

# Chapter 6 : Functional Data Fusion Application to Supervised Learning

John Klein

Lille1 Université - CRISyAL UMR CNRS 9189

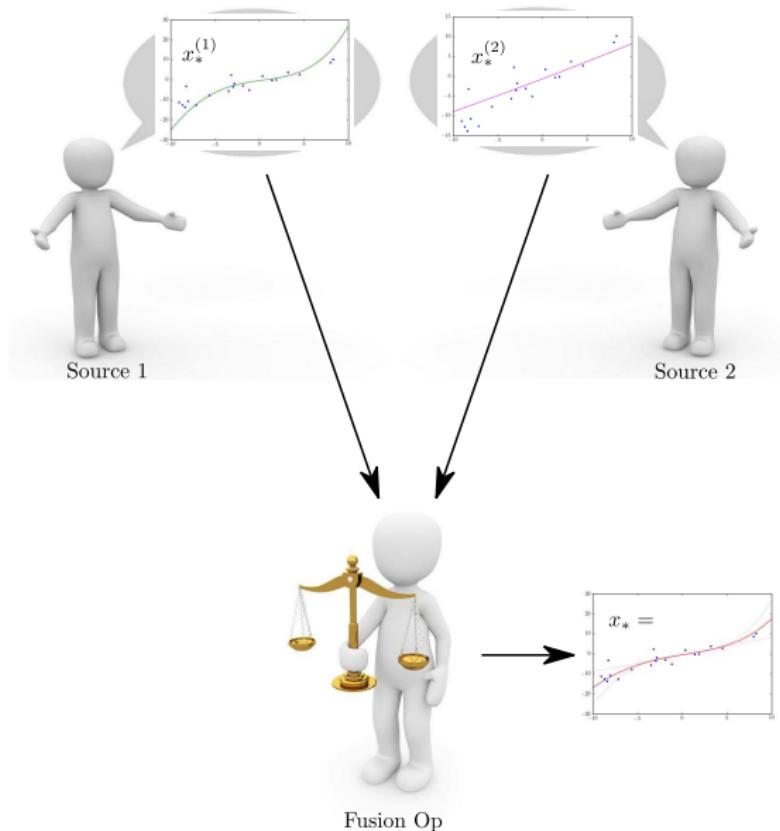


# Chapter organization

- 1 Mixture models
- 2 Ensemble methods
- 3 Bayesian methods
- 4 Statistical aggregation theory

- In this chapter, each datum delivered by a **source** is a **function**.
- We focus on **regression** or **classification** functions in a **supervised learning** paradigm.

## Probabilistic data fusion setting :



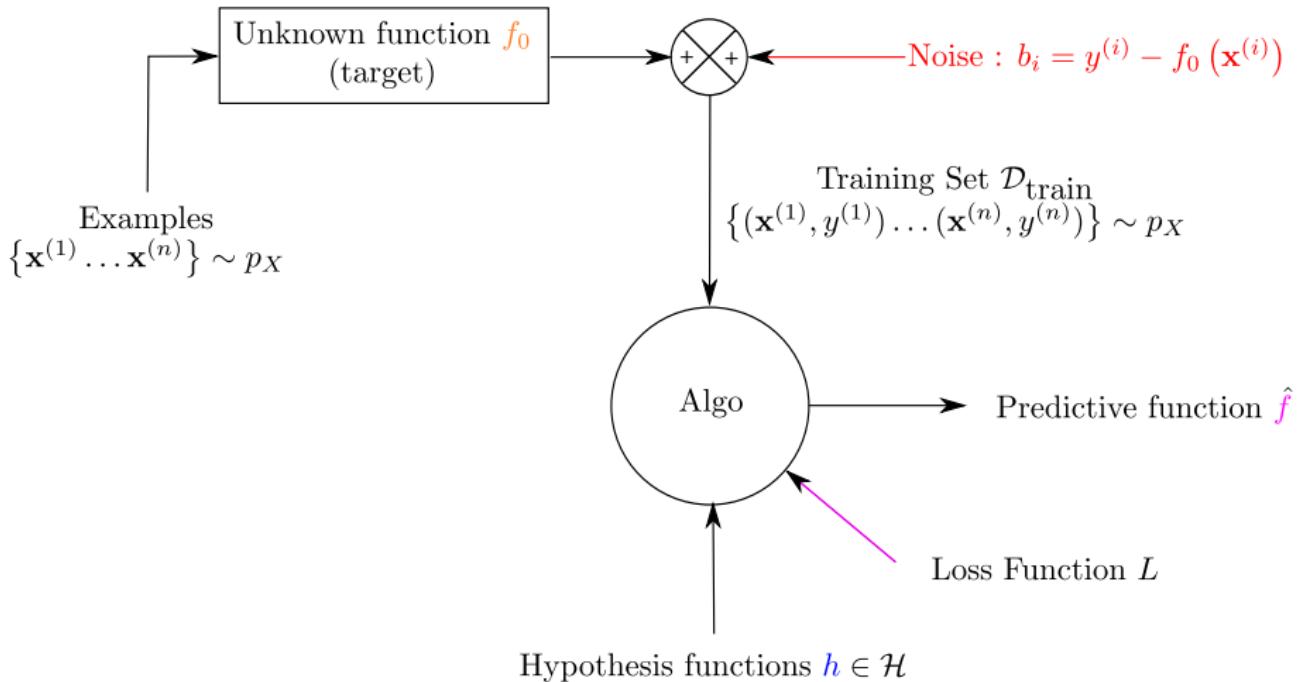
## Functional data fusion :

- Let  $\mathcal{L}_2$  denote the space of square integrable functions.
- Let us give a formal definition of functional data fusion problems :

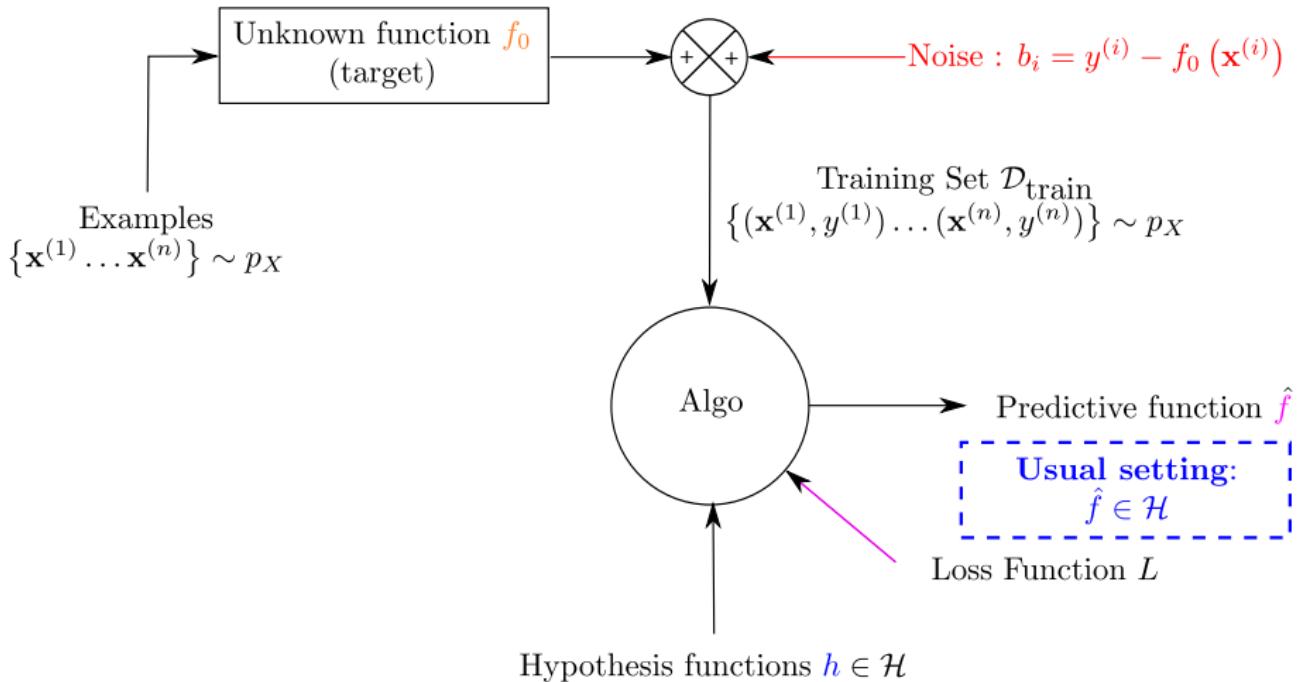
### Definition

**Functional data fusion** is a subclass of data fusion where advocacies live in  $\mathbb{X} = \mathcal{L}_2$ .

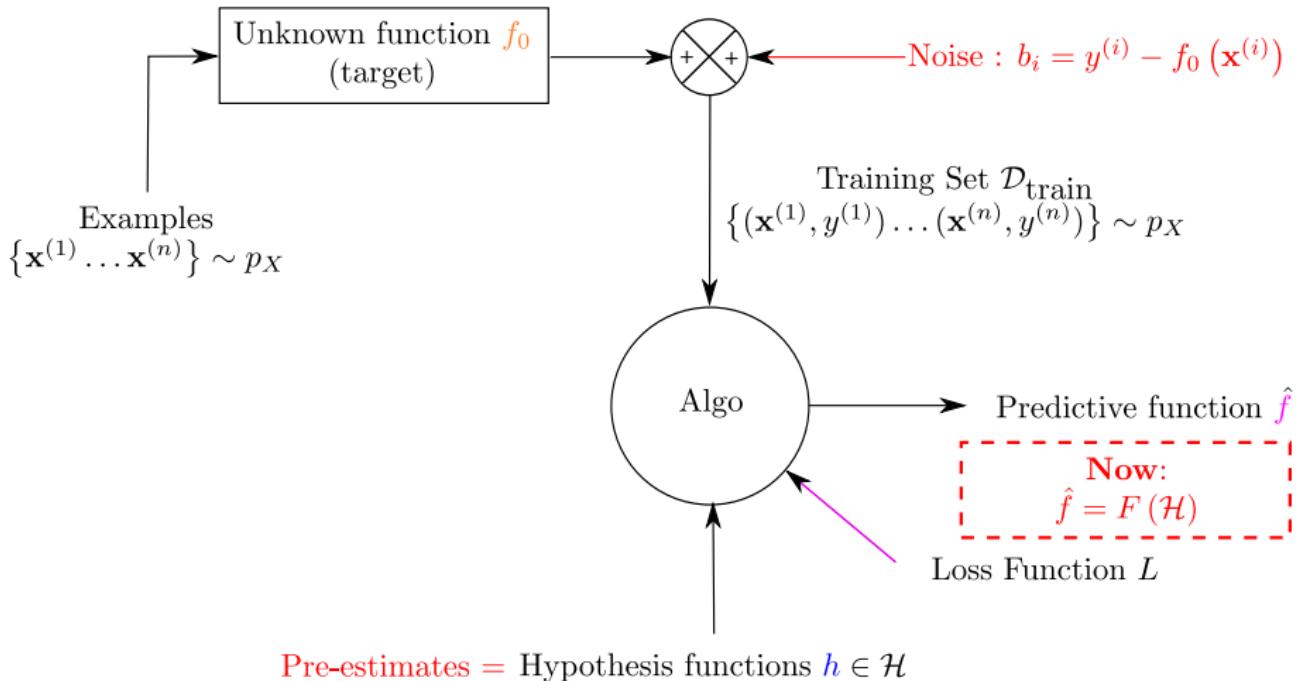
## Supervised Learning : the learning diagram



## Supervised Learning : the learning diagram



## Supervised Learning : the learning diagram



## New notations :

Before	Now
$x$ (solution)	$f_0$ (solution)
$x_*$ (aggregate)	$\hat{f}$ (aggregated estimate)
$x_*^{(i)}$ (advocacy)	$\hat{f}_i$ (pre-estimates)

## Supervised Learning : important notions

- Train Error (empirical risk) :

$$Err_{\text{train}} = \sum_{i=1}^n L\left(\hat{f}\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) = \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}_{\text{train}}} L\left(\hat{f}\left(\mathbf{x}^{(i)}\right), y^{(i)}\right).$$

- Test Error :

$$Err_{\text{test}} = \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}_{\text{test}}} L\left(\hat{f}\left(\mathbf{x}^{(i)}\right), y^{(i)}\right).$$

- Generalization Error :

$$Err_{\text{gen}} = \mathbb{E}_{(X, Y)} \left[ L\left(\hat{f}\left(\cdot\right), \cdot\right) \right].$$

## Supervised Learning : important notions

- Train Error (empirical risk) :

$$Err_{\text{train}} = \sum_{i=1}^n L\left(\hat{f}\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) = \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}_{\text{train}}} L\left(\hat{f}\left(\mathbf{x}^{(i)}\right), y^{(i)}\right).$$

- Test Error :

$$Err_{\text{test}} = \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}_{\text{test}}} L\left(\hat{f}\left(\mathbf{x}^{(i)}\right), y^{(i)}\right).$$

- Generalization Error :

$$Err_{\text{gen}} = \mathbb{E}_{(X, Y)} \left[ L\left(\hat{f}\left(\cdot\right), \cdot\right) \right].$$

## Supervised Learning : important notions

- Train Error (empirical risk) :

$$Err_{\text{train}} = \sum_{i=1}^n L\left(\hat{f}\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) = \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}_{\text{train}}} L\left(\hat{f}\left(\mathbf{x}^{(i)}\right), y^{(i)}\right).$$

- Test Error :

$$Err_{\text{test}} = \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}_{\text{test}}} L\left(\hat{f}\left(\mathbf{x}^{(i)}\right), y^{(i)}\right).$$

- Generalization Error :

$$Err_{\text{gen}} = \mathbb{E}_{(X, Y)} \left[ L\left(\hat{f}\left(\cdot\right), \cdot\right) \right].$$

## Supervised Learning : important notions

- $Err_{\text{train}}$  is a **biaised** estimator of  $Err_{\text{gen}}$ .
- $Err_{\text{test}}$  is a **unbiaised** estimator of  $Err_{\text{gen}}$ .
- If  $\hat{f}$  is learnt from  $\mathcal{D}_{\text{train}}$ , then the genuine **generalization** error is :

$$Err_{\text{gen}} = \mathbb{E}_{\mathcal{D}} \left[ \mathbb{E}_{(X,Y)} \left[ L \left( \hat{f}(\cdot), \cdot \right) \right] \right].$$

- This definition of  $Err_{\text{gen}}$  is also known as the **risk**.
- It features the ability for the **algorithm** to generalize correctly in this situation while the previous definition features **only** the ability of the **output predictor**  $\hat{f}$  to generalize.

## Supervised Learning : important notions

- $Err_{\text{train}}$  is a **biaised** estimator of  $Err_{\text{gen}}$ .
- $Err_{\text{test}}$  is a **unbiaised** estimator of  $Err_{\text{gen}}$ .
- If  $\hat{f}$  is learnt from  $\mathcal{D}_{\text{train}}$ , then the genuine **generalization** error is :

$$Err_{\text{gen}} = \mathbb{E}_{\mathcal{D}} \left[ \mathbb{E}_{(X,Y)} \left[ L \left( \hat{f}(\cdot), \cdot \right) \right] \right].$$

- This definition of  $Err_{\text{gen}}$  is also known as the **risk**.
- It features the ability for the **algorithm** to generalize correctly in this situation while the previous definition features **only** the ability of the **output predictor**  $\hat{f}$  to generalize.

