

Radial Basis Function Networks

Sargur Srihari

Topics

- Basis Functions
- Radial Basis Functions
- Gaussian Basis Functions
- Nadaraya Watson Kernel Regression Model
- Decision Tree Initialization of RBF

Basis Functions

- Summary of Linear Regression Models

- Polynomial Regression with one variable

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots = \sum w_i x^i$$

- Simple Linear Regression with D variables

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_D x_D = \mathbf{w}^T \mathbf{x}$$

In one-dimensional case

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x$$

which is a straight line

- Linear Regression with Basis functions $f_j(\mathbf{x})$

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$

- There are now M parameters instead of D parameters
- What form should the basis functions take?

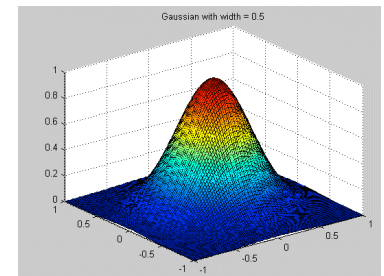
Radial Basis Functions

- A radial basis function depends only on the radial distance (Euclidean) from the origin

$$f(\mathbf{x}) = f(\|\mathbf{x}\|)$$

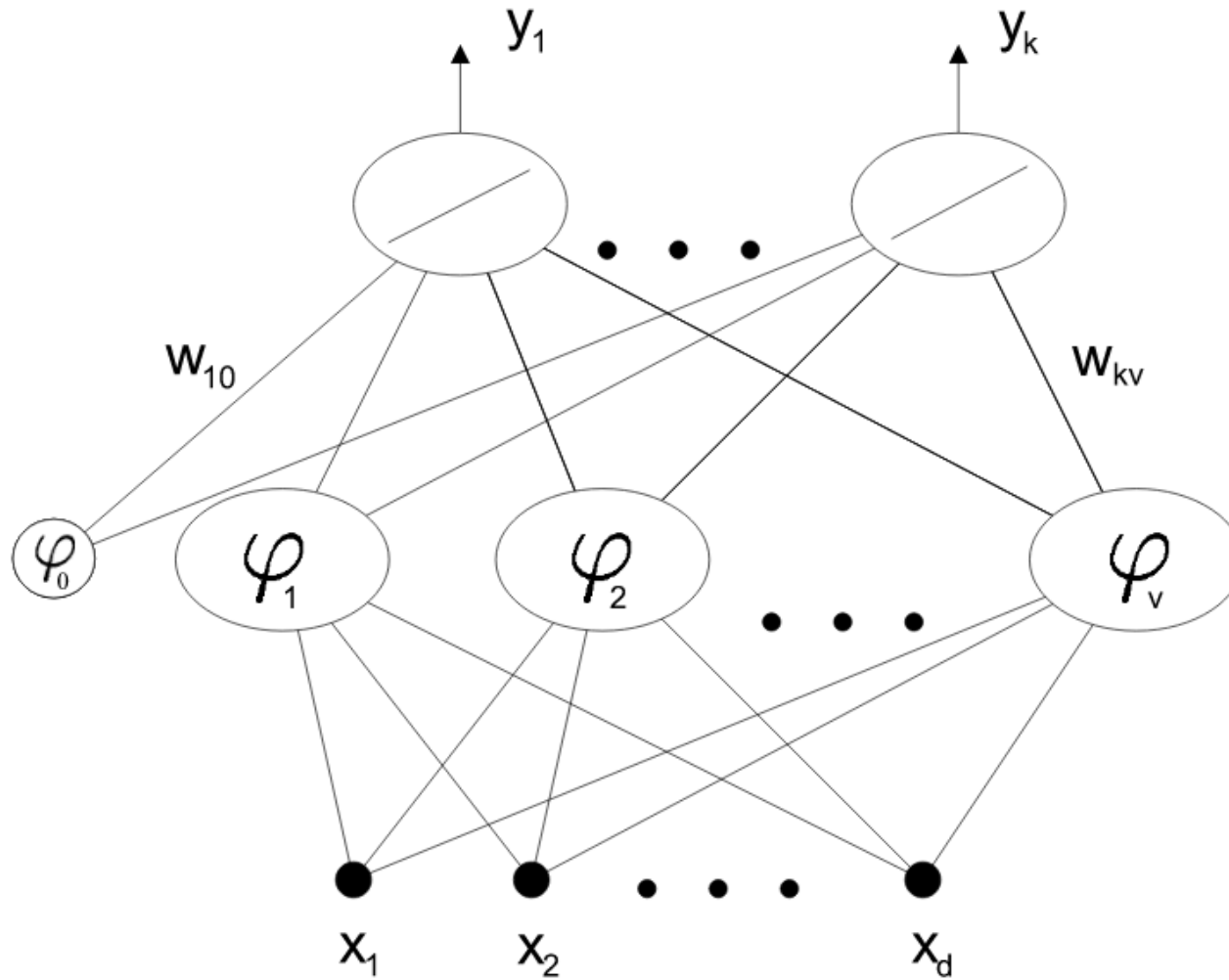
- If the basis function is centered at m_j th

