

# Basis Functions

Volker Tresp  
Summer 2015

*I am an AI optimist. We've got a lot of work in machine learning, which is sort of the polite term for AI nowadays because it got so broad that it's not that well defined.*

Bill Gates (Scientific American Interview, 2004)

*"If you invent a breakthrough in artificial intelligence, so machines can learn," Mr. Gates responded, "that is worth 10 Microsofts." (Quoted in NY Times, Monday March 3, 2004)*

# Amazon Europe Machine Learning Team Coming To Berlin!

Posted by Victoria Nicholl on Fri, 18/01/2013 - 14:37

Amazon is building a European Machine Learning (ML) team in Berlin! Machine Learning Scientists at Amazon are technical leaders who develop planet-scale platforms for machine learning on the cloud, assist the benchmarking and future development of existing machine learning applications across Amazon, and help develop novel and infinitely-scalable applications that optimize Amazon's systems using cutting edge quantitative techniques. The ML team innovates algorithms that model patterns within data to drive automated decisions at scale in all corners of the company, including our eCommerce site and subsidiaries, Amazon Web Services, Seller & Buyer Services and Digital Media including Kindle. Amazon was one of the first companies to build eCommerce customer recommendations, fraud detection, and product search using machine learning innovations. Being part of the Machine Learning team at Amazon is one of the most exciting machine learning job opportunities in the world today. If you have deep technical knowhow in Machine Learning, know how to deliver, are deeply technical, highly innovative and long for the opportunity to build solutions to challenging problems that directly affect millions of people: there may be no better place than Amazon for you to impact the world!

If you are interested send your CV to [strategic-recruiting@amazon.com](mailto:strategic-recruiting@amazon.com).

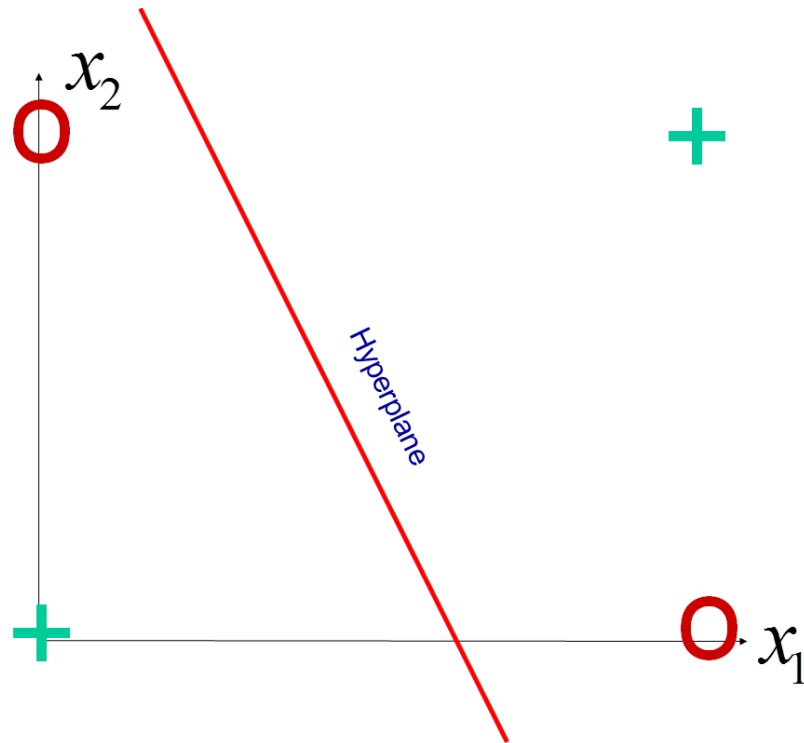
## Nonlinear Mappings and Nonlinear Classifiers

- Regression:
  - Linearity is often a good assumption when many inputs influence the output
  - Some natural laws are (approximately) linear  $F = ma$
  - But in general, it is rather unlikely that a true function is linear
- Classification:
  - Similarly, it is often not reasonable to assume that the classification boundaries are linear hyper planes

## Trick

- We simply transform the input into a high-dimensional space where the regression/classification is again linear!
- Other view: let's define appropriate features
- Other view: let's define appropriate basis functions

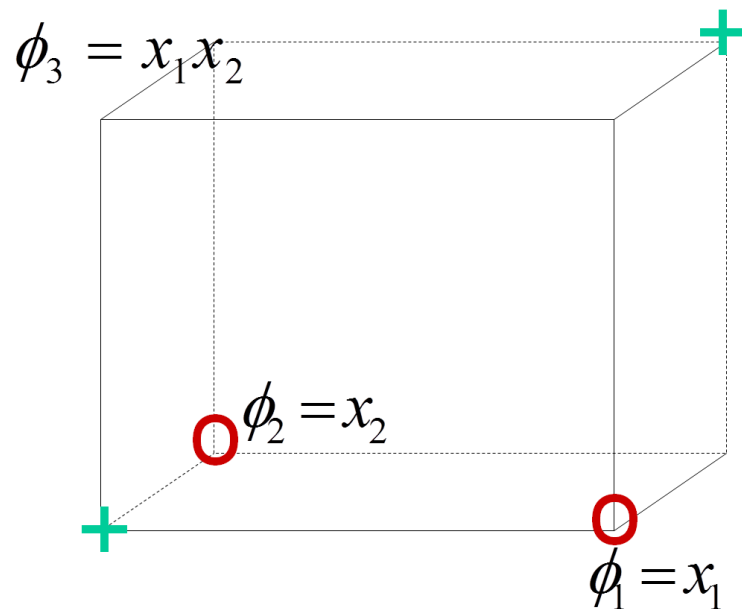
XOR is not linearly separable



## Trick: Let's Add Basis Functions

- Linear Model: input vector:  $1, x_1, x_2$
- Let's consider  $x_1x_2$  in addition
- The interaction term  $x_1x_2$  couples two inputs nonlinearly

With a Third Input  $z_3 = x_1x_2$  the XOR Becomes Linearly Separable

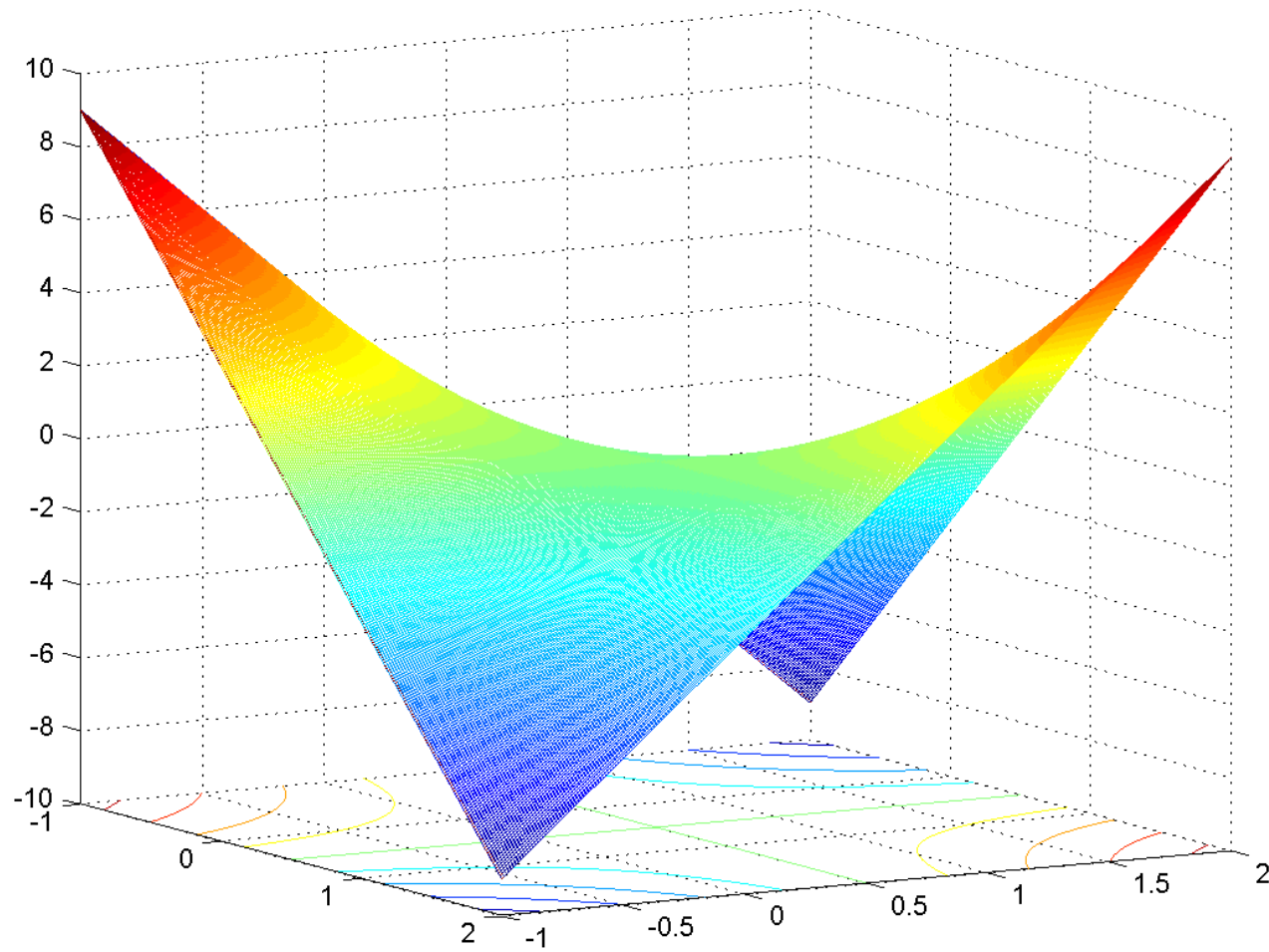


$$f(\mathbf{x}) = 1 - 2x_1 - 2x_2 + 4x_1x_2 = \phi_1(x) - 2\phi_2(x) - 2\phi_3(x) + 4\phi_4(x)$$

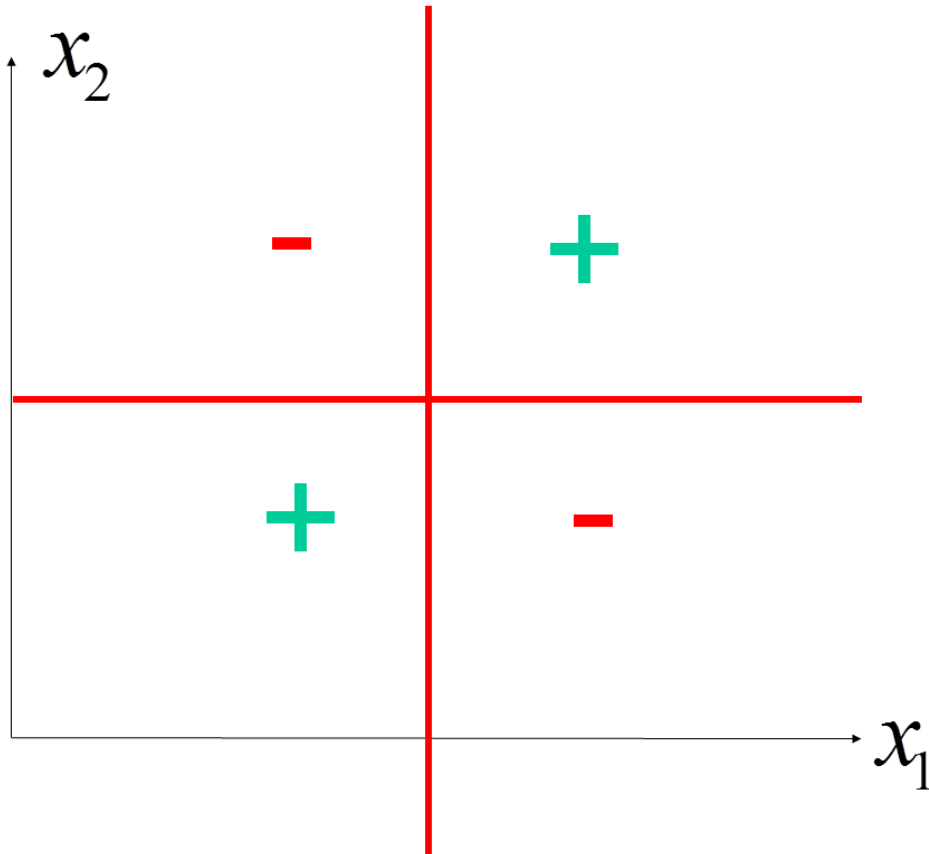
with  $\phi_1(x) = 1, \phi_2(x) = x_1, \phi_3(x) = x_2, \phi_4(x) = x_1x_2$



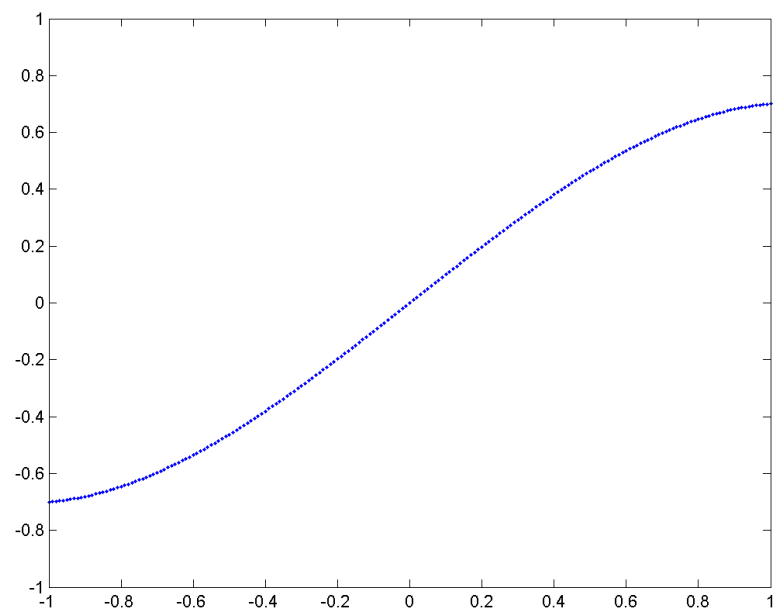
$$f(\mathbf{x}) = 1 - 2x_1 - 2x_2 + 4x_1x_2$$



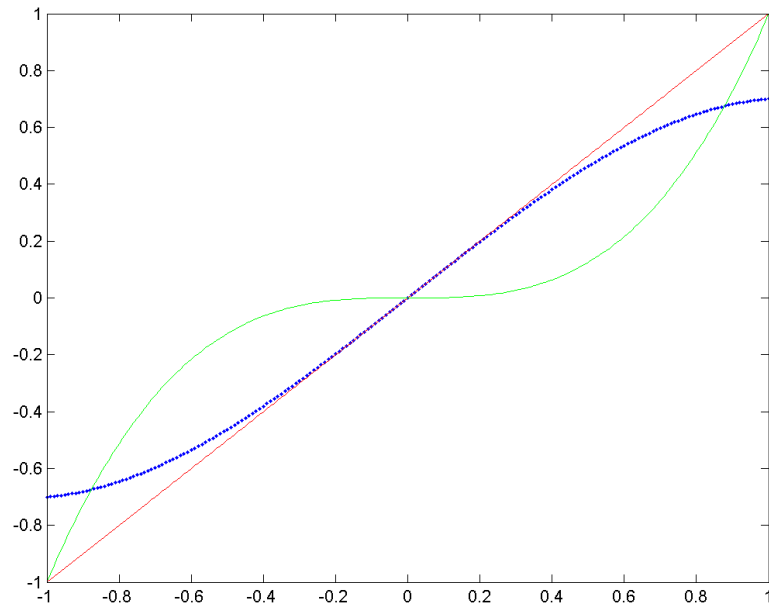
## Separating Planes



## A Nonlinear Function



$$f(x) = x - 0.3x^3$$



Basis functions  $\phi_1(x) = 1, \phi_2(x) = x, \phi_3(x) = x^2, \phi_4(x) = x^3$  und  $\mathbf{w} = (0, 1, 0, -0.3)$

## Basic Idea

- The simple idea: in addition to the original inputs, we add inputs that are calculated as deterministic functions of the existing inputs and treat them as additional inputs
- Example: Polynomial Basis Functions

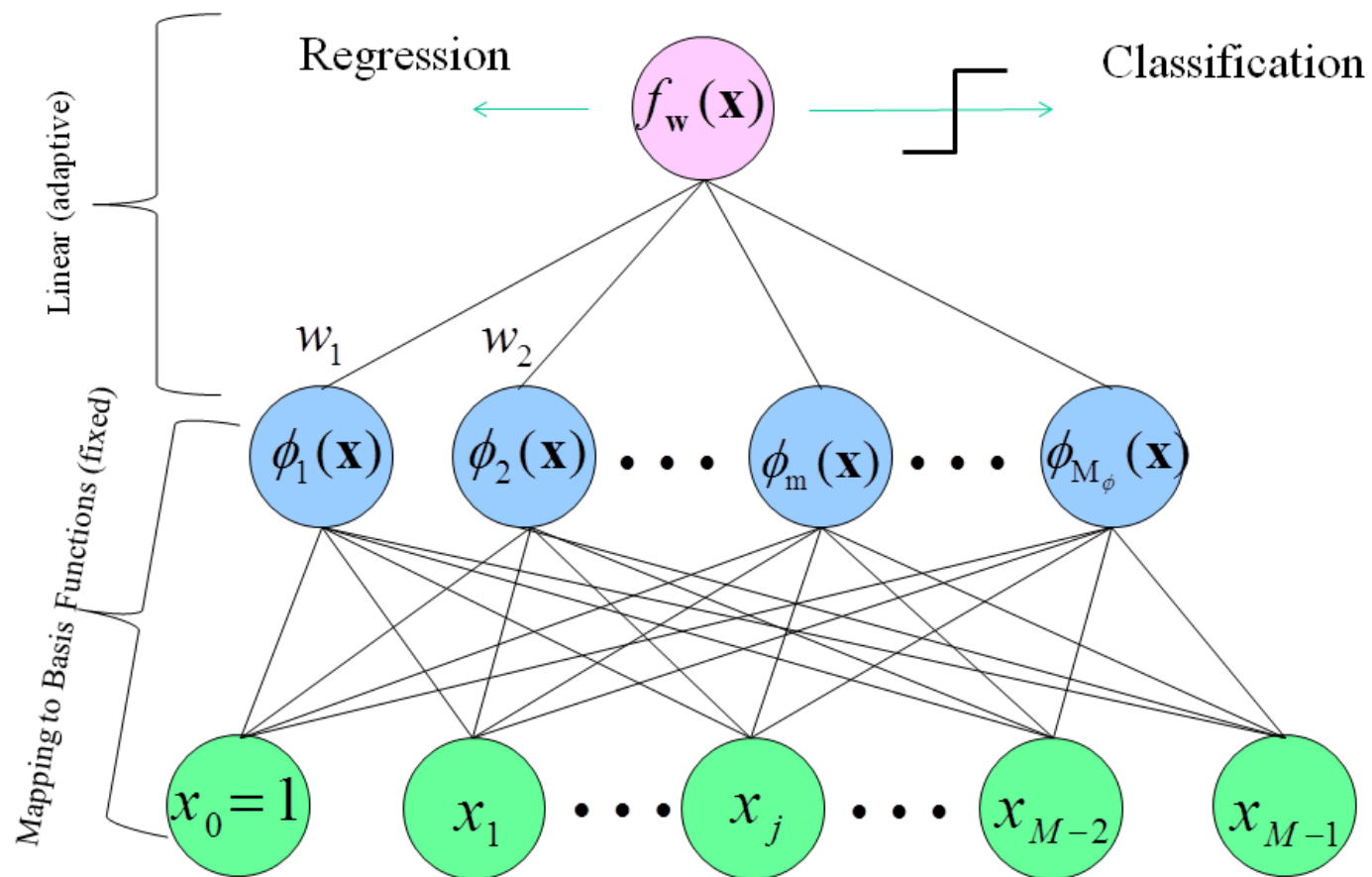
$$\{1, x_1, x_2, x_3, x_1x_2, x_1x_3, x_2x_3, x_1^2, x_2^2, x_3^2\}$$

