \mathbf{x}')) \exp\left(-\frac{\sigma^2}{2}(4\pi^2 \mathbf{s} \cdot \mathbf{s})\right) d\mathbf{s} \quad (6.17)$$

$$k(\mathbf{x} - \mathbf{x}') = \frac{1}{(2\pi\sigma^2)^{D/2}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{x} - \mathbf{x}')^\top (\mathbf{x} - \mathbf{x}')\right) \quad (6.18)$$

as shown by Yuille and Grzywacz [1989]. This is the squared exponential covariance function that we have seen earlier

6.2.2 Obtaining the Regularized Solution

The representer theorem tells us the general form of the solution to eq. (6.8). We now consider a specific functional

$$J[f] = \frac{1}{2} \|f\|_{\mathcal{H}}^2 + \frac{1}{2\sigma_n^2} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 \quad (6.19)$$

which uses a squared error data-fit term (corresponding to the negative log likelihood of a Gaussian noise model with variance σ_n^2).

The solution $f(\mathbf{x}) = \sum_{i=1}^N \alpha_i k(\mathbf{x}, \mathbf{x}_i)$ that minimizes eq. (6.19) was called a regularization network regularization network in Poggio and Girosi [1990].

Substituting $f(\mathbf{x}) = \sum_{i=1}^N \alpha_i k(\mathbf{x}, \mathbf{x}_i)$ and using $\langle k(\cdot, \mathbf{x}_i), k(\cdot, \mathbf{x}_j) \rangle_{\mathcal{H}} = k(\mathbf{x}_i, \mathbf{x}_j)$ we obtain

$$\begin{aligned} J[\boldsymbol{\alpha}] &= \frac{1}{2} \boldsymbol{\alpha}^\top K \boldsymbol{\alpha} + \frac{1}{2\sigma_n^2} \|\mathbf{y} - K \boldsymbol{\alpha}\|^2 \\ &= \frac{1}{2} \boldsymbol{\alpha}^\top (K + \frac{1}{\sigma_n^2} K^2) \boldsymbol{\alpha} - \frac{1}{\sigma_n} \mathbf{y}^\top K \boldsymbol{\alpha} + \frac{1}{2\sigma_n^2} \mathbf{y}^\top \mathbf{y} \quad (6.20) \end{aligned}$$

Minimizing J by differentiating w.r.t. the vector of coefficients $\boldsymbol{\alpha}$ we obtain $\hat{\boldsymbol{\alpha}} = (K + \sigma_n^2 I)^{-1} \mathbf{y}$, so that the prediction for a test point x_* is $\hat{f}(x_*) = k(x_*)^\top (K + \sigma_n^2 I)^{-1} \mathbf{y}$. This should look very familiar-it is exactly the form of the predictive mean obtained in eq. (2.23).

6.2.3 The Relationship of the Regularization View to Gaussian Process Prediction

The regularization method returns $\hat{f} = \operatorname{argmin}_f J[f]$. For a Gaussian process predictor we obtain a posterior distribution over functions.

In fact we shall see in this section that \hat{f} can be viewed as the maximum a posteriori (MAP) function under the posterior.

Following Szeliski [1987] and Poggio and Girosi [1990] we consider

$$\exp(-J[f]) = \exp\left(-\frac{\lambda}{2}\|Pf\|^2\right) \times \exp(-Q(\mathbf{y}, \mathbf{f})) \quad (6.21)$$

- The first term on the RHS is a Gaussian process prior on f
- The second is proportional to the likelihood
- As \hat{f} is the minimizer of $J[f]$, it is the MAP function.

To get some intuition for the Gaussian process prior, imagine $f(\mathbf{x})$ being represented on a grid in \mathbf{x} -space, so that f is now an (infinite dimensional) vector \mathbf{f} . Thus we obtain

$$\|Pf\|^2 \simeq \sum_{m=0}^M a_m (D_m \mathbf{f})^\top (D_m \mathbf{f}) = \mathbf{f}^\top \left(\sum_m a_m D_m^\top D_m \right) \mathbf{f}$$

where D_m is an appropriate finite-difference approximation of the differential operator O_m . Observe that this prior term is a quadratic form in \mathbf{f} .