## Supervised Learning : important notions

- $Err_{\text{train}}$  is a **biaised** estimator of  $Err_{\text{gen}}$ .
- $Err_{\text{test}}$  is a **unbiaised** estimator of  $Err_{\text{gen}}$ .
- If  $\hat{f}$  is learnt from  $\mathcal{D}_{\text{train}}$ , then the genuine **generalization** error is :

$$Err_{\text{gen}} = \mathbb{E}_{\mathcal{D}} \left[ \mathbb{E}_{(X,Y)} \left[ L \left( \hat{f} \left( . \right), . \right) \right] \right].$$

- This definition of  $Err_{\text{gen}}$  is also known as the **risk**.
- It features the ability for the **algorithm** to generalize correctly in this situation while the previous definition features **only** the ability of the **output predictor**  $\hat{f}$  to generalize.

## Supervised Learning : important notions

- $Err_{\text{train}}$  is a **biaised** estimator of  $Err_{\text{gen}}$ .
- $Err_{\text{test}}$  is a **unbiaised** estimator of  $Err_{\text{gen}}$ .
- If  $\hat{f}$  is learnt from  $\mathcal{D}_{\text{train}}$ , then the genuine **generalization** error is :

$$Err_{\text{gen}} = \mathbb{E}_{\mathcal{D}} \left[ \mathbb{E}_{(X,Y)} \left[ L \left( \hat{f} \left( . \right), . \right) \right] \right].$$

- This definition of  $Err_{\text{gen}}$  is also known as the **risk**.
- It features the ability for the **algorithm** to generalize correctly in this situation while the previous definition features **only** the ability of the **output predictor**  $\hat{f}$  to generalize.

## Supervised Learning : important notions

- $Err_{\text{train}}$  is a **biaised** estimator of  $Err_{\text{gen}}$ .
- $Err_{\text{test}}$  is a **unbiaised** estimator of  $Err_{\text{gen}}$ .
- If  $\hat{f}$  is learnt from  $\mathcal{D}_{\text{train}}$ , then the genuine **generalization** error is :

$$Err_{\text{gen}} = \mathbb{E}_{\mathcal{D}} \left[ \mathbb{E}_{(X,Y)} \left[ L \left( \hat{f}(\cdot), \cdot \right) \right] \right].$$

- This definition of  $Err_{\text{gen}}$  is also known as the **risk**.
- It features the ability for the **algorithm** to generalize correctly in this situation while the previous definition features **only** the ability of the **output predictor**  $\hat{f}$  to generalize.

# Chapter organization

- 1 Mixture models
- 2 Ensemble methods
- 3 Bayesian methods
- 4 Statistical aggregation theory

## Mixture models : definition

- A mixture model is a convex combination of posterior (predictive) distributions  $p_{Y|X;\theta_i}(y|x)$ :

$$p_{Y|X;\Theta}(y|x) = \sum_{i=1}^{N_s} \pi_i p_{Y|X;\theta_i}(y|x)$$

- The combined distribution depends on the sources parameters  $\theta_i$  and the combination parameters  $\pi_i$ :

$$\Theta = \{\theta_1, \dots, \theta_{N_s}, \pi_1, \dots, \pi_{N_s}\}.$$

- We have  $\sum_{i=1}^{N_s} \pi_i = 1$  and  $\pi_i \geq 0, \forall i$ .

## Mixture models : definition

- A mixture model is a convex combination of posterior (predictive) distributions  $p_{Y|X;\theta_i}(y|x)$ :

$$p_{Y|X;\Theta}(y|x) = \sum_{i=1}^{N_s} \pi_i p_{Y|X;\theta_i}(y|x)$$

- The combined distribution depends on the sources parameters  $\theta_i$  and the combination parameters  $\pi_i$ :

$$\Theta = \{\theta_1, \dots, \theta_{N_s}, \pi_1, \dots, \pi_{N_s}\}.$$

- We have  $\sum_{i=1}^{N_s} \pi_i = 1$  and  $\pi_i \geq 0, \forall i$ .

## Mixture models : definition

- A mixture model is a convex combination of posterior (predictive) distributions  $p_{Y|X;\theta_i}(y|x)$ :

$$p_{Y|X;\Theta}(y|x) = \sum_{i=1}^{N_s} \pi_i p_{Y|X;\theta_i}(y|x)$$

- The combined distribution depends on the sources parameters  $\theta_i$  and the combination parameters  $\pi_i$ :

$$\Theta = \{\theta_1, \dots, \theta_{N_s}, \pi_1, \dots, \pi_{N_s}\}.$$

- We have  $\sum_{i=1}^{N_s} \pi_i = 1$  and  $\pi_i \geq 0, \forall i$ .

## Mixture models : definition

- For a classification task, each pre-estimate is given by

$$\hat{f}_i(\mathbf{x}) = \arg \max_y p_{Y|\mathbf{X};\boldsymbol{\theta}_i}(y|\mathbf{x}).$$

- Consequently, the aggregated estimate is obtained by a weighted vote.
- All parameters can be jointly learnt from  $\mathcal{D}_{\text{train}}$  using EM.
- Non-probabilistic pre-estimates (SVM,  $k$ -NN) cannot be combined this way.

## Mixture models : definition

- For a classification task, each pre-estimate is given by

$$\hat{f}_i(\mathbf{x}) = \arg \max_y p_{Y|\mathbf{X};\theta_i}(y|\mathbf{x}).$$

- Consequently, the aggregated estimate is obtained by a weighted vote.
- All parameters can be jointly learnt from  $\mathcal{D}_{\text{train}}$  using EM.
- Non-probabilistic pre-estimates (SVM,  $k$ -NN) cannot be combined this way.

## Mixture models : definition

- For a classification task, each pre-estimate is given by

$$\hat{f}_i(\mathbf{x}) = \arg \max_y p_{Y|\mathbf{X};\boldsymbol{\theta}_i}(y|\mathbf{x}).$$

- Consequently, the aggregated estimate is obtained by a weighted vote.
- All parameters can be jointly learnt from  $\mathcal{D}_{\text{train}}$  using EM.
- Non-probabilistic pre-estimates (SVM,  $k$ -NN) cannot be combined this way.

## Mixture models : definition

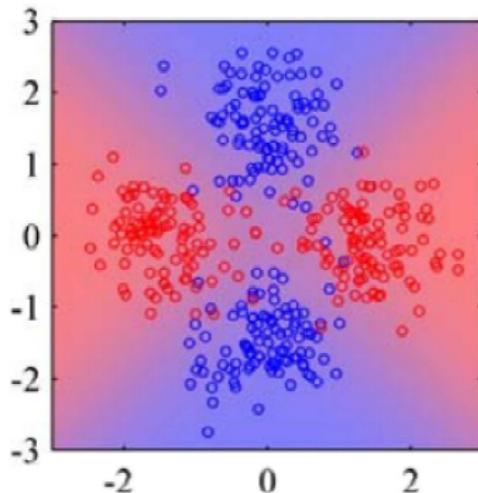
- For a classification task, each pre-estimate is given by

$$\hat{f}_i(\mathbf{x}) = \arg \max_y p_{Y|\mathbf{X};\boldsymbol{\theta}_i}(y|\mathbf{x}).$$

- Consequently, the aggregated estimate is obtained by a weighted vote.
- All parameters can be jointly learnt from  $\mathcal{D}_{\text{train}}$  using EM.
- Non-probabilistic pre-estimates (SVM,  $k$ -NN) cannot be combined this way.

## Mixture models : Example - Mixture of LogRegs

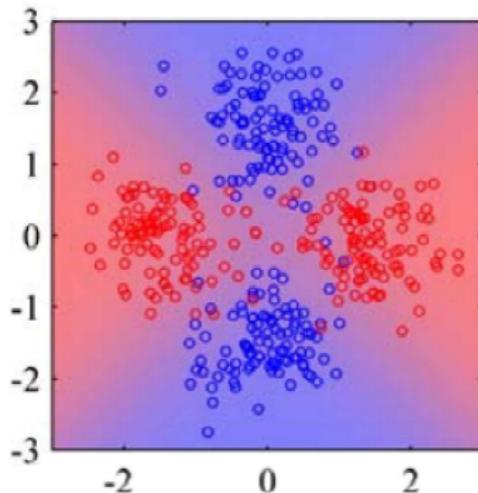
- Impossible to fit a linear model to this dataset :



- Let's try to fit a mixture of 2 logistic regressions.

## Mixture models : Example - Mixture of LogRegs

- Impossible to fit a linear model to this dataset :



- Let's try to fit a mixture of 2 logistic regressions.

## Mixture models : Example - Mixture of LogRegs

- Outputs  $y$  are **binary** variables  $\{0; 1\}$  and **inputs**  $x$  are **vectors** in  $\mathbb{R}^d$ .
- The posterior is :

$$\begin{aligned} p_{Y|X;\Theta}(y|x) &= \pi_1(1 - \hat{y}_1)^{1-y} \hat{y}_1^y + \pi_2(1 - \hat{y}_2)^{1-y} \hat{y}_2^y, \\ &= \begin{cases} \pi_1 \hat{y}_1 + \pi_2 \hat{y}_2 & \text{if } y = 1 \\ \pi_1(1 - \hat{y}_1) + \pi_2(1 - \hat{y}_2) & \text{if } y = 0 \end{cases} \end{aligned}$$

with  $\hat{y}_1 = \text{sgm}(\boldsymbol{\theta}_1^T \cdot \mathbf{x}_+)$  and  $\hat{y}_2 = \text{sgm}(\boldsymbol{\theta}_2^T \cdot \mathbf{x}_+)$  the logistic outputs and

$$\mathbf{x}_+^{(i)} = \begin{bmatrix} x_1^{(i)} \\ \vdots \\ x_d^{(i)} \\ 1 \end{bmatrix}$$

## Mixture models : Example - Mixture of LogRegs

- Outputs  $y$  are **binary** variables  $\{0; 1\}$  and **inputs**  $\mathbf{x}$  are **vectors** in  $\mathbb{R}^d$ .
- The posterior is :

$$\begin{aligned} p_{Y|X;\Theta}(y|\mathbf{x}) &= \pi_1 (1 - \hat{y}_1)^{1-y} \hat{y}_1^y + \pi_2 (1 - \hat{y}_2)^{1-y} \hat{y}_2^y, \\ &= \begin{cases} \pi_1 \hat{y}_1 + \pi_2 \hat{y}_2 & \text{if } y = 1 \\ \pi_1 (1 - \hat{y}_1) + \pi_2 (1 - \hat{y}_2) & \text{if } y = 0 \end{cases}, \end{aligned}$$

with  $\hat{y}_1 = \text{sgm}(\boldsymbol{\theta}_1^T \cdot \mathbf{x}_+)$  and  $\hat{y}_2 = \text{sgm}(\boldsymbol{\theta}_2^T \cdot \mathbf{x}_+)$  the logistic outputs and

$$\mathbf{x}_+^{(i)} = \begin{bmatrix} x_1^{(i)} \\ \vdots \\ x_d^{(i)} \\ 1 \end{bmatrix}$$

## Mixture models : Example - Mixture of LogRegs

- The likelihood is

$$\begin{aligned}
 \mathcal{L}(\Theta) &= p(\text{data}|\Theta), \\
 &= \prod_{i=1}^n p(\text{datum number } i|\Theta), \\
 &= \prod_{i=1}^n p\left(y^{(i)}|\mathbf{x}^{(i)}, \Theta\right), \\
 &= \prod_{i=1}^n \pi_1 \left(1 - \hat{y}_1^{(i)}\right)^{1-y^{(i)}} \left(\hat{y}_1^{(i)}\right)^{y^{(i)}} + \pi_2 \left(1 - \hat{y}_2^{(i)}\right)^{1-y^{(i)}} \left(\hat{y}_2^{(i)}\right)^{y^{(i)}}.
 \end{aligned}$$

## Mixture models : Example - Mixture of LogRegs

- Let us introduce **latent variables**  $z^{(i)}$  standing for the fact that example  $\mathbf{x}^{(i)}$  was generated by mixture **component n°1**.
- The **complete data**<sup>1</sup> likelihood is then :

$$\begin{aligned}\mathcal{L}_{\text{comp}}(\Theta) &= \prod_{i=1}^n p(y^{(i)}, z^{(i)} | \mathbf{x}^{(i)}, \Theta), \\ &= \prod_{i=1}^n \prod_{k=1}^2 \left( \pi_k \left(1 - \hat{y}_k^{(i)}\right)^{1-y^{(i)}} \left(\hat{y}_k^{(i)}\right)^{y^{(i)}} \right)^{\mathbb{I}_{k-1}(1-z^{(i)})}.\end{aligned}$$

- A **fake multiplication** appears because now each point is concerned with only one component of the mixture !

## Mixture models : Example - Mixture of LogRegs

- Let us introduce **latent variables**  $z^{(i)}$  standing for the fact that example  $x^{(i)}$  was generated by mixture **component n°1**.
- The **complete data**<sup>1</sup> likelihood is then :

$$\begin{aligned}\mathcal{L}_{\text{comp}}(\Theta) &= \prod_{i=1}^n p\left(y^{(i)}, z^{(i)} | \mathbf{x}^{(i)}, \Theta\right), \\ &= \prod_{i=1}^n \prod_{k=1}^2 \left( \pi_k \left(1 - \hat{y}_k^{(i)}\right)^{1-y^{(i)}} \left(\hat{y}_k^{(i)}\right)^{y^{(i)}} \right)^{\mathbb{I}_{k-1}(1-z^{(i)})}.\end{aligned}$$

- A **fake multiplication** appears because now each point is concerned with only one component of the mixture!

## Mixture models : Example - Mixture of LogRegs

- Let us introduce **latent variables**  $z^{(i)}$  standing for the fact that example  $\mathbf{x}^{(i)}$  was generated by mixture **component n°1**.
- The **complete data**<sup>1</sup> likelihood is then :

$$\begin{aligned}\mathcal{L}_{\text{comp}}(\Theta) &= \prod_{i=1}^n p\left(y^{(i)}, z^{(i)} | \mathbf{x}^{(i)}, \Theta\right), \\ &= \prod_{i=1}^n \prod_{k=1}^2 \left( \pi_k \left(1 - \hat{y}_k^{(i)}\right)^{1-y^{(i)}} \left(\hat{y}_k^{(i)}\right)^{y^{(i)}} \right)^{\mathbb{I}_{k-1}(1-z^{(i)})}.\end{aligned}$$

- A **fake multiplication** appears because now each point is concerned with only one component of the mixture!

