

Let consider a general mixture model. We observe  $N$  random variables  $x_1, x_2, \dots, x_N$  which are independently and identically distributed with  $x_i \sim f_{\boldsymbol{\pi}}(x_i)$  where  $f_{\boldsymbol{\pi}}$  is given by:

$$f_{\boldsymbol{\pi}}(x) = \sum_{j=1}^K \pi_j f_j(x) \quad (1)$$

Let suppose that each component density  $f_i$  is known, but not necessarily Gaussian. We will focus on the estimation of the weights vector  $\boldsymbol{\pi} \in \mathbb{R}^K$  and assume that this vector is sparse. We will focus our study on the performance of the Maximum Likelihood Estimator (et excess risk ?). From a rewriting of the loglikelihood we can define  $\Phi_N(\boldsymbol{\pi})$  as following:

$$\Phi_N(\boldsymbol{\pi}) = -\frac{1}{N} \sum_{i=1}^N \log f_{\boldsymbol{\pi}}(x_i) \quad (2)$$

We can rewrite the minimization problem:

$$\hat{\boldsymbol{\pi}} \in \arg \min_{\boldsymbol{\pi} \in \Pi} \{\Phi_N(\boldsymbol{\pi})\}, \quad \Pi = \{\boldsymbol{\pi} \in [0, 1]^K : \sum_{j=1}^K \pi_j = 1\} \quad (3)$$

For theoretical objectives, we will make the following assumption:

**Hypothesis 1.** *All realizations are not probably unlikely to be observed. Therefore  $\exists m > 0$  such that  $f_{\boldsymbol{\pi}}(x) \geq m$  for all  $x \in \{x_1, \dots, x_N\}$ .*

Let denote  $M = \max_{x \in \{x_1, \dots, x_N\}, j \in [K]} \{f_j(x)\}$ , since  $\boldsymbol{\pi}^T \mathbf{1} = 1$  then  $f_{\boldsymbol{\pi}} \leq M$ . Therefore,  $\forall \boldsymbol{\pi} \in \Pi, f_{\boldsymbol{\pi}}(x_1, \dots, x_N) \in [m, M]$ . We have the following lemma:

**Lemma 1.** *Under hypothesis 1,  $\Phi_N$  is Lipschitz-smooth and strongly convex.*

*Proof.* For each  $i \in [K]$ ,  $g_{x_i}(\boldsymbol{\pi}) = f_{\boldsymbol{\pi}}(x_i)$  is a linear function defined on the convex compact set  $\Pi$  and it's image is the interval  $[m, M]$  where  $m > 0$ . We will prove that  $-\log$  is strongly convex on  $[m, M]$ .

$$\forall x \in [m, M], \frac{1}{M^2} \leq \frac{d^2(-\log)}{dx^2}(x) = \frac{1}{x^2} \leq \frac{1}{m^2}. \quad (4)$$

The first inequality proves the  $1/M^2$ -strong convexity of  $-\log$ , the second proves that it is  $1/m^2$ -Lipschitz smooth. The sum of strongly convex functions is strongly convex. Therefore,  $\Phi_N$  is strongly convex.  $\square$

With these nice property under assumption 1, the minimization problem can be rewritten as follows:

$$\hat{\boldsymbol{\pi}} \in \arg \min_{\boldsymbol{\pi} \in \Pi} \{\Phi_N(\boldsymbol{\pi})\}, \quad \Pi = \{\boldsymbol{\pi} \in [0, 1]^K : \boldsymbol{\pi}^T \mathbf{1} = 1, \forall i \in [N], \sum_{j=1}^K \pi_j f_j(x_i) \geq m\} \quad (5)$$

24 In this work, we will study different loss function:  $\|\hat{\pi} - \pi^*\|_1$ ,  $\|\hat{\pi} - \pi^*\|_2$  and some  
 25  $\text{dist}(f_{\hat{\pi}}, f_{\pi^*})$  (donner un exemple). It turns out that this problem is close to the  
 26 regression with random design in the context of transductive learning [Bellec et al.](#)  
 27 (2016) since we do not observe the true cluster labels in our problem. We can  
 28 consider  $\Phi_N$  as a function of two random variable  $X_i$  and  $\pi$ :

$$\Phi_N(\pi) = \frac{1}{N} \sum_{i=1}^N \varphi(x_i, \pi) \quad (6)$$

29 In this setting  $\varphi(.,.)$  (in our problem it is  $-\log(f(.))$ ) is strongly convex and Lips-  
 30 chitz smooth. We will recall some interesting results for our work on regression with  
 31 random design.

32

33 Let consider the following trace regression model :

$$Y_i = \text{tr}(X_i^T \mathbf{B}^*) + \xi_i \quad i = 1, \dots, N \quad (7)$$

34 with  $B^* \in \mathbb{R}^{p \times q}$  and let assume that  $\text{rank}(B^*)$  is small. Let denote  $\sigma = [\sigma_1, \dots, \sigma_p]$   
 35 the singular values of  $B^*$ . The rank of this matrix is given by  $\|\sigma\|_0$ . Unfortunately,  
 36 the  $L_0$  norm is not convex, we tackle this problem by considering the convex  $L_1$   
 37 norm  $\|\sigma\|_1$ . Assume the constraint  $\sigma^T \mathbf{1} = 1$ , then according to [Koltchinskii et al.](#)  
 38 (2016) (quel theoreme ?) an empirical risk minimization method or a Maximum  
 39 Likelihood Estimator with this constraint leads to a sparse estimator  $\hat{\mathbf{B}}$ .