- Basis functions  $\{\phi_m(\mathbf{x})\}_{m=1}^{M_\phi}$
- In the example:

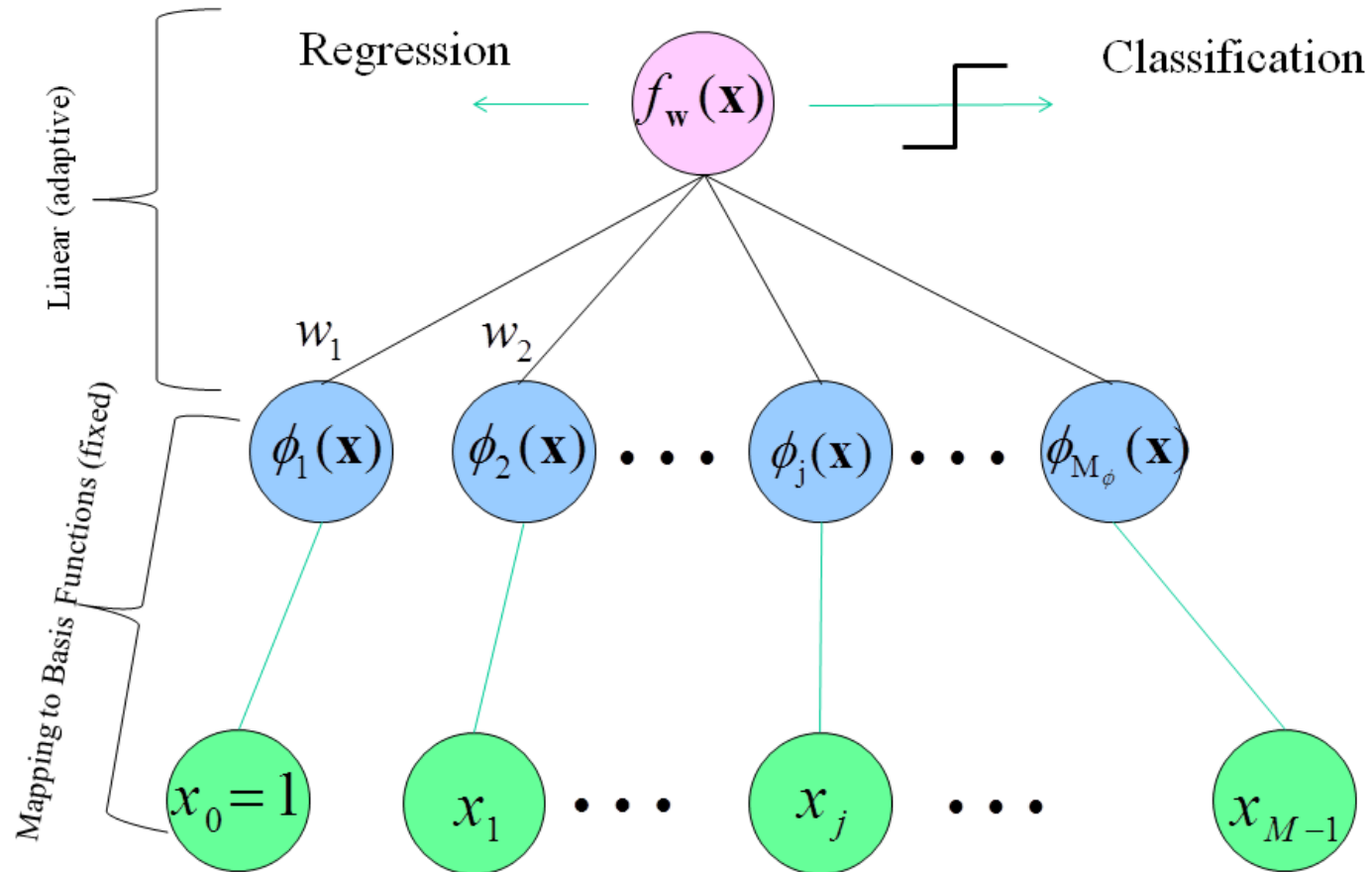
$$\phi_1(\mathbf{x}) = 1 \quad \phi_2(\mathbf{x}) = x_1 \quad \phi_6(\mathbf{x}) = x_1x_3 \quad \dots$$

- Independent of the choice of basis functions, the regression parameters are calculated using the well-known equations for linear regression

# Network of Basis Functions



## Network of Linear Basis Functions



## Review: Penalized LS for Linear Regression

- Multiple Linear Regression:

$$f_{\mathbf{w}}(\mathbf{x}) = w_0 + \sum_{j=1}^{M-1} w_j x_j = \mathbf{x}^T \mathbf{w}$$

- Regularized cost function

$$\text{cost}^{pen}(\mathbf{w}) = \sum_{i=1}^N (y_i - f_{\mathbf{w}}(\mathbf{x}_i))^2 + \lambda \sum_{i=0}^{M-1} w_i^2$$

- Die penalized LS-Solution

$$\hat{\mathbf{w}}_{pen} = \left( \mathbf{X}^T \mathbf{X} + \lambda I \right)^{-1} \mathbf{X}^T \mathbf{y} \quad \text{with} \quad \mathbf{X} = \begin{pmatrix} x_{1,0} & \dots & x_{1,M-1} \\ \dots & \dots & \dots \\ x_{N,0} & \dots & x_{N,M-1} \end{pmatrix}$$



## Regression with Basis Functions

- Model with basis functions:

$$f_{\mathbf{w}}(\mathbf{x}) = \sum_{m=1}^{M_{\phi}} w_m \phi_m(\mathbf{x})$$

- Regularized cost function with only basis function weights as free parameters (version 1)

$$\text{cost}^{pen}(\mathbf{w}) = \sum_{i=1}^N \left( y_i - \sum_{m=1}^{M_{\phi}} w_m \phi_m(\mathbf{x}_i) \right)^2 + \lambda \sum_{m=1}^{M_{\phi}} w_m^2$$

- The penalized least-squares solution

$$\hat{\mathbf{w}}_{pen} = \left( \Phi^T \Phi + \lambda I \right)^{-1} \Phi^T \mathbf{y}$$

with

$$\Phi = \begin{pmatrix} \phi_1(\mathbf{x}_1) & \dots & \phi_{M_\phi}(\mathbf{x}_1) \\ \dots & \dots & \dots \\ \phi_1(\mathbf{x}_N) & \dots & \phi_{M_\phi}(\mathbf{x}_N) \end{pmatrix}$$

## Nonlinear Models for Regression and Classification

- Regression:

$$f_{\mathbf{w}}(\mathbf{x}) = \sum_{m=1}^{M_{\phi}} w_m \phi_m(\mathbf{x})$$

As discussed, the weights can be calculated via penalized LS

- Classification:

$$\hat{y} = \text{sign}(f_{\mathbf{w}}(\mathbf{x})) = \text{sign} \left( \sum_{m=1}^{M_{\phi}} w_m \phi_m(\mathbf{x}) \right)$$

The Perceptron learning rules can be applied, if we replace  $1, x_{i,1}, x_{i,2}, \dots$  with  $\phi_1(\mathbf{x}_i), \phi_2(\mathbf{x}_i), \dots$

## Which Basis Functions?

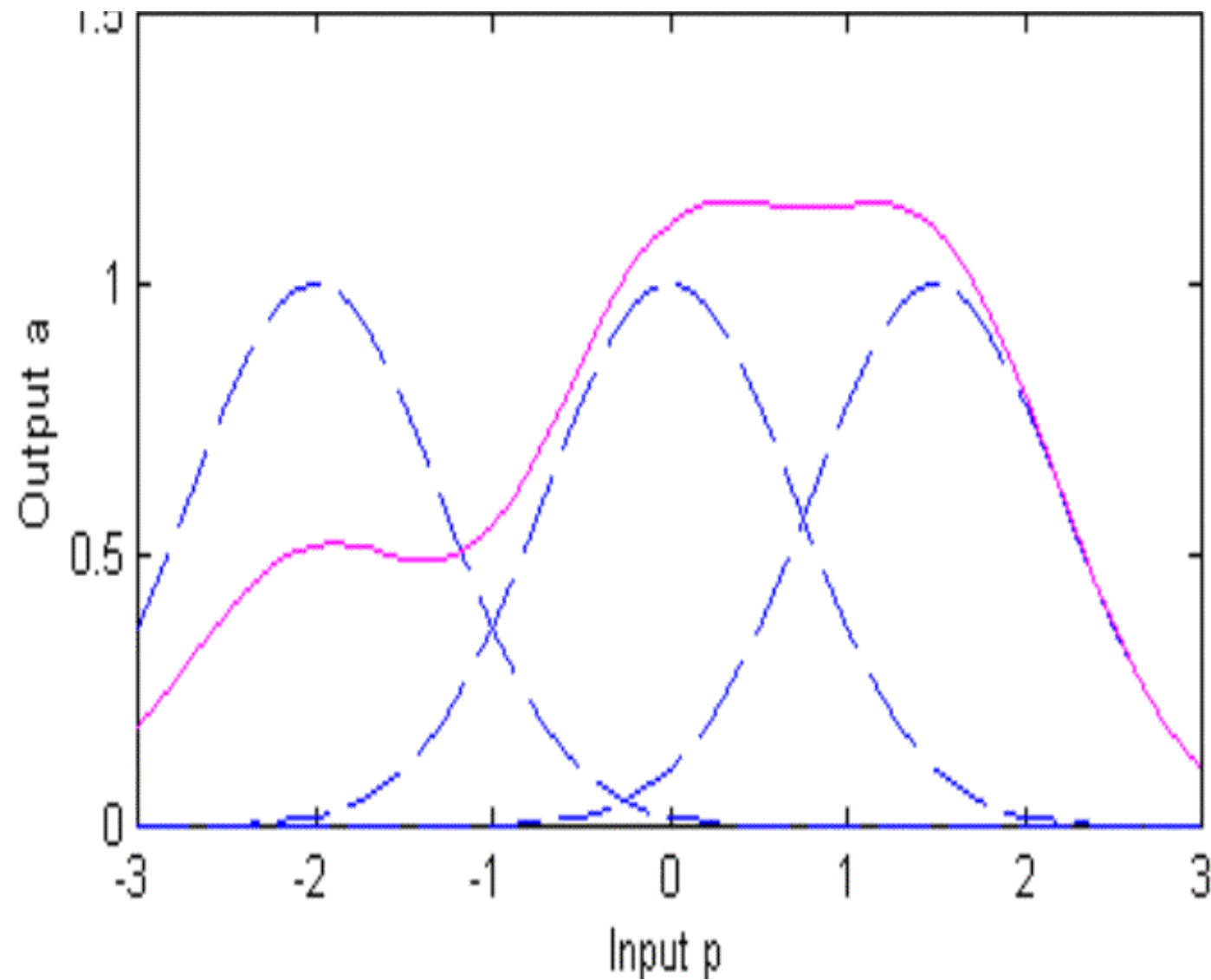
- The challenge is to find problem specific basis functions which are able to effectively model the true mapping, resp. that make the classes linearly separable; in other words we assume that the true dependency  $f(\mathbf{x})$  can be modelled by at least one of the functions  $f_{\mathbf{w}}(\mathbf{x})$  that can be represented by a linear combination of the basis functions, i.e., by one function in the function class under consideration
- If we include too few basis functions or unsuitable basis functions, we might not be able to model the true dependency
- If we include too many basis functions, we need many data points to fit all the unknown parameters (This sound very plausible, although we will see in the lecture on kernels that it is possible to work with an infinite number of basis functions)