## Mixture models : Example - Mixture of LogRegs

- E step : one can show that

$$\mathbb{E}_{z^{(1)}, \dots, z^{(n)}} [\log \mathcal{L}_{\text{comp}} (\Theta)] = \sum_{i=1}^n \sum_{k=1}^2 \gamma_k^{(i)} \left[ \log (\pi_k) + (1 - y^{(i)}) \log (1 - \hat{y}_k^{(i)}) \right.$$

$$\left. + y^{(i)} \log (\hat{y}_k^{(i)}) \right],$$

$$\gamma_k^{(i)} = \frac{\pi_k (1 - \hat{y}_k^{(i)})^{1-y^{(i)}} (\hat{y}_k^{(i)})^{y^{(i)}}}{\sum_{k'} \pi_{k'} (1 - \hat{y}_{k'}^{(i)})^{1-y^{(i)}} (\hat{y}_{k'}^{(i)})^{y^{(i)}}}.$$

- M step : parameters  $\theta_i$  need to be estimated using a gradient ascent (Newton's method) while mixing weights are given by :

$$\pi_k = \frac{1}{n} \sum_{i=1}^n \gamma_k^{(i)}.$$

## Mixture models : Example - Mixture of LogRegs

- E step : one can show that

$$\mathbb{E}_{z^{(1)}, \dots, z^{(n)}} [\log \mathcal{L}_{\text{comp}}(\Theta)] = \sum_{i=1}^n \sum_{k=1}^2 \gamma_k^{(i)} \left[ \log(\pi_k) + (1 - y^{(i)}) \log(1 - \hat{y}_k^{(i)}) \right.$$

$$\left. + y^{(i)} \log(\hat{y}_k^{(i)}) \right],$$

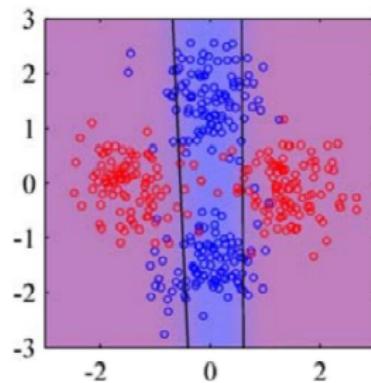
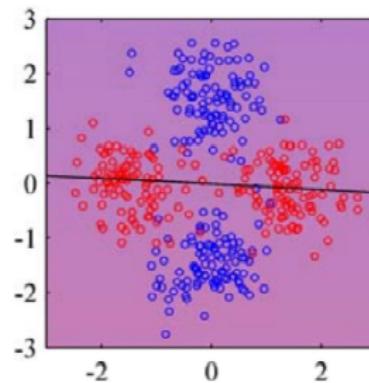
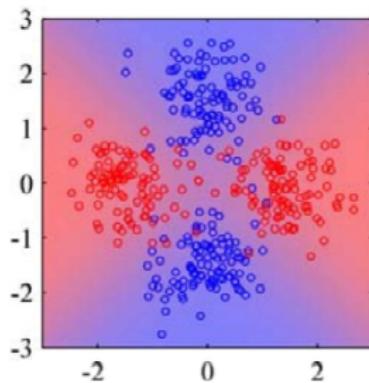
$$\gamma_k^{(i)} = \frac{\pi_k (1 - \hat{y}_k^{(i)})^{1-y^{(i)}} (\hat{y}_k^{(i)})^{y^{(i)}}}{\sum_{k'} \pi_{k'} (1 - \hat{y}_{k'}^{(i)})^{1-y^{(i)}} (\hat{y}_{k'}^{(i)})^{y^{(i)}}}.$$

- M step : parameters  $\theta_i$  need to be estimated using a gradient ascent (Newton's method) while mixing weights are given by :

$$\pi_k = \frac{1}{n} \sum_{i=1}^n \gamma_k^{(i)}.$$

## Mixture models : Example - Mixture of LogRegs

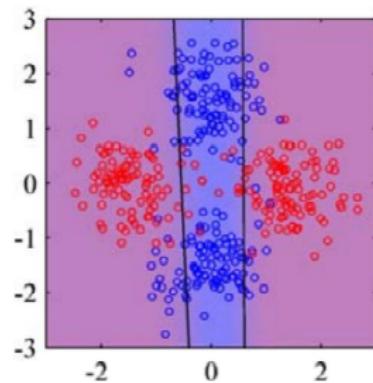
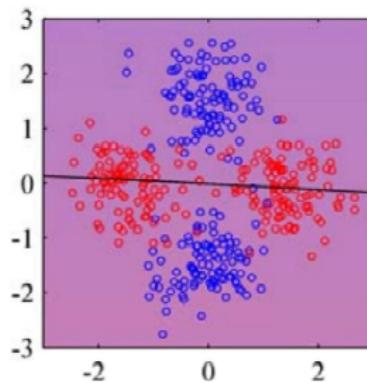
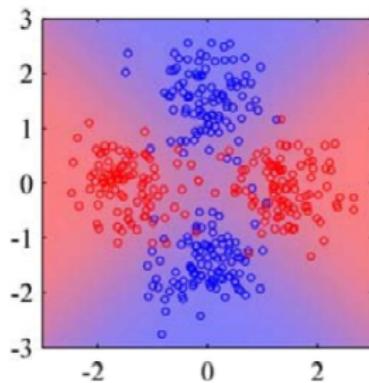
- The fit result is



- See [Bishop 14.5.2] for more details.

## Mixture models : Example - Mixture of LogRegs

- The fit result is



- See [Bishop 14.5.2] for more details.

## Mixture of Experts : definition

- The fact that pre-estimates are functions of  $\mathbf{x}$  is not exploited in the previous model.
- Mixtures of experts generalize this model by making mixing weights input-dependent :

$$\pi_k(\mathbf{x}) = \text{smax}(\mathbf{V}^T \cdot \mathbf{x}).$$

- In this context, mixing weights are called gating functions and each pre-estimate is called an expert.

## Mixture of Experts : definition

- The fact that pre-estimates are functions of  $\mathbf{x}$  is not exploited in the previous model.
- Mixtures of experts generalize this model by making mixing weights input-dependent :

$$\pi_k(\mathbf{x}) = \text{smax}(\mathbf{V}^T \cdot \mathbf{x}).$$

- In this context, mixing weights are called gating functions and each pre-estimate is called an expert.

## Mixture of Experts : definition

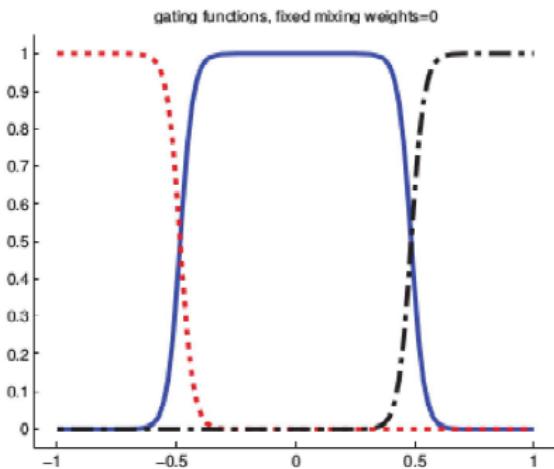
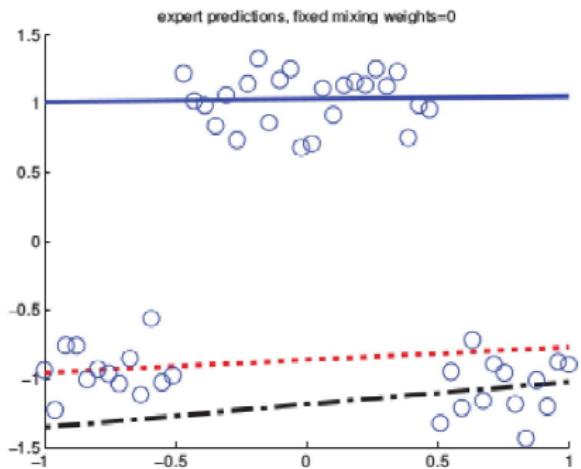
- The fact that pre-estimates are functions of  $\mathbf{x}$  is not exploited in the previous model.
- Mixtures of experts generalize this model by making mixing weights input-dependent :

$$\pi_k(\mathbf{x}) = \text{smax}(\mathbf{V}^T \cdot \mathbf{x}).$$

- In this context, mixing weights are called gating functions and each pre-estimate is called an expert.

## Mixture of Experts : Example - Mixture of Linear Regressors

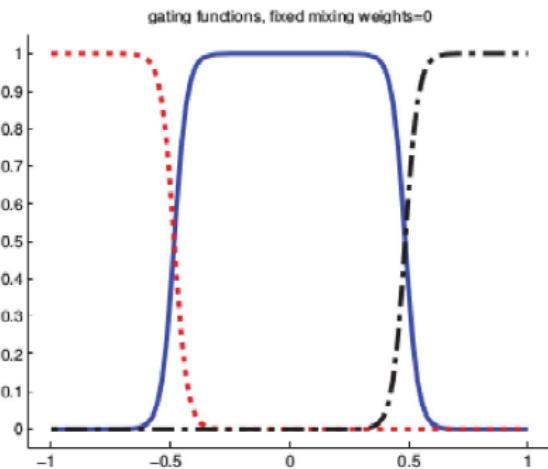
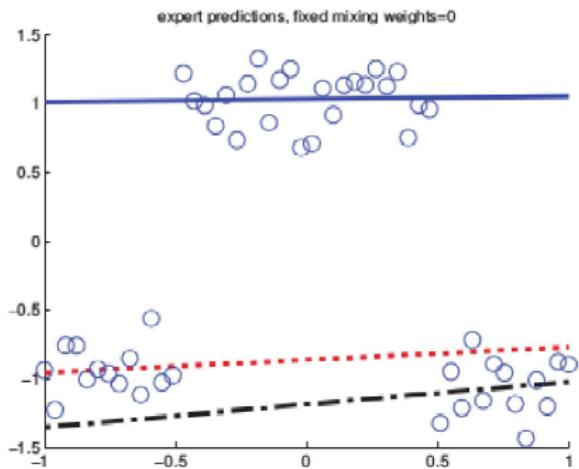
- Obviously, a linear regression is a good model for subsets of the following data :



- The right figure shows trained **gating functions** for each regressor.

## Mixture of Experts : Example - Mixture of Linear Regressors

- Obviously, a linear regression is a good model for subsets of the following data :



- The right figure shows trained **gating functions** for each regressor.

## Mixture of Experts : Example - Mixture of Linear Regressors

- E step : one can show that

$$\begin{aligned} \mathbb{E}_{z^{(1)}, \dots, z^{(n)}} [\log \mathcal{L}_{\text{comp}} (\Theta)] &= \sum_{i=1}^n \sum_{k=1}^2 \gamma_k^{(i)} \left[ \log \left( \pi_k^{(i)} \right) - \frac{(y^{(i)} - \theta_k^T \cdot \mathbf{x}^{(i)})^2}{2\sigma_k^2} \right], \\ \gamma_k^{(i)} &= \frac{\pi_k^{(i)} \times \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{(y^{(i)} - \theta_k^T \cdot \mathbf{x}^{(i)})^2}{2\sigma_k^2}}}{\sum_{k'} \pi_{k'}^{(i)} \times \frac{1}{\sqrt{2\pi}\sigma_{k'}} e^{-\frac{(y^{(i)} - \theta_{k'}^T \cdot \mathbf{x}^{(i)})^2}{2\sigma_{k'}^2}}}, \\ \pi_k^{(i)} &= \text{smax} \left( \mathbf{V}^T \cdot \mathbf{x}^{(i)} \right). \end{aligned}$$

- M step : there is a closed-form MLE solution for parameters  $\theta_k$  and  $\sigma_k$ .  $\mathbf{V}$  is estimated by gradient ascent (Newton's method).

## Mixture of Experts : Example - Mixture of Linear Regressors

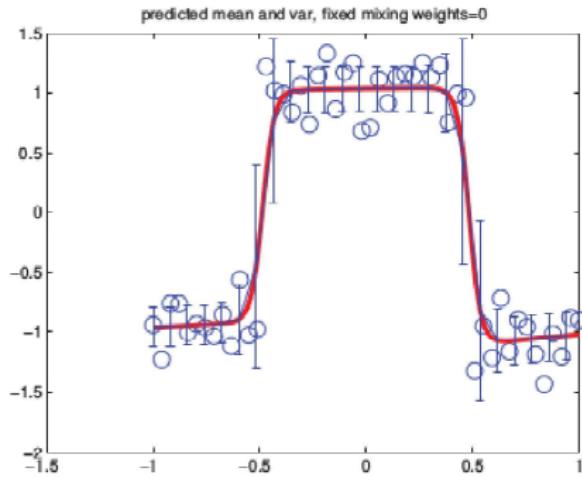
- E step : one can show that

$$\begin{aligned} \mathbb{E}_{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(n)}} [\log \mathcal{L}_{\text{comp}} (\Theta)] &= \sum_{i=1}^n \sum_{k=1}^2 \gamma_k^{(i)} \left[ \log \left( \pi_k^{(i)} \right) - \frac{(y^{(i)} - \boldsymbol{\theta}_k^T \cdot \mathbf{x}^{(i)})^2}{2\sigma_k^2} \right], \\ \gamma_k^{(i)} &= \frac{\pi_k^{(i)} \times \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{(y^{(i)} - \boldsymbol{\theta}_k^T \cdot \mathbf{x}^{(i)})^2}{2\sigma_k^2}}}{\sum_{k'} \pi_{k'}^{(i)} \times \frac{1}{\sqrt{2\pi}\sigma_{k'}} e^{-\frac{(y^{(i)} - \boldsymbol{\theta}_{k'}^T \cdot \mathbf{x}^{(i)})^2}{2\sigma_{k'}^2}}}, \\ \pi_k^{(i)} &= \text{smax} \left( \mathbf{V}^T \cdot \mathbf{x}^{(i)} \right). \end{aligned}$$

- M step : there is a closed-form MLE solution for parameters  $\boldsymbol{\theta}_k$  and  $\sigma_k$ .  $\mathbf{V}$  is estimated by gradient ascent (Newton's method).

## Mixture of Experts : Example - Mixture of Linear Regressors

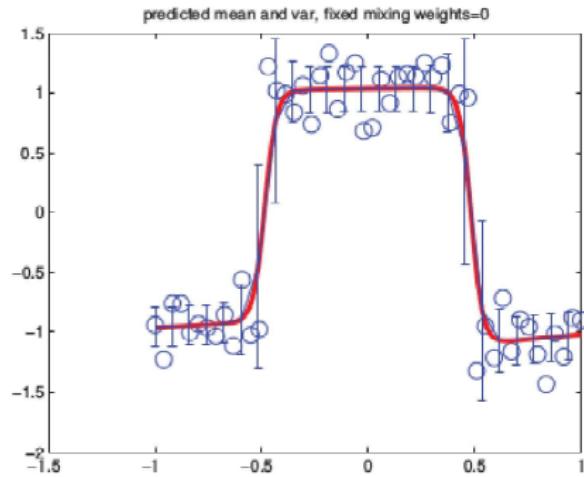
- The fit result is



- See [Murphy 2012 - 11.4.3] for more details.

## Mixture of Experts : Example - Mixture of Linear Regressors

- The fit result is



- See [Murphy 2012 - 11.4.3] for more details.

# Chapter organization

- 1 Mixture models
- 2 Ensemble methods
- 3 Bayesian methods
- 4 Statistical aggregation theory

## Ensemble methods : definition

- Ensemble methods are **homogeneous** classifier combination methods.
- This means that **pre-estimates** are trained using the **same algorithm** under different circumstances : different **hyperparameters**, different datasets, etc.

## Ensemble methods : definition

- Ensemble methods are **homogeneous** classifier combination methods.
- This means that **pre-estimates** are trained using the **same algorithm** under different circumstances : different **hyperparameters**, different datasets, etc.

## Ensemble methods : bagging

- ➊ Generate a dataset  $\mathcal{D}_{\text{boot}}$  from  $\mathcal{D}_{\text{train}}$  (unif. sampling with replacement).
- ➋ Train your algorithm on  $\mathcal{D}_{\text{boot}}$ .
- ➌ After  $m$  such trainings, combine using majority voting (classifiers) or average (regressors).

