

Linear Models for Regression

Alberto Paccanaro
EMAp – FGV

www.paccanarolab.org

Material and images in these slides are from (or adapted from):
C. Bishop, Pattern Recognition and Machine Learning, Springer, 2006

The Regression problem

Predict the value of **continuous** *target* variables t given the value of a D -dimensional vector \mathbf{x} of *input* variables.

Given N observations $\{\mathbf{x}_n\}$, where $n = 1, \dots, N$, together with corresponding target values $\{t_n\}$, the goal is to predict the value of t for a new value of \mathbf{x} .

The Roadmap

We model the predictive distribution $p(t|\mathbf{x})$ expressing our uncertainty about the value of t for each value of \mathbf{x} .

- learn $p(t|\mathbf{x})$ by minimizing a loss function.
- common choice of loss function for real-valued variables: Sum of Squared Errors – **consequence of maximizing likelihood under the assumption of a Gaussian noise distribution.**
 - ➔ the optimal solution is given by the conditional expectation of t (from decision theory)

Linear Models: linear functions of the adjustable parameters

IMPORTANT: can be nonlinear with respect to the input variables (e.g. the polynomial we saw earlier)

- simplest form: models are also linear functions of the input variables.
- *basis functions*: nonlinear functions of the input variables, of which we take linear combinations

Regression by linear combination of basis functions

- Simplest linear model for regression (often just called *linear regression*):

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1x_1 + \dots + w_Dx_D = w_0 + \sum_{i=1}^D w_ix_i$$

- linear combinations of ~~fixed~~ ~~nonlinear~~ Basis Functions, $\phi_j(\mathbf{x})$:

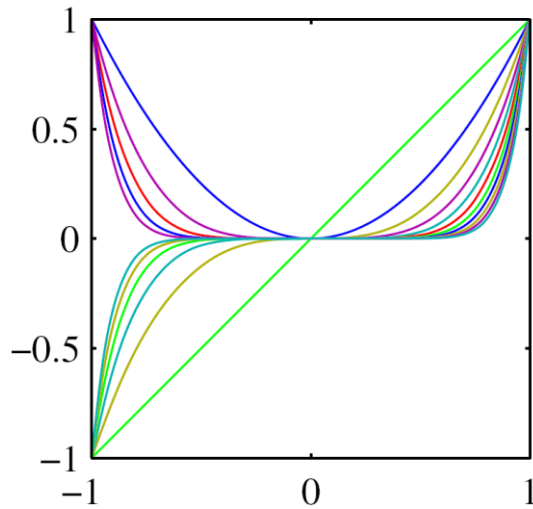
$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j\phi_j(\mathbf{x})$$

Nonlinear basis functions
allow $y(\mathbf{x}, \mathbf{w})$ to be nonlinear !

- w_0 , called **bias** parameter, can be included in the sum, by defining one extra basis function $\phi_0(\mathbf{x}) = 1$

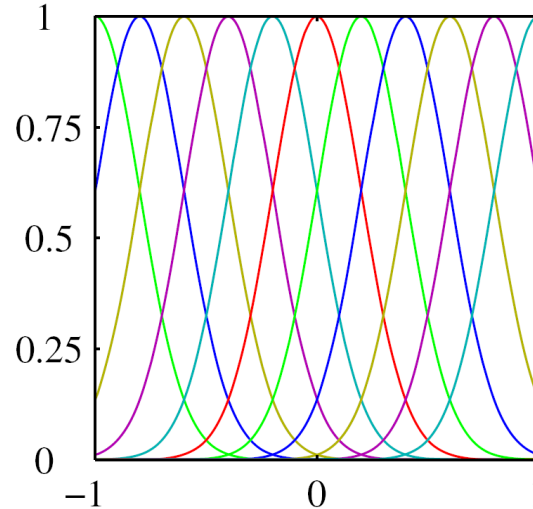
$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j\phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

Classic basis functions



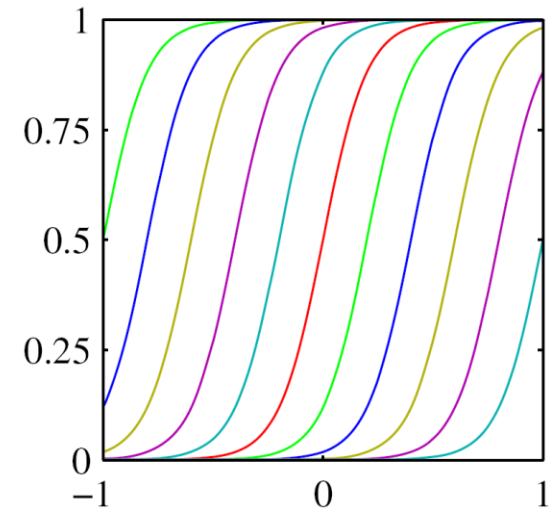
Polynomial

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



Gaussian

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



Sigmoidal

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

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

μ_i control the locations of the basis functions
 s controls the scale of the basis function

Maximum likelihood and least squares

(it's similar to what we did earlier...)

