

Kernels

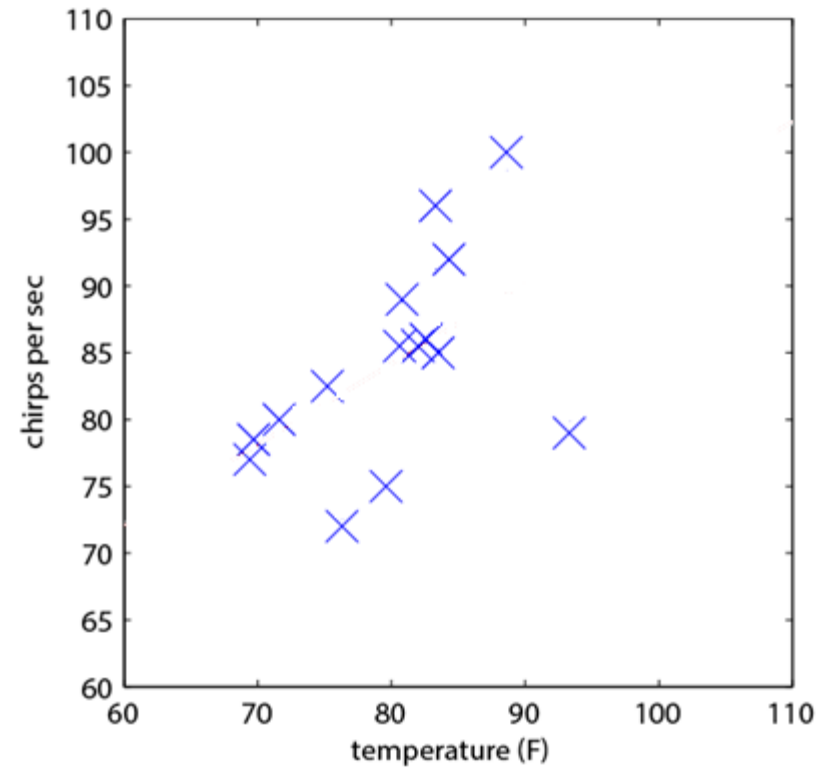
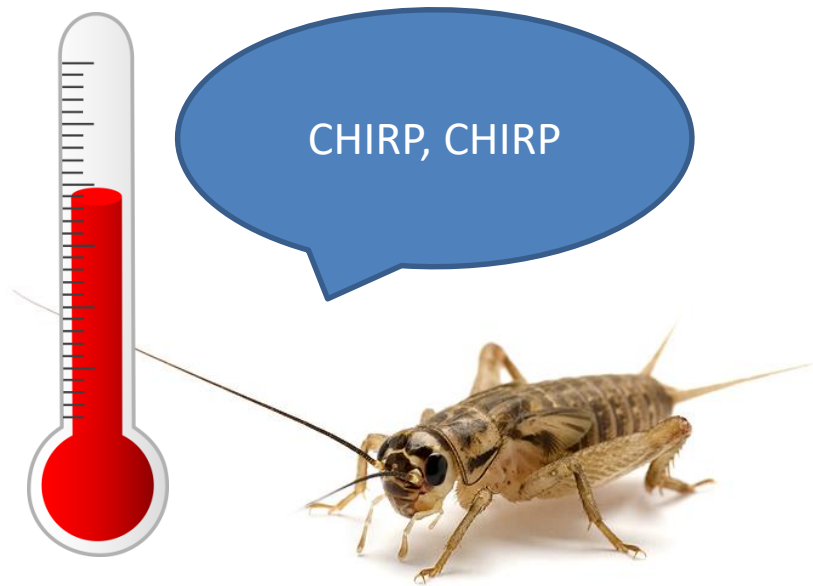
Lecture given by Grady Jensen

Slides created by Sebastian Urban

Notation

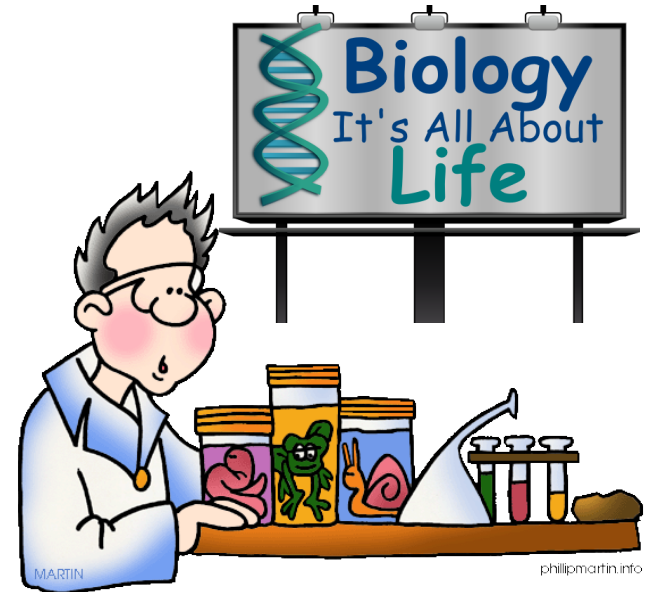
| | | |
|---------------|---------------|--|
| \mathcal{D} | \mathcal{D} | dataset with N samples |
| x_n | $x^{(n)}$ | n-th input sample |
| z_n | $y^{(n)}$ | target output for n-th sample |
| Z | \mathbf{y} | $= (y^{(1)}, y^{(2)}, \dots, y^{(N)})^T$ |
| $y(x)$ | $y(x)$ | output computed for input x by model |
| $y(x_n)$ | $Y^{(n)}$ | output computed for n-th input sample by model |
| | \mathbf{Y} | $= (Y^{(1)}, Y^{(2)}, \dots, Y^{(N)})^T$ |
| W | \mathbf{w} | weights (model parameters) |

Crickets

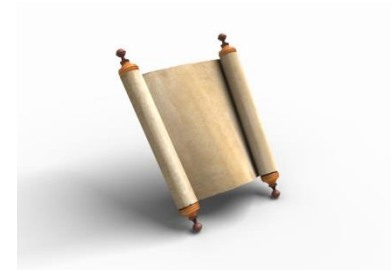


Chirps per second versus temperature.

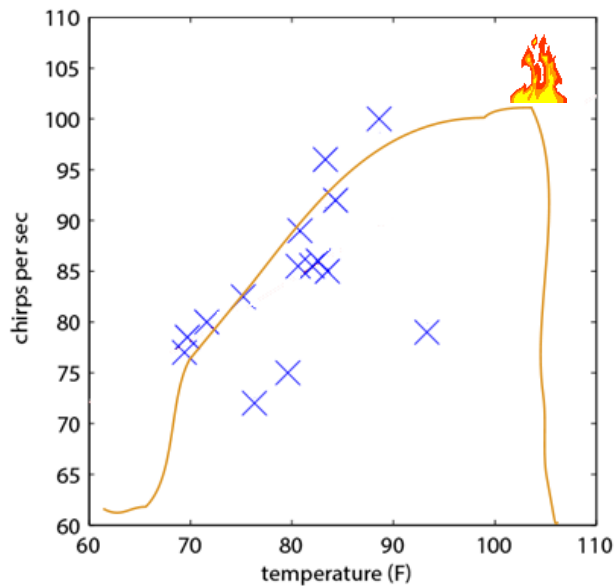
In an ideal world...



Cricket-Chirp model: $y = f_w(x)$



fit model parameters w

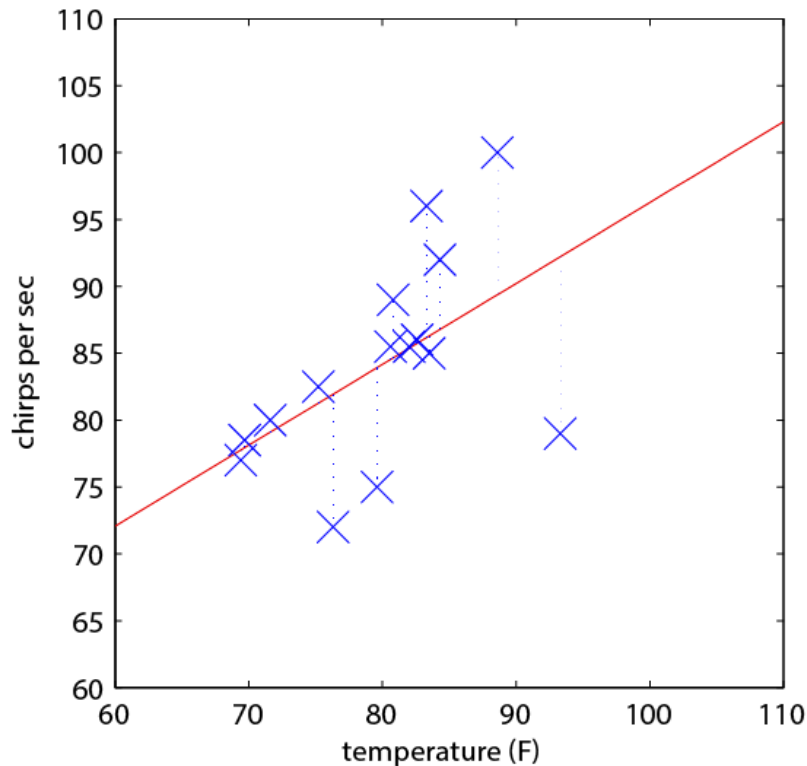




Suppose that a biologist is not available.

Can we make predictions without a scientific model?

Linear regression



Chirps per second versus temperature.