## Ensemble methods : bagging

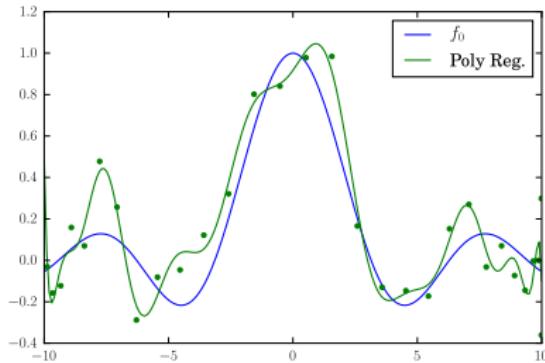
- ➊ Generate a dataset  $\mathcal{D}_{\text{boot}}$  from  $\mathcal{D}_{\text{train}}$  (unif. sampling with replacement).
- ➋ Train your algorithm on  $\mathcal{D}_{\text{boot}}$ .
- ➌ After  $m$  such trainings, combine using majority voting (classifiers) or average (regressors).

## Ensemble methods : bagging

- ➊ Generate a dataset  $\mathcal{D}_{\text{boot}}$  from  $\mathcal{D}_{\text{train}}$  (unif. sampling with replacement).
- ➋ Train your **algorithm** on  $\mathcal{D}_{\text{boot}}$ .
- ➌ After  $m$  such trainings, **combine** using majority voting (classifiers) or average (regressors).

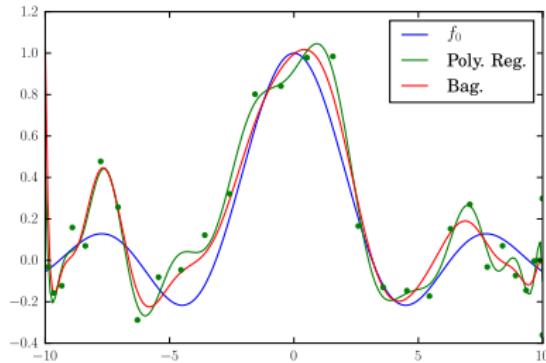
## Ensemble methods : bagging - Regression example

- Suppose one fits a 20 degree polynomial on this noisy data :



## Ensemble methods : bagging - Regression example

- Bagging will help a bit to mitigate overfitting



## Ensemble methods : bagging - Why should it work ?

- Let  $\epsilon_k$  denote the (signed) error function between the true function  $f_0$  and one of my pre-estimate  $f_k$  that was trained on a bootstrap sample :

$$\epsilon_k(x) = f_k(x) - f_0(x).$$

- For regression problems, the usual loss function is the quadratic one.
- The generalization error for the predictor  $f_k(x)$  is thus :

$$Err_{\text{gen}}^{(k)} = \mathbb{E}_X \left[ (f_k(x) - f_0(x))^2 \right] = \mathbb{E}_X \left[ \epsilon_k(x)^2 \right]$$

## Ensemble methods : bagging - Why should it work ?

- Let  $\epsilon_k$  denote the (signed) error function between the true function  $f_0$  and one of my pre-estimate  $f_k$  that was trained on a bootstrap sample :

$$\epsilon_k(x) = f_k(x) - f_0(x).$$

- For regression problems, the usual loss function is the quadratic one.
- The generalization error for the predictor  $f_k(x)$  is thus :

$$Err_{\text{gen}}^{(k)} = \mathbb{E}_X \left[ (f_k(x) - f_0(x))^2 \right] = \mathbb{E}_X \left[ \epsilon_k(x)^2 \right]$$

## Ensemble methods : bagging - Why should it work ?

- Let  $\epsilon_k$  denote the (signed) error function between the true function  $f_0$  and one of my pre-estimate  $f_k$  that was trained on a bootstrap sample :

$$\epsilon_k(x) = f_k(x) - f_0(x).$$

- For regression problems, the usual loss function is the quadratic one.
- The generalization error for the predictor  $f_k(x)$  is thus :

$$Err_{\text{gen}}^{(k)} = \mathbb{E}_X \left[ (f_k(x) - f_0(x))^2 \right] = \mathbb{E}_X \left[ \epsilon_k(x)^2 \right]$$

## Ensemble methods : bagging - Why should it work ?

- Now, the generalization error for the bagging ensemble writes :

$$\text{Err}_{\text{gen}}^{\text{ens}} = \mathbb{E}_X \left[ \left( \frac{1}{m} \sum_{k=1}^m f_k(x) - f_0(x) \right)^2 \right] = \mathbb{E}_X \left[ \left( \frac{1}{m} \sum_{k=1}^m \epsilon_k(x) \right)^2 \right]$$

- If  $\mathbb{E}_X [\epsilon_k] = 0$  ( $\forall k$ ) and  $\mathbb{E}_X [\epsilon_k \epsilon_{k'}] = 0$  ( $\forall k \neq k'$ ), then

$$\text{Err}_{\text{gen}}^{\text{ens}} = \frac{1}{m} \times \text{Err}_{\text{gen}}^{\text{ave}} \quad \text{with} \quad \text{Err}_{\text{gen}}^{\text{ave}} = \frac{1}{m} \sum_{k=1}^m \mathbb{E}_X [\epsilon_k^2] = \frac{1}{m} \sum_{k=1}^m \text{Err}_{\text{gen}}^{(k)}$$

- reduced error !

## Ensemble methods : bagging - Why should it work ?

- Now, the generalization error for the bagging ensemble writes :

$$\text{Err}_{\text{gen}}^{\text{ens}} = \mathbb{E}_X \left[ \left( \frac{1}{m} \sum_{k=1}^m f_k(x) - f_0(x) \right)^2 \right] = \mathbb{E}_X \left[ \left( \frac{1}{m} \sum_{k=1}^m \epsilon_k(x) \right)^2 \right]$$

- If  $\mathbb{E}_X [\epsilon_k] = 0$  ( $\forall k$ ) and  $\mathbb{E}_X [\epsilon_k \epsilon_{k'}] = 0$  ( $\forall k \neq k'$ ), then

$$\text{Err}_{\text{gen}}^{\text{ens}} = \frac{1}{m} \times \text{Err}_{\text{gen}}^{\text{ave}} \quad \text{with} \quad \text{Err}_{\text{gen}}^{\text{ave}} = \frac{1}{m} \sum_{k=1}^m \mathbb{E}_X [\epsilon_k^2] = \frac{1}{m} \sum_{k=1}^m \text{Err}_{\text{gen}}^{(k)}$$

- reduced error !

## Ensemble methods : bagging - Why should it work ?

- Now, the generalization error for the bagging ensemble writes :

$$\text{Err}_{\text{gen}}^{\text{ens}} = \mathbb{E}_X \left[ \left( \frac{1}{m} \sum_{k=1}^m f_k(x) - f_0(x) \right)^2 \right] = \mathbb{E}_X \left[ \left( \frac{1}{m} \sum_{k=1}^m \epsilon_k(x) \right)^2 \right]$$

- If  $\mathbb{E}_X [\epsilon_k] = 0$  ( $\forall k$ ) and  $\mathbb{E}_X [\epsilon_k \epsilon_{k'}] = 0$  ( $\forall k \neq k'$ ), then

$$\text{Err}_{\text{gen}}^{\text{ens}} = \frac{1}{m} \times \text{Err}_{\text{gen}}^{\text{ave}} \quad \text{with} \quad \text{Err}_{\text{gen}}^{\text{ave}} = \frac{1}{m} \sum_{k=1}^m \mathbb{E}_X [\epsilon_k^2] = \frac{1}{m} \sum_{k=1}^m \text{Err}_{\text{gen}}^{(k)}$$

- reduced error !

## Ensemble methods : boosting

- Use bagging in case of overfitting ( $\mathcal{H}$  is big  $\rightarrow$  high variance in trained estimates).
- Use boosting in case of underfitting ( $\mathcal{H}$  is small  $\rightarrow$  low variance in trained estimates but high bias).

## Ensemble methods : boosting

- Use bagging in case of overfitting ( $\mathcal{H}$  is big  $\rightarrow$  high variance in trained estimates).
- Use boosting in case of underfitting ( $\mathcal{H}$  is small  $\rightarrow$  low variance in trained estimates but high bias).

## Ensemble methods : boosting - procedure

Sequentially train weak classifiers as :

- ① Train classifier  $f_k$  on the weighted training set

$$\mathcal{D}_{\text{train}} \times \left\{ w_k^{(1)}, \dots, w_k^{(n)} \right\}.$$

- ② Evaluate the classification uncertainty on each data point and update weights accordingly.
- ③ Update classifier mixing coefficients  $\alpha_k$

Finally, return (in the binary classification case) :

$$\hat{f} = \text{sign} \left( \sum_{k=1}^m \alpha_k f_k \right).$$

## Ensemble methods : boosting - procedure

Sequentially train weak classifiers as :

- ① Train classifier  $f_k$  on the weighted training set

$$\mathcal{D}_{\text{train}} \times \left\{ w_k^{(1)}, \dots, w_k^{(n)} \right\}.$$

- ② Evaluate the classification uncertainty on each data point and update weights accordingly.
- ③ Update classifier mixing coefficients  $\alpha_k$

Finally, return (in the binary classification case) :

$$\hat{f} = \text{sign} \left( \sum_{k=1}^m \alpha_k f_k \right).$$

## Ensemble methods : boosting - procedure

Sequentially train weak classifiers as :

- ① Train classifier  $f_k$  on the weighted training set

$$\mathcal{D}_{\text{train}} \times \left\{ w_k^{(1)}, \dots, w_k^{(n)} \right\}.$$

- ② Evaluate the classification uncertainty on each data point and update weights accordingly.
- ③ Update classifier mixing coefficients  $\alpha_k$

Finally, return (in the binary classification case) :

$$\hat{f} = \text{sign} \left( \sum_{k=1}^m \alpha_k f_k \right).$$

## Ensemble methods : boosting - procedure

Sequentially train weak classifiers as :

- ① Train classifier  $f_k$  on the weighted training set

$$\mathcal{D}_{\text{train}} \times \left\{ w_k^{(1)}, \dots, w_k^{(n)} \right\}.$$

- ② Evaluate the classification uncertainty on each data point and update weights accordingly.
- ③ Update classifier mixing coefficients  $\alpha_k$

Finally, return (in the binary classification case) :

$$\hat{f} = \text{sign} \left( \sum_{k=1}^m \alpha_k f_k \right).$$

## Ensemble methods : boosting - Why should it work ?

- In classification, the standard loss function is the **0-1 loss**.
- Boosting aims at minimizing a weighted **0-1 loss** over the training set :

$$J_k = \sum_{i=1}^n w_k^{(i)} \left( 1 - \mathbb{I}_{y^{(i)}} \left( f_k \left( \mathbf{x}^{(i)} \right) \right) \right).$$

- The boosted ensemble, however, minimizes the **exponential loss** :

$$J_{\text{ens}} = \sum_{i=1}^n \exp \left( -y^{(i)} \frac{1}{2} \sum_{k=1}^m \alpha_k f_k \left( \mathbf{x}^{(i)} \right) \right).$$

- The first loss is derived from the ensemble one when the minimization focusses on  $f_k$  only.

## Ensemble methods : boosting - Why should it work ?

- In classification, the standard loss function is the **0-1 loss**.
- Boosting aims at minimizing a **weighted 0-1 loss** over the training set :

$$J_k = \sum_{i=1}^n w_k^{(i)} \left( 1 - \mathbb{I}_{y^{(i)}} \left( f_k \left( \mathbf{x}^{(i)} \right) \right) \right).$$

- The boosted ensemble, however, minimizes the **exponential loss** :

$$J_{\text{ens}} = \sum_{i=1}^n \exp \left( -y^{(i)} \frac{1}{2} \sum_{k=1}^m \alpha_k f_k \left( \mathbf{x}^{(i)} \right) \right).$$

- The first loss is derived from the ensemble one when the minimization focusses on  $f_k$  only.

## Ensemble methods : boosting - Why should it work ?

- In classification, the standard loss function is the **0-1 loss**.
- Boosting aims at minimizing a **weighted 0-1 loss** over the training set :

$$J_k = \sum_{i=1}^n w_k^{(i)} \left( 1 - \mathbb{I}_{y^{(i)}} \left( f_k \left( \mathbf{x}^{(i)} \right) \right) \right).$$

- The boosted ensemble, however, minimizes the **exponential loss** :

$$J_{\text{ens}} = \sum_{i=1}^n \exp \left( -y^{(i)} \frac{1}{2} \sum_{k=1}^m \alpha_k f_k \left( \mathbf{x}^{(i)} \right) \right).$$

- The first loss is derived from the ensemble one when the minimization focusses on  $f_k$  only.

## Ensemble methods : boosting - Why should it work ?

- In classification, the standard loss function is the **0-1 loss**.
- Boosting aims at minimizing a **weighted 0-1 loss** over the training set :

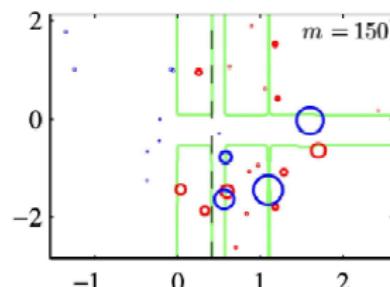
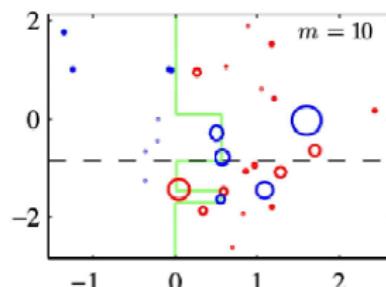
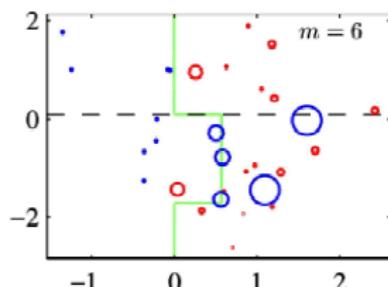
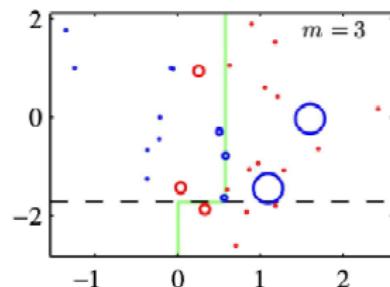
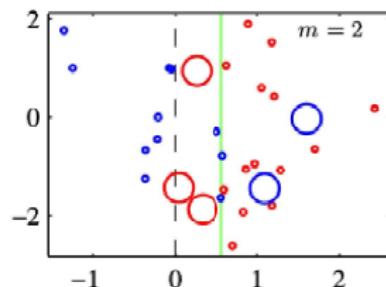
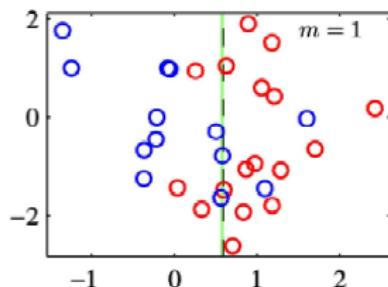
$$J_k = \sum_{i=1}^n w_k^{(i)} \left( 1 - \mathbb{I}_{y^{(i)}} \left( f_k \left( \mathbf{x}^{(i)} \right) \right) \right).$$

- The boosted ensemble, however, minimizes the **exponential loss** :

$$J_{\text{ens}} = \sum_{i=1}^n \exp \left( -y^{(i)} \frac{1}{2} \sum_{k=1}^m \alpha_k f_k \left( \mathbf{x}^{(i)} \right) \right).$$

- The first loss is derived from the ensemble one when the minimization focusses on  $f_k$  only.

