

Multivariate classification

Balázs Kégl

LAL / University of Paris-Sud

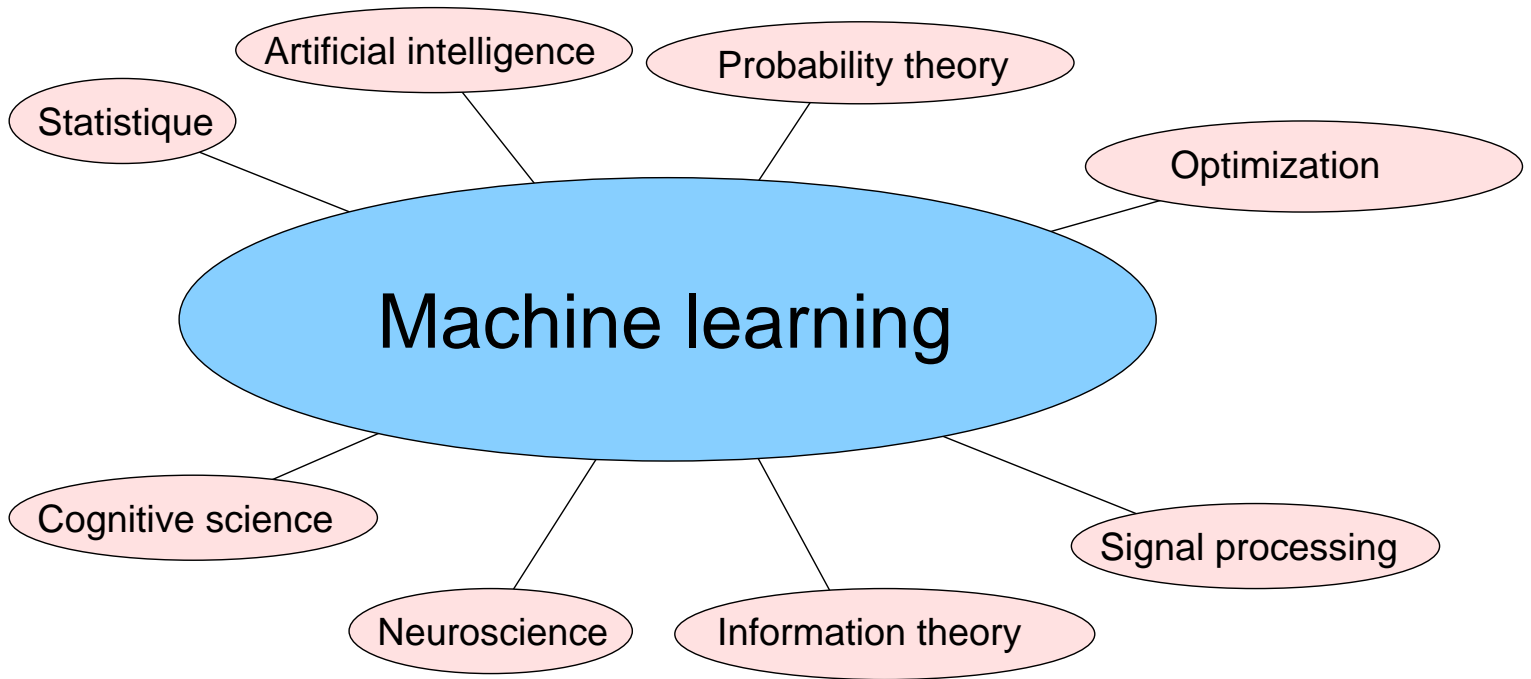
School of Statistics

20 May, 2010

Machine learning

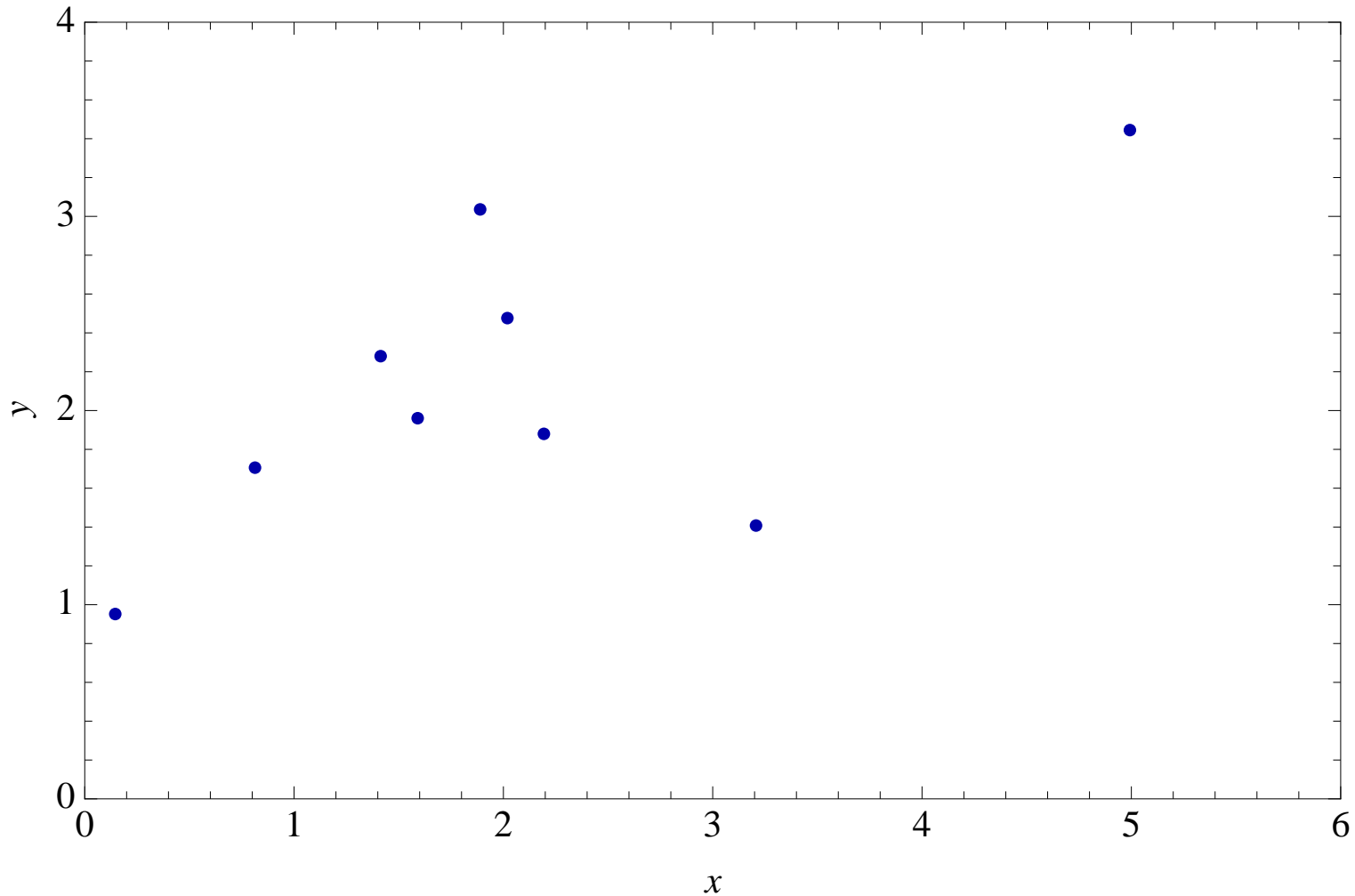
- Can be considered as a **sub-domain of statistics**
 - Often **non-parametric**, **high-dimensional** spaces, **large data sets** → **computational** issues (optimization) become as important as statistical issues (capacity control)
 - Most often associated to **multivariate discrimination (classification)**
- Most known **algorithms**
 - **Classification**: neural network, boosting, support vector machine
 - **Regression**: neural network, Gaussian process
 - **Density estimator**: mixture of Gaussians
 - **Clustering, dimensionality reduction**: k-means, spectral clustering, local linear embedding, ISOMAP, nonlinear PCA, kernel PCA

Machine learning at the crossroads

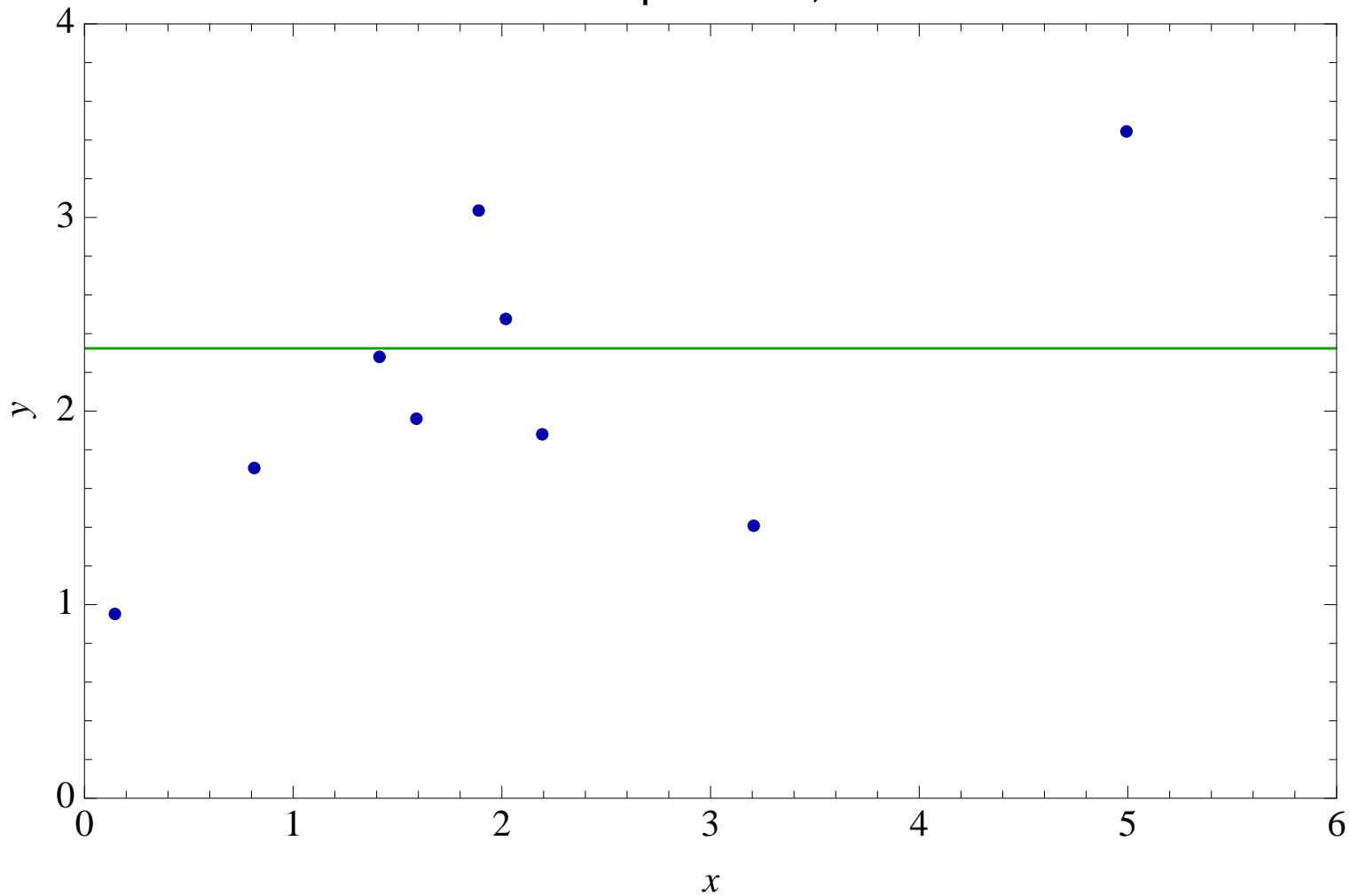


- Non-parametric fitting: two simple examples
- The formal model for classification, learning principles
- Classification algorithms (from a user's point of view)
 - perceptron, neural networks (NN)
 - AdaBoost
 - the Support Vector Machine (SVM)

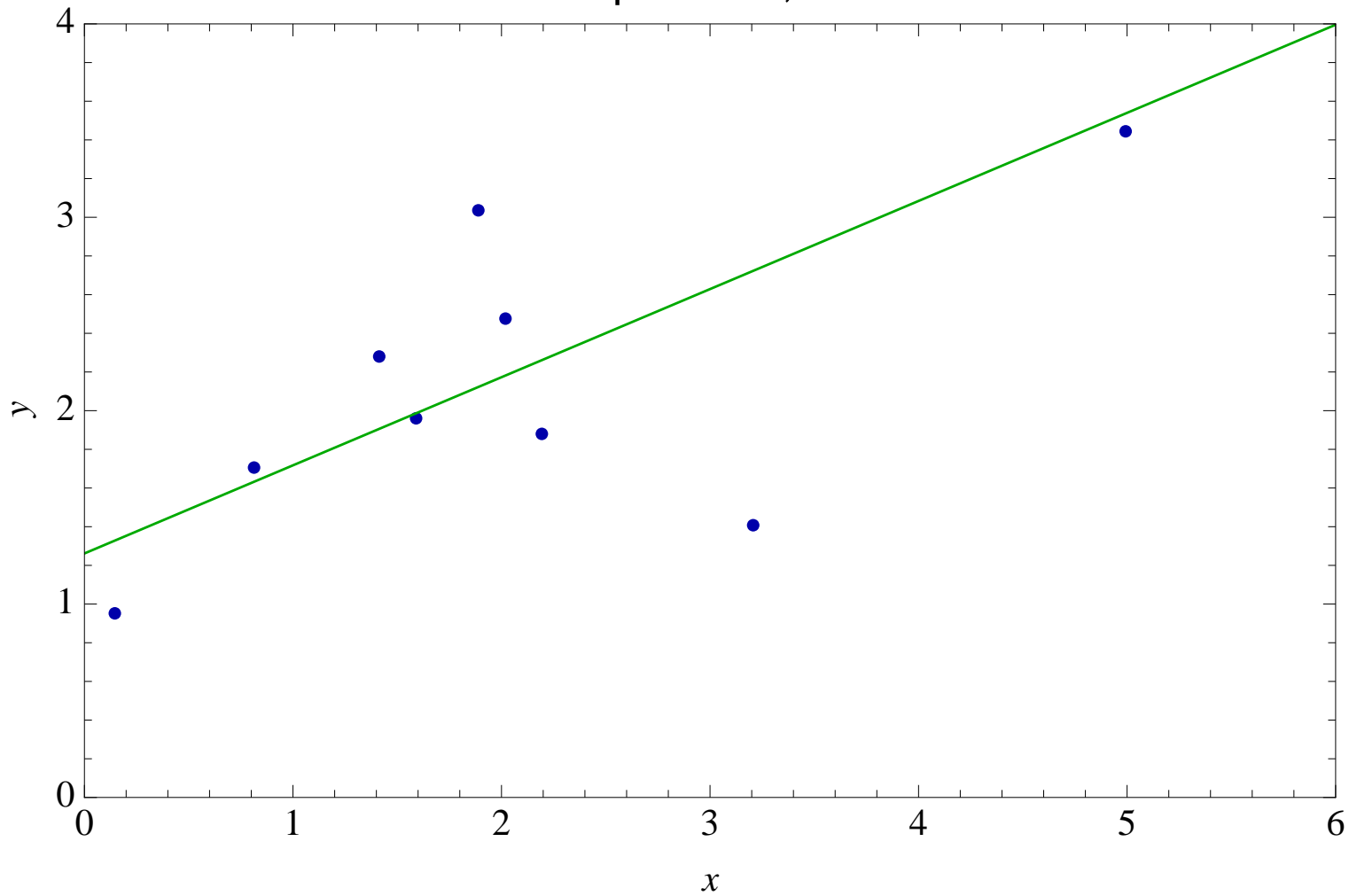
Data generated from an unknown function with unknown noise



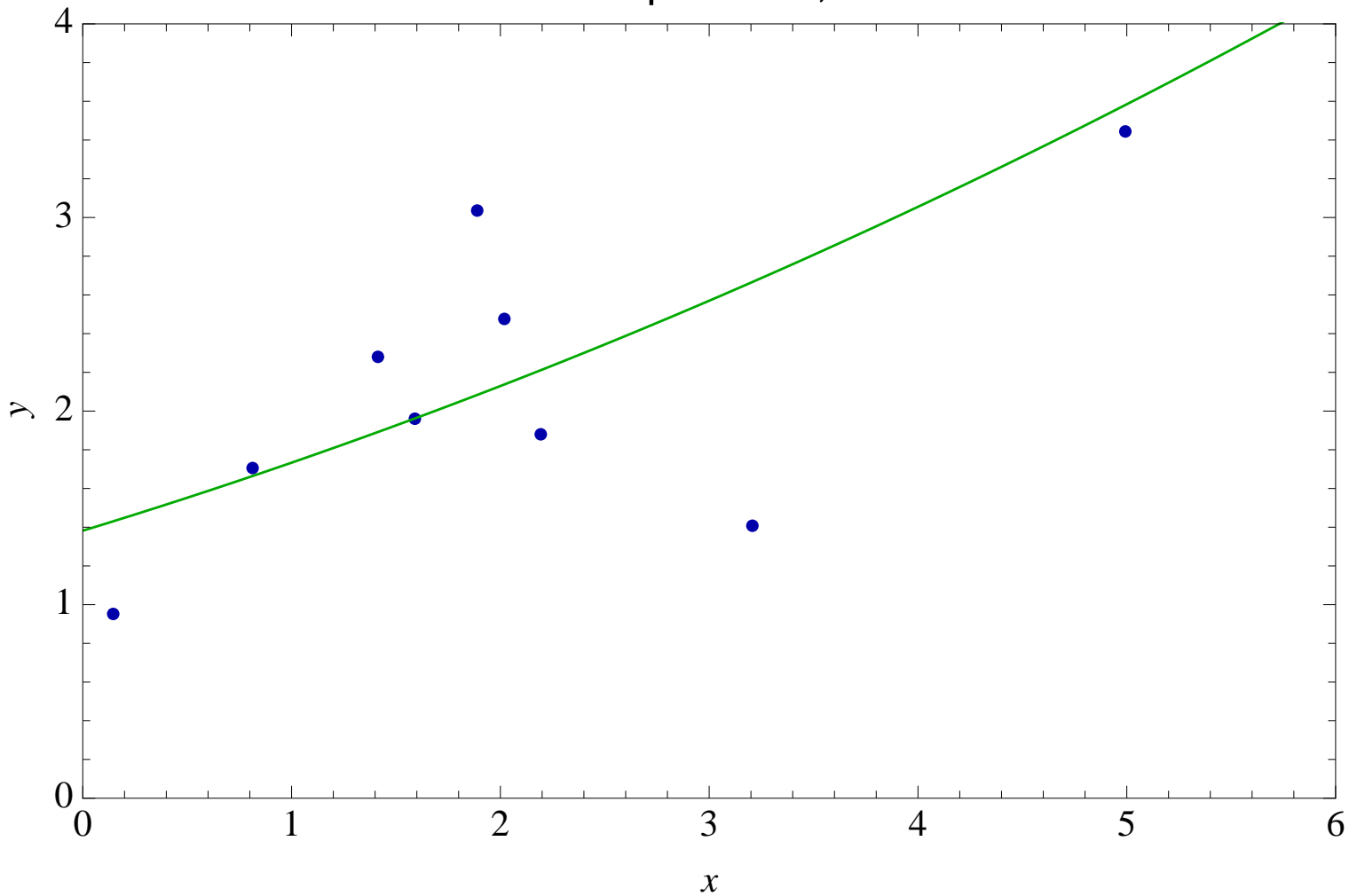
Constant least squares fit, RMSE = 0.915



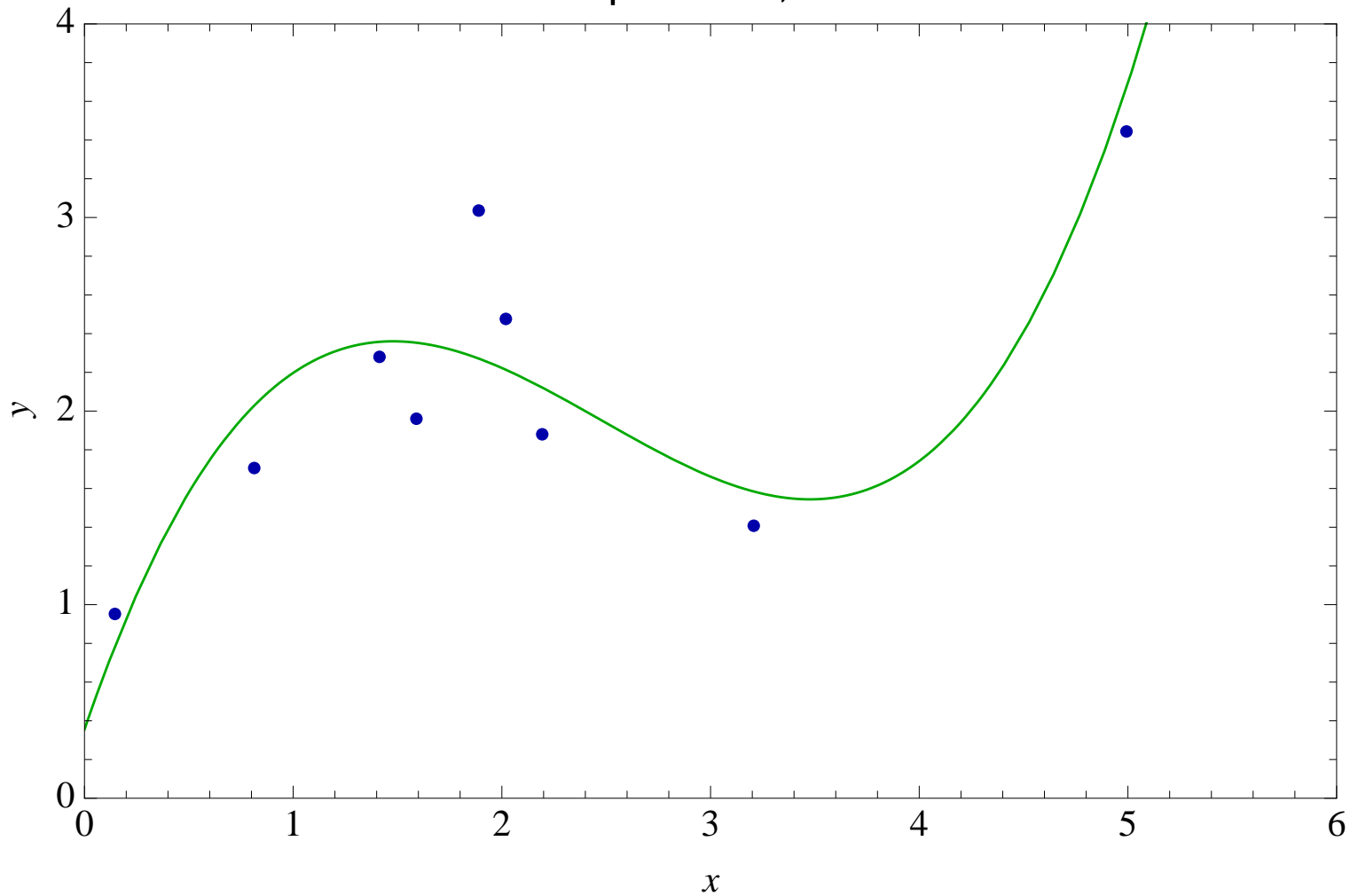
Linear least squares fit, RMSE = 0.581



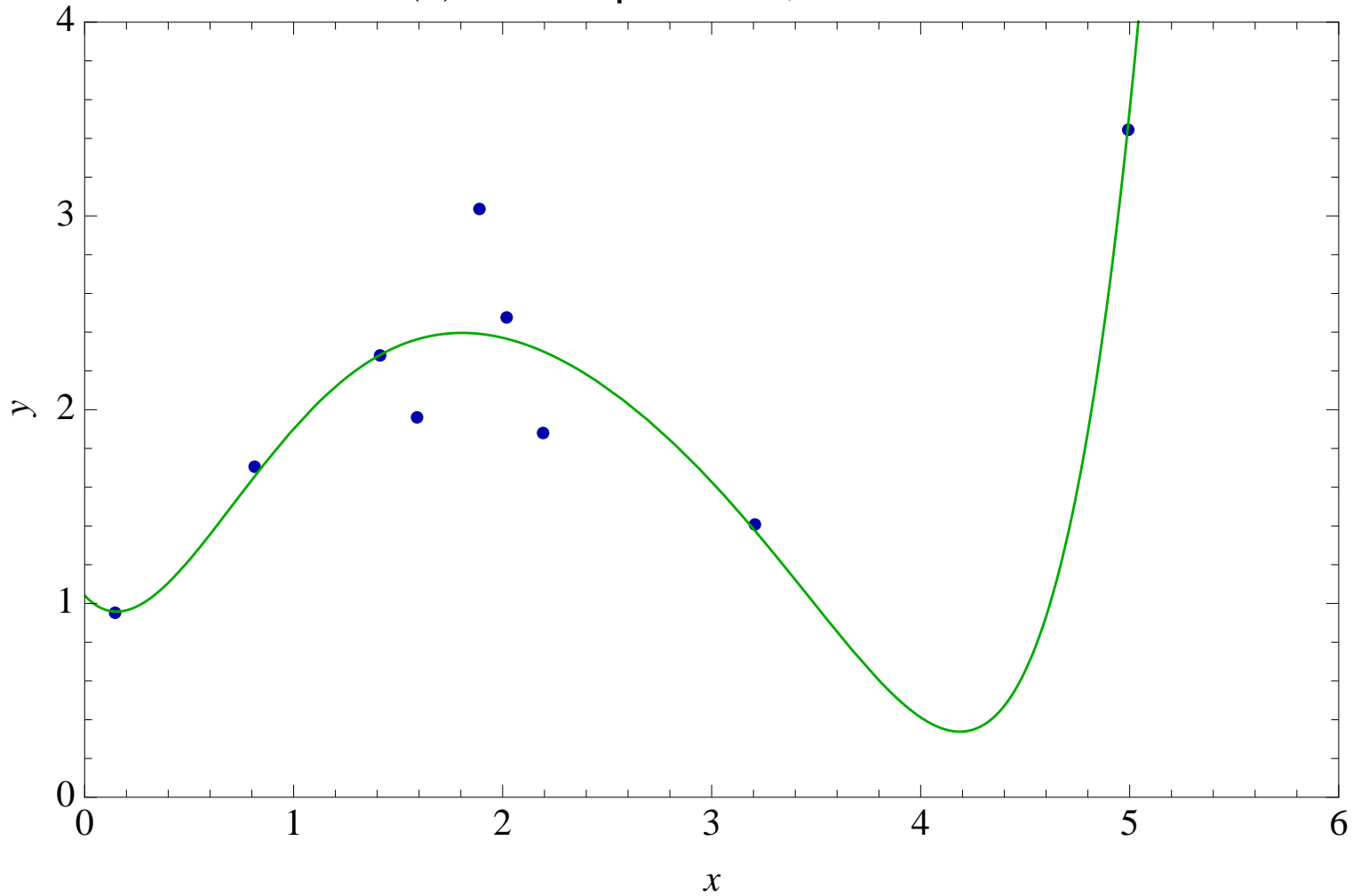
Quadratic least squares fit, RMSE = 0.579



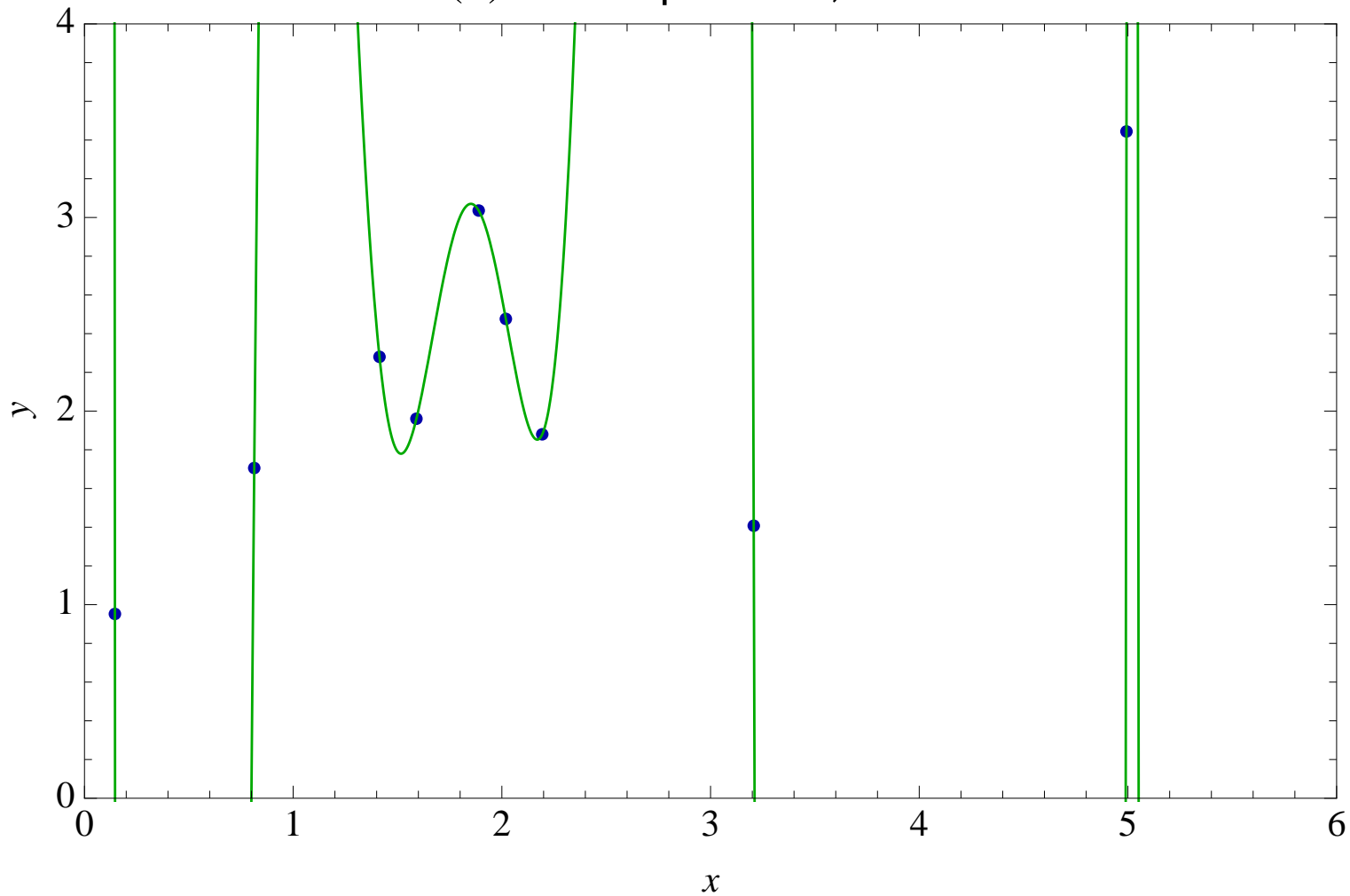
Cubic least squares fit, RMSE = 0.339



Poli(6) least squares fit, RMSE = 0.278



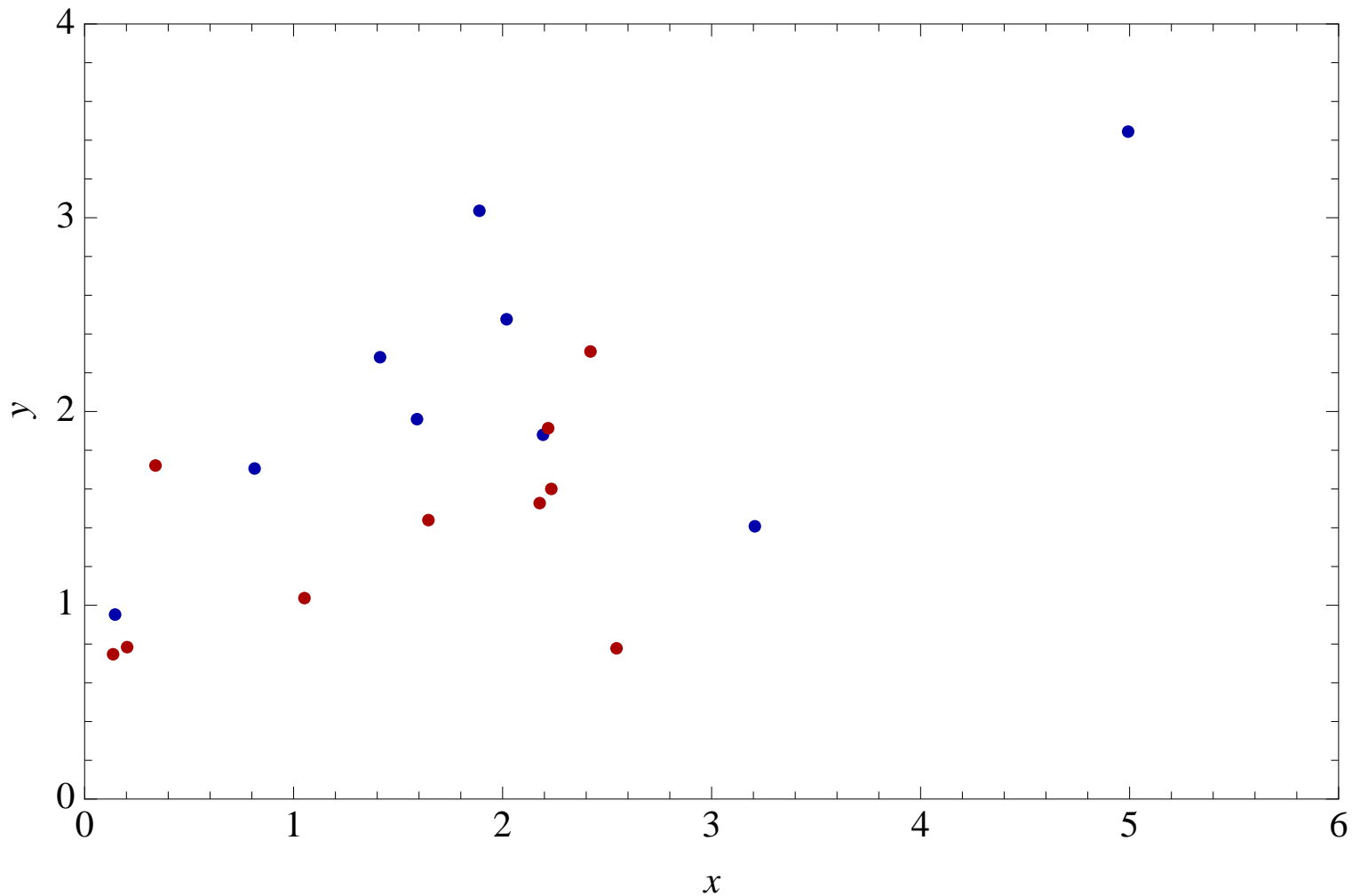
Poli(9) least squares fit, RMSE = 0



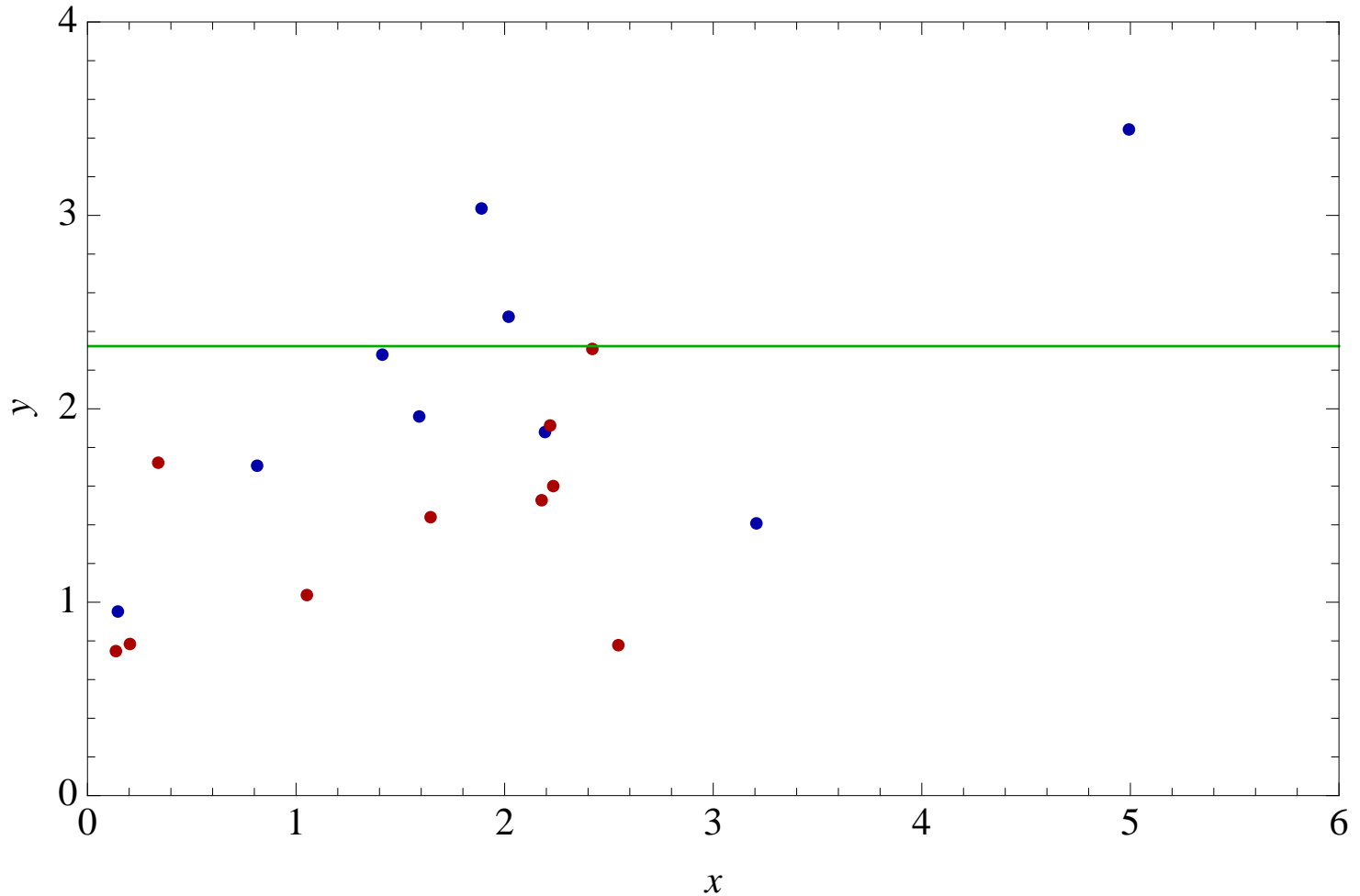
Non-parametric fitting

- Minimizing the **training** cost (RMSE) does not work if the function **class is not fixed beforehand**
- If you **do not know** the correct function **class**, you should **not fix it**
- Capacity control, regularization
 - **theoretical** results, **data-independent**, based on **sample size** and **measures of complexity**
 - use **independent (test) set**

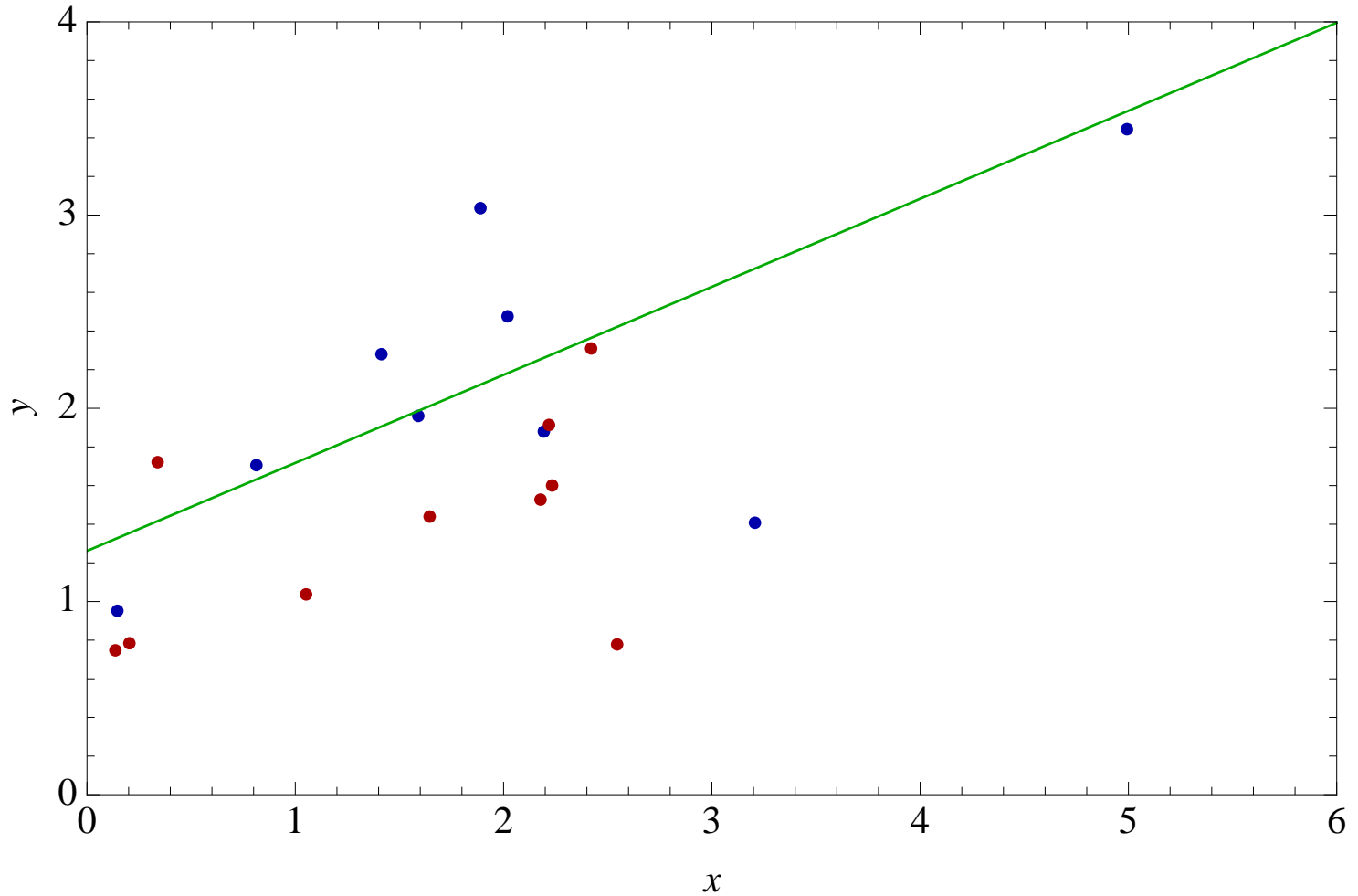
Data generated from an unknown function with unknown noise



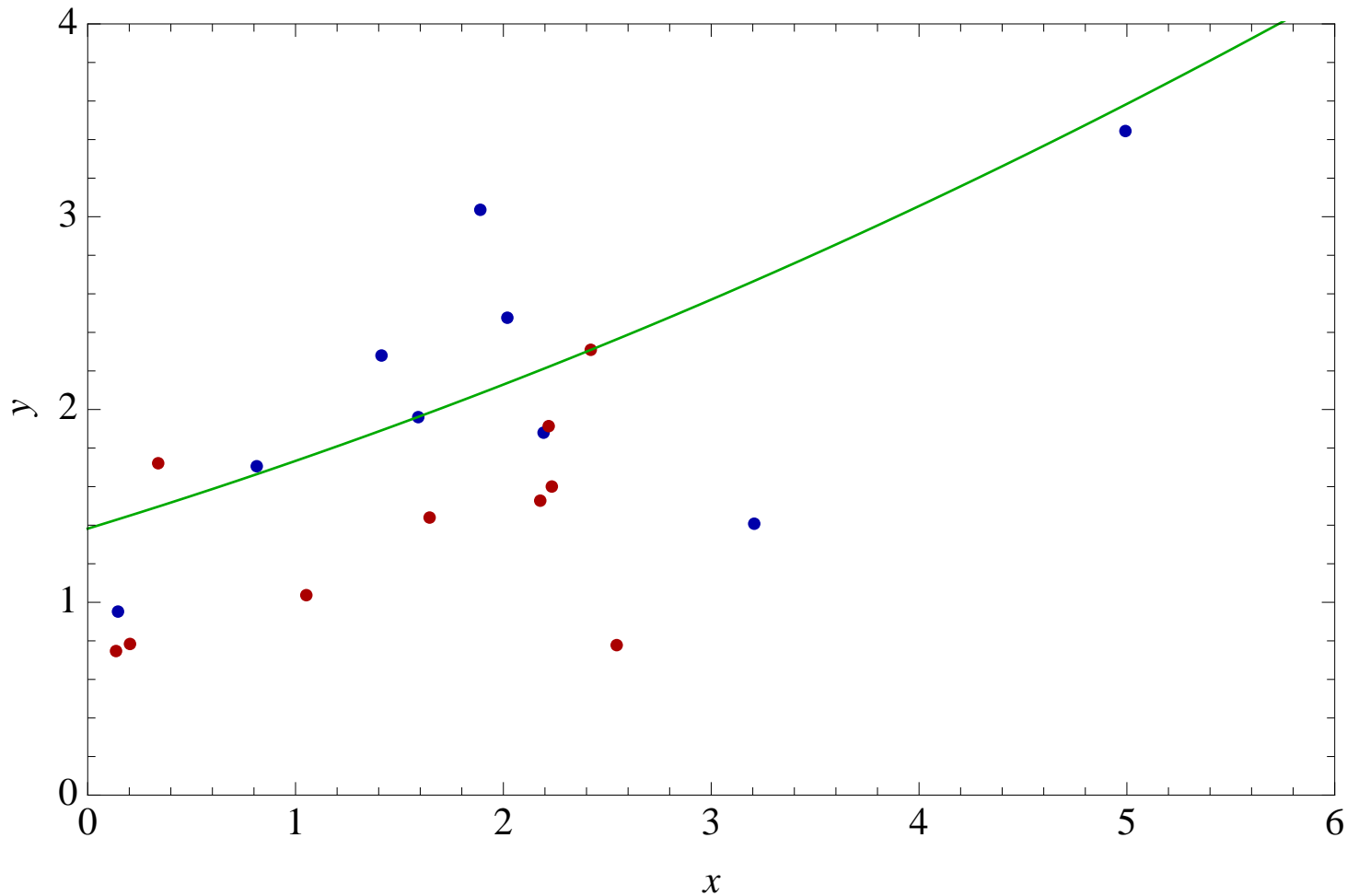
Const. least squares fit, **training** RMSE = 0.915, **test** RMSE = 1.067