Given the dataset

$$\mathcal{D} = \{(x^{(n)}, y^{(n)}), n = 1, \dots, N\}$$

we want to fit the linear model

$$y(x) = ax + b.$$

We define the *feature map*

$$\boldsymbol{\phi}: \mathbb{R} \rightarrow \mathbb{R}^2, \quad \boldsymbol{\phi}(x) = (x, 1)^T$$

and can formulate our linear model as

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

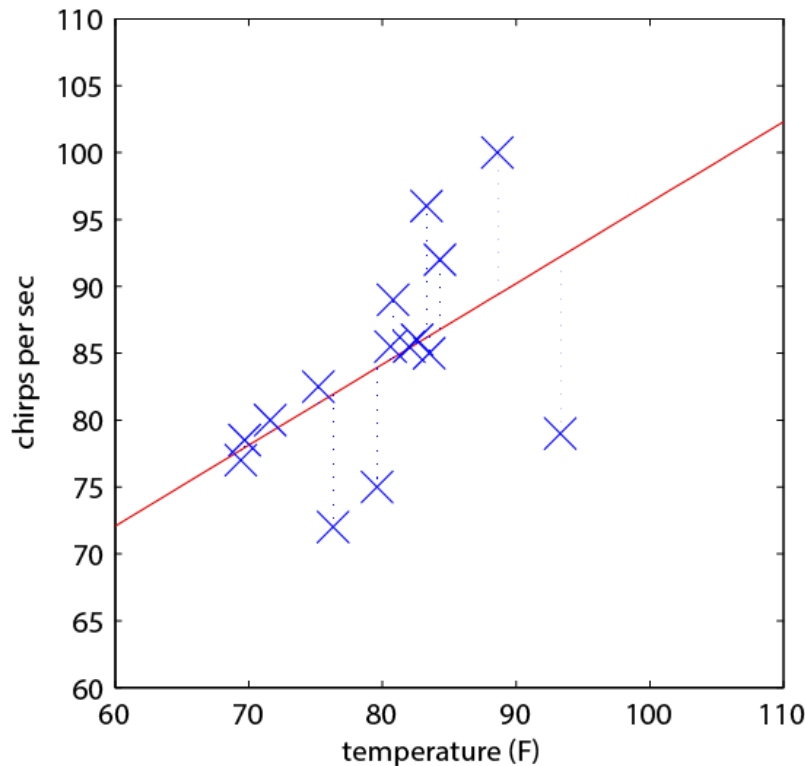
with $\mathbf{w} = (a, b)^T$.

Using the squared prediction error we obtain the objective function

$$E(\mathbf{w}) = \sum_{n=1}^N [y^{(n)} - \mathbf{w}^T \boldsymbol{\phi}^{(n)}]^2$$

where $\boldsymbol{\phi}^{(n)} := \boldsymbol{\phi}(x^{(n)})$.

Linear regression in matrix form



Chirps per second versus temperature.

In matrix form the error function is

$$E(\mathbf{w}) = (\mathbf{y}^T - \mathbf{w}^T \Phi)^T (\mathbf{y}^T - \mathbf{w}^T \Phi)$$

with $\mathbf{y}_n = y^{(n)}$ and $\Phi_{ij} = \phi_i(\mathbf{x}^{(j)})$.

Minimum error is obtained by

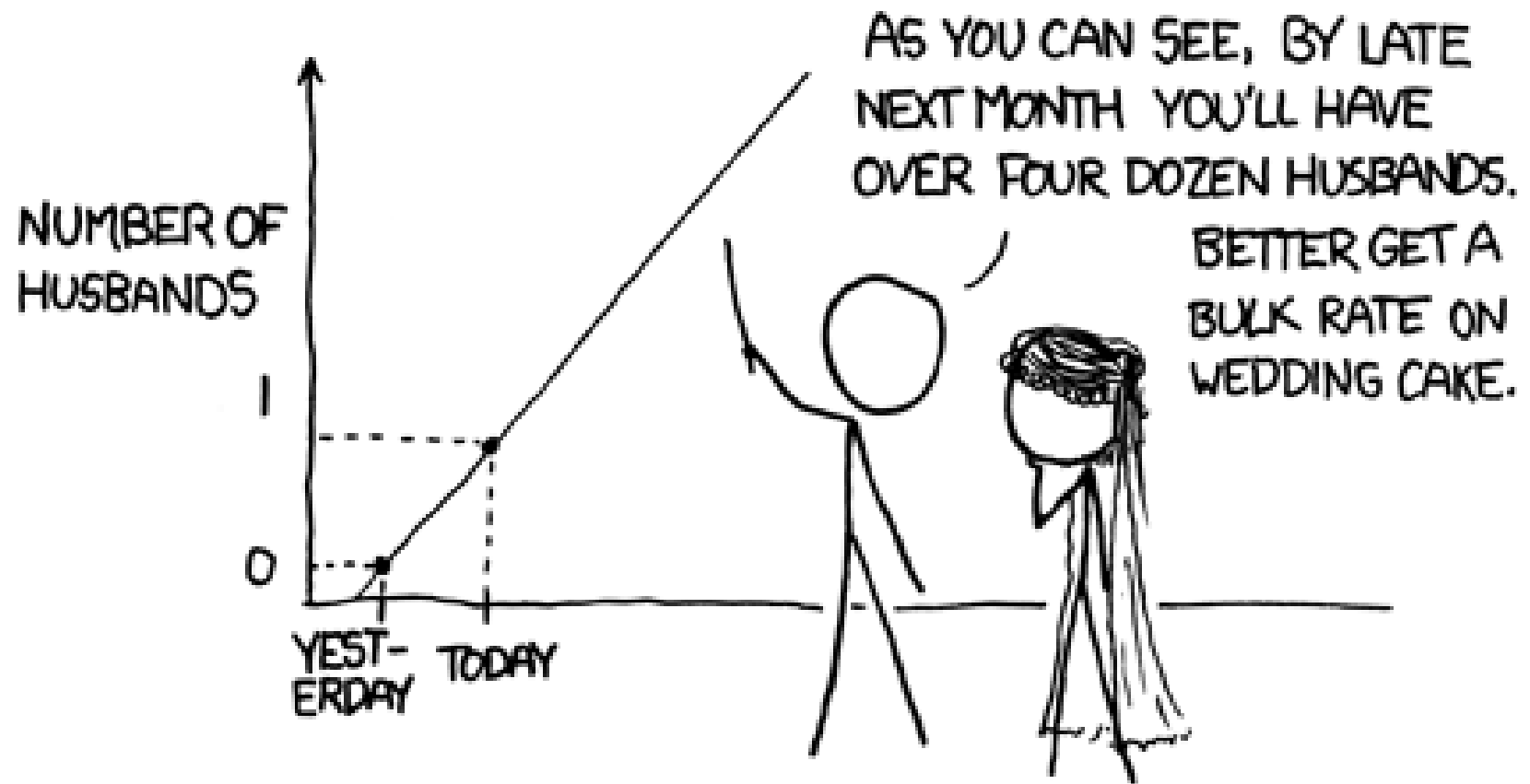
$$\mathbf{w} = (\Phi \Phi^T)^{-1} \Phi \mathbf{y}$$

where $(\Phi \Phi^T)^{-1} \Phi$ is called the Moore-Penrose pseudo-inverse of Φ .

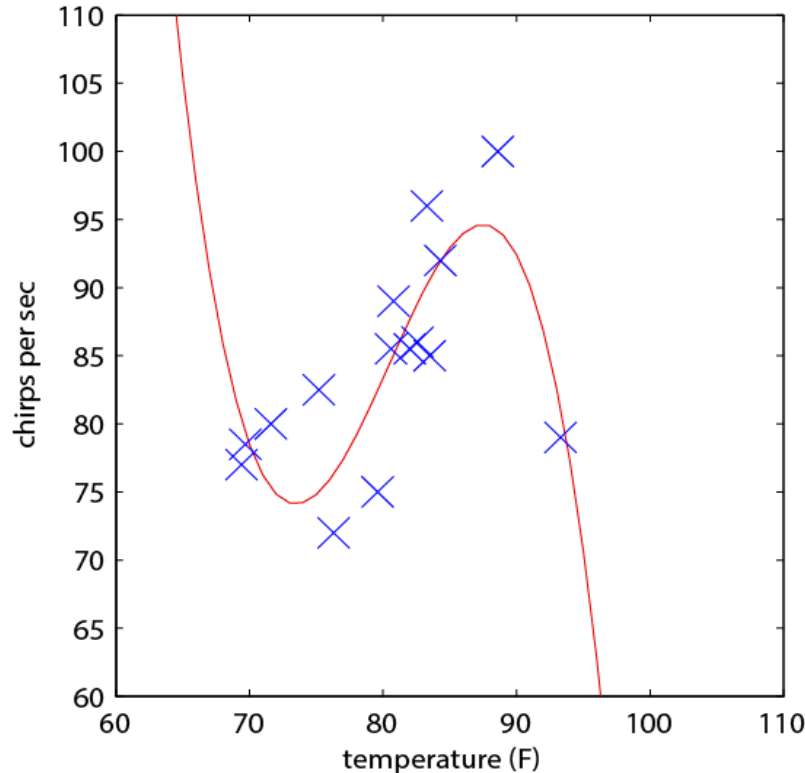
$(\Phi \Phi^T)$ is a $M \times M$ matrix where M is the number of basis functions.

$\Rightarrow O(M^3)$ operations to calculate $(\Phi \Phi^T)^{-1}$.

MY HOBBY: EXTRAPOLATING



More complex basis functions



Chirps per second versus temperature.

We can also use a more complex model, for example a 3th order polynomial

$$y(x) = w_1 + w_2x + w_3x^2 + w_4x^3 .$$

This yields the feature map

$$\phi(x) = \begin{pmatrix} 1 \\ x \\ x^2 \\ x^3 \end{pmatrix} .$$

Quality of solution? **Worse.**

Model complexity? **Higher.**

Why? Model does not match data.

Dealing with unknown basis functions

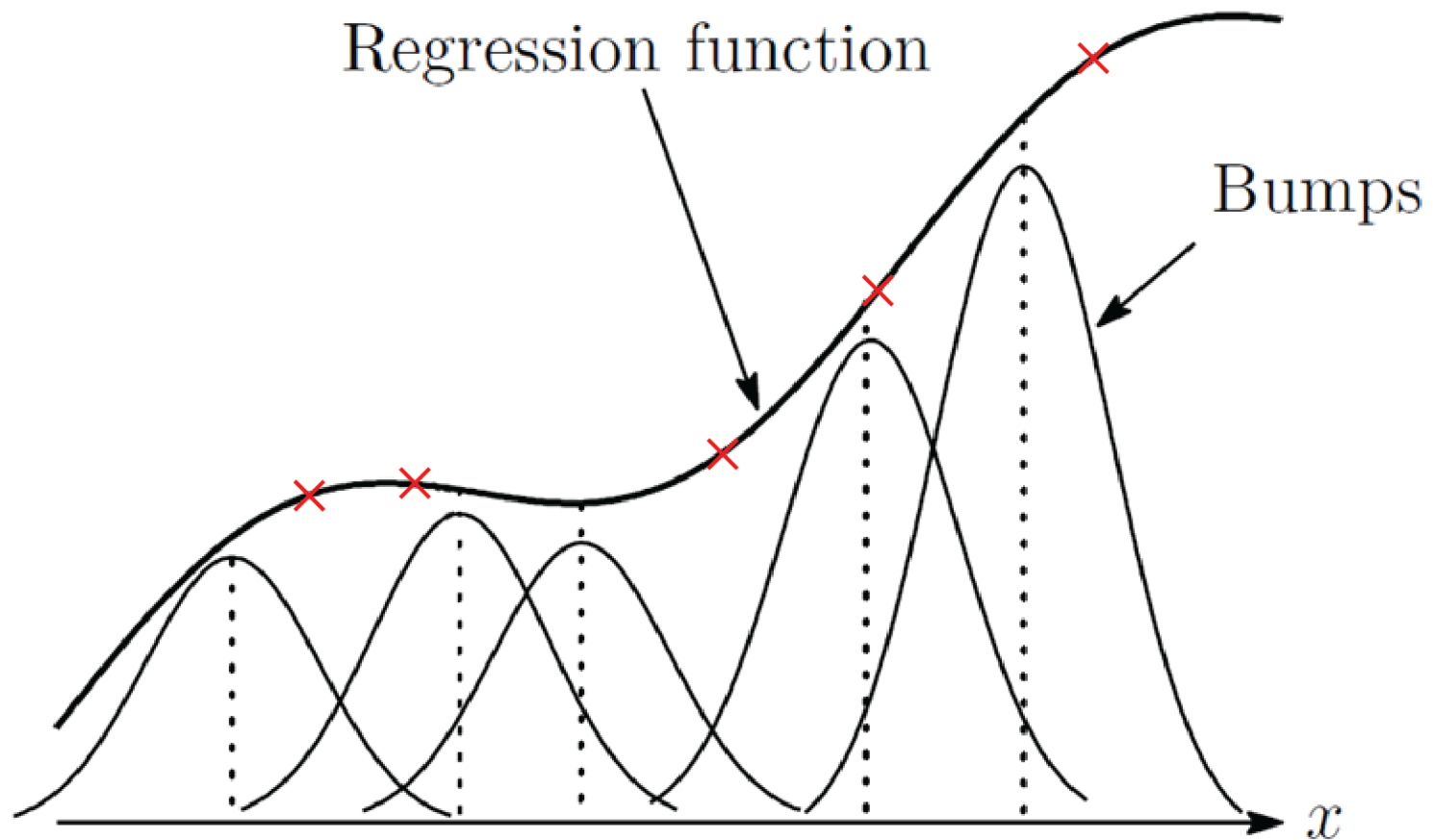
One way to deal with overfitting is by using regularization. However it will not help much if the basis functions do not match the data.

So, what should we do if we do not know the model?

Is there a **universal set of basis functions** that can approximate arbitrary functions reasonably well?

Linear regression using “bump” functions

Approximate the function with a sum of bump-shaped functions.



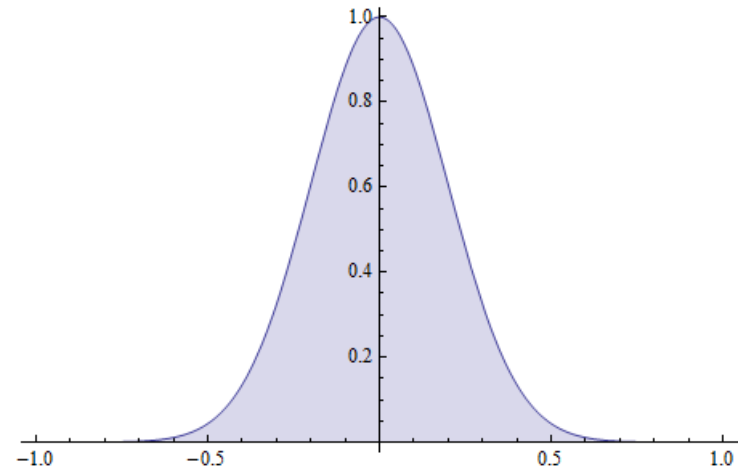
Radial basis functions (RBFs)

Gaussian Radial Basis Functions (RBFs):

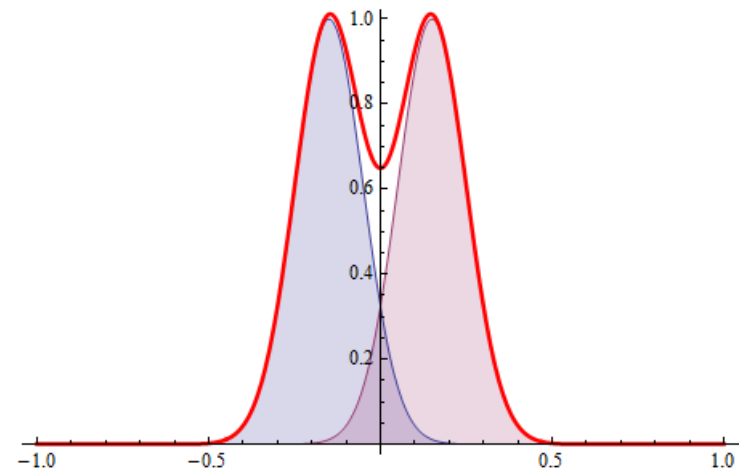
$$\phi_i(\mathbf{x}) = \exp\left(-\frac{(\mathbf{x} - \mathbf{m}^{(i)})^2}{2\alpha^2}\right)$$

The center is at $\mathbf{m}^{(i)}$ and α determines the width of the bump.

Advantage: They decrease smoothly to zero and do not show oscillatory behavior unlike higher order polynomials.

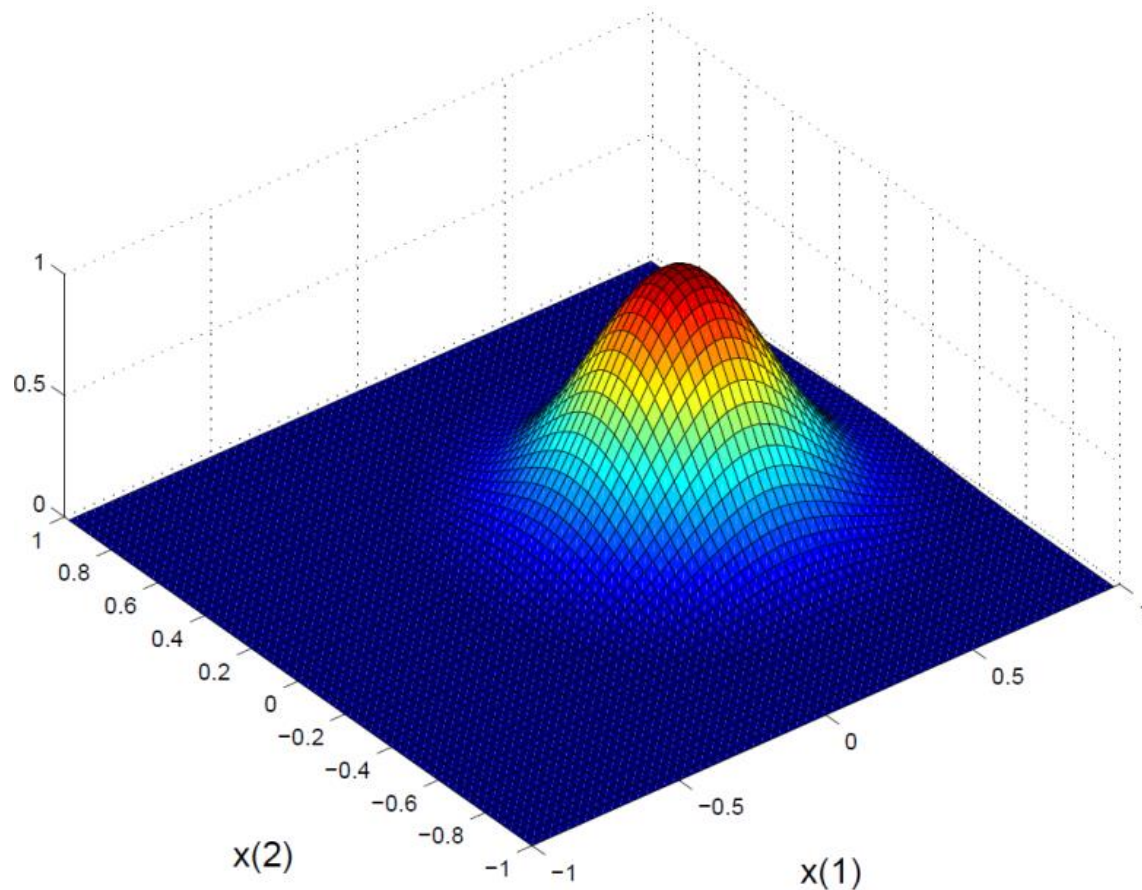


$$\alpha = 0.2, m = 0$$



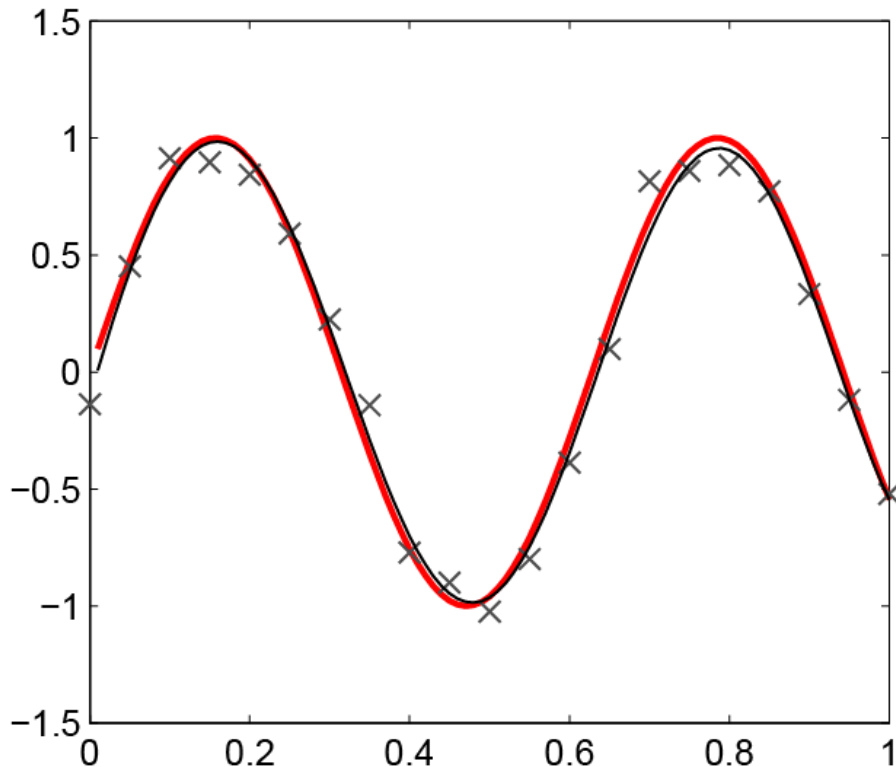
$$\alpha = 0.1, m = -0.15; m = +0.15$$

Gaussian radial basis functions (RBFs)



Gaussian RBF in two dimensional space.

Linear regression using RBFs



Approximation of $\sin(10x)$ by
16 Gaussian RBFs.

Black line: original function
Red line: regression

Use

$$\phi_i(\mathbf{x}) = \exp\left(-\frac{(\mathbf{x} - \mathbf{m}^{(i)})^2}{2\alpha^2}\right)$$

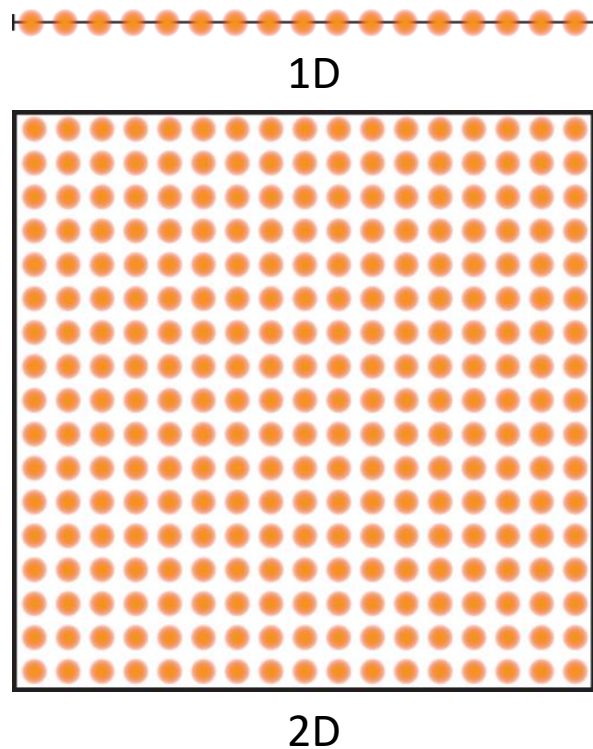
with $i = 1, \dots, 16$ as basis functions.

Apply linear regression.

In one dimension this seems to work well, but does it scale to higher dimensions?

Curse of dimensionality

To cover the data with a constant discretization level (number of RBFs per unit volume) the number of basis functions and weights grows exponentially with the number of dimensions.



| Dimensions | # of basis functions / weights |
|------------|--------------------------------|
| 1 | 16 |
| 2 | $16^2 = 256$ |
| 3 | $16^3 = 4096$ |
| 4 | $16^4 = 65\,536$ |
| 5 | $16^5 = 1\,048\,576$ |
| \vdots | \vdots |
| 10 | $16^{10} \approx 10^{12}$ |

Is there a way to reduce the number of required weights?

DUAL REPRESENTATION

Notation

| | | |
|---------------|---------------|--|
| \mathcal{D} | \mathcal{D} | dataset with N samples |
| x_n | $x^{(n)}$ | n-th input sample |
| z_n | $y^{(n)}$ | target output for n-th sample |
| Z | \mathbf{y} | $= (y^{(1)}, y^{(2)}, \dots, y^{(N)})^T$ |
| $y(x)$ | $y(x)$ | output computed for input x by model |
| $y(x_n)$ | $Y^{(n)}$ | output computed for n-th input sample by model |
| Y | \mathbf{Y} | $= (Y^{(1)}, Y^{(2)}, \dots, Y^{(N)})^T$ |
| W | \mathbf{w} | weights (model parameters) |

Row space and null space of a matrix

Let $A \in \mathbb{R}^{N \times M}$ be a matrix and let $\mathbf{A}_i = (A_{i1}, \dots, A_{iM})$ be its i th row.

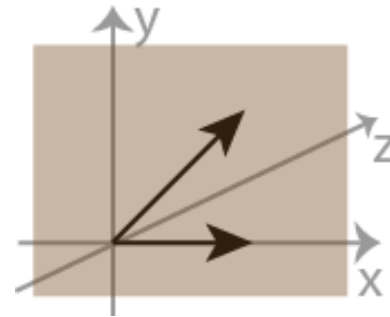
The **row space** \mathcal{A} of A is the set of all possible linear combinations of its rows, $\mathcal{A} = \text{span}(\{\mathbf{A}_1, \dots, \mathbf{A}_n\})$.

The **null space** $\ker(A)$ of A is the orthogonal complement of \mathcal{A} , $\ker(A) = \mathcal{A}^\perp = \{\mathbf{x} \in \mathbb{R}^M : A\mathbf{x} = \mathbf{0}\}$.

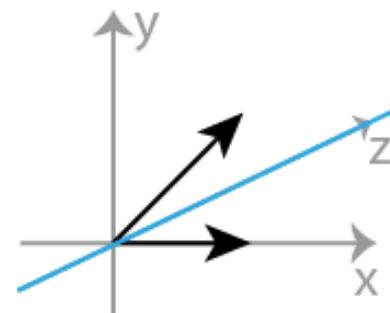
Every element \mathbf{v} in the null space of A is orthogonal to all rows of A , $\mathbf{A}_i \cdot \mathbf{v} = \mathbf{0} \quad \forall i \in \{1, \dots, N\}$.

$$A = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \end{pmatrix}$$

$$\mathcal{A} = \{c_1(1,0,0) + c_2(1,1,0), \mathbf{c} \in \mathbb{R}^2\}$$



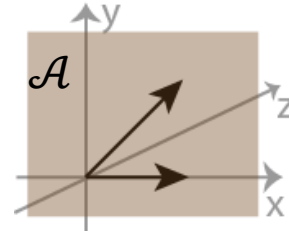
$$\ker(A) = \mathcal{A}^\perp = \{c(0,0,1), c \in \mathbb{R}\}$$



Orthogonal Complement and Decomposition

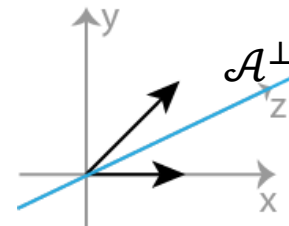
Let \mathcal{A} be a subspace of \mathbb{R}^N .

$$\mathcal{A} = \{c_1(1,0,0) + c_2(1,1,0), c \in \mathbb{R}^2\}$$



The **orthogonal complement** \mathcal{A}^\perp of \mathcal{A} is the set of all vectors that are orthogonal to all elements of \mathcal{A} ,

$$\mathcal{A}^\perp = \{c(0,0,1), c \in \mathbb{R}\}$$



$$\mathcal{A}^\perp = \{\mathbf{z} \in \mathbb{R}^N : \mathbf{x}^T \mathbf{z} = 0 \quad \forall \mathbf{x} \in \mathcal{A}\}.$$

Orthogonal decomposition: Any $\mathbf{w} \in \mathbb{R}^N$ can be written uniquely in the form

$$\mathbf{w} = \hat{\mathbf{w}} + \mathbf{z}$$

with $\hat{\mathbf{w}} \in \mathcal{A}$ and $\mathbf{z} \in \mathcal{A}^\perp$.

$$\mathbf{w} = (2,3,5)$$

$$\hat{\mathbf{w}} = (2,3,0)$$

$$\mathbf{z} = (0,0,5)$$

Null space of training set

Let $\mathcal{X} = \{\mathbf{x}^{(n)}, n = 1, \dots, N\}$ be the set of all training points.