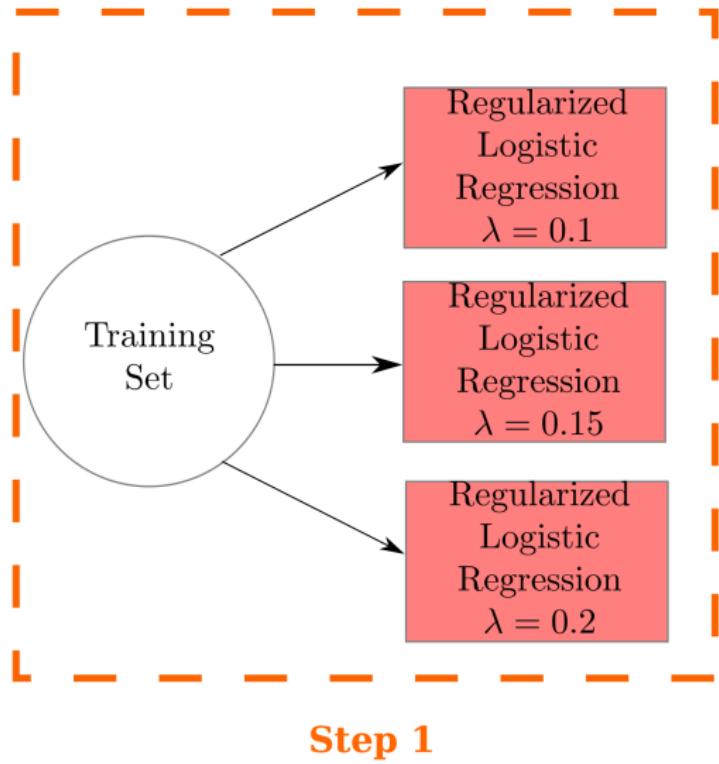
## Ensemble methods : boosting - Illustration



[Bishop 2006]

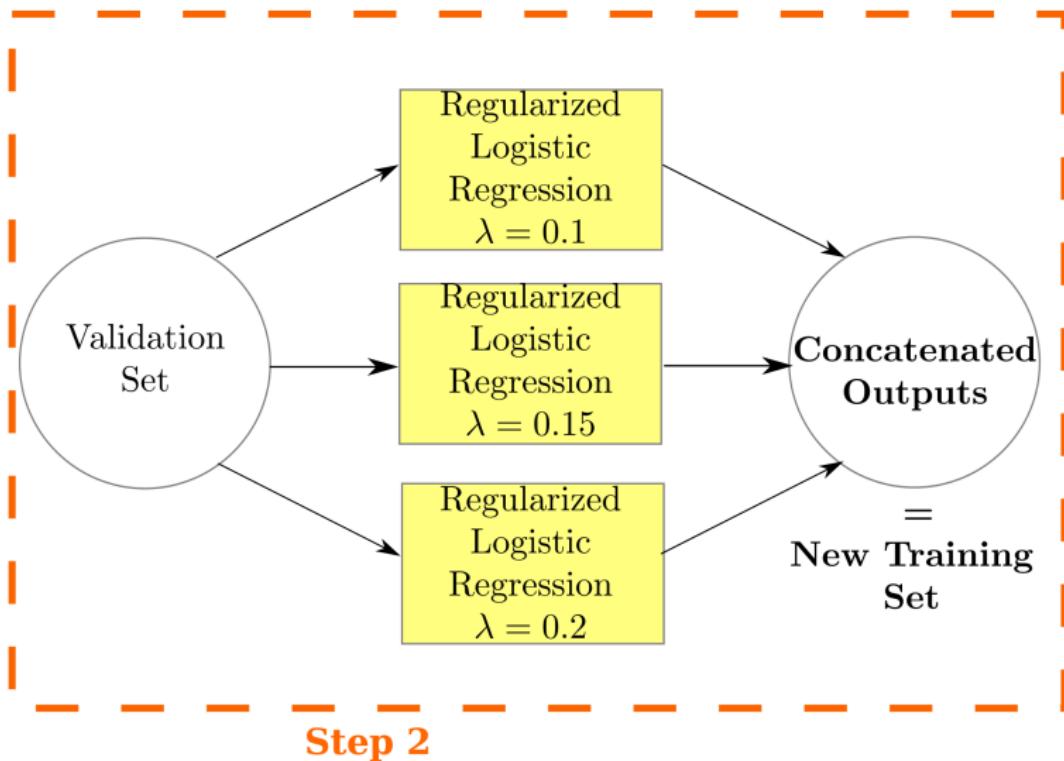
## Ensemble methods : stacking - Procedure

Train each member of the ensemble.



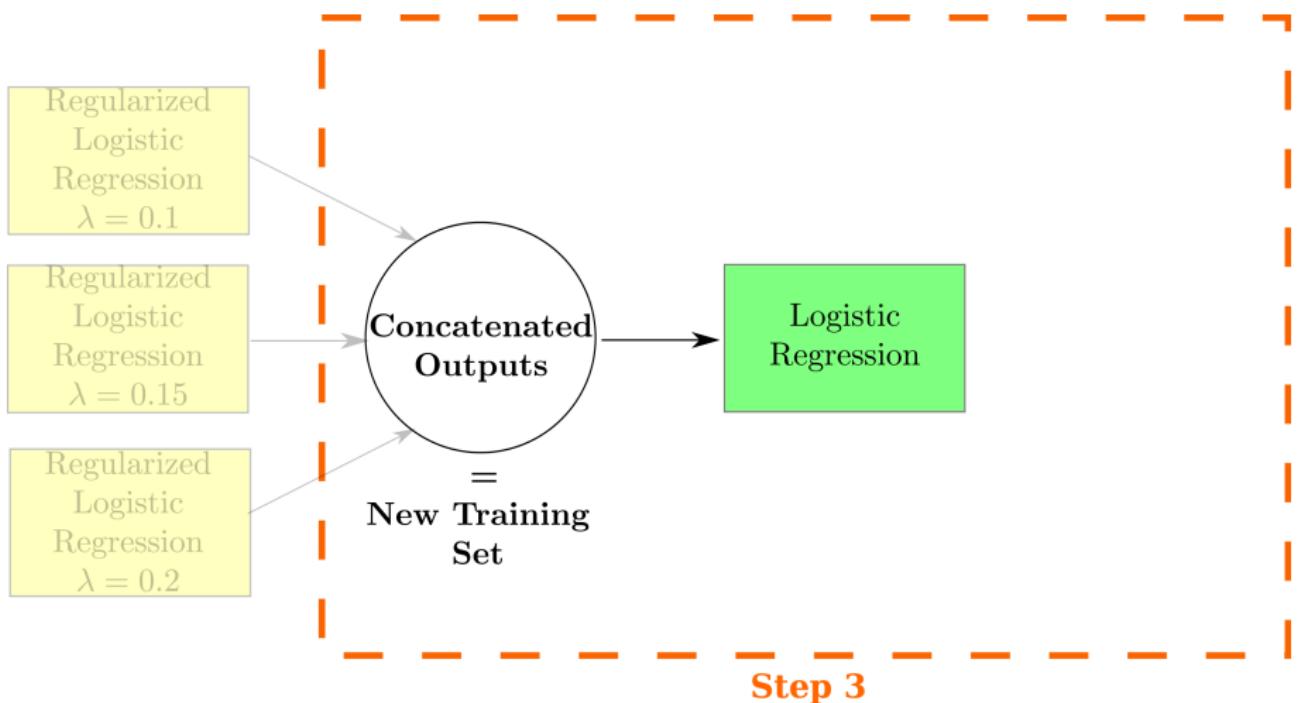
## Ensemble methods : stacking - Procedure

Generate a new training set for the fusion op.



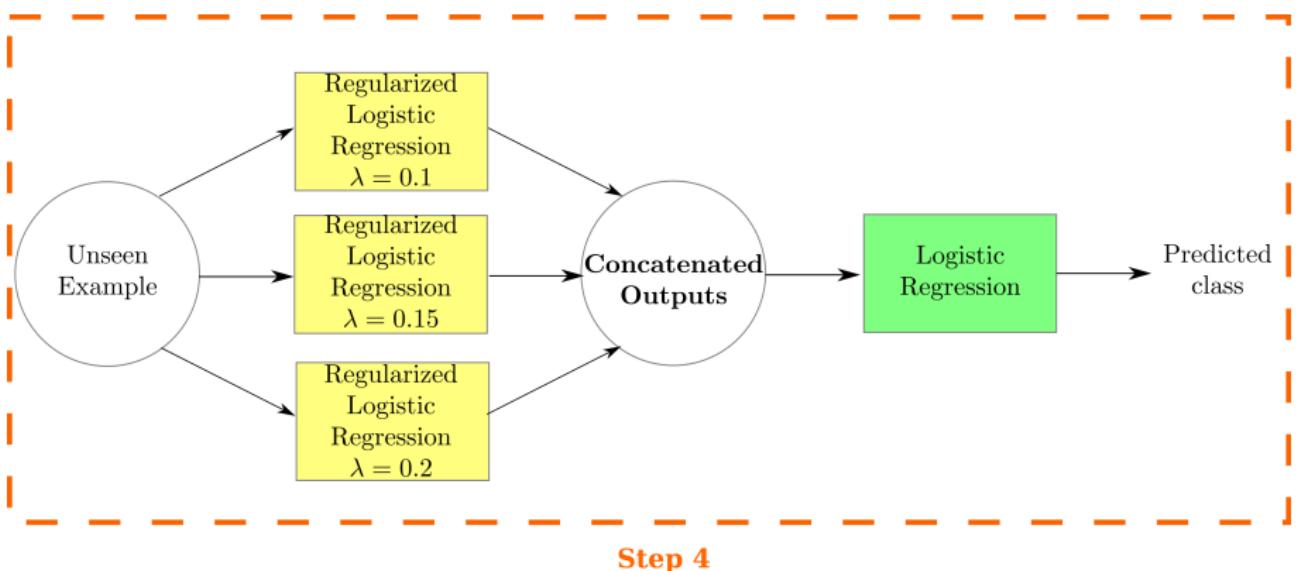
## Ensemble methods : stacking - Procedure

Train the fusion op.



## Ensemble methods : stacking - Procedure

Use it all.



## Ensemble methods : stacking - Comments

- Stacking can be used for heterogeneous classifiers too.
- Unlike mixture models, stacking need private data for training the fusion op.
- Solution : repeat the procedure in a cross validation (CV) fashion.
- Supposing one performs leave-one-out CV (LOOCV) and let  $f_k^{(-i)}$  denote the  $k^{\text{th}}$  training on the whole training set except data point  $x^{(i)}$ , then in step 3, one needs to minimize

$$J(\mathbf{w}) = \sum_{i=1}^n L \left( y^{(i)}, \operatorname{sgm} \left( \sum_{k=1}^m w_k f_k^{(-i)} \right) \right).$$

## Ensemble methods : stacking - Comments

- Stacking can be used for heterogeneous classifiers too.
- Unlike mixture models, stacking need private data for training the fusion op.
- Solution : repeat the procedure in a cross validation (CV) fashion.
- Supposing one performs leave-one-out CV (LOOCV) and let  $f_k^{(-i)}$  denote the  $k^{\text{th}}$  training on the whole training set except data point  $x^{(i)}$ , then in step 3, one needs to minimize

$$J(\mathbf{w}) = \sum_{i=1}^n L\left(y^{(i)}, \operatorname{sgm}\left(\sum_{k=1}^m w_k f_k^{(-i)}\right)\right).$$

## Ensemble methods : stacking - Comments

- Stacking can be used for heterogeneous classifiers too.
- Unlike mixture models, stacking need private data for training the fusion op.
- Solution : repeat the procedure in a cross validation (CV) fashion.
- Supposing one performs leave-one-out CV (LOOCV) and let  $f_k^{(-i)}$  denote the  $k^{\text{th}}$  training on the whole training set except data point  $x^{(i)}$ , then in step 3, one needs to minimize

$$J(\mathbf{w}) = \sum_{i=1}^n L \left( y^{(i)}, \operatorname{sgm} \left( \sum_{k=1}^m w_k f_k^{(-i)} \right) \right).$$

## Ensemble methods : stacking - Comments

- Stacking can be used for **heterogeneous** classifiers too.
- Unlike **mixture models**, stacking need **private data** for training the fusion op.
- **Solution** : repeat the procedure in a **cross validation (CV)** fashion.
- Supposing one performs leave-one-out CV (**LOOCV**) and let  $f_k^{(-i)}$  denote the  $k^{\text{th}}$  training on the whole training set **except** data point  $\mathbf{x}^{(i)}$ , then in step 3, one needs to minimize

$$J(\mathbf{w}) = \sum_{i=1}^n L\left(y^{(i)}, \operatorname{sgm}\left(\sum_{k=1}^m w_k f_k^{(-i)}\right)\right).$$

## Ensemble methods : Error-correcting output codes

- Suppose examples  $x$  need to be sorted into  $\ell > 2$  classes.
- One can assign a **code** encoded as a bit-vector  $\mathbf{b}$  for each class.
- For example :

Class	$b_1$	$b_2$	$b_3$	$b_4$	$b_5$	$b_6$
a	1	1	0	0	0	0
b	1	0	1	1	0	0
c	0	0	0	0	1	1

- This encryption uses more bits than necessary.
- Hamming distance : nbr. of unequal bits

## Ensemble methods : Error-correcting output codes

- Suppose examples  $x$  need to be sorted into  $\ell > 2$  classes.
- One can assign a **code** encoded as a **bit-vector  $b$**  for each class.
- For example :

Class	$b_1$	$b_2$	$b_3$	$b_4$	$b_5$	$b_6$
a	1	1	0	0	0	0
b	1	0	1	1	0	0
c	0	0	0	0	1	1

- This encryption uses more bits than necessary.
- Hamming distance : nbr. of unequal bits

## Ensemble methods : Error-correcting output codes

- Suppose examples  $x$  need to be sorted into  $\ell > 2$  classes.
- One can assign a **code** encoded as a **bit-vector**  $\mathbf{b}$  for each class.
- For **example** :

Class	$b_1$	$b_2$	$b_3$	$b_4$	$b_5$	$b_6$
a	1	1	0	0	0	0
b	1	0	1	1	0	0
c	0	0	0	0	1	1

- This encryption uses more bits than necessary.
- Hamming distance : nbr. of unequal bits

## Ensemble methods : Error-correcting output codes

- Suppose examples  $x$  need to be sorted into  $\ell > 2$  classes.
- One can assign a **code** encoded as a **bit-vector**  $\mathbf{b}$  for each class.
- For **example** :

Class	$b_1$	$b_2$	$b_3$	$b_4$	$b_5$	$b_6$
a	1	1	0	0	0	0
b	1	0	1	1	0	0
c	0	0	0	0	1	1

- This encryption uses more bits than necessary.
- Hamming distance : nbr. of unequal bits

## Ensemble methods : Error-correcting output codes

- Suppose examples  $x$  need to be sorted into  $\ell > 2$  classes.
- One can assign a **code** encoded as a **bit-vector**  $\mathbf{b}$  for each class.
- For **example** :

Class	$b_1$	$b_2$	$b_3$	$b_4$	$b_5$	$b_6$
a	1	1	0	0	0	0
b	1	0	1	1	0	0
c	0	0	0	0	1	1

- This encryption uses more bits than necessary.
- Hamming distance : nbr. of unequal bits

## Ensemble methods : Error-correcting output codes

- The minimal Hamming distance between each pair of codes is 3.
- Suppose we train  $m = 6$  different classifiers allowing to estimate probabilities  $p_{b_i}(x)$  (ex : LogReg).
- A class is decided for example  $x$  according to the rule :

$$\arg \min_y \sum_{i=1}^m |b_i(y) - \hat{p}_{b_i}(x)|$$

- This method is more robust to bit-flips, i.e. classification errors.