## Radial Basis Function (RBF)

- We already have learned about polynomial basis functions
- Another class are radial basis functions (RBF). Typical representatives are Gaussian basis functions

$$\phi_j(\mathbf{x}) = \exp \left( -\frac{1}{2s_j^2} \|\mathbf{x} - \mathbf{c}_j\|^2 \right)$$

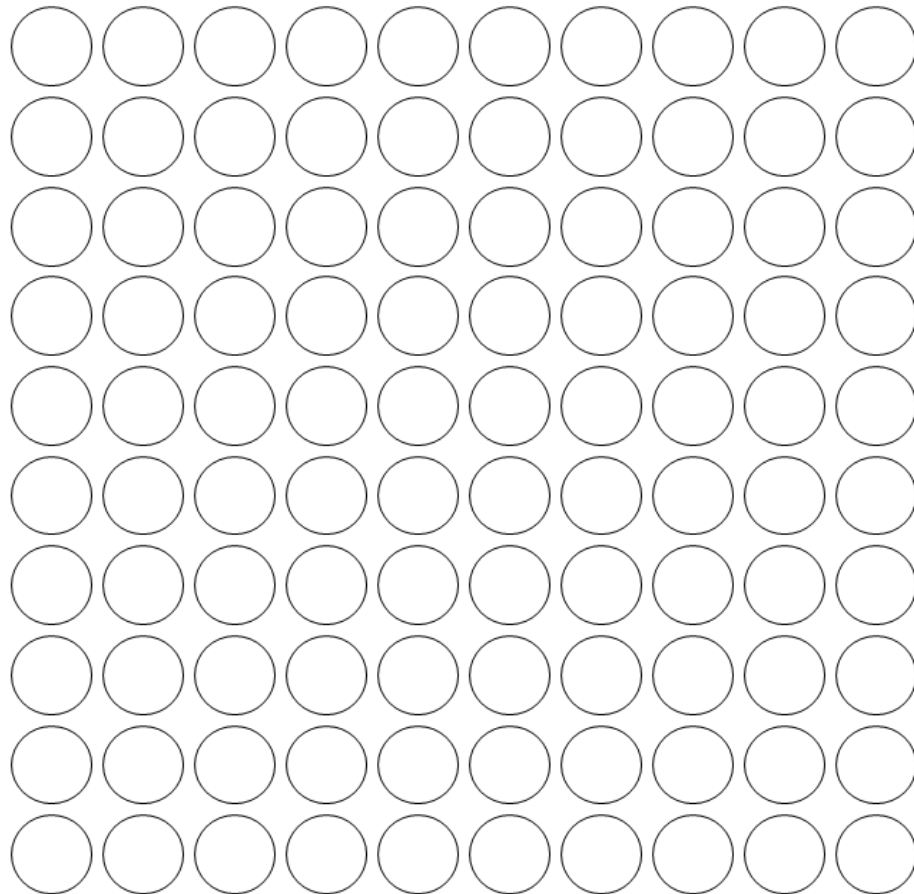
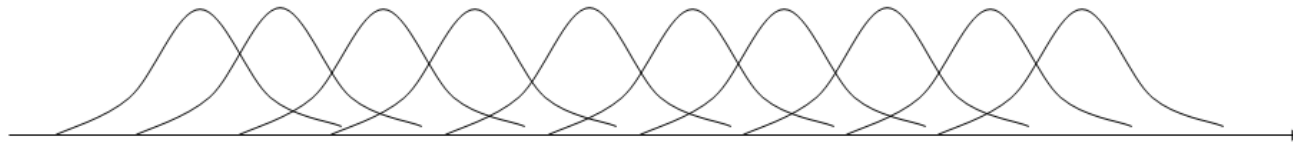
Three RBFs (blue) form  $f(x)$  (pink)



## Optimal Basis Functions

- So far all seems to be too simple
- Here is the catch: the number of “sensible” basis functions increases exponential with the number of inputs
- If I am willing to use  $K$  RBF-basis functions “per dimension”. then I need  $K^M$  RBFs in  $M$  dimensions
- We get a similar exponential increase for polynomial basis functions; the number of polynomial basis functions of a given degree increases quickly with the number of dimensions  $(x^2)$ ;  $(x^2, y^2, xy)$ ;  $(x^2, y^2, z^2, xy, xz, yz), \dots$
- *The most important challenge: How can I get a small number of relevant basis functions, i.e., a small number of basis functions that define a function class that contains the true function (true dependency)  $f(\mathbf{x})$ ?*

10 RBFs in one dimension



100 RBFs in  
two dimensions



## Strategy: Stepwise Increase of Model Class Complexity

- Start with a model class which is too simple and then incrementally add complexity
- First we only work with the original inputs and form a linear model
- Then we stepwise add basis functions that improve the model significantly
- For example we explore all quadratic basis functions. We include the quadratic basis function that mostly decreases the training cost; then we explore the remaining basis functions and, again, include the basis function that mostly decreases the training cost, and so on
- Examples: Polynomklassifikatoren (OCR, J. Schürmann, AEG)
  - Pixel-based image features (e.g., of hand written digits)
  - Dimensional reduction via PCA (see later lecture)
  - Start with a linear classifier and add polynomials that significantly increase performance
  - Apply a linear classifier

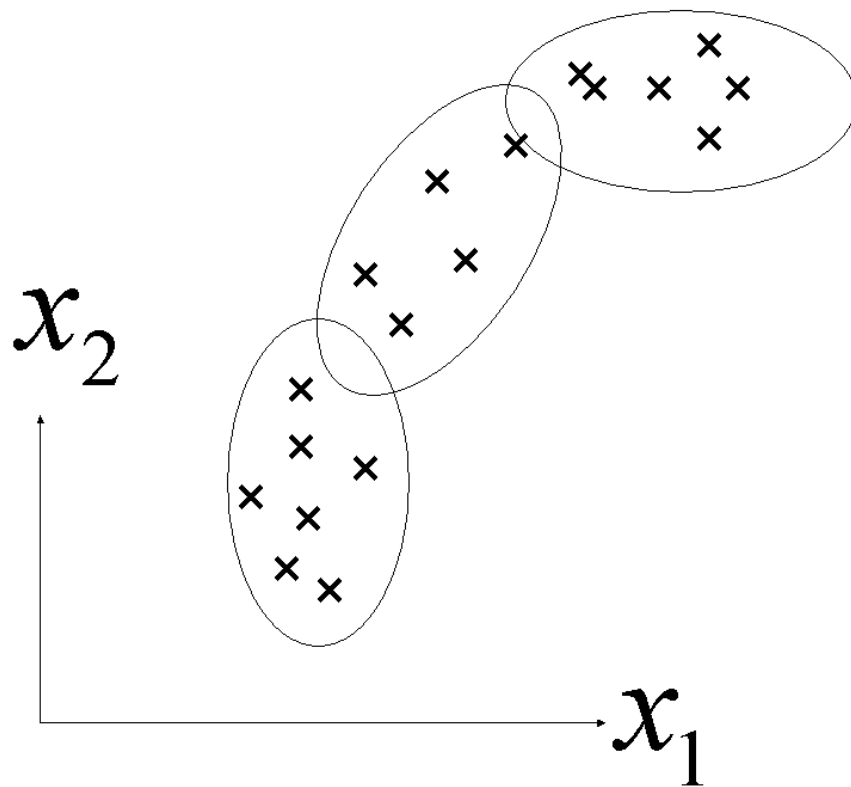
## Strategy: Stepwise Decrease of Model Class Complexity (Model Pruning)

- Start with a model class which is too complex and then incrementally decrease complexity
- First start with many basis functions
- Then we stepwise remove basis functions that increase the training cost the least
- A stepwise procedure is not optimal. Better: what is the best subset of  $K$  basis functions. Unfortunately, this problem is NP-hard