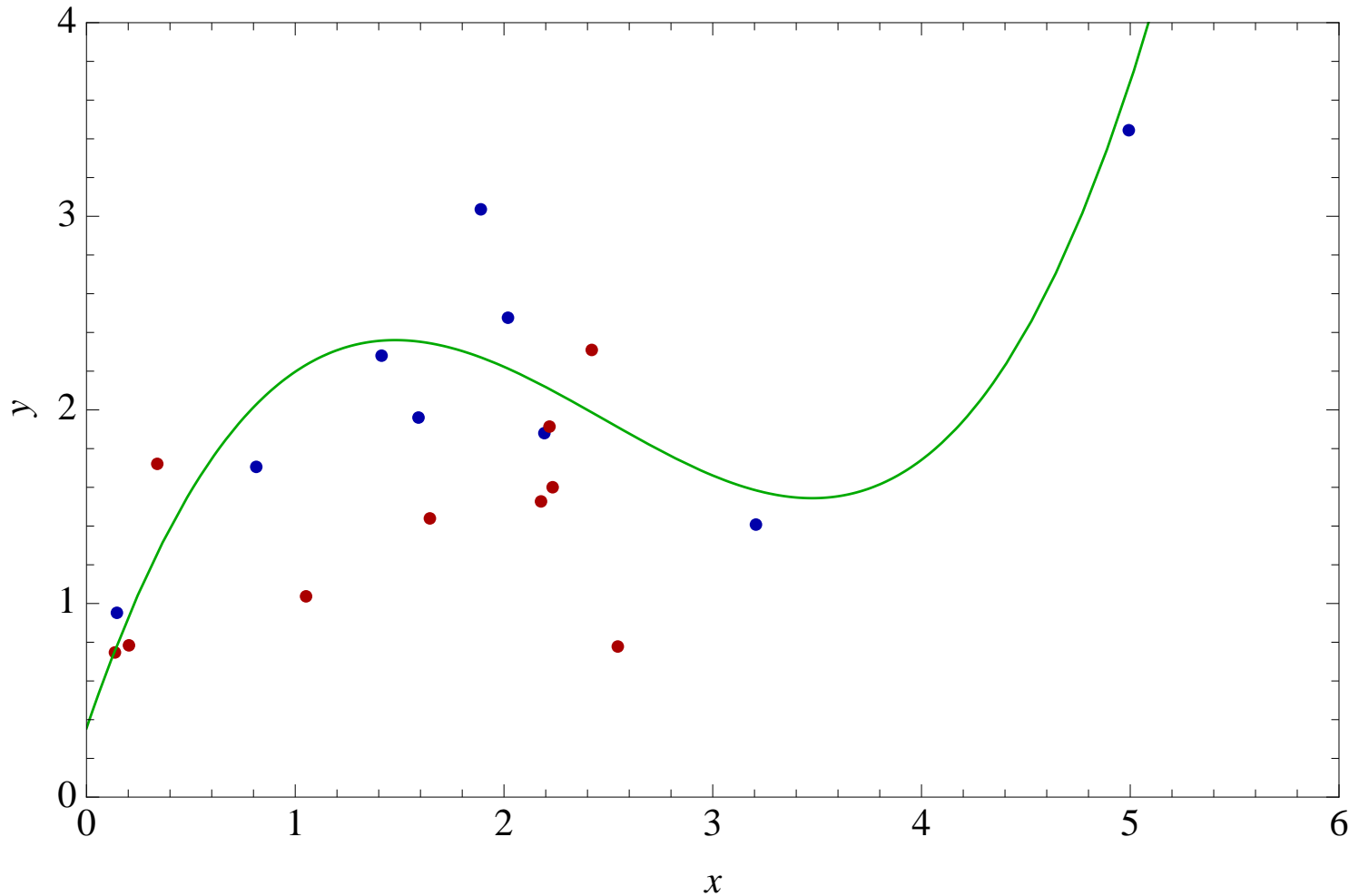
Linear least squares fit, **training** RMSE = 0.581, **test** RMSE = 0.734



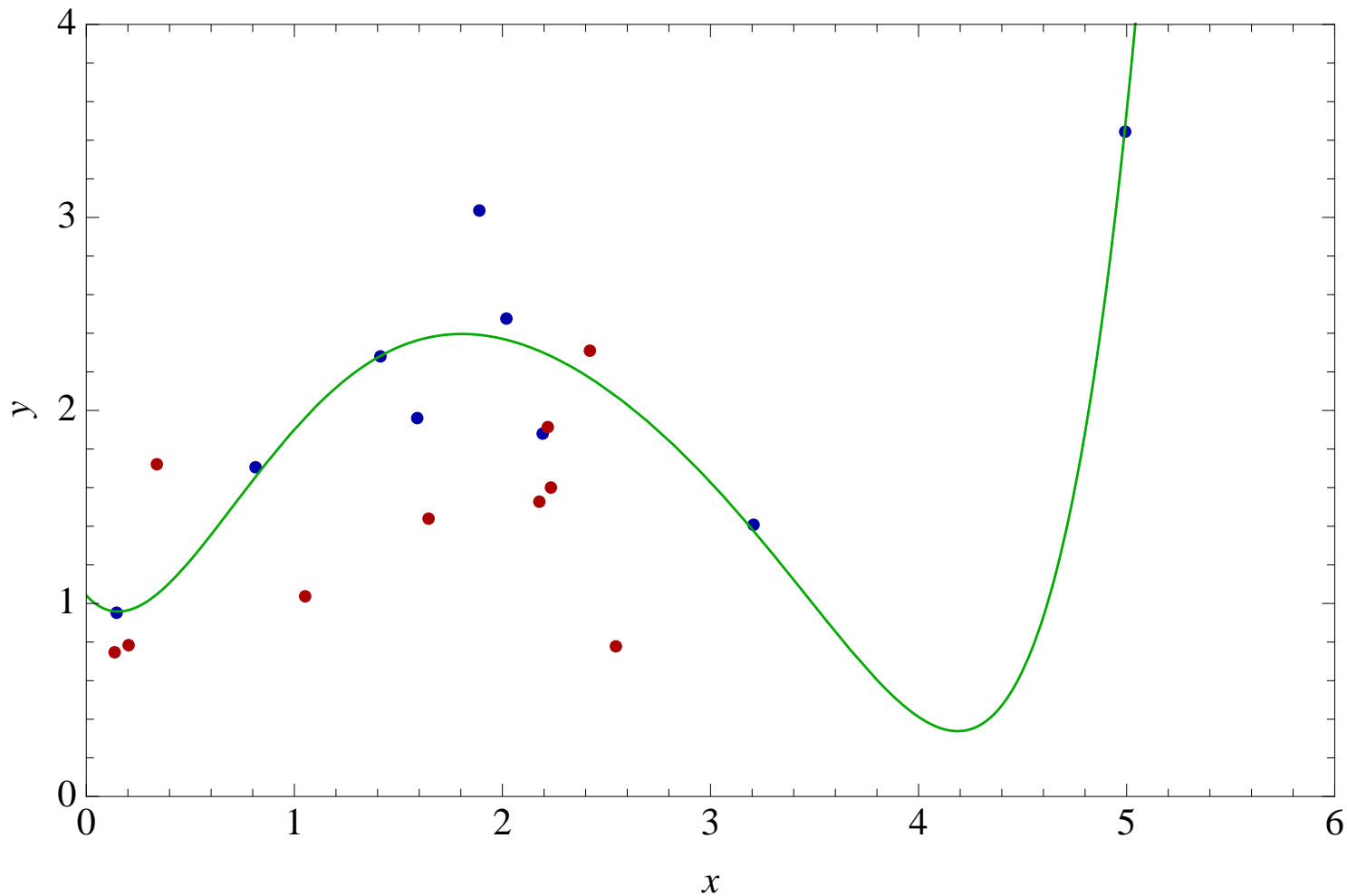
Quadr. least squares fit, **training** RMSE = 0.579, **test** RMSE = 0.723



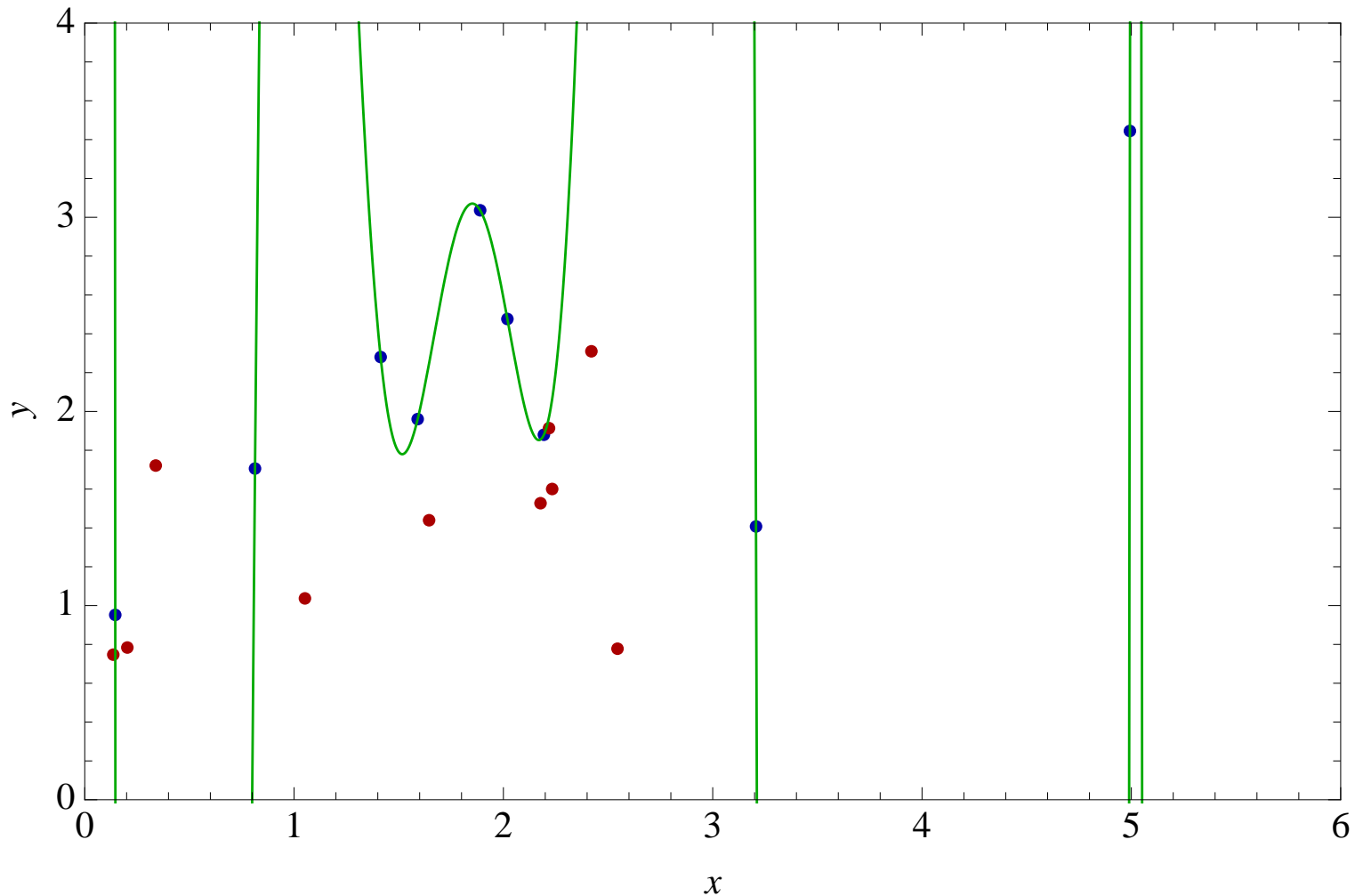
Cubic least squares fit, **training** RMSE = 0.339, **test** RMSE = 0.672



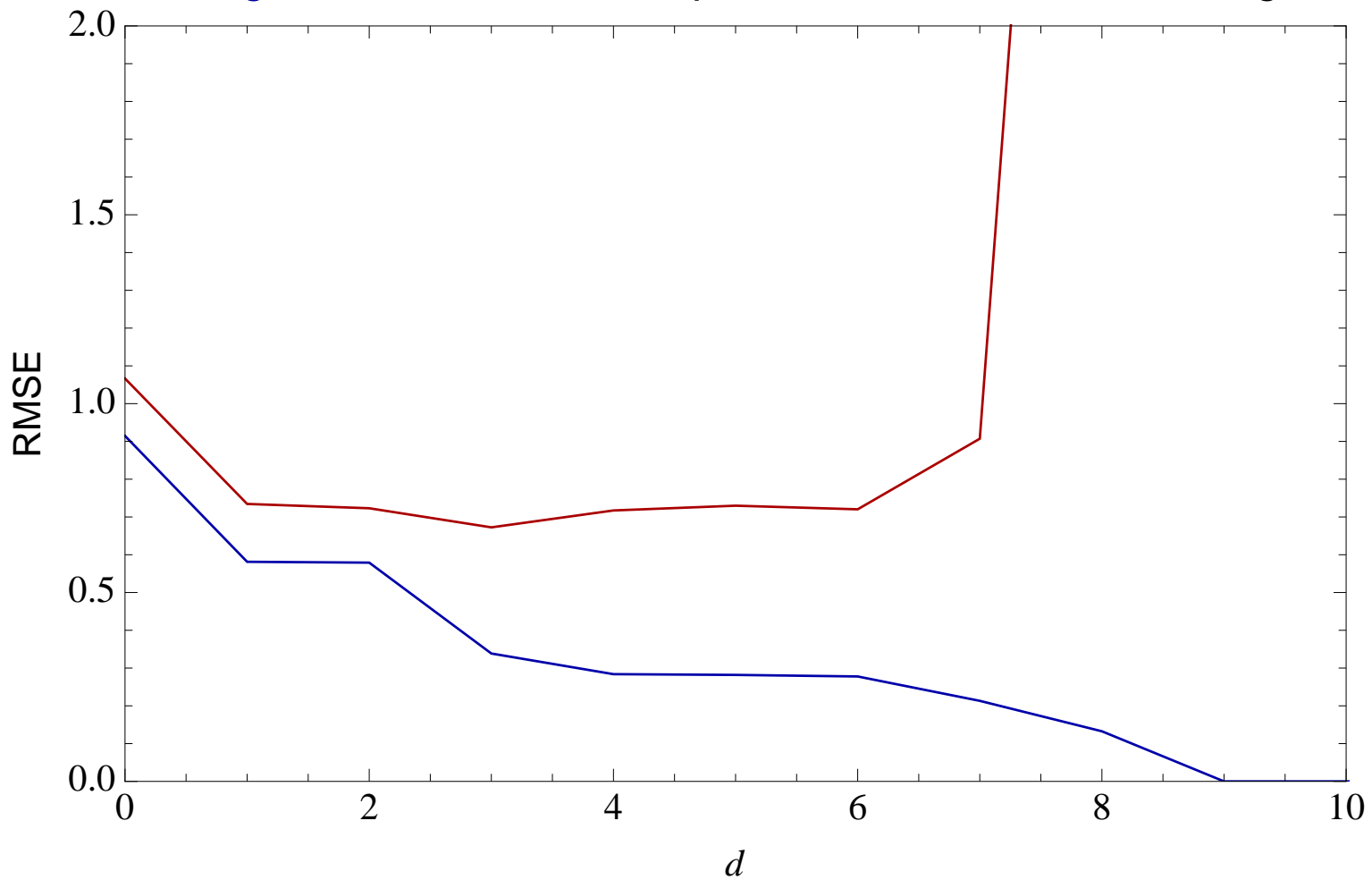
Poli(6) least squares fit, **training** RMSE = 0.278, **test** RMSE = 0.72



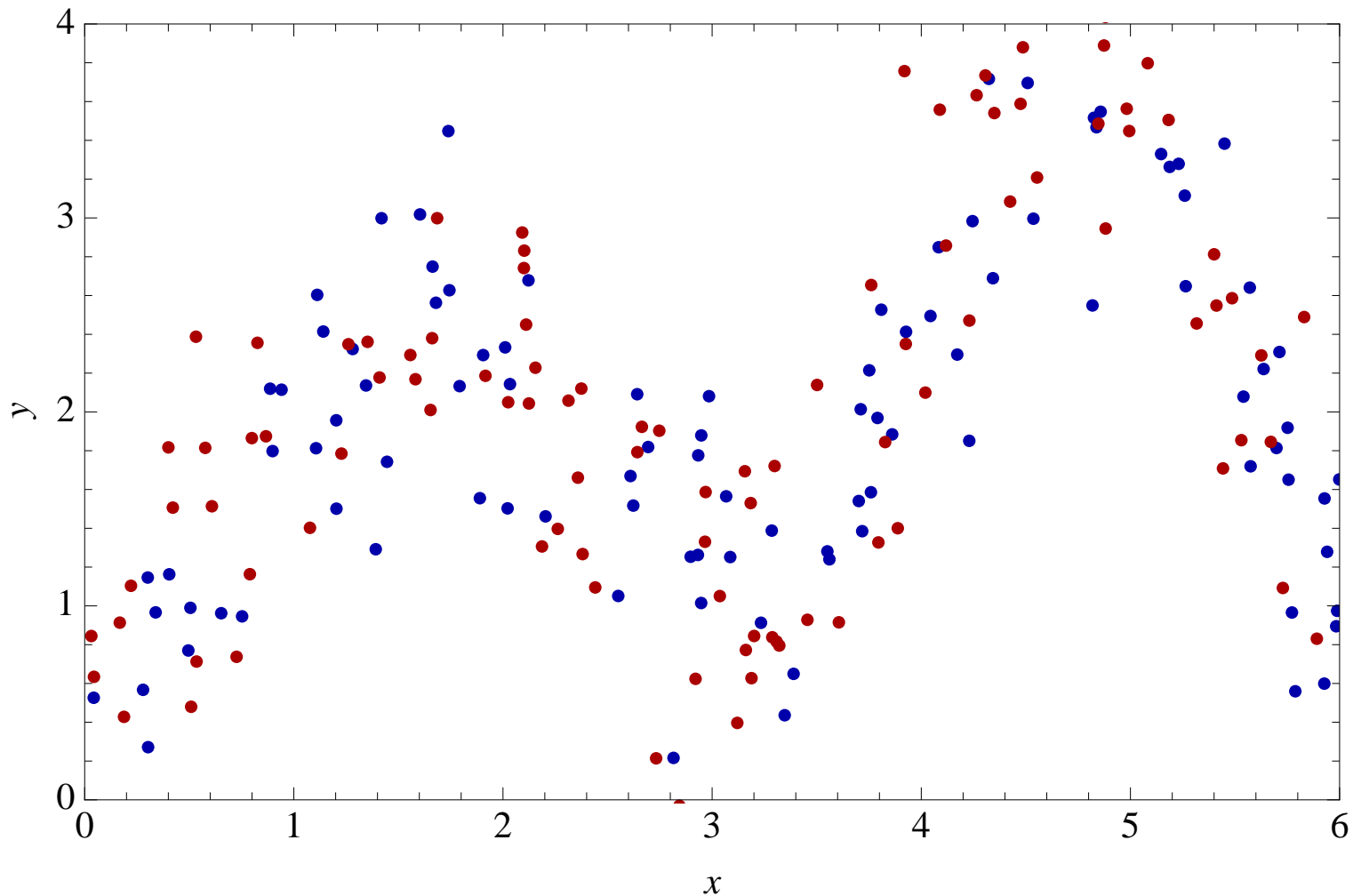
Poli(9) least squares fit, **training** RMSE = 0, **test** RMSE = 46.424



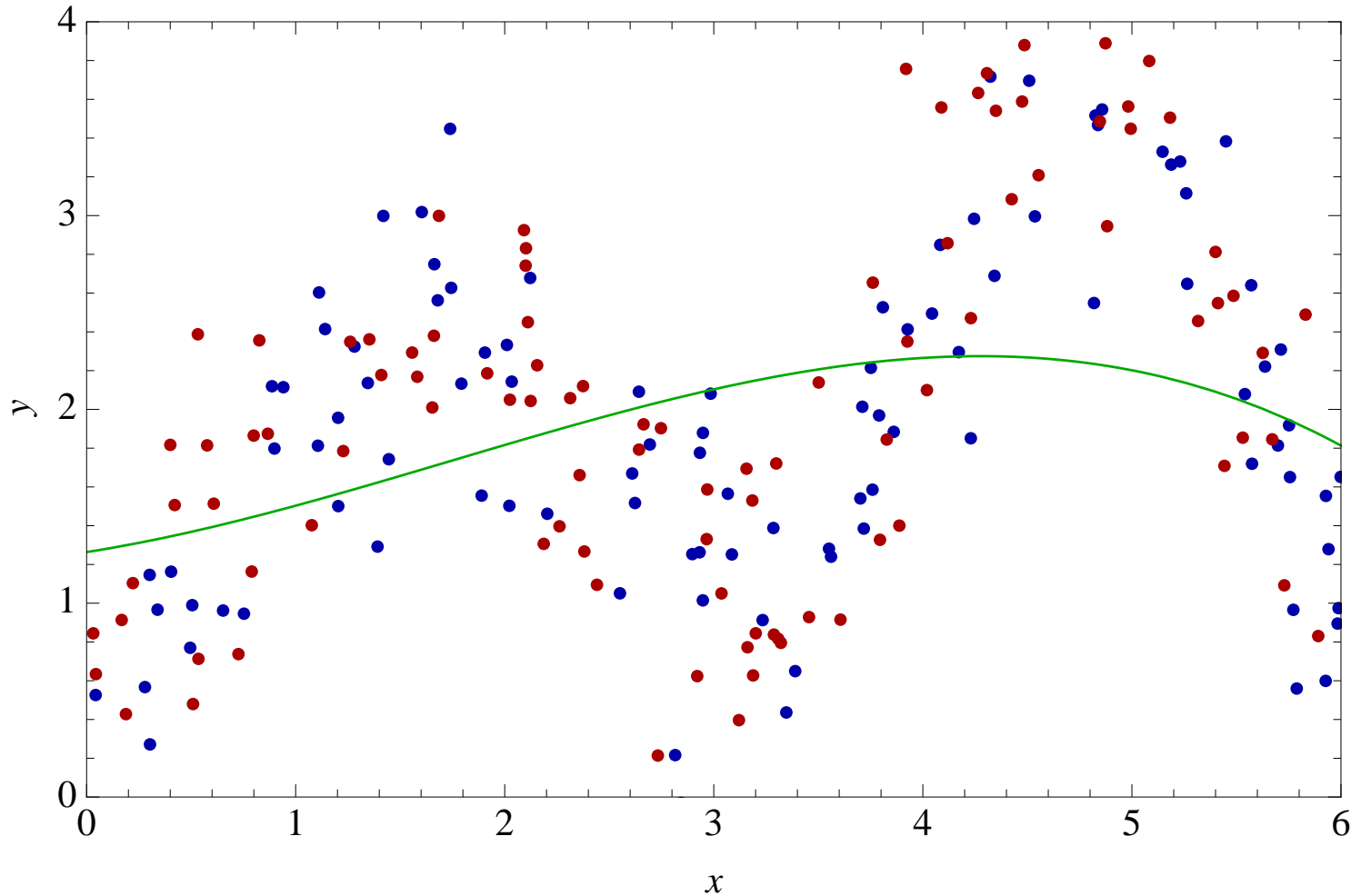
Training and test RMSE's for polinomial fits of different degrees