The MAP relationship

- ❶ when $Q(y, f)$ is quadratic (corresponding to a Gaussian likelihood);
- ❷ when $Q(y, f)$ is not quadratic but convex
- ❸ when $Q(y, f)$ is not convex.

The MAP relationship

- 1 In case 1 the posterior mean function can be obtained exactly, and the posterior is Gaussian. As the mean of a Gaussian is also its mode this is the MAP solution.
- 2 In case 2 we have seen in chapter 3 for classification problems using the logistic, probit or softmax response functions that $Q(y, f)$ is convex.
- 3 In case 3 there will be more than one local minimum of $J[f]$ under the regularization approach. One could check these minima to find the deepest one. However, in this case the argument for MAP is rather weak (especially if there are multiple optima of similar depth) and suggests the need for a fully Bayesian treatment.

While the regularization solution gives a part of the Gaussian process solution, there are the following limitations:

- 1 It does not characterize the uncertainty in the predictions, nor does it handle well multimodality in the posterior.
- 2 The analysis is focussed at approximating the first level of Bayesian inference, concerning predictions for f . It is not usually extended to the next level, e.g. to the computation of the marginal likelihood. The marginal likelihood is very useful for setting any parameters of the covariance function, and for model comparison (see chapter 5).

In addition, we find the specification of smoothness via the penalties on derivatives to be not very intuitive. The regularization viewpoint can be thought of as directly specifying the inverse covariance rather than the covariance.

As marginalization is achieved for a Gaussian distribution directly from the covariance (and not the inverse covariance) it seems more natural to us to specify the covariance function. Also, while non-stationary covariance functions can be obtained from the regularization viewpoint, e.g. by replacing the Lebesgue measure in eq. (6.10) with a non-uniform measure $\mu(\mathbf{x})$, calculation of the corresponding covariance function can then be very difficult.

Section 3

6.3 Spline Models

6.3 Spline Models

In section 6.2 we discussed regularizers which had $a_0 > 0$ in eq. (6.12). We now consider the case when $a_0 = 0$; in particular we consider the regularizer to be of the form $\|O^m f\|^2$, as defined in eq. (6.10). In this case polynomials of degree up to $m - 1$ are in the null space of the regularization operator, in that they are not penalized at all.

In the case that $\mathcal{X} = \mathfrak{R}^D$ we can again use Fourier techniques to obtain the Green's function G corresponding to the Euler-Lagrange equation $(-1)^m \nabla^{2m} G(\mathbf{x}) = \delta(\mathbf{x})$. The result, as shown by Duchon [1977] and Meinguet [1979] is

$$G(\mathbf{x} - \mathbf{x}') = \begin{cases} c_{m,D} |\mathbf{X} - \mathbf{x}|^{2m-D} \log |\mathbf{x} - \mathbf{x}'| & \text{if } 2m > D \text{ and even} \\ c_{m,D} |\mathbf{X} - \mathbf{x}|^{2m-D} & \text{otherwise} \end{cases} \quad (6.22)$$

where $c_{m,D}$ is a constant (Wahba [1990, p. 31] gives the explicit form). Note that the constraint $2m > D$ has to be imposed to avoid having a Green's function that is singular at the origin.

Explicit calculation of the Green's function for other domains \mathcal{X} is sometimes possible.

Because of the null space, a minimizer of the regularization functional has the form

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i G(\mathbf{x}, \mathbf{x}_i) + \sum_{j=1}^k \beta_j h_j(\mathbf{x}) \quad (6.23)$$

where $h_1(\mathbf{x}), \dots, h_k(\mathbf{x})$ are polynomials that span the null space.