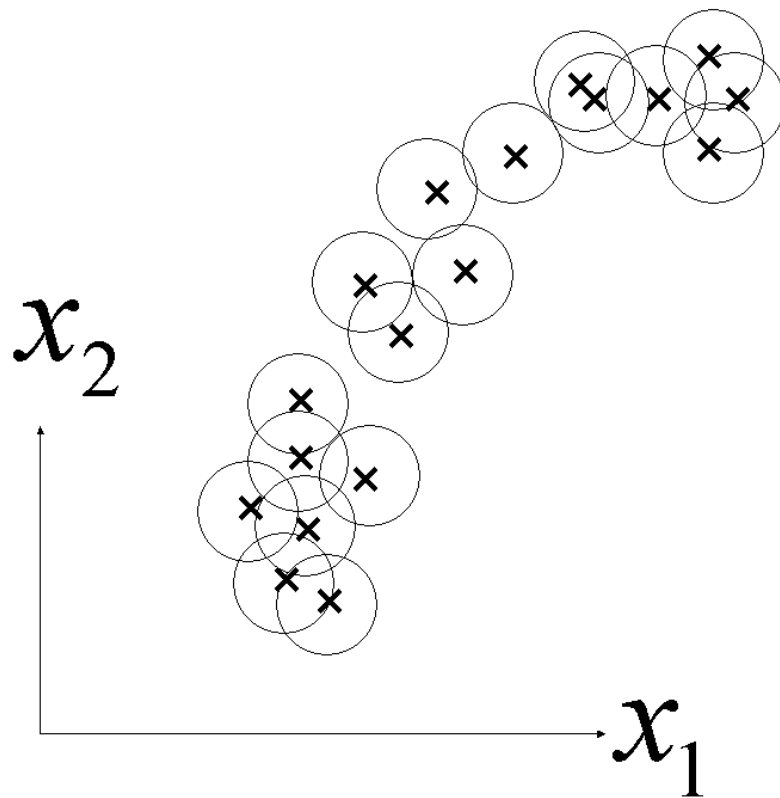
## Model Selection: RBFs

- Sometimes it is sensible to first group (cluster) data in input space and to then use the cluster centers as positions for the Gaussian basis functions
- The widths of the Gaussian basis functions might be derived from the variances of the data in the cluster
- An alternative is to use one RBF per data point. The centers of the RBFs are simply the data points themselves and the widths are determined via some heuristics (or via cross validation, see later lecture)

## RBFs via Clustering



## One Basis Function per Data Point



## Application-Specific Features

- Often the basis functions can be derived from sensible application features
  - Given an image with  $256 \times 256 = 65536$  pixels. The pixels form the input vector for a linear classifier. This representation would not work well for face recognition
  - With fewer than 100 appropriate features one can achieve very good results (example: PCA features, see later lecture)
- The definition of suitable features for documents, images, gene sequences, ... is a very active research area
- If the feature extraction already delivers many features, it is likely that a linear model will solve the problem and no additional basis functions need to be calculated
- This is quite remarkable: learning problems can become simpler in high-dimensions, in apparent contradiction to the famous “curse of dimensionality” (Bellman) (although there still is the other “curse of dimensionality” since the number of required basis functions might increase exponentially with the number of inputs! )

## Interpretation of Systems with Fixed Basis Functions

- The best way to think about models with fixed basis functions is that they implement a form of prior knowledge: we make the assumption that the true function can be modelled by the set of weighted basis function
- The data then favors certain members of the function class
- In the lecture on kernel systems we will see that the set of basis functions can be translated in assuming certain correlations between (mostly near-by) function values, implementing a smoothness prior

# Appendix: Detour on Function Spaces



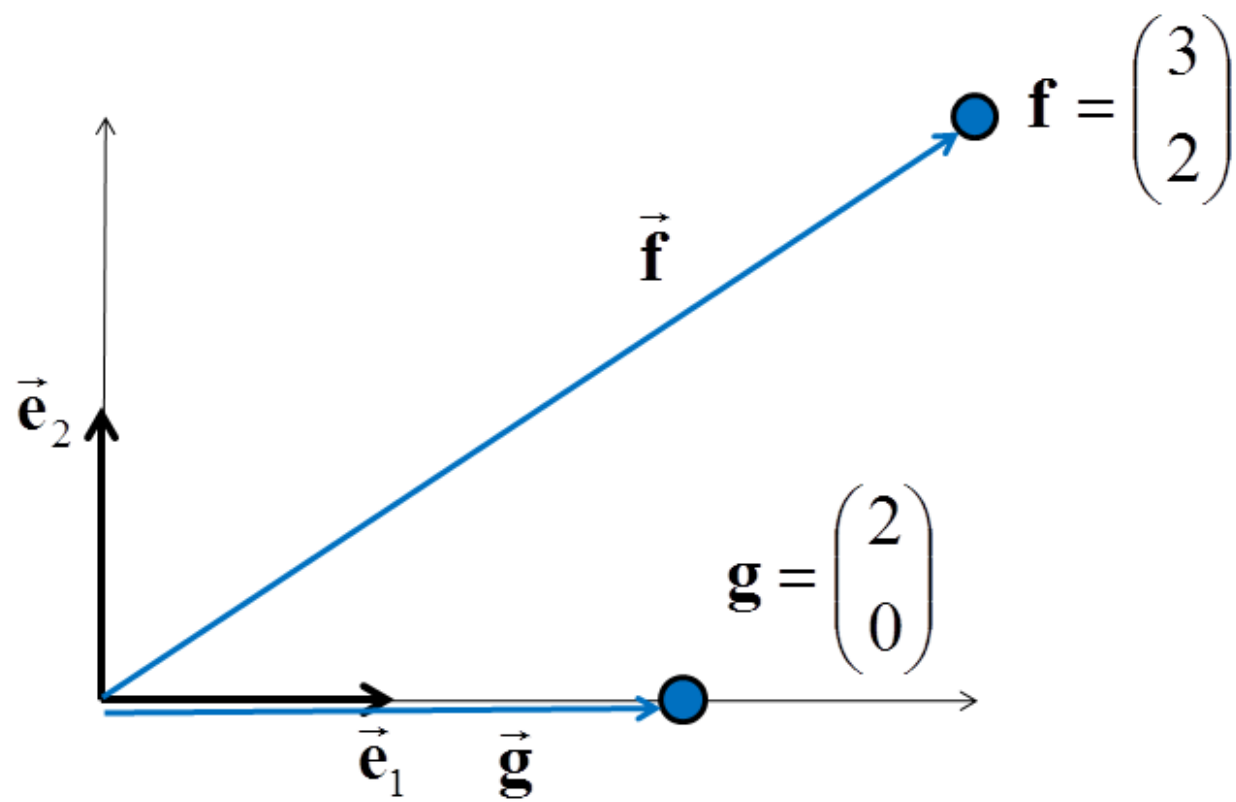
# Vectors

## Vectors

- To describe a vector  $\vec{\mathbf{f}}$  we need basis vectors  $\vec{\mathbf{e}}_i$  that define the orthogonal unit vectors in a coordinate system (i.e., the standard/natural/canonical coordinate system) and then we can write  $\vec{\mathbf{f}} = \sum_j f_j \vec{\mathbf{e}}_j$
- Orthonormality of basis vectors:  $\langle \vec{\mathbf{e}}_j, \vec{\mathbf{e}}_{j'} \rangle_{\mathbf{e}} = \delta_{j,j'}$  (Kronecker-Delta)
- The coordinates of a vector in a coordinate system are defined by the inner product of the vector with the basis vectors  $f_j = \langle \vec{\mathbf{e}}_j, \vec{\mathbf{f}} \rangle_{\mathbf{e}}$
- The inner product of two vectors is then

$$\langle \vec{\mathbf{f}}, \vec{\mathbf{g}} \rangle_{\mathbf{e}} = \sum_j \sum_{j'} f_j g_{j'} \langle \vec{\mathbf{e}}_j, \vec{\mathbf{e}}_{j'} \rangle_{\mathbf{e}} = \sum_j f_j g_j = \vec{\mathbf{f}} \cdot \vec{\mathbf{g}} = \mathbf{f}^T \mathbf{g}$$

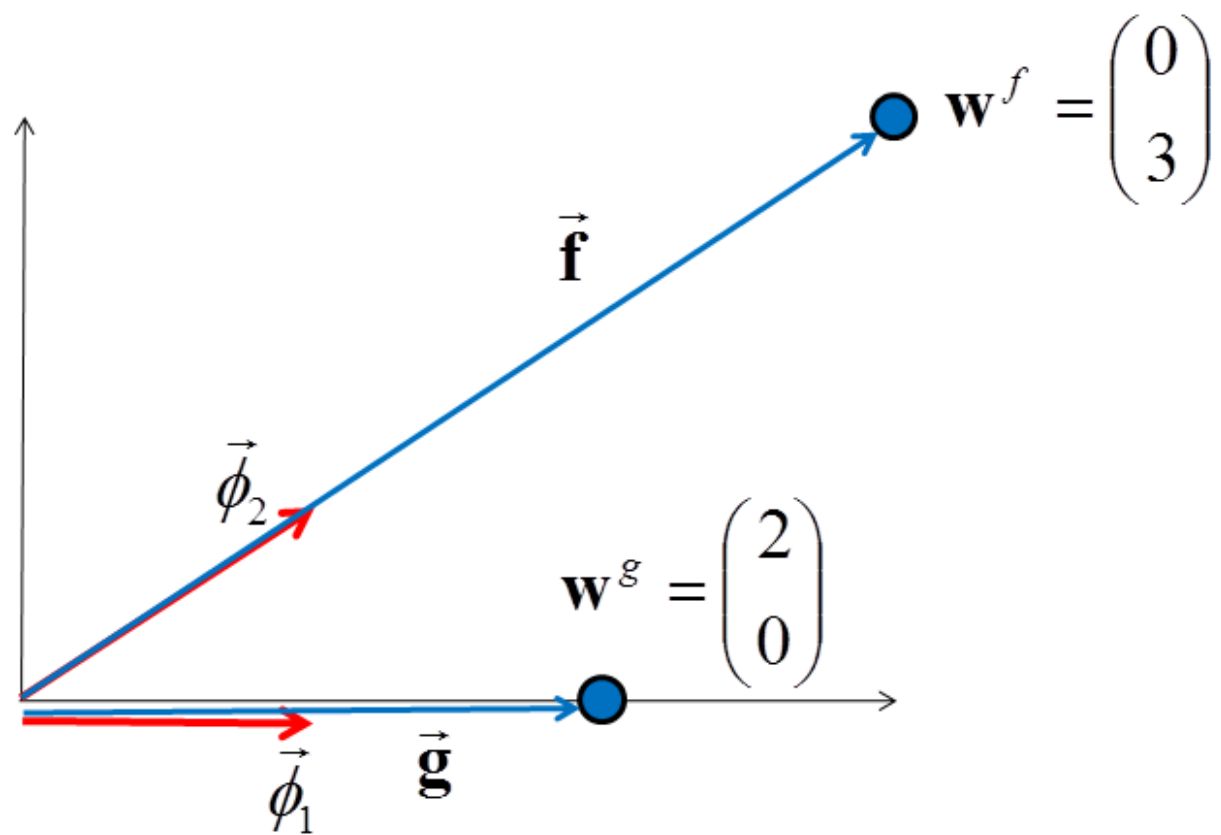
Thus the inner product as the dot product in the standard system.