Consider the predictions of the model

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

on the training set, i.e. $\mathbf{x} \in \mathcal{X}$.

Using orthogonal decomposition we write

$$\mathbf{w} = \hat{\mathbf{w}} + \mathbf{z}$$

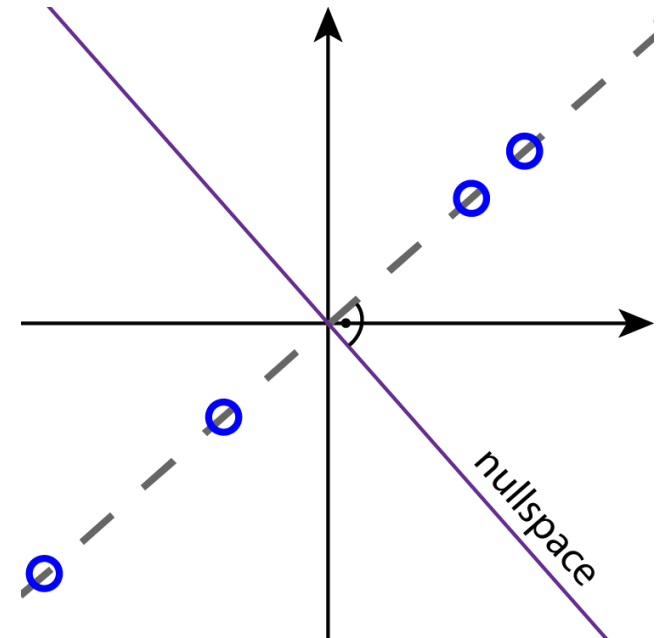
with $\hat{\mathbf{w}} \in \text{span}(\mathcal{X})$ and $\mathbf{z} \in \mathcal{X}^\perp$.

The prediction becomes

$$y(\mathbf{x}) = (\hat{\mathbf{w}} + \mathbf{z})^T \mathbf{x} = \hat{\mathbf{w}}^T \mathbf{x} + \mathbf{z}^T \mathbf{x} = \hat{\mathbf{w}}^T \mathbf{x}.$$

Hence we can assume that

$$\mathbf{w} \in \text{span}(\mathcal{X}).$$



Dual representation

The weight vector \mathbf{w} is thus a linear combination of the training samples \mathcal{X} ,

$$\mathbf{w} = \sum_{n=1}^N a_n \mathbf{x}^{(n)}.$$

The parameters $\mathbf{a} = (a_1, \dots, a_N)$ are called the *dual parameters*.

This also applies when using basis functions,

$$\mathbf{w} = \sum_{n=1}^N a_n \boldsymbol{\phi}(\mathbf{x}^{(n)}).$$

There are **as many dual parameters as training samples**. Their number is independent of the number of basis functions.

Predictions in dual representation

$$\phi^{(m)} = \phi(\mathbf{x}^{(m)})$$

Substituting \mathbf{w} back into the linear regression $y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$ yields

$$y(\mathbf{x}) = \sum_{n=1}^N a_n (\phi^{(n)})^T \phi(\mathbf{x}) = \sum_{n=1}^N a_n K(\mathbf{x}^{(n)}, \mathbf{x}).$$

with the **kernel** function of ϕ ,

$$K(\mathbf{x}, \mathbf{y}) := \phi(\mathbf{x})^T \phi(\mathbf{y}).$$

Using the **Gram matrix**

$$\mathbf{K}_{mn} := K(\mathbf{x}^{(m)}, \mathbf{x}^{(n)}) = \phi(\mathbf{x}^{(m)})^T \phi(\mathbf{x}^{(n)})$$

the predictions on the training set are

$$Y^{(m)} = y(\mathbf{x}^{(m)}) = \sum_{n=1}^N a_n K(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) = \sum_{n=1}^N a_n \mathbf{K}_{mn}$$

$$\mathbf{Y} = \mathbf{a}^T \mathbf{K} \quad \text{with} \quad \mathbf{Y} = (Y^{(1)}, Y^{(2)}, \dots, Y^{(N)})^T.$$

Solution of dual representation

Primal representation:

$$E(\mathbf{w}) = \|\mathbf{y}^T - \mathbf{w}^T \Phi\|_2^2$$

where Φ is the matrix of feature vectors.

Solution:

$$\mathbf{w} = (\Phi \Phi^T)^{-1} \Phi \mathbf{y} = \Phi^+ \mathbf{y}$$

Complexity is $O(M^3)$ where M is the *number of basis functions*.

Predictions:

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

Dual representation:

$$E(\mathbf{a}) = \|\mathbf{y}^T - \mathbf{a}^T \mathbf{K}\|_2^2$$

where \mathbf{K} is the gram matrix (also called kernel matrix).

Solution:

$$\mathbf{a} = (\mathbf{K} \mathbf{K}^T)^{-1} \mathbf{K} \mathbf{y} = \mathbf{K}^+ \mathbf{y}$$

Complexity is $O(N^3)$ where N is the *number of training samples*.

Predictions:

$$y(\mathbf{x}) = \sum_{n=1}^N a_n K(\mathbf{x}^{(n)}, \mathbf{x})$$

“Component a_n weighted with similarity to training sample $\mathbf{x}^{(n)}$.”

Why is the dual representation useful?

$$y(\mathbf{x}) = \sum_{n=1}^N a_n K(\mathbf{x}^{(n)}, \mathbf{x})$$

N is the *number of training samples*.

Lots of basis functions and moderate number of training samples:

⇒ dual representation saves lots of parameters.

Example: Kernel of RBF basis (*not* RBF kernel)

Gaussian RBF basis functions:

$$\phi_i(\mathbf{x}) = \exp\left(-\frac{(\mathbf{x} - \mathbf{m}^{(i)})^2}{2\alpha^2}\right), \quad i = 1, \dots, M$$

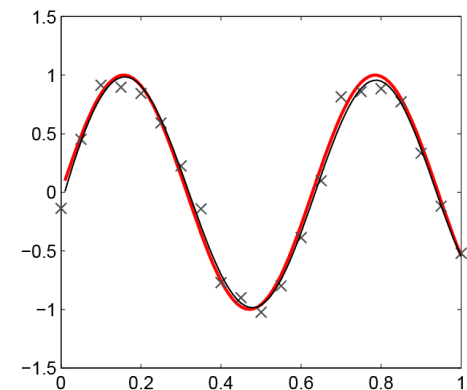
Calculate kernel:

$$\begin{aligned} K(\mathbf{x}, \mathbf{y}) &= \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{y}) = \sum_{i=1}^M \phi_i(\mathbf{x}) \phi_i(\mathbf{y}) \\ &= \exp\left(-\frac{\mathbf{x}^2 + \mathbf{y}^2}{2\alpha^2}\right) \sum_{i=1}^M \exp\left(\frac{\mathbf{m}^{(i)}(\mathbf{x} + \mathbf{y}) - \mathbf{m}^{(i)2}}{\alpha^2}\right) \end{aligned}$$

Linear regression:

$$y(\mathbf{x}) = \sum_{n=1}^N a_n K(\mathbf{x}^{(n)}, \mathbf{x})$$

where N is number of training samples.



Example: Kernel of polynomial basis

The polynomial basis of order D

$$\phi_i(x) = x^{i-1}, \quad i = 1, \dots, D$$

induces the kernel

$$K(x, y) = \boldsymbol{\phi}(x)^T \boldsymbol{\phi}(y) = \sum_{i=1}^D x^{i-1} y^{i-1} = \sum_{i=0}^{D-1} (xy)^i = \frac{1 - (xy)^D}{1 - xy}$$

and the Gram matrix

$$K_{nm} = K(x^{(n)}, x^{(m)}) = \frac{1 - (x^{(n)} x^{(m)})^D}{1 - x^{(n)} x^{(m)}}.$$

Evaluating the kernel is much cheaper than calculating the mapping $\boldsymbol{\phi}$ into feature space and the scalar product explicitly.

PROPERTIES OF KERNELS

Notation

| | | |
|---------------|---------------|--|
| \mathcal{D} | \mathcal{D} | dataset with N samples |
| x_n | $x^{(n)}$ | n-th input sample |
| z_n | $y^{(n)}$ | target output for n-th sample |
| Z | \mathbf{y} | $= (y^{(1)}, y^{(2)}, \dots, y^{(N)})^T$ |
| $y(x)$ | $y(x)$ | output computed for input x by model |
| $y(x_n)$ | $Y^{(n)}$ | output computed for n-th input sample by model |
| Y | \mathbf{Y} | $= (Y^{(1)}, Y^{(2)}, \dots, Y^{(N)})^T$ |
| W | \mathbf{w} | weights (model parameters) |

Properties of a kernel

$$\Phi_{ij} = \phi_i(x^{(j)})$$

Calculate $\Phi^T \Phi$ component-wise,

$$\begin{aligned} (\Phi^T \Phi)_{mn} &= \sum_k (\Phi^T)_{mk} \Phi_{kn} = \sum_k \Phi_{km} \Phi_{kn} \\ &= \phi(x^{(m)})^T \phi(x^{(n)}) = K_{mn} . \end{aligned}$$

Thus we have

$$K = \Phi^T \Phi$$

with **K symmetric**.