$$f_j(\mathbf{x}) = h(\|\mathbf{x} - m_j\|)$$



- We look at radial basis functions centered at the data points \mathbf{x}_n , $n = 1, \dots, N$

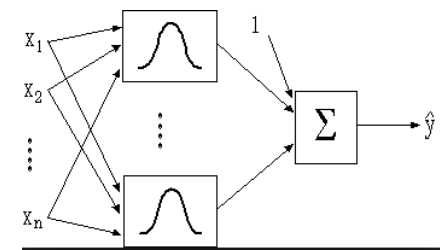
An RBF Network



History of Radial Basis Functions

- Introduced for exact function interpolation
 - Given set of input vectors $\mathbf{x}_1, \dots, \mathbf{x}_N$ and target values t_1, \dots, t_N
 - Goal is to find a smooth function $f(\mathbf{x})$ that fits every target value exactly so that $f(\mathbf{x}_n) = t_n$ for $n=1, \dots, N$
- Achieved by expressing $f(\mathbf{x})$ as a linear combination of radial basis functions, one per data point

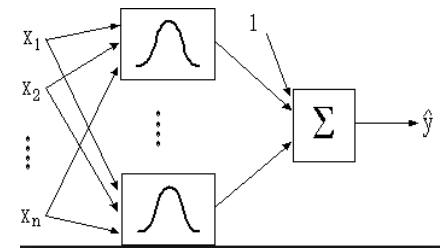
$$f(\mathbf{x}) = \sum_{n=1}^N w_n h(\|\mathbf{x} - \mathbf{x}_n\|)$$



History of Radial Basis Functions

- Values of coefficients w_n are found by least squares
 - Very large number of weights
- Since there are same number of coefficients as constraints, result is a function that fits target values exactly
 - When target values are noisy exact interpolation is undesirable

$$f(\mathbf{x}) = \sum_{n=1}^N w_n h(\|\mathbf{x} - \mathbf{x}_n\|)$$



Another Motivation for Radial Basis Functions

- Interpolation when Input rather than target is noisy
 - If noise on input variable \mathbf{x} is described by variable $\boldsymbol{\xi}$ having distribution $\nu(\boldsymbol{\xi})$ then sum of squared error function becomes

$$E = \frac{1}{2} \sum_{n=1}^N \int \{y(\mathbf{x}_n + \boldsymbol{\xi}) - t_n\}^2 \nu(\boldsymbol{\xi}) d\boldsymbol{\xi}$$

ν : pronounced “nu”
for noise

Another Motivation for Radial Basis Functions

- Using calculus of variations, optimize wrt $y(\mathbf{x})$ to give

$$y(\mathbf{x}) = \sum_{n=1}^N t_n h(\mathbf{x} - \mathbf{x}_n)$$

where the basis functions are given by

$$h(\mathbf{x} - \mathbf{x}_n) = \frac{v(\mathbf{x} - \mathbf{x}_n)}{\sum_{n=1}^N v(\mathbf{x} - \mathbf{x}_n)}$$

