

INM431 Machine Learning

Artur S. d'Avila Garcez

a.garcez@city.ac.uk

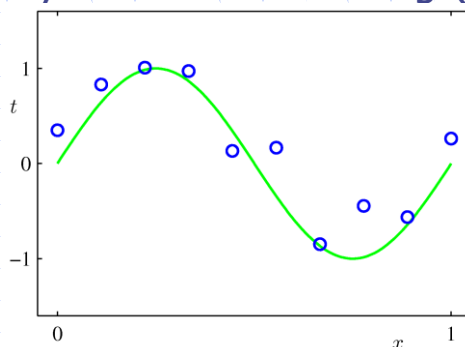
<http://www.staff.city.ac.uk/~aag/>

Content

Regression models

Linear models for regression

E.g.: Polynomial Curve Fitting (revisited)



$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

Linear models for regression (cont.)

Generally:

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^M w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

where $\phi(x)$ are known as *basis functions*

Typically, $\phi_0(x)=1$, so that w_0 acts as a bias

In the simplest case, we use linear basis functions...

Linear regression

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_n x_n$$

given a set of n-dimensional data points

$$\mathbf{x} = (x_1, \dots, x_n).$$

This model is a linear function of the parameters \mathbf{w} AND of the input variables x_i

In the general case, this is extended to consider linear combinations of fixed nonlinear functions of the input variables

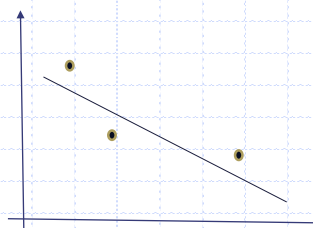
© Artur Garcez

Example

When 2 highway patrol cars are deployed, the average speed on the freeway is 75mph. When 4 patrol cars are deployed the average speed is 45mph. When 10 patrol cars are deployed the average speed is 35mph. Using linear (linear) regression and least squares, what will be the average speed when 5 cars are deployed?

$$E = \frac{1}{2} \sum_i (t_i - (w_1 x_i + w_0))^2$$

$$\frac{\partial E}{\partial W} = 0$$



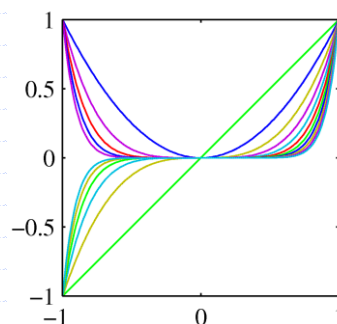
© Artur Garcez

Possible choices of basis function

Polynomial basis functions:

$$\phi_j(x) = x^j.$$

These are global; a small change in x affect all basis functions.

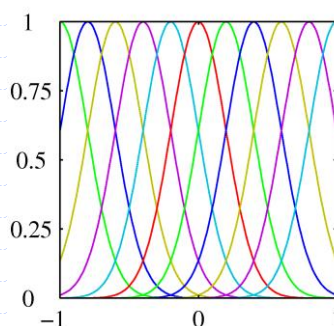


Possible choices of basis function (cont.)

Gaussian basis functions:

$$\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\}$$

These are local; a small change in x only affect nearby basis functions...



Note: s in equation above is *standard deviation*, i.e. σ from previous slides

Possible choices of basis function (cont.)

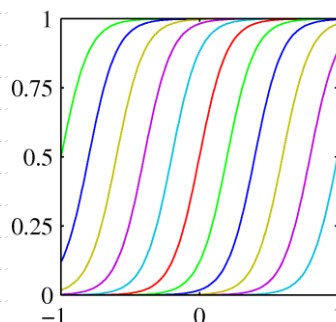
Sigmoidal basis functions:

$$\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right)$$

where

$$\sigma(a) = \frac{1}{1 + \exp(-a)}.$$

Also these are local; a small change in x only affect nearby basis functions.



Note: s in equation above is *standard deviation*, i.e. σ from previous slides; σ in equation above is sigmoid function, i.e. not *standard deviation*...

Using Maximum Likelihood and Least Squares

Assume observations from a deterministic function with added Gaussian noise:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon \quad \text{where} \quad p(\epsilon|\beta) = \mathcal{N}(\epsilon|0, \beta^{-1})$$

which is the same as saying:

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1}).$$

Given observed inputs $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ and targets $\mathbf{t} = [t_1, \dots, t_N]^T$

We obtain the likelihood function:

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n|\mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}).$$

Using Maximum Likelihood and Least Squares (cont.)

Taking the logarithm, we get

$$\begin{aligned}\ln p(\mathbf{t}|\mathbf{w}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w})\end{aligned}$$

where

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2$$

is the sum-of-squares error.

Maximum Likelihood and Least Squares...

Computing the gradient and setting it to zero yields

$$\nabla_{\mathbf{w}} \ln p(\mathbf{t}|\mathbf{w}, \beta) = \beta \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\} \phi(\mathbf{x}_n)^T = \mathbf{0}.$$

Solving for \mathbf{w} , we get

$$\mathbf{w}_{\text{ML}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

where

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}.$$