## Ensemble methods : Error-correcting output codes

- The minimal Hamming distance between each pair of codes is 3.
- Suppose we train  $m = 6$  different classifiers allowing to estimate probabilities  $p_{b_i}(\mathbf{x})$  (ex : LogReg).
- A class is decided for example  $\mathbf{x}$  according to the rule :

$$\arg \min_y \sum_{i=1}^m |b_i(y) - \hat{p}_{b_i}(\mathbf{x})|$$

- This method is more robust to bit-flips, i.e. classification errors.

## Ensemble methods : Error-correcting output codes

- The minimal Hamming distance between each pair of codes is 3.
- Suppose we train  $m = 6$  different classifiers allowing to estimate probabilities  $p_{b_i}(\mathbf{x})$  (ex : LogReg).
- A class is decided for example  $\mathbf{x}$  according to the rule :

$$\arg \min_y \sum_{i=1}^m |b_i(y) - \hat{p}_{b_i}(\mathbf{x})|$$

- This method is more robust to bit-flips, i.e. classification errors.

## Ensemble methods : Error-correcting output codes

- The minimal Hamming distance between each pair of codes is 3.
- Suppose we train  $m = 6$  different classifiers allowing to estimate probabilities  $p_{b_i}(\mathbf{x})$  (ex : LogReg).
- A class is decided for example  $\mathbf{x}$  according to the rule :

$$\arg \min_y \sum_{i=1}^m |b_i(y) - \hat{p}_{b_i}(\mathbf{x})|$$

- This method is more robust to bit-flips, i.e. classification errors.

# Chapter organization

- 1 Mixture models
- 2 Ensemble methods
- 3 Bayesian methods
- 4 Statistical aggregation theory

## Bayesian Learning :

- A catch sentence for this chapter could be :  
« Why use **only one** classifier when I can use **many** ? »
- With **Bayesian learning**, this would become :  
« Why use **only one** classifier when I can use **infinitely many** ? »
- Let us see under which circumstances such a result can be achieved.

## Bayesian Learning :

- A catch sentence for this chapter could be :  
« Why use **only one** classifier when I can use **many** ? »
- With **Bayesian learning**, this would become :  
« Why use **only one** classifier when I can use **infinitely many** ? »
- Let us see under which circumstances such a result can be achieved.

## Bayesian Learning :

- A catch sentence for this chapter could be :  
« Why use **only one** classifier when I can use **many** ? »
- With **Bayesian learning**, this would become :  
« Why use **only one** classifier when I can use **infinitely many** ? »
- Let us see under which circumstances such a result can be achieved.

## Bayesian Learning :- starting point

- Most of learning algorithms translate into an optimization problem of the following kind :

$$\arg \min_{\theta} \text{DataFit}(\theta) + \text{Regularizer}(\theta).$$

- In this setting, each  $f \in \mathcal{H}$  is in bijective correspondance with a given  $\theta \in \Theta$ .
- Almost all such algorithms have an equivalent probabilistic formulation :

$$\arg \max_{\theta} \text{Likelihood}(\theta) \times \text{Prior}(\theta).$$

## Bayesian Learning :- starting point

- Most of learning algorithms translate into an optimization problem of the following kind :

$$\arg \min_{\theta} \text{DataFit}(\theta) + \text{Regularizer}(\theta).$$

- In this setting, each  $f \in \mathcal{H}$  is in bijective correspondance with a given  $\theta \in \Theta$ .
- Almost all such algorithms have an equivalent probabilistic formulation :

$$\arg \max_{\theta} \text{Likelihood}(\theta) \times \text{Prior}(\theta).$$

## Bayesian Learning :- starting point

- Most of learning algorithms translate into an optimization problem of the following kind :

$$\arg \min_{\theta} \text{DataFit}(\theta) + \text{Regularizer}(\theta).$$

- In this setting, each  $f \in \mathcal{H}$  is in bijective correspondance with a given  $\theta \in \Theta$ .
- Almost all such algorithms have an equivalent probabilistic formulation :

$$\arg \max_{\theta} \text{Likelihood}(\theta) \times \text{Prior}(\theta).$$

## Bayesian Learning :- Linear regression example

- Suppose we are trying to predict the **selling price**  $y$  of a house.
- For each house, we collected data like **surface**, **previous buying price**, **GPS coordinates**, etc.
- These **features** are concatenated into a vector  $\mathbf{x}$ ;
- We need to **learn** the function  $f_0$  mapping vectors  $\mathbf{x}$  to  $y$ .
- We believe a **linear combination** of the features shoud be a relevant model :

$$y = \boldsymbol{\theta}^T \cdot \mathbf{x}.$$

## Bayesian Learning :- Linear regression example

- Suppose we are trying to predict the **selling price**  $y$  of a house.
- For each house, we collected data like **surface**, **previous buying price**, **GPS coordinates**, etc.
- These **features** are concatenated into a vector  $\mathbf{x}$ ;
- We need to **learn** the function  $f_0$  mapping vectors  $\mathbf{x}$  to  $y$ .
- We believe a **linear combination** of the features shoud be a relevant model :

$$y = \boldsymbol{\theta}^T \cdot \mathbf{x}.$$

## Bayesian Learning :- Linear regression example

- Suppose we are trying to predict the **selling price**  $y$  of a house.
- For each house, we collected data like **surface**, **previous buying price**, **GPS coordinates**, etc.
- These **features** are concatenated into a **vector  $x$** ;
- We need to **learn** the **function**  $f_0$  mapping vectors  $x$  to  $y$ .
- We believe a **linear combination** of the features shoud be a relevant model :

$$y = \theta^T \cdot x.$$

## Bayesian Learning :- Linear regression example

- Suppose we are trying to predict the **selling price**  $y$  of a house.
- For each house, we collected data like **surface**, **previous buying price**, **GPS coordinates**, etc.
- These **features** are concatenated into a **vector  $x$** ;
- We need to **learn** the **function**  $f_0$  mapping vectors  $x$  to  $y$ .
- We believe a **linear combination** of the features should be a relevant model :

$$y = \theta^T \cdot x.$$

## Bayesian Learning :- Linear regression example

- Suppose we are trying to predict the **selling price**  $y$  of a house.
- For each house, we collected data like **surface**, **previous buying price**, **GPS coordinates**, etc.
- These **features** are concatenated into a **vector  $x$** ;
- We need to **learn** the **function**  $f_0$  mapping vectors  $x$  to  $y$ .
- We believe a **linear combination** of the features shoud be a relevant model :

$$y = \theta^T \cdot x.$$

## Bayesian Learning :- Linear regression example

- Yet we also believe that this **linear combination** is just an **approximation** of  $f_0$  and therefore we go for a probabilistic formulation :

$$Y \sim \mathcal{N}(\boldsymbol{\theta}^T \cdot \mathbf{x}, \sigma)$$

- Now the **likelihood** is given by :

$$\text{Likelihood}(\boldsymbol{\theta}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y^{(i)} - \boldsymbol{\theta}^T \cdot \mathbf{x}^{(i)})^2}{2\sigma^2}}$$

- For simplicity, we assume the noise variance  $\sigma^2$  is known.

## Bayesian Learning :- Linear regression example

- Yet we also believe that this **linear combination** is just an **approximation** of  $f_0$  and therefore we go for a probabilistic formulation :

$$Y \sim \mathcal{N}(\boldsymbol{\theta}^T \cdot \mathbf{x}, \sigma)$$

- Now the **likelihood** is given by :

$$\text{Likelihood}(\boldsymbol{\theta}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y^{(i)} - \boldsymbol{\theta}^T \cdot \mathbf{x}^{(i)})^2}{2\sigma^2}}$$

- For simplicity, we assume the noise variance  $\sigma^2$  is known.

## Bayesian Learning :- Linear regression example

- Yet we also believe that this **linear combination** is just an **approximation** of  $f_0$  and therefore we go for a probabilistic formulation :

$$Y \sim \mathcal{N}(\boldsymbol{\theta}^T \cdot \mathbf{x}, \sigma)$$

- Now the **likelihood** is given by :

$$\text{Likelihood}(\boldsymbol{\theta}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y^{(i)} - \boldsymbol{\theta}^T \cdot \mathbf{x}^{(i)})^2}{2\sigma^2}}$$

- For simplicity, we assume the noise variance  $\sigma^2$  is known.

## Bayesian Learning :- Linear regression example

- We already have some knowledge on what values of  $\theta$  are more likely before seeing any datum :

$$\text{Prior}(\theta) = \frac{1}{(2\pi)^{\frac{d}{2}} \det(V_0)^{\frac{1}{2}}} e^{-\frac{1}{2}(\theta - \theta_0)^T \cdot V_0^{-1} (\theta - \theta_0)}$$

- After seeing data  $\mathcal{D}$ , our knowledge is given by the following posterior distribution

$$p(\theta | \mathcal{D}, \theta_0, V_0) \propto \text{Likelihood}(\theta) \times \text{Prior}(\theta).$$

## Bayesian Learning :- Linear regression example

- We already have some knowledge on what values of  $\theta$  are more likely before seeing any datum :

$$\text{Prior}(\theta) = \frac{1}{(2\pi)^{\frac{d}{2}} \det(V_0)^{\frac{1}{2}}} e^{-\frac{1}{2}(\theta - \theta_0)^T \cdot V_0^{-1} (\theta - \theta_0)}$$

- After seeing data  $\mathcal{D}$ , our knowledge is given by the following posterior distribution

$$p(\theta | \mathcal{D}, \theta_0, V_0) \propto \text{Likelihood}(\theta) \times \text{Prior}(\theta).$$

## Bayesian Learning :- Linear regression example

- If the prior parameters are such that  $\theta_0 = \mathbf{0}$  and  $V_0 = \tau^2 I$ , applying  $-\log$  leads to the following cost function (up to an additive constant)

$$J(\boldsymbol{\theta}) = \underbrace{\sum_{i=1}^n \frac{(y^{(i)} - \boldsymbol{\theta}^T \cdot \mathbf{x}^{(i)})^2}{2\sigma^2}}_{\text{Least Squares}} + \underbrace{\frac{1}{\tau^2} \|\boldsymbol{\theta}\|_2}_{\text{Ridge Reg.}}$$

## Bayesian Learning :- Linear regression example

- Going back to probabilities, one can show that the posterior  $p(\boldsymbol{\theta}|\mathcal{D}, \boldsymbol{\theta}_0, \mathbf{V}_0)$  is also Gaussian, in which case our prior is conjugate<sup>2</sup>.

$$p(\boldsymbol{\theta}|\mathcal{D}, \boldsymbol{\theta}_0, \mathbf{V}_0) \sim \mathcal{N}(\boldsymbol{\theta}_n, \mathbf{V}_n), \quad (1)$$

$$\boldsymbol{\theta}_n = \mathbf{V}_n \left( \mathbf{V}_0^{-1} \cdot \boldsymbol{\theta}_0 + \frac{1}{\sigma^2} \mathbf{X}^T \cdot \mathbf{y} \right), \quad (2)$$

$$\mathbf{V}_n = \left( \mathbf{V}_0^{-1} + \frac{1}{\sigma^2} \mathbf{X}^T \cdot \mathbf{X} \right)^{-1}, \quad (3)$$

with

$$\mathbf{X} = \begin{pmatrix} \vdots & & \\ - (\mathbf{x}^{(1)})^T & - \\ \vdots & & \\ - (\mathbf{x}^{(n)})^T & - \end{pmatrix} \quad \text{and} \quad \mathbf{y} = \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{pmatrix}.$$

- Under conjugacy, learning boils down to updating the prior parameters and the updates are easy to compute.

## Bayesian Learning :- Linear regression example

- No fusion for now .. just Bayesian statistics !
- As learners, what we really need is the posterior predictive  $p(y|x, \mathcal{D})$ .
- The expectation of this distribution is our proxy for  $f_0$  and allows to make a prediction for the selling price of a house whose features are the entries of the unseen example  $x$ .

## Bayesian Learning :- Linear regression example

- No fusion for now .. just Bayesian statistics !
- As learners, what we really need is the **posterior predictive**  $p(y|x, \mathcal{D})$ .
- The expectation of this distribution is our proxy for  $f_0$  and allows to make a **prediction** for the selling price of a house whose features are the entries of the unseen example  $x$ .

## Bayesian Learning :- Linear regression example

- No fusion for now .. just Bayesian statistics !
- As learners, what we really need is the **posterior predictive**  $p(y|x, \mathcal{D})$ .
- The expectation of this distribution is our proxy for  $f_0$  and allows to make a **prediction** for the selling price of a house whose features are the entries of the unseen example  $x$ .

## Bayesian Learning :- Linear regression example

- Observe that the predictive distribution is **free of un-observed parameter conditioning**... because we **marginalized them out** :

$$p(y|\mathbf{x}, \mathcal{D}) = \int_{\Theta} p(y|\mathbf{x}, \mathcal{D}, \boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathbf{x}, \mathcal{D}) d\boldsymbol{\theta} \quad (4)$$

- The above calculus is the **weighted combination** of an **infinity** of regressors !
- The weights depend on the ability of each regressor to fit well the data.

## Bayesian Learning :- Linear regression example

- Observe that the predictive distribution is **free of un-observed parameter conditioning**... because we **marginalized** them out :

$$p(y|\mathbf{x}, \mathcal{D}) = \int_{\Theta} p(y|\mathbf{x}, \mathcal{D}, \boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathbf{x}, \mathcal{D}) d\boldsymbol{\theta} \quad (4)$$

- The above calculus is the **weighted combination** of an **infinity** of regressors !
- The weights depend on the ability of each regressor to fit well the data.

## Bayesian Learning :- Linear regression example

- Observe that the predictive distribution is **free of un-observed parameter conditioning**... because we **marginalized** them out :

$$p(y|\mathbf{x}, \mathcal{D}) = \int_{\Theta} p(y|\mathbf{x}, \mathcal{D}, \boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathbf{x}, \mathcal{D}) d\boldsymbol{\theta} \quad (4)$$

- The above calculus is the **weighted combination** of an **infinity** of regressors !
- The weights depend on the ability of each regressor to fit well the data.

## Bayesian Learning : - Linear regression example

- In our linear regression case, we have

$$p(y|\mathbf{x}, \mathcal{D}) = \int_{\Theta} p(y|\mathbf{x}, \mathcal{D}, \boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathbf{x}, \mathcal{D}) d\boldsymbol{\theta}, \quad (5)$$

$$= \int_{\Theta} p(y|\mathbf{x}, \mathcal{D}, \boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathcal{D}) d\boldsymbol{\theta}, \quad (6)$$

$$= \int_{\Theta} G\left(y; \boldsymbol{\theta}^T \cdot \mathbf{x}, \sigma^2\right) G\left(\boldsymbol{\theta}; \boldsymbol{\theta}_n, \mathbf{V}_n\right) d\boldsymbol{\theta}, \quad (7)$$

with  $G$  the Gaussian density function.

