


The Kernel Trick

Yingzhen Li

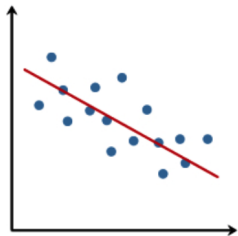
Department of Computing
Imperial College London

@liyzhen2
yingzhen.li@imperial.ac.uk

November 26, 2021

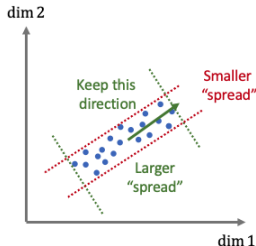
Linear models in ML

With linear features



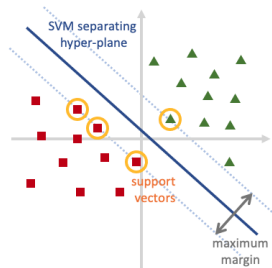
Linear regression

$$f(x, \theta) = \theta^\top x$$



PCA

$$\max \mathbb{V}[\mathbf{b}^\top \mathbf{x}]$$

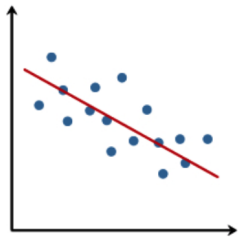


SVM

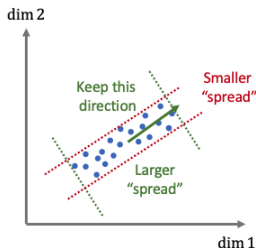
$$f(x) = \text{sgn}(\mathbf{w}^\top \mathbf{x} + b)$$

Linear models in ML

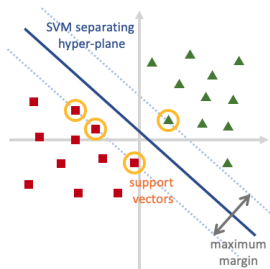
With **non-linear** features



Linear regression
 $f(x, \theta) = \theta^\top \phi(x)$



PCA
 $\max \mathbb{V}[\mathbf{b}^\top \phi(x)]$



SVM
 $f(x) = \text{sgn}(\mathbf{w}^\top \phi(x) + b)$

How to choose feature mapping $\phi(x)$?

Choice of feature map

The choice of feature map $\phi(\cdot)$: crucial for performance:

- Want a flexible feature map
 - example: sine waves with different frequency and phase
 - often achieved by using large $p > \dim(\mathbf{x})$

Choice of feature map

The choice of feature map $\phi(\cdot)$: crucial for performance:

- ▶ Want a flexible feature map
 - ▶ example: sine waves with different frequency and phase
 - ▶ often achieved by using large $p > \dim(\mathbf{x})$
- ▶ Can we make $p \rightarrow +\infty$?
- ▶ What if the output domain of $\phi(\cdot)$ is not even \mathbb{R}^p ?

Choice of feature map

The choice of feature map $\phi(\cdot)$: crucial for performance:

- ▶ Want a flexible feature map
 - ▶ example: sine waves with different frequency and phase
 - ▶ often achieved by using large $p > \dim(\mathbf{x})$
- ▶ Can we make $p \rightarrow +\infty$?
- ▶ What if the output domain of $\phi(\cdot)$ is not even \mathbb{R}^p ?

Yes! Using the **kernel trick**:

compute $\phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$ directly without explicitly evaluating $\phi(\mathbf{x}_i), \phi(\mathbf{x}_j)$

Solution of non-linear ridge regression

Fitting non-linear ridge regression models:

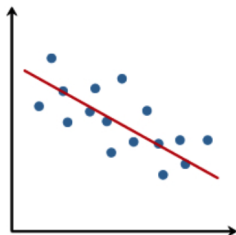
- ▶ Dataset: $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$,
 $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D}$,
 $\mathbf{y} = [y_1, \dots, y_N]^\top \in \mathbb{R}^{N \times 1}$
- ▶ The non-linear regression model:

$$f(\mathbf{x}, \boldsymbol{\theta}) = \boldsymbol{\phi}(\mathbf{x})^\top \boldsymbol{\theta}, \quad \boldsymbol{\phi}(\cdot) : \mathbb{R}^D \rightarrow \mathbb{R}^p$$

$$y = f(\mathbf{x}, \boldsymbol{\theta}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

- ▶ Solution: $\boldsymbol{\Phi} = [\boldsymbol{\phi}(\mathbf{x}_1), \dots, \boldsymbol{\phi}(\mathbf{x}_N)]^\top \in \mathbb{R}^{N \times p}$

$$\boldsymbol{\theta}^* = (\sigma^2 \lambda \mathbf{I} + \boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \mathbf{y}$$



$$f(\mathbf{x}, \boldsymbol{\theta}) = \boldsymbol{\phi}(\mathbf{x})^\top \boldsymbol{\theta}$$

The kernel trick

Solution of non-linear ridge regression:

$$\boldsymbol{\theta}^* = (\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}$$

Constructing predictions for new inputs $\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_M$ with features

$\hat{\Phi} = [\phi(\hat{\mathbf{x}}_1), \dots, \phi(\hat{\mathbf{x}}_M)]^\top$:

$$\hat{\mathbf{y}} = \hat{\Phi} \boldsymbol{\theta}^* = \hat{\Phi} (\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}$$

The kernel trick

Solution of non-linear ridge regression:

$$\boldsymbol{\theta}^* = (\sigma^2 \lambda \mathbf{I} + \boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \mathbf{y}$$

Constructing predictions for new inputs $\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_M$ with features $\hat{\boldsymbol{\Phi}} = [\phi(\hat{\mathbf{x}}_1), \dots, \phi(\hat{\mathbf{x}}_M)]^\top$:

$$\hat{\mathbf{y}} = \hat{\boldsymbol{\Phi}} \boldsymbol{\theta}^* = \hat{\boldsymbol{\Phi}} (\sigma^2 \lambda \mathbf{I} + \boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \mathbf{y}$$

► Notice the following identities:

$$(\sigma^2 \lambda \mathbf{I} + \boldsymbol{\Phi}^\top \boldsymbol{\Phi}) \boldsymbol{\Phi}^\top = \sigma^2 \lambda \boldsymbol{\Phi}^\top + \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \boldsymbol{\Phi}^\top = \boldsymbol{\Phi}^\top (\sigma^2 \lambda \mathbf{I} + \boldsymbol{\Phi} \boldsymbol{\Phi}^\top)$$

The kernel trick

Solution of non-linear ridge regression:

$$\boldsymbol{\theta}^* = (\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}$$

Constructing predictions for new inputs $\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_M$ with features

$\hat{\Phi} = [\phi(\hat{\mathbf{x}}_1), \dots, \phi(\hat{\mathbf{x}}_M)]^\top$:

$$\hat{\mathbf{y}} = \hat{\Phi} \boldsymbol{\theta}^* = \hat{\Phi} (\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}$$

► Notice the following identities:

$$(\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi) \Phi^\top = \sigma^2 \lambda \Phi^\top + \Phi^\top \Phi \Phi^\top = \Phi^\top (\sigma^2 \lambda \mathbf{I} + \Phi \Phi^\top)$$

$$\Rightarrow \Phi^\top (\sigma^2 \lambda \mathbf{I} + \Phi \Phi^\top)^{-1} = (\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top$$

The kernel trick

Solution of non-linear ridge regression:

$$\boldsymbol{\theta}^* = (\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}$$

Constructing predictions for new inputs $\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_M$ with features

$$\hat{\Phi} = [\phi(\hat{\mathbf{x}}_1), \dots, \phi(\hat{\mathbf{x}}_M)]^\top:$$

$$\hat{\mathbf{y}} = \hat{\Phi} \boldsymbol{\theta}^* = \hat{\Phi} (\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}$$

► Notice the following identities:

$$(\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi) \Phi^\top = \sigma^2 \lambda \Phi^\top + \Phi^\top \Phi \Phi^\top = \Phi^\top (\sigma^2 \lambda \mathbf{I} + \Phi \Phi^\top)$$

$$\Rightarrow \Phi^\top (\sigma^2 \lambda \mathbf{I} + \Phi \Phi^\top)^{-1} = (\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top$$

$$\Rightarrow \hat{\Phi} \Phi^\top (\sigma^2 \lambda \mathbf{I} + \Phi \Phi^\top)^{-1} \mathbf{y} = \hat{\Phi} (\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y} = \hat{\Phi} \boldsymbol{\theta}^*$$

The kernel trick

Constructing predictions for new inputs $\hat{\mathbf{X}}$ with features $\hat{\Phi}$:

$$\text{Parametric regression: } \hat{\mathbf{y}} = \hat{\Phi} \boldsymbol{\theta}^* = \hat{\Phi} (\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}$$

$$\Leftrightarrow \text{Kernel regression: } \hat{\mathbf{y}} = \hat{\Phi} \Phi^\top (\sigma^2 \lambda \mathbf{I} + \Phi \Phi^\top)^{-1} \mathbf{y}$$

The kernel trick

Constructing predictions for new inputs $\hat{\mathbf{X}}$ with features $\hat{\Phi}$:

$$\text{Parametric regression: } \hat{\mathbf{y}} = \hat{\Phi}\boldsymbol{\theta}^* = \hat{\Phi}(\sigma^2\lambda\mathbf{I} + \Phi^\top\Phi)^{-1}\Phi^\top\mathbf{y}$$

$$\Leftrightarrow \text{Kernel regression: } \hat{\mathbf{y}} = \hat{\Phi}\Phi^\top(\sigma^2\lambda\mathbf{I} + \Phi\Phi^\top)^{-1}\mathbf{y}$$

Advantages of the kernel regression view:

- ▶ Efficient in $N \ll p$ scheme:
 - ▶ $\Phi^\top\Phi \in \mathbb{R}^{p \times p} \Rightarrow$ matrix inversion needs $\mathcal{O}(p^3)$ time
 - ▶ $\Phi\Phi^\top \in \mathbb{R}^{N \times N} \Rightarrow$ matrix inversion needs $\mathcal{O}(N^3)$ time

The kernel trick

Constructing predictions for new inputs $\hat{\mathbf{X}}$ with features $\hat{\Phi}$:

$$\text{Parametric regression: } \hat{\mathbf{y}} = \hat{\Phi} \boldsymbol{\theta}^* = \hat{\Phi} (\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}$$

$$\Leftrightarrow \text{Kernel regression: } \hat{\mathbf{y}} = \mathbf{K}_{\hat{\mathbf{X}}\mathbf{X}} (\sigma^2 \lambda \mathbf{I} + \mathbf{K}_{\mathbf{X}\mathbf{X}})^{-1} \mathbf{y}$$

Advantages of the kernel regression view:

- ▶ No need to explicitly define $\phi(\cdot)$:
 - ▶ The entries of $\mathbf{K}_{\mathbf{X}\mathbf{X}} := \Phi \Phi^\top$: $\mathbf{K}_{\mathbf{X}\mathbf{X}}[i, j] := \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) := \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$
 - ▶ The entries of $\mathbf{K}_{\hat{\mathbf{X}}\mathbf{X}} := \hat{\Phi} \Phi^\top$: $\mathbf{K}_{\hat{\mathbf{X}}\mathbf{X}}[k, l] := \phi(\hat{\mathbf{x}}_k)^\top \phi(\mathbf{x}_l) := \mathcal{K}(\hat{\mathbf{x}}_k, \mathbf{x}_l)$
 - ▶ Enables the usage of $p = \infty$ -dim features
(or even non- \mathbb{R}^p features)

The kernel trick

Constructing predictions for new inputs $\hat{\mathbf{X}}$ with features $\hat{\Phi}$:

$$\text{Parametric regression: } \hat{\mathbf{y}} = \hat{\Phi} \boldsymbol{\theta}^* = \hat{\Phi} (\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}$$

$$\Leftrightarrow \text{Kernel regression: } \hat{\mathbf{y}} = \mathbf{K}_{\hat{\mathbf{X}}\mathbf{X}} (\sigma^2 \lambda \mathbf{I} + \mathbf{K}_{\mathbf{X}\mathbf{X}})^{-1} \mathbf{y}$$

Advantages of the parametric regression view:

- ▶ Easy to design $\phi(\cdot)$ with prior knowledge
- ▶ Better interpretability for non-linear features $\phi(\cdot)$
- ▶ (Mark will argue both points for the kernel view as well :)

Solution of non-linear PCA

Fitting non-linear PCA:

- ▶ Dataset: $\mathcal{D} = \mathbf{X}$,
 $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D}$
- ▶ The non-linear PCA model:

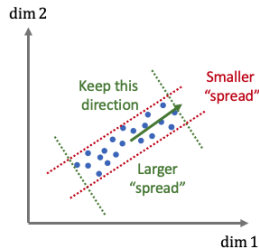
$$\mathbf{z}_n := \mathbf{B}^\top \phi(\mathbf{x}_n), \quad \phi(\cdot) : \mathbb{R}^D \rightarrow \mathbb{R}^p$$

$$\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{p \times M}$$

is an orthonormal basis

- ▶ Solution: the leading M eigenvector of

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top$$



$$\max \mathbb{V}[\mathbf{b}^\top \phi(\mathbf{x})]$$

The kernel trick

Solving the following eigenvector computation problem:

$$\lambda_m \mathbf{b}_m = \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top \mathbf{b}_m$$

The kernel trick

Solving the following eigenvector computation problem:

$$\lambda_m \mathbf{b}_m = \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top \mathbf{b}_m$$

Define $\alpha_{mn} = \frac{1}{N\lambda_m} \phi(\mathbf{x}_n)^\top \mathbf{b}_m$:

$$\mathbf{b}_m = \sum_{n=1}^N \alpha_{mn} \phi(\mathbf{x}_n)$$

The kernel trick

Solving the following eigenvector computation problem:

$$\lambda_m \mathbf{b}_m = \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top \mathbf{b}_m$$

Define $\alpha_{mn} = \frac{1}{N\lambda_m} \phi(\mathbf{x}_n)^\top \mathbf{b}_m$:

$$\mathbf{b}_m = \sum_{n=1}^N \alpha_{mn} \phi(\mathbf{x}_n)$$

The objective of PCA for \mathbf{b}_m :

$$\mathbf{b}_m^\top \mathbf{S} \mathbf{b}_m = \mathbf{b}_m^\top \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top \mathbf{b}_m = \frac{1}{N} \sum_{n=1}^N \|\phi(\mathbf{x}_n)^\top \mathbf{b}_m\|_2^2$$

The kernel trick

Solving the following eigenvector computation problem:

$$\lambda_m \mathbf{b}_m = \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top \mathbf{b}_m$$

Define $\alpha_{mn} = \frac{1}{N\lambda_m} \phi(\mathbf{x}_n)^\top \mathbf{b}_m$:

$$\mathbf{b}_m = \sum_{n=1}^N \alpha_{mn} \phi(\mathbf{x}_n)$$

Plugging-in the above reparam. of \mathbf{b}_m :

$$\phi(\mathbf{x}_n)^\top \mathbf{b}_m = \sum_{n'=1}^N \alpha_{mn'} \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_{n'}) = \sum_{n'=1}^N \alpha_{mn'} \mathcal{K}(\mathbf{x}_n, \mathbf{x}_{n'})$$

\Rightarrow we can write $\mathbf{b}_m^\top \mathbf{S} \mathbf{b}_m$ as a function of $\{\alpha_{mn}\}$ and $\{\mathcal{K}(\mathbf{x}_n, \mathbf{x}_{n'})\}$

The kernel trick

Solving the following eigenvector computation problem:

$$\lambda_m \mathbf{b}_m = \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top \mathbf{b}_m$$

Define $\alpha_{mn} = \frac{1}{N\lambda_m} \phi(\mathbf{x}_n)^\top \mathbf{b}_m$:

$$\mathbf{b}_m = \sum_{n=1}^N \alpha_{mn} \phi(\mathbf{x}_n)$$

Constraints for \mathbf{b}_m to satisfy:

$$\|\mathbf{b}_m\|_2^2 = 1 \quad \Rightarrow \quad \sum_{n,n'} \alpha_{mn} \alpha_{mn'} \underbrace{\phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_{n'})}_{=\mathcal{K}(\mathbf{x}_n, \mathbf{x}_{n'})} = 1$$

$$\mathbf{b}_i^\top \mathbf{b}_j = 0, i \neq j \quad \Rightarrow \quad \sum_{n,n'} \alpha_{in} \alpha_{jn'} \underbrace{\phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_{n'})}_{=\mathcal{K}(\mathbf{x}_n, \mathbf{x}_{n'})} = 0$$

The kernel trick

Solving the following eigenvector computation problem:

$$\lambda_m \mathbf{b}_m = \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top \mathbf{b}_m$$

Define $\alpha_{mn} = \frac{1}{N\lambda_m} \phi(\mathbf{x}_n)^\top \mathbf{b}_m$:

$$\mathbf{b}_m = \sum_{n=1}^N \alpha_{mn} \phi(\mathbf{x}_n)$$

Rewrite the objective & constraints using $\{\alpha_{mn}\}$ and $\{\mathcal{K}(\mathbf{x}_n, \mathbf{x}_{n'})\}$
 \Rightarrow solve the constrained opt. for $\{\alpha_{mn}\}$ which implicitly define \mathbf{b}_m^*

The kernel trick

Solving the following eigenvector computation problem:

$$\lambda_m \mathbf{b}_m = \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top \mathbf{b}_m$$

Define $\alpha_{mn} = \frac{1}{N\lambda_m} \phi(\mathbf{x}_n)^\top \mathbf{b}_m$:

$$\mathbf{b}_m = \sum_{n=1}^N \alpha_{mn} \phi(\mathbf{x}_n)$$

Rewrite the objective & constraints using $\{\alpha_{mn}\}$ and $\{\mathcal{K}(\mathbf{x}_n, \mathbf{x}_{n'})\}$
 \Rightarrow solve the constrained opt. for $\{\alpha_{mn}\}$ which implicitly define \mathbf{b}_m^*

Prediction:

$$\hat{z}_m := \mathbf{b}_m^\top \phi(\hat{\mathbf{x}}) = \sum_{n=1}^N \alpha_{mn} \underbrace{\phi(\mathbf{x}_n)^\top \phi(\hat{\mathbf{x}})}_{=\mathcal{K}(\mathbf{x}_n, \hat{\mathbf{x}})}$$

Constructing kernels

What is a kernel function?

- ▶ In previous non-linear regression example:

$$\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \lambda^{-1} \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$$

- ▶ Can also directly define $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$ without explicitly define $\phi(\mathbf{x})$:
 - ▶ Non-negative: $\mathcal{K}(\mathbf{x}, \mathbf{x}) \geq 0$ for any \mathbf{x}
 - ▶ Symmetric: $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \mathcal{K}(\mathbf{x}', \mathbf{x})$
 - ▶ Positive semi-definite (PSD):
for any $N > 0$, $\mathbf{X} \in \mathbb{R}^{N \times D}$, $\mathbf{K}_{\mathbf{X}\mathbf{X}}$ is a PSD matrix.