The exact values of the coefficients α and β for a specific problem can be obtained in an analogous manner to the derivation in section 6.2.2;

The solution is equivalent to that given in eq. (2.42). To gain some more insight into the form of the Green's function we consider the equation $(-1)^m \nabla^{2m} G(\mathbf{x}) = \delta(\mathbf{x})$ in Fourier space, leading to $\tilde{G}(\mathbf{s}) = (4\pi^2 \mathbf{s} \cdot \mathbf{s})^{-m}$. $\tilde{G}(\mathbf{s})$ plays a role like that of the power spectrum in eq. (6.13), but notice that $\int \tilde{G}(\mathbf{s}) d\mathbf{s}$ is infinite, which would imply that the corresponding process has infinite variance.

The problem is of course that the null space is unpenalized; for example any arbitrary constant function can be added to f without changing the regularizer.

Because of the null space we have seen that one cannot obtain a simple connection between the spline solution and a corresponding Gaussian process problem.

However, by introducing the notion of an intrinsic random function (IRF) one can define a generalized covariance;

The basic idea is to consider linear combinations of $f(\mathbf{x})$ of the form $g(\mathbf{x}) = \sum_{i=1}^k a_i f(\mathbf{x} + \boldsymbol{\delta}_i)$ for which $g(\mathbf{x})$ is second-order stationary and where $(h_j(\boldsymbol{\delta}_1), \dots, h_j(\boldsymbol{\delta}_k))\mathbf{a} = 0$ for $j = 1, \dots, k$. A careful description of the equivalence of spline and IRF prediction is given in Kent and Mardia [1994].

The power-law form of $\tilde{G}(\mathbf{s}) = (4\pi^2 \mathbf{s} \cdot \mathbf{s})^{-m}$ means that there is no characteristic length-scale for random functions drawn from this (improper) prior.

Some authors argue that the lack of a characteristic length-scale is appealing. This may sometimes be the case, but if we believe there is an appropriate length-scale (or set of length-scales) for a given problem but this is unknown in advance, we would argue that a hierarchical Bayesian formulation of the problem (as described in chapter 5) would be more appropriate.

Splines were originally introduced for one-dimensional interpolation and smoothing problems, and then generalized to the multivariate setting.

Schoenspline interpolation berg [1964] considered the problem of finding the function that minimizes $\int_a^b (f^{(m)}(\mathbf{x}))^2 dx$ (6.24), where $f^{(m)}$ denotes the m th derivative of f , subject to the interpolation constraints $f(x_i) = f_i$, $x_i \in (a, b)$ for $i = 1, \dots, n$ and for f in an appropriate natural polynomial Sobolev space. He showed that the solution is the natural polynomial spline, which is a piecewise polynomial of order $2m - 1$ in each interval $[x_i, x_{i+1}]$, $i = 1, \dots, n - 1$, and of order $m - 1$ in the two outermost intervals.

The pieces are joined so that the solution has $2m - 2$ continuous derivatives. Schoenberg also proved that the solution to the univariate smoothing problem (see eq. (6.19)) is a natural polynomial spline.

A common choice is $m = 2$, leading to the cubic spline. One possible way of writing this solution is $f(\mathbf{x}) =$

$$\sum_{j=0}^1 \beta_j x^j + \sum_{i=1}^n \alpha_i (x - x_i)_+^3, \quad (\mathbf{x})_+ = \begin{cases} (x)_+ & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases} \quad (6.25)$$

It turns out that the coefficients α and β can be computed in time $\mathcal{O}(n)$ using an algorithm due to Reinsch.

Splines were first used in regression problems. However, by using generalized linear modelling [**McCullagh and Nelder, 1983**] they can be extended to classification problems and other non-Gaussian likelihoods, as we did for GP classification in section 3.3.

* 6.3.1 A 1-d Gaussian Process Spline Construction

In this section we will further clarify the relationship between splines and Gaussian processes by giving a GP construction for the solution of the univariate cubic spline smoothing problem whose cost functional is

$$\sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \int_0^1 (f''(x))^2 dx \quad (6.26)$$

where the observed data are

$\{(x_i, y_i) | i = 1, \dots, n, 0 < x_1 < \dots < x_n < 1\}$ and λ is a smoothing parameter controlling the trade-off between the first term, the data-fit, and the second term, the regularizer, or complexity penalty.

