

Préliminaires

- Balázs Kégl (prononcé “**Bolage**”)
- balazs.kegl@gmail.com
- **Arrêtez-moi** si qqchose n’est pas clair!
- **Posez des questions!**

- Introduction
- Exemple : OCR – MNIST
- Une solution : AdaBoost
- Extraction des traits : filtres de Haar

Algorithmes “classiques” vs. d’apprentissage

- Approche classique
 - **description formelle** des contraintes de l’entrée et de la sortie souhaitée
 - **compréhension du problème** computationnel
 - design d’une solution algorithmique **basée sur ces connaissances**
- Problèmes
 - **Connaissances incomplètes**
 - Algorithme trop **coûteux**

Algorithmes “classiques” vs. d’apprentissage

- Approche d’apprentissage
 - données (exemples) de forme (entrée, sortie)
 - compréhension **partielle** du problème : connaissances a-priori
 - apprendre : **chercher dans un ensemble** de fonctions

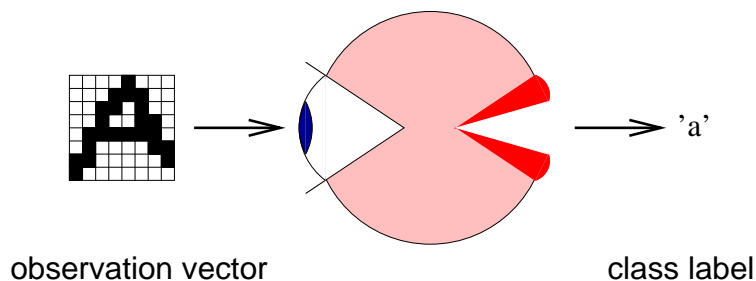
Catégories d'apprentissage automatique

- Classification
 - classifier le nouvel exemple
- Régression
 - faire une prédiction à partir du nouvel exemple
- Estimation de densité
 - dire si le nouvel exemple ressemble aux exemples déjà vus
- (Modèles graphiques)

Exemple de classification

- OCR = Optical character recognition
 - **entrée** : images de caractères (chiffres) manuscrites
 - **sortie** : l'étiquette ou classe
 - **but** : design d'une **fonction de classification** :

$$f : \{\text{image}\} \rightarrow \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$



Exemple de classification

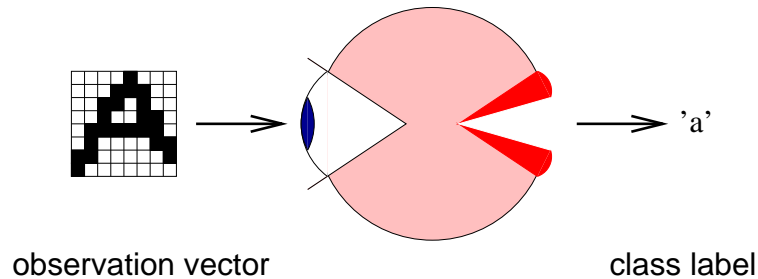
- MNIST (yann.lecun.com/exdb/mnist/)



Exemple de classification

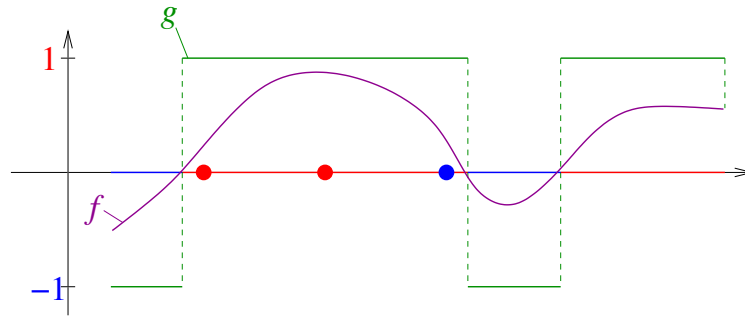
- MNIST (yann.lecun.com/exdb/mnist/)
 - extrait de la base de données NIST
 - images **grayscale**, **28 × 28**
 - **60000** entraînement, **10000** test

Le modèle d'apprentissage supervisé



- vecteur d'observation: $\mathbf{x} \in \mathbb{R}^d$
- étiquette de classe: $y \in \{-1, 1\}$ (ou $y \in \{1, \dots, K\}$)
- classifieur: $g : \mathbb{R}^d \rightarrow \{-1, 1\}$

Le modèle d'apprentissage supervisé



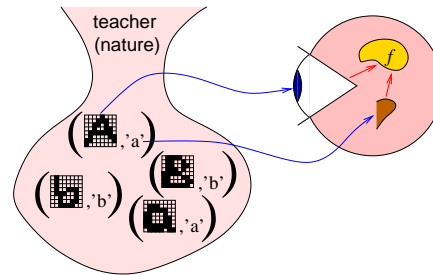
- Fonction discriminante: $f : \mathbb{R}^d \rightarrow [-1, 1]$

- \longrightarrow fonction de classification

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

- frontière de décision: $\{\mathbf{x} : f(\mathbf{x}) = 0\}$

Le modèle d'apprentissage supervisé



- Apprentissage par l'expérience, avec un superviseur

- échantillon d'entraînement: $D_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$

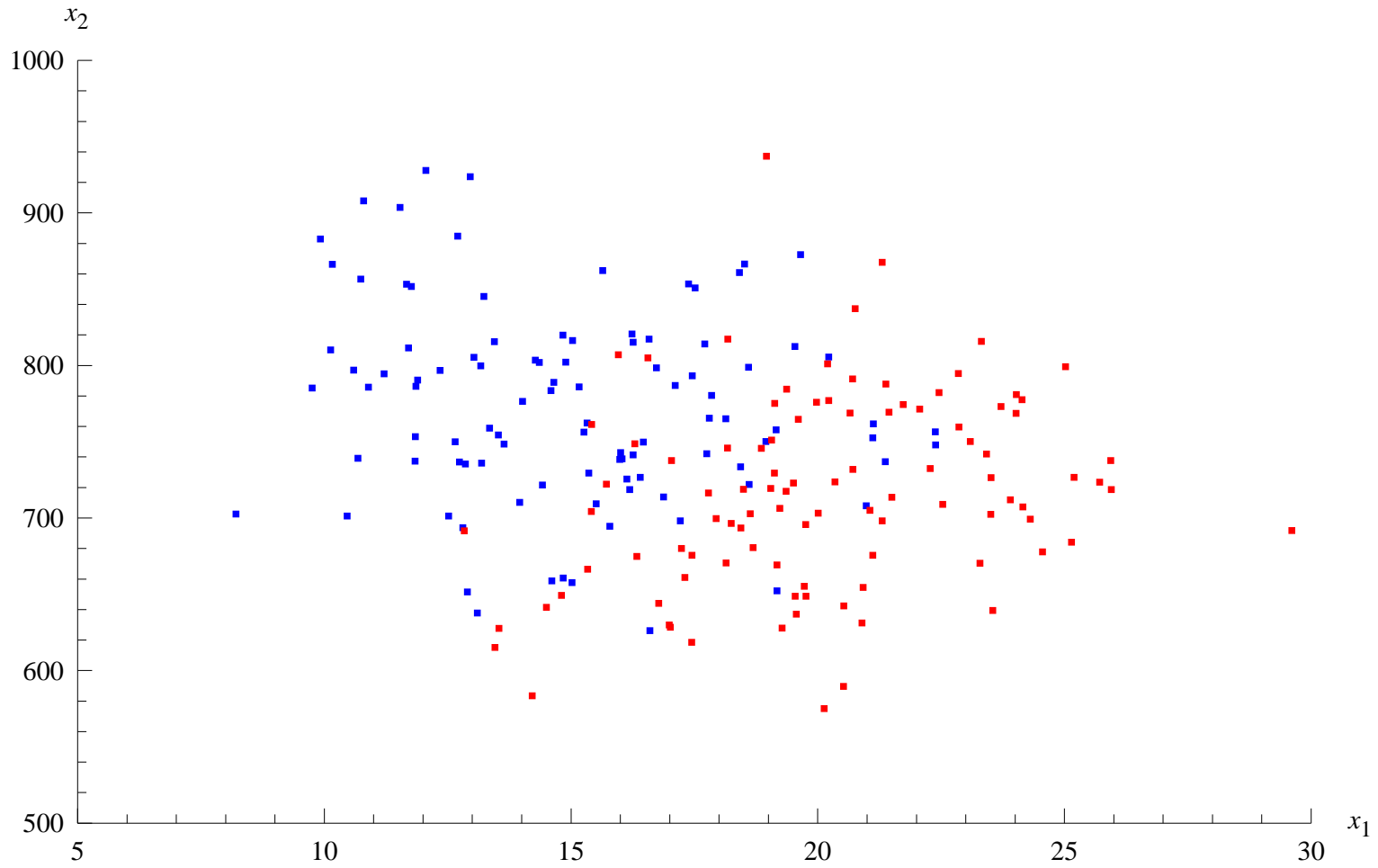
- ensemble des fonctions: \mathcal{F}

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

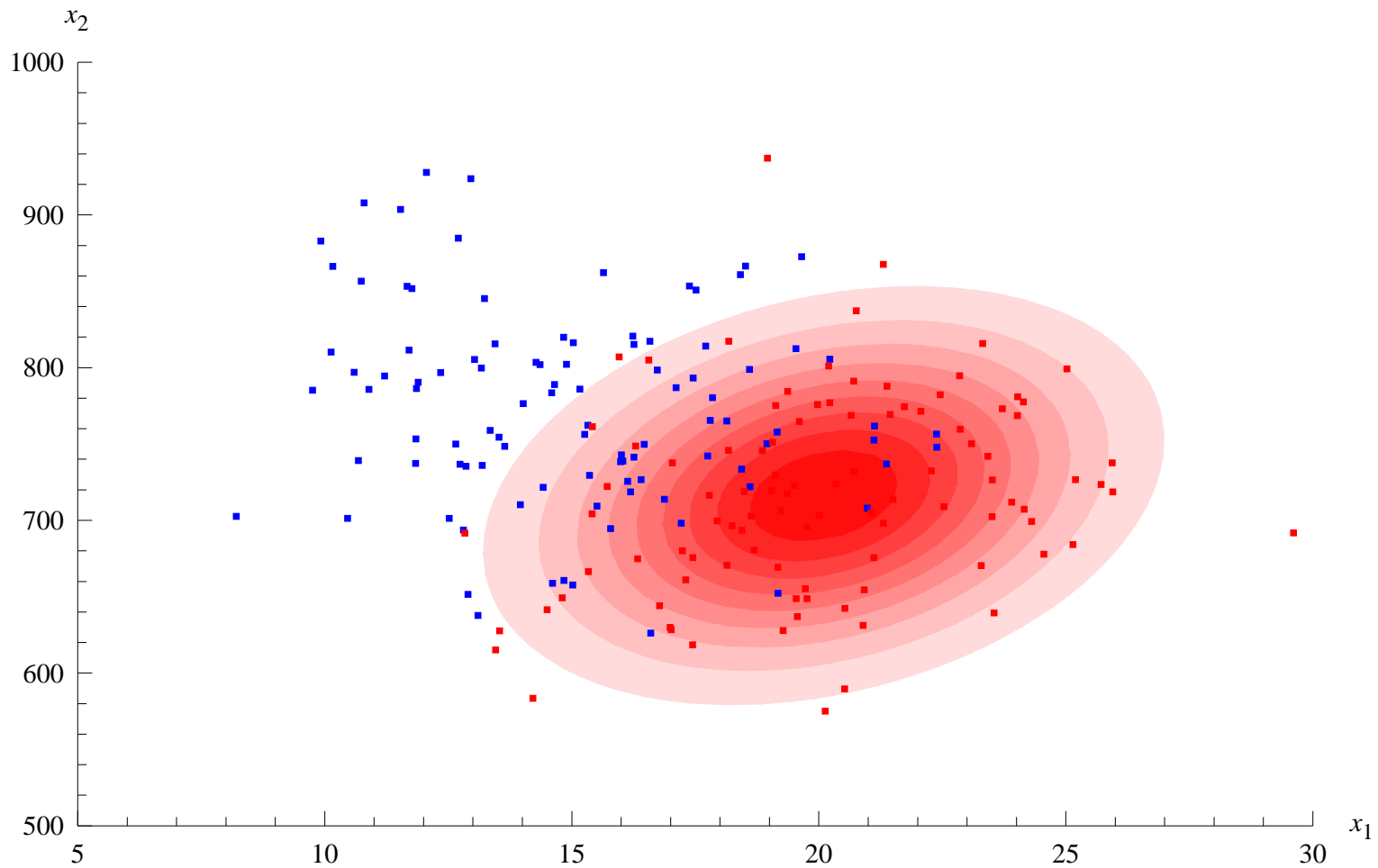
$$\text{ALGO}(D_n) \mapsto f$$

- but: petite erreur de généralisation $P[g(\mathbf{X}) \neq Y] = P[f(\mathbf{X})Y \leq 0]$

Data for two-class classification problem



2-D Gaussian fit for class 2



- Terminology

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

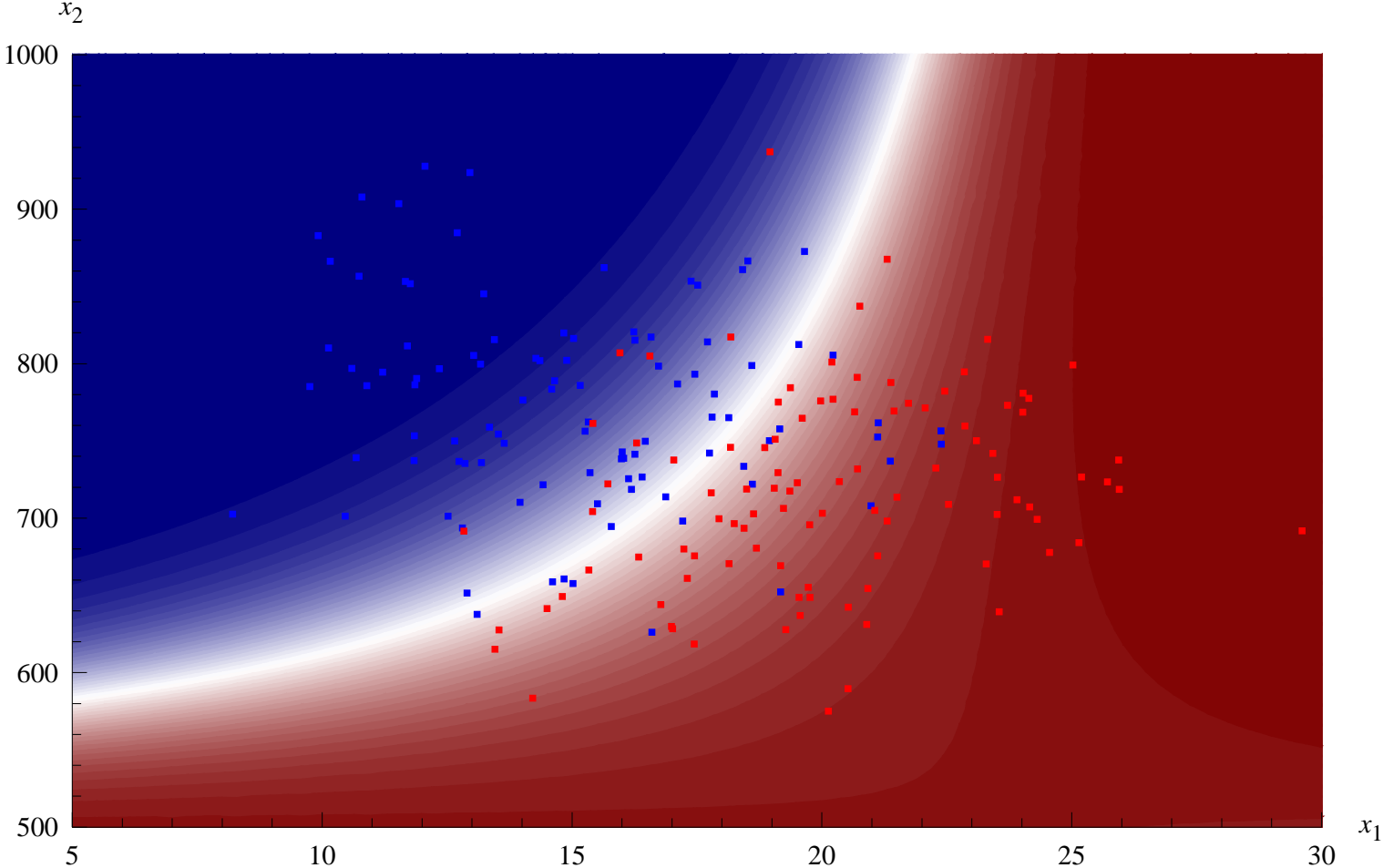
- **Bayes** theorem:

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

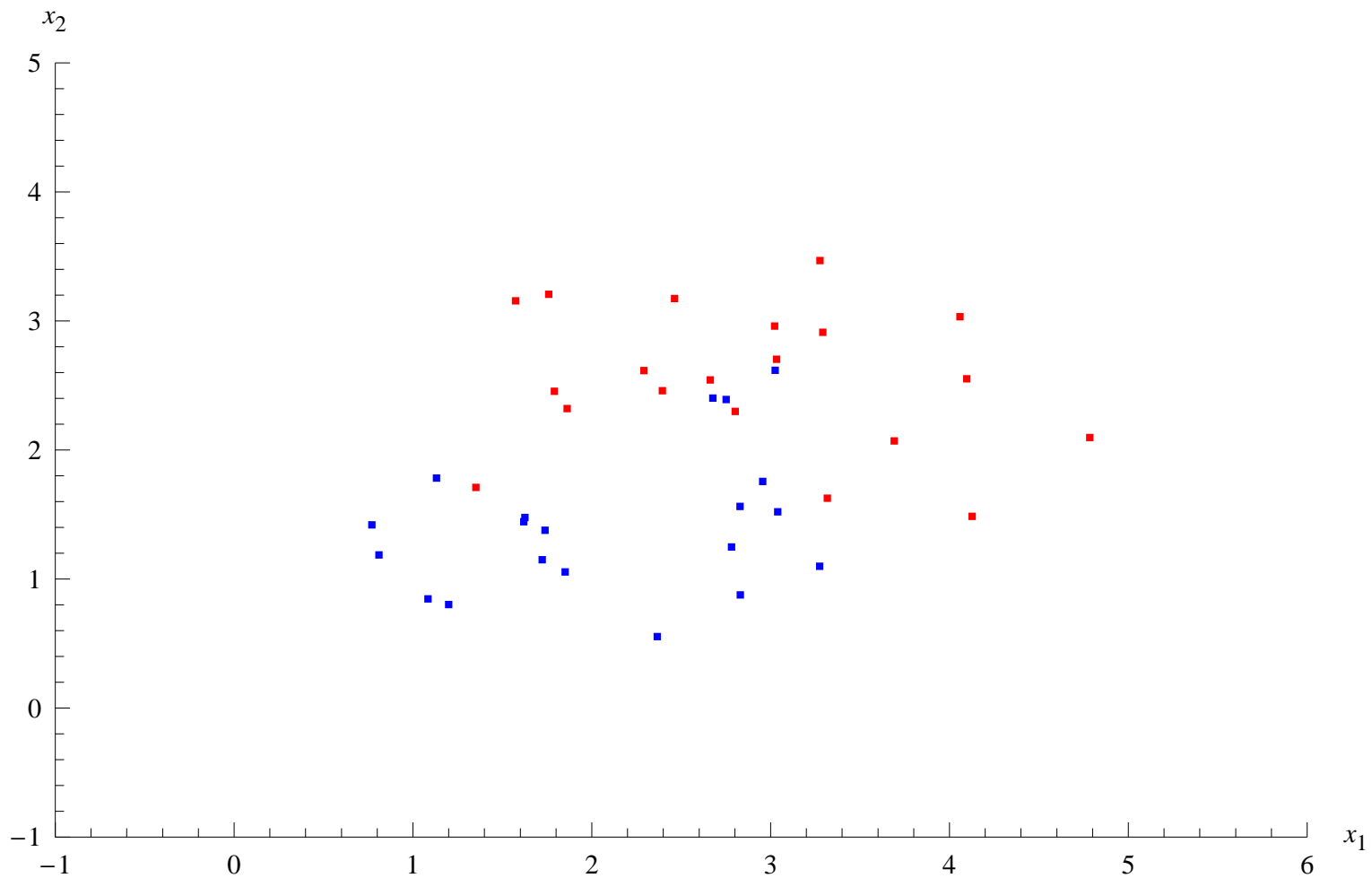
- **Decision:**

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

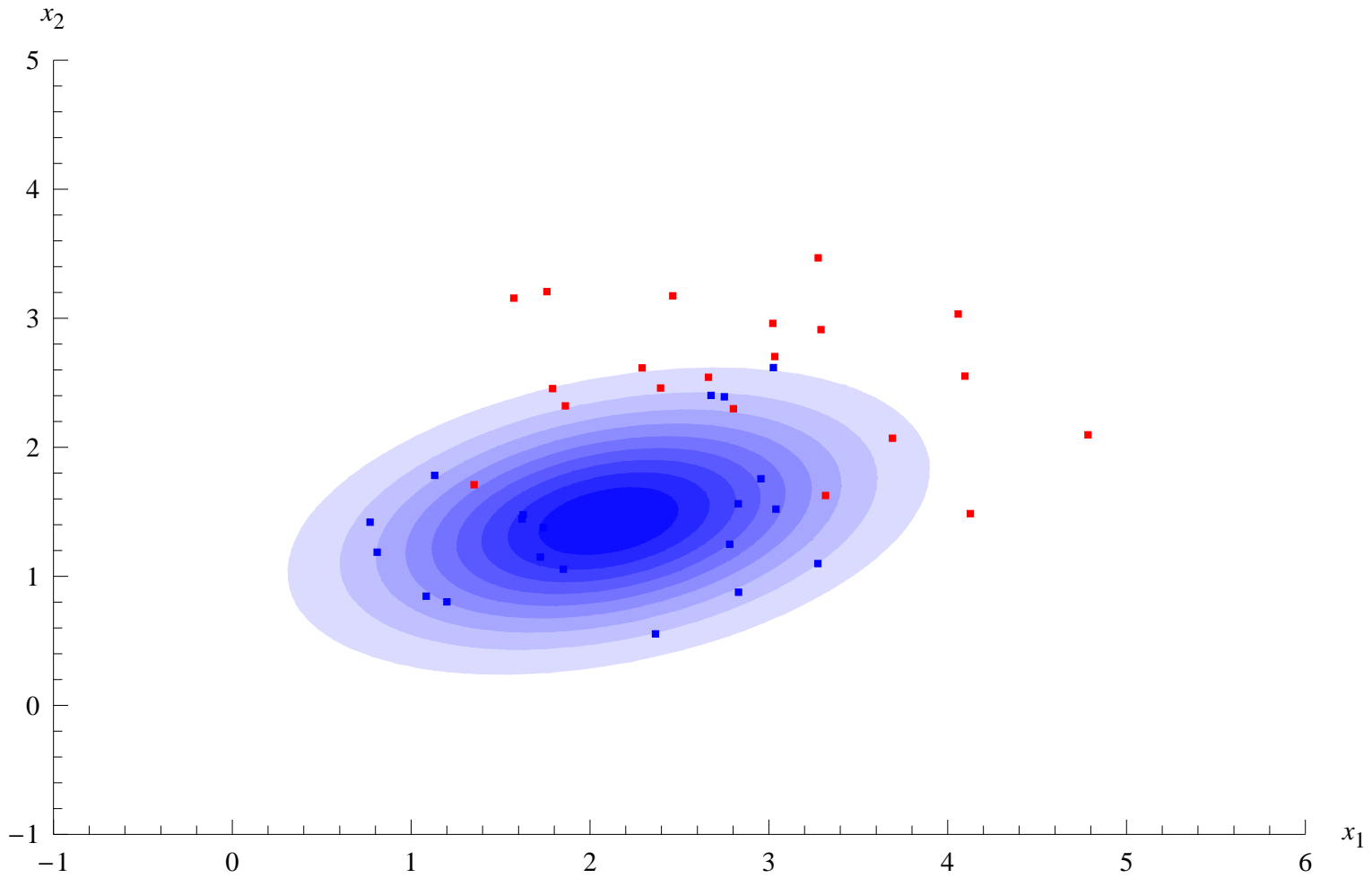
Discriminant function with Gaussian fits



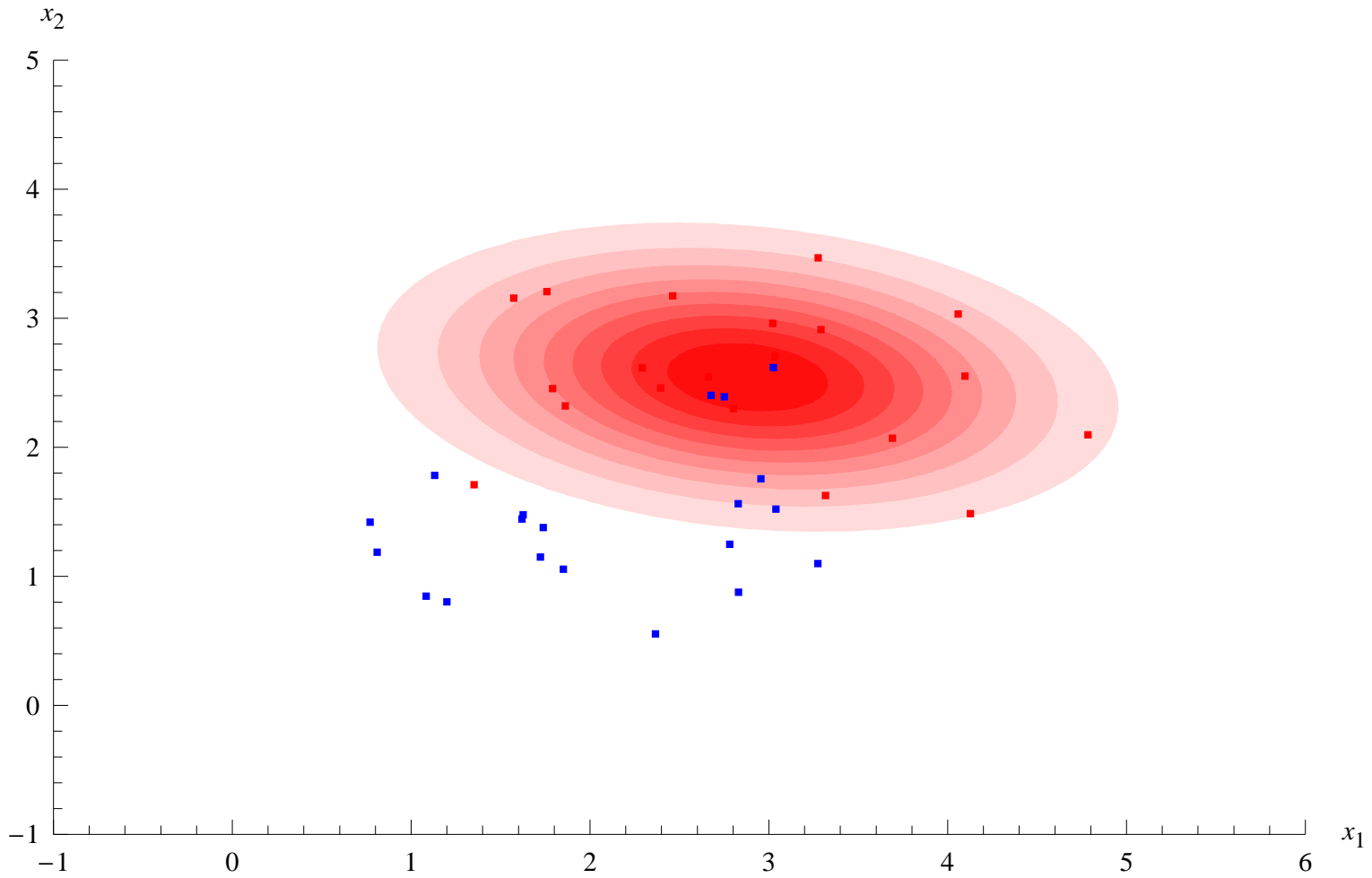
'Two Moons' data for two-class classification problem



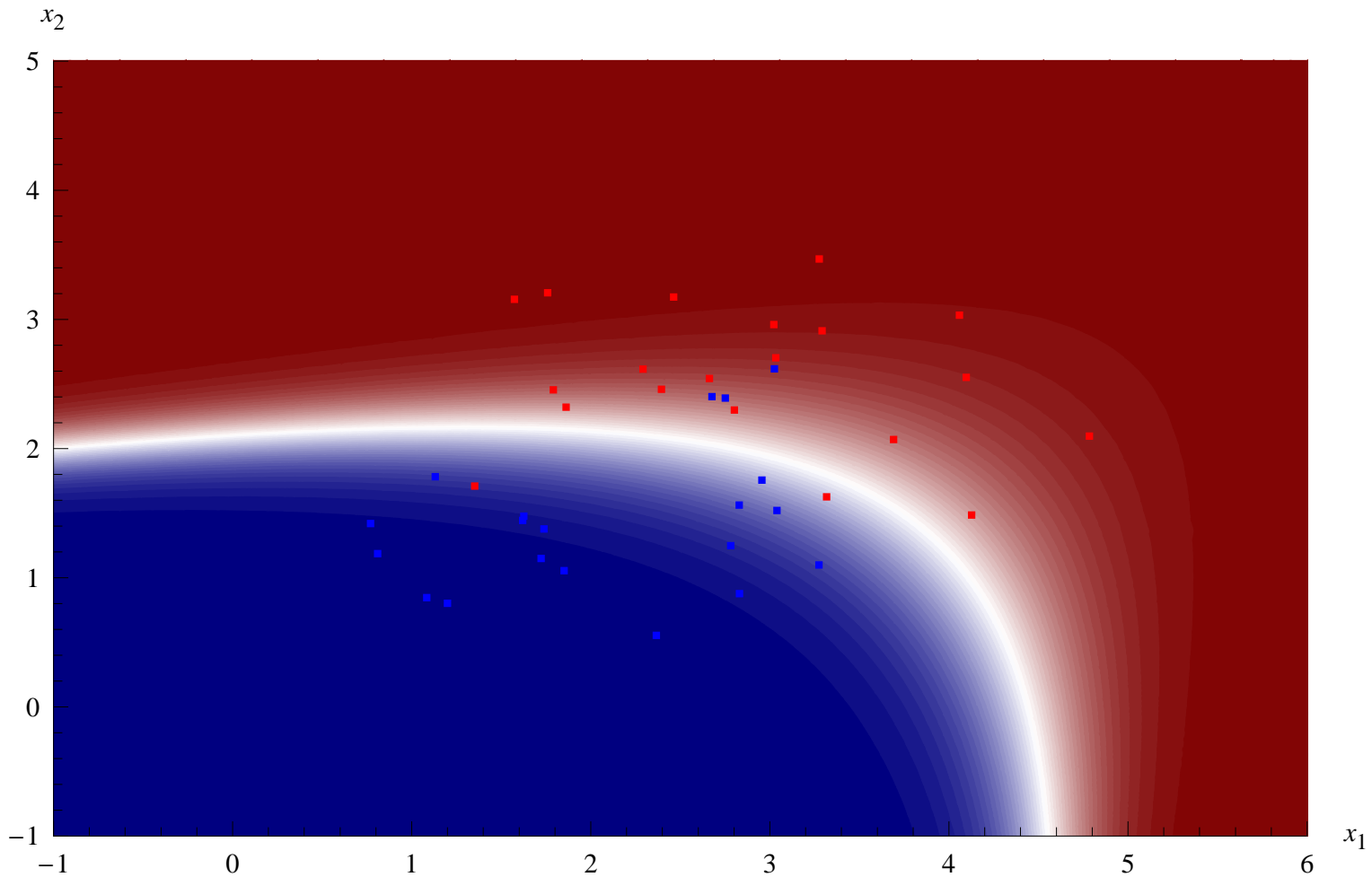
2-D Gaussian fit for class 1

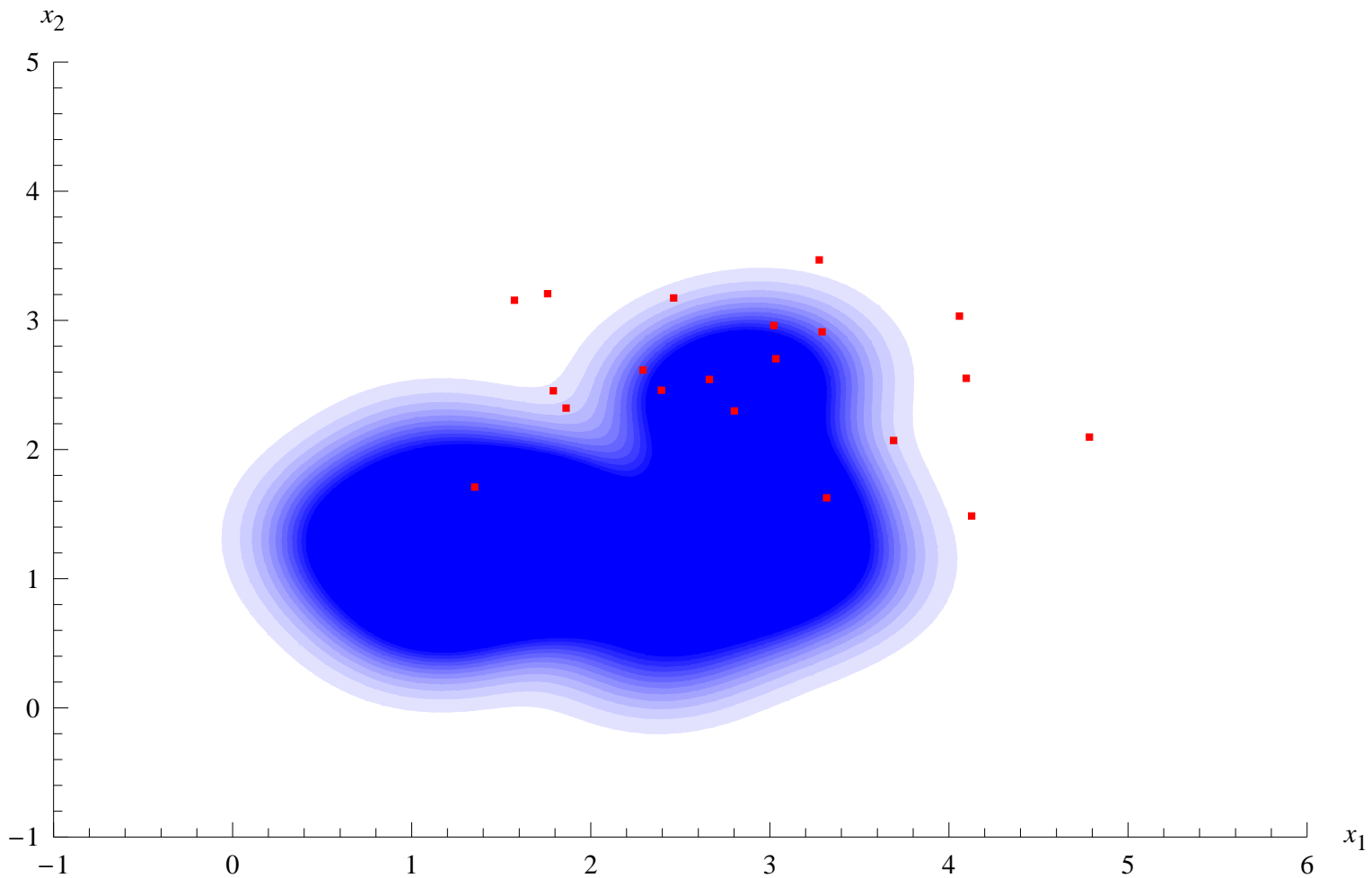


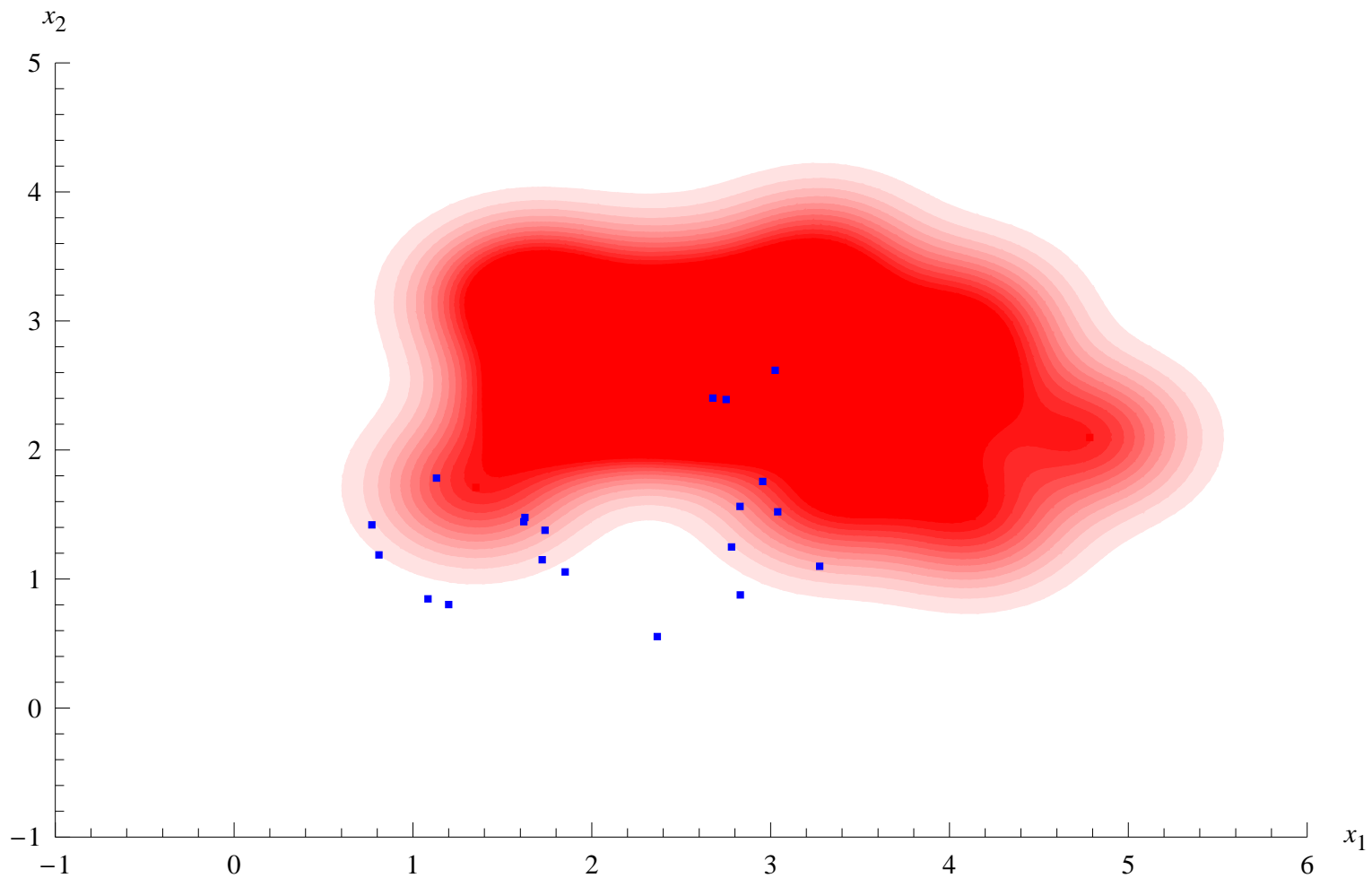
2-D Gaussian fit for class 2

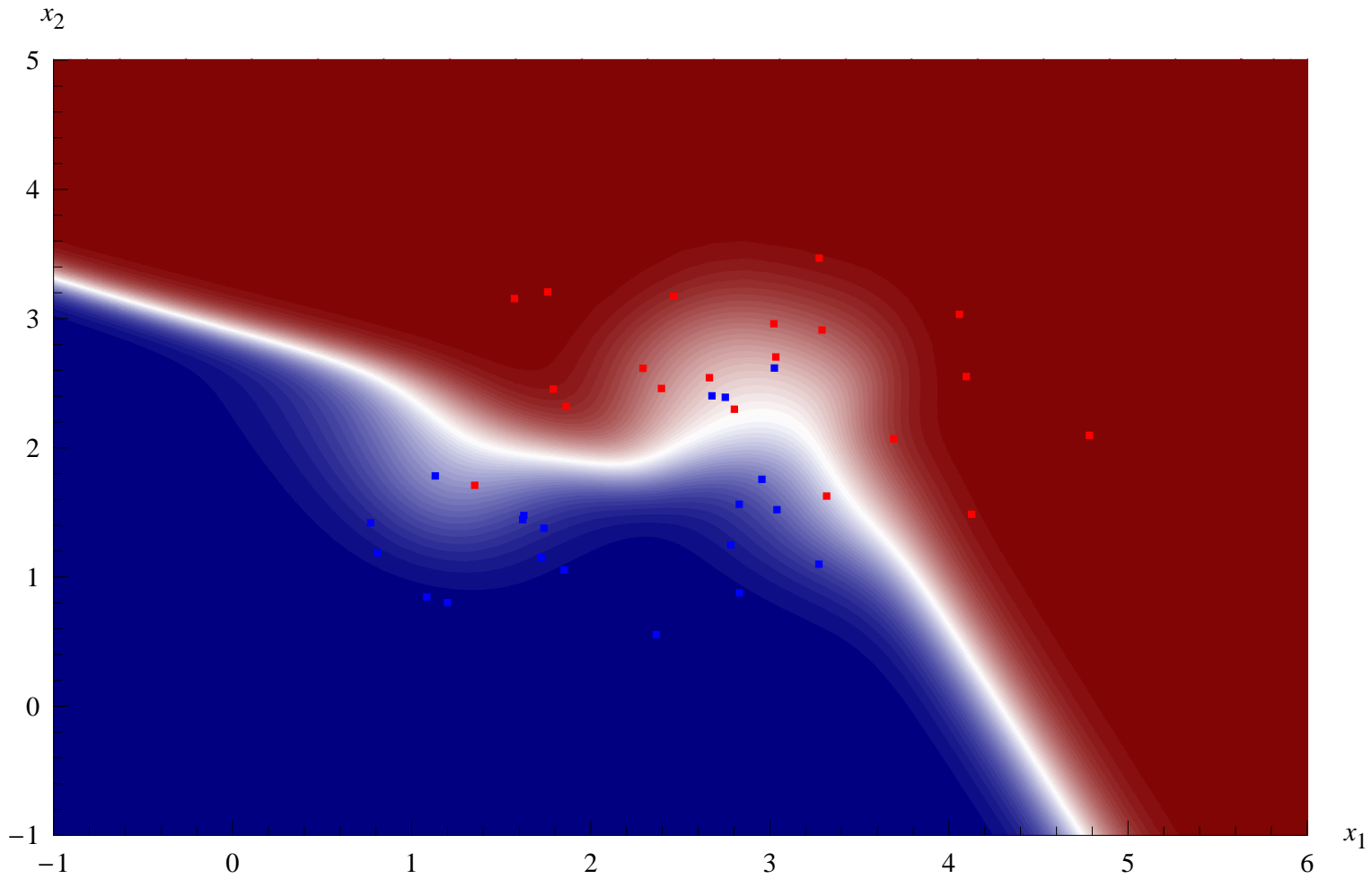


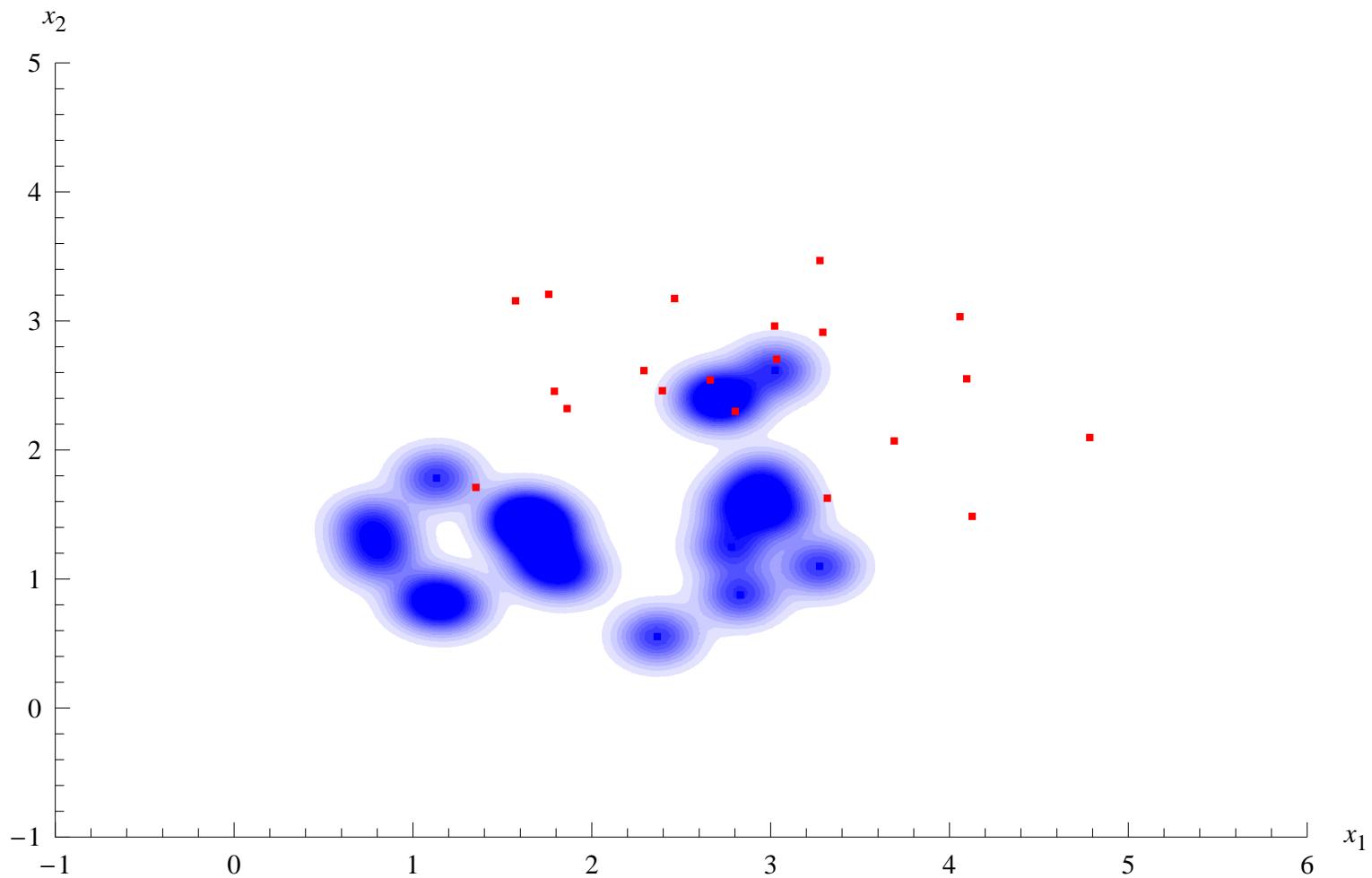
Discriminant function with Gaussian fits

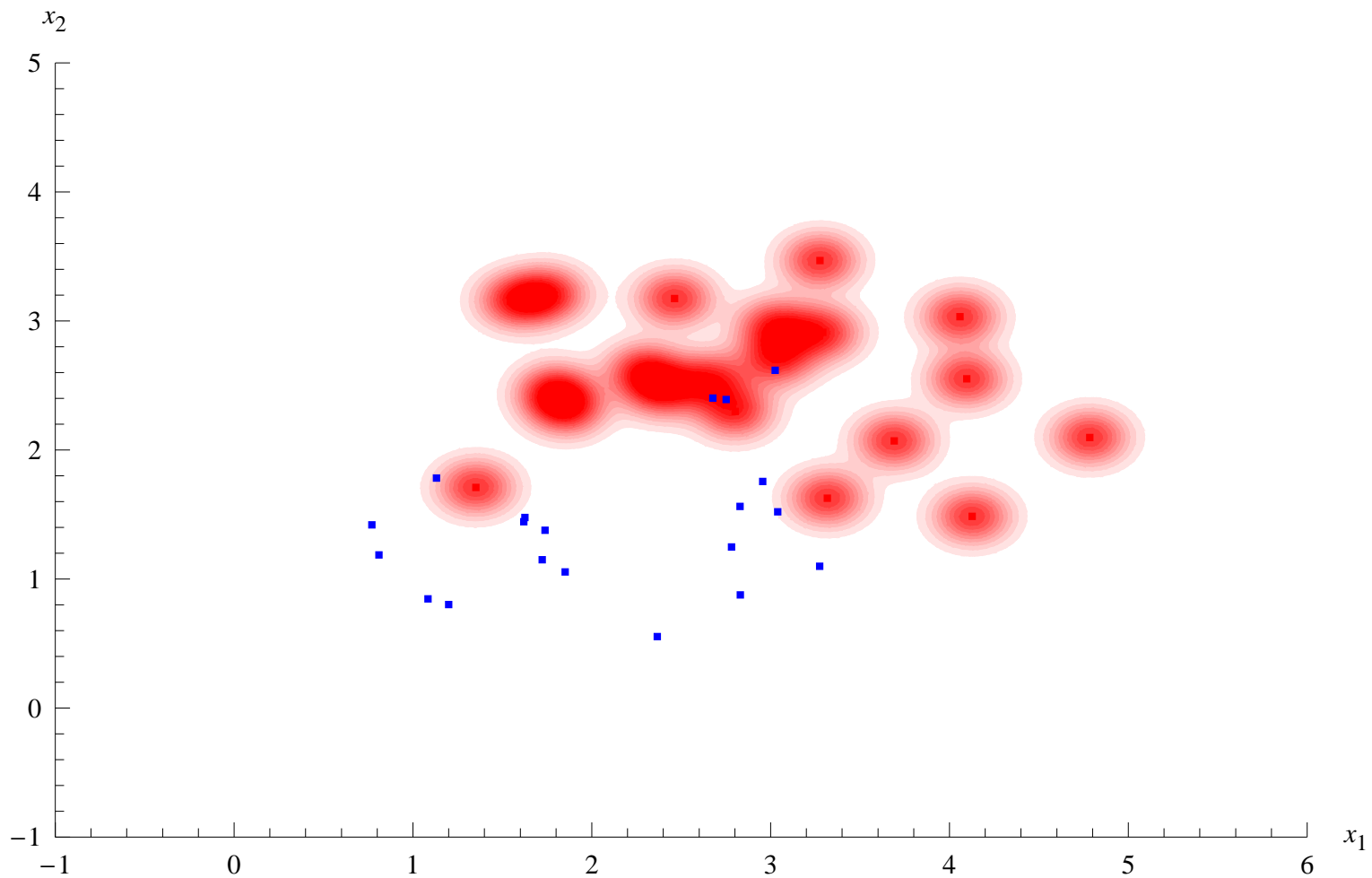


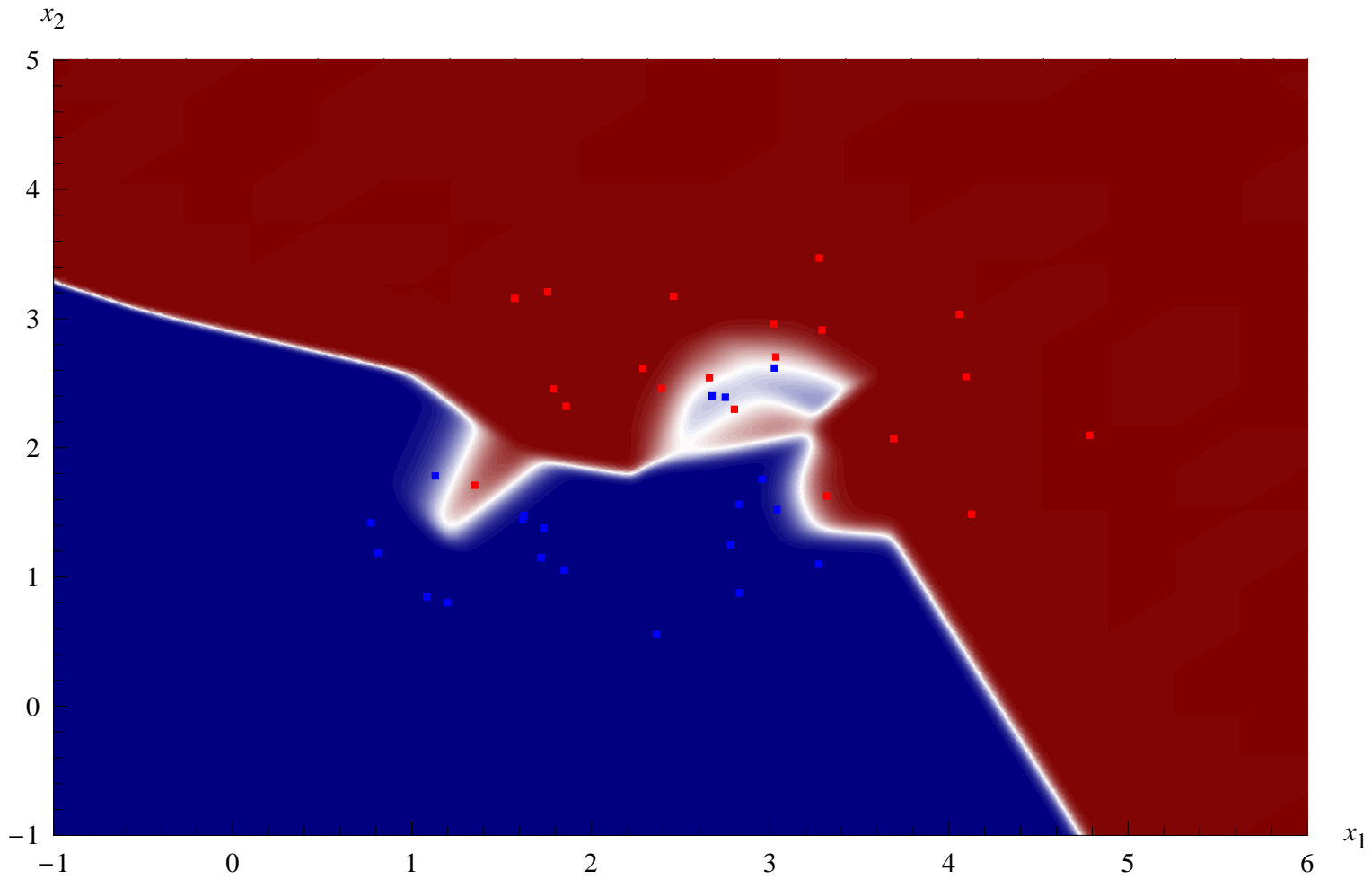
2-D Parzen fit for class 1, $h = 0.12$ 

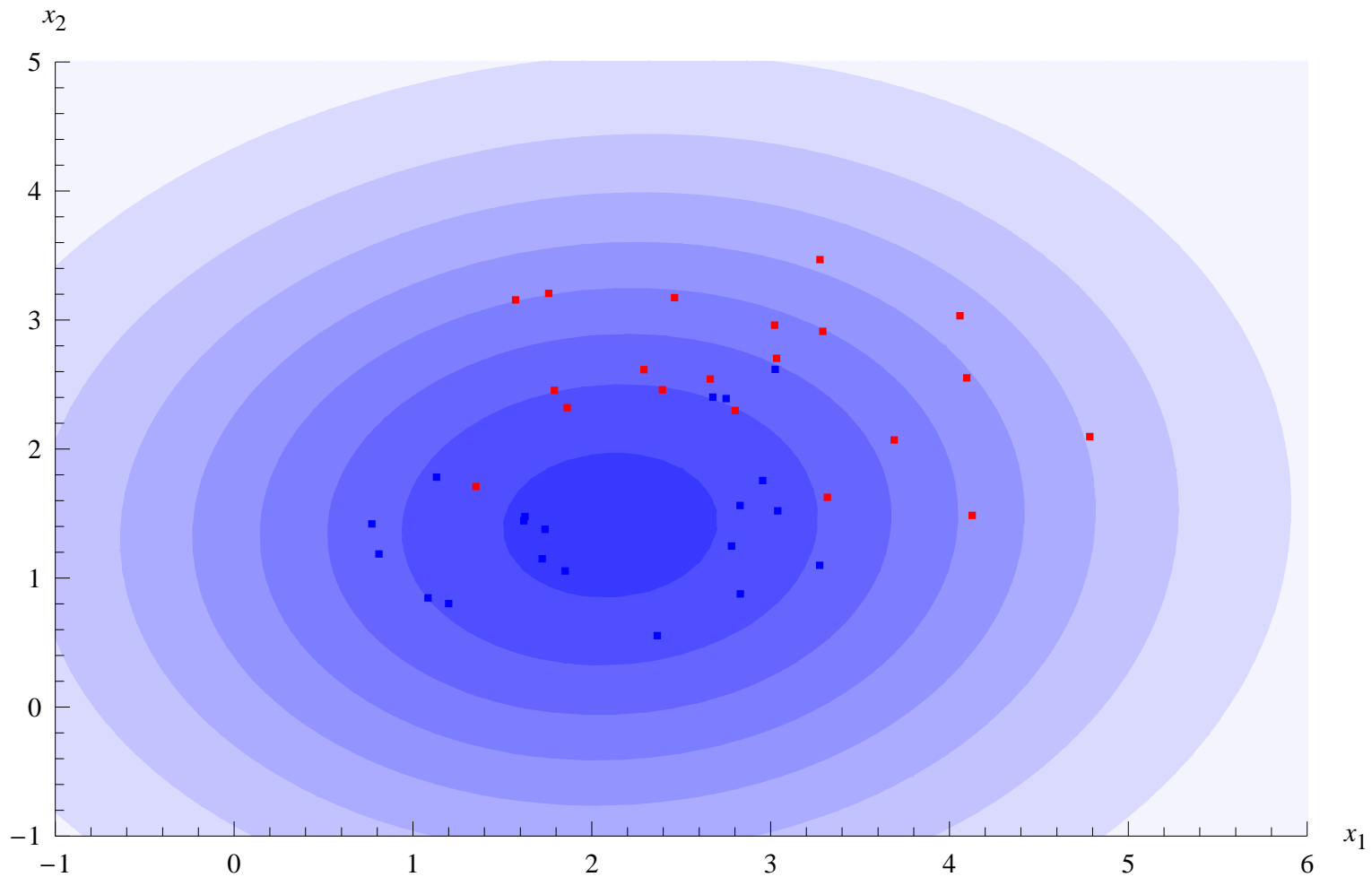
2-D Parzen fit for class 2, $h = 0.12$ 

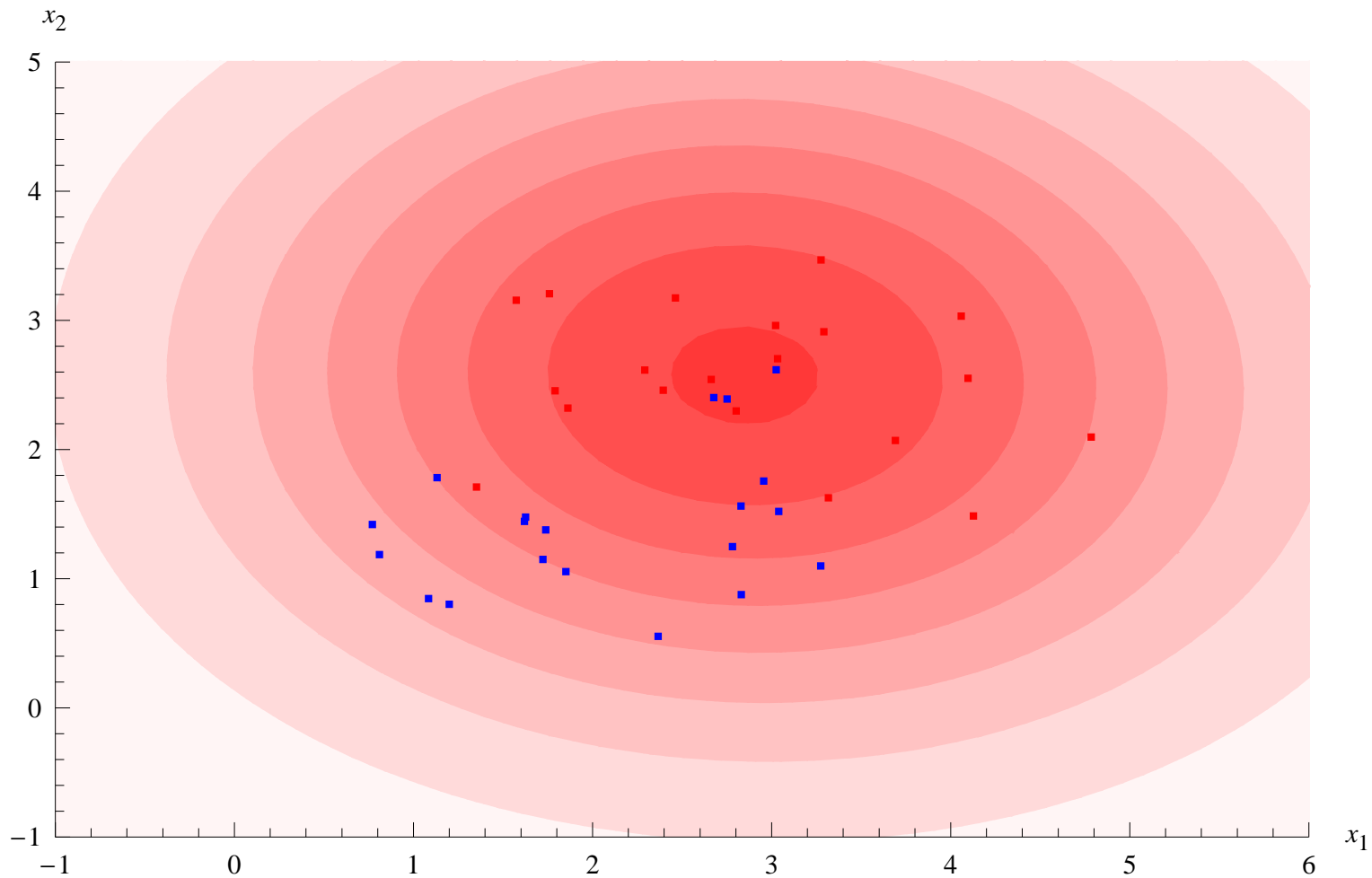


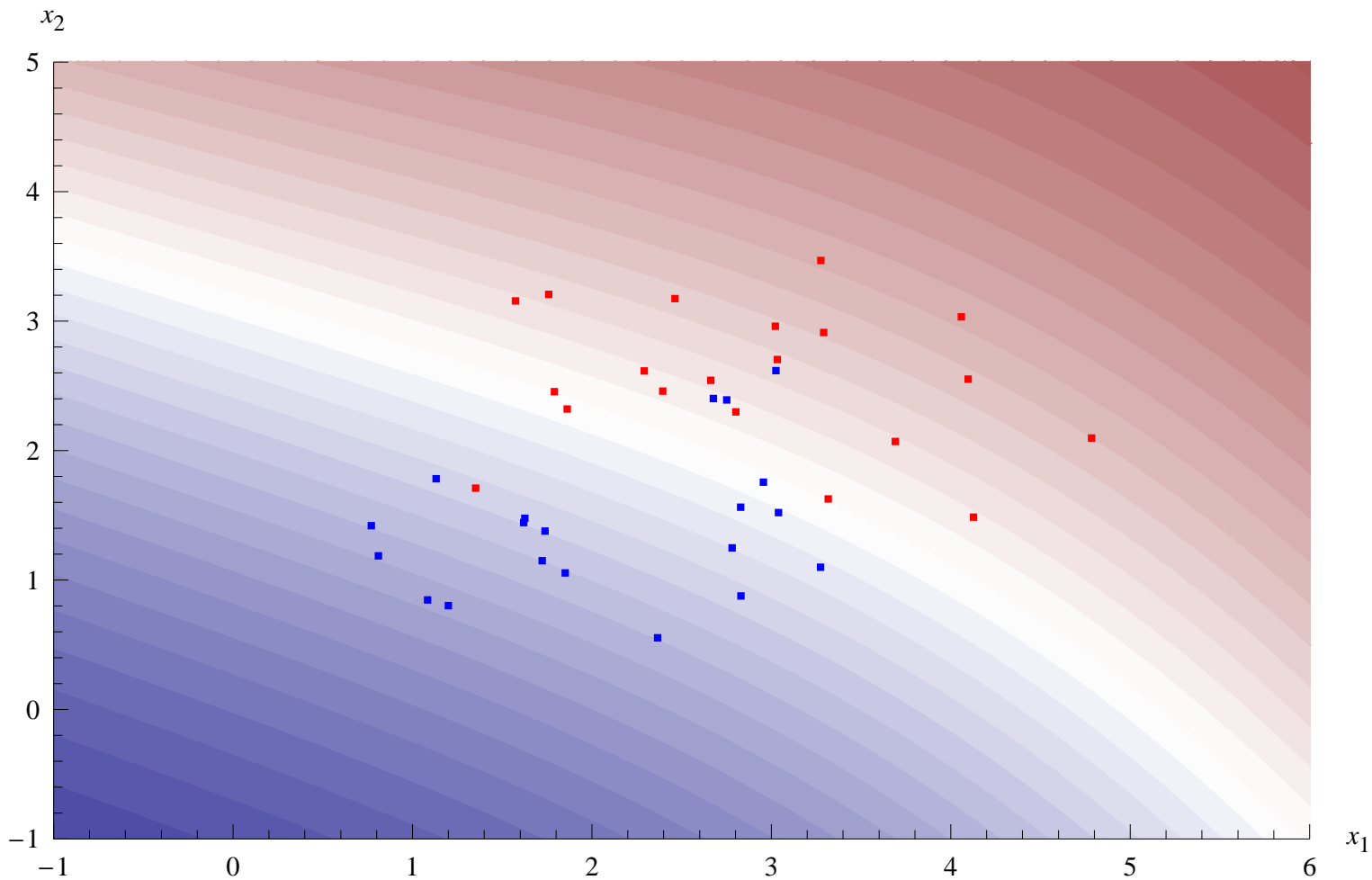
2-D Parzen fit for class 1, $h = 0.02$ 

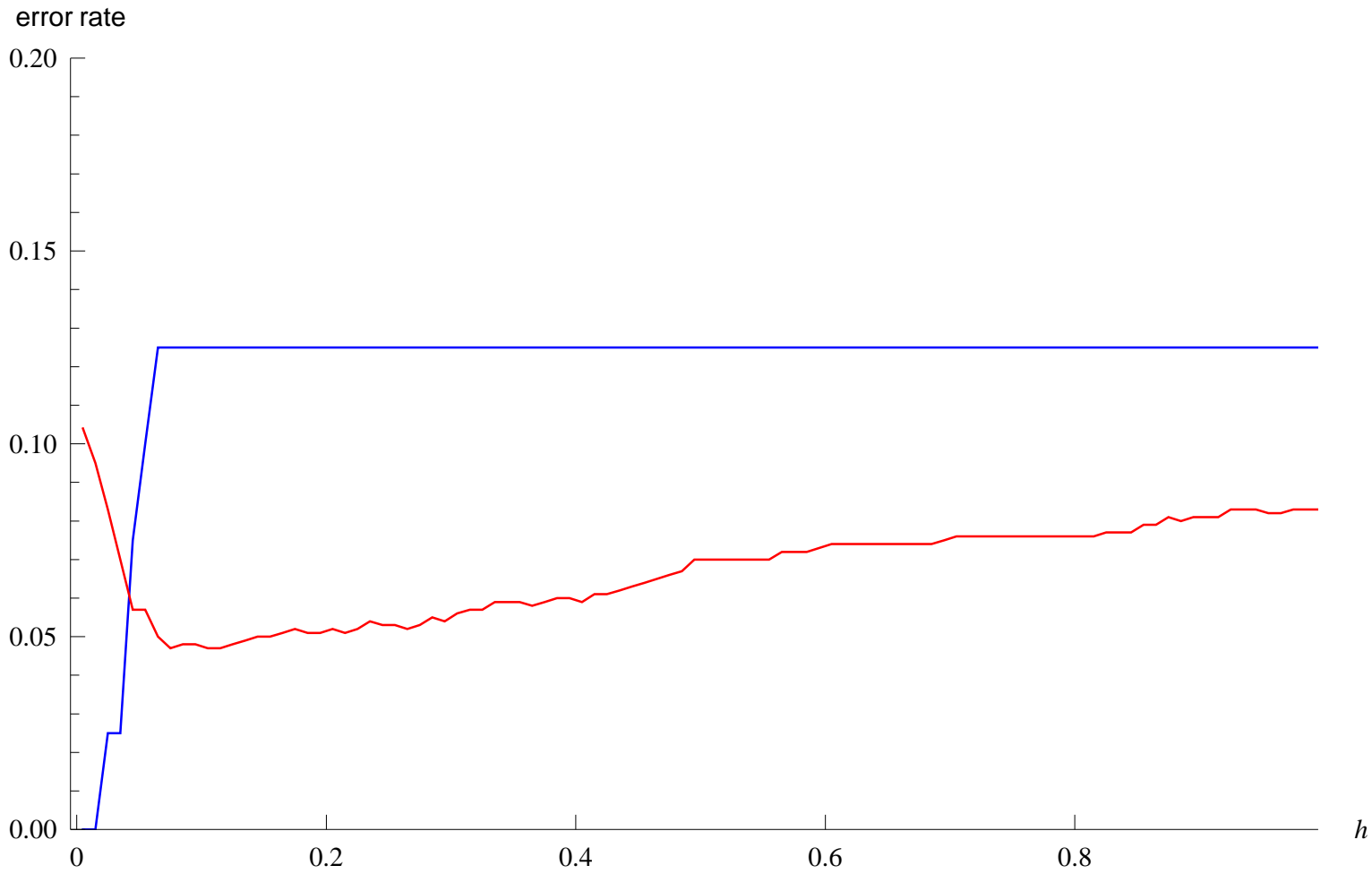
2-D Parzen fit for class 2, $h = 0.02$ 



2-D Parzen fit for class 1, $h = 3$ 

2-D Parzen fit for class 2, $h = 3$ 

Discriminant function with Parzen fits, $h = 3$ 



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

- Algorithmes

- 1958 : **Perceptron** [Rosenblatt, '58] – [Minsky–Papert '69]
- 1986 : **Réseaux de neurones** et l'algorithme de **rétro-propagation** [Rumelhart–Hinton–Williams, '86]
- 1995: **Machines à vecteurs de support** [Boser–Guyon–Vapnik, '92], [Cortes–Vapnik, '95]
- 1997: **boosting, AdaBoost** [Freund, '95], [Freund–Schapire, '97]

- Classification **linéaire**
 - Perceptron
 - Machine à vecteurs de support (SVM) linéaire
- Classification **non-linéaire**
 - Perceptron multi-couche (réseaux de neurones)
 - Plus proche voisin
 - Arbres de décision
 - SVM
 - **Boosting**

Le plus proche voisin

- Le plus proche voisin:

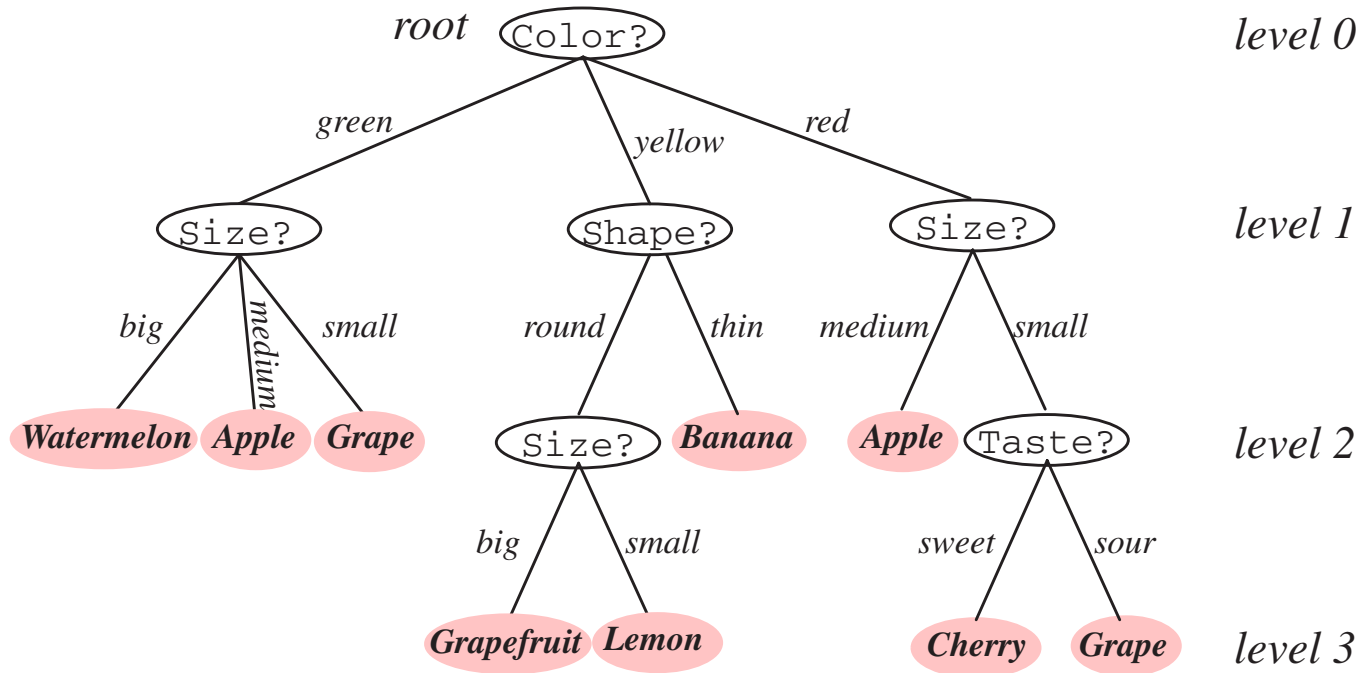
$$f(\mathbf{x}) = y_{i_{\mathbf{x}}}$$

ou

$$i_{\mathbf{x}} = \arg \min_{i:1,\dots,n} d(\mathbf{x}, \mathbf{x}_i)$$

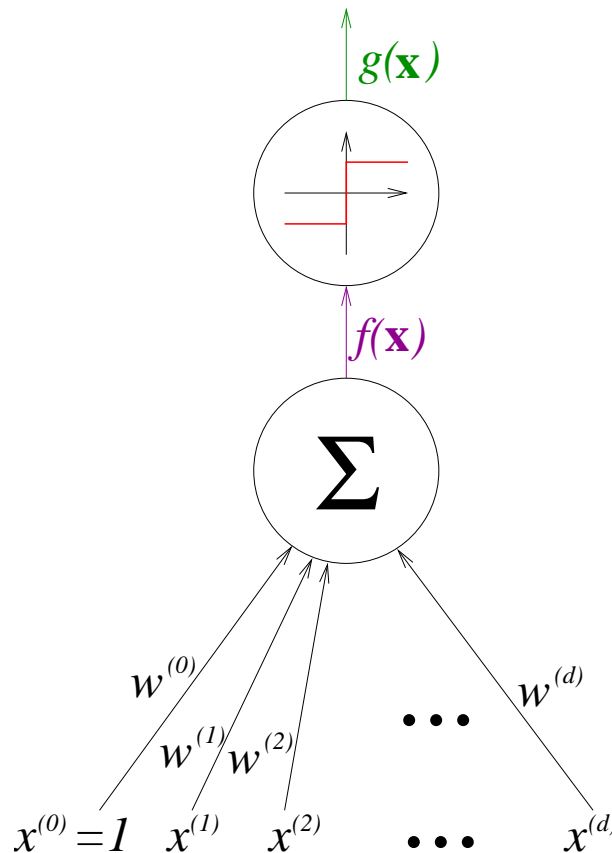
- très simple
- problèmes algorithmiques et statistiques

Arbres de décision



Le perceptron

- Fonctions discriminantes **linéaires** : $f(\mathbf{x}) = \sum_{i=0}^d w^{(i)} \cdot x^{(i)} = \langle \mathbf{w}, \mathbf{x} \rangle$



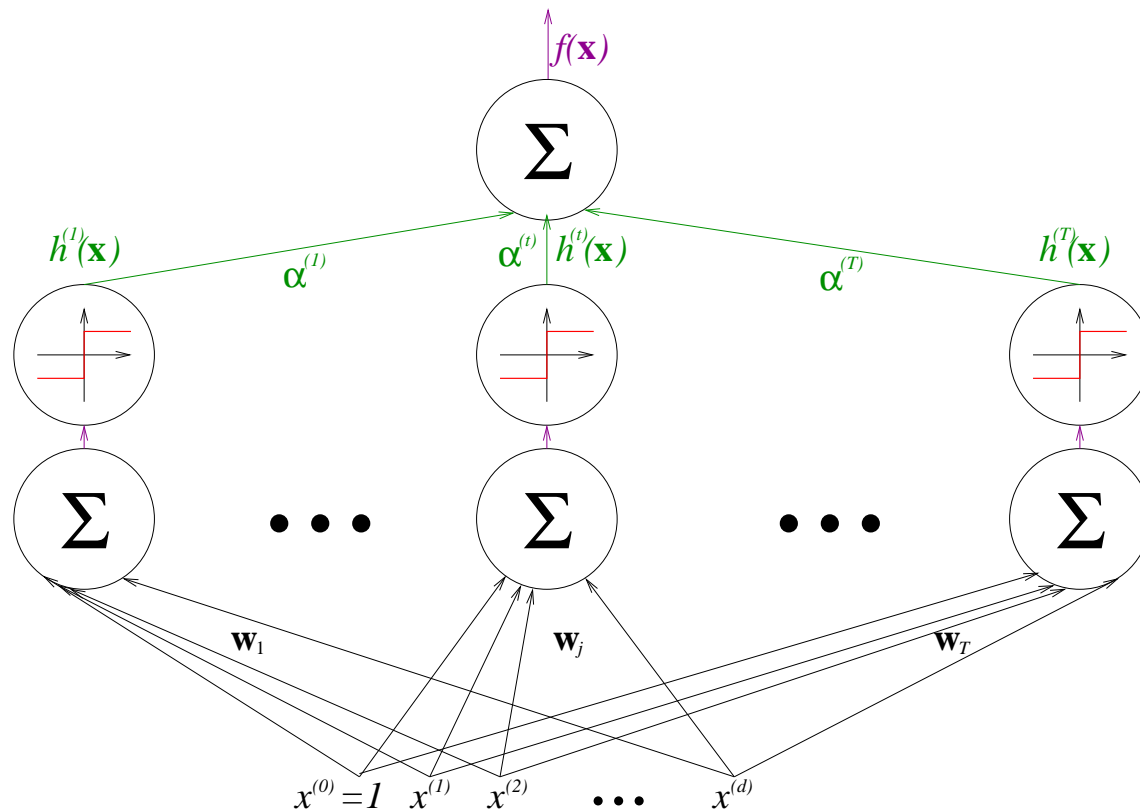
- Modèle :

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

- $h^{(j)} : \mathbb{R}^d \rightarrow [-1, 1]$
 - **classifieurs**/fonctions discriminantes **simples**, **traits**, **experts**
- $\alpha^{(j)} \in \mathbb{R}^+$
 - poids de l'expert $h^{(j)}$ dans le **vote final**

Le perceptron multi-couche

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



- Modèle :

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

- $I_{sv} \subset \{1, \dots, n\}$ est l'ensemble de **vecteurs de support**
- $K(\cdot, \cdot)$ est une fonction de **similarité** (noyau)
- But : une **frontière équidistante** des classes
- “plus proche voisin sophistiqué”

- Modèle :

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

- Noyaux :

- $K(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle \longrightarrow f(\mathbf{x})$ est linéaire
- $K(\mathbf{x}, \mathbf{x}') = (1 + \langle \mathbf{x}, \mathbf{x}' \rangle)^d \longrightarrow f(\mathbf{x})$ est un polynôme de degré d
- $K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2) \longrightarrow f(\mathbf{x})$ est un mélange de gaussiens

- Modèle :

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

- aucune restriction sur la forme de $h^{(j)}(\mathbf{x})$
- souvent des “decision stumps” :

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

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

- L'intuition de l'algorithme
 - ajouter un expert (stump) à la fois
 - ajouter le meilleur sur les points d'entraînement mal-classifiés par des experts précédents
 - mettre le poids de l'expert choisi proportionnel à son erreur
 - à la fin, retourner le vote pondéré

- **Pondération sur les points** d'entraînement w_1, \dots, w_n
 - pondération **normalisée** : $\sum_{i=1}^n w_i = 1$
 - initialisée **uniformément** : $\mathbf{w} = (1/n, \dots, 1/n)$
 - si \mathbf{x}_i est **mal-classifié** par $h^{(j)}$, **augmenter** w_i
 - sinon, **diminuer** w_i
 - au fur et à mesure les **points "difficiles"** obtiennent des poids élevés

- Choix d'expert

- experts **binaires** : $h^{(j)} : \mathcal{X} \rightarrow \{-1, 1\}$

- erreur **pondérée**

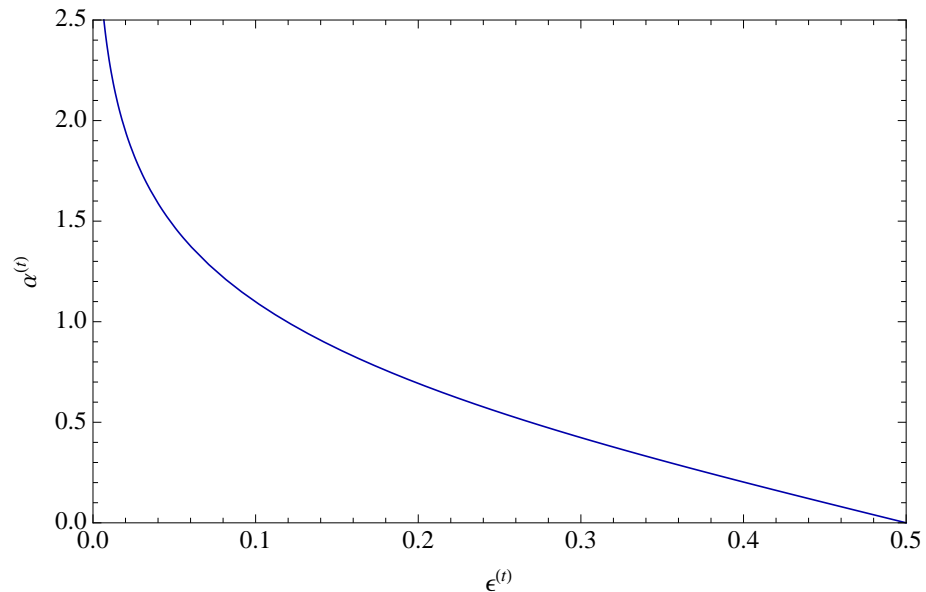
$$\varepsilon = \varepsilon(h) = \sum_{i=1}^n w_i \mathbb{I}\{h(\mathbf{x}_i) \neq y_i\}$$

- choisir l'expert qui **minimise l'erreur pondérée**

$$h^{(t)} = \arg \min_h \varepsilon(h)$$

- Choix d'expert
 - mettre son poids à

$$\alpha^{(t)} = \frac{1}{2} \ln \frac{1 - \epsilon}{\epsilon}$$



- Re-pondération des points

- si $h(\mathbf{x}_i) \neq y_i$ alors $w_i \leftarrow w_i \frac{1}{2\varepsilon}$

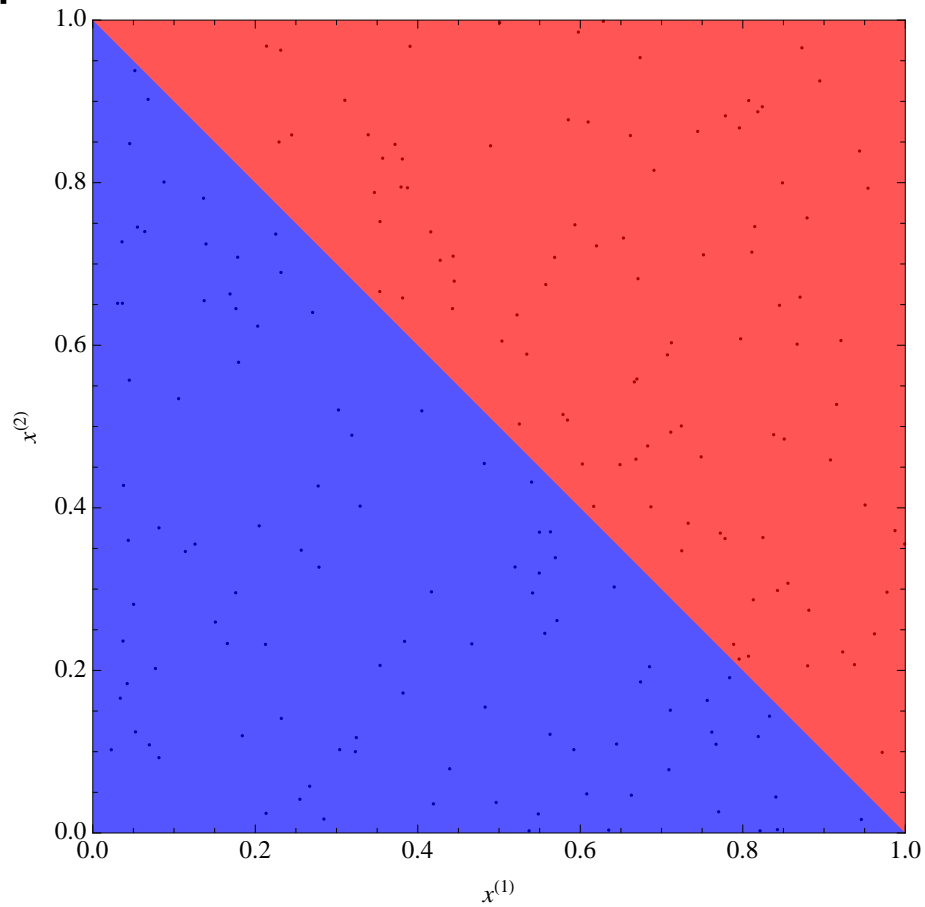
- si $h(\mathbf{x}_i) = y_i$ alors $w_i \leftarrow w_i \frac{1}{2(1-\varepsilon)}$

ADABOOST($D_n, \text{BASE}(D_n, \mathbf{w}), T$)

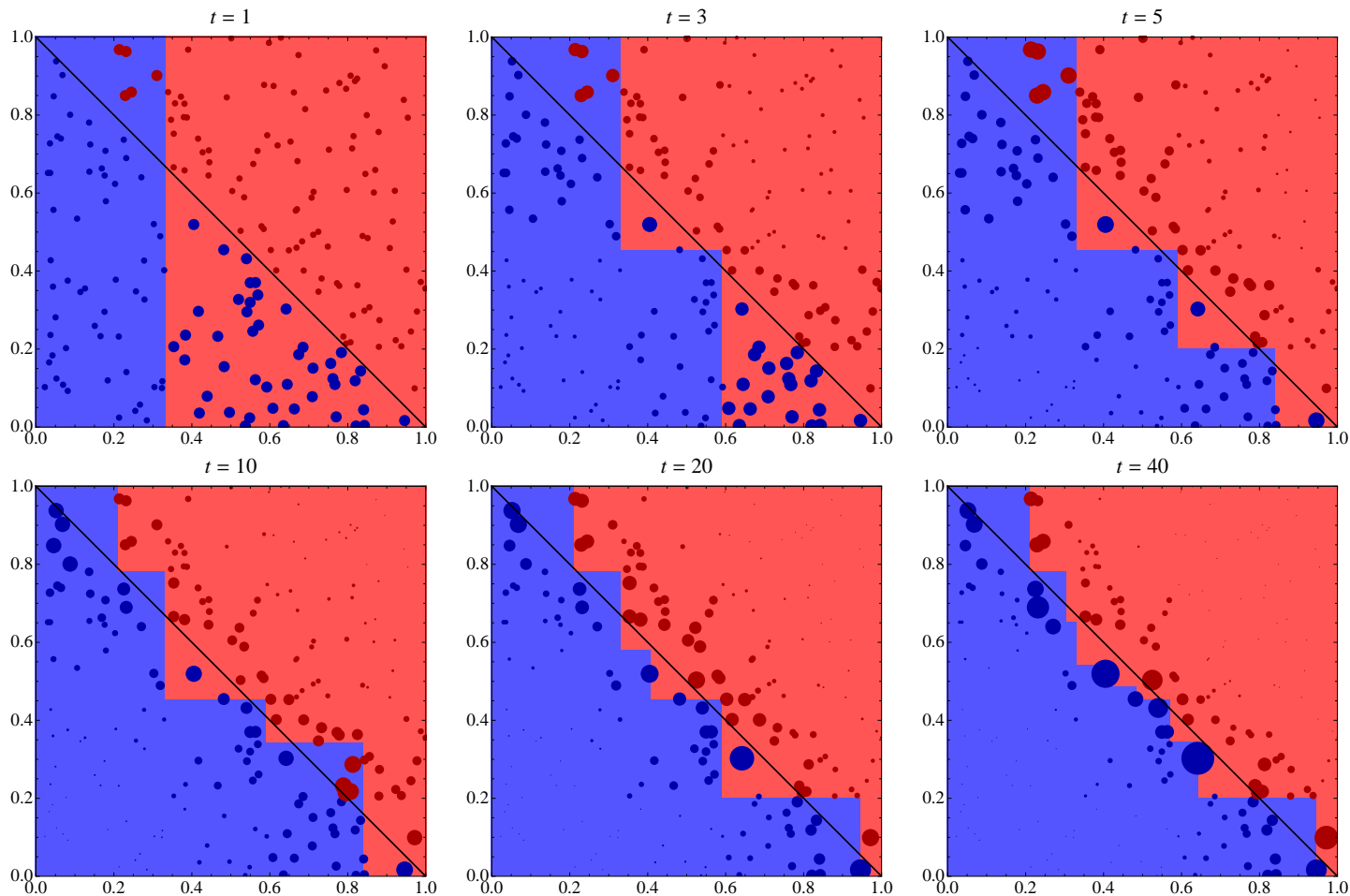
- 1 $\mathbf{w}^{(1)} \leftarrow (1/n, \dots, 1/n)$ \triangleright *poids initiaux*
- 2 **pour** $t \leftarrow 1$ **à** T
- 3 $h^{(t)} \leftarrow \text{BASE}(D_n, \mathbf{w})$ \triangleright *choix d'expert*
- 4 $\epsilon^{(t)} \leftarrow \sum_{h^{(t)}(\mathbf{x}_i) \neq y_i} w_i$ \triangleright *erreur pondérée*
- 5 $\alpha^{(t)} \leftarrow \frac{1}{2} \ln \left(\frac{1 - \epsilon^{(t)}}{\epsilon^{(t)}} \right)$ \triangleright *poids de $h^{(t)}$*
- 6 **pour** $i \leftarrow 1$ **à** n \triangleright *re-pondération des points*
- 7 **si** $h^{(t)}(\mathbf{x}_i) \neq y_i$ **alors**
- 8 $w_i^{(t+1)} \leftarrow \frac{w_i^{(t)}}{2\epsilon^{(t)}}$
- 9 **sinon**
- 10 $w_i^{(t+1)} \leftarrow \frac{w_i^{(t)}}{2(1 - \epsilon^{(t)})}$
- 11 **retourner** $f^{(T)}(\cdot) = \sum_{t=1}^T \alpha^{(t)} h^{(t)}(\cdot)$

AdaBoost

- Données:

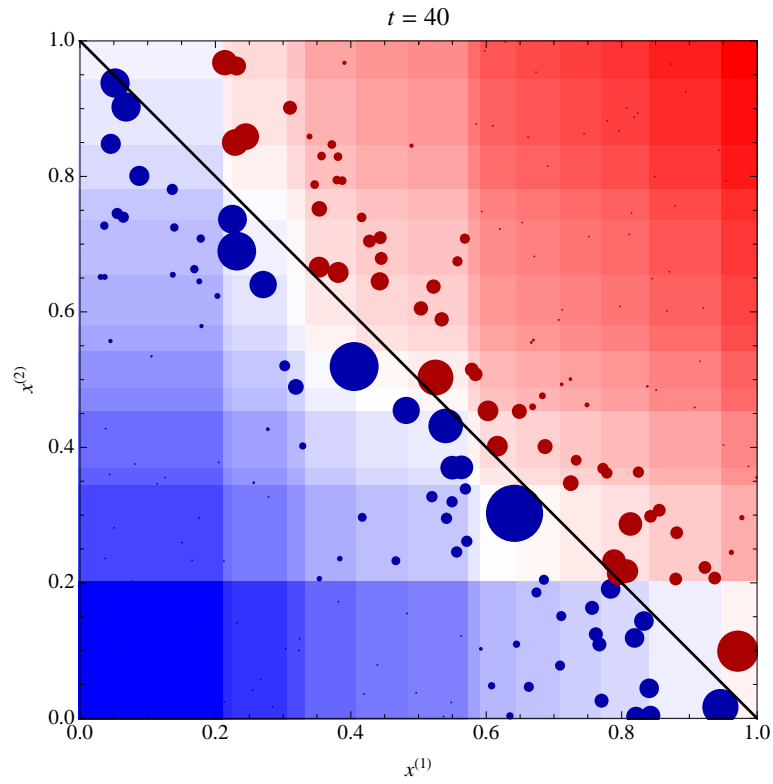


AdaBoost



AdaBoost

- Marges:



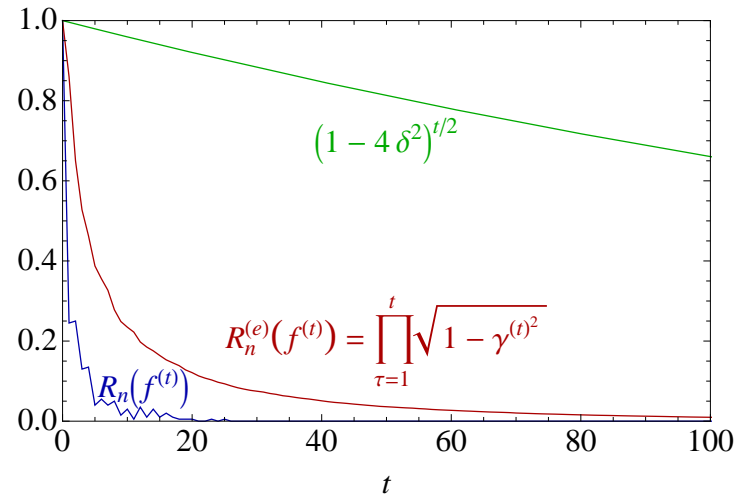
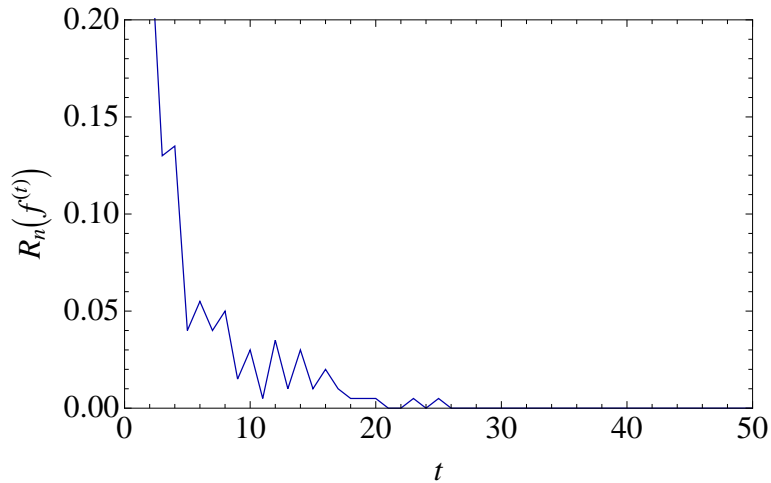
- Convergence

- soit $\varepsilon \leq \frac{1}{2} - \delta$ dans toutes les itérations : h est un peu meilleur qu'une décision aléatoire
- l'erreur d'entraînement

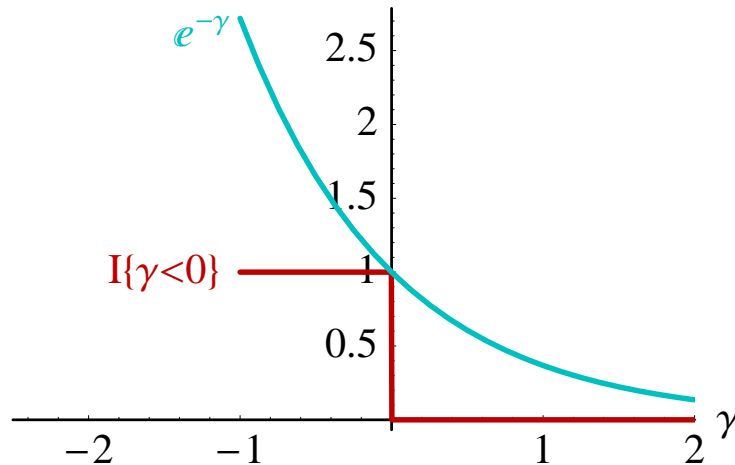
$$\widehat{R}(f) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{f(\mathbf{x}_i)y_i < 0\}$$

est 0 après $\left\lceil \frac{\ln n}{2\delta^2} \right\rceil + 1$ itérations

- Convergence



- Justification des formules :
 - minimiser un **coût exponentiel** sur la **marge** $\gamma = f(\mathbf{x})y$
 - **descente de gradient** dans l'**espace fonctionnel** des experts



- Extension 1 : experts ternaires : $h^{(j)} : \mathcal{X} \rightarrow \{-1, 0, 1\}$
 - experts peuvent s'abstenir
 - experts locaux
 - experts spécialistes
- Extension 2 : experts avec confiances : $h^{(j)} : \mathcal{X} \rightarrow [-1, 1]$
 - $\text{signe}(h^{(j)}(\mathbf{x}))$ est la classification de \mathbf{x} par l'expert
 - $|h^{(j)}(\mathbf{x})|$ est la confiance de l'expert

- Approches **générales**
 - one-against-all
 - one-against-one
- **Fonction discriminante** de multi-classe: $\mathbf{f} : \mathbf{x} \rightarrow \mathbb{R}^K$
 - notation : $\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_K(\mathbf{x}))$
- **Classifieur** de multi-classe:

$$g(\mathbf{x}) = \arg \max_{\ell \in \{1, \dots, K\}} f_{\ell}(\mathbf{x})$$

Multi-classe

- **Experts** binaires de multi-classe : $\mathbf{h} : \mathbf{x} \rightarrow \{-1, 1\}^K$

- Convertir un **expert binaire** ϕ à multi-classe:

$$\mathbf{h}(\mathbf{x}) = \mathbf{v} \cdot \phi(\mathbf{x})$$

- $\mathbf{v} \in \{-1, 1\}^K$ est le **vecteur de votes**

- **Étiquettes** de multi-classe : $\mathbf{y} = \{-1, \dots, -1, 1, -1, \dots, -1\}$

$$y_\ell = \begin{cases} 1 & \text{si } \ell \text{ est la vraie classe } \ell(\mathbf{x}), \\ -1 & \text{sinon} \end{cases}$$

- Erreurs de multi-classe :

$$\widehat{R}(\mathbf{f}) = \frac{1}{n} \sum_{i=1}^n \mathbb{I} \left\{ \ell(\mathbf{x}_i) \neq \widehat{\ell}(\mathbf{x}_i) \right\}$$

$$\widehat{R}(\mathbf{f}, \mathbf{w}^{(1)}) = \sum_{i=1}^n \sum_{\ell=1}^K w_{i,\ell}^{(1)} \mathbb{I} \{ f_{\ell}(\mathbf{x}_i) y_{\ell} < 0 \}$$

- par exemple

$$w_{i,\ell}^{(1)} = \begin{cases} \frac{1}{2n} & \text{si } \ell \text{ est la vraie classe (si } y_{i,\ell} = 1), \\ \frac{1}{2n(K-1)} & \text{sinon (si } y_{i,\ell} = -1). \end{cases}$$

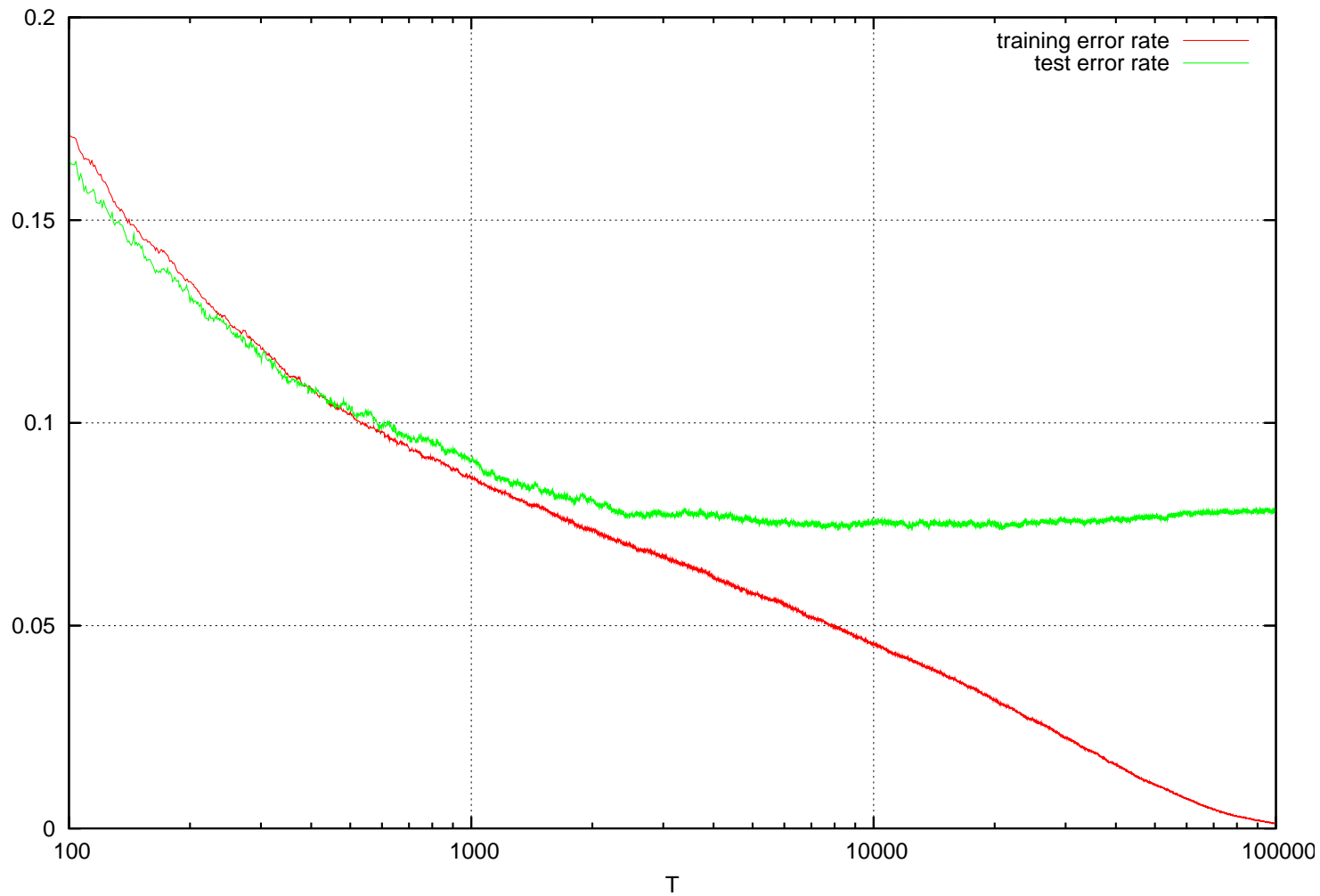
- AdaBoost

- pondération sur les points d'entraînement **et sur les classes** :

$$\{w_{i,\ell}\}_{i=1,\dots,n}^{\ell=1,\dots,K}$$

- $w_{i,\ell}$ est **initialisé** à $w_{i,\ell}^{(1)}$
- trouver le meilleur stump $\phi(\mathbf{x})$ et le meilleur vecteur de votes \mathbf{v} **reste efficace** ($O(nK)$)

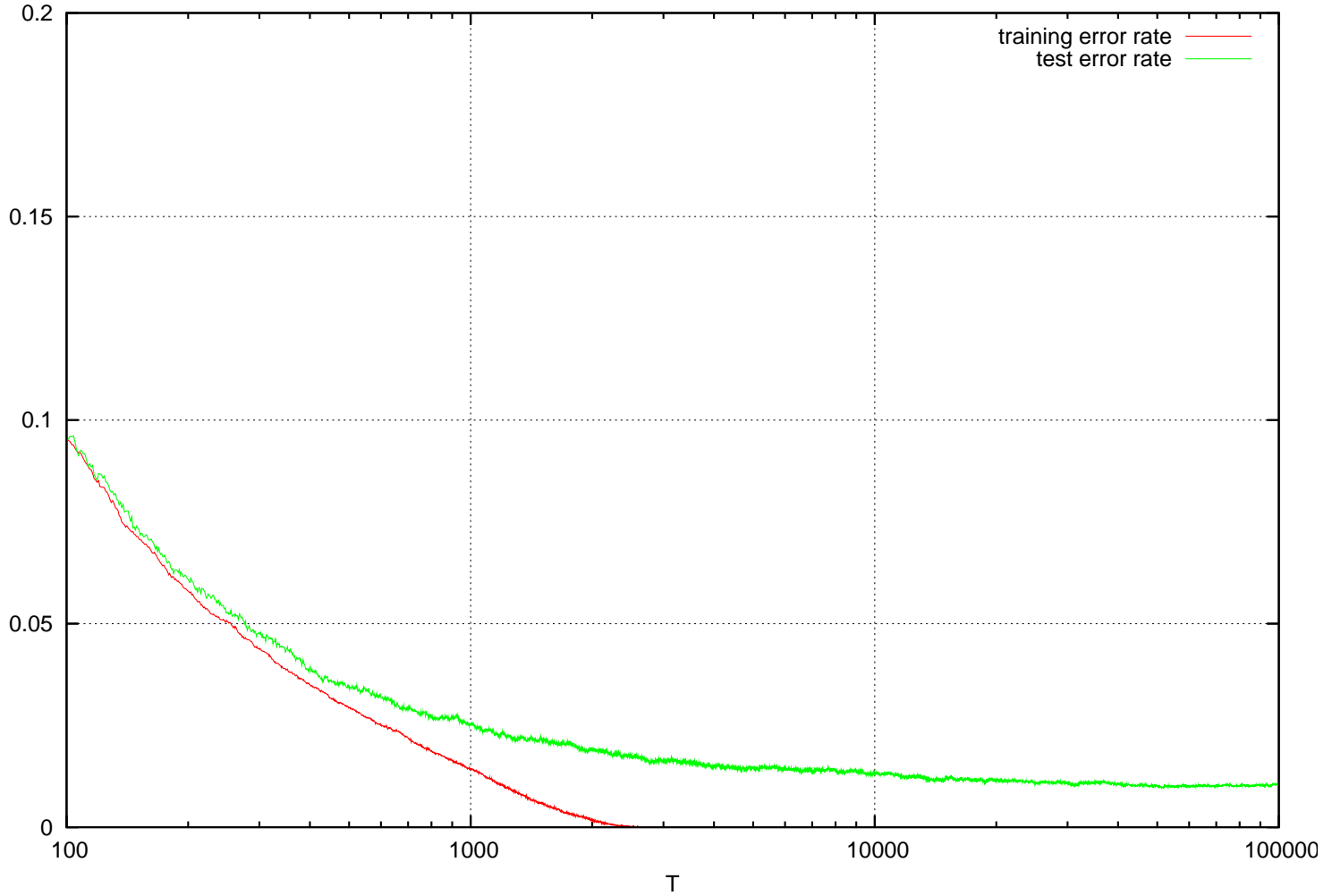
MNIST learning curves with stumps



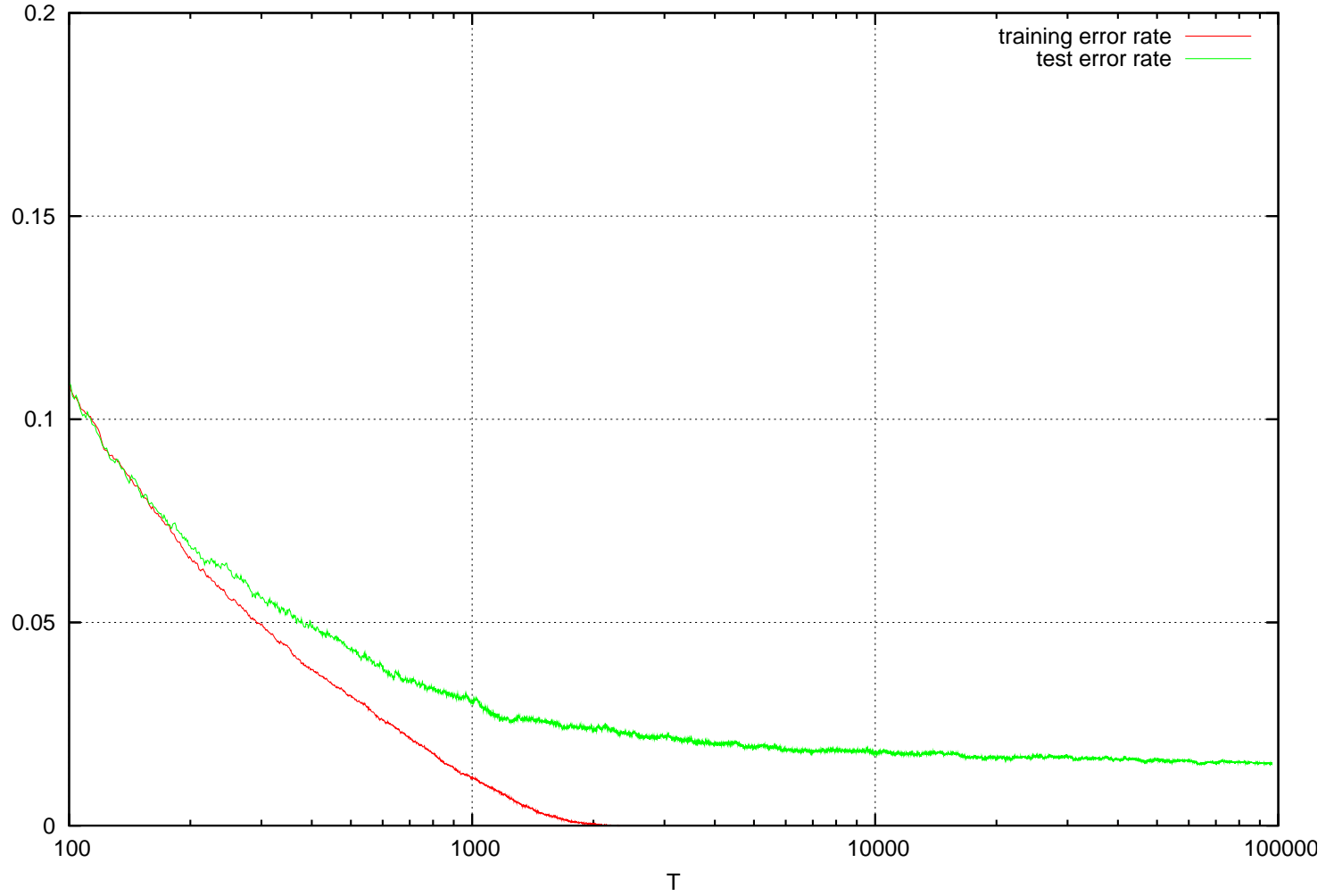
Meilleur experts

- Experts **génériques** plus puissants
 - arbres
 - produits
- Experts **spécifiques** aux applications
 - **filtres**, **extraction des traits**

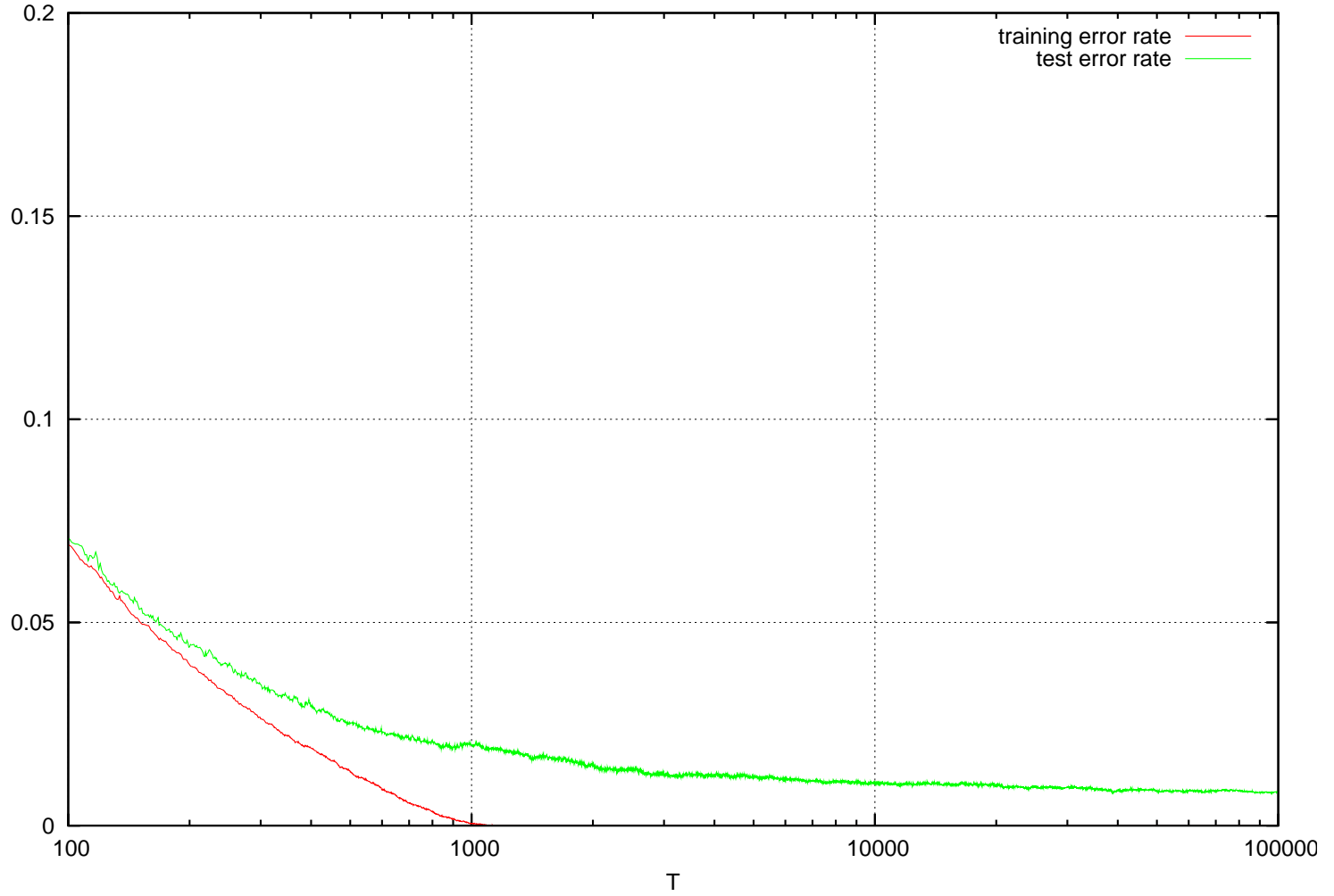
MNIST learning curves with stumps over Haar filters



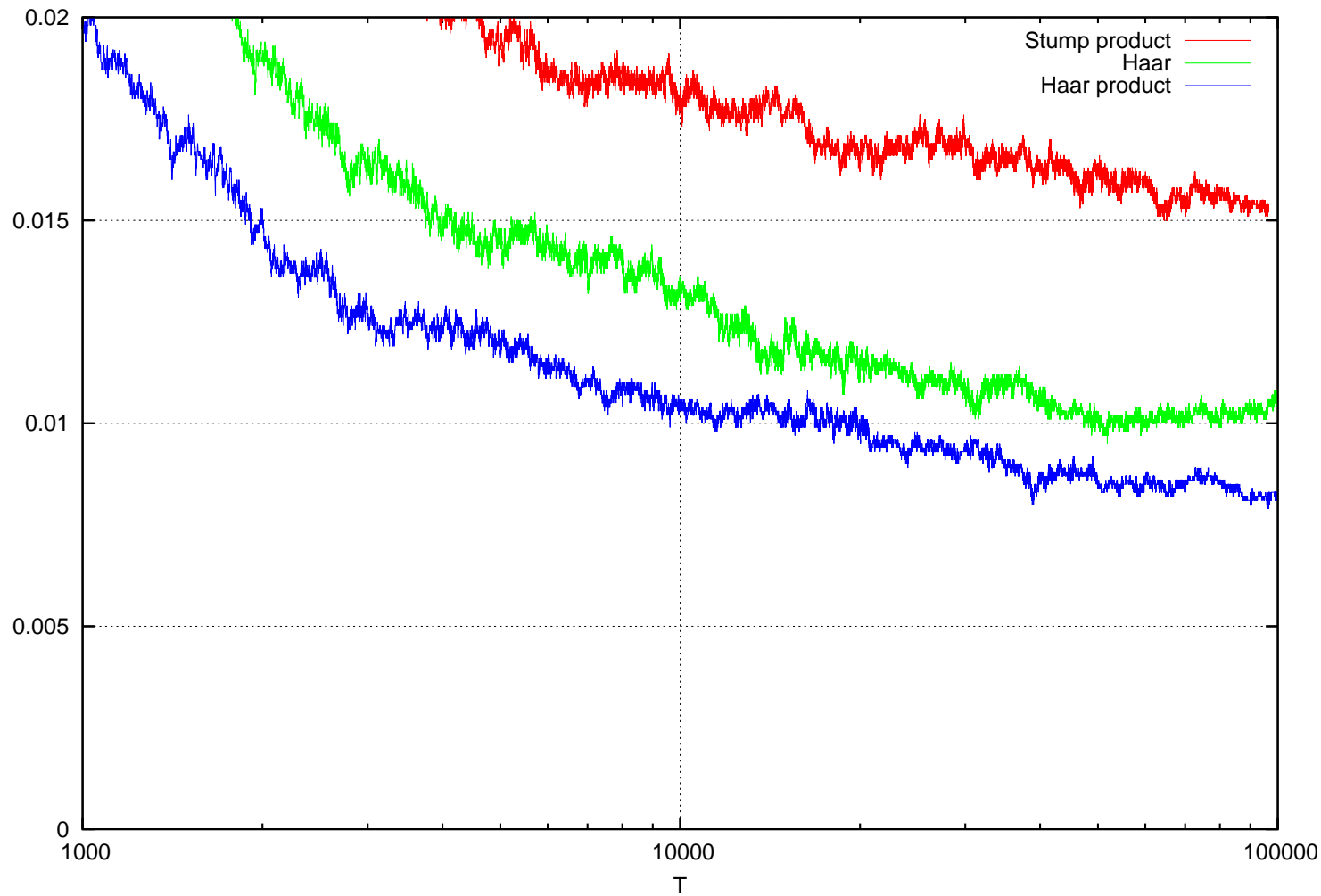
MNIST learning curves with products of stumps



MNIST learning curves with products of stumps over Haar filters



MNIST learning curves



- Avantages

- apprentissage **extrêmement simple**, pas de descente de gradient, pas d'optimisation numérique compliquée
- rapide
- **interprétation** intuitive : **vote pondéré des experts**
- le choix du **pool d'experts** "capsule" les connaissances à priori
- **aucune restriction** sur la forme des experts
- extension naturelle à la **classification multi-classe**

- Désavantages

- pas très bon sur des étiquettes **bruitées**
- l'extension à la **régression** est un peu forcée