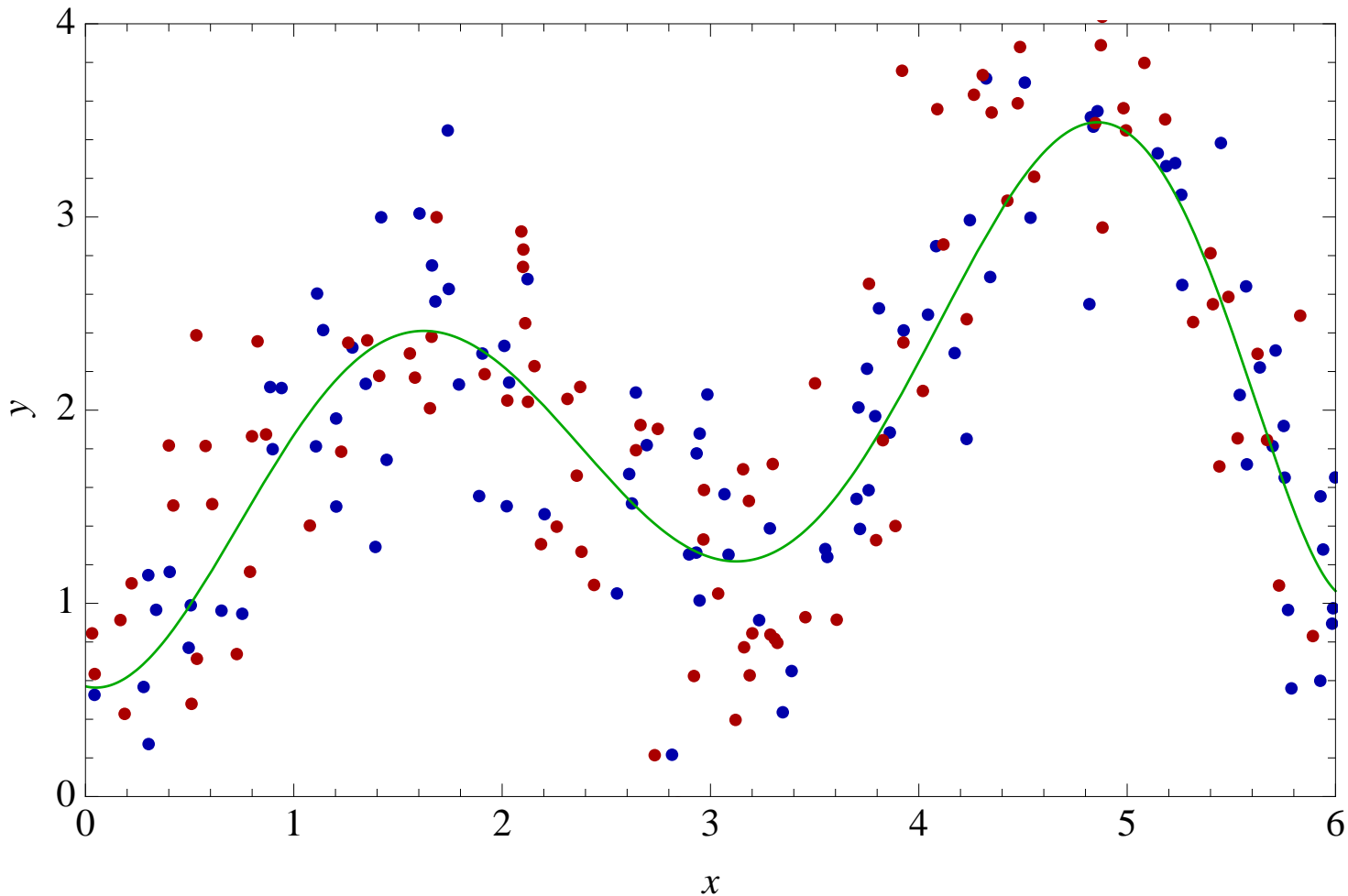
Data generated from an unknown function with unknown noise



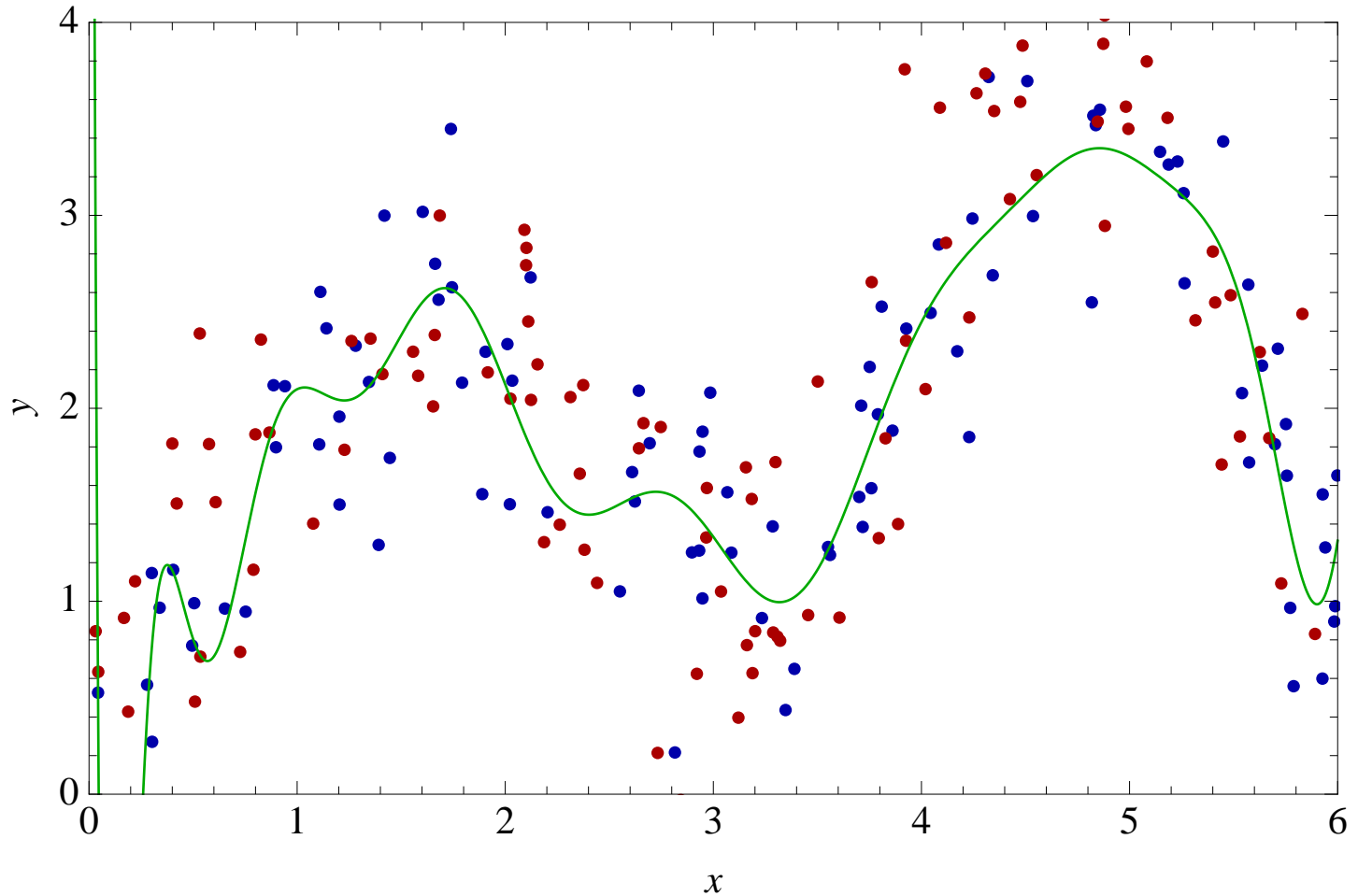
Cubic least squares fit, **training** RMSE = 0.793, **test** RMSE = 0.913



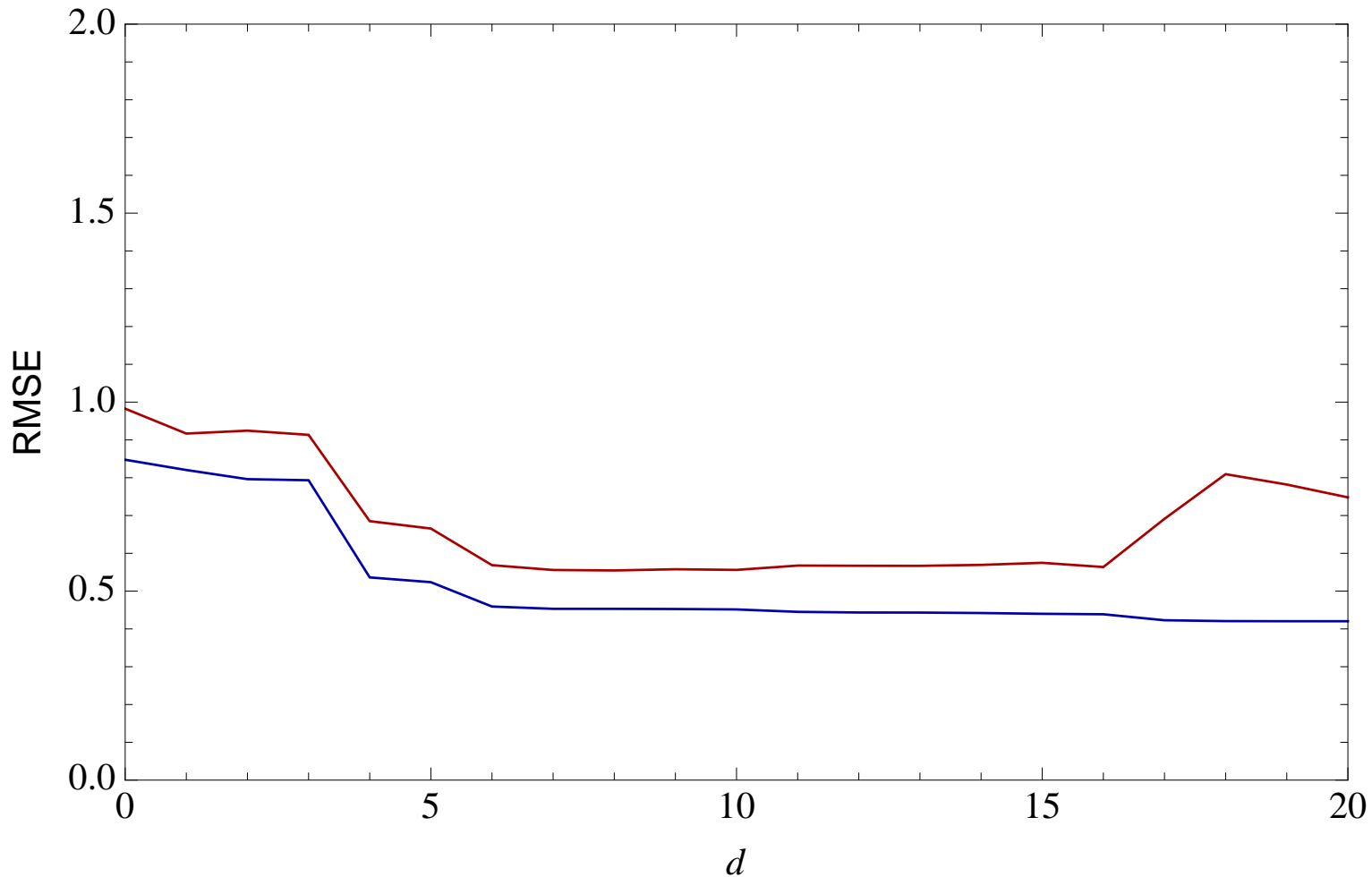
Poli(8) least squares fit, **training** RMSE = 0.453, **test** RMSE = 0.555



Poli(18) least squares fit, **training** RMSE = 0.421, **test** RMSE = 0.809



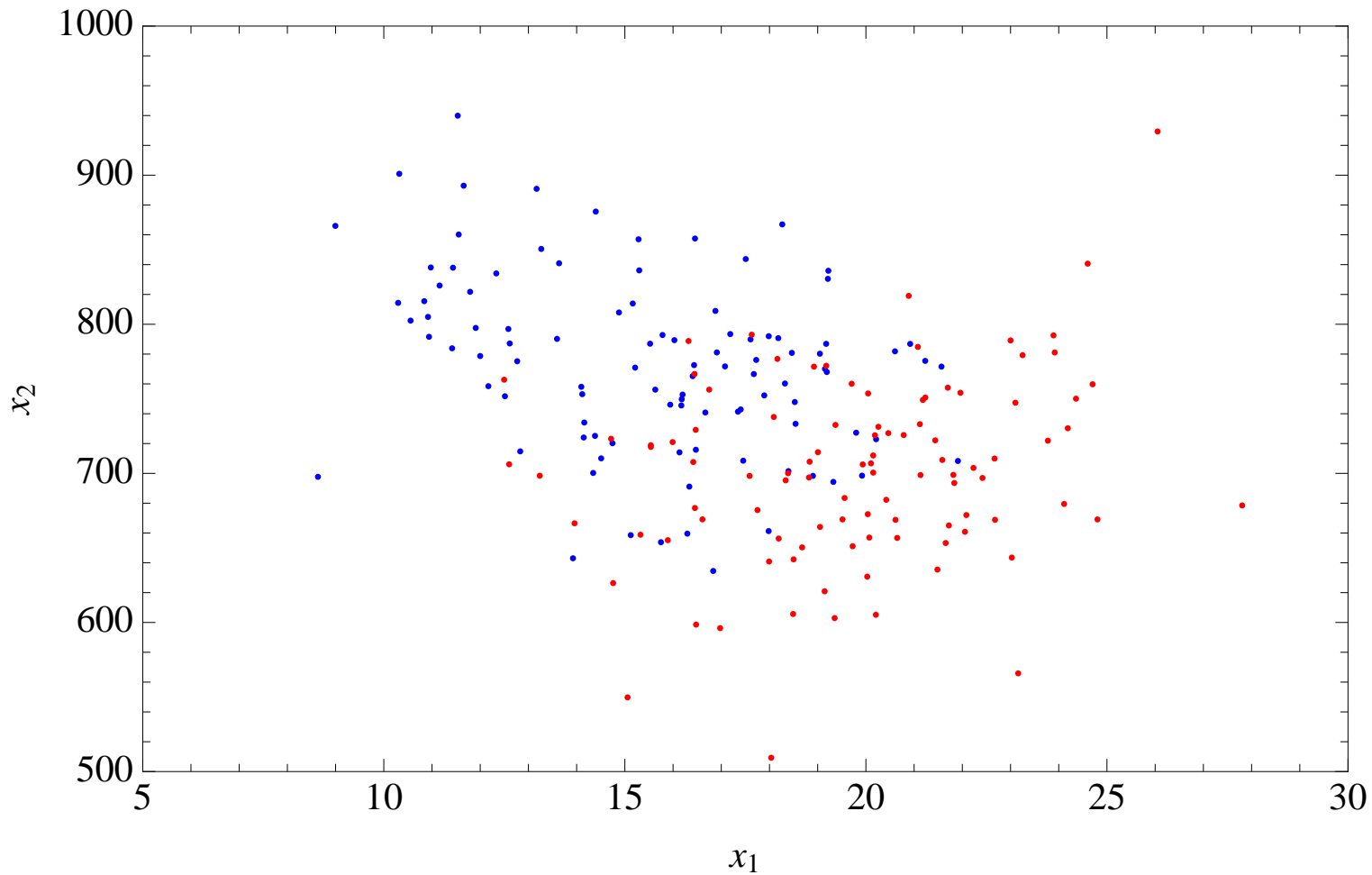
Training and test RMSE's for polinomial fits of different degrees



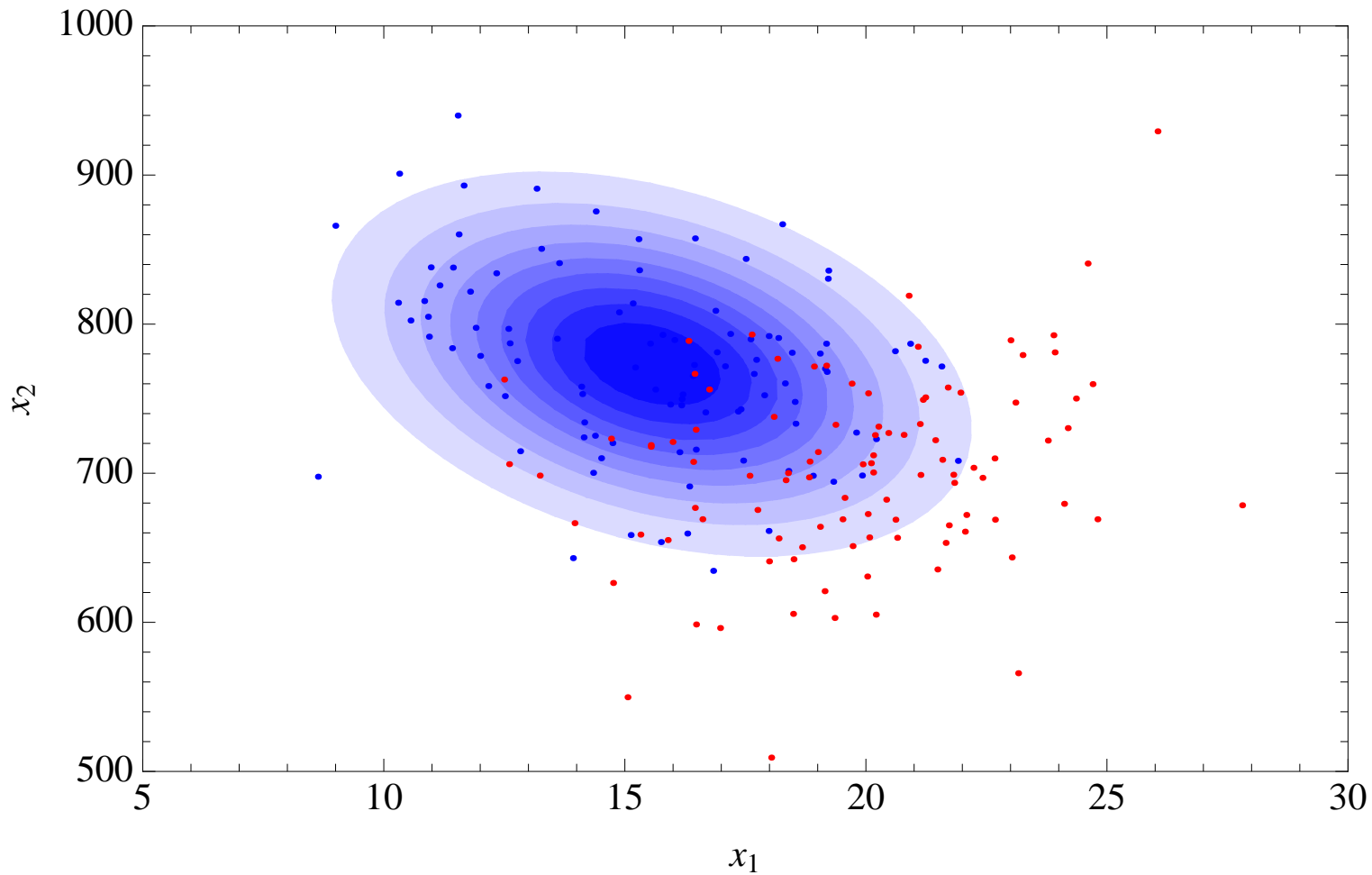
Non-parametric fitting

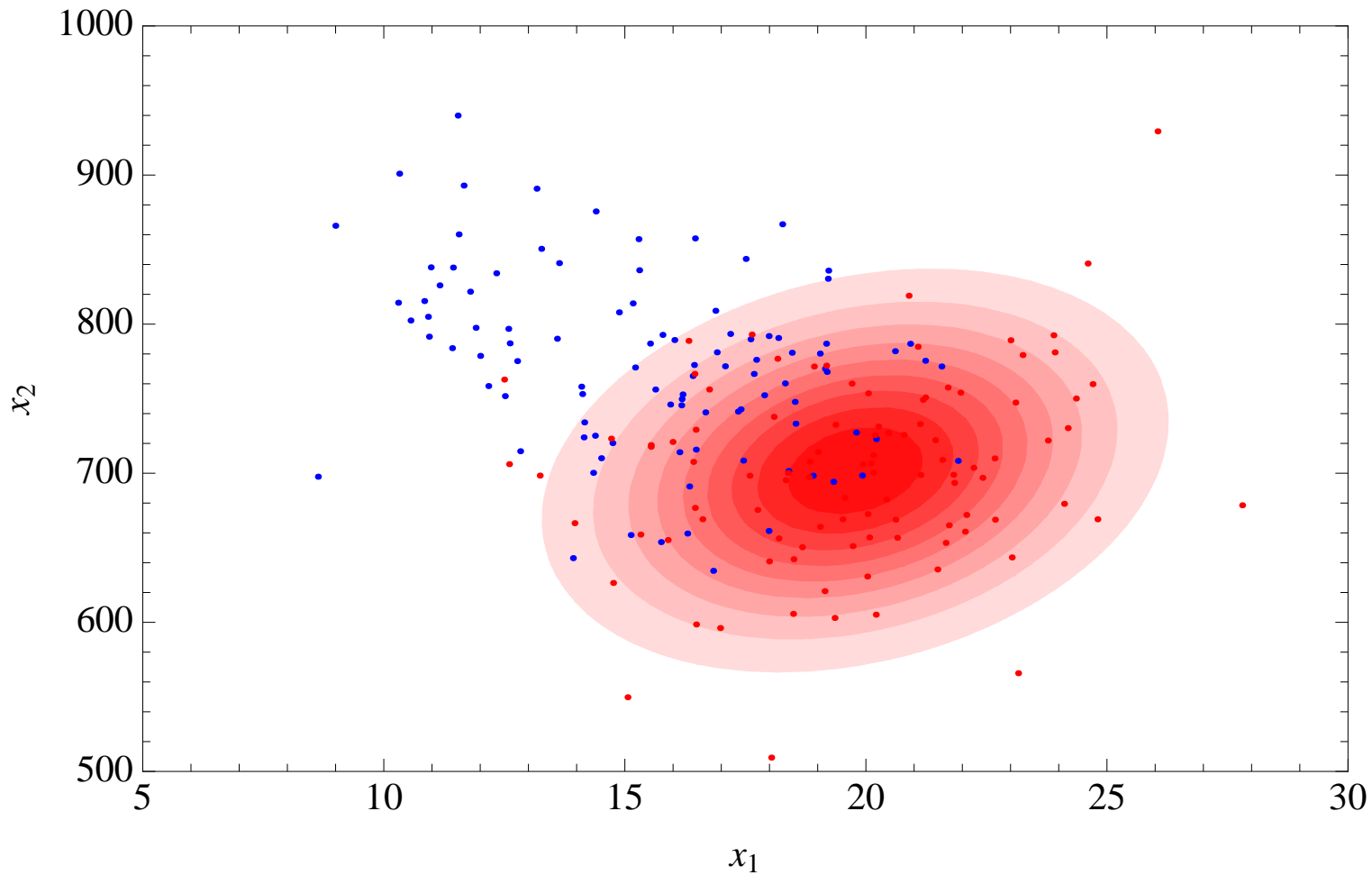
- Capacity control, regularization
 - trade-off between **approximation error** and **estimation error**
 - **complexity** grows with **data size**
 - no need to correctly guess the function class