- Finally, one can show that

$$y|\mathbf{x}, \mathcal{D} \sim \mathcal{N}\left(\boldsymbol{\theta}_n^T \cdot \mathbf{x}, \sigma_n\right), \quad (8)$$

$$\sigma_n^2 = \sigma^2 + \mathbf{x}^T \cdot \mathbf{V}_n \cdot \mathbf{x}. \quad (9)$$

## Bayesian Learning : - Linear regression example

- In our linear regression case, we have

$$p(y|\mathbf{x}, \mathcal{D}) = \int_{\Theta} p(y|\mathbf{x}, \mathcal{D}, \boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathbf{x}, \mathcal{D}) d\boldsymbol{\theta}, \quad (5)$$

$$= \int_{\Theta} p(y|\mathbf{x}, \mathcal{D}, \boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathcal{D}) d\boldsymbol{\theta}, \quad (6)$$

$$= \int_{\Theta} G\left(y; \boldsymbol{\theta}^T \cdot \mathbf{x}, \sigma^2\right) G\left(\boldsymbol{\theta}; \boldsymbol{\theta}_n, \mathbf{V}_n\right) d\boldsymbol{\theta}, \quad (7)$$

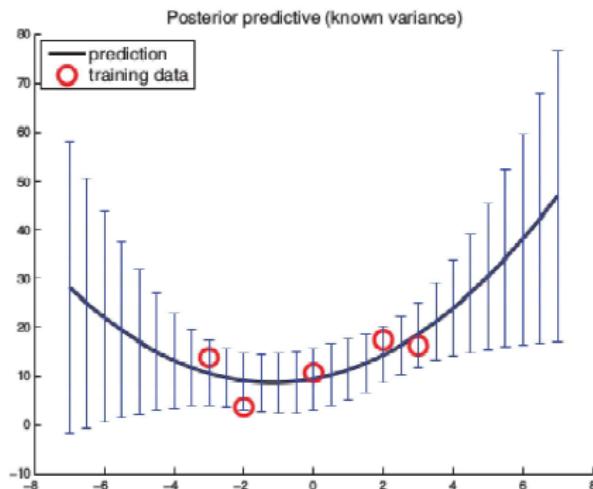
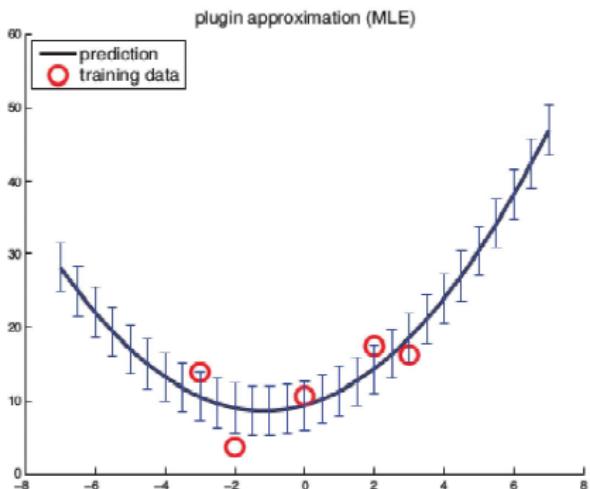
with  $G$  the Gaussian density function.

- Finally, one can show that

$$y|\mathbf{x}, \mathcal{D} \sim \mathcal{N}\left(\boldsymbol{\theta}_n^T \cdot \mathbf{x}, \sigma_n\right), \quad (8)$$

$$\sigma_n^2 = \sigma^2 + \mathbf{x}^T \cdot \mathbf{V}_n \cdot \mathbf{x}. \quad (9)$$

## Bayesian Learning : - Illustration (polynomial reg.)



[Murphy 2012 - 7.6]

## Bayesian Learning :- Comments

- The **posterior predictive** is not always known in closed form → use Monte-Carlo to approximate the **marginalization**.
- Have we really gotten rid of all the parameters ?
- No, we are still conditioning w.r.t.  $\theta_0 = \mathbf{0}$  and  $\mathbf{V}_0$ .
- They can be marginalized out too by introducing a distribution for them called a **hyperprior**. This setting is known as **hierarchical Bayes**.

## Bayesian Learning :- Comments

- The **posterior predictive** is not always known in closed form → use Monte-Carlo to approximate the **marginalization**.
- Have we really gotten rid of all the parameters ?
- No, we are still conditioning w.r.t.  $\theta_0 = \mathbf{0}$  and  $\mathbf{V}_0$ .
- They can be marginalized out too by introducing a distribution for them called a **hyperprior**. This setting is known as **hierarchical Bayes**.

## Bayesian Learning :- Comments

- The **posterior predictive** is not always known in closed form → use Monte-Carlo to approximate the **marginalization**.
- Have we really gotten rid of all the parameters ?
- No, we are still conditioning w.r.t.  $\theta_0 = \mathbf{0}$  and  $\mathbf{V}_0$ .
- They can be marginalized out too by introducing a distribution for them called a **hyperprior**. This setting is known as **hierarchical Bayes**.

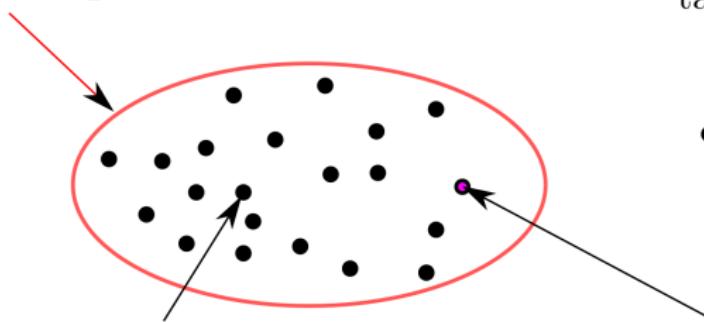
## Bayesian Learning :- Comments

- The **posterior predictive** is not always known in closed form → use Monte-Carlo to approximate the **marginalization**.
- Have we really gotten rid of all the parameters ?
- No, we are still conditioning w.r.t.  $\theta_0 = \mathbf{0}$  and  $\mathbf{V}_0$ .
- They can be marginalized out too by introducing a distribution for them called a **hyperprior**. This setting is known as **hierarchical Bayes**.

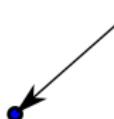
## Bayesian Model Averaging : let's start with model selection

Example : Polynomial regression with small degree  $q = 1$

model  $\mathcal{H}_1$



target function  $f_0$



Usual setting:

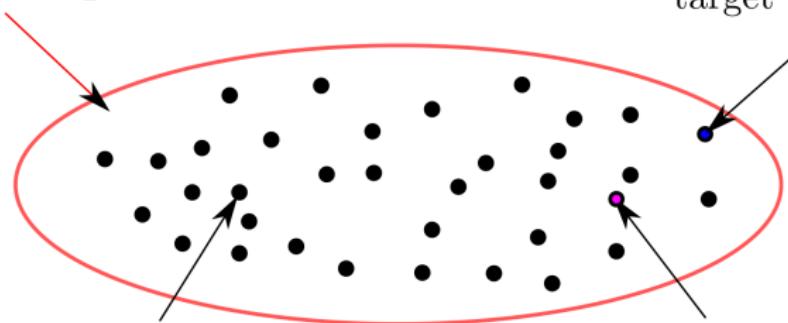
$$\hat{f} \in \mathcal{H}_1$$

chosen hypothesis  $h$   
= learnt estimate  $\hat{f}$

## Bayesian Model Averaging : let's start with model selection

Example : Polynomial regression with higher degree  $q = 2$

model  $\mathcal{H}_2$



target function  $f_0$

Usual setting:

$$\hat{f} \in \mathcal{H}_2$$

hypothesis  $h$

chosen hypothesis  $h$   
= learnt estimate  $\hat{f}$

## Bayesian Model Averaging :

Example : Polynomial regression with degree  $q$

- In model **selection**, the candidate value for  $q$  is sought using, for example, CV.

In general, it could be obtained as

$$q^* = \arg \max_{q \in \mathbb{N}^*} p(q|\mathcal{D}).$$

- In model **averaging**, several values for  $q$  are considered as relevant candidates. We are now elaborating the predictive posterior as

$$\begin{aligned} p(y|\mathbf{x}, \mathcal{D}) &= \sum_{q \in \mathbb{N}^*} p(y|\mathbf{x}, \mathcal{D}, q) p(q|\mathbf{x}, \mathcal{D}), \\ &= \sum_{q \in \mathbb{N}^*} p(y|\mathbf{x}, \mathcal{D}, q) p(q|\mathcal{D}). \end{aligned}$$

## Bayesian Model Averaging :

Example : Polynomial regression with degree  $q$

- In model **selection**, the candidate value for  $q$  is sought using, for example, CV.

In general, it could be obtained as

$$q^* = \arg \max_{q \in \mathbb{N}^*} p(q|\mathcal{D}).$$

- In model **averaging**, several values for  $q$  are considered as relevant candidates. We are now elaborating the predictive posterior as

$$\begin{aligned} p(y|\mathbf{x}, \mathcal{D}) &= \sum_{q \in \mathbb{N}^*} p(y|\mathbf{x}, \mathcal{D}, q) p(q|\mathbf{x}, \mathcal{D}), \\ &= \sum_{q \in \mathbb{N}^*} p(y|\mathbf{x}, \mathcal{D}, q) p(q|\mathcal{D}). \end{aligned}$$

## Bayesian Model Averaging : comments

- BMA only works for probabilistic models allowing to determine both  $p(\mathbf{q}|\mathcal{D})$  and  $p(y|\mathbf{x}, \mathcal{D}, \mathbf{q})$ .
- Its philosophy is close to **hierarchical Bayes** in the sense that each hyperprior parameter choice can be regarded as a given model.

## Bayesian Model Averaging : comments

- BMA only works for probabilistic models allowing to determine both  $p(\mathbf{q}|\mathcal{D})$  and  $p(y|\mathbf{x}, \mathcal{D}, \mathbf{q})$ .
- Its philosophy is close to **hierarchical Bayes** in the sense that each hyperprior parameter choice can be regarded as a given model.

# Chapter organization

- 1 Mixture models
- 2 Ensemble methods
- 3 Bayesian methods
- 4 Statistical aggregation theory

## Statistical aggregation theory :

- The **statistical theory of aggregation** formalizes the intuitions behind the **fusion of prediction functions** in a **supervised learning context**.
- Its main **contribution** are generalization performance warranties.

## Statistical aggregation theory :

- The **statistical theory of aggregation** formalizes the intuitions behind the **fusion of prediction functions** in a **supervised learning context**.
- Its main **contribution** are generalization performance warranties.

## Statistical aggregation theory : problem statement

- For simplicity, suppose the hypothesis set<sup>3</sup>  $\mathcal{H}$  is finite :  $|\mathcal{H}| = m < \infty$ .
- Let  $\mathcal{H}_\Theta$  denote the following space :

$$\mathcal{H}_\Theta = \left\{ f_\theta = \sum_{j=1}^m \theta_j f_j : \theta \in \Theta \right\}.$$

- $\mathcal{H}_\Theta$  is the functional vector space spanned by the members of  $\mathcal{H}$ .

---

3. often called dictionary in this framework.

## Statistical aggregation theory : problem statement

- For simplicity, suppose the hypothesis set<sup>3</sup>  $\mathcal{H}$  is finite :  $|\mathcal{H}| = m < \infty$ .
- Let  $\mathcal{H}_\Theta$  denote the following space :

$$\mathcal{H}_\Theta = \left\{ f_\theta = \sum_{j=1}^m \theta_j f_j : \theta \in \Theta \right\}.$$

- $\mathcal{H}_\Theta$  is the functional vector space spanned by the members of  $\mathcal{H}$ .

---

3. often called dictionary in this framework.

## Statistical aggregation theory : problem statement

- For simplicity, suppose the hypothesis set<sup>3</sup>  $\mathcal{H}$  is finite :  $|\mathcal{H}| = m < \infty$ .
- Let  $\mathcal{H}_\Theta$  denote the following space :

$$\mathcal{H}_\Theta = \left\{ f_\theta = \sum_{j=1}^m \theta_j f_j : \theta \in \Theta \right\}.$$

- $\mathcal{H}_\Theta$  is the functional vector space spanned by the members of  $\mathcal{H}$ .

---

3. often called dictionary in this framework.

## Statistical aggregation theory : problem statement

Depending the  $\Theta \subseteq \mathbb{R}^m$ , several sub-classes of problems are obtained :

- **Model selection** :  $\Theta$  is the set of canonical base vectors, or **one-hot** vectors.

$$\boldsymbol{\theta}^T = [0 \ \dots \ 0 \ 1 \ 0 \ \dots \ 0]$$

- **Convex aggregation** :  $\Theta$  is the canonical simplex of  $\mathbb{R}^m$ ,

$$\sum_{j=1}^m \theta_j = 1 \text{ and } \theta_j \geq 0.$$

- **Linear aggregation** :  $\Theta = \mathbb{R}^m$ .

- **Sparse linear aggregation** :  $\Theta = \{\boldsymbol{\theta} \in \mathbb{R}^m \text{ s.t. } \|\boldsymbol{\theta}\|_0 < r\}$ .

## Statistical aggregation theory : problem statement

Depending the  $\Theta \subseteq \mathbb{R}^m$ , several sub-classes of problems are obtained :

- **Model selection** :  $\Theta$  is the set of canonical base vectors, or **one-hot** vectors.

$$\boldsymbol{\theta}^T = [0 \ \dots \ 0 \ 1 \ 0 \ \dots \ 0]$$

- **Convex aggregation** :  $\Theta$  is the **canonical simplex** of  $\mathbb{R}^m$ ,

$$\sum_{j=1}^m \theta_j = 1 \text{ and } \theta_j \geq 0.$$

- **Linear aggregation** :  $\Theta = \mathbb{R}^m$ .

- **Sparse linear aggregation** :  $\Theta = \{\boldsymbol{\theta} \in \mathbb{R}^m \text{ s.t. } \|\boldsymbol{\theta}\|_0 < r\}$ .

## Statistical aggregation theory : problem statement

Depending the  $\Theta \subseteq \mathbb{R}^m$ , several sub-classes of problems are obtained :

- **Model selection** :  $\Theta$  is the set of canonical base vectors, or one-hot vectors.

$$\boldsymbol{\theta}^T = [0 \quad \dots \quad 0 \quad 1 \quad 0 \quad \dots \quad 0]$$

- **Convex aggregation** :  $\Theta$  is the canonical simplex of  $\mathbb{R}^m$ ,

$$\sum_{j=1}^m \theta_j = 1 \text{ and } \theta_j \geq 0.$$

- **Linear aggregation** :  $\Theta = \mathbb{R}^m$ .

- **Sparse linear aggregation** :  $\Theta = \{\boldsymbol{\theta} \in \mathbb{R}^m \text{ s.t. } \|\boldsymbol{\theta}\|_0 < r\}$ .

## Statistical aggregation theory : problem statement

Depending the  $\Theta \subseteq \mathbb{R}^m$ , several sub-classes of problems are obtained :

- **Model selection** :  $\Theta$  is the set of canonical base vectors, or one-hot vectors.

$$\theta^T = [0 \quad \dots \quad 0 \quad 1 \quad 0 \quad \dots \quad 0]$$

- **Convex aggregation** :  $\Theta$  is the canonical simplex of  $\mathbb{R}^m$ ,

$$\sum_{j=1}^m \theta_j = 1 \text{ and } \theta_j \geq 0.$$

- **Linear aggregation** :  $\Theta = \mathbb{R}^m$ .

- **Sparse linear aggregation** :  $\Theta = \{\theta \in \mathbb{R}^m \text{ s.t. } \|\theta\|_0 < r\}$ .

**Statistical aggregation theory** : problem statement

Let  $R(f_{\theta})$  denote the **generalization error (risk)** of the pre-estimator  $f_{\theta}$  :

$$R(f_{\theta}) = \mathbb{E}_{X,Y} [L(f_{\theta}(x), y)].$$

### Definition

$\theta^*$  is an **oracle** for the aggregation problem if

$$R(f_{\theta^*}) = \inf_{\theta \in \Theta} R(f_{\theta}).$$

## Statistical aggregation theory : problem statement

What are we looking for ?