- There is a basis function centered on every data point
- Known as Nadaraya-Watson Model

Normalized Basis Functions

If noise $\mathbf{v}(\mathbf{x})$ is isotropic so that it is only a function of $\|\mathbf{x}\|$ then the basis functions are radial

Functions are normalized

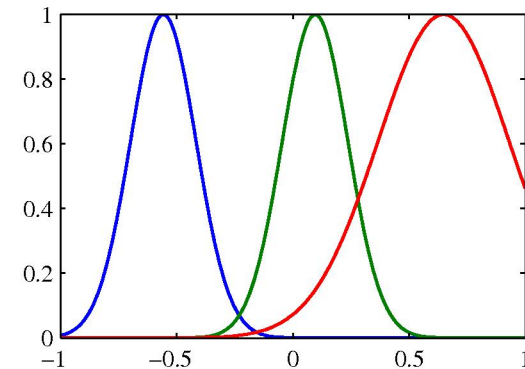
$$h(\mathbf{x} - \mathbf{x}_n) = \frac{\mathbf{v}(\mathbf{x} - \mathbf{x}_n)}{\sum_{n=1}^N \mathbf{v}(\mathbf{x} - \mathbf{x}_n)}$$

so that

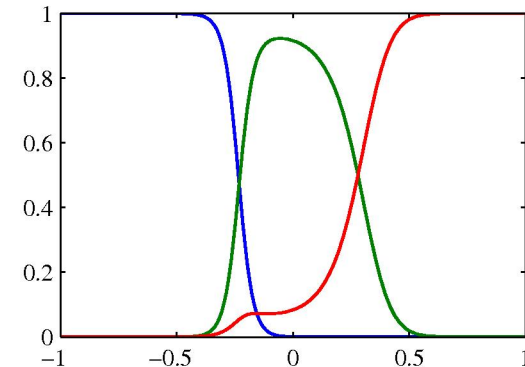
$$\sum_n h(\mathbf{x} - \mathbf{x}_n) = 1 \quad \text{for any value of } \mathbf{x} \quad \rightarrow$$

Normalization is useful in regions of input space where all basis functions are small

Gaussian Basis Functions



Normalized Basis Functions



$h(\mathbf{x} - \mathbf{x}_n)$ is called a kernel function since we use it with every sample to determine value at \mathbf{x}

Factorization into basis functions is not used

Computational Efficiency

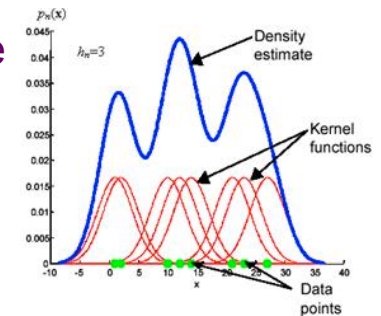
- Since there is one basis function per data point corresponding computational model will be costly to evaluate for new data points
- Methods for choosing a smaller set
 - Use a random subset of data points
 - Orthogonal least squares
 - At each step the next data point to be chosen as basis function center is the one that gives the greatest reduction in least-squared error
 - Clustering algorithms which give basis function centers that no longer coincide with training data points

Nadaraya-Watson Regression Model

- More general formulation than before
 - Proposed in 1964. Nadaraya: a Russian statistician
- Begin with Parzen estimation to derive kernel function
 - Given a training set $\{\mathbf{x}_n, t_n\}$ the joint distribution of two variables is

$$p(\mathbf{x}, t) = \frac{1}{N} \sum_{n=1}^N f(\mathbf{x} - \mathbf{x}_n, t - t_n)$$

Parzen window estimate
With one variable $p(x)$



- where $f(\mathbf{x}, t)$ is the component density function
- there is one such component centered at each data point

- Starting $\int_{-\infty}^{\infty} f(\mathbf{x}, t) t dt = 0 = \int t p(t|\mathbf{x}) dt$ we derive the kernel form