Assume t is given by:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon \quad \epsilon \text{ gaussian random variable}$$

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

Given data set of N inputs \mathbf{X}

Likelihood function (a function of \mathbf{w} and β)

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

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

where: $E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2$ is the sum-of-squares error function

Gradient of the log likelihood function

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

Setting the gradient to 0:

$$0 = \sum_{n=1}^N t_n \phi(\mathbf{x}_n)^T - \mathbf{w}^T \left(\sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \right)$$

Solving for \mathbf{w} we obtain

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

Moore-Penrose
pseudo-inverse of Φ

the *normal equations* for
the least squares problem

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

Design Matrix

$$\Phi_{nj} = \phi_j(\mathbf{x}_n)$$

Important connection to help us understand...

The Moore-Penrose pseudo-inverse of the matrix Φ , Φ^\dagger can be regarded as a generalization of the notion of matrix inverse to nonsquare matrices

$$\Phi^\dagger \equiv (\Phi^T \Phi)^{-1} \Phi^T$$

If Φ is square and invertible then using the property $(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$ we see that $\Phi^\dagger \equiv \Phi^{-1}$.

Sequential/online learning

Batch techniques can be computationally costly for large datasets.

Sequential algorithms: data points are considered one at a time. Model parameters updated after each presentation.

Stochastic gradient descent/ sequential gradient descent

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n$$

E_n is the error for point n

For sum-of-squares error function:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \eta(t_n - \mathbf{w}^{(\tau)\top} \phi_n) \phi_n$$

Least-Mean-Squares or the LMS algorithm

$$\phi_n = \phi(\mathbf{x}_n)$$

Regularized Least Squares

Idea: adding a regularization term to an error function to control over-fitting

$$E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$$

1. Weight decay (in machine learning)

$$E_W(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

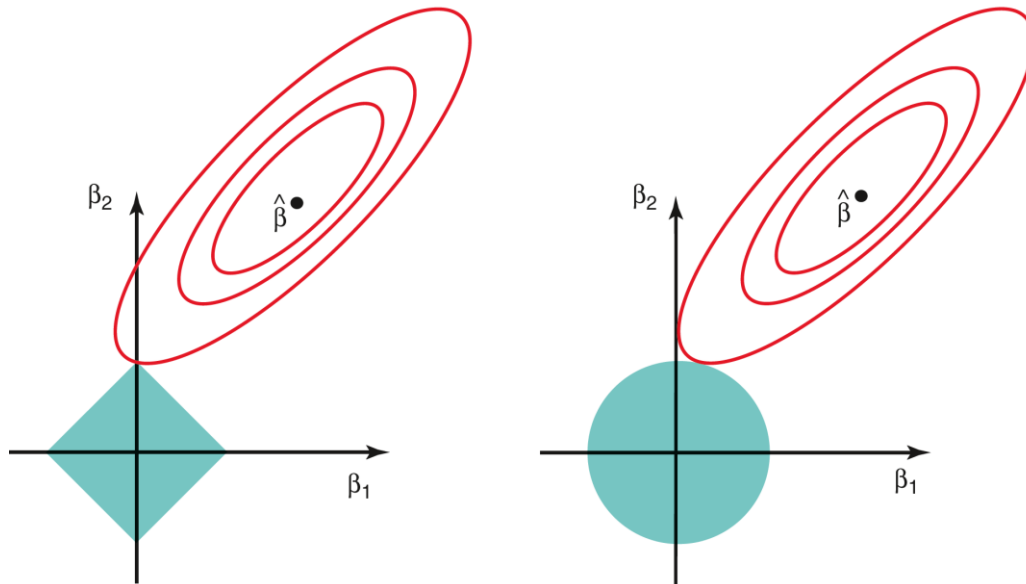
$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \quad (\text{ridge regression in statistics})$$

$$\mathbf{w} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \quad \text{normal equations with weight decay}$$

2. Lasso (Least Absolute Shrinkage and Selection Operator)

$$E_W(\mathbf{w}) = \frac{\lambda}{2} \sum_{j=1}^M |w_j|$$

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j| \quad \text{Lasso regression}$$



some of the coefficients w_j become 0, leading to a sparse model

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q \quad \text{General regularization form}$$

3. Elastic Net

$$E_W(\mathbf{w}) = \frac{\lambda_1}{2} \sum_{j=1}^M |w_j| + \frac{\lambda_2}{2} \mathbf{w}^T \mathbf{w}$$

combines Lasso and Ridge regularization

Linear regression for multiple outputs

Predict $K > 1$ target variables – target vector \mathbf{t}

The generalization is straightforward

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

$$p(\mathbf{t}|\mathbf{x}, \mathbf{W}, \beta) = \mathcal{N}(\mathbf{t}|\mathbf{W}^T \phi(\mathbf{x}), \beta^{-1}\mathbf{I})$$

$$\begin{aligned} \ln p(\mathbf{T}|\mathbf{X}, \mathbf{W}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(\mathbf{t}_n|\mathbf{W}^T \phi(\mathbf{x}_n), \beta^{-1}\mathbf{I}) \\ &= \frac{NK}{2} \ln \left(\frac{\beta}{2\pi} \right) - \frac{\beta}{2} \sum_{n=1}^N \|\mathbf{t}_n - \mathbf{W}^T \phi(\mathbf{x}_n)\|^2 \end{aligned}$$

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

$$\mathbf{w}_k = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}_k = \Phi^\dagger \mathbf{t}_k$$