Constructing kernels

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp\left(-\frac{1}{2\ell^2} \|\mathbf{x} - \mathbf{x}'\|_2^2\right)$$

$\sigma_f = 1, \ell = 0.25$

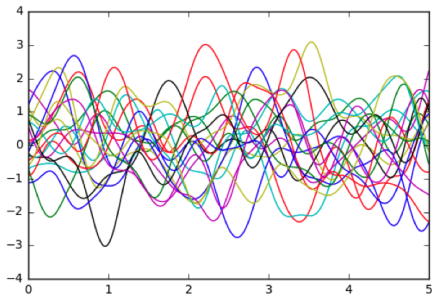


Figure from Richard Wilkinson's GPSS 2019 lecture

Constructing kernels

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp\left(-\frac{1}{2\ell^2} \|\mathbf{x} - \mathbf{x}'\|_2^2\right)$$

$\sigma_f = 1, \ell = 4$

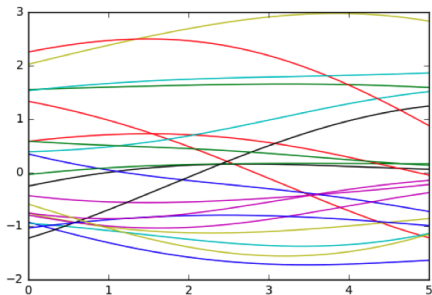


Figure from Richard Wilkinson's GPSS 2019 lecture

Constructing kernels

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp\left(-\frac{1}{2\ell^2} \|\mathbf{x} - \mathbf{x}'\|_2^2\right)$$

$\sigma_f = 1, \ell = 1$

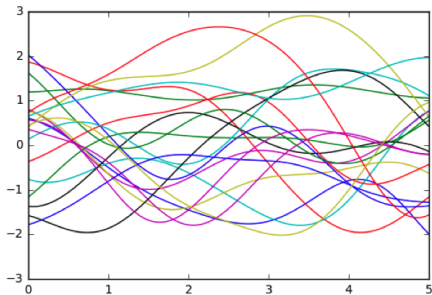


Figure from Richard Wilkinson's GPSS 2019 lecture

Constructing kernels

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp\left(-\frac{1}{2\ell^2} \|\mathbf{x} - \mathbf{x}'\|_2^2\right)$$

$$\sigma_f = 10, \ell = 1$$

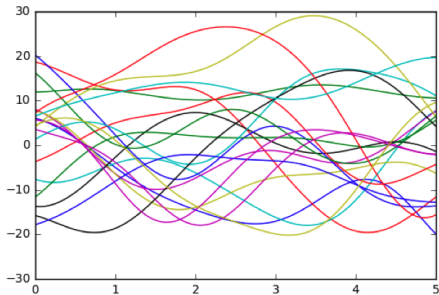
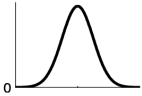
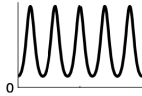

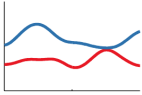
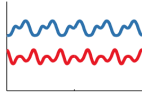
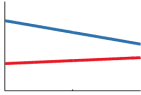


Figure from Richard Wilkinson's GPSS 2019 lecture

Constructing kernels

Kernel name:	Squared-exp (SE)	Periodic (Per)	Linear (Lin)
$k(x, x') =$	$\sigma_f^2 \exp\left(-\frac{(x-x')^2}{2\ell^2}\right)$	$\sigma_f^2 \exp\left(-\frac{2}{\ell^2} \sin^2\left(\pi \frac{x-x'}{p}\right)\right)$	$\sigma_f^2(x-c)(x'-c)$
Plot of $k(x, x')$:			
Functions $f(x)$ sampled from GP prior:			
Type of structure:	local variation	repeating structure	linear functions

See David Duvenaud's PhD thesis, Chapter 2

Summary

The kernel trick

- ▶ $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}')$ for some feature mapping $\phi(\cdot)$
 - ▶ As defining the “inner product” in some feature space (might not be Euclidean and might be infinite-dimensional)
- ▶ Useful when model prediction can be computed using inner products between features
- ▶ Examples:
 - ▶ Non-linear ridge regression \Rightarrow Kernel ridge regression
 - ▶ PCA with non-linear features \Rightarrow Kernel PCA
 - ▶ SVM (read MML book chapter 12)
 - ▶ Gaussian processes (“Probabilistic Inference” course in spring)