$$\langle \vec{f}, \vec{g} \rangle_e = \mathbf{f}^T \mathbf{g} = 6$$

## A Vector in A Non-standard Coordinate System

- We consider new basis vectors  $\vec{\phi}_m$  that define the unit vectors in a second coordinate system, i.e., the non-standard system, and the coordinates of a vector  $w_m$ , and  $\vec{f} = \sum_m w_m \vec{\phi}_m$
- We now define **another inner product** by  $\langle \vec{\phi}_m, \vec{\phi}_{m'} \rangle_\phi = \delta_{m,m'}$ , i.e. as a dot product in the non-standard coordinate system
- The coordinates of a vector in a coordinate system are defined by the inner product of the vector with the basis vectors  $w_m = \langle \vec{\phi}_m, \vec{f} \rangle_\phi$
- The inner product of two vectors is then  $\langle \vec{f}, \vec{g} \rangle_\phi = \sum_m w_m^f w_m^g$ . This would be the dot product in the space where the  $\vec{\phi}_m$ -basis are orthonormal



$$\langle \vec{\mathbf{f}}, \vec{\mathbf{g}} \rangle_{\phi} = (\mathbf{w}^f)^T \mathbf{w}^g = 0$$

## Making the Link

- Typically we would know the representation of the non-standard basis in the standard system

$$\vec{\phi}_{m,j} = \left\langle \vec{e}_j, \vec{\phi}_m \right\rangle_{\mathbf{e}}$$

- Thus

$$f_j = \left\langle \vec{e}_j, \vec{\mathbf{f}} \right\rangle_{\vec{\mathbf{e}}} = \sum_m w_m \left\langle \vec{e}_j, \vec{\phi}_m \right\rangle_{\mathbf{e}} = \sum_m w_m \vec{\phi}_{m,j}$$

- Also,

$$\mathbf{f} = \Phi \mathbf{w} \quad \text{with} \quad (\Phi)_{m,j} = \phi_{m,j}$$

- ... and

$$\left\langle \vec{\mathbf{f}}, \vec{\mathbf{g}} \right\rangle_{\vec{\mathbf{e}}} = \sum_{m,m'} w_m^f w_{m'}^g \left\langle \vec{\phi}_m, \vec{\phi}_{m'} \right\rangle_{\mathbf{e}}$$

- Note, that in general, unless the non-standard set of basis vectors is also orthonormal in the standard system,  $\langle \vec{\mathbf{f}}, \vec{\mathbf{g}} \rangle_{\mathbf{e}} \neq \langle \vec{\mathbf{f}}, \vec{\mathbf{g}} \rangle_{\phi}$

## A Kernel Identity

Consider the *kernel* vector

$$\vec{\mathbf{k}}_{j'} = \sum_m \phi_{m,j'} \vec{\phi}_m \quad \text{such that} \quad \mathbf{k}_{j',j} = \sum_m \phi_{m,j'} \phi_{m,j}$$

Somewhat unusually, the weights are derived from the basis functions. There is one kernel vector for each dimension. We form the inner product with  $\vec{\mathbf{f}}$  and get

$$\langle \vec{\mathbf{f}}, \vec{\mathbf{k}}_{j'} \rangle_{\phi} = \sum_m w_m \phi_{m,j'} = f_{j'} = \langle \vec{\mathbf{f}}, \vec{\mathbf{e}}_{j'} \rangle_{\mathbf{e}}$$

Thus I can calculate  $f_{j'}$ , i.e., the  $j'$ -th component of  $\vec{\mathbf{f}}$  in the standard system, by calculating the dot product of  $\vec{\mathbf{f}}$  with  $\vec{\mathbf{k}}_{j'}$  in the non-standard system



## The Kernel Matrix and the Covariance Matrix link Inner Products

- The columns in the kernel matrix  $K$  contains the kernel vector coefficients and can be written as  $K = \Phi\Phi^T$  with  $(\Phi)_{m,j} = \phi_{m,j}$ ; in contrast, the covariance matrix is defined as  $C = \Phi^T\Phi$
- Starting with the inner product in standard space,

$$\langle \vec{f}, \vec{g} \rangle_e = (\Phi \mathbf{w}^f)^T \Phi \mathbf{w}^g = (\mathbf{w}_f)^T (C \mathbf{w}^g) = \langle \vec{f}, \vec{h} \rangle_\phi$$

and  $\vec{h}$  is a vector with coefficients  $\mathbf{w}^h = C \mathbf{w}^g$  in the non-standard space. Thus an inner (dot) product in the standard system can be written as an inner (dot) product in the non-standard system. Here, I need to know  $C$ .

- Starting with the inner product in the non-standard space, when the inverses exist,

$$\langle \vec{f}, \vec{g} \rangle_\phi = \left( (\Phi^{-1}) \mathbf{f} \right)^T \Phi^{-1} \mathbf{g} = \mathbf{f}^T \left( (\Phi \Phi^T)^{-1} \mathbf{g} \right) = \mathbf{f}^T K^{-1} \mathbf{g} = \langle \vec{f}, \vec{n} \rangle_e$$

where  $\vec{\mathbf{m}}$  is a vector with coefficients  $\mathbf{n} = K^{-1}\mathbf{g}$  in standard space. Thus an inner (dot) product in the non-standard system can be written as an inner (dot) product in the standard system; I need to know  $K$ .

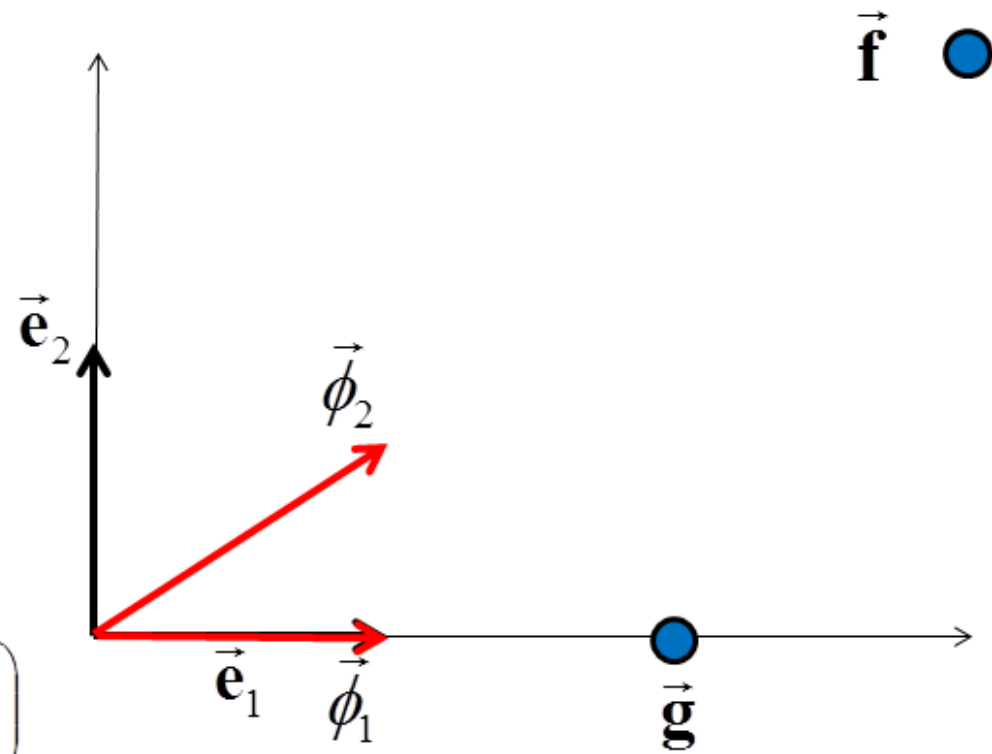
$$\vec{\mathbf{f}} = 3\vec{\mathbf{e}}_1 + 2\vec{\mathbf{e}}_2 = 0\vec{\phi}_1 + 3\vec{\phi}_2$$

$$\vec{\mathbf{g}} = 2\vec{\mathbf{e}}_1 + 0\vec{\mathbf{e}}_2 = 2\vec{\phi}_1 + 0\vec{\phi}_2$$

$$\langle \vec{\mathbf{f}}, \vec{\mathbf{g}} \rangle_e = 6 \quad \langle \vec{\mathbf{f}}, \vec{\mathbf{g}} \rangle_\phi = 0$$

$$\Phi = \begin{pmatrix} 1 & 1 \\ 0 & 2/3 \end{pmatrix} \quad C = \begin{pmatrix} 1 & 1 \\ 1 & 13/9 \end{pmatrix}$$

$$K = \begin{pmatrix} 2 & 2/3 \\ 2/3 & 4/9 \end{pmatrix} \quad K^{-1} = \begin{pmatrix} 1 & -1.5 \\ -1.5 & 4.5 \end{pmatrix}$$



$$\langle \vec{\mathbf{f}}, \vec{\mathbf{g}} \rangle_e = \langle \vec{\mathbf{f}}, \vec{\mathbf{h}} \rangle_\phi = (\mathbf{w}^f)^T C \mathbf{w}^g = \begin{pmatrix} 0 \\ 3 \end{pmatrix}^T \begin{pmatrix} 1 & 0 \\ 1 & 13/9 \end{pmatrix} \begin{pmatrix} 2 \\ 0 \end{pmatrix} = 6$$

$$\langle \vec{\mathbf{f}}, \vec{\mathbf{g}} \rangle_\phi = \langle \vec{\mathbf{f}}, \vec{\mathbf{m}} \rangle_e = \mathbf{f}^T K^{-1} \mathbf{g} = \begin{pmatrix} 3 \\ 2 \end{pmatrix}^T \begin{pmatrix} 1 & -1.5 \\ -1.5 & 4.5 \end{pmatrix} \begin{pmatrix} 2 \\ 0 \end{pmatrix} = 0$$