'Two Gaussians' data for 2-class classification

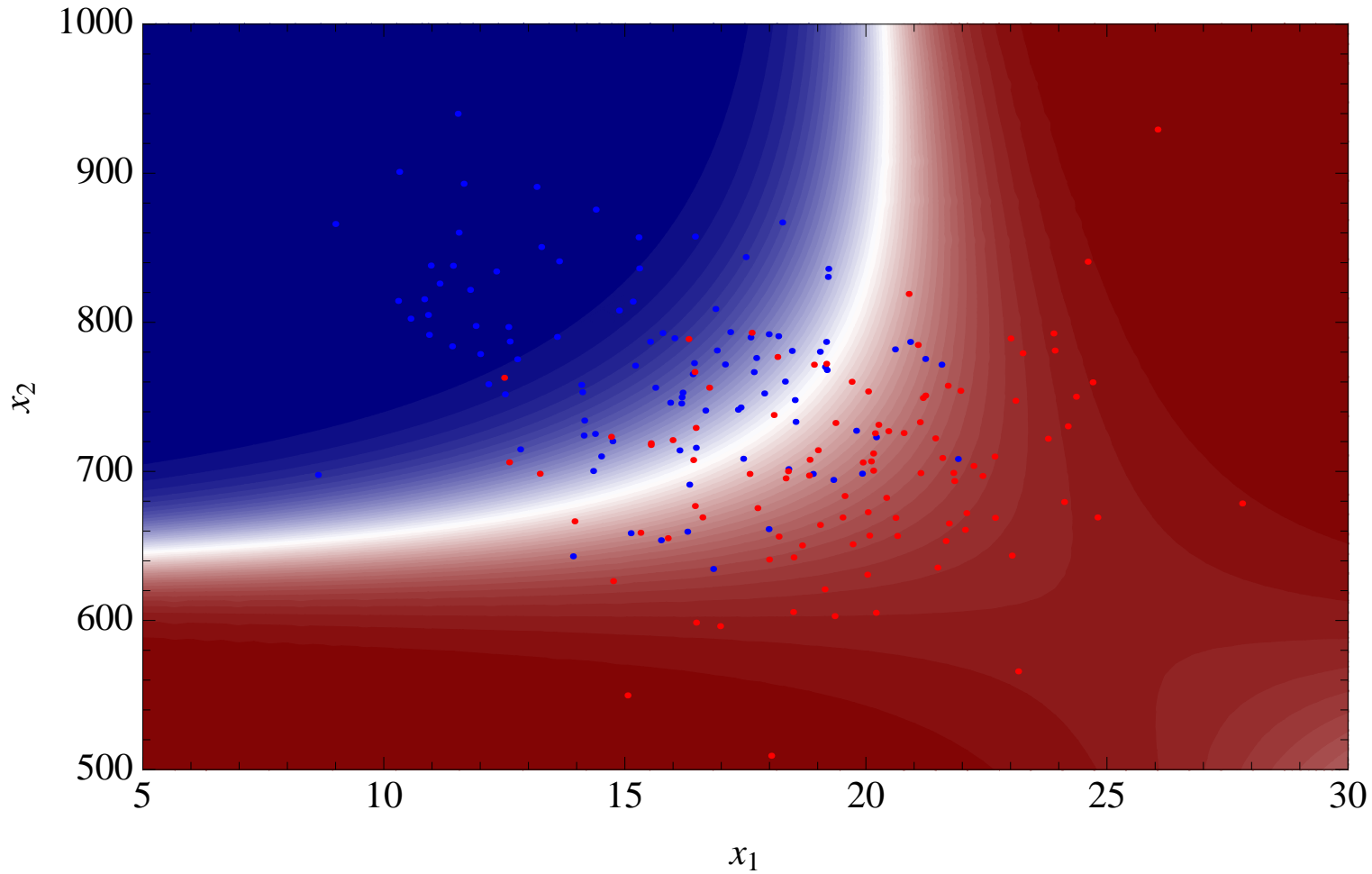


2-D Gaussian fit for class 1



2-D Gaussian fit for **class -1**

Discriminant function with Gaussian density fits



- Terminology

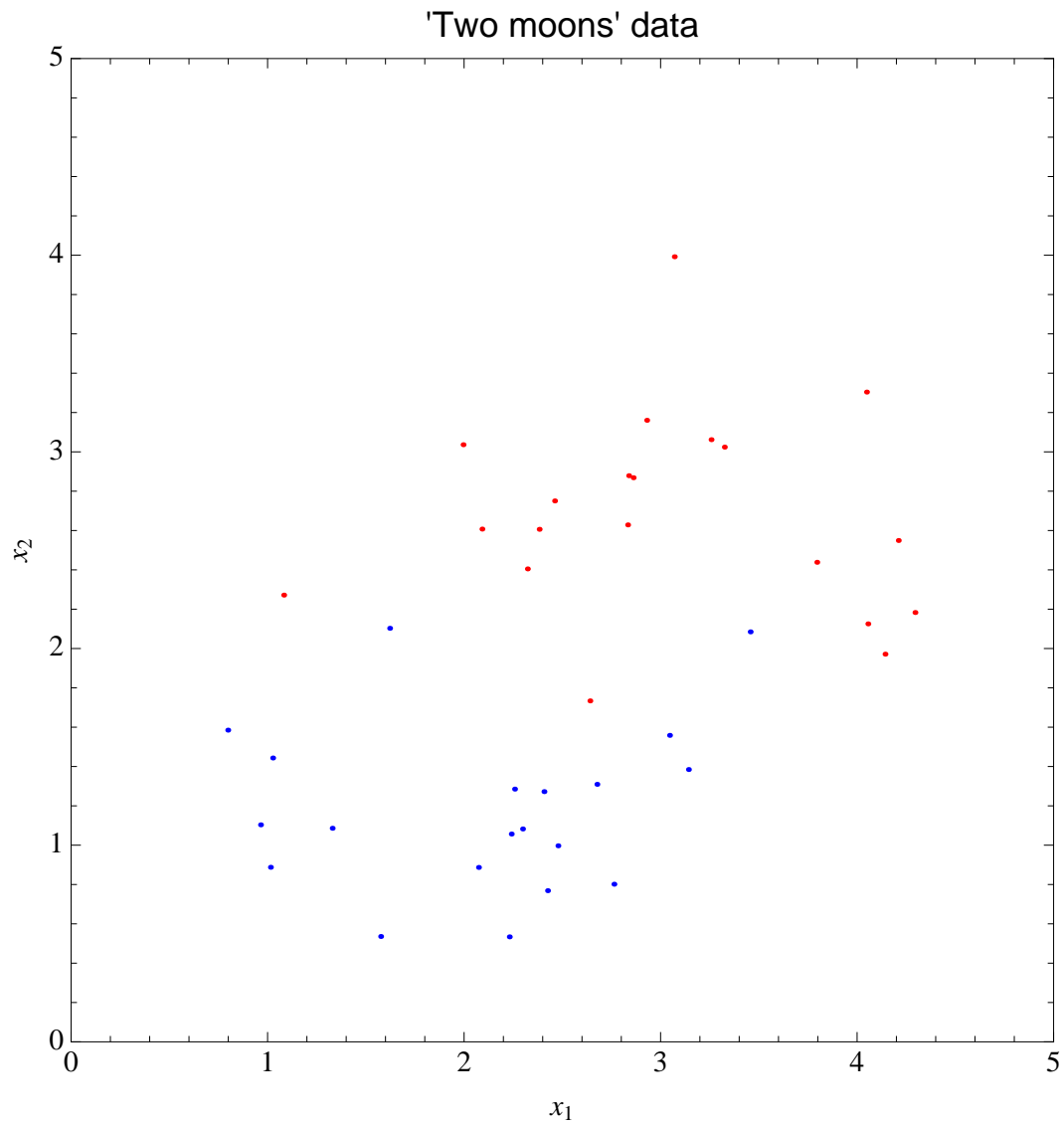
- **Conditional densities:** $p(\mathbf{x}|Y = 1)$, $p(\mathbf{x}|Y = -1)$
- **Prior** probabilities: $P(Y = 1)$, $P(Y = -1)$
- **Posterior** probabilities: $P(Y = 1|\mathbf{x})$, $P(Y = -1|\mathbf{x})$

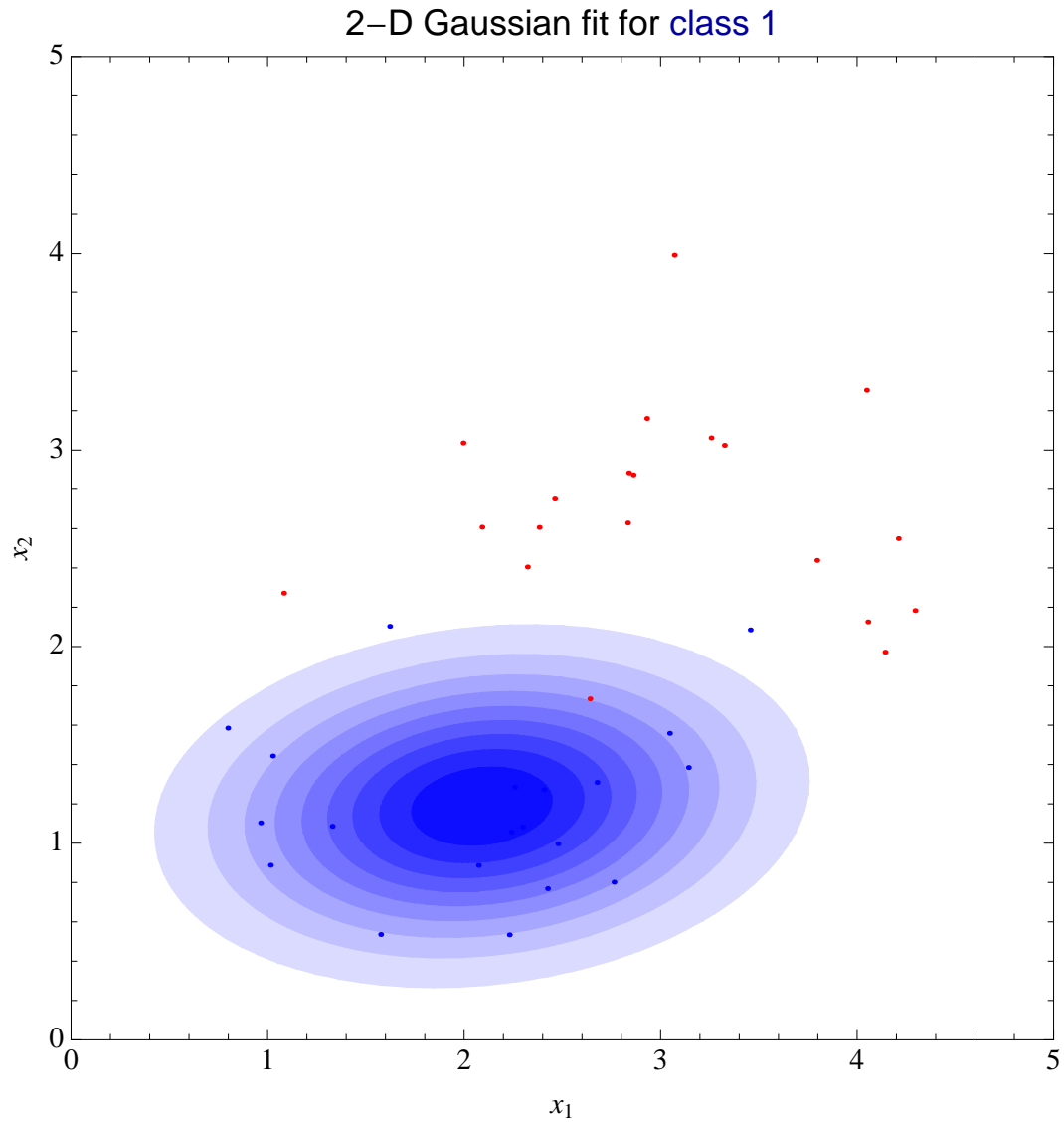
- **Bayes** theorem:

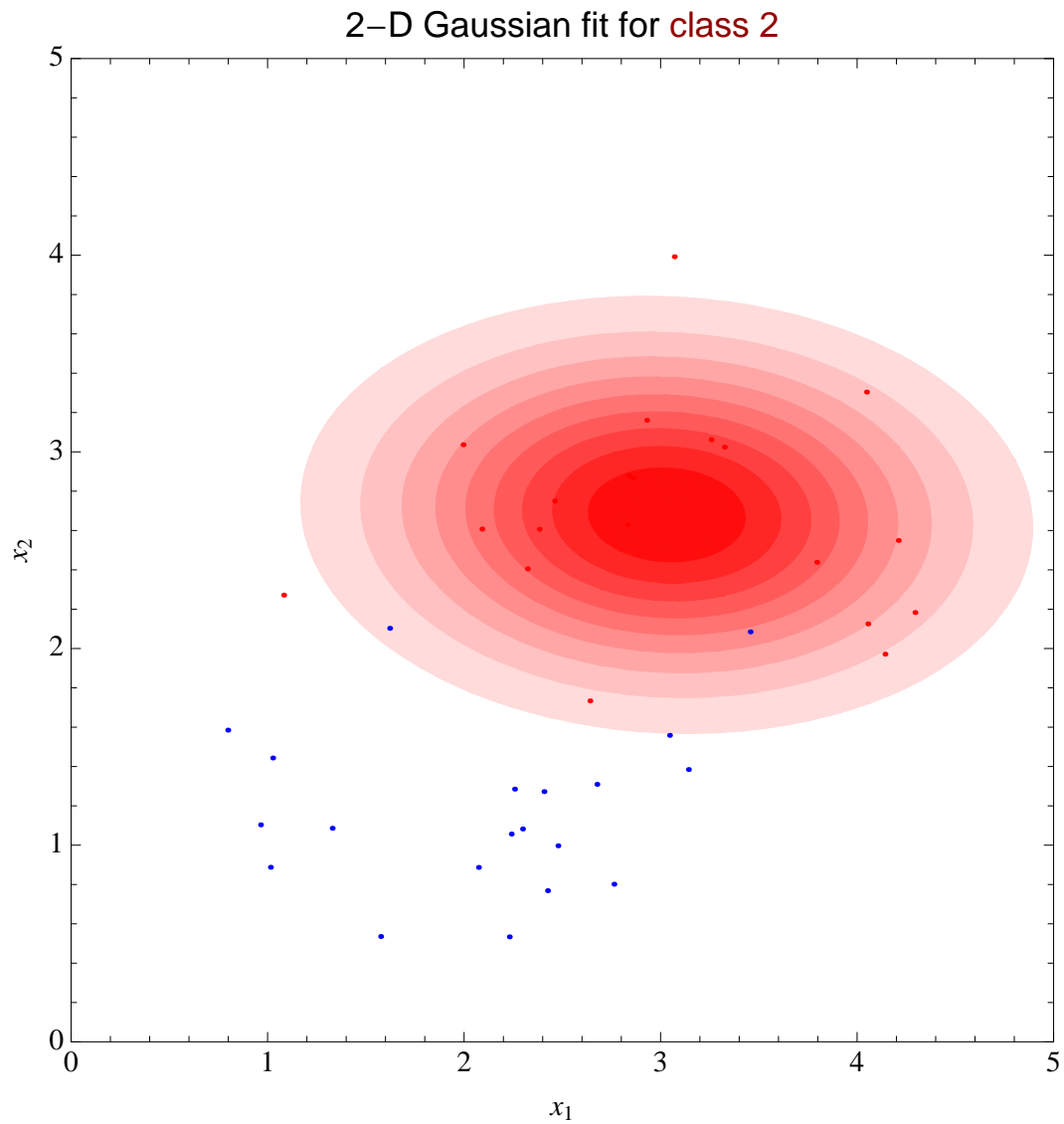
$$P(Y = 1|\mathbf{x}) = \frac{p(\mathbf{x}|Y = 1)P(Y = 1)}{p(\mathbf{x})} \propto p(\mathbf{x}|Y = 1)P(Y = 1)$$

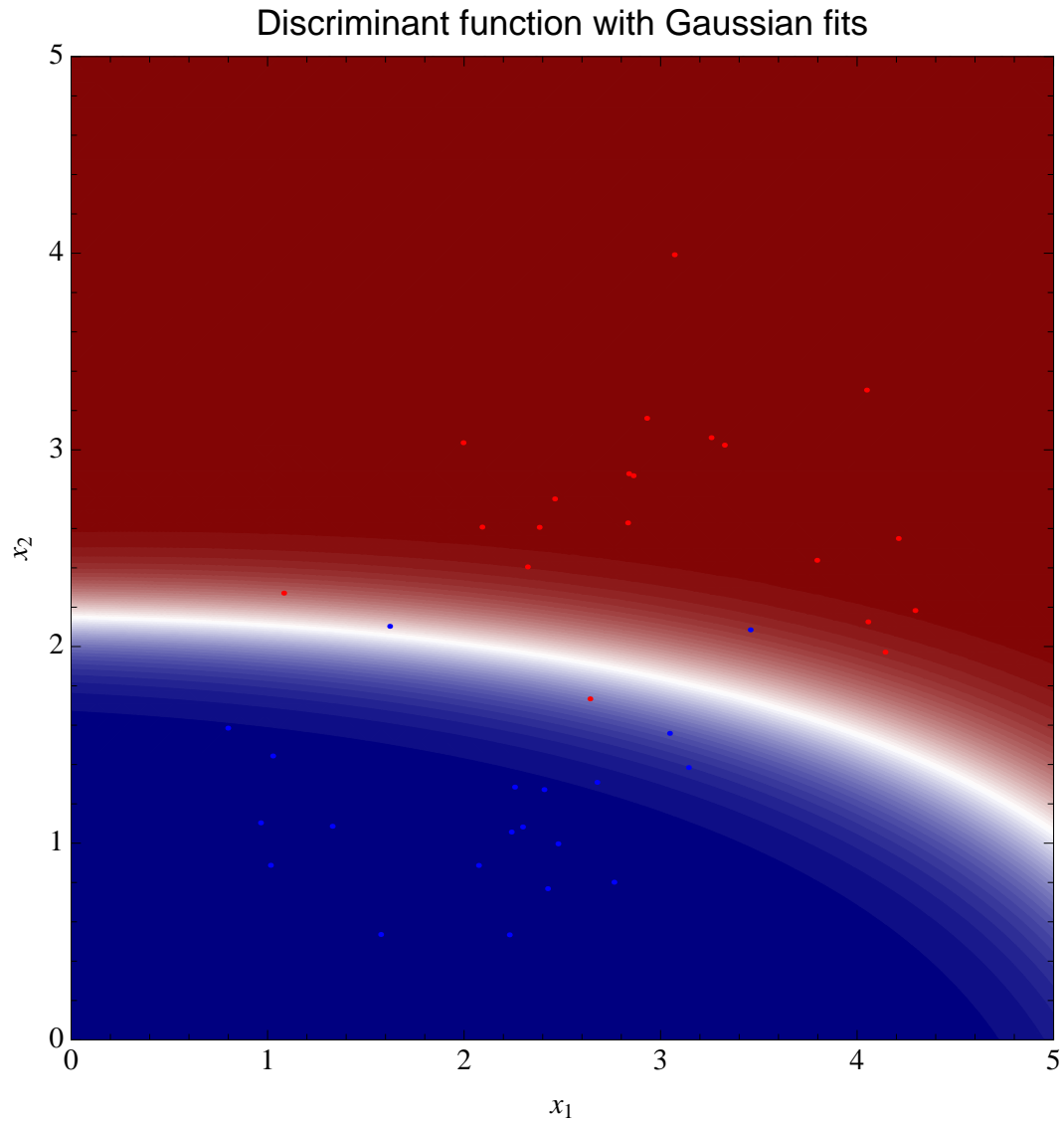
- **Decision:**

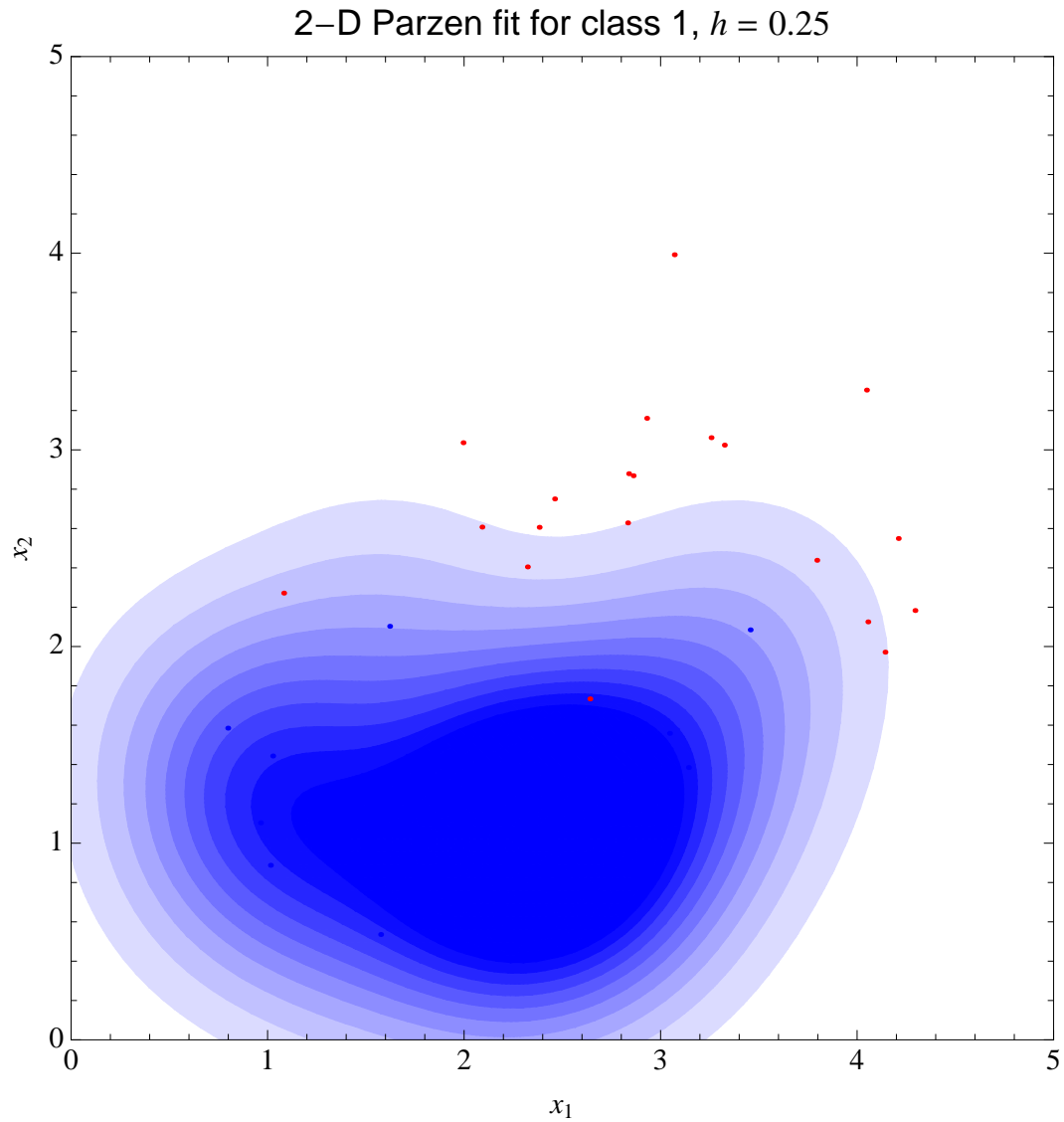
$$g(\mathbf{x}) = \begin{cases} 1 & \text{if } \frac{p(\mathbf{x}|Y=1)P(Y=1)}{p(\mathbf{x}|Y=-1)P(Y=-1)} > 1, \\ -1 & \text{otherwise.} \end{cases}$$

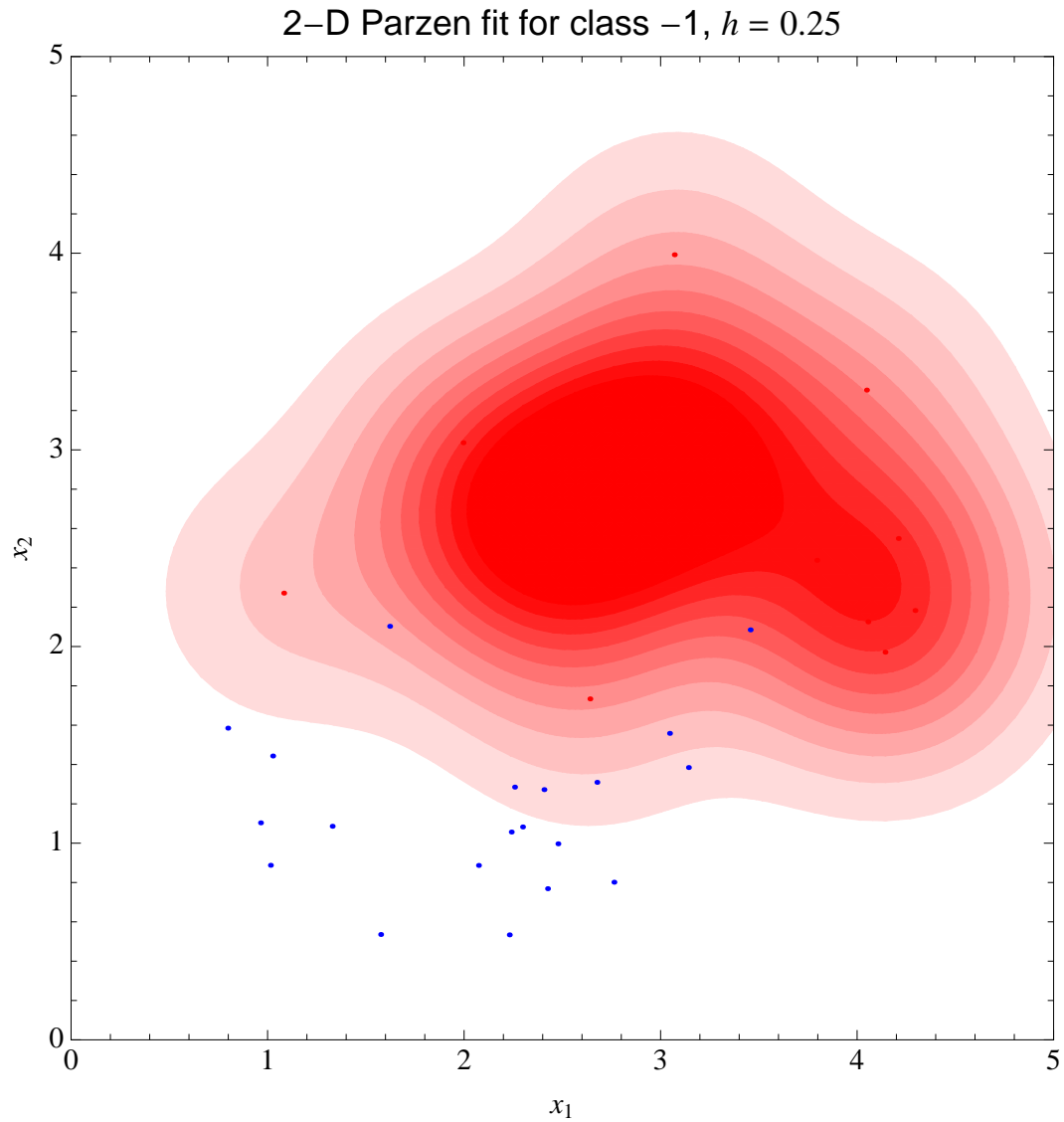


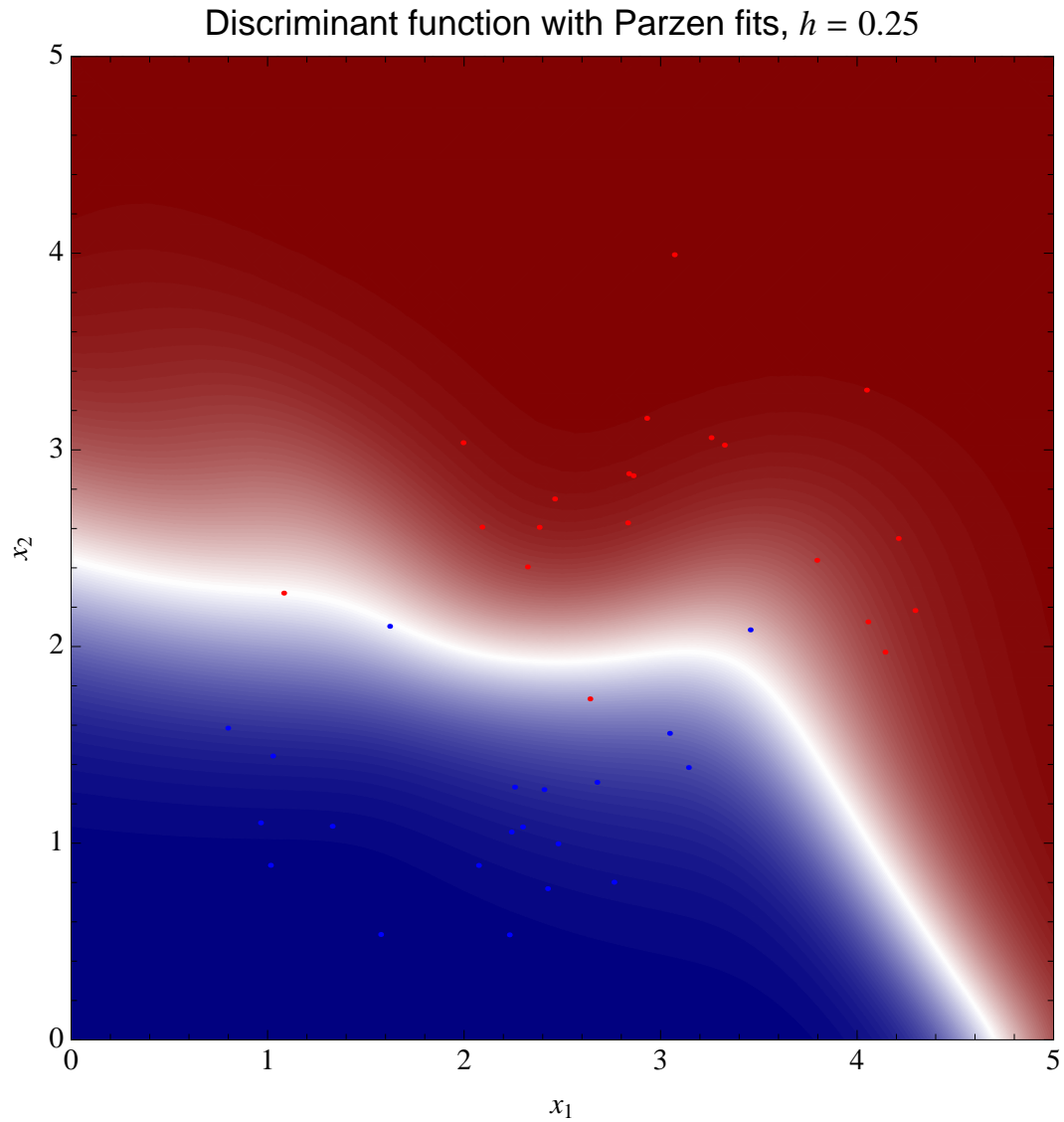


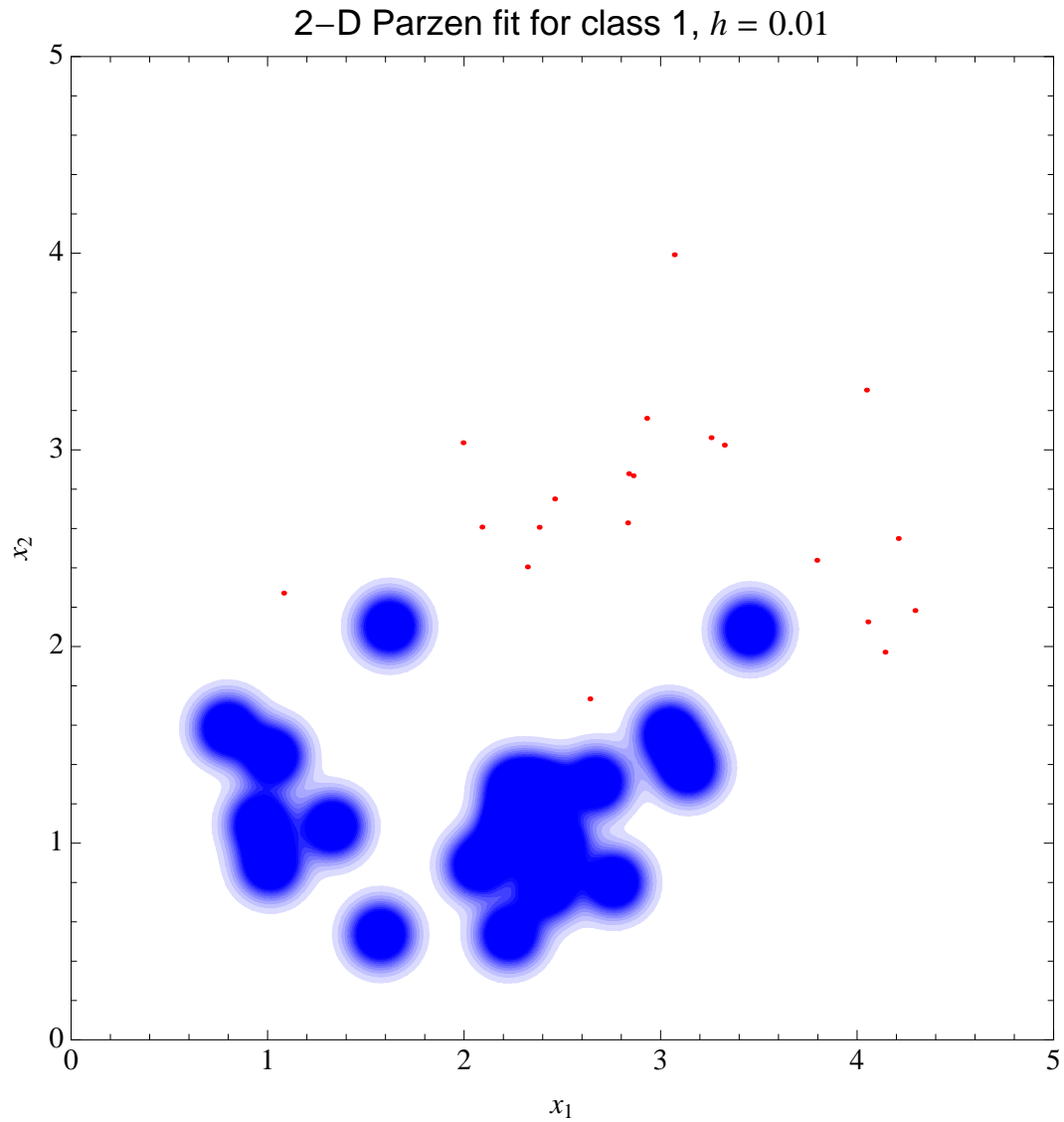


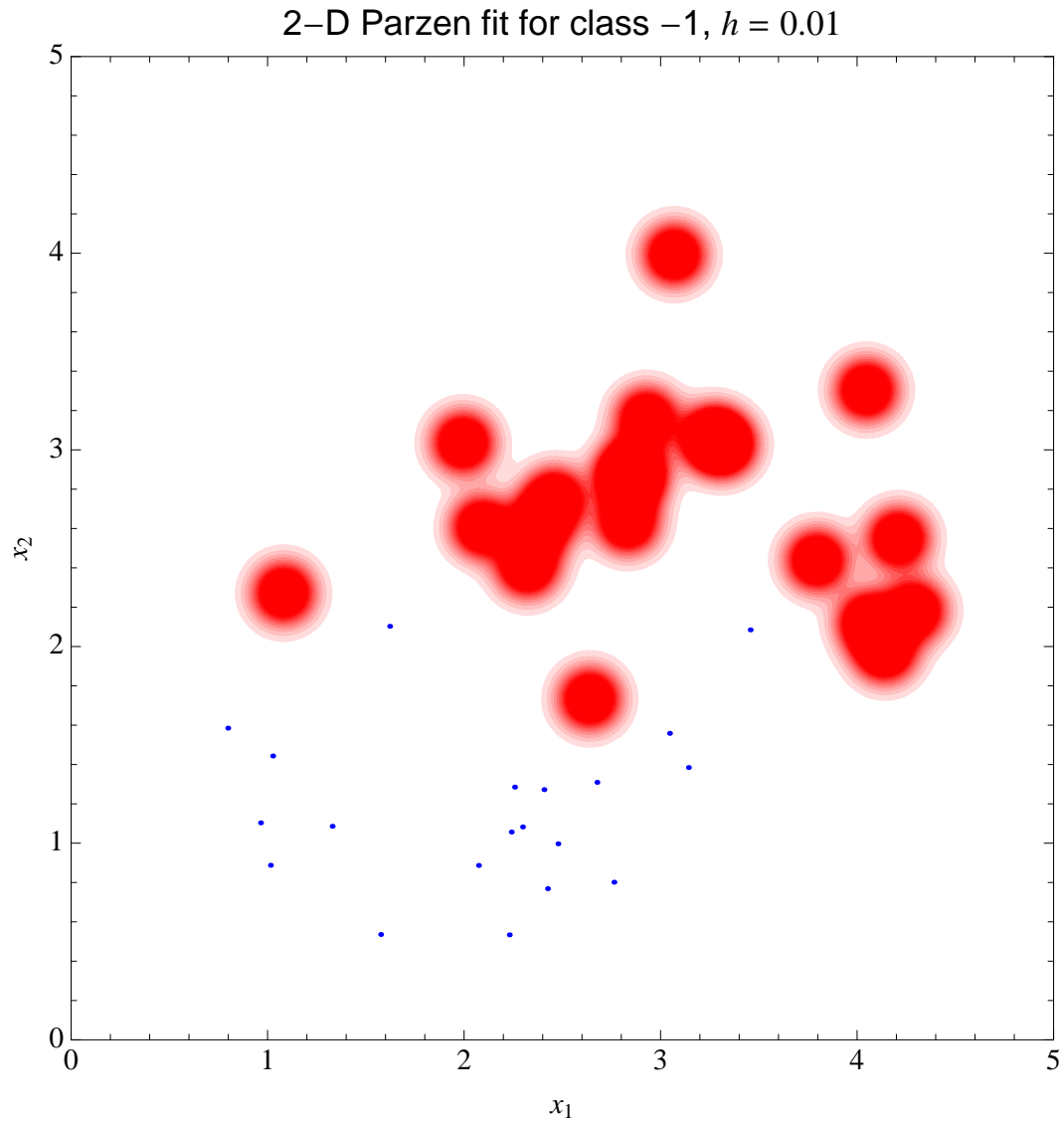


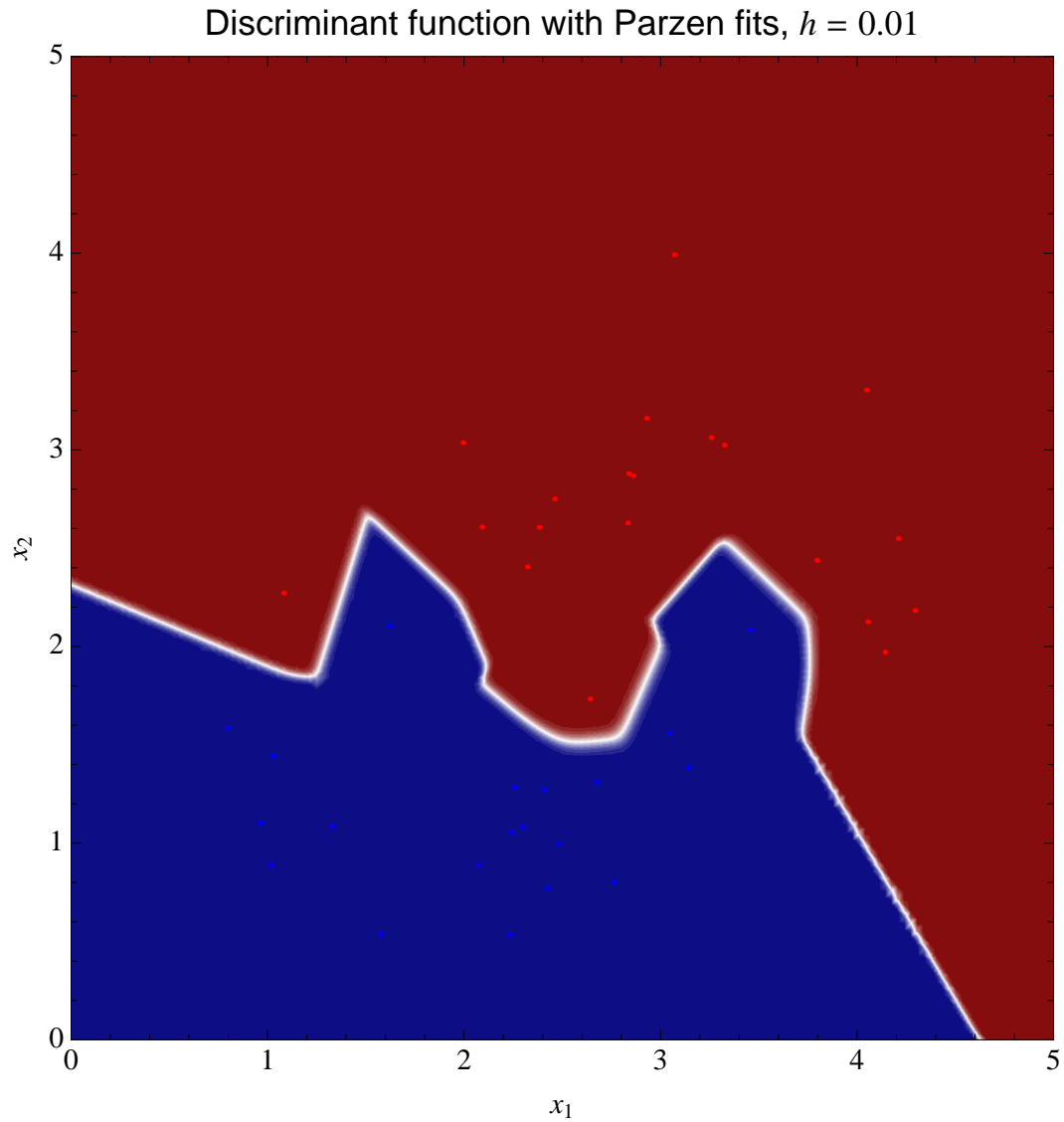


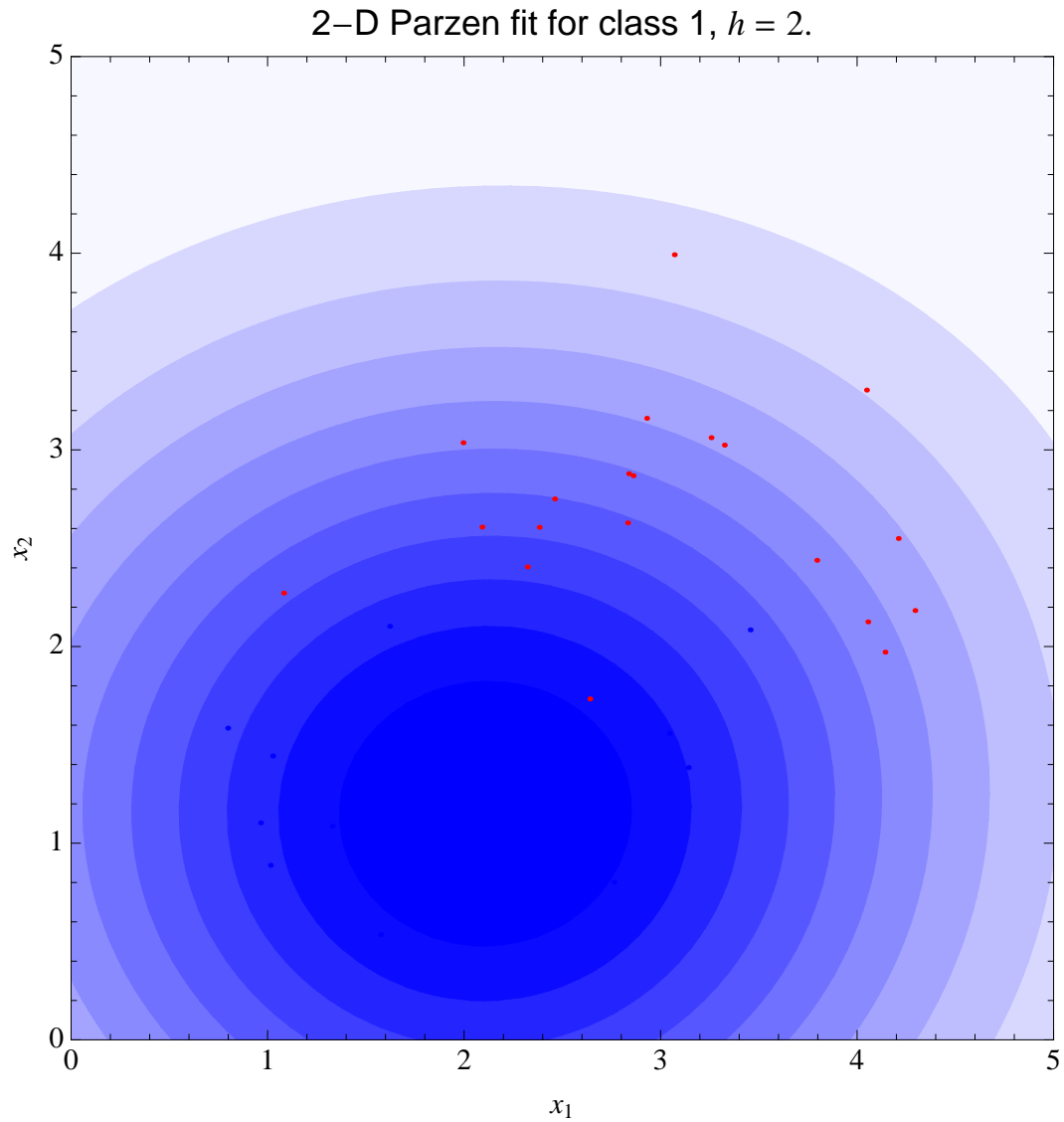


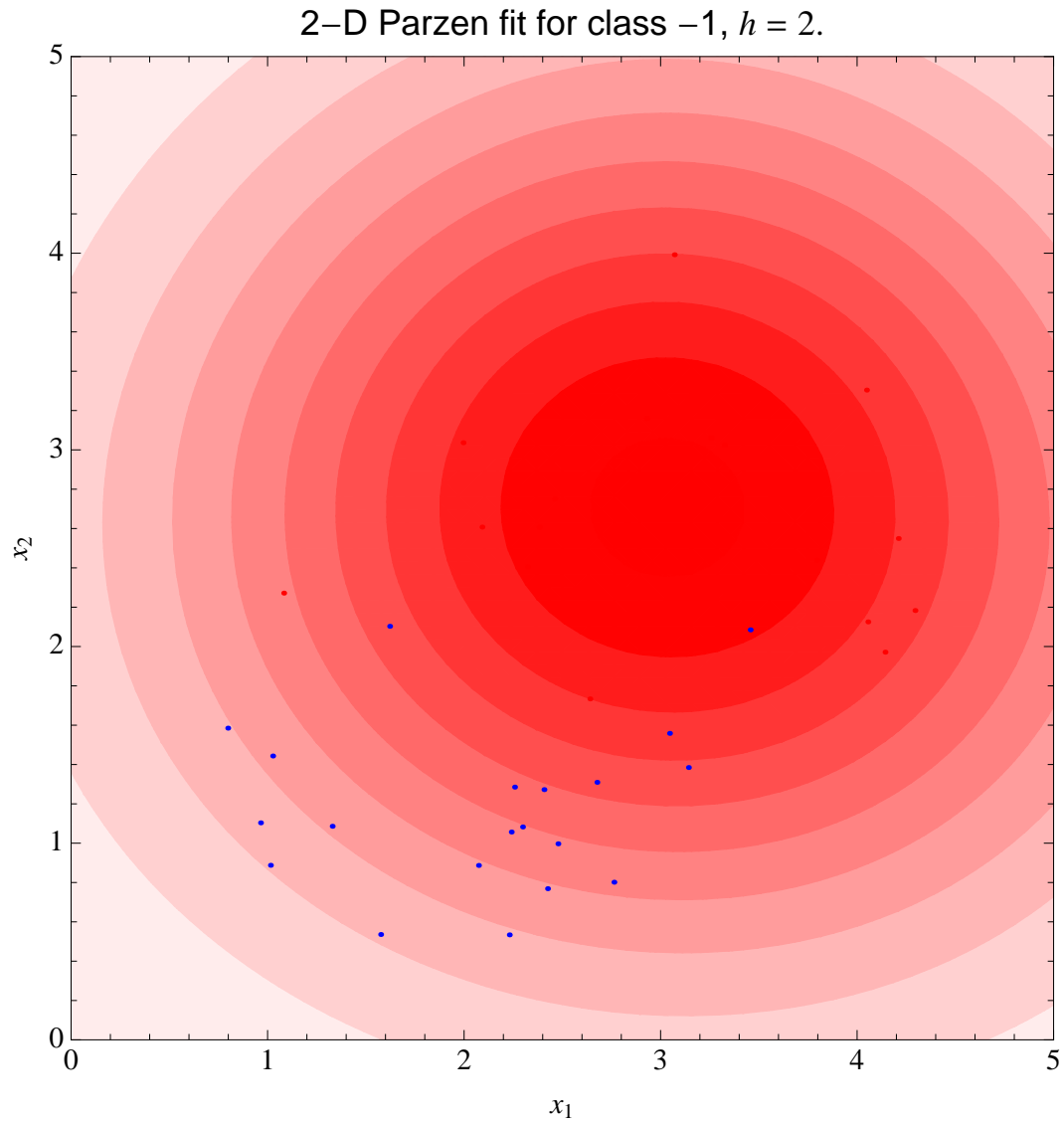


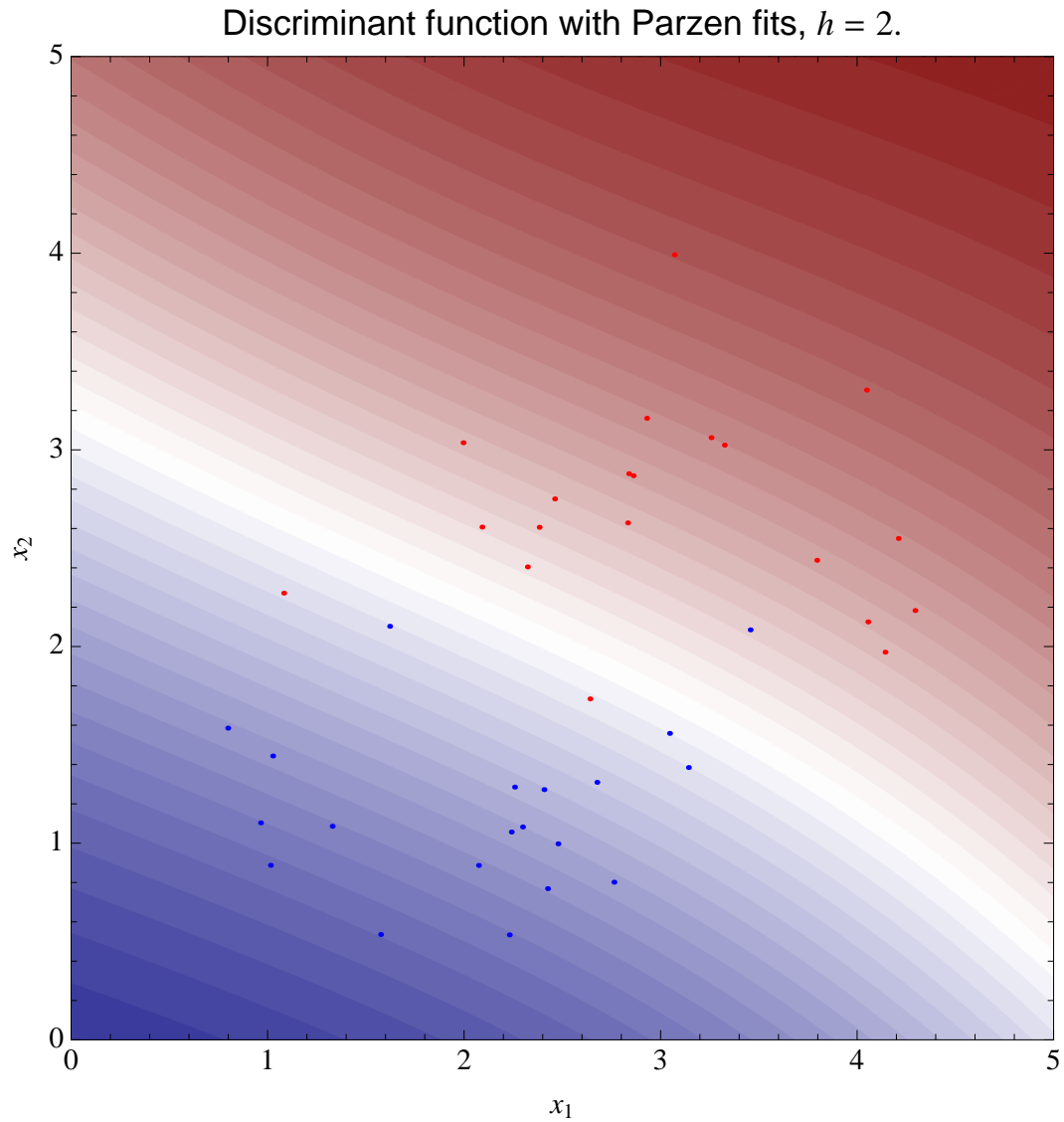




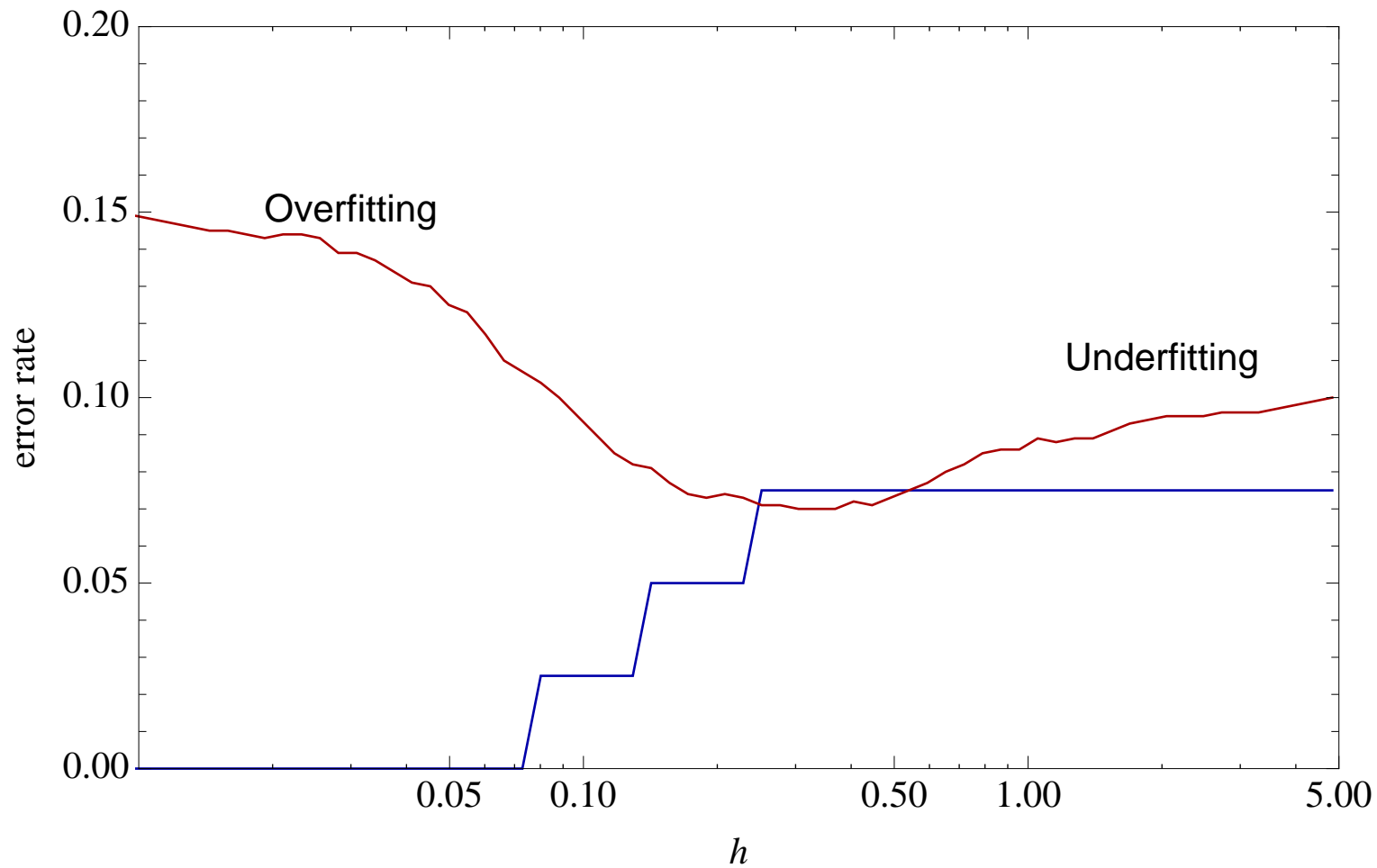








Training and test error rates



- Capacity/complexity control becomes a real issue in **high-dimensional** spaces
 - in a 10000-dimensional space a **linear** function has **10000 parameters!**
 - we **need a lot of training points** to “fill” the space
- Examples
 - images, music, language, text, bioinfo (genetics, proteomics)
- We prefer to **learn the discriminant function directly**, without going through the density estimation step

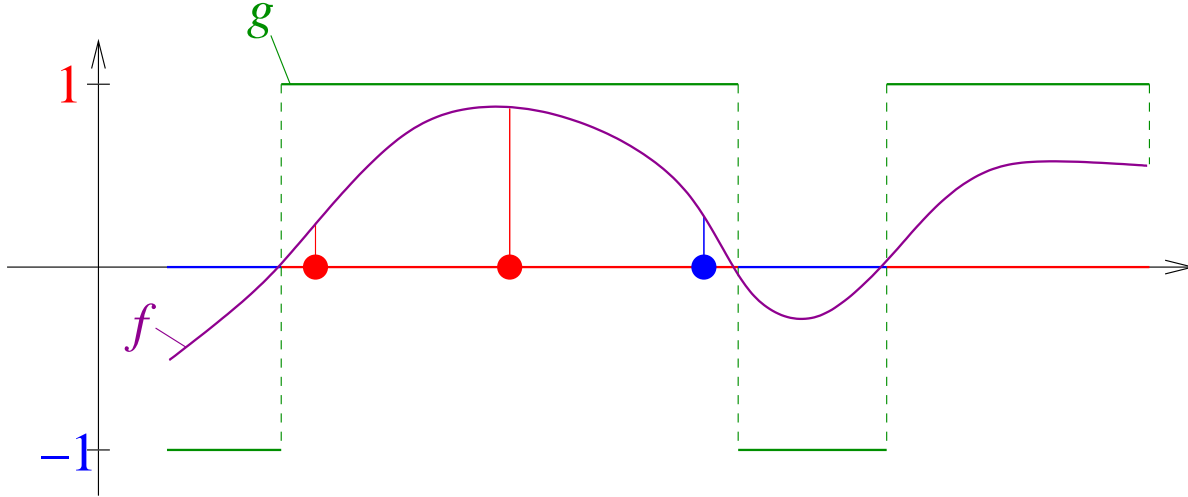
- **Observation** vector: $\mathbf{x} \in \mathbb{R}^d$
- **Class label**: $y \in \{-1, 1\}$ (or $y \in \{1, \dots, K\}$)
- **Classifier**: $g : \mathbb{R}^d \rightarrow \{-1, 1\}$
- **Discriminant** function: $f : \mathbb{R}^d \rightarrow [-1, 1]$

- \longrightarrow classifier

$$g(\mathbf{x}) = \begin{cases} 1, & \text{if } f(\mathbf{x}) \geq 0, \\ -1, & \text{if } f(\mathbf{x}) < 0 \end{cases}$$

- **decision boundary**: $\{\mathbf{x} : f(\mathbf{x}) = 0\}$

The classification model



- Discriminant function: $f : \mathbb{R}^d \rightarrow [-1, 1]$

- \longrightarrow classifier

$$g(\mathbf{x}) = \begin{cases} 1, & \text{if } f(\mathbf{x}) \geq 0, \\ -1, & \text{if } f(\mathbf{x}) < 0 \end{cases}$$

- decision boundary: $\{\mathbf{x} : f(\mathbf{x}) = 0\}$

- Learning from data

- training set : $D_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$

- function class : \mathcal{F}

- learning algorithm : $\text{ALGO} : (\mathbb{R}^d \times \{-1, 1\})^n \rightarrow \mathcal{F}$

$$\text{ALGO}(D_n) \mapsto f$$

- goal: small generalization error $R(g) = P[g(\mathbf{X}) \neq Y] = P[f(\mathbf{X})Y \leq 0]$

- learning principle: minimize the training error

$$\hat{R}(g) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{g(\mathbf{x}_i) \neq y_i\}$$

The classification model

- goal: small **generalization error**

$$R(g) = P[g(\mathbf{X}) \neq Y]$$

- Learning principle: minimize the **training error**

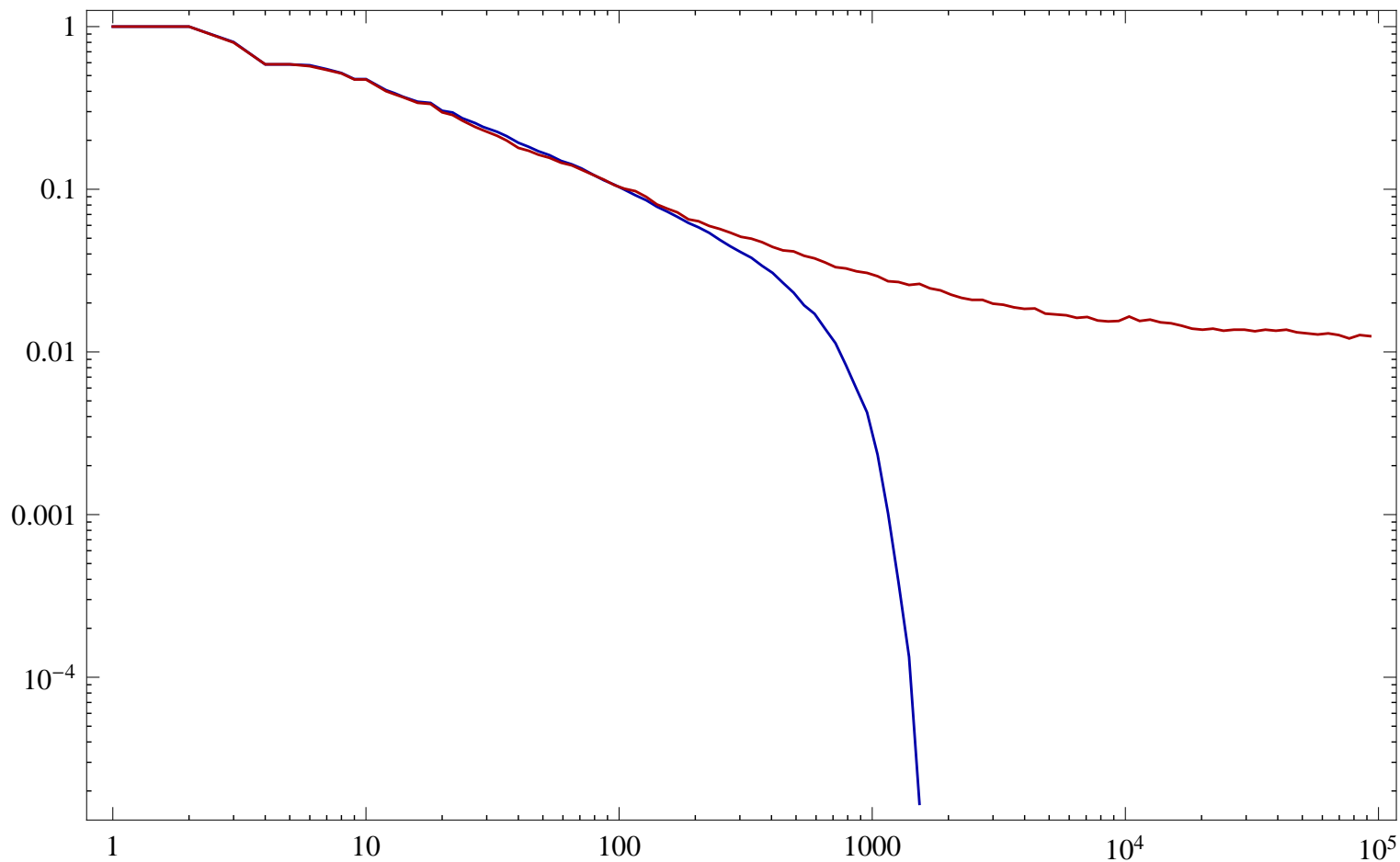
$$\hat{R}(g) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{g(\mathbf{x}_i) \neq y_i\}$$

- Generalization error **bounds**

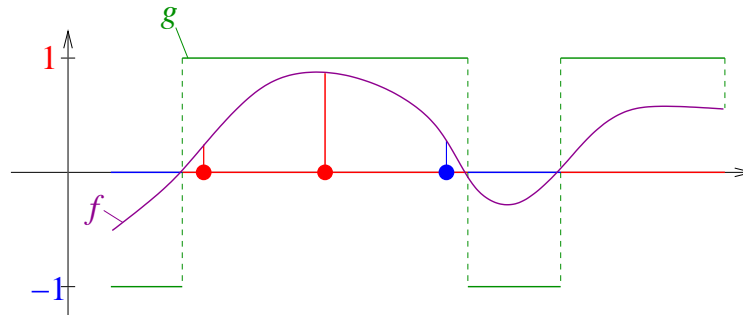
$$\text{“ } R(g) \leq \hat{R}(g) + O\left(\frac{\dim_{\text{VC}}(\mathcal{F}) \log n}{n}\right) \text{”}$$

The classification model

AdaBoost on MNIST: **training** and **test** error rates



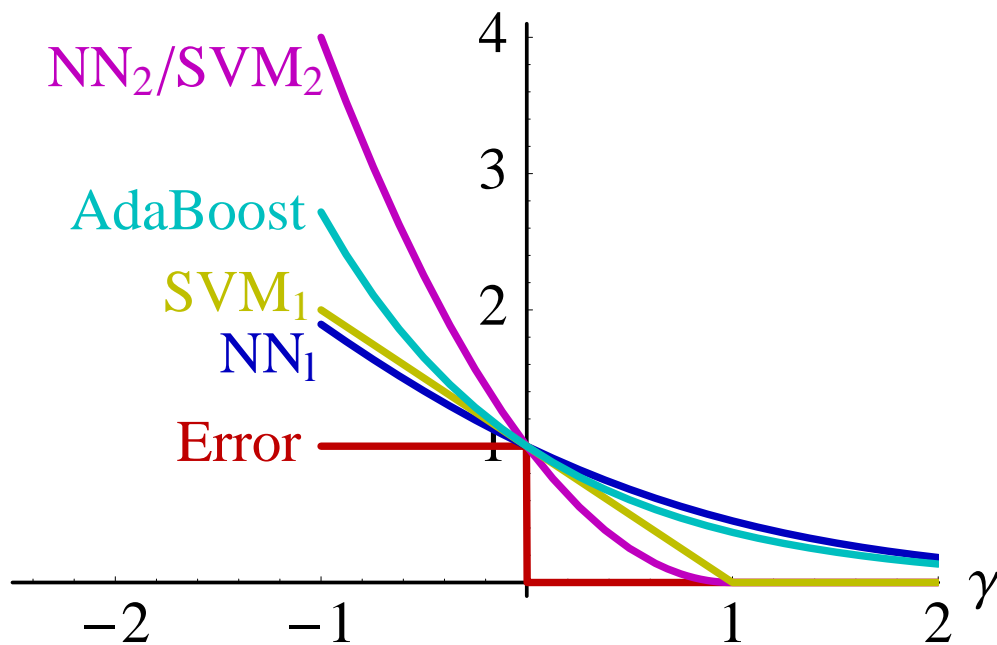
The classification model



- **Margin**: $\gamma = y \cdot f(\mathbf{x})$
 - **classification error** \equiv negative margin
 - the **magnitude** of a positive margin quantifies the **confidence**
 - learning principle: minimize a **smooth loss** function over the **margin**

$$\widehat{R}_\gamma(f) = \frac{1}{n} \sum_{i=1}^n L(f(\mathbf{x}_i)y_i)$$

- Margin loss functions

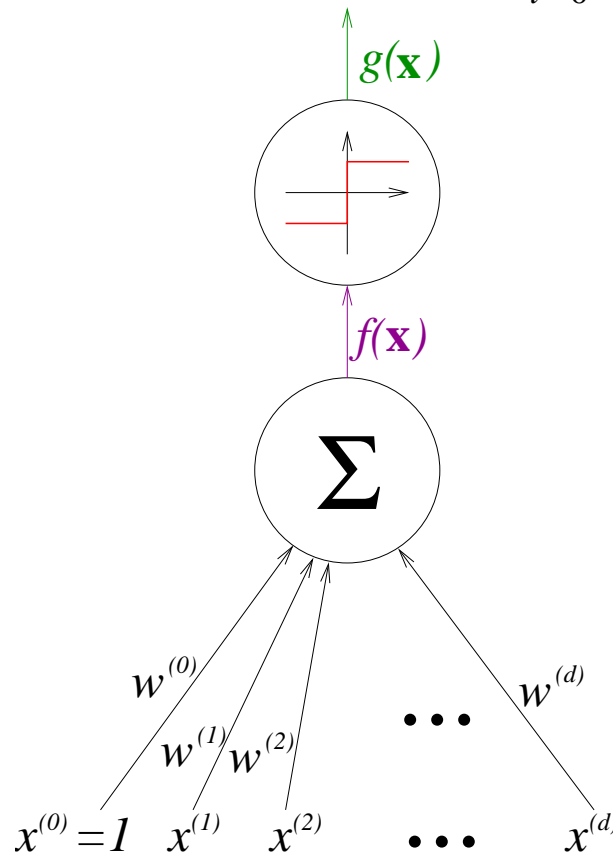


- Algorithms

- 1958: **Perceptron** [Rosenblatt, '58] – [Minsky–Papert '69]
- 1986: **Multilayer perceptrons (neural networks)** and the **back-propagation** algorithm [Rumelhart–Hinton–Williams, '86]
- 1995: **Support vector machines** [Boser–Guyon–Vapnik, '92], [Cortes–Vapnik, '95]
- 1997: **boosting, AdaBoost** [Freund, '95], [Freund–Schapire, '97]

The perceptron

- **Linear** discriminant functions: $f(\mathbf{x}) = \sum_{i=0}^d w^{(i)} \cdot x^{(i)} = \langle \mathbf{w}, \mathbf{x} \rangle$



- **Linear** discriminant functions: $f(\mathbf{x}) = \sum_{i=0}^d w^{(i)} \cdot x^{(i)} = \langle \mathbf{w}, \mathbf{x} \rangle$
- Algorithm
 - simple iterative **error correction**
 - **convergence** if the data is **linearly separable**
 - **oscillation** for **linearly non-separable** data

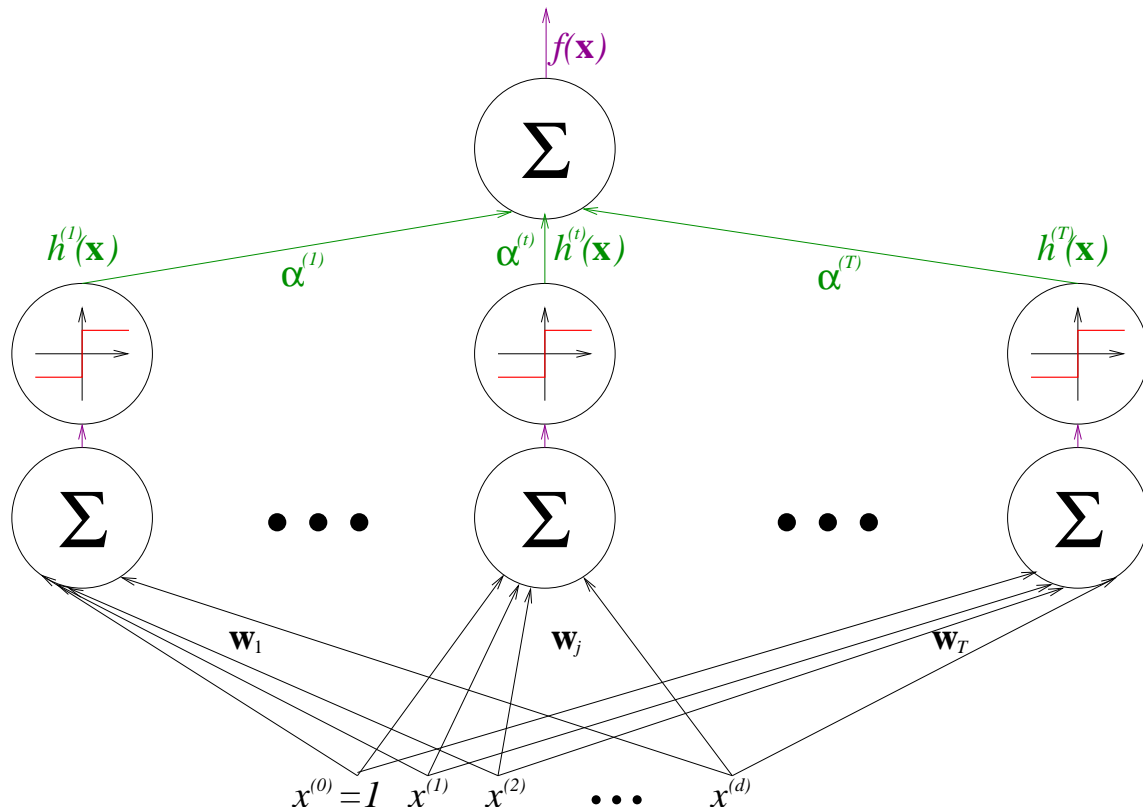
- Model:

$$f(\mathbf{x}) = \sum_{j=1}^N \alpha^{(j)} h^{(j)}(\mathbf{x})$$

- $h^{(j)} : \mathbb{R}^d \rightarrow [-1, 1]$
 - simple classifiers/discriminant functions, features, experts
- $\alpha^{(j)} \in \mathbb{R}^+$
 - weight of the expert $h^{(j)}$ in the final vote

Multilayer perceptron (neural net)

- Model: $f(\mathbf{x}) = \sum_{j=1}^N \alpha^{(j)} \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle)$



- Model: $f(\mathbf{x}) = \sum_{j=1}^N \alpha^{(j)} \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle)$
- Algorithm:
 - gradient descent optimization
 - differentiable error functions \rightarrow margin loss
 - differentiable activation function σ : the sigmoid
 - local minima, “engineering”, parameters to tune
 - fast, works well if well-tuned
 - versatile: multi-class classification, regression, density estimation

- Model:

$$f(\mathbf{x}) = \sum_{j=1}^N \alpha^{(j)} h^{(j)}(\mathbf{x})$$

- no restriction on the form of $h^{(j)}(\mathbf{x})$
- often “decision stumps” :

$$h_{\ell, \theta}(\mathbf{x}) = \begin{cases} +1 & \text{if } x^{(\ell)} \geq \theta, \\ -1 & \text{otherwise} \end{cases}$$

where $\mathbf{x} = (x^{(1)}, \dots, x^{(d)})$

- Algorithm
 - **extremely simple** learning, limited parameter tuning
 - **fast**
 - the choice of the **pool of experts** captures the a-priori knowledge
 - **no restriction** on the form of the simple classifiers
 - **multi-class** classification is natural, **regression** is not

- Model:

$$f(\mathbf{x}) = \sum_{j \in I_{sv}} \alpha^{(j)} y_j K(\mathbf{x}_j, \mathbf{x})$$

- $I_{sv} \subset \{1, \dots, n\}$ is the set of support vectors
- $K(\cdot, \cdot)$ is a similarity function (kernel)

- Kernel:

- $K(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle \longrightarrow f(\mathbf{x})$ is linear
- $K(\mathbf{x}, \mathbf{x}') = (1 + \langle \mathbf{x}, \mathbf{x}' \rangle)^d \longrightarrow f(\mathbf{x})$ is a polynomial of degree d
- $K(\mathbf{x}, \mathbf{x}') = \exp(-1/h \|\mathbf{x} - \mathbf{x}'\|^2) \longrightarrow f(\mathbf{x})$ is a Gaussian mixture (\rightarrow Parzen)

- Model:

$$f(\mathbf{x}) = \sum_{j \in I_{sv}} \alpha^{(j)} y_j K(\mathbf{x}_j, \mathbf{x})$$

- Algorithm:

- goal: **classification boundary** equidistant from classes
- “sophisticated nearest neighbor”
- slow and complex **quadratic programming** optimization
- **turn-key** algorithm, very **limited** parameter **tuning**
- no multi-class, no regression

Machines à vecteurs de support

- Idée 1: séparation linéaire avec une **marge maximale**

- marge **fonctionnelle**:

$$\gamma_i = f(\mathbf{x}_i)y_i = (\mathbf{w}^t \mathbf{x}_i + w_0)y_i$$

- marge **géométrique**:

$$\gamma_i^{(g)} = \frac{1}{\|\mathbf{w}\|} (\mathbf{w}^t \mathbf{x}_i + w_0)y_i = \frac{1}{\|\mathbf{w}\|} \gamma_i$$

- Optimisations équivalentes

- maximiser la marge **géométrique**

- maximiser la marge **fonctionnelle** sous la **contrainte** $\|\mathbf{w}\| = 1$

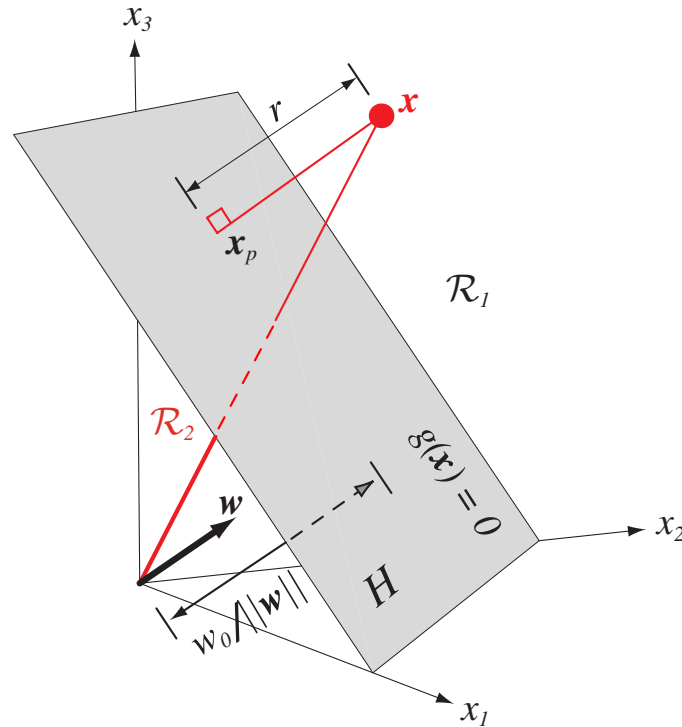
- minimiser $\|\mathbf{w}\|$ sous la **contrainte** $\gamma_i \geq 1$

- Géométrie – deux classes

- r = distance algébrique de \mathbf{x} et H :

$$\begin{aligned}\mathbf{x} &= \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \\ g(\mathbf{x}) &= \mathbf{w}^t \mathbf{x} + w_0 = r \|\mathbf{w}\| \\ r &= \frac{g(\mathbf{x})}{\|\mathbf{w}\|}\end{aligned}$$

- Géométrie – deux classes



- Problème d'optimisation:
 - soit $D_n = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n))$ linéairement séparable
 - minimiser $\|\mathbf{w}\|^2 = \mathbf{w}^t \mathbf{w}$
 - sous les contraintes $\gamma_i = (\mathbf{w}^t \mathbf{x}_i + w_0) y_i \geq 1, \quad i = 1, \dots, n$
- Résultat:
 - l'hyperplan $(\mathbf{w}^t \mathbf{x} + w_0 = 0)$ avec une marge géométrique $\frac{1}{\|\mathbf{w}\|}$ maximale

Machines à vecteurs de support

- Problème **primal** – théorème de **Kuhn-Tucker**
 - minimiser par rapport à \mathbf{w} et w_0 et maximiser par rapport à $\boldsymbol{\alpha}$:

$$\begin{aligned} L(\mathbf{w}, w_0, \boldsymbol{\alpha}) &= \frac{1}{2} \mathbf{w}^t \mathbf{w} - \sum_{i=1}^n \alpha_i (\gamma_i - 1) \\ &= \frac{1}{2} \mathbf{w}^t \mathbf{w} - \sum_{i=1}^n \alpha_i ((\mathbf{w}^t \mathbf{x}_i + w_0) y_i - 1) \end{aligned}$$

- sous les contraintes $\alpha_i \geq 0$, $i = 1, \dots, n$

Machines à vecteurs de support

- Optimisation:

- les gradients:

$$\frac{\partial L(\mathbf{w}, w_0, \boldsymbol{\alpha})}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^n \alpha_i \mathbf{x}_i y_i = \mathbf{0}$$

$$\frac{\partial L(\mathbf{w}, w_0, \boldsymbol{\alpha})}{\partial w_0} = \sum_{i=1}^n \alpha_i y_i = 0$$

- resubstitution: maximiser par rapport à $\boldsymbol{\alpha}$ (problème dual):

$$\begin{aligned} W(\boldsymbol{\alpha}) &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^t \mathbf{x}_j - \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^t \mathbf{x}_j + \sum_{i=1}^n \alpha_i \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^t \mathbf{x}_j \end{aligned}$$

- sous les contraintes $\alpha_i \geq 0$, $i = 1, \dots, n$ et $\sum_{i=1}^n \alpha_i y_i = 0$

- La solution

- $\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i$

- $w_0^* = -\frac{1}{2} \left(\max_{y_i=-1} \mathbf{w}^* \mathbf{x}_i + \min_{y_i=1} \mathbf{w}^* \mathbf{x}_i \right)$

Machines à vecteurs de support

- La **structure** de la solution
 - \mathbf{w}^* est une **combinaison linéaire des points** d'entraînement
 - théorème de **Kuhn-Tucker**:

$$\alpha_i^* ((\mathbf{w}^{*t} \mathbf{x}_i + w_0^*) y_i - 1) = 0, \quad i = 1, \dots, n$$

- si $\gamma_i > 1$ alors $\alpha_i^* = 0$
- si $\alpha_i^* > 0$ alors $\gamma_i = 1$: **vecteurs de support**
- $\mathbf{w}^* = \sum_{i \in sv} \alpha_i^* y_i \mathbf{x}_i$
- $f^*(\mathbf{x}) = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i^t \mathbf{x} + w_0^* = \sum_{i \in sv} \alpha_i^* y_i \mathbf{x}_i^t \mathbf{x} + w_0^*$

- La structure de la solution

- pour $j \in sv$

$$\gamma_j = y_j f^*(\mathbf{x}_j) = y_j \left(\sum_{i \in sv} \alpha_i^* y_i \mathbf{x}_i^t \mathbf{x}_j + w_0^* \right) = 1$$

- alors

$$\begin{aligned} \mathbf{w}^{*t} \mathbf{w}^* &= \sum_{i=1}^n \sum_{j=1}^n \alpha_i^* \alpha_j^* y_i y_j \mathbf{x}_i^t \mathbf{x}_j = \sum_{j \in sv} \alpha_j^* y_j \sum_{i \in sv} \alpha_i^* y_i \mathbf{x}_i^t \mathbf{x}_j = \sum_{j \in sv} \alpha_j^* (1 - y_j w_0^*) \\ &= \sum_{j \in sv} \alpha_j^* \end{aligned}$$

- la marge maximale:

$$\gamma^* = \frac{1}{\|\mathbf{w}\|} = \left(\sum_{i \in sv} \alpha_i^* \right)^{-1/2}$$

Machines à vecteurs de support

- Idée 2: le **noyau**

- l'optimisation: $W(\boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^t \mathbf{x}_j$

- la fonction optimale: $f^*(\mathbf{x}) = \sum_{i \in sv} \alpha_i^* y_i \mathbf{x}_i^t \mathbf{x} + w_0^*$

- Remplacer $\mathbf{x}^t \mathbf{x}'$ par une **fonction de noyaux** $K(\mathbf{x}, \mathbf{x}')$

- l'optimisation: $W(\boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$

- la fonction optimale: $f^*(\mathbf{x}) = \sum_{i \in sv} \alpha_i^* y_i K(\mathbf{x}_i, \mathbf{x}) + w_0^*$

- matrice de **Gram**: $\mathbf{G}_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j)$

- équivalent à une **transformation non-linéaire** dans une espace de **traits** de dimension très élevée

- Exemples de noyaux

- linéaire: $K^\ell(\mathbf{x}, \mathbf{x}') = \mathbf{x}^t \mathbf{x}'$

- polynômial: $K_d^p(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^t \mathbf{x}' + 1)^d$

- base radiale: $K^r(\mathbf{x}, \mathbf{x}') = K(\|\mathbf{x} - \mathbf{x}'\|)$

- gaussien: $K_\gamma^g(\mathbf{x}, \mathbf{x}') = e^{-\gamma \|\mathbf{x} - \mathbf{x}'\|^2}$

- réseaux de neurones: $K^\sigma(\mathbf{x}, \mathbf{x}') = \sigma(s\mathbf{x}^t \mathbf{x}' + c)$

- sigmoïde: $K_{s,c}^t(\mathbf{x}, \mathbf{x}') = \tanh(s\mathbf{x}^t \mathbf{x}' + c)$

Machines à vecteurs de support

- Discrimination linéaire **généralisée**

- linéaire:

$$f^*(\mathbf{x}) = \mathbf{w}^{*t} \mathbf{x} + w_0^* = \sum_{j=1}^d w_j^* x^{(j)} + w_0^* = \sum_{i \in sv} \alpha_i^* y_i \mathbf{x}_i^t \mathbf{x} + w_0^*$$

- généralisée:

$$f^*(\mathbf{x}) = \sum_{j=1}^D w_j^* \phi^{(j)}(\mathbf{x}) + w_0^* = \sum_{i \in sv} \alpha_i^* y_i \sum_{j=1}^D \phi^{(j)}(\mathbf{x}_i) \phi^{(j)}(\mathbf{x}) + w_0^*$$

- produit scalaire dans l'**espace de traits**:

$$\phi^t(\mathbf{x}) \phi(\mathbf{x}') = \sum_{j=1}^D \phi^{(j)}(\mathbf{x}) \phi^{(j)}(\mathbf{x}')$$

Machines à vecteurs de support

- Discrimination linéaire généralisée

- exemple:

$$\begin{aligned}\phi^{(j)}(\mathbf{x}) &= \sqrt{2}x^{(j)}, \quad j = 1, \dots, d, \\ \phi^{(id+j)}(\mathbf{x}) &= x^{(i)}x^{(j)}, \quad i, j = 1, \dots, d \\ \phi^{(d^2)}(\mathbf{x}) &= 1\end{aligned}$$

- produit scalaire dans l'espace des traits:

$$\begin{aligned}\phi^t(\mathbf{x})\phi(\mathbf{x}') &= \sum_{j=1}^{d^2} \phi^{(j)}(\mathbf{x})\phi^{(j)}(\mathbf{x}') \\ &= \sum_{j=1}^d \sqrt{2}x^{(j)}\sqrt{2}x'^{(j)} + \sum_{i=1}^d \sum_{j=1}^d x^{(i)}x^{(j)}x'^{(i)}x'^{(j)} + 1 \\ &= 2\mathbf{x}^t\mathbf{x}' + (\mathbf{x}^t\mathbf{x}')^2 + 1 \\ &= (\mathbf{x}^t\mathbf{x}' + 1)^2 = K_2^p(\mathbf{x}, \mathbf{x}')\end{aligned}$$

- Caractérisation des noyaux

- $K(\mathbf{x}, \mathbf{x}') = \boldsymbol{\phi}^t(\mathbf{x})\boldsymbol{\phi}(\mathbf{x}') = \sum_{j=1}^D \phi^{(j)}(\mathbf{x})\phi^{(j)}(\mathbf{x}')$

- commutativité: $K(\mathbf{x}, \mathbf{x}') = K(\mathbf{x}', \mathbf{x})$

- inégalité de Cauchy-Schwartz:

$$\begin{aligned} K(\mathbf{x}, \mathbf{x}')^2 &= (\boldsymbol{\phi}^t(\mathbf{x})\boldsymbol{\phi}(\mathbf{x}'))^2 \\ &\leq \|\boldsymbol{\phi}(\mathbf{x})\|^2 \|\boldsymbol{\phi}(\mathbf{x}')\|^2 \\ &= \boldsymbol{\phi}^t(\mathbf{x})\boldsymbol{\phi}(\mathbf{x}) \cdot \boldsymbol{\phi}^t(\mathbf{x}')\boldsymbol{\phi}(\mathbf{x}') \\ &= K(\mathbf{x}, \mathbf{x})K(\mathbf{x}', \mathbf{x}') \end{aligned}$$

- Théorème de Mercer:

- $\phi^t(\mathbf{x})\phi(\mathbf{x}') = K(\mathbf{x}, \mathbf{x}')$

- Conditions

- diagonaliser la matrice de Gram: $\mathbf{G}_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{V}\mathbf{\Lambda}\mathbf{V}$

- valeurs propres: λ_t , vecteurs propres: \mathbf{v}_t

- condition suffisante: \mathbf{G} est positif semi-défini ($\forall t : \lambda_t > 0$) pour tout $\mathcal{X}_n = (\mathbf{x}_1, \dots, \mathbf{x}_n)$

- D peut être ∞ !!!

Machines à vecteurs de support

- Idée 3: les **variables d'écart** (slack variables)
 - cas **non-séparable**
- Permettre des **erreurs**:
 - minimiser $\|\mathbf{w}\|^2 = \mathbf{w}^t \mathbf{w}$
 - sous les **contraintes** $\gamma_i = (\mathbf{w}^t \mathbf{x}_i + w_0)y_i \geq 1 - \xi_i, \quad i = 1, \dots, n$
 - où $\xi_i \geq 0, \quad i = 1, \dots, n$
 - minimiser l'erreur \equiv minimiser le **nombre de points avec $\xi_i > 0$** :
NP-difficile

Machines à vecteurs de support

- **Problème soluble 1**: minimiser $\|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i^2$
 - on peut supprimer la contrainte de positivité des ξ_i
 - C est un **hyper-paramètre** réglé par la validation croisée (par exemple)

Machines à vecteurs de support

- Problème primal

- minimiser par rapport à \mathbf{w} , ξ et w_0 et maximiser par rapport à α :

$$L(\mathbf{w}, w_0, \xi, \alpha) = \frac{1}{2} \mathbf{w}^t \mathbf{w} + \frac{C}{2} \sum_{i=1}^n \xi_i^2 - \sum_{i=1}^n \alpha_i [(\mathbf{w}^t \mathbf{x}_i + w_0) y_i - 1 + \xi_i]$$

- sous les contraintes $\alpha_i \geq 0$, $i = 1, \dots, n$
- les gradients:

$$\frac{\partial L(\mathbf{w}, w_0, \alpha)}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^n \alpha_i \mathbf{x}_i y_i = \mathbf{0}$$

$$\frac{\partial L(\mathbf{w}, w_0, \alpha)}{\partial \xi} = C \xi - \alpha = \mathbf{0}$$

$$\frac{\partial L(\mathbf{w}, w_0, \alpha)}{\partial w_0} = \sum_{i=1}^n \alpha_i y_i = 0$$

Machines à vecteurs de support

- Problème dual

- resubstitution: maximiser par rapport à α :

$$\begin{aligned}
 W(\boldsymbol{\alpha}) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) + \frac{1}{2C} \boldsymbol{\alpha}^t \boldsymbol{\alpha} - \frac{1}{C} \boldsymbol{\alpha}^t \boldsymbol{\alpha} \\
 &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) - \frac{1}{2C} \boldsymbol{\alpha}^t \boldsymbol{\alpha} \\
 &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \left(K(\mathbf{x}_i, \mathbf{x}_j) + \frac{1}{C} \delta_{i,j} \right)
 \end{aligned}$$

- sous les contraintes $\alpha_i \geq 0$, $i = 1, \dots, n$ et $\sum_{i=1}^n \alpha_i y_i = 0$

- La solution

- $\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i = \sum_{i \in sv} \alpha_i^* y_i \mathbf{x}_i$

- $f^*(\mathbf{x}) = \sum_{i=1}^n \alpha_i^* y_i K(\mathbf{x}_i, \mathbf{x}) + w_0^* = \sum_{i \in sv} \alpha_i^* y_i K(\mathbf{x}_i, \mathbf{x}) + w_0^*$

- w_0^* est choisi tel que $\gamma_i = f(\mathbf{x}_i) y_i \geq 1 - \xi_i = 1 - \frac{\alpha_i^*}{C}$

- la marge obtenue: $\gamma = \left(\sum_{i \in sv} \alpha_i^* - \frac{1}{C} \alpha^{*t} \alpha^* \right)^{-1/2}$

Machines à vecteurs de support

- **Problème soluble 2**: minimiser $\|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i$

- **Problème primal**

- minimiser par rapport à \mathbf{w} , $\boldsymbol{\xi}$ et w_0 et maximiser par rapport à $\boldsymbol{\alpha}$:

$$L(\mathbf{w}, w_0, \boldsymbol{\xi}, \boldsymbol{\alpha}, \mathbf{r}) = \frac{1}{2} \mathbf{w}^t \mathbf{w} + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i [(\mathbf{w}^t \mathbf{x}_i + w_0) y_i - 1 + \xi_i] - \sum_{i=1}^n r_i \xi_i$$

- sous les contraintes $\alpha_i \geq 0, r_i \geq 0 \quad i = 1, \dots, n$

- Problème dual

- les gradients:

$$\frac{\partial L(\mathbf{w}, w_0, \boldsymbol{\alpha})}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^n \alpha_i \mathbf{x}_i y_i = \mathbf{0}$$

$$\frac{\partial L(\mathbf{w}, w_0, \boldsymbol{\alpha})}{\partial \xi} = C - \alpha_i - r_i = 0$$

$$\frac{\partial L(\mathbf{w}, w_0, \boldsymbol{\alpha})}{\partial w_0} = \sum_{i=1}^n \alpha_i y_i = 0$$

Machines à vecteurs de support

- Problème dual

- resubstitution: maximiser par rapport à α :

$$W(\boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

- sous les contraintes $\alpha_i \geq 0, r_i \geq 0, C - \alpha_i - r_i = 0 \quad i = 1, \dots, n$ et

$$\sum_{i=1}^n \alpha_i y_i = 0$$

- \equiv sous les contraintes $0 \leq \alpha_i \leq C$ (contraintes de boîte)

- La solution

- $\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i = \sum_{i \in sv} \alpha_i^* y_i \mathbf{x}_i$

- $f^*(\mathbf{x}) = \sum_{i=1}^n \alpha_i^* y_i K(\mathbf{x}_i, \mathbf{x}) + w_0^* = \sum_{i \in sv} \alpha_i^* y_i K(\mathbf{x}_i, \mathbf{x}) + w_0^*$

- w_0^* est choisi tel que $\gamma_i = f(\mathbf{x}_i) y_i = 1$ pour tous $i : 0 < \alpha^* < C$

- la marge obtenue: $\gamma = \left(\sum_{i,j \in sv} \alpha_i^* \alpha_j^* y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) \right)^{-1/2}$

- If you can **build a generative model** for the class densities, do it
- All algorithms are **optimized for classification, not hypothesis testing**
 - doesn't mean they don't work, but they are probably not optimal
- None of the margin-based algorithms like **label noise**
- Non-parametric methods also exist for other objectives (e.g., **regression, density estimation**)

- `mloss.org`
- `www.cs.waikato.ac.nz/ml/weka`
- `multiboost.org`
- `svmlight.joachims.org`
- `www.csie.ntu.edu.tw/~cjlin/libsvm`
- `www.torch.ch`

Thank you!