K is **positive semi-definite**, that is $\mathbf{b}^T \mathbf{K} \mathbf{b} \geq 0$ for any \mathbf{b} , since

$$0 \leq |\Phi \mathbf{b}|_2^2 = \mathbf{b}^T \Phi^T \Phi \mathbf{b} = \mathbf{b}^T \mathbf{K} \mathbf{b} \quad \text{for any } \mathbf{b}.$$

When is a function a kernel?

Mercer's theorem (for finite input spaces): Consider a *finite* input space $\mathcal{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$ with $K(\mathbf{x}^{(n)}, \mathbf{x}^{(m)})$ a function on \mathcal{X} .

Then $K(\mathbf{x}^{(n)}, \mathbf{x}^{(m)})$ is a *kernel function*, that is a scalar product in a feature space, if and only if K is symmetric and the matrix $\mathbf{K}_{nm} = K(\mathbf{x}^{(n)}, \mathbf{x}^{(m)})$ is positive semi-definite.

Proof.

\mathbf{K} symmetric \Rightarrow The eigenvalue decomposition has the form $\mathbf{K} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$, where $\mathbf{\Lambda}$ is a diagonal matrix containing the eigenvalues $\lambda_t = \Lambda_{tt}$ of \mathbf{K} , and \mathbf{V} is an orthogonal matrix containing the eigenvectors $\mathbf{V}_{\cdot t}$.

\mathbf{K} is positive semi-definite \Rightarrow All eigenvalues are non-negative, $\lambda_t \geq 0$.

Define the feature map

$$\phi_t(\mathbf{x}^{(n)}) = \sqrt{\lambda_t} \mathbf{V}_{nt}$$

and see that it corresponds to the kernel by calculating the scalar product

$$\phi(\mathbf{x}^{(n)})^T \phi(\mathbf{x}^{(m)}) = \sum_{t=1}^N \mathbf{V}_{nt} \lambda_t \mathbf{V}_{mt} = (\mathbf{V}\mathbf{\Lambda}\mathbf{V}^T)_{nm} = \mathbf{K}_{nm} .$$



Making kernels from kernels

Use the following rules to construct more complex kernels from simple ones.

Let K_1 and K_2 be kernels on $\mathcal{X} \subseteq \mathbb{R}^n$, then the following functions are kernels:

1. $K(\mathbf{x}, \mathbf{y}) = K_1(\mathbf{x}, \mathbf{y}) + K_2(\mathbf{x}, \mathbf{y})$,
2. $K(\mathbf{x}, \mathbf{y}) = a K_1(\mathbf{x}, \mathbf{y})$ for $a > 0$,
3. $K(\mathbf{x}, \mathbf{y}) = K_1(\mathbf{x}, \mathbf{y})K_2(\mathbf{x}, \mathbf{y})$,
4. $K(\mathbf{x}, \mathbf{y}) = K_3(\boldsymbol{\phi}(\mathbf{x}), \boldsymbol{\phi}(\mathbf{y}))$ for K_3 kernel on \mathbb{R}^m and $\boldsymbol{\phi}: \mathcal{X} \rightarrow \mathbb{R}^m$,
5. $K(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{B} \mathbf{y}$ for \mathbf{B} symmetric and positive semi-definite $n \times n$ matrix.

Proofs.

Fix a finite set of points $\{\mathbf{x}_1, \dots, \mathbf{x}_l\}$ and let $\mathbf{K}, \mathbf{K}_1, \mathbf{K}_2$ be the corresponding Gram matrices of the kernel functions K, K_1, K_2 .

Obviously the resulting \mathbf{K} is symmetric for every case.

It remains to be shown that \mathbf{K} is positive semi-definite, that is

$$\mathbf{z}^T \mathbf{K} \mathbf{z} \geq 0 \quad \text{for all } \mathbf{z}.$$

Making kernels from kernels

1. $K(\mathbf{x}, \mathbf{y}) = K_1(\mathbf{x}, \mathbf{y}) + K_2(\mathbf{x}, \mathbf{y})$

We have

$$\mathbf{z}^T (\mathbf{K}_1 + \mathbf{K}_2) \mathbf{z} = \mathbf{z}^T \mathbf{K}_1 \mathbf{z} + \mathbf{z}^T \mathbf{K}_2 \mathbf{z} \geq 0$$

since K_1 and K_2 are positive semi-definite.

2. $K(\mathbf{x}, \mathbf{y}) = a K_1(\mathbf{x}, \mathbf{y})$ for $a > 0$

Similarly, we have

$$\mathbf{z}^T a \mathbf{K}_1 \mathbf{z} = a \mathbf{z}^T \mathbf{K}_1 \mathbf{z} \geq 0 .$$

Making kernels from kernels

3. $K(\mathbf{x}, \mathbf{y}) = K_1(\mathbf{x}, \mathbf{y})K_2(\mathbf{x}, \mathbf{y})$

We explicitly construct a feature space corresponding to $K(\mathbf{x}, \mathbf{y})$,

$$\begin{aligned} K(\mathbf{x}, \mathbf{y}) &= K_1(\mathbf{x}, \mathbf{y})K_2(\mathbf{x}, \mathbf{y}) = \boldsymbol{\phi}_1(\mathbf{x})^T \boldsymbol{\phi}_1(\mathbf{y}) \boldsymbol{\phi}_2(\mathbf{x})^T \boldsymbol{\phi}_2(\mathbf{y}) \\ &= \sum_i \phi_1(\mathbf{x})_i \phi_1(\mathbf{y})_i \sum_j \phi_2(\mathbf{x})_j \phi_2(\mathbf{y})_j \\ &= \sum_{i,j} \phi_1(\mathbf{x})_i \phi_2(\mathbf{x})_j \phi_1(\mathbf{y})_i \phi_2(\mathbf{y})_j = \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{y}) \end{aligned}$$

with $\boldsymbol{\phi}(\mathbf{x}) = \boldsymbol{\phi}_1(\mathbf{x}) \otimes \boldsymbol{\phi}_2(\mathbf{x})$ (Kronecker product).

Making kernels from kernels

4. $K(\mathbf{x}, \mathbf{y}) = K_3(\boldsymbol{\phi}(\mathbf{x}), \boldsymbol{\phi}(\mathbf{y}))$ for K_3 kernel on \mathbb{R}^m and $\boldsymbol{\phi}: \mathcal{X} \rightarrow \mathbb{R}^m$

Since K_3 is a kernel for all input values there is nothing to prove.

5. $K(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{B} \mathbf{y}$ for \mathbf{B} symmetric and positive semi-definite $n \times n$ matrix

A positive-definite matrix has exactly one square root $\mathbf{B}^{1/2}$ with the property $\mathbf{B}^{1/2} \mathbf{B}^{1/2} = \mathbf{B}$.

Define the feature map $\boldsymbol{\phi}(\mathbf{x}) = \mathbf{B}^{1/2} \mathbf{x}$ and see that $\boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{y}) = (\mathbf{B}^{1/2} \mathbf{x})^T \mathbf{B}^{1/2} \mathbf{y} = \mathbf{x}^T \mathbf{B}^{T/2} \mathbf{B}^{1/2} \mathbf{x} = \mathbf{x}^T \mathbf{B} \mathbf{y} = K(\mathbf{x}, \mathbf{y})$.



EXAMPLE AND SUMMARY

Notation

| | | |
|---------------|---------------|--|
| \mathcal{D} | \mathcal{D} | dataset with N samples |
| x_n | $x^{(n)}$ | n-th input sample |
| z_n | $y^{(n)}$ | target output for n-th sample |
| Z | \mathbf{y} | $= (y^{(1)}, y^{(2)}, \dots, y^{(N)})^T$ |
| $y(x)$ | $y(x)$ | output computed for input x by model |
| $y(x_n)$ | $Y^{(n)}$ | output computed for n-th input sample by model |
| Y | \mathbf{Y} | $= (Y^{(1)}, Y^{(2)}, \dots, Y^{(N)})^T$ |
| W | \mathbf{w} | weights (model parameters) |

Example: Gaussian kernel

Linear regression:

$$y(\mathbf{x}) = \sum_{n=1}^N a_n K(\mathbf{x}^{(n)}, \mathbf{x})$$

„Similarity“ between data points \mathbf{x} and \mathbf{y} is calculated by the **Gaussian kernel**,

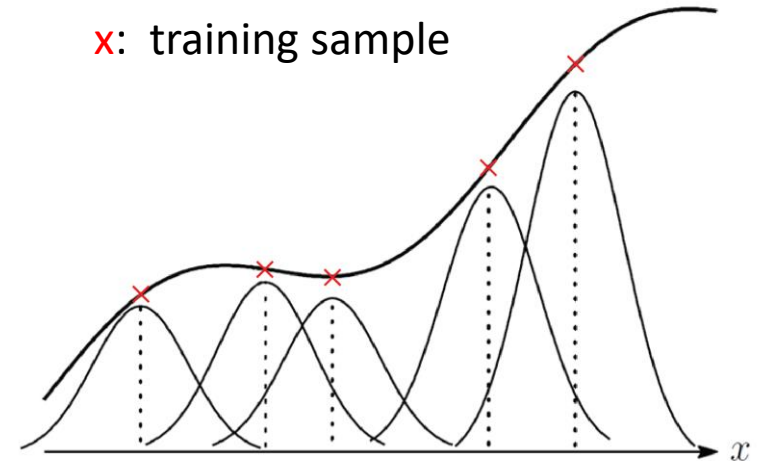
$$K(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{|\mathbf{x} - \mathbf{y}|^2}{2\sigma^2}\right),$$

where σ is a scale parameter.

