



LINEAR REGRESSION FOR NONLINEAR PROBLEMS

CS6140

Predrag Radivojac
KHOURY COLLEGE OF COMPUTER SCIENCES
NORTHEASTERN UNIVERSITY

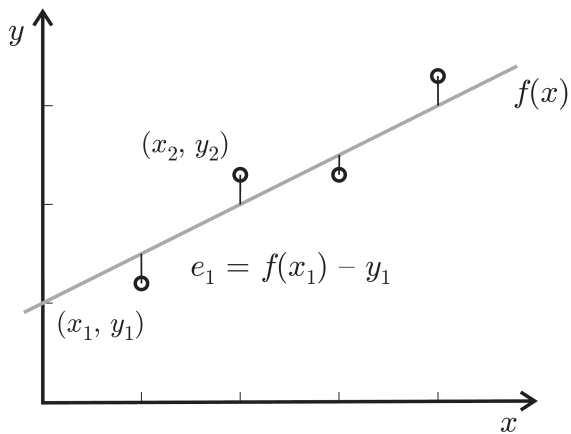
Spring 2021

LINEAR REGRESSION

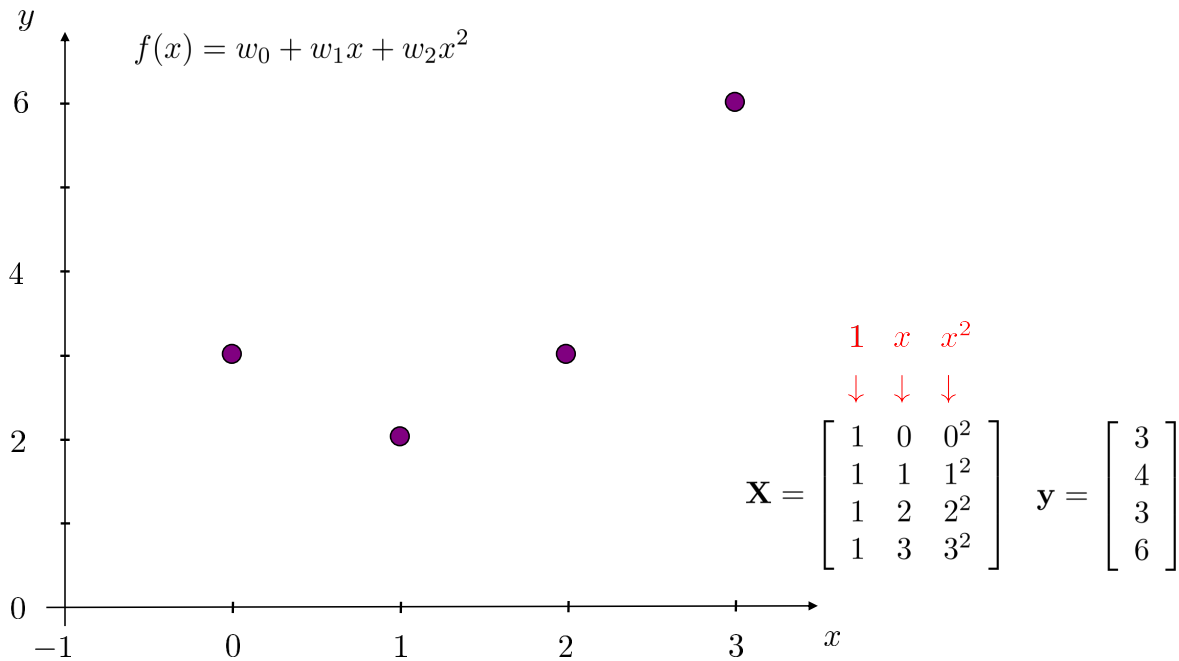
Given: a set of observations $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, $(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$

Objective: find best linear approximator $f(\mathbf{x}) = w_0 + \sum_{j=1}^d w_j x_j$

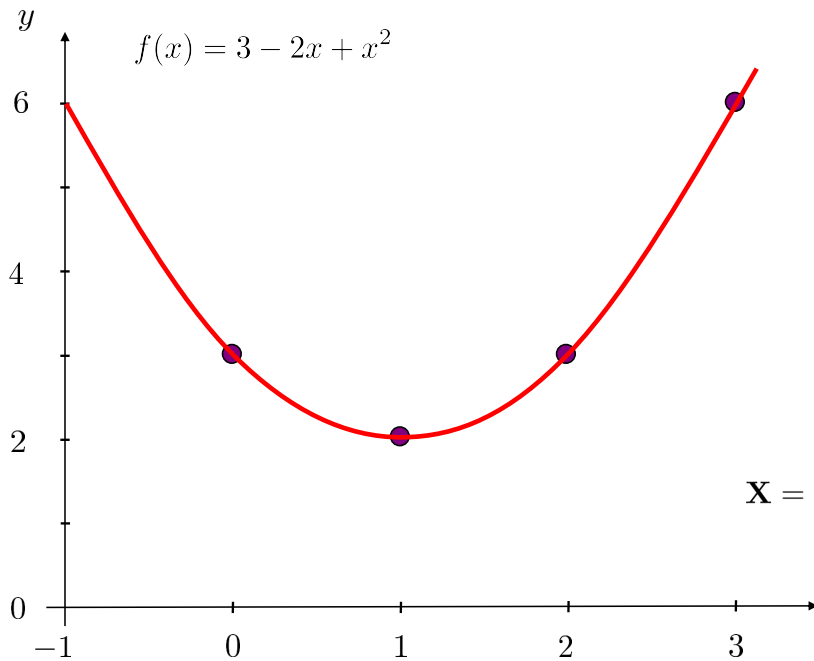
$\mathcal{X} = \mathbb{R}$, $\mathcal{Y} = \mathbb{R}$



LEARNING POLYNOMIAL FUNCTIONS



LEARNING POLYNOMIAL FUNCTIONS



$$\mathbf{X} = \begin{matrix} & \begin{matrix} 1 & x & x^2 \\ \downarrow & \downarrow & \downarrow \end{matrix} \\ \begin{bmatrix} 1 & 0 & 0^2 \\ 1 & 1 & 1^2 \\ 1 & 2 & 2^2 \\ 1 & 3 & 3^2 \end{bmatrix} & \mathbf{y} = \begin{bmatrix} 3 \\ 4 \\ 3 \\ 6 \end{bmatrix} \end{matrix}$$

$$e_1^2 + e_2^2 + e_3^2 + e_4^2 = 0$$

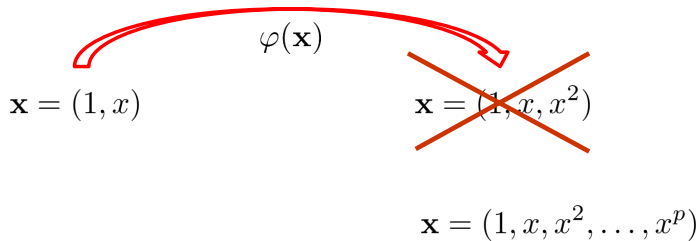
LET'S REFLECT FOR A MOMENT

Q: Have we learned a non-linear function using linear regression?

A: Yup!

Q: How did it happen?

A: We mapped the data into a higher-dimension using a nonlinear transformation.



p = degree of the polynomial

SUMMARY

Original data:

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{bmatrix}$$



Intermediate data:

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ 1 & x_3 & x_3^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{bmatrix}$$




Still the same solution:

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

WHAT IF THE DATA IS MULTIVARIATE?

$\varphi(\mathbf{x})$


$$\mathbf{X} = \begin{matrix} \mathbf{x}_2^T \\ \mathbf{X} \end{matrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1d} \\ 1 & x_{21} & x_{22} & \cdots & x_{2d} \\ 1 & x_{31} & x_{32} & \cdots & x_{3d} \\ \vdots & & & \ddots & \vdots \\ 1 & x_{n1} & & & x_{nd} \end{bmatrix}$$
$$\Phi = \begin{bmatrix} 1 & \phi_1(\mathbf{x}_1) & \phi_2(\mathbf{x}_1) & \cdots & \phi_p(\mathbf{x}_1) \\ 1 & \phi_1(\mathbf{x}_2) & \phi_2(\mathbf{x}_2) & \cdots & \phi_p(\mathbf{x}_2) \\ 1 & \phi_1(\mathbf{x}_3) & \phi_2(\mathbf{x}_3) & \cdots & \phi_p(\mathbf{x}_3) \\ \vdots & & & \ddots & \vdots \\ 1 & \phi_1(\mathbf{x}_n) & & & \phi_p(\mathbf{x}_n) \end{bmatrix}$$

$$\mathbf{w}^* = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}.$$

WHAT IF THE DATA IS MULTIVARIATE?

1. Cluster data into p clusters.
2. Pick p points $\mathbf{c}_1 \dots \mathbf{c}_p$; e.g., examples or centers
3. Make a transformation for each input example

Centers: $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_p$

Parameters: $\sigma_1, \sigma_2, \dots, \sigma_p$

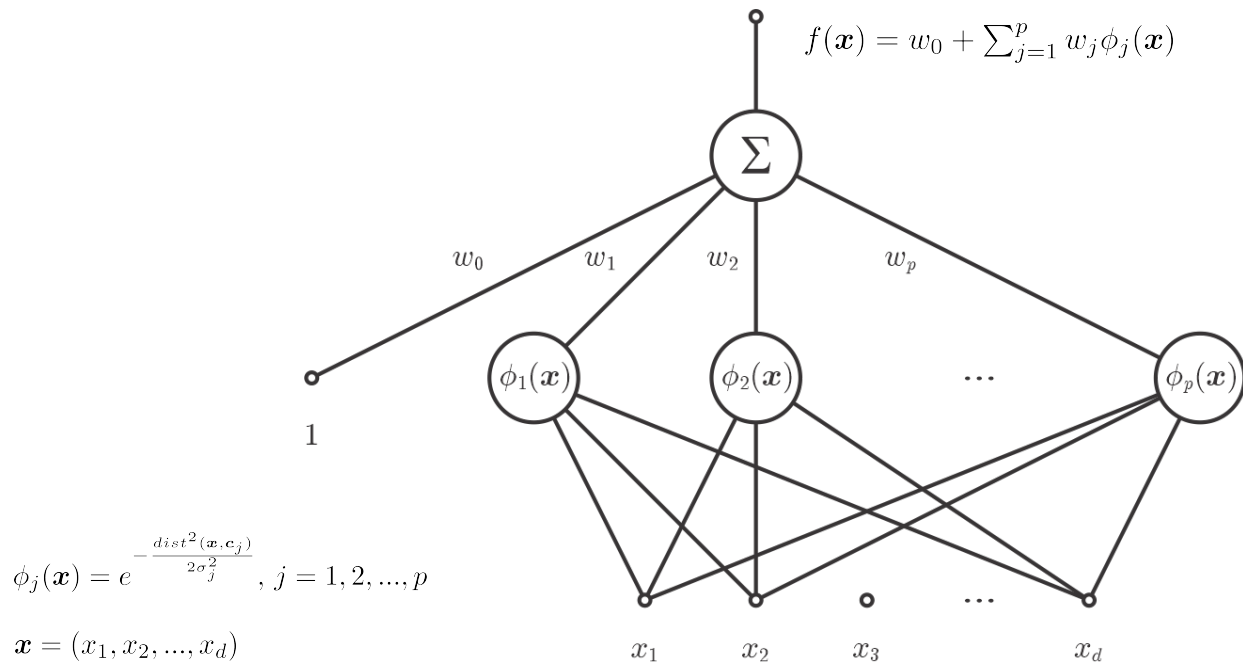
$$\phi_j(\mathbf{x}) = e^{-\frac{(\mathbf{x} - \mathbf{c}_j)^T (\mathbf{x} - \mathbf{c}_j)}{2\sigma_j^2}}$$



Radial basis function

$$j = 1, 2, \dots, p$$

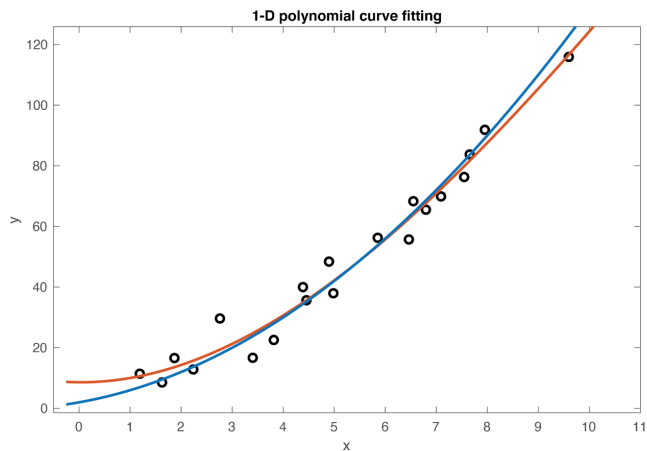
RADIAL BASIS FUNCTION NETWORK



EXAMPLE: POLYNOMIAL CURVE FITTING

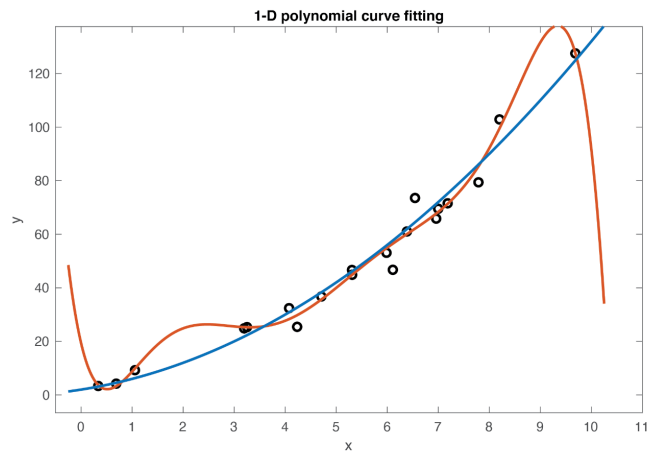
$$y = 2 + 3x + x^2 + \epsilon, \text{ where } \epsilon \sim \mathcal{N}(0, 25)$$

$p = 3$



$$\hat{\mathbf{w}} = (13, -2.6, 2.0 - 0.1)$$

$p = 7$



$$\hat{\mathbf{w}} = (19.4, -80.9, 123.5, -70.3, 19.8, -2.9, 0.2, 0.0)$$

REGULARIZATION

Idea: modify the objective function

$$\text{Objective} = \underbrace{\text{original objective}}_{\text{sum of squared errors}} + \underbrace{\text{regularization term}}_{\text{a function of } \mathbf{w}}$$

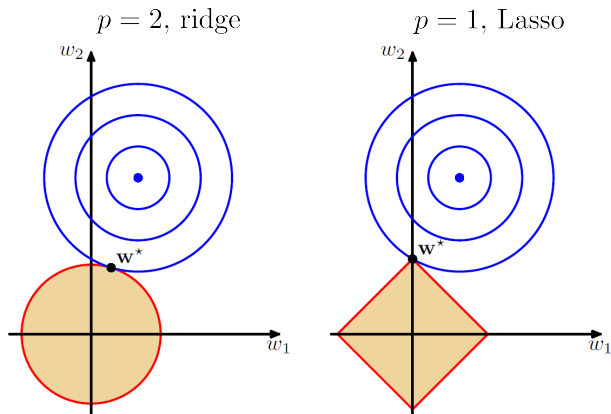
Ridge regression:

$$\text{Objective} = \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \sum_{j=1}^d w_j^2$$

REGULARIZATION

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left\{ \sum_{i=1}^n (y_i - w_0 - \sum_{j=1}^d w_j x_{ij})^2 \right\}$$

subject to $\sum_{j=1}^d |w_j|^p \leq t$



Picture from Bishop's textbook (Chapter 3).

MAP ESTIMATION

$$p(\mathbf{w}|\mathcal{D}) \propto p(\mathcal{D}|\mathbf{w}) \cdot p(\mathbf{w})$$

$$p(\mathcal{D}|\mathbf{w}) = \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n e^{-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w})}$$

likelihood

$$p(\mathbf{w}) = \frac{1}{(2\pi\alpha^2)^{(d+1)/2}} e^{-\frac{1}{2\alpha^2}\mathbf{w}^T\mathbf{w}}$$

prior

$$\mathbf{w}_{\text{MAP}} = \arg \min_{\mathbf{w}} \left\{ (\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w}) + \frac{\sigma^2}{\alpha^2}\mathbf{w}^T\mathbf{w} \right\}$$

PERFORMANCE OF REGRESSION MODELS

Given: training set \mathcal{D} , large test set \mathcal{T} , and model $f(x)$ trained on \mathcal{D}

Goal: estimate accuracy of $f(x)$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2$$

$f(x)$ = prediction on x
 y = observed target for x
 n = number of examples in \mathcal{T}
 \bar{y} = mean of the target

$$R^2 = 1 - \frac{\sum_{i=1}^n (f(x_i) - y_i)^2}{\sum_{i=1}^n (\bar{y} - y_i)^2}$$

$$\text{MSE} \in [0, \infty)$$

$$R^2 \in (-\infty, 1]$$

R^2 = percentage of variance “explained” by $f(x)$. Target mean is the trivial predictor.

PERFORMANCE OF REGRESSION MODELS

Given: training set \mathcal{D} , large test set \mathcal{T} , and model $f(x)$ trained on \mathcal{D}

Goal: estimate accuracy of $f(x)$

Pearson's correlation:

$$\rho = \frac{\sum_{i=1}^n (f_i - \bar{f})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (f_i - \bar{f})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

f_i = prediction $f(x_i)$

y_i = observed target for x_i

n = number of examples in \mathcal{T}

\bar{y} = mean of the target

\bar{f} = mean of the predictions

Kendall's tau:

$$\tau_K = \frac{2}{n(n-1)} \sum_{i < j} \text{sign}(f(x_i) - f(x_j)) \cdot \text{sign}(y_i - y_j)$$

$$\rho \in [-1, 1]$$

$$\tau_K \in [0, 1]$$