- Where we need to note that component densities have zero mean

Derivation of Nadaraya-Watson Regression

- Desired Regression function

$$\begin{aligned} y(\mathbf{x}) &= \mathbb{E}[t|\mathbf{x}] = \int_{-\infty}^{\infty} tp(t|\mathbf{x}) dt \\ &= \frac{\int tp(\mathbf{x}, t) dt}{\int p(\mathbf{x}, t) dt} \end{aligned}$$

- Using Parzen window estimate and change of variables

$$\begin{aligned} y(\mathbf{x}) &= \frac{\sum_n g(\mathbf{x} - \mathbf{x}_n)t_n}{\sum_m g(\mathbf{x} - \mathbf{x}_m)} \\ &= \sum_n k(\mathbf{x}, \mathbf{x}_n)t_n \end{aligned}$$

- Where k is a normalized kernel

$$k(\mathbf{x}, \mathbf{x}_n) = \frac{g(\mathbf{x} - \mathbf{x}_n)}{\sum_m g(\mathbf{x} - \mathbf{x}_m)} \quad g(\mathbf{x}) = \int_{-\infty}^{\infty} f(\mathbf{x}, t) dt.$$

Example of Nadaraya-Watson Regression

Regression function is

$$y(x) = \sum_n k(x, x_n) t_n$$

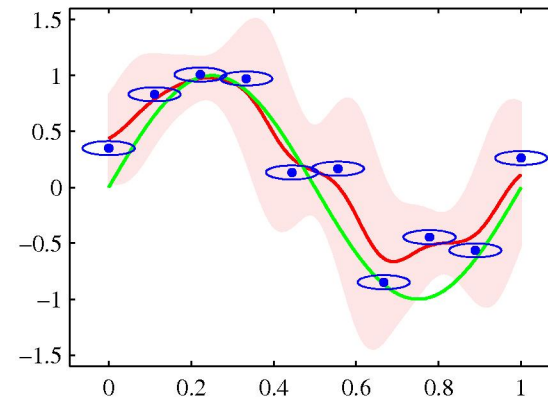
Where the kernel function is

$$k(x, x_n) = \frac{g(x - x_n)}{\sum_m g(x - x_m)}$$

and we have defined

$$g(x) = \int_{-\infty}^{\infty} f(x, t) dt$$

and $f(x, t)$ is a zero-mean component density function centered at each data point



Green: Original sine function

Blue: Data points

Red: Resulting regression function

Pink region: Two std deviation region for distribution $p(t|x)$

Blue ellipses: one std deviation contour

Different scales on x and y axes make circles elongated

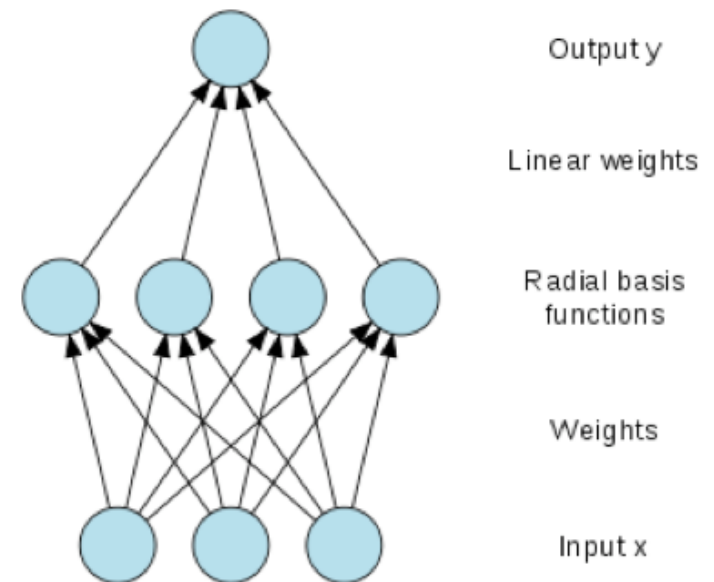
Radial Basis Function Network

- A neural network that uses RBFs as activation function

$$\varphi(\mathbf{x}) = \sum_{i=1}^N a_i \rho(\|\mathbf{x} - \mathbf{c}_i\|)$$

- In Nadaraya-Watson

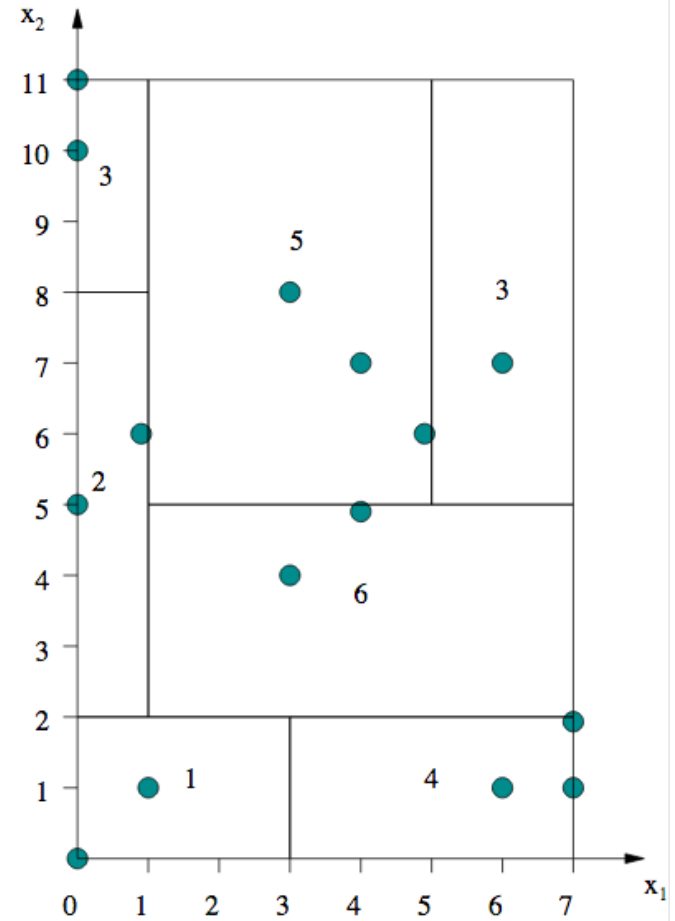
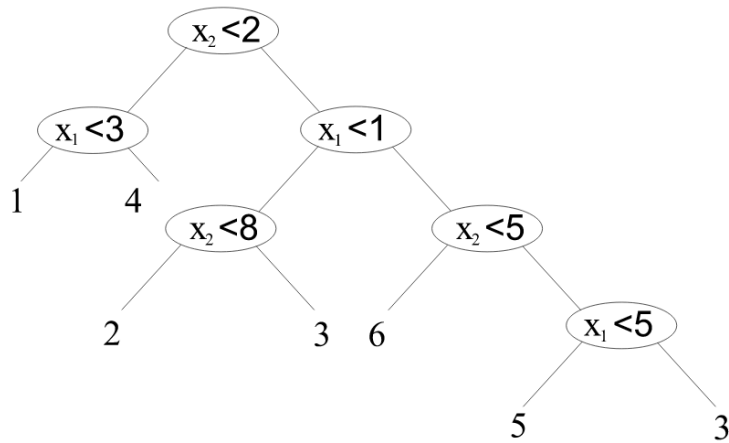
- Weights a_i are target values
- r is component density (Gaussian)
- Centers \mathbf{c}_i are samples