This puts a bump shaped function on every training sample,

$$y(\mathbf{x}) = \sum_{n=1}^N a_n \exp\left(-\frac{|\mathbf{x} - \mathbf{x}^{(n)}|^2}{2\sigma^2}\right).$$

Due to the fast decline of the exponential function, data points far away do not influence the result.



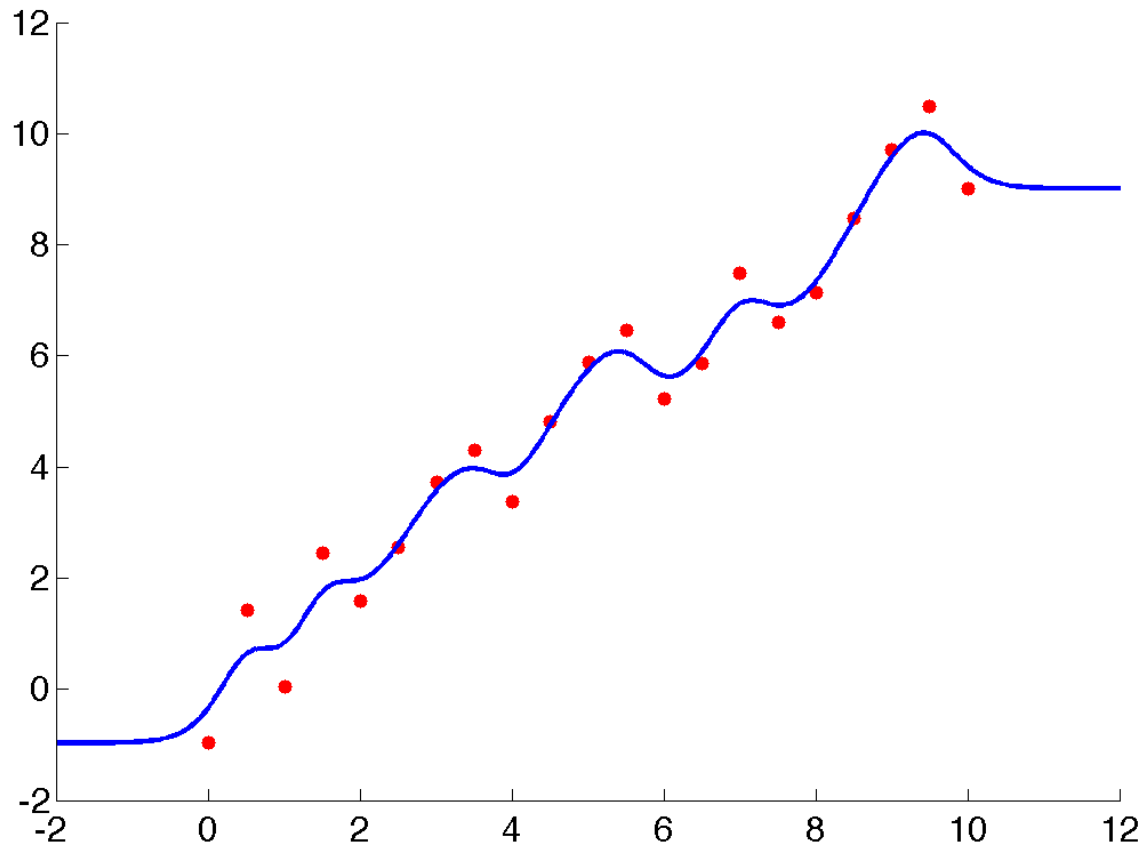
Is this really a kernel?

Yes, use rules from „Making kernels from kernels“ to prove it.

Corresponding feature space?

Infinite dimensional, will be shown in an exercise.

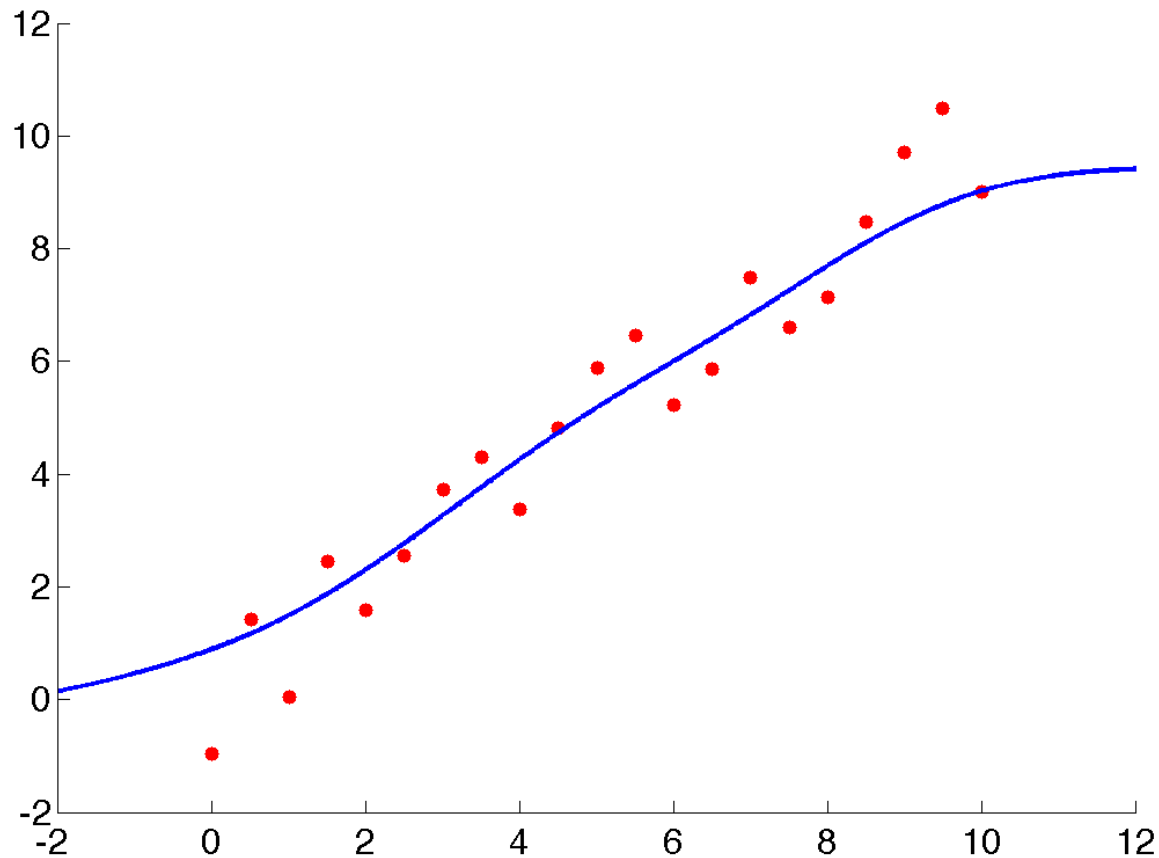
Example: Gaussian kernel



<https://alliance.seas.upenn.edu/~cis520/wiki/index.php?n=Lectures.LocalLearning>

Linear regression using a Gaussian kernel with $\sigma = 0.5$.
Generating function was linear.

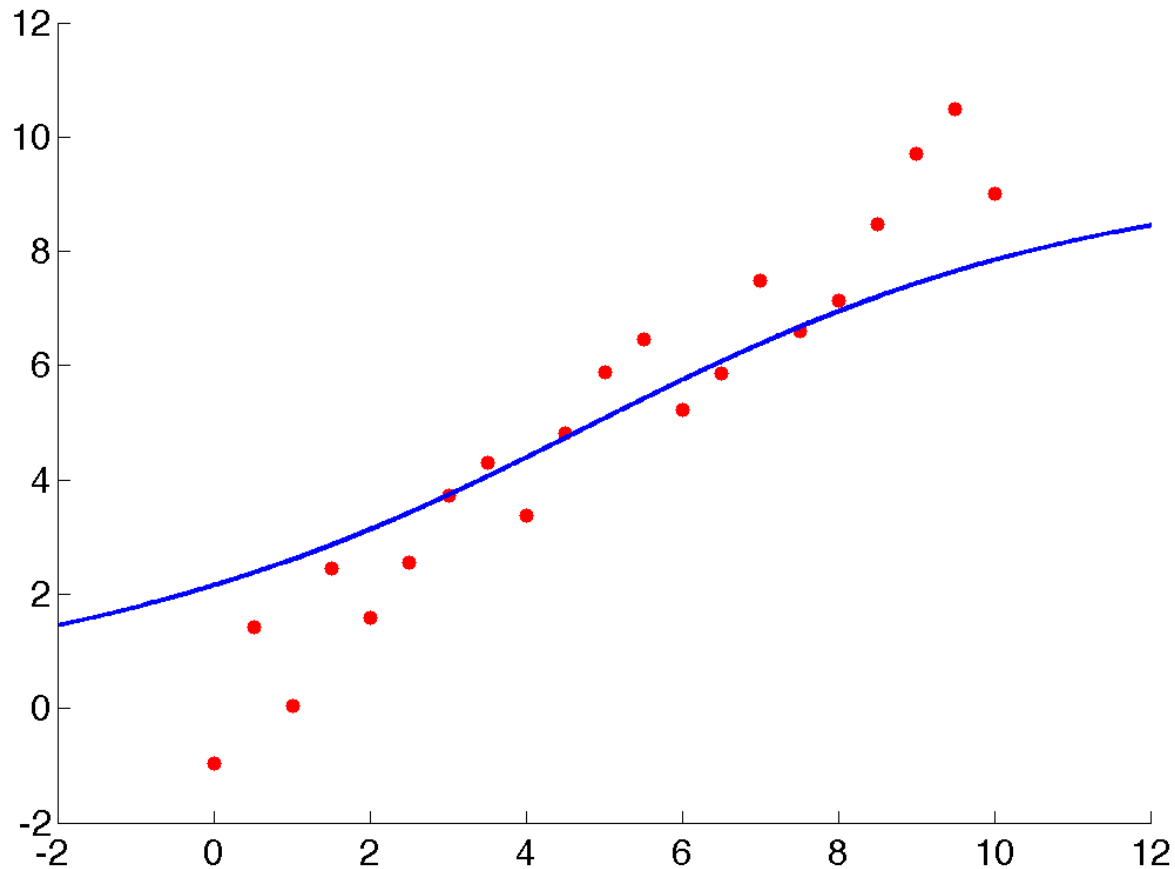
Example: Gaussian kernel



<https://alliance.seas.upenn.edu/~cis520/wiki/index.php?n=Lectures.LocalLearning>

Linear regression using a Gaussian kernel with $\sigma = 2.0$.
Generating function was linear.