# Functions

## Functions are Vectors in a Standard Basis System

- Functions are just like vectors in a vector space
- Consider as basis functions in the standard space the Dirac delta functions

$$\delta(\mathbf{x} - \mathbf{x}') = \begin{cases} +\infty, & \mathbf{x} = \mathbf{x}' \\ 0, & \mathbf{x} \neq \mathbf{x}' \end{cases}$$

with

$$\int_{-\infty}^{\infty} \delta(\mathbf{x} - \mathbf{x}') d\mathbf{x} = 1.$$

- Then the inner product in which the delta-functions are orthonormal is defined as

$$\langle \delta(\mathbf{x} - \mathbf{x}_i), \delta(\mathbf{x} - \mathbf{x}_{i'}) \rangle_{\delta} = \delta_{\mathbf{x}_i, \mathbf{x}_{i'}}$$

- Furthermore we have

$$f(\mathbf{x}_i) = \langle \delta(\mathbf{x} - \mathbf{x}_i), f \rangle_{\delta} = \int f(\mathbf{x}) \delta(\mathbf{x} - \mathbf{x}_i) d\mathbf{x}$$

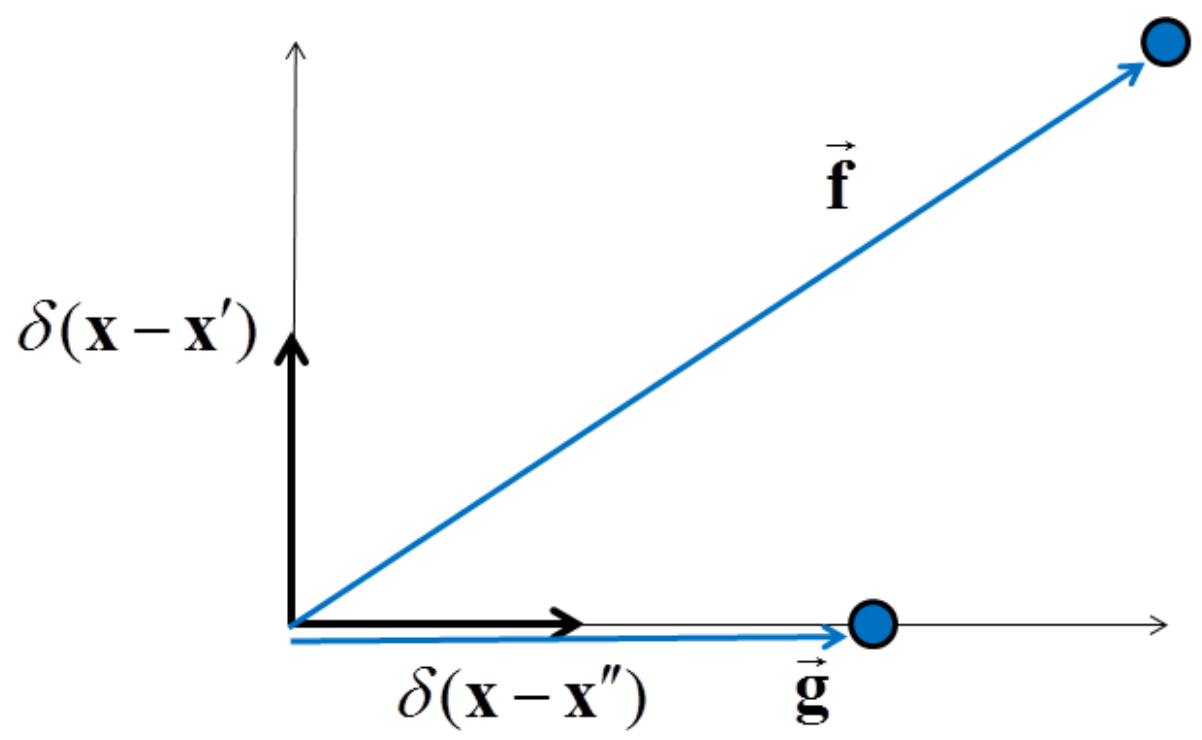
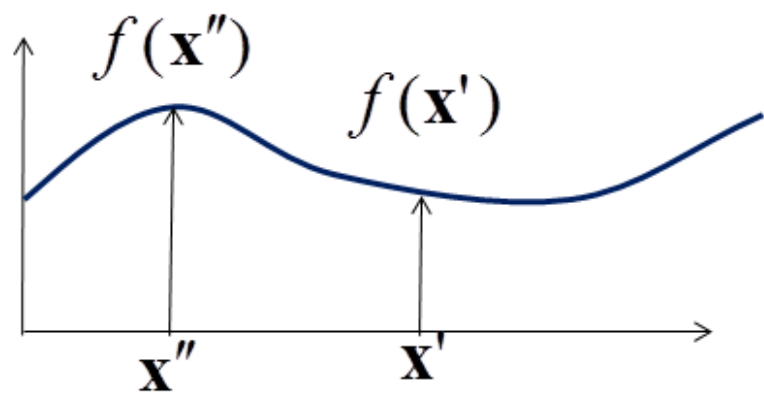
and

$$\vec{\mathbf{f}} = \int f(\mathbf{x}') \delta(\mathbf{x} - \mathbf{x}') d\mathbf{x}'$$

- In this coordinate system (if the integrals exists),

$$\langle \vec{\mathbf{f}}, \vec{\mathbf{g}} \rangle_{\delta} = \int f(\mathbf{x}) g(\mathbf{x}) d\mathbf{x}$$

Thus the inner product is a dot product in delta-function space



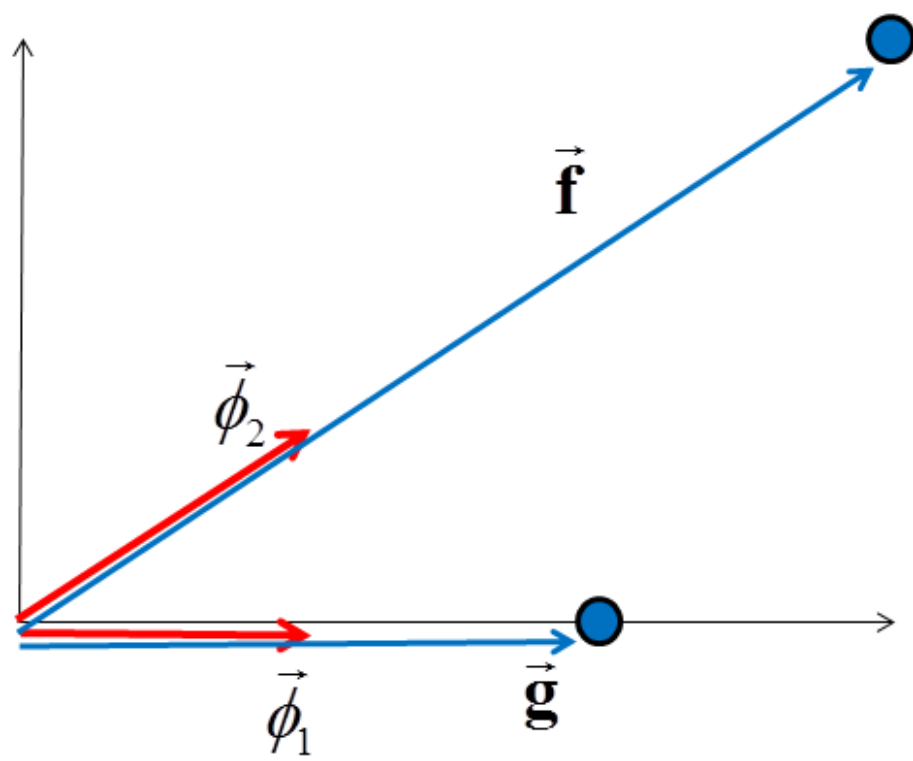
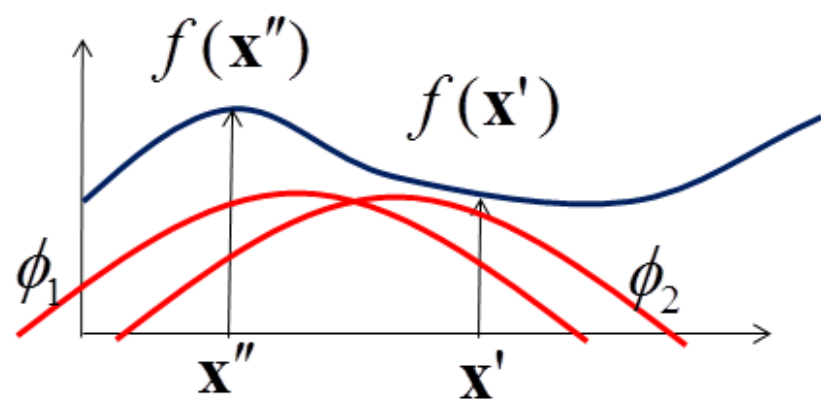
## Functions are Vectors in a Non-standard Basis Systems

- Now let's look at a function described in another coordinate system  $\vec{\mathbf{f}} = \sum_m w_m \vec{\phi}_m$
- We now define **another inner product** by  $\langle \vec{\phi}_m, \vec{\phi}_{m'} \rangle_\phi = \delta_{m,m'}$ : again the basis functions are orthonormal in the non-standard coordinate system
- The coordinates are simply the weights:  $\langle \vec{\phi}_m, \vec{\mathbf{f}} \rangle_\phi = w_m$
- In this coordinate system,

$$\langle \vec{\mathbf{f}}, \vec{\mathbf{g}} \rangle_\phi = \sum_m w_m^f w_m^g$$

Thus the inner product is a dot product in  $\phi$ -space





## Making the Link

- Typically we would know the representation of a basis function in the standard system

$$\phi_m(\mathbf{x}') = \left\langle \delta(\mathbf{x} - \mathbf{x}'), \vec{\phi}_m \right\rangle_\delta$$

