

# Feature Selection in Finite Mixture of Regression Models with Diverging Number of Parameters

Abbas Khalili

Department of Mathematics and Statistics,

McGill University, Montreal

(Joint work with Shili Lin, The Ohio State University.)

June 10, 2011.

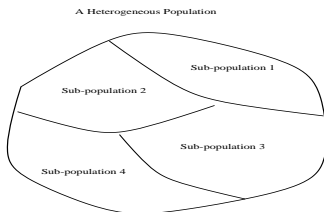
- Introduction
- Finite mixture of regression (FMR) models with diverging number of parameters.
- Feature selection problem in FMR models with diverging number of parameters.
- A new feature selection method in FMR models and its statistical properties.
- Analysis of the Parkinson's disease data.

- Recent advancements in medical and other fields of scientific research have allowed scientists to collect data of unprecedented size and complexity.
- A common statistical problem in such applications is:  
To model a response variable  $Y$  as a function of a small subset of features  $\mathbf{x} = (x_1, x_2, \dots, x_p)^\top$ .
- This is often referred to as a *feature (or variable) selection problem*.

- Sometimes, in applications, many variables are introduced to reduce possible modeling biases, but the number of variables a model can accommodate is often limited by the amount of data available.
- In other words, the number of variables considered depends on the sample size, i.e.  $p_n$ , which reflects the estimability of the parametric model.

- The problem becomes even more complex when the population under study is made up of *subpopulations* and the relationship between  $Y$  and  $x$  varies across the subpopulations.
- *Finite mixture of regression (FMR) models* provide a flexible statistical tool in studying such relationships.

- In some applications the population under study is heterogeneous, i.e. it is made of sub-populations:



- Each sub-population calls for its own regression modelling between  $Y$  and  $x$ .
- FMR models provide a natural way to model **unobserved** heterogeneity in such populations.

## Example: Parkinson disease (PD)



- Parkinson disease (PD) is a neurological disorder.
- Over one million people in North America have PD.
- No cure available for PD, although current medication are effective in controlling its symptoms.
- Early **diagnostic** and **monitoring** is crucial in controlling the disease and in improving the life quality for the patients.

- The main symptoms of PD are tremor, rigidity and other general movement disorders, as well as vocal impairment.
- Tracking PD symptoms involves various physical examinations performed by trained clinical staff.
- A medical rater subjectively assesses the ability of a patient in performing certain tasks. The physical tests are mapