

Outline and References

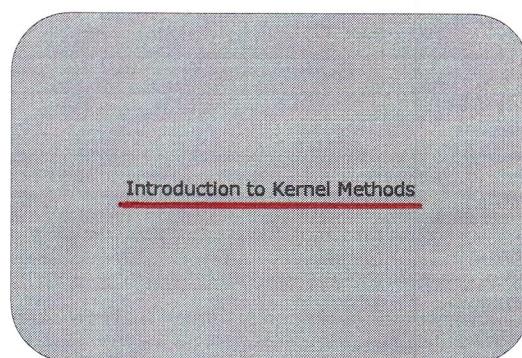
- Outline
 - ▶ Introduction [PRML 6]
 - ▶ Kernel Ridge Regression [PRML 6.1]
 - ▶ Kernel Design [PRML 6.2]
 - ▶ Kernel Regression [PRML 6.3, 2.5.1]
 - ▶ Gaussian Processes [PRML 6.4]

- References
 - ▶ These slides are based on material of prof. Marcello Restelli
 - ▶ [Pattern Recognition and Machine Learning, Bishop \[PRML\]](#)



Machine Learning - Daniele Loiacono

2

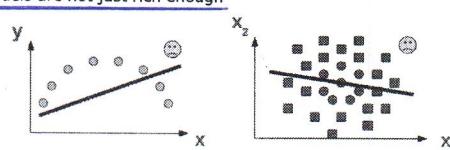


Machine Learning - Daniele Loiacono

3

Feature Mapping

- Often we want to capture nonlinear patterns in the data
 - ▶ Nonlinear Regression: input-output relationship may not be linear
 - ▶ Nonlinear Classification: Classes may not be separable by a linear boundary
- Linear models are not just rich enough

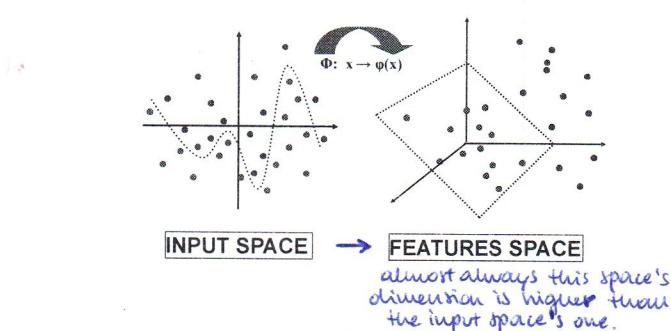


- Kernel methods allow to make linear models work in nonlinear settings by mapping data to higher dimensions where it exhibits linear patterns

Machine Learning - Daniele Loiacono

5

Feature Mapping (2)

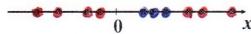


Machine Learning - Daniele Loiacono

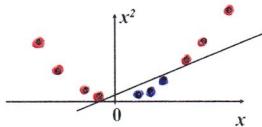
6

Example: 1D Binary Classification

- Let consider this binary classification problem for which **no linear separator** exists:

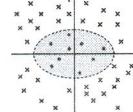


- Now let map the input space (single variable x) to a feature space with **two features**: $x \rightarrow \{x, x^2\}$
- Data is now linear separable:

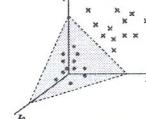


Example: 2D Binary Classification

- Let consider the following binary classification problem for which **no linear separator** exists in the input space $x = \{x_1, x_2\}$



- We can make data linearly separable by mapping the input space to a suitable feature space: $x = \{x_1, x_2\} \rightarrow z = \{x_1^2, \sqrt{2}x_1x_2, x_2^2\}$



Why kernel methods?

- Let consider a **quadratic mapping** for a problem with M input variables:

$$x = \{x_1, \dots, x_M\} \rightarrow \phi(x) = \{x_1^2, x_2^2, \dots, x_M^2, x_1x_2, x_2x_3, \dots, x_1x_M, \dots, x_{M-1}x_M\}$$

- Curse of dimensionality!** The number of features grows significantly with the number of input variable and the mapping become quickly **computationally unfeasible**

- Kernels methods** deal with this issue:
 - they **don't require to explicitly compute** the feature mapping
 - they are **expensive** but computationally **feasible**

Kernel Functions

- The **kernel function** is defined as the **scalar product** between the feature vectors of two data samples:

$$k(x, x') = \phi(x)^T \phi(x')$$

- Kernel function is **symmetric**: $k(x, x') = k(x', x)$
- Kernel function can be interpreted as a **similarity measure** between x and x'
- Very large feature vectors (even non finite ones) might result in an easy to compute kernel function
- Special class of kernels:
 - Stationary kernels**: $k(x, x') = k(x - x')$
 - Homogeneous kernels** (or **radial basis functions**): $k(x, x') = k(\|x - x'\|)$

Kernel Trick

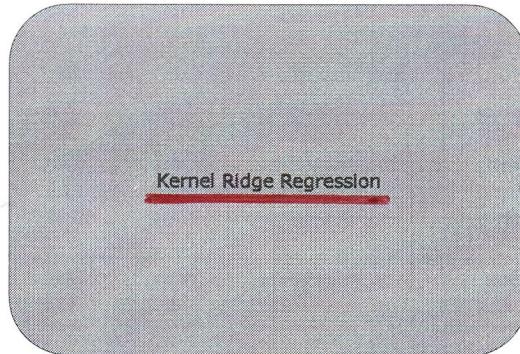
- What kernel function is used for?
 - It is possible to rework the representation of linear models to replace all the terms that involve $\phi(x)$ with other terms that involve only $k(x, \cdot)$
 - In other words, the output of linear model can be computed only on the basis of the similarities between data samples (computed with the kernel function)
- This approach, called **kernel trick**, is used in several learning algorithms
 - Ridge Regression
 - K-NN Regression
 - Perceptron
 - (Nonlinear) PCA
 - Support Vector Machines
 - ...

we may think that we solved nothing since we still have to compute $\phi(\cdot)^T \phi(\cdot)$ (and so we have to det. $\phi(\cdot)$). However, it's possible to obtain $k(\cdot, \cdot)$ even without computing the scalar product !! In fact, we'll try to construct a kernel function s.t. it respects some properties. At that point it'll be still possible to represent the kernel function as $\phi(\cdot)^T \phi(\cdot)$ but we won't be interested.

Kernel Trick

- What kernel function is used for?
 - It is possible to rework the representation of linear models to replace all the terms that involve $\phi(x)$ with other terms that involve only $k(x_i)$
 - In other words, the output of linear model can be computed only on the basis of the similarities between data samples (computed with the kernel function)
- This approach, called **kernel trick**, is used in several learning algorithms
 - Ridge Regression
 - K-NN Regression
 - Perceptron
 - (Nonlinear) PCA
 - Gaussian Processes
 - Support Vector Machines
 - ...

Machine Learning - Daniele Lolliano 12



Machine Learning - Daniele Lolliano 13

Dual Representation

- Let's go back to the loss function used for the **ridge regression**: (residual sum of squares + regulariz. term)

$$L(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \phi(\mathbf{x}_n) - t_n)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} = \frac{1}{2} (\mathbf{t} - \Phi \mathbf{w})^T (\mathbf{t} - \Phi \mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- To solve it, we set to zero the gradient of L with respect to \mathbf{w} :

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \lambda \mathbf{w} - \Phi^T (\mathbf{t} - \Phi \mathbf{w}) = 0$$

$$\Phi = \begin{bmatrix} \Phi(\mathbf{x}_1) \\ \vdots \\ \Phi(\mathbf{x}_n) \end{bmatrix} : \text{we're not able to compute it}$$

- Now, instead of solving it for \mathbf{w} , let do a variable change:

$$\mathbf{w} = \Phi^T \lambda^{-1} (\mathbf{t} - \Phi \mathbf{w}) = \Phi^T \mathbf{a}$$

$$\mathbf{a} = \lambda^{-1} (\mathbf{t} - \Phi \mathbf{w})$$

Machine Learning - Daniele Lolliano 14

Dual Representation (2)

- Now we can replace \mathbf{w} in the gradient:

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \lambda \mathbf{w} - \Phi^T (\mathbf{t} - \Phi \mathbf{w}) = 0$$

$$\rightarrow \Phi^T (\lambda \mathbf{a} - (\mathbf{t} - \Phi \Phi^T \mathbf{a})) = 0$$

$$\rightarrow \Phi \Phi^T \mathbf{a} + \lambda \mathbf{a} = \mathbf{t}$$

$$\rightarrow \mathbf{a} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{t}$$

$$\mathbf{K} = \Phi \Phi^T$$

Gram Matrix

$$\Phi \in \mathbb{R}^{N \times M}, \quad \mathbf{K} \in \mathbb{R}^{N \times N}$$

Machine Learning - Daniele Lolliano 15

Gram Matrix and Kernel Function

- The **Gram matrix** is a $N \times N$ matrix, where each element is the inner product between the feature vectors:

$$K = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & \dots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

$$k(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$$

- K matrix represents the **similarities** between each pair of samples in the **training data**

Machine Learning - Daniele Lolliano 16

Prediction Function

- How can we compute the prediction using the dual representation?

$$\begin{aligned} \mathbf{w} &= \Phi^T \mathbf{a} \\ y(\mathbf{x}) &= \boxed{\mathbf{w}^T \phi(\mathbf{x})} = \boxed{\mathbf{a}^T \Phi \phi(\mathbf{x})} = \mathbf{k}(\mathbf{x})^T (K + \lambda \mathbf{I}_N)^{-1} \mathbf{t} \\ \mathbf{a} &= (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{t} \end{aligned}$$

$\Phi = \begin{bmatrix} \underline{\Phi}(\underline{\mathbf{x}}_1) \\ \vdots \\ \underline{\Phi}(\underline{\mathbf{x}}_n) \end{bmatrix} \Rightarrow \underline{\Phi} \cdot \underline{\Phi}(\underline{\mathbf{x}}) = \begin{bmatrix} \underline{\Phi}(\underline{\mathbf{x}}_1) \\ \vdots \\ \underline{\Phi}(\underline{\mathbf{x}}_n) \end{bmatrix} \cdot \underline{\Phi}(\underline{\mathbf{x}}) = \begin{bmatrix} \underline{\Phi}(\underline{\mathbf{x}})^T \cdot \underline{\Phi}(\underline{\mathbf{x}}) \\ \vdots \\ \underline{\Phi}(\underline{\mathbf{x}}_n)^T \cdot \underline{\Phi}(\underline{\mathbf{x}}) \end{bmatrix} = \begin{bmatrix} k_n(\underline{\mathbf{x}}_1, \underline{\mathbf{x}}) \\ \vdots \\ k_n(\underline{\mathbf{x}}_n, \underline{\mathbf{x}}) \end{bmatrix} = \underline{\mathbf{k}}(\underline{\mathbf{x}})$

Machine Learning - Daniele Lolli

17

Prediction Function

- How can we compute the prediction using the dual representation?

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) = \mathbf{a}^T \Phi \phi(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T (K + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$$

► Where $\mathbf{k}(\mathbf{x})$ is such that $k_n(\mathbf{x}) = k(\mathbf{x}_n, \mathbf{x}) \forall \mathbf{x}_n \in \mathcal{D}$

- Accordingly the prediction is computed as the linear combination of the target values of the samples in the training set

Machine Learning - Daniele Lolli

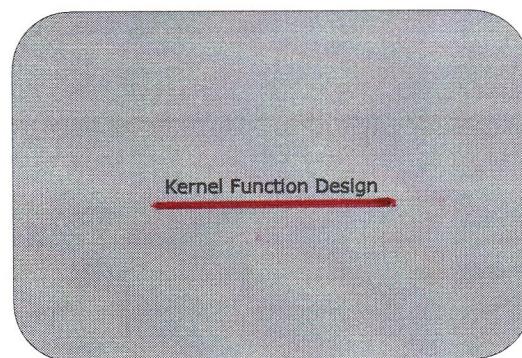
18

Original vs Dual Representation?

- Original representation
 - Requires to compute the inverse of $(\Phi^T \Phi + \lambda I_M)$ which is a $M \times M$ matrix
 - Is computationally convenient when M is rather small
- Dual representation
 - Requires to compute the inverse of $(K + \lambda I_N)$ which is a $N \times N$ matrix
 - Is computationally convenient when M is very large or even infinite
 - Does not require to explicitly compute Φ , making it possible to apply this approach also to complex type of data (e.g., graphs, sets, strings, text, etc.)
 - The similarity between data samples (i.e., the kernel function) is generally both less expensive to compute and easy to design than computing Φ

Machine Learning - Daniele Lolli

19



Machine Learning - Daniele Lolli

20

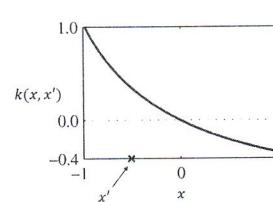
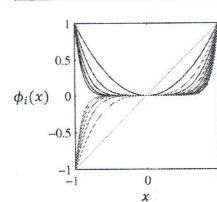
Design of Kernels

- We defined the kernel function as the dot product in the feature space:

$$k(x, x') = \phi(x)^T \phi(x') = \sum_{i=1}^M \phi_i(x) \phi_i(x')$$

- Examples (1D input space)

If we restrict to this then it's all useless (because we're passing through $\Phi(\cdot)$ to calculate $k(\cdot, \cdot)$)



Machine Learning - Daniele Lolli

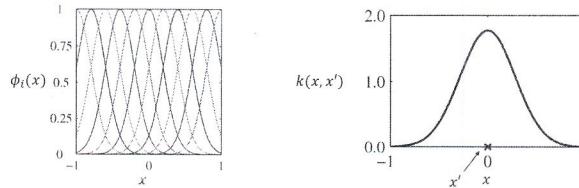
21

Design of Kernels

- We defined the kernel function as the dot product in the feature space:

$$k(x, x') = \phi(x)^T \phi(x') = \sum_{i=1}^M \phi_i(x) \phi_i(x')$$

- Examples (1D input space)

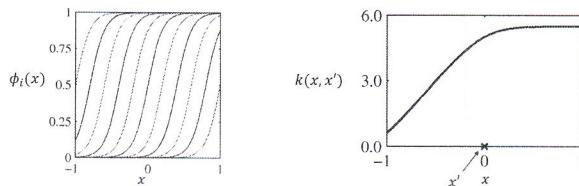


Design of Kernels

- We defined the kernel function as the dot product in the feature space:

$$k(x, x') = \phi(x)^T \phi(x') = \sum_{i=1}^M \phi_i(x) \phi_i(x')$$

- Examples (1D input space)



Design of Kernels: alternative methods

- We do not necessarily have to compute the kernel function starting from the feature space, as we do not want to explicitly compute the feature vectors
- There are two major alternatives to design a kernel function:
 - Directly design the kernel functions from scratch
 - Design from existing kernel functions by applying a set of rules
- In general, we must make sure that the designed kernel functions is valid, that is it correspond to a scalar product into any feature space
- Mercer Theorem:** Any continuous, symmetric, positive semi-definite kernel function $k(x, x')$ can be expressed as a dot product in a high-dimensional space
 - Necessary and sufficient condition for a function $k(x, x')$ to be a valid kernel is that Gram matrix K is positive semi-definite for all possible choice of $\mathcal{D} = \{x_i\}$
 - It means $x^T K x > 0$ for any non zero real vector x , i.e., $\sum_j K_{ij} x_i x_j$ for any real numbers x_i and x_j

Kernel direct design: an example

- Let consider the following kernel function: $k(x, x') = (x^T x')^2$
- Is it a valid kernel function?
- Let check it in a two dimensional space:

$$\begin{aligned} k(x, x') &= (x^T x')^2 = (x_1 x'_1 + x_2 x'_2)^2 = x_1^2 x'^2_1 + 2x_1 x'_1 x_2 x'_2 + x_2^2 x'^2_2 \\ &= (x_1^2, \sqrt{2} x_1 x_2, x_2^2) (x'^2_1, \sqrt{2} x'_1 x'_2, x'^2_2)^T = \phi(x)^T \phi(x') \end{aligned}$$

- It does correspond to the scalar product in a feature space with only second order terms (also notice that is less computationally expensive to compute)
- To get also constant and linear terms, we can define $k(x, x') = (x^T x' + c)^2$
- To get all the terms up to degree z , we can define $k(x, x') = (x^T x' + c)^z$

Rules to design valid kernels

- Given valid kernels $k_1(x, x')$ and $k_2(x, x')$ the following rules can be applied to design a new valid kernel:
 - $k(x, x') = ck_1(x, x')$, where $c > 0$ is a constant
 - $k(x, x') = f(x)k_1(x, x')f(x')$, where $f(\cdot)$ is any function
 - $k(x, x') = q(k_1(x, x'))$, where $q(\cdot)$ is a polynomial with non-negative coefficients
 - $k(x, x') = \exp(k_1(x, x'))$
 - $k(x, x') = k_1(x, x') + k_2(x, x')$
 - $k(x, x') = k_1(x, x')k_2(x, x')$
 - $k(x, x') = k_3(\phi(x), \phi(x'))$, where $\phi(x)$ maps x to \mathbb{R}^M and $k_3(\cdot, \cdot)$ is a valid kernel in \mathbb{R}^M
 - $k(x, x') = x^T A x'$, where A is a symmetric semidefinite matrix
 - $k(x, x') = k_a(x_a, x'_a) + k_b(x_b, x'_b)$
 - $k(x, x') = k_a(x_a, x'_a)k_b(x_b, x'_b)$

Where $x = [x_a] \cup [x_b]$ are two subsets no necessarily disjoints of variables and k_a, k_b are valid kernels

Gaussian Kernel

- This is a commonly used kernel:

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / 2\sigma^2)$$

- We can check it is a valid kernel, expanding the square:

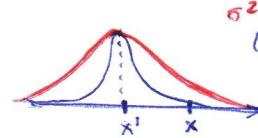
$$\|\mathbf{x} - \mathbf{x}'\|^2 = \mathbf{x}^T \mathbf{x} + \mathbf{x}'^T \mathbf{x}' - 2\mathbf{x}^T \mathbf{x}'$$

$$\Rightarrow k(\mathbf{x}, \mathbf{x}') = \exp(-\mathbf{x}^T \mathbf{x} / 2\sigma^2) \exp(-\mathbf{x}'^T \mathbf{x}' / 2\sigma^2) \exp(-\mathbf{x}^T \mathbf{x}' / 2\sigma^2) \quad \text{Rule 2 and 4}$$

- The feature space corresponding to Gaussian kernel has infinite dimension

- We can extend Gaussian Kernel by replacing $\mathbf{x}'^T \mathbf{x}'$ with a nonlinear kernel $\kappa(\mathbf{x}, \mathbf{x}')$:

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2\sigma^2}(\kappa(\mathbf{x}, \mathbf{x}) + \kappa(\mathbf{x}', \mathbf{x}') - 2\kappa(\mathbf{x}, \mathbf{x}'))\right)$$



$\sigma^2 > \sigma'$

(Note: σ^2 is also called BANDWIDTH)

this kernel corresponds to the dot product of two infinite-dim vectors (dim(features space) = ∞)
Basically kernels allow us to avoid to deal with feature spaces.

The original problem of choosing the right feature space now becomes the problem of choosing the right kernel

Kernels for Symbolic Data

- Kernel methods can be extended also to inputs different from real vectors, such as graphs, sets, strings, texts, etc.
- In fact, the kernel function represents a measure of the similarity between two samples
- A common kernel used for set is:

$$k(A_1, A_2) = \frac{2|A_1 \cap A_2|}{|A_1| + |A_2|} \quad \text{number of elements that appear in both sets}$$

Kernels Based on Generative Models

- It is also possible to define a kernel function based on probability distribution
- Given a generative model $p(\mathbf{x})$ we can define a kernel as:

$$k(\mathbf{x}, \mathbf{x}') = p(\mathbf{x})p(\mathbf{x}')$$

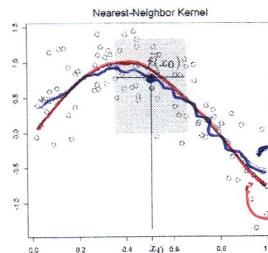
- It is a valid kernel, since it corresponds to the inner product in the one dimensional feature space defined mapping \mathbf{x} to $p(\mathbf{x})$

Kernel Regression

k-NN Regression

- The k Nearest Neighbour (k-NN) can be applied to solve a regression problem by averaging the K nearest samples in the training data:

$$\hat{f}(\mathbf{x}) = \frac{1}{k} \sum_{\mathbf{x}_i \in N_k(\mathbf{x})} t_i$$



obtained regression

true function

What is a problem? There is a lot of noise; moving a little bit makes the function change a lot

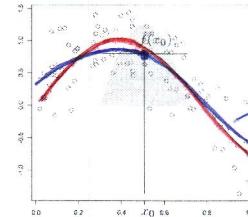
To mitigate it we can introduce a kernel function

Nadaraya-Watson Model

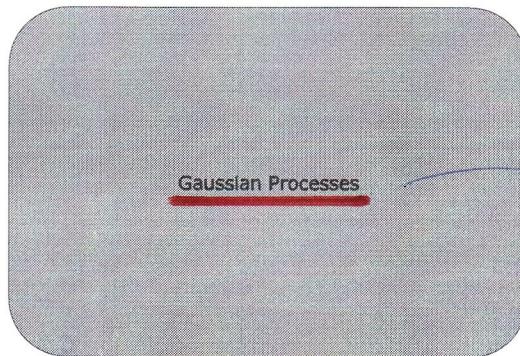
- In k-NN regression the model output is generally very noisy due to the discontinuity of neighborhood averages
- Nadaraya-Watson model (or kernel regression) deal with this issue by using kernel function to compute a weighted average of samples:

$$\hat{f}(\mathbf{x}) = \frac{\sum_{i=1}^N k(\mathbf{x}, \mathbf{x}_i) t_i}{\sum_{i=1}^N k(\mathbf{x}, \mathbf{x}_i)}$$

- Typical choices for kernels are:
 - Epanechnikov Kernel (bounded support)
 - Gaussian Kernel (infinite support)



this time the obtained regression function is more smooth



Linear Regression Revisited

- Let's start from the same assumptions used in Bayesian Linear Regression

$$y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x})$$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | 0, \tau^2 \mathbf{I})$$

linear model with a prior distribution over the weights

- Now compute the prior distribution of the outputs of the regression function:

$$\mathbf{y} = \Phi \mathbf{w} \rightarrow p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mu, \mathbf{S})$$

$$\mu = \mathbb{E}[\mathbf{y}] = \Phi \mathbb{E}[\mathbf{w}] = \mathbf{0}$$

$$\mathbf{S} = \text{cov}[\mathbf{y}] = \mathbb{E}[\mathbf{y}\mathbf{y}^T] = \Phi \mathbb{E}[\mathbf{w}\mathbf{w}^T] \Phi^T = \tau^2 \Phi \Phi^T = \mathbf{K}$$

Gram matrix

Now, instead of modelling the distribution of the weights, we try to model the outputs of the regression function (the "y's"). (To do it we start from the assumption that we did on the weights: $\mathbf{w} \sim \mathcal{N}(0, \tau^2 \mathbf{I})$)

Note that we're including τ^2

Gaussian Processes and Gram Matrix

- In general, a Gaussian Process is defined as a distribution probability over a function $y(\mathbf{x})$ such that the set of values $y(\mathbf{x}_i)$ – for an arbitrary $\{\mathbf{x}_i\}$ – jointly have a Gaussian Distribution.

- In our specific case,

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | 0, \mathbf{K})$$

- where \mathbf{K} is the Gram matrix defined as:

$$K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m) = \tau^2 \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$$

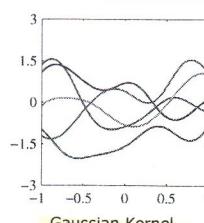
- This gives a probabilistic interpretation of the Kernel function as:

$$k(\mathbf{x}_n, \mathbf{x}_m) = \mathbb{E}[y(\mathbf{x}_n) y(\mathbf{x}_m)]$$

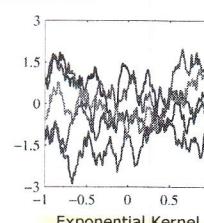
the kernel function computed between two samples represents how much the output of the model will be correlated in the two samples. This is why the kernel function can be interpreted as similarity. The kernel shapes how correlated will be the output of the model for each pair of points.

Kernel Design

- We can apply the usual approaches to design the kernels
- Two families of kernels typically used with GP are:



$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_2^2 / 2\sigma^2)$$



$$k(\mathbf{x}, \mathbf{x}') = \exp(-\theta |\mathbf{x} - \mathbf{x}'|)$$

longer σ^2
smoother the function

- To apply GP to regression we must take into account also noise on the target:

$$t_n = y(\mathbf{x}_n) + \epsilon_n$$

- Assuming the noise Gaussian and independent on each point, the joint distribution of the targets will be:

$$p(\mathbf{t}|\mathbf{y}) = \mathcal{N}(\mathbf{t}|\mathbf{y}, \sigma^2 \mathbf{I})$$

- We can compute the marginal distribution $p(\mathbf{t})$ from $p(\mathbf{t}|\mathbf{y})$ and $p(\mathbf{y})$ as:

$$p(\mathbf{t}) = \int p(\mathbf{t}|\mathbf{y})p(\mathbf{y})d\mathbf{y}$$

$$\mathcal{N}(\mathbf{t}|\mathbf{y}, \sigma^2 \mathbf{I}) \quad \mathcal{N}(\mathbf{y}|0, \mathbf{K})$$

early since both are gaussian.
The result will be a gaussian with
the mean which is the product of
the two means and the variance
which is the sum of the two variances

Gaussian Processes for Regression

- To apply GP to regression we must take into account also noise on the target:

$$t_n = y(\mathbf{x}_n) + \epsilon_n$$

- Assuming the noise Gaussian and independent on each point, the joint distribution of the targets will be:

$$p(\mathbf{t}|\mathbf{y}) = \mathcal{N}(\mathbf{t}|\mathbf{y}, \sigma^2 \mathbf{I})$$

- We can compute the marginal distribution $p(\mathbf{t})$ from $p(\mathbf{t}|\mathbf{y})$ and $p(\mathbf{y})$ as:

$$p(\mathbf{t}) = \int p(\mathbf{t}|\mathbf{y})p(\mathbf{y})d\mathbf{y} = \mathcal{N}(\mathbf{t}|0, \mathbf{C})$$

► Being $p(\mathbf{t}|\mathbf{y})$ and $p(\mathbf{y})$ their covariances simply add:

$$C(\mathbf{x}_n, \mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) + \sigma^2 \delta_{nm}$$

$$\mathbf{C} = \mathbf{K} + \sigma^2 \mathbf{I}$$

We now want to use this
distribution to make prediction

Making Prediction with GP

- We observe $\{\mathbf{x}, \mathbf{t}\}$ and want to make a prediction on a new point \mathbf{x}_{N+1}

- First, we need to compute the joint distribution of \mathbf{t}_{N+1} (that will be still Gaussian):

$$p(\mathbf{t}_{N+1}) = \mathcal{N}(\mathbf{t}_{N+1}|0, \mathbf{C}_{N+1})$$

$$\mathbf{t}_{N+1} = (t_1, t_2, \dots, t_{N+1})^T$$

$$\mathbf{C}_{N+1} = \begin{pmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^T & c \end{pmatrix}$$

► \mathbf{k} is a vector $k(\mathbf{x}_i, \mathbf{x}_{N+1}) \quad \forall i = 1, \dots, N$

$$c = k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) + \sigma^2$$

target that we want to predict (t_{N+1})

Making Prediction with GP (2)

- ... now we can use the joint distribution to derive the conditioned distribution that will be used in practice to make the prediction:

$$p(t_{N+1}|\mathbf{t}) = \mathcal{N}(m, \sigma^2)$$

► where:

$$m(\mathbf{x}_{N+1}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t}$$

$$\sigma^2(\mathbf{x}_{N+1}) = c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k} \quad \rightarrow \text{uncertainty of our prediction}$$

- Computational cost is:

- $O(N^3)$ for inversion of \mathbf{C} (required only once for training)
- $O(N)$ for computing m
- $O(N^2)$ for computing σ^2

Kernel design for regression

- Popular choice is $k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left(-\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|_2^2\right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$

- How to choose the kernel hyperparameters?

► Grid Search (very expensive)

► Maximization of marginal likelihood over the training set using gradient optimization

