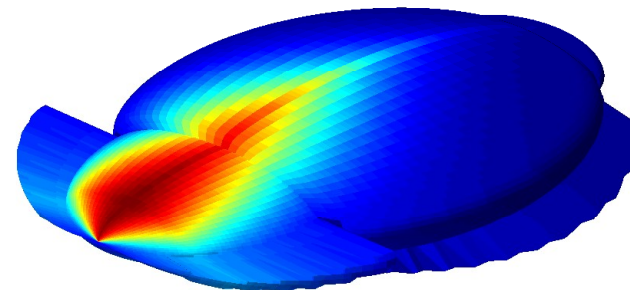
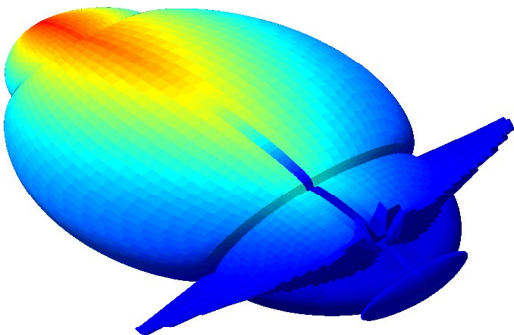


Lecture 9 Approximation



The Approximation Problem

Input values $\bar{x}^q \in X \subset \mathbb{R}^{n+1} \quad q = 1, \dots, M$

$$\bar{x}^q = \begin{bmatrix} x_0^q & x_1^q & x_2^q & \dots & x_n^q \end{bmatrix}^T$$

Output values $\bar{y}^q \in Y \subset \mathbb{R}^{l+1} \quad q = 1, \dots, M$

$$\bar{y}^q = \begin{bmatrix} y_1^q & y_2^q & \dots & y_l^q \end{bmatrix}^T$$

Training Set with real-valued inputs and outputs

Fit a model: $\hat{y} = f(\bar{x}; \bar{w})$

\bar{w} are the parameters of the model

to minimize the loss function

$$J(\bar{w}) = \sum_k J^q(\bar{w}) = \sum_k e(\hat{y}^q, \bar{y}^q; \bar{w})$$

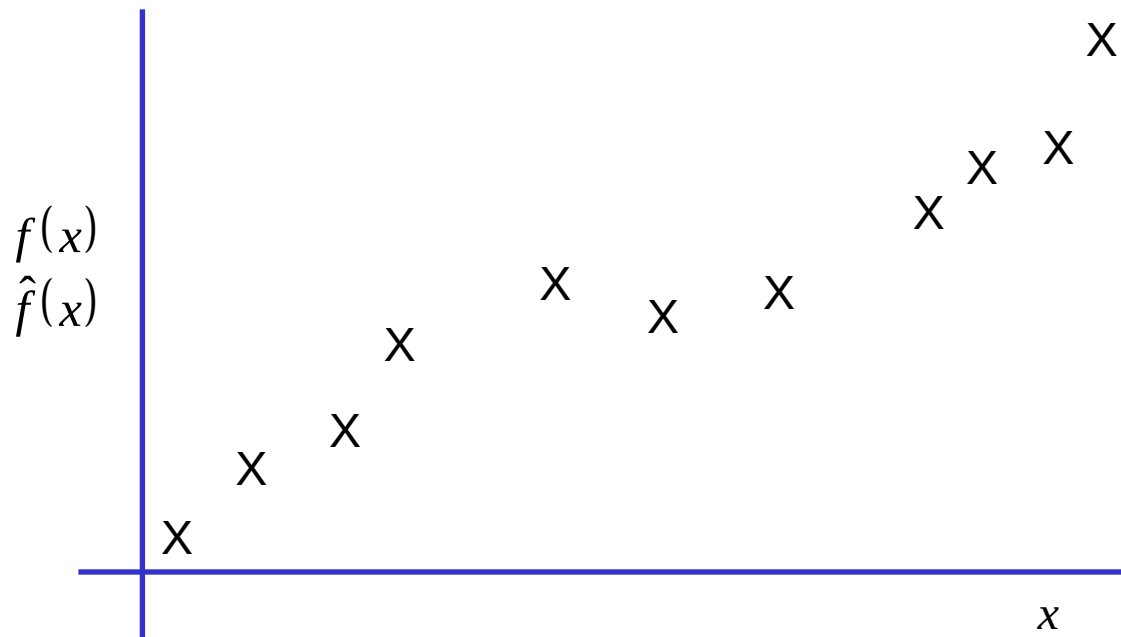
$e(\hat{y}, \bar{y}; w)$ is an error function

by adapting the parameters \bar{w}

Approximation = curve-fitting / regression

Approximation = curve fitting using a finite data set

interpolation → generalization
extrapolation →



x = available (sample)
data points

The Problem with Approximation

A finite number of points are given in real space

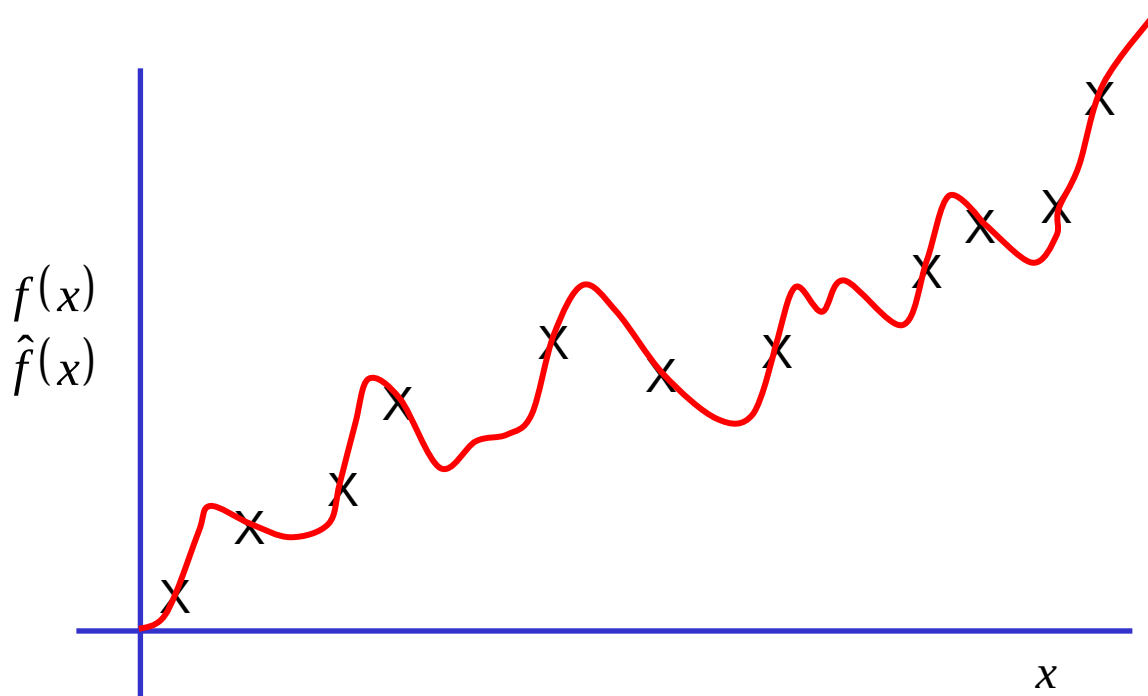
⇒ An infinite number of functions are compatible with it.

How should we choose the best?

Question to ask: Which function is most likely to be correct for unknown data?

Approximation = curve fitting using a finite data set

interpolation →
extrapolation → generalization



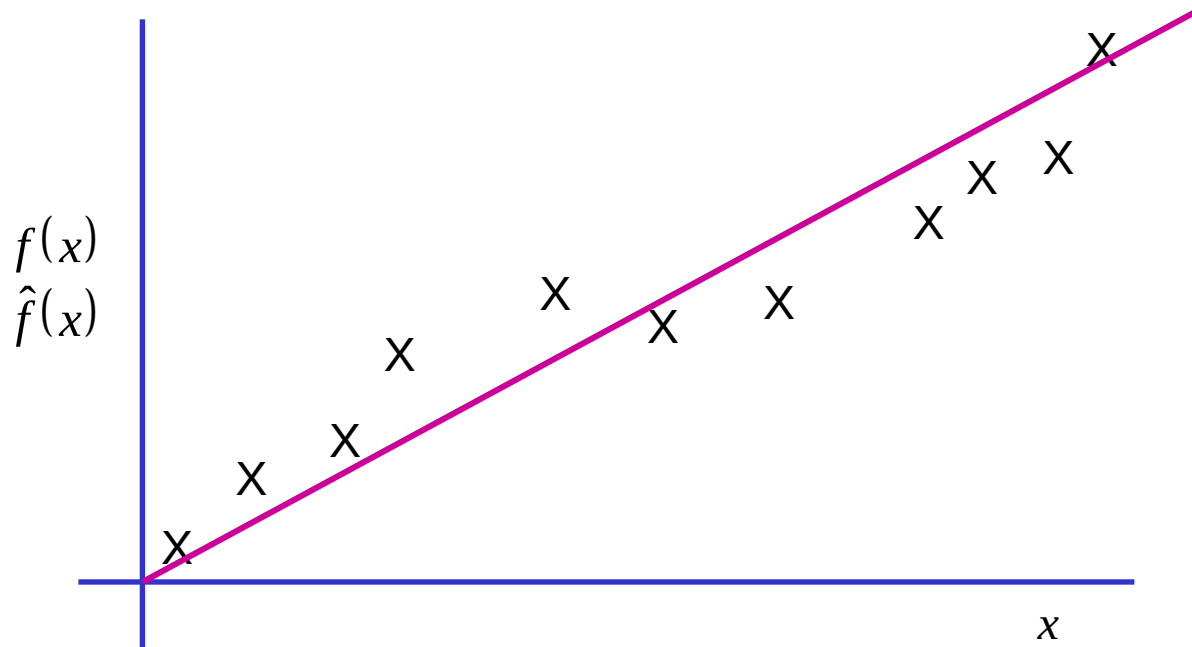
x = available (sample)
data points

Is this a good
approximation?

Why not?

Approximation = curve fitting using a finite data set

interpolation →
extrapolation → generalization



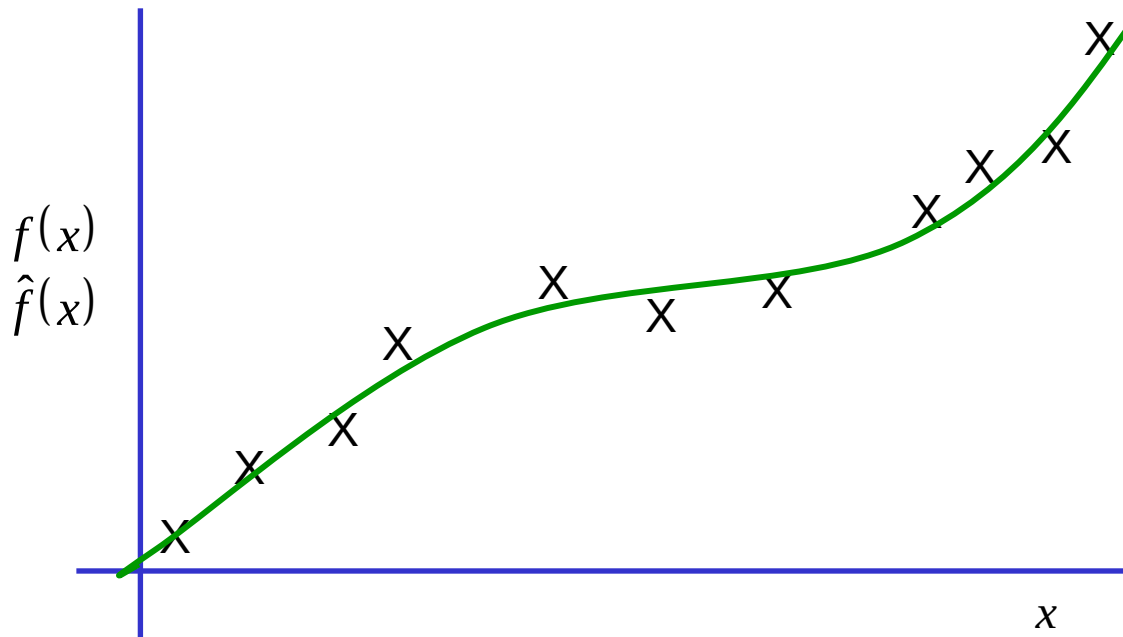
x = available (sample)
data points

What about this
one?

Better?

Approximation = curve fitting using a finite data set

interpolation →
extrapolation → generalization



x = available (sample)
data points

And this one?

Probably the best!

The Problem with Approximation

A finite number of points are given in real space
⇒ An infinite number of functions are compatible with it.

How should we choose the best?

Question to ask: Which function is most likely to be correct for unknown data?

Answer:

“One should not increase, beyond what is necessary, the number of entities required to explain anything” William of Ockham → Occam’s Razor

“Everything should be made as simple as possible, but no simpler ”
Attributed to Albert Einstein

The simplest adequate model is likelier to be correct.

Why are feedforward networks approximators?

Consider a 1 hidden-layer net with n inputs and one output.
The network can be written as:

$$\hat{y}_i = f_i \left[\sum_{j=0}^m w_{ij} f_j \left(\sum_{k=0}^n w_{jk} x_k \right) \right] \quad \text{where } i \text{ is the only output neuron}$$

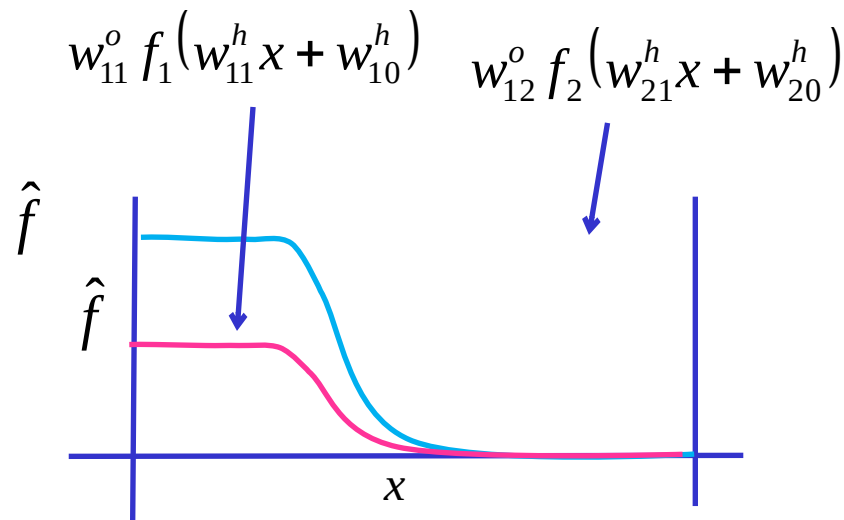
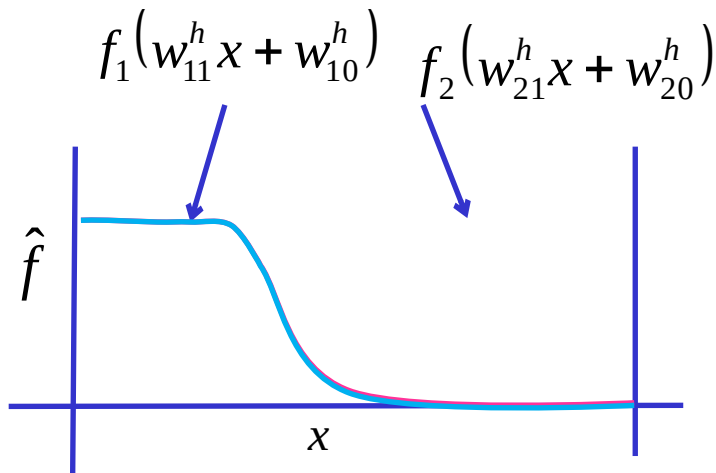
Assume $f_i(u) = u$ (linear output) $\Rightarrow \hat{y}_i = \sum_j w_{ij} f_j \left(\sum_k w_{jk} x_k \right)$

Thus $f_j(u)$ serve as **basis functions**

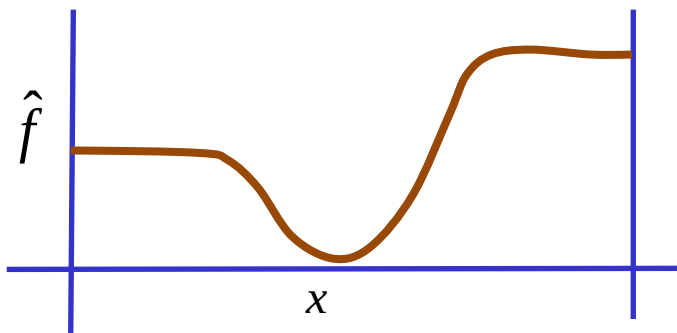
i.e. a weighted sum of several nonlinear functions, f_j , defined over the input space

Recall the sinusoidal basis $f(x) = \sum_i a_i \cos(w_i x + \theta_i)$ = Fourier series approximation

Let $n = 1, m = 2$

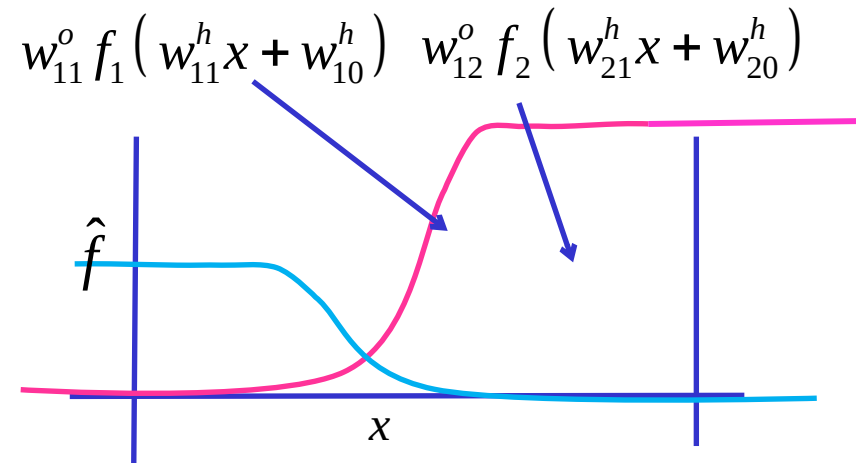
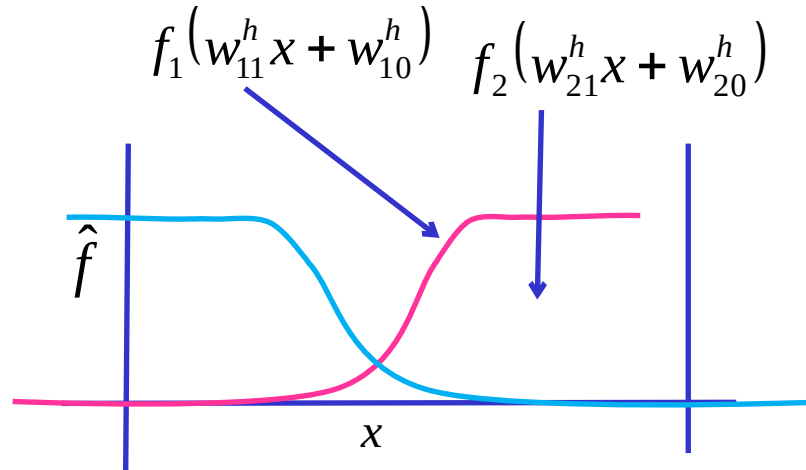


$$w_{11}^o f_1(w_{11}^h x + w_{10}^h) + w_{12}^o f_2(w_{21}^h x + w_{20}^h)$$

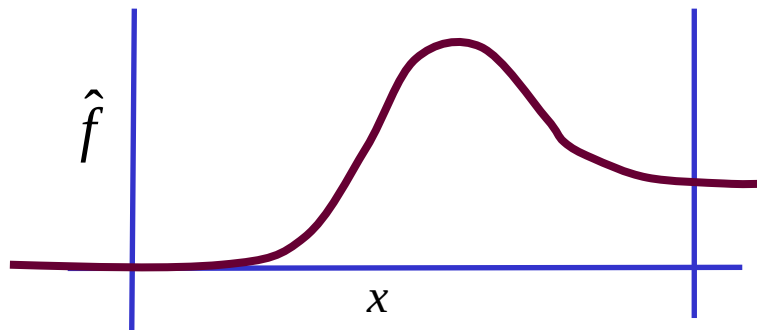


By adding several scaled and possibly reversed sigmoids, we can form pretty much any function. This remains true even if $n > 1$

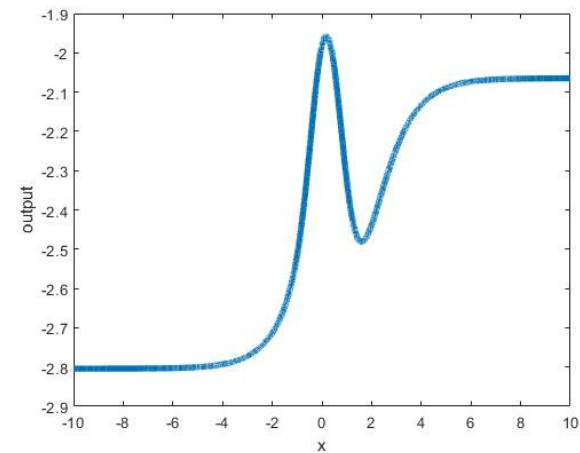
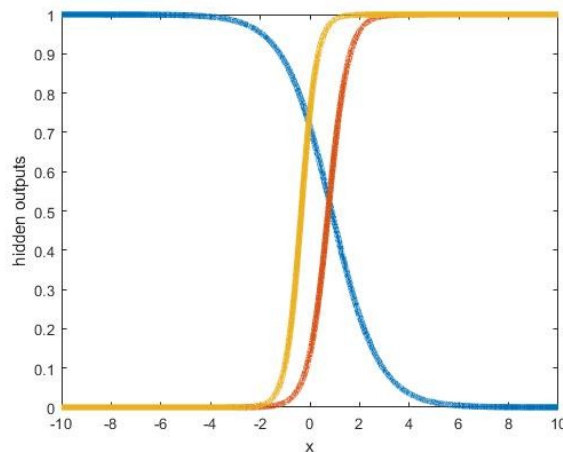
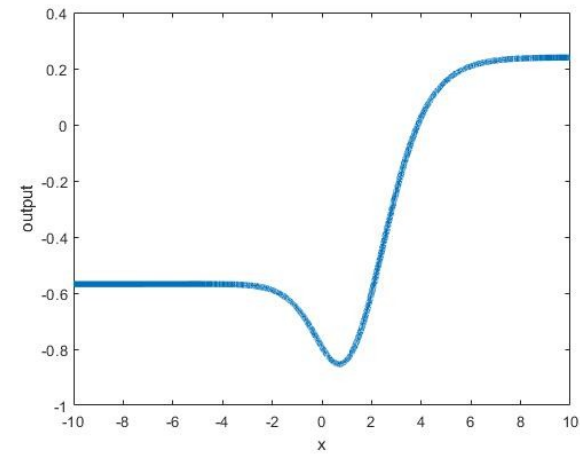
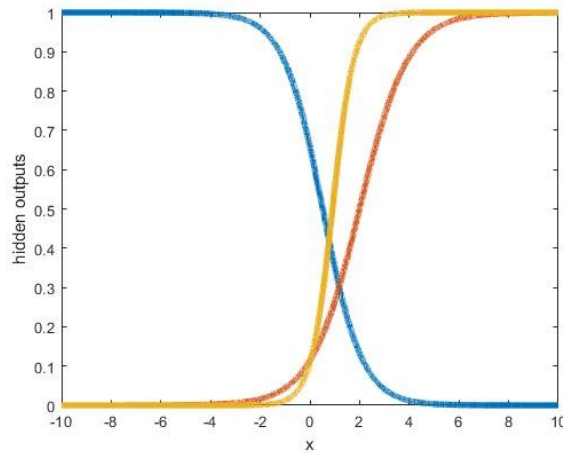
Let $n = 1, m = 2$