- An aggregated estimate  $\hat{\theta}$  close to the oracle  $\theta^*$ .
- An aggregated estimate  $\hat{\theta}$  learnt from private data  $\mathcal{D}_{\text{val}}$ .

Close in what sense ?

- In the following probabilistic sense :

$$\mathbb{E}_{\mathcal{D}_{\text{val}}} [R(f_{\hat{\theta}})] \leq C \inf_{\theta \in \Theta} \{R(f_{\theta}) + \Delta_{n,m}(\theta)\},$$

with constant  $C \geq 1$  and rate  $\Delta_{n,m} \rightarrow 0$  as  $n \rightarrow \infty$ .

**Assymptotical warranty** : when  $n$  is big, the risk is not worse (in average) than  $C$  times the oracle risk.

## Statistical aggregation theory : problem statement

What are we looking for ?

- An aggregated estimate  $\hat{\theta}$  close to the oracle  $\theta^*$ .
- An aggregated estimate  $\hat{\theta}$  learnt from private data  $\mathcal{D}_{\text{val}}$ .

Close in what sense ?

- In the following probabilistic sense :

$$\mathbb{E}_{\mathcal{D}_{\text{val}}} [R(f_{\hat{\theta}})] \leq C \inf_{\theta \in \Theta} \{R(f_{\theta}) + \Delta_{n,m}(\theta)\},$$

with constant  $C \geq 1$  and rate  $\Delta_{n,m} \rightarrow 0$  as  $n \rightarrow \infty$ .

Assymptotical warranty : when  $n$  is big, the risk is not worse (in average) than  $C$  times the oracle risk.

## Statistical aggregation theory : problem statement

What are we looking for ?

- An aggregated estimate  $\hat{\theta}$  close to the oracle  $\theta^*$ .
- An aggregated estimate  $\hat{\theta}$  learnt from private data  $\mathcal{D}_{\text{val}}$ .

Close in what sense ?

- In the following probabilistic sense :

$$\mathbb{E}_{\mathcal{D}_{\text{val}}} [R(f_{\hat{\theta}})] \leq C \inf_{\theta \in \Theta} \{R(f_{\theta}) + \Delta_{n,m}(\theta)\},$$

with constant  $C \geq 1$  and rate  $\Delta_{n,m} \rightarrow 0$  as  $n \rightarrow \infty$ .

Assymptotical warranty : when  $n$  is big, the risk is not worse (in average) than  $C$  times the oracle risk.

## Statistical aggregation theory : problem statement

What are we looking for ?

- An aggregated estimate  $\hat{\theta}$  close to the oracle  $\theta^*$ .
- An aggregated estimate  $\hat{\theta}$  learnt from private data  $\mathcal{D}_{\text{val}}$ .

Close in what sense ?

- In the following probabilistic sense :

$$\mathbb{E}_{\mathcal{D}_{\text{val}}} [R(\hat{f}_\theta)] \leq C \inf_{\theta \in \Theta} \{R(f_\theta) + \Delta_{n,m}(\theta)\},$$

with constant  $C \geq 1$  and rate  $\Delta_{n,m} \rightarrow 0$  as  $n \rightarrow \infty$ .

**Assymptotical warranty** : when  $n$  is big, the risk is not worse (in average) than  $C$  times the oracle risk.

## Statistical aggregation theory : problem analysis

- There is a bias-variance tradeoff binding approximation quality ( $\Delta_{n,m}$ ) and the size of the parameter space  $\Theta$ .
- In the  $C = 1$  case, oracle inequalities are said to be sharp.

## Statistical aggregation theory : problem analysis

- There is a bias-variance tradeoff binding approximation quality ( $\Delta_{n,m}$ ) and the size of the parameter space  $\Theta$ .
- In the  $C = 1$  case, oracle inequalities are said to be sharp.

## Statistical aggregation theory : problem analysis

- New question : for a given  $\Theta$ , what is the best rate  $\Delta_{n,m}$  ?

### Definition

Let  $\mathcal{F}$  denote a functional space.

Optimal rates or **minimax** rates are a sequence  $\Phi_{n,m}$  such that

- (i) For any pre-estimates  $f_1, \dots, f_m$  and any target  $f_0$ , there is an aggregate  $f_{\tilde{\theta}}$  s.t.

$$R(f_{\tilde{\theta}}) - R(f_0) \leq c_1 \Phi_{n,m}$$

- (ii) For any estimator  $f_{\tilde{\theta}}$  trained on  $\mathcal{D}_{\text{val}}$ ,

$$\sup_{f_0 \in \mathcal{F}} \{R(f_{\tilde{\theta}}) - R(f_0)\} \geq c_2 \Phi_{n,m}.$$

$c_1$  and  $c_2$  are positive constants. Note that we may have  $f_{\tilde{\theta}} \notin \mathcal{H}_{\Theta}$ .

## Statistical aggregation theory : problem analysis

- New question : for a given  $\Theta$ , what is the best rate  $\Delta_{n,m}$  ?

### Definition

Let  $\mathcal{F}$  denote a functional space.

Optimal rates or **minimax** rates are a sequence  $\Phi_{n,m}$  such that

- (i) For any pre-estimates  $f_1, \dots, f_m$  and any target  $f_0$ , there is an aggregate  $f_{\hat{\theta}}$  s.t.

$$R(f_{\hat{\theta}}) - R(f_0) \leq c_1 \Phi_{n,m}$$

- (ii) For any estimator  $f_{\tilde{\theta}}$  trained on  $\mathcal{D}_{\text{val}}$ ,

$$\sup_{f_0 \in \mathcal{F}} \{R(f_{\tilde{\theta}}) - R(f_0)\} \geq c_2 \Phi_{n,m}.$$

$c_1$  and  $c_2$  are positive constants. Note that we may have  $f_{\tilde{\theta}} \notin \mathcal{H}_{\Theta}$ .

## Statistical aggregation theory : problem analysis

- New question : for a given  $\Theta$ , what is the best rate  $\Delta_{n,m}$  ?

### Definition

Let  $\mathcal{F}$  denote a functional space.

Optimal rates or **minimax** rates are a sequence  $\Phi_{n,m}$  such that

- (i) For any pre-estimates  $f_1, \dots, f_m$  and any target  $f_0$ , there is an aggregate  $f_{\hat{\theta}}$  s.t.

$$R(f_{\hat{\theta}}) - R(f_0) \leq c_1 \Phi_{n,m}$$

- (ii) For any estimator  $f_{\tilde{\theta}}$  trained on  $\mathcal{D}_{\text{val}}$ ,

$$\sup_{f_0 \in \mathcal{F}} \{R(f_{\tilde{\theta}}) - R(f_0)\} \geq c_2 \Phi_{n,m}.$$

$c_1$  and  $c_2$  are positive constants. Note that we may have  $f_{\tilde{\theta}} \notin \mathcal{H}_{\Theta}$ .

## Statistical aggregation theory : problem analysis

- New question : for a given  $\Theta$ , what is the best rate  $\Delta_{n,m}$  ?

### Definition

Let  $\mathcal{F}$  denote a functional space.

Optimal rates or **minimax** rates are a sequence  $\Phi_{n,m}$  such that

- (i) For any pre-estimates  $f_1, \dots, f_m$  and any target  $f_0$ , there is an aggregate  $f_{\hat{\theta}}$  s.t.

$$R(f_{\hat{\theta}}) - R(f_0) \leq c_1 \Phi_{n,m}$$

- (ii) For any estimator  $f_{\tilde{\theta}}$  trained on  $\mathcal{D}_{\text{val}}$ ,

$$\sup_{f_0 \in \mathcal{F}} \{R(f_{\tilde{\theta}}) - R(f_0)\} \geq c_2 \Phi_{n,m}.$$

$c_1$  and  $c_2$  are positive constants. Note that we may have  $f_{\tilde{\theta}} \notin \mathcal{H}_{\Theta}$ .

## Statistical aggregation theory : problem analysis

Example : for linear aggregation of regressors (under additive Gaussian noise), we have

$$\Phi_{n,m} = \min \left\{ 1; \frac{m}{n} \right\}.$$

- The dimensionality of  $\Theta$  increases w.r.t.  $m$ .
- The more data, the better the generalization.

## Statistical aggregation theory : problem analysis

Example : for linear aggregation of regressors (under additive Gaussian noise), we have

$$\Phi_{n,m} = \min \left\{ 1; \frac{m}{n} \right\}.$$

- The dimensionality of  $\Theta$  increases w.r.t.  $m$ .
- The more data, the better the generalization.

# Statistical aggregation theory : some solutions

## Penalized approaches :

- Since the distribution  $p_{X,Y}$  is **unknown**, explicit risk computation is **impossible**.

- A **biased** proxy for it is the **empirical risk**

$$\begin{aligned} R_{\text{emp}}(f_{\theta}) &= \sum_{i=1}^n L\left(y^{(i)}, f_{\theta}\left(\mathbf{x}^{(i)}\right)\right), \\ &= \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}_{\text{val}}} L\left(y^{(i)}, \theta_1 f_1\left(\mathbf{x}^{(i)}\right) + \dots + \theta_m f_m\left(\mathbf{x}^{(i)}\right)\right). \end{aligned}$$

- Keep in mind that for each  $f_k$ , the training error is

$$Err_{\text{train}} = \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}_{\text{train}}} L\left(y^{(i)}, f_k\left(\mathbf{x}^{(i)}\right)\right).$$

# Statistical aggregation theory : some solutions

## Penalized approaches :

- Since the distribution  $p_{X,Y}$  is **unknown**, explicit risk computation is **impossible**.
- A **biased** proxy for it is the **empirical risk**

$$\begin{aligned} R_{\text{emp}}(f_{\theta}) &= \sum_{i=1}^n L\left(y^{(i)}, f_{\theta}\left(\mathbf{x}^{(i)}\right)\right), \\ &= \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}_{\text{val}}} L\left(y^{(i)}, \theta_1 f_1\left(\mathbf{x}^{(i)}\right) + \dots + \theta_m f_m\left(\mathbf{x}^{(i)}\right)\right). \end{aligned}$$

- Keep in mind that for each  $f_k$ , the training error is

$$Err_{\text{train}} = \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}_{\text{train}}} L\left(y^{(i)}, f_k\left(\mathbf{x}^{(i)}\right)\right).$$

# Statistical aggregation theory : some solutions

## Penalized approaches :

- Since the distribution  $p_{X,Y}$  is **unknown**, explicit risk computation is **impossible**.
- A **biased** proxy for it is the **empirical risk**

$$\begin{aligned} R_{\text{emp}}(f_{\theta}) &= \sum_{i=1}^n L\left(y^{(i)}, f_{\theta}\left(\mathbf{x}^{(i)}\right)\right), \\ &= \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}_{\text{val}}} L\left(y^{(i)}, \theta_1 f_1\left(\mathbf{x}^{(i)}\right) + \dots + \theta_m f_m\left(\mathbf{x}^{(i)}\right)\right). \end{aligned}$$

- Keep in mind that for each  $f_k$ , the training error is

$$Err_{\text{train}} = \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}_{\text{train}}} L\left(y^{(i)}, f_k\left(\mathbf{x}^{(i)}\right)\right).$$

# Statistical aggregation theory : some solutions

## Penalized approaches :

- A possibility to learn  $\hat{\theta}$  is to minimize the empirical risk

$$\hat{\theta} = \arg \min_{\theta} R_{\text{emp}}(f_{\theta}).$$

- This is doomed to overfit, and therefore we solve

$$\hat{\theta} = \arg \min_{\theta} R_{\text{emp}}(f_{\theta}) + \text{Regularizer}(\theta).$$

- This penalized problem mimics the oracle bound, yet the regularizer does not depend on  $n$  !
- By carefully choosing Regularizer( $\theta$ ) (with a  $\|\theta\|_0$  term in it), then minimax rates are achieved for linear, convex and model selection aggregation problems.

# Statistical aggregation theory : some solutions

## Penalized approaches :

- A possibility to learn  $\hat{\theta}$  is to minimize the empirical risk

$$\hat{\theta} = \arg \min_{\theta} R_{\text{emp}}(f_{\theta}).$$

- This is doomed to **overfit**, and therefore we solve

$$\hat{\theta} = \arg \min_{\theta} R_{\text{emp}}(f_{\theta}) + \text{Regularizer}(\theta).$$

- This penalized problem mimics the **oracle bound**, yet the **regularizer** does not depend on  $n$  !
- By carefully choosing  $\text{Regularizer}(\theta)$  (with a  $\|\theta\|_0$  term in it), then minimax rates are achieved for linear, convex and model selection aggregation problems.

# Statistical aggregation theory : some solutions

## Penalized approaches :

- A possibility to learn  $\hat{\theta}$  is to minimize the empirical risk

$$\hat{\theta} = \arg \min_{\theta} R_{\text{emp}}(f_{\theta}).$$

- This is doomed to **overfit**, and therefore we solve

$$\hat{\theta} = \arg \min_{\theta} R_{\text{emp}}(f_{\theta}) + \text{Regularizer}(\theta).$$

- This penalized problem mimics the **oracle bound**, yet the **regularizer** does not depend on  $n$  !
- By carefully choosing  $\text{Regularizer}(\theta)$  (with a  $\|\theta\|_0$  term in it), then minimax rates are achieved for linear, convex and model selection aggregation problems.

## Statistical aggregation theory : some solutions

### Penalized approaches :

- A possibility to learn  $\hat{\theta}$  is to minimize the empirical risk

$$\hat{\theta} = \arg \min_{\theta} R_{\text{emp}}(f_{\theta}).$$

- This is doomed to **overfit**, and therefore we solve

$$\hat{\theta} = \arg \min_{\theta} R_{\text{emp}}(f_{\theta}) + \text{Regularizer}(\theta).$$

- This penalized problem mimics the **oracle bound**, yet the **regularizer** does not depend on  $n$  !
- By carefully choosing  $\text{Regularizer}(\theta)$  (with a  $\|\theta\|_0$  term in it), then minimax rates are achieved for linear, convex and model selection aggregation problems.

## Statistical aggregation theory : some solutions

### Exponential weights :

- In the convex aggregation problem, weights  $[\theta_1 \dots \theta_m]$  can be regarded as a probability distribution.
- By choosing a regularizer involving the Kullback-Leibler divergence between  $\theta$  and a prior  $\pi$  and by replacing the empirical risk with the surrogate

$$\sum_{j=1}^m \theta_k R_{\text{emp}}(f_k),$$

we obtain

$$\hat{\theta} = \text{smax} \left( n \begin{bmatrix} \pi_1 R_{\text{emp}}(f_1) \\ \vdots \\ \pi_m R_{\text{emp}}(f_m) \end{bmatrix} \right).$$

## Statistical aggregation theory : some solutions

### Exponential weights :

- In the convex aggregation problem, weights  $[\theta_1 \dots \theta_m]$  can be regarded as a probability distribution.
- By choosing a regularizer involving the Kullback-Leibler divergence between  $\theta$  and a prior  $\pi$  and by replacing the empirical risk with the surrogate

$$\sum_{j=1}^m \theta_k R_{\text{emp}}(f_k),$$

we obtain

$$\hat{\theta} = \text{smax} \left( n \begin{bmatrix} \pi_1 R_{\text{emp}}(f_1) \\ \vdots \\ \pi_m R_{\text{emp}}(f_m) \end{bmatrix} \right).$$

## Statistical aggregation theory : our local expert

- Benjamin Guedj (INRIA / Painlevé - MODAL) [homepage](#)
- Internship position.

## Statistical aggregation theory : our local expert

- Benjamin Guedj (INRIA / Painlevé - MODAL) [homepage](#)
- Internship position.