Sapienza University of Rome

Master in Artificial Intelligence and Robotics

# Machine Learning

A.Y. 2025/2026

Prof. Luca Iocchi

# 8. Linear models for regression

Luca Iocchi

# Overview

- Linear models for regression
- Maximum likelihood and Least squares
- Sequential learning
- Regularization

*References*
C. Bishop. Pattern Recognition and Machine Learning. Sect. 3.1

# Linear Models for Regression

Learning a function $f : X \to Y$, with

- $X \subseteq \mathbb{R}^d$
- $Y = \mathbb{R}$

from data set $D = \{(x_n, t_n)_{n=1}^{N}\}$

# Linear Models for Regression

Define a model $y(x; w)$ with parameters $w$ to approximate the target function $f$.
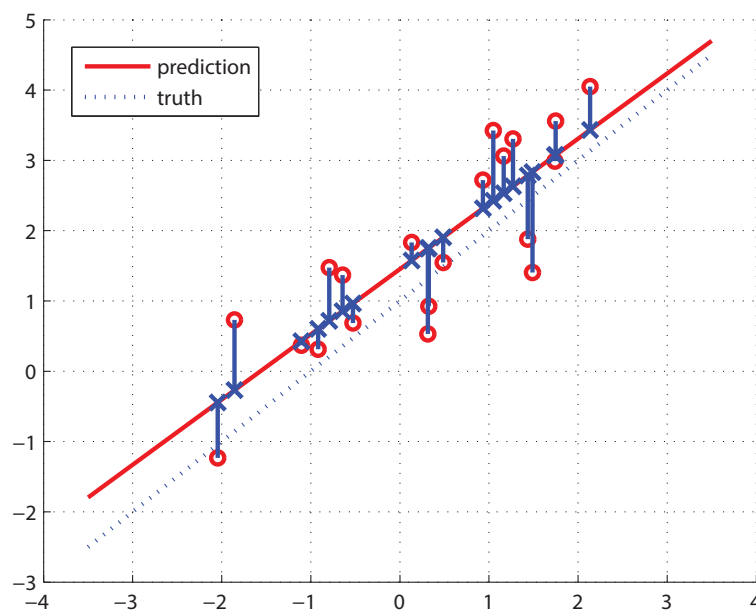
Linear model for linear functions

$$y(x; w) = w_0 + w_1 x_1 + \ldots + w_d x_d = w^T x$$

$$\text{with } x = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} \text{ and } w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}$$

# Example: 2D line fitting

$$y = w_0 + w_1 x_1$$

# Linear Models for Regression

Linear Basis Function Models
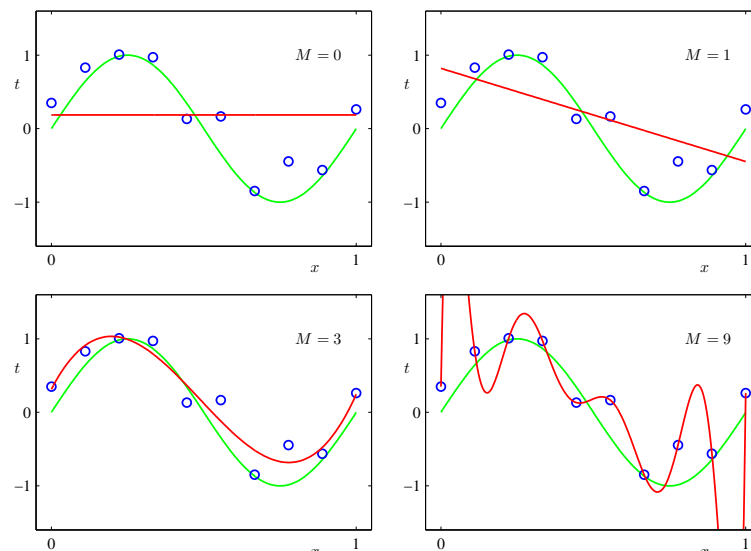
Using nonlinear functions of input variables:

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

$$\text{with } w = \begin{bmatrix} w_0 \\ \vdots \\ w_M \end{bmatrix}, \ \phi(x) = \begin{bmatrix} \phi_0(x) \\ \vdots \\ \phi_M(x) \end{bmatrix}, \text{ and } \phi_0(x) = 1.$$

- Still linear in the parameters w!
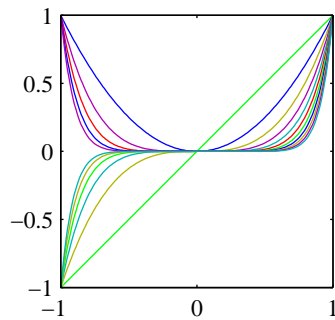
# Example: Polynomial curve fitting

$$y = w_0 + w_1 x + w_2 x^2 + \ldots + w_M x^M = \sum_{j=0}^{M} w_j x^j$$
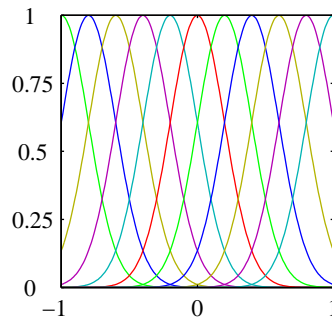


**Warning: overfitting!!!**

# Linear Regression Basis Functions
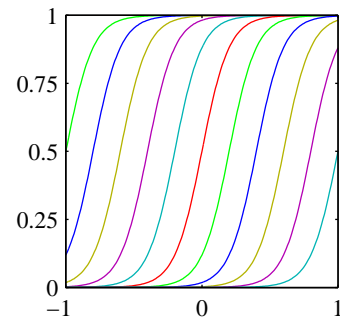
Examples of basis functions



Polynomial          Radial          Sigmoid / Tanh

# Linear Regression - Algorithms

**Maximum likelihood and least squares**

Target value $t$ is given by $y(\mathrm{x}; \mathrm{w})$ affected by additive noise $\epsilon$

$$t = y(\mathrm{x}; \mathrm{w}) + \epsilon$$

Assume Gaussian noise $P(\epsilon|\beta) = \mathcal{N}(\epsilon|0, \beta^{-1})$, with precision (inverse variance) $\beta$.

We have:

$$P(t|\mathrm{x}, \mathrm{w}, \beta) = \mathcal{N}(t|y(\mathrm{x}; \mathrm{w}), \beta^{-1})$$

# Linear Regression - Algorithms

Assume observations independent and identically distributed (i.i.d.)

We seek the maximum of the likelihood function:

$$P(\{t_1, \ldots, t_N\}|\mathsf{x}_1, \ldots, \mathsf{x}_N, \mathsf{w}, \beta) = \prod_{n=1}^{N} \mathcal{N}(t_n|\mathsf{w}^T \phi(\mathsf{x}_n), \beta^{-1}).$$

or equivalently:

$$\ln P(\{t_1, \ldots, t_N\}|\mathsf{x}_1, \ldots, \mathsf{x}_N, \mathsf{w}, \beta) = \sum_{n=1}^{N} \ln \mathcal{N}(t_n|\mathsf{w}^T \phi(\mathsf{x}_n), \beta^{-1})$$

$$= -\beta \underbrace{\frac{1}{2} \sum_{n=1}^{N} [t_n - \mathsf{w}^T \phi(\mathsf{x}_n)]^2}_{E_D(\mathsf{w})} - \frac{N}{2} \ln(2\pi\beta^{-1}).$$

# Linear Regression - Algorithms

Maximum likelihood (zero-mean Gaussian noise assumption)

$$\operatorname*{argmax}_{\mathsf{w}} P(\{t_1, \ldots, t_N\}|\mathsf{x}_1, \ldots, \mathsf{x}_N, \mathsf{w}, \beta)$$

corresponds to least square error minimization

$$\operatorname*{argmin}_{\mathsf{w}} E_D(\mathsf{w}) = \operatorname*{argmin}_{\mathsf{w}} \frac{1}{2} \sum_{n=1}^{N} [t_n - \mathsf{w}^T \phi(\mathsf{x}_n)]^2$$

# Linear Regression - Algorithms

Note:

$$E_D(w) = \frac{1}{2}(t - \Phi w)^T (t - \Phi w),$$

$$\text{with } t = \begin{bmatrix} t_1 \\ \vdots \\ t_N \end{bmatrix} \text{ and } \Phi = \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \cdots & \phi_M(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \cdots & \phi_M(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \cdots & \phi_M(x_N) \end{bmatrix}.$$

Optimality condition:

$$\nabla E_D = 0 \iff \Phi^T \Phi w = \Phi^T t.$$

Hence:

$$w_{ML} = \underbrace{(\Phi^T \Phi)^{-1} \Phi^T}_{\Phi^\dagger: \text{ pseudo-inverse}} t.$$

# Linear Regression - Algorithms

## Sequential Learning

Stochastic gradient descent algorithm:

$$\hat{w} \leftarrow \hat{w} - \eta \nabla E_S$$

$\eta$: learning rate parameter

$\nabla E_S$: estimation of the gradient of $E$ wrt w on a random subset $S \subset D$.

Algorithm converges for suitable small values of $\eta$.

# Linear Regression - Regularization

Regularization is a technique to control over-fitting.

$$\underset{\mathsf{w}}{\operatorname{argmin}}\ E_D(\mathsf{w}) + \lambda E_W(\mathsf{w})$$

with $\lambda > 0$ being the regularization factor

A common choice:

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

Other choices:

$$E_W(\mathsf{w}) = \sum_{j=0}^{M} |w_j|^q.$$

# Linear Regression - Regularization

$$\underset{\mathsf{w}}{\operatorname{argmin}}\ E_D(\mathsf{w})$$

# Linear Regression - Regularization

$$\operatorname*{argmin}_{\mathsf{w}}\ E_D(\mathsf{w}) + \lambda \frac{1}{2}\mathsf{w}^T\mathsf{w}$$

# Linear Regression - Multiple outputs

y: vector with $K$ components

$$y(\mathsf{x}; \mathsf{W}) = \mathsf{W}^T \phi(\mathsf{x})$$

Target variable T, with $\mathsf{t}_n$ vector of $K$ output values for input $\mathsf{x}_n$

$$\ln P(\mathsf{T}|\mathsf{X}, \mathsf{W}, \beta) = \sum_{n=1}^{N} \ln \mathcal{N}(\mathsf{t}_n|\mathsf{W}^T\phi(\mathsf{x}_n), \beta^{-1}\mathsf{I})$$

Similarly as before we obtain:

$$\mathsf{W}_{ML} = (\Phi^T\Phi)^{-1}\Phi^T\mathsf{T}.$$