40

41 Therefore, it might be interesting to compare this result with our problem [5](#)

## 42 1 Error Bound

43 Using the strong convexity property, we have:

$$(\nabla \Phi_n(\pi^*) - \nabla \Phi_n(\hat{\pi}))^T (\pi^* - \hat{\pi}) \geq \frac{1}{M^2} \|\pi^* - \hat{\pi}\|^2 \quad (8)$$

44 By definition of the estimator  $\hat{\pi}$  we have  $\nabla \Phi_n(\hat{\pi}) = 0$  therefore, we develop  $\nabla \Phi_n(\pi^*)^T (\pi^* -$   
 45  $\hat{\pi})$ ; for  $l \in [K]$ :

$$[\nabla \Phi_n(\pi^*)]_l = -\frac{1}{N} \sum_{i=1}^N \frac{f_l(x_i)}{\sum_{j=1}^K \pi_j^* f_j(x_i)} \quad (9)$$

46 Therefore:

$$\nabla \Phi_n(\pi^*)^T(\pi^* - \hat{\pi}) = -\frac{1}{N} \sum_{l=1}^K \sum_{i=1}^N \frac{f_l(x_i)(\pi_l^* - \hat{\pi}_l)}{\sum_{j=1}^K \pi_j^* f_j(x_i)} \quad (10)$$

$$= -\frac{1}{N} \sum_{i=1}^N \left( \frac{\sum_{l=1}^K \pi_l^* f_l(x_i)}{\sum_{j=1}^K \pi_j^* f_j(x_i)} - \frac{\sum_{l=1}^K \hat{\pi}_l f_l(x_i)}{\sum_{j=1}^K \pi_j^* f_j(x_i)} \right) \quad (11)$$

$$= \frac{1}{N} \sum_{i=1}^N \left( \frac{f_{\hat{\pi}}(x_i)}{f_{\pi^*}(x_i)} - 1 \right) = \frac{1}{N} \sum_{i=1}^N Z_i \quad (12)$$

47 The idea is to use the Bernstein inequality,  $Z_1, \dots, Z_N$  are independent real-valued  
 48 random variables, we need to prove that there exist a constant  $b$  such that  $\mathbb{E}[Z_i^2] \leq \infty$   
 49 and  $|Z_i - \mathbb{E}Z_i| \leq b$ .

$$\mathbb{E}[Z_i^2] = \mathbb{E}\left[\frac{f_{\hat{\pi}}^2}{f_{\pi^*}^2} - 2\frac{f_{\hat{\pi}}}{f_{\pi^*}} + 1\right] = \mathbb{E}\left[\frac{f_{\hat{\pi}}^2}{f_{\pi^*}^2}\right] - 1 \quad (13)$$

50 **Note:** On l'assume pour l'instant, c'est surement le cas si le support de  $\hat{\pi}$  est inclu  
 51 dans celui de  $\pi^*$ , en effet, soit  $\hat{S}$  le support de  $\hat{\pi}$  et  $S^*$  le support de  $\pi^*$ , alors, pour  
 52 tout  $x \in \mathbb{R}^p$  on note  $q = \arg \max_{l \in \hat{S}} f_l(x)$  et donc

$$\frac{\sum_{l \in \hat{S}} \hat{\pi}_l f_l(x)}{\sum_{l \in S^*} \pi_l^* f_l(x)} \leq \frac{f_q(x)}{\pi_q^* f_q(x)} = \frac{1}{\pi_q^*} \quad (14)$$

53 donc

$$\mathbb{E}\left[\frac{f_{\hat{\pi}}^2}{f_{\pi^*}^2}\right] = \int_{\mathbb{R}^p} \frac{f_{\hat{\pi}}(x)^2}{f_{\pi^*}(x)^2} f_{\pi^*}(x) dx \leq \max_{l \in \hat{S}} (\pi_l^*)^{-2} \quad (15)$$

54 il faut donc etudier le support de l'estimateur

55 We can use now the Bernstein inequality (à completer):

$$|\bar{Z}_N - \mathbb{E}[\bar{Z}_N]| \leq \sigma_N \left( \frac{2 \log(2/\delta)}{N} \right)^{1/2} \left[ 1 + \frac{b}{6N\sigma_N} \left( \frac{2 \log(2/\delta)}{N} \right)^{1/2} \right] \quad (16)$$

56 Since  $\mathbb{E}[Z_i] = 0$ , we have:

$$\|\pi^* - \hat{\pi}\|^2 \leq M^2 \sigma_N \left( \frac{2 \log(2/\delta)}{N} \right)^{1/2} \left[ 1 + \frac{b}{6N\sigma_N} \left( \frac{2 \log(2/\delta)}{N} \right)^{1/2} \right] \quad (17)$$

## 57 References

- 58 P. C. Bellec, A. S. Dalalyan, E. Grappin, and Q. Paris. On the prediction loss of  
 59 the lasso in the partially labeled setting. 2016.
- 60 V. Koltchinskii, K. Lounici, and A. B. Tsybakov. Nuclear-norm penalization and  
 61 optimal rates for noisy low-rank matrix completion. *Annals of Statistics*, 2016.