Recall that the solution is a piece-wise polynomial as in eq. (6.25).

Following Wahba [1978], we consider the random function

$$g(x) = \sum_{j=0}^1 \beta_j x^j + f(x) \quad (6.27)$$

where $\beta_j \sim \mathcal{N}(0, \sigma^2 I)$ and $f(x)$ is a Gaussian process with covariance

To complete the analogue of the regularizer in eq. (6.26), we need to remove any penalty on polynomial terms in the null space by making the prior vague, i.e. by taking the limit $\sigma_\beta^2 \rightarrow \infty$. Notice that the covariance has the form of contributions from explicit basis functions, $\mathbf{h}(\mathbf{x}) = (1, x)^\top$ and a regular covariance function $k_{sp}(x, x')$, a problem which we have already studied in section 2.7.

Indeed we have computed the limit where the prior becomes vague $\sigma_\beta^2 \rightarrow \infty$, the result is given in eq. (2.42).

Plugging into the mean equation from eq. (2.42), we get the predictive mean $\bar{f}(x_*) = \mathbf{k}(x_*)^\top K_y^{-1}(\mathbf{y} - H^\top \bar{\beta}) + \mathbf{h}(x_*)^\top \bar{\beta}$ (6.29), where K_y is the covariance matrix corresponding to $\sigma_f^2 k_{sp}(x_i, x_j) + \sigma_n^2 \delta_{ij}$ evaluated at the training points, H is the matrix that collects the $\mathbf{h}(x_i)$ vectors at all training points, and $\bar{\beta} = (HK_y^{-1}H^\top)^{-1}HK_y^{-1}\mathbf{y}$ is given below eq. (2.42).

It is not difficult to show that this predictive mean function is a piecewise cubic polynomial, since the elements of $k(x_*)$ are piecewise cubic polynomials.

So far k_{sp} has been produced rather mysteriously “from the hat”;

Shepp [1966] defined the l -fold integrated Wiener process as

$$W_l(\mathbf{x}) = \int_0^1 \frac{(x-u)_+^l}{l!} Z(u) du, \quad l = 0, 1, \dots \quad (6.30)$$

where $Z(u)$ denotes the Gaussian white noise process with covariance $\delta(u - u')$. Note that W_0 is the standard Wiener process.

It is easy to show that $k_{sp}(x, x')$ is the covariance of the once-integrated Wiener process by writing $W_1(x)$ and $W_1(x')$ using eq. (6.30) and taking the expectation using the covariance of the white noise process.

Note that W_l is the solution to the stochastic differential equation (SDE) $X^{(l+1)} = Z$; see Appendix B for further details on SDEs.

Thus for the cubic spline we set $l = 1$ to obtain the SDE $\mathbf{x}'' = Z$, corresponding to the regularizer $\int (f''(x))^2 dx$.

Consider the family of random functions given by

$$f_N(x) = \frac{1}{\sqrt{N}} \sum_{i=0}^{N-1} \gamma_i (x - \frac{i}{N})_+ \quad (6.31)$$

where γ is a vector of parameters with $\gamma \sim \mathcal{N}(\mathbf{0}, I)$. Note that the sum has the form of evenly spaced “ramps” whose magnitudes are given by the entries in the γ vector.

$$E[f_N(x)f_N(x')] = \frac{1}{N} \sum_{i=0}^{N-1} (x - \frac{i}{N})_+ (x' - \frac{i}{N})_+ \quad (6.32)$$

Taking the limit $N \rightarrow \infty$, we obtain eq. (6.28), a derivation which is also found in [Vapnik, 1998, sec. 11.6].

Notice that the covariance function k_{sp} given in eq. (6.28) corresponds to a Gaussian process which is MS continuous but only once MS differentiable.

Thus samples from the prior will be quite “rough”, although (as noted in section 6.1) the posterior mean, eq. (6.25), is smoother.