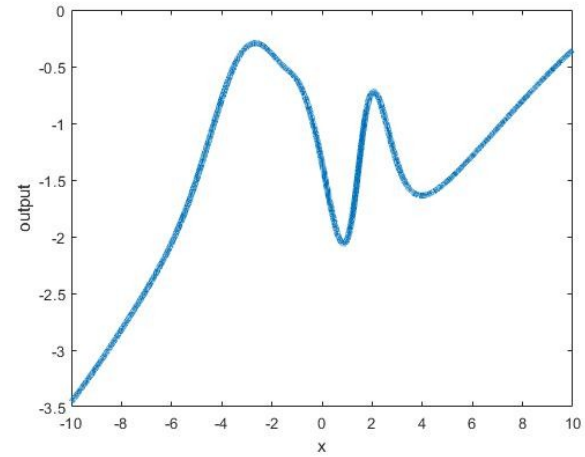
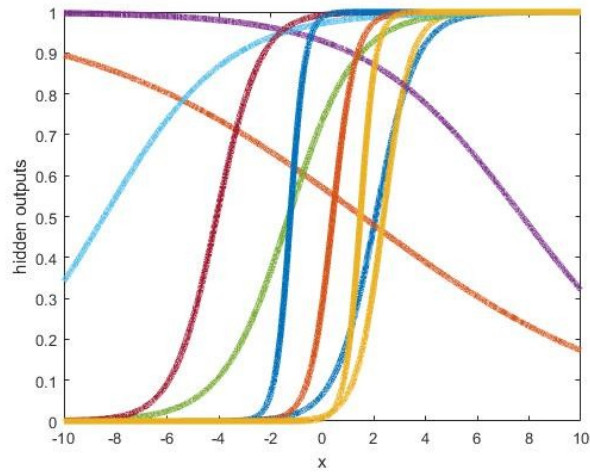
$$w_{11}^o f_1(w_{11}^h x + w_{10}^h) - w_{12}^o f_2(w_{21}^h x + w_{20}^h)$$



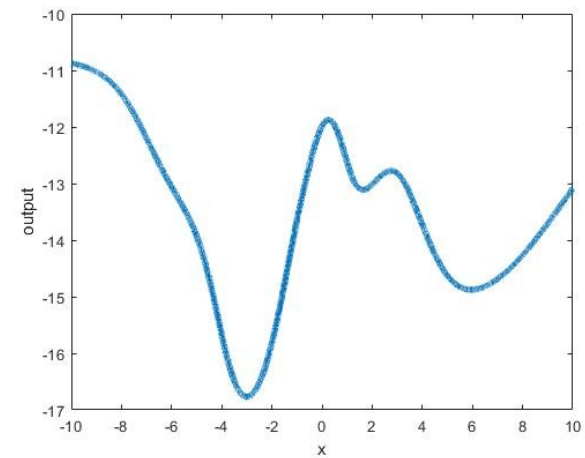
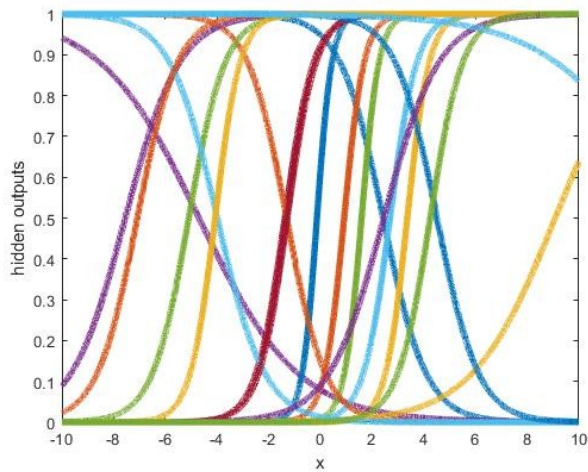
Two Real Examples with 3 Hidden Neurons



10
Hidden



20
Hidden



The Universal Approximation Theorem

(Cybenko; Funahashi; Hornik, Stinchcombe & White)

Let $f(\cdot)$ be a non-constant, bounded and monotonically increasing continuous function.

Let $I_n \equiv [0,1]^n$

Let $C(I_n)$ be the space of continuous functions on I_n

Then, given any $\Phi \in C(I_n)$ and $\varepsilon > 0$

\exists integer m and real constants a_j, b_j, w_{jk}
 $j = 1, \dots, m$
 $k = 1, \dots, n$

such that we can find:

$$\hat{\Phi}(x_1, \dots, x_n) = \sum_{j=1}^m a_j f\left(\sum_{k=1}^n w_{jk} x_k + b_j\right)$$

i.e. any cont. function on I_n can be approximated to arbitrary precision using one hidden layer of sigmoid neurons

and $\left| \hat{\Phi}(x_1, \dots, x_n) - \Phi(x_1, \dots, x_n) \right| < \varepsilon \quad \forall x_1, \dots, x_n \in I_n$

Applications for Approximation

Modeling:

Given data points from an input-output system

$$\left\{ (x^k, y^k) \right\}$$

determine $F : x^k \rightarrow y^k$

e.g.

$F(\text{heart rate, blood pressure}) = \text{autonomic nervous system function}$

$F(\text{externally monitored plant variables}) = \text{internal variables.}$

or, for a dynamical system

$$\dot{x} = F(x)$$

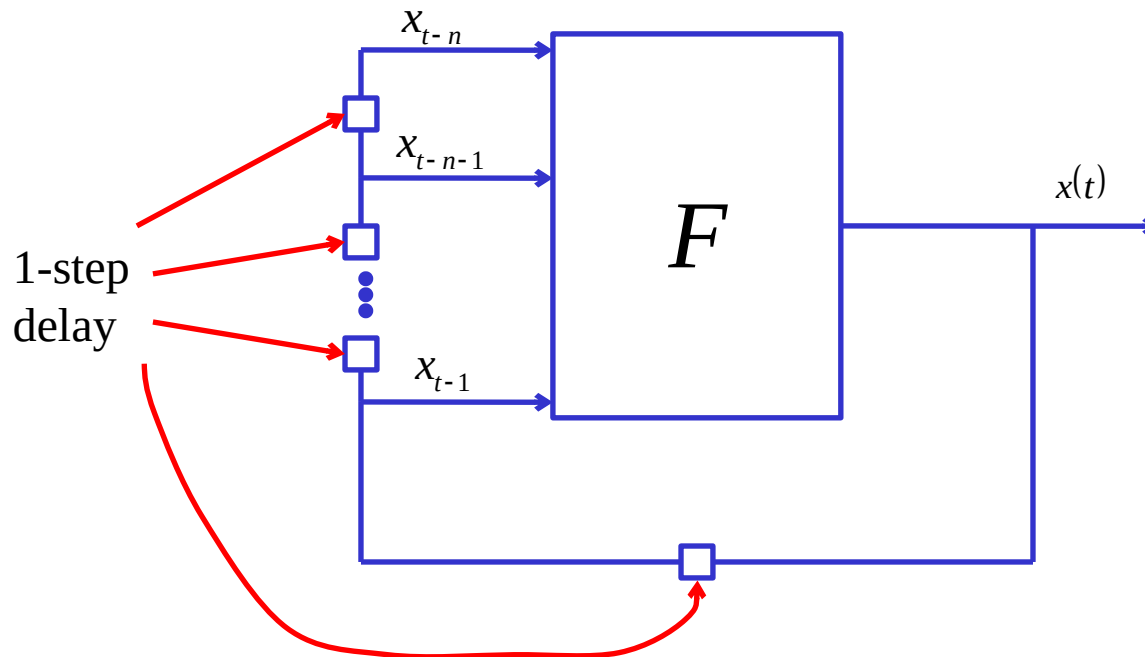
use data points $\left\{ (x^k, \dot{x}^k) \right\}$ to find $\hat{F}(x) = \dot{x}$

Prediction:

Given a time series x_t

hypothesize $x_t = F(x_{t-1}, \dots, x_{t-m}; \mathbf{w})$

Then use approximation to find F .
This works also when the system
has external inputs



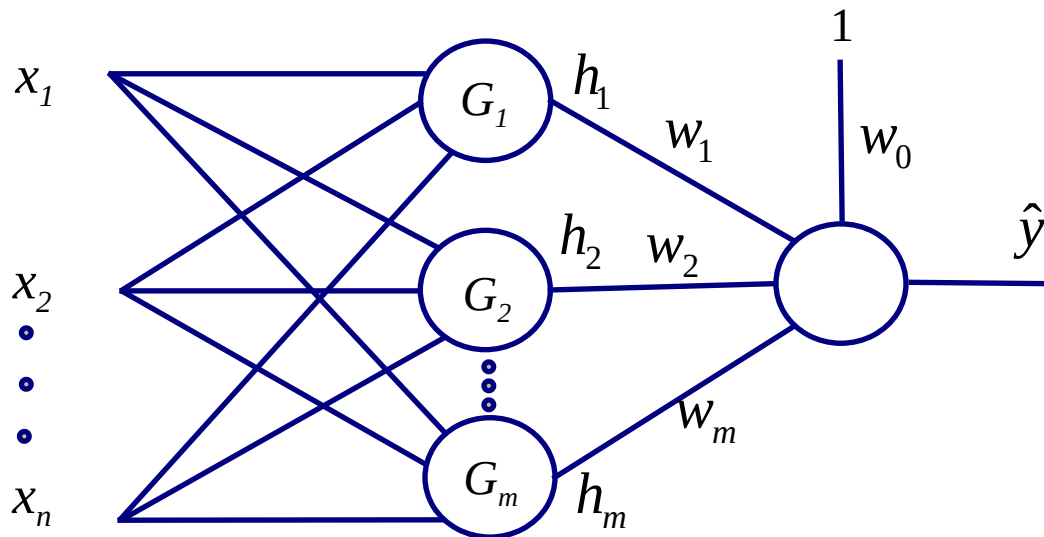
Compare with a linear
autoregressive model:

$$x_t = \sum_{k=1}^n \alpha_i x_{t-k}$$

$$x_t = F(x_{t-1}, \dots, x_{t-m}; \mathbf{w}) = \sum_{j=1}^m w_{ij}^o f \left(\sum_{k=1}^n w_{jk}^h x_{t-k} + w_{j0}^h \right)$$

Radial Basis Function Networks

Radial Basis Function (RBF) Networks



$$\bar{x} = [x_1 \quad x_2 \quad . \quad . \quad . \quad x_n]^T \quad n\text{-dim}$$

$$\bar{w} = [w_0 \quad w_1 \quad . \quad . \quad . \quad w_m]^T \quad (m+1)\text{-dim}$$

$$\bar{h} = [h_0 \quad h_1 \quad . \quad . \quad . \quad h_m]^T \quad (m+1)\text{-dim}$$

$$\bar{C} = [C_1 \quad C_2 \quad . \quad . \quad . \quad C_n] \quad n\text{-dim}$$