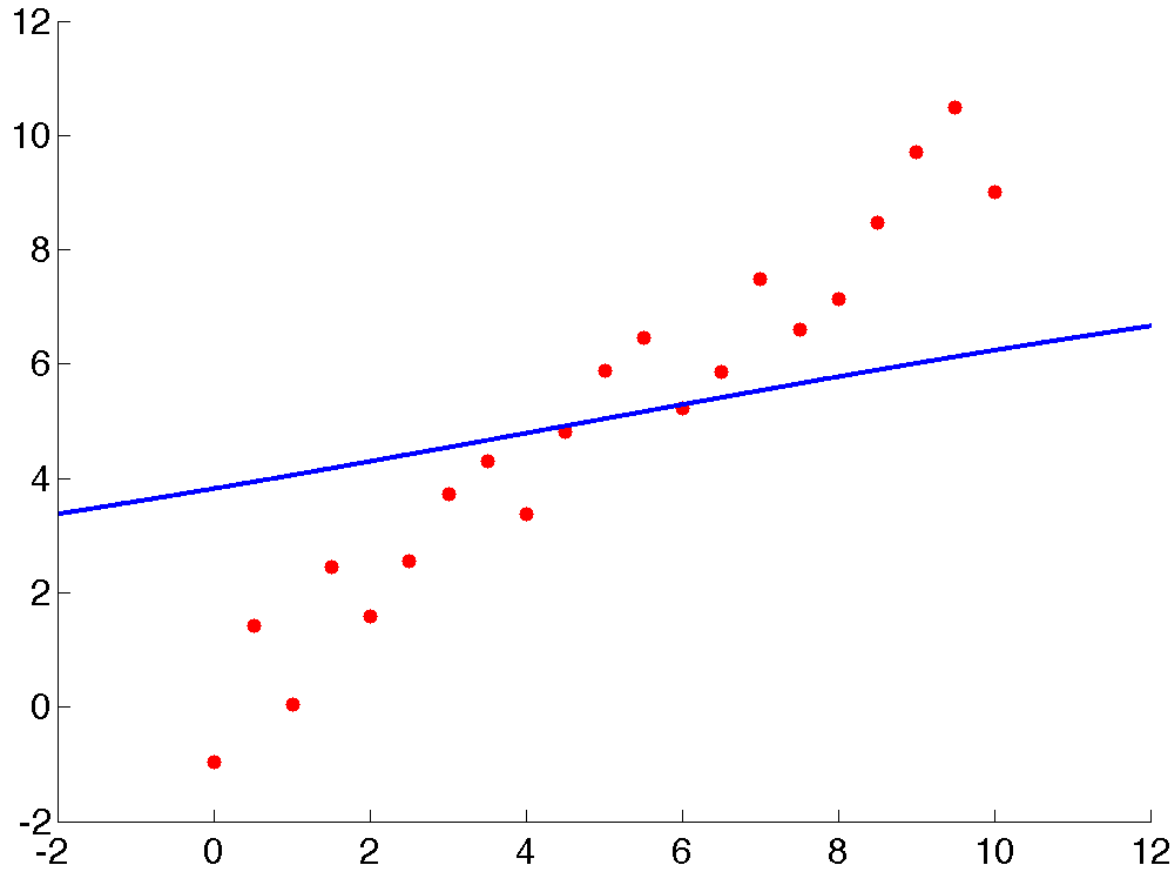
Example: Gaussian kernel



<https://alliance.seas.upenn.edu/~cis520/wiki/index.php?n=Lectures.LocalLearning>

Linear regression using a Gaussian kernel with $\sigma = 4.0$.
Generating function was linear.

Example: Gaussian kernel



<https://alliance.seas.upenn.edu/~cis520/wiki/index.php?n=Lectures.LocalLearning>

Linear regression using a Gaussian kernel with $\sigma = 8.0$.
Generating function was linear.

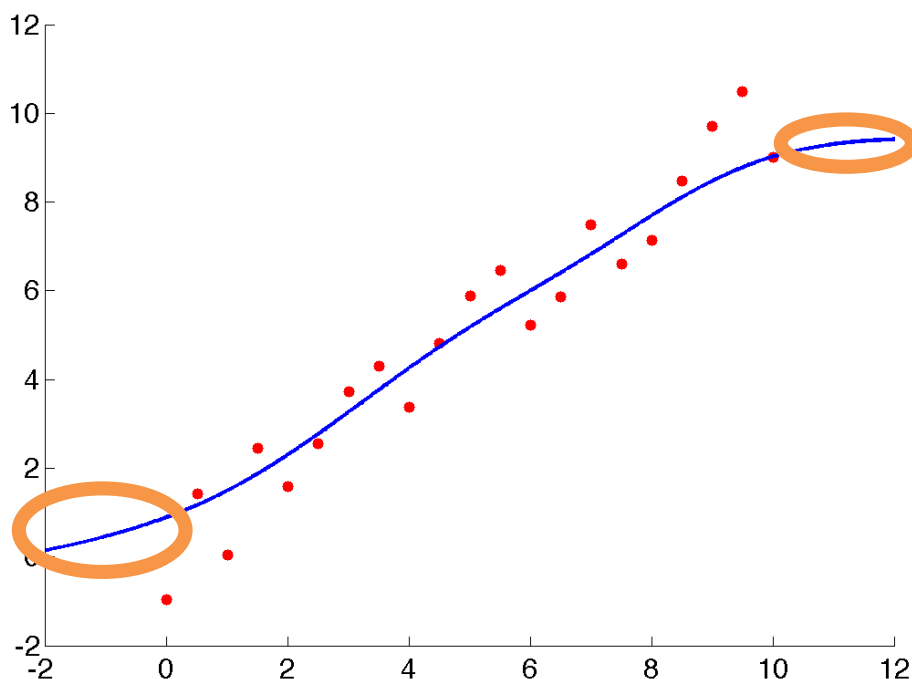
Example: Gaussian kernel

$$K(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{|\mathbf{x} - \mathbf{y}|^2}{2\sigma^2}\right)$$

The quality of the result is very sensitive to the choice of the variance σ .

Use cross-validation to choose the right value.

Combining model-based and model-free approaches



If there is a feature space $\phi(\mathbf{x})$ that partly explains the data, we can add the kernel induced by that feature space, $\phi(\mathbf{x})^T \phi(\mathbf{y})$, to a generic kernel $K_2(\mathbf{x}, \mathbf{y})$ to obtain a better fit.

$$K(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^T \phi(\mathbf{y}) + K_2(\mathbf{x}, \mathbf{y})$$

Here, we could combine a linear model, $\phi(\mathbf{x}) = \mathbf{x}$, with the Gaussian kernel $K_2(\mathbf{x}, \mathbf{y})$.

Problem: Gaussian kernel cannot extrapolate, since

$$\sum_n a_n \exp\left(-\frac{|\mathbf{x} - \mathbf{x}^{(n)}|^2}{2\sigma^2}\right) \rightarrow 0$$

further away from the training set.

Linear model: computes general trend of data.

Gaussian kernel: captures local variations.

Commonly used kernels

- ▶ Linear: $K(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$
- ▶ Polynomial: $K(\mathbf{x}, \mathbf{y}) = ((\mathbf{x}^T \mathbf{y}) + c)^d$
- ▶ Gaussian: $K(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{|\mathbf{x}-\mathbf{y}|^2}{2\sigma^2}\right)$
- ▶ Sometimes functions are used as kernels that are not positive semi-definite. Usually this works but it can lead to strange results.

Applications of kernels

- ▶ Regression
 - ▶ Gaussian Processes (GPs)
 - ▶ Classification (especially Support Vector Machines)
 - ▶ Principal Component Analysis (PCA)
 - ▶ *Every algorithm with a scalar product.*
-
- ▶ Kernels can also be defined on domains other than \mathbb{R}^N , for example strings. This allows the above methods to be applied on new domains.

Summary

- ▶ Kernels compute the value of the scalar product in a high dimensional feature space without explicitly computing the features.
- ▶ They can be used in any model that depends (or that can be rewritten so that it depends) on the scalar product of the training samples.
- ▶ Not every function is a kernel.
- ▶ Use kernel construction rules to prove that a function is a valid kernel.

References

C. Bishop

Pattern Recognition and Machine Learning

Springer-Verlag 2006

N. Cristianini, J. Shawe-Taylor

*An Introduction to Support Vector Machines and other
kernel-based learning methods*

Cambridge University Press 2000