- Thus

$$f(\mathbf{x}') = \left\langle \delta(\mathbf{x} - \mathbf{x}'), \vec{\mathbf{f}} \right\rangle_\delta = \sum_m w_m \left\langle \delta(\mathbf{x} - \mathbf{x}'), \vec{\phi}_m \right\rangle_\delta = \sum_m w_m \phi_m(\mathbf{x}')$$

$$\left\langle \vec{\mathbf{f}}, \vec{\mathbf{g}} \right\rangle_\delta = \sum_{m,m'} w_m^f w_{m'}^g \left\langle \vec{\phi}_m, \vec{\phi}_{m'} \right\rangle_\delta$$

- Note, that in general:  $\left\langle \vec{\mathbf{f}}, \vec{\mathbf{g}} \right\rangle_\delta \neq \left\langle \vec{\mathbf{f}}, \vec{\mathbf{g}} \right\rangle_\phi$ , unless the non-standard basis functions are orthonormal in the standard system as well, like Fourier basis functions

## Rewriting the Cost Function using Kernels

- Consider the function  $\vec{f} = \sum_m w_m \phi_m$  and the function

$$\mathbf{k}_{x'} = \sum_m \phi_m(\mathbf{x}') \vec{\phi}_m \quad \mathbf{k}_{x'}(x) = \sum_m \phi_m(\mathbf{x}') \phi_m(x)$$

As before, the weights of the non-standard vector are derived from the basis functions.  $\vec{\mathbf{k}}_{x'}$  is called a **kernel function** and is called the *reproducing kernel for the Hilbert space*. We form the inner product and get

$$\langle \vec{f}, \vec{\mathbf{k}}_{x'} \rangle_{\phi} = \sum_m w_m \phi_m(\mathbf{x}') = f(\mathbf{x}') = \langle \vec{f}, \delta(\mathbf{x} - \mathbf{x}') \rangle_{\delta}$$

- With all of this, we can write our cost function in terms of functions using only inner products in  $\phi$ -space as (version 2)

$$\text{cost}^{\text{pen}}(\vec{f}) = \sum_{i=1}^N \left( y_i - \langle \vec{f}, \vec{\mathbf{k}}_{x_i} \rangle_{\phi} \right)^2 + \lambda \langle \vec{f}, \vec{f} \rangle_{\phi}$$

## The Kernel Operator links Inner Products

- Similar to the vector case we can link the inner products via the kernels. The multiplications with matrices then correspond to operators (“infinite matrices”).

Let  $\mathcal{O}_K(\vec{g})(\mathbf{x})$  be an operator that transforms a function into another function defined by

$$\mathcal{O}_K(\vec{g})(\mathbf{x}) = \int k_{\mathbf{x}}(\mathbf{x}')g(\mathbf{x}')d\mathbf{x}'$$

- Starting with the inner product in the non-standard space, when the inverses exist,

$$\langle \vec{\mathbf{f}}, \vec{\mathbf{g}} \rangle_{\phi} = \langle \vec{\mathbf{f}}, \vec{\mathbf{n}} \rangle_{\delta}$$

with the function  $\vec{\mathbf{n}}$  defined in standard space as  $\vec{\mathbf{n}} = \mathcal{O}_K^{-1}(\vec{g})$

- With all of this, we can write our cost function in terms of functions using only inner products in  $\delta$ -space as (version 3)

$$\text{cost}^{pen}(\vec{\mathbf{f}}) = \sum_{i=1}^N \left( y_i - \left\langle \vec{\mathbf{f}}, \delta(\mathbf{x} - \mathbf{x}_i) \right\rangle_{\delta} \right)^2 + \lambda \left\langle \vec{\mathbf{f}}, \vec{\mathbf{n}} \right\rangle_{\delta}$$

## Dimensionality of the Function Space

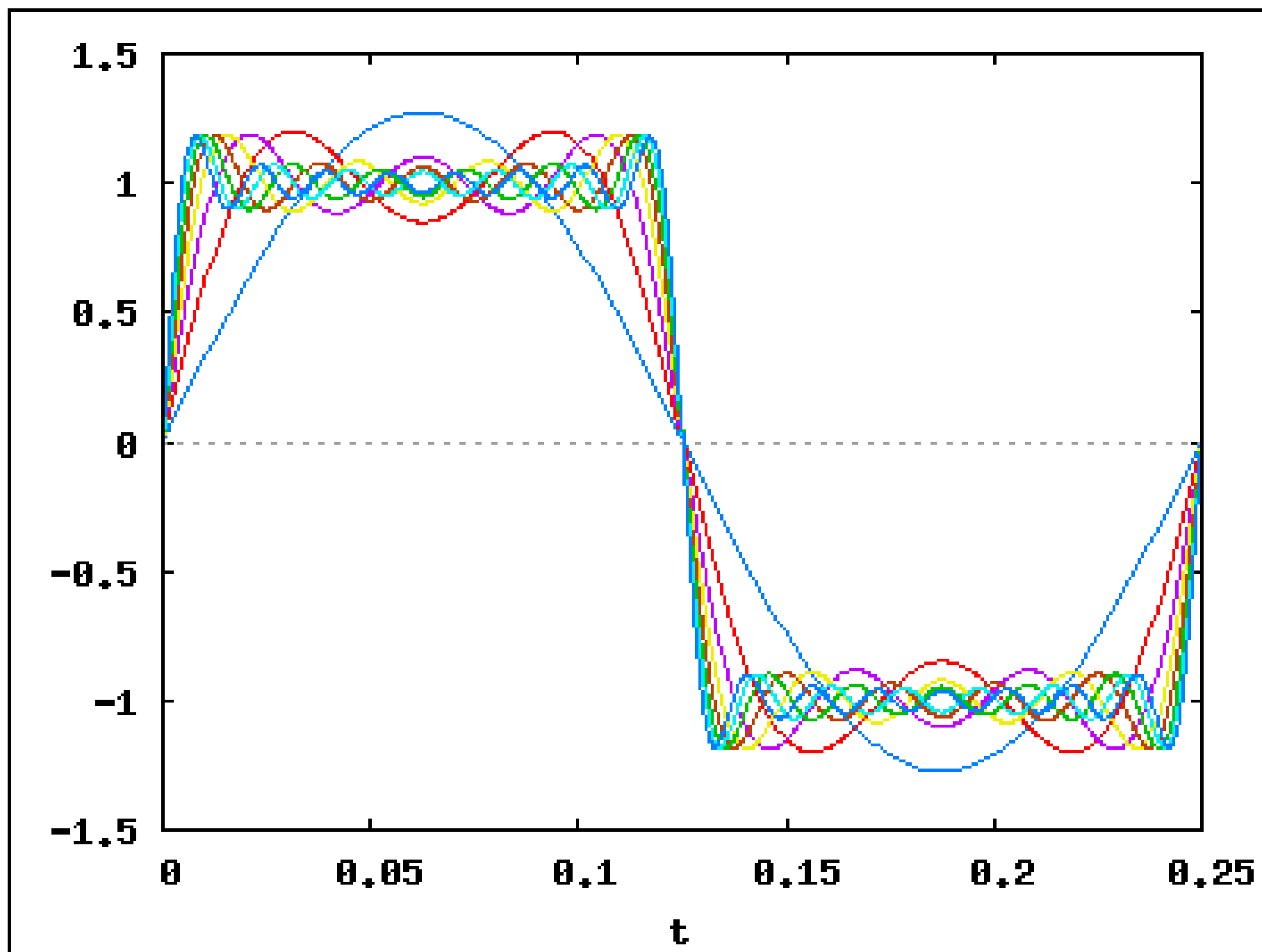
- In a kernel approach (see lecture on kernels) both the standard basis functions and the non-standard set of basis functions might describe the same **infinite** space
- When we work with fixed basis functions the standard space typically still has infinite dimensions ( $f(\mathbf{x})$  is defined for any  $\mathbf{x}$ ), but the non-standard set of basis functions describes an  $M_\Phi$ -dimensional subspace (assuming no degeneracy)
- Thus, not surprisingly, any function that can be described in the non-standard basis function system can be described in the standard system but not vice versa

## Example: Fourier Basis Functions

- A common set of basis functions (typically in 1-D or 2-D) are Fourier basis functions  $\phi_{c,\omega_i}(x) = \cos(\omega_i x)$ ,  $\phi_{s,\omega_i}(x) = \sin(\omega_i x)$  defined for periodic functions or functions defined in an interval
- Thus we can write  $f(x) = \sum_i w_{c,i} \cos(\omega_i x) + w_{s,i} \sin(\omega_i x)$ ; the weights  $w_{c,i}$  and the  $w_{s,i}$  are the Fourier coefficients
- The basis functions are orthogonal in the basis function space, but also in the standard space

$$\langle \vec{\phi}_{\omega_i}, \vec{\phi}_{\omega_j} \rangle_{\phi} = \langle \vec{\phi}_{\omega_i}, \vec{\phi}_{\omega_j} \rangle_{\delta} = \delta_{i,j}$$

This is the reason why the Fourier coefficients can be calculated very rapidly (Fourier transform, Fast Fourier Transform (FFT))





## Other Orthogonal Basis Functions

- You might have heard of the classical **orthogonal polynomials** (Jacobi polynomials, Laguerre polynomials, Hermite polynomials, and their special cases Gegenbauer polynomials, Chebyshev polynomials and Legendre polynomials)
- **Wavelets basis functions** are local both in location space and frequency space and are used in signal processing, image processing and compressed sensing. Some wavelets form a complete, orthonormal set of basis functions
- In general, orthogonal basis functions are only used in low-dimensional problems (signals, images) and are not often used in machine learning