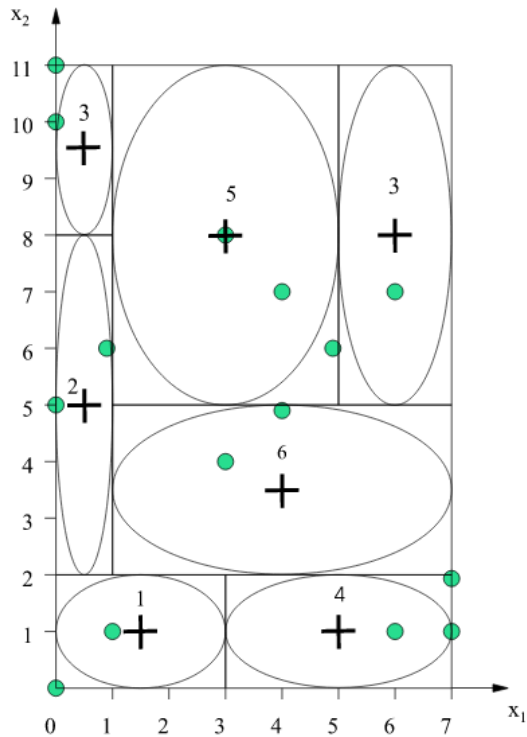
Speeding-up RBFs

- More flexible forms of Gaussian components can be used
 - Model $p(\mathbf{x}, t)$ using a Gaussian mixture model
 - Trained using EM
- Number of components in the mixture model can be smaller than the number of training samples
 - Faster evaluation for test data points
 - But increased training time

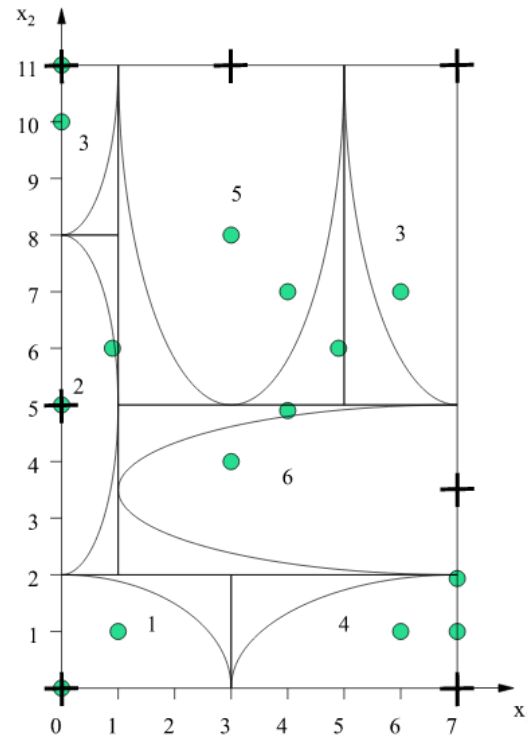
Decision Tree Initialization



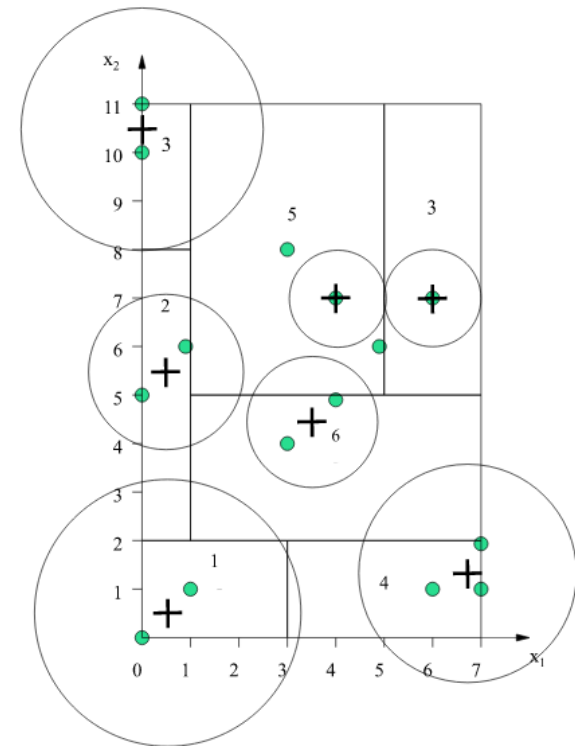
Three Initialization Methods



(a) C1 + S1



(b) C2 + S2



(c) C3 + S3

References

- <http://mccormickml.com/2013/08/15/radial-basis-function-network-rbfn-tutorial/>
- <http://www.informatik.uni-ulm.de/ni/staff/CDietrich/publications/nnw2000.pdf>