$$\|z\| = \sqrt{\sum_k z_k^2}$$

$$h_j = G_j \left(\|\bar{x} - \bar{C}_j\|^2 \right)$$

where

$$G_j(u) = e^{\frac{-u}{2\sigma_j^2}}$$

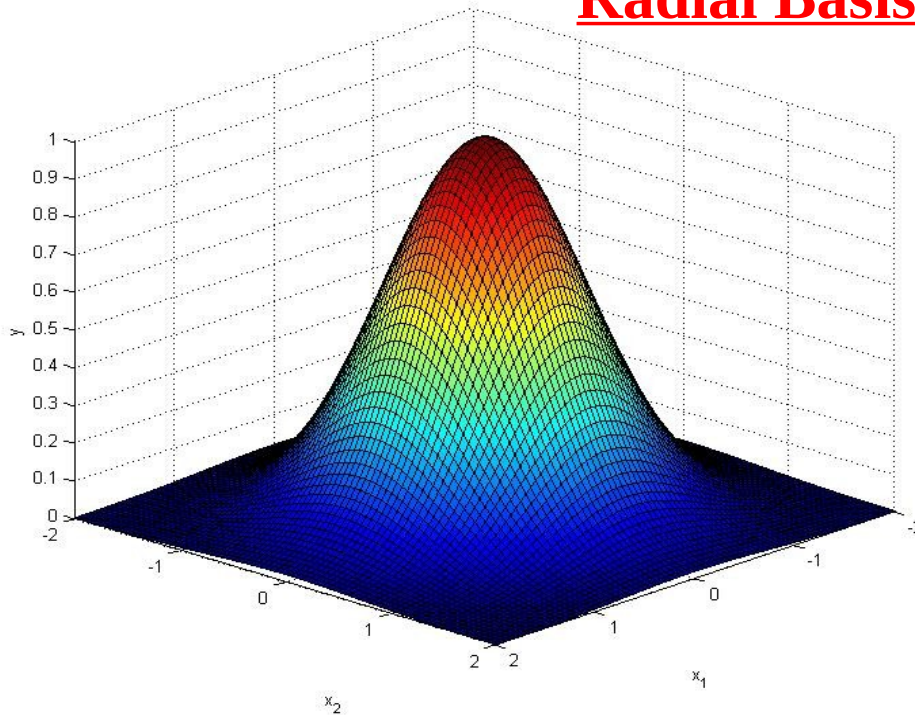
Gaussian

and

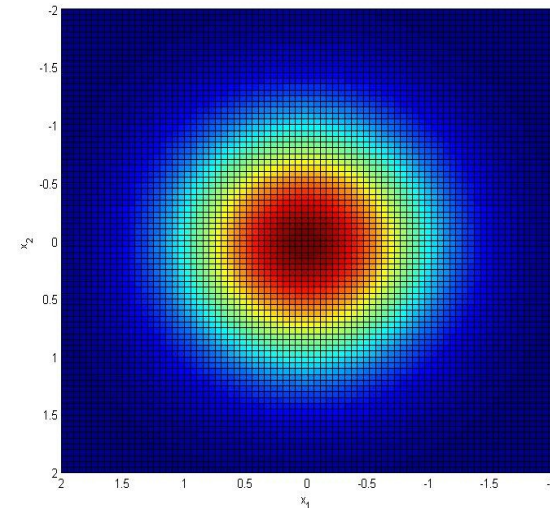
$$\hat{y} = \bar{w}^T \bar{h}$$

linear

Radial Basis Functions



since σ_j is a scalar, the contours on $x_1 - x_2$ are circular (equal variance)



A more general form would be to use
$$h_j = e^{-\left|(\bar{x} - \bar{c}_j)^T \Sigma^{-1} (\bar{x} - \bar{c}_j)\right|}$$

where Σ is a symmetric $n \times n$ matrix with non-negative diagonal values

One can use $\Sigma^{-1} = K^T K$ where K is any real $n \times n$ matrix – weighting matrix.

Training RBF Networks

Option 1:

- Choose fixed $\bar{C}_j, \sigma_j, \forall_j$
- Train \bar{w}

Given training set

$$X = \{\bar{x}^1, \bar{x}^2, \dots, \bar{x}^N\}$$

$$Y = \{y^1, y^2, \dots, y^N\}$$

define

$$H = \{\bar{h}^1, \bar{h}^2, \dots, \bar{h}^N\}$$

$$\text{where } \bar{h}^q = \begin{bmatrix} 1 & h_1(\bar{x}^q) & \dots & h_m(\bar{x}^q) \end{bmatrix}^T$$

Now training \bar{w} is finding weights

for the linear estimate $y = \bar{w}^T \bar{h}$

using the training set $H \rightarrow Y$



Use LMS

(or use other linear estimation methods)

Note: σ_j s could be identical or different.

Option 2:

- Adapt centers \bar{C}_j through self-organization
- Fix $\sigma_j \quad \forall_j$
- Obtain H using the \bar{C}_j, σ_j
- Train \bar{w} using LMS etc.

Notes: \bar{C}_j 's can be obtained by
k-means clustering or some
other clustering method

σ_j can be identical for each
 j or chosen heuristically

we'll do this later

Option 3:

Treat $\sigma_j, \bar{C}_{jk}, w_j$ as
parameters in an optimization
problem and train them all by
gradient descent or some other
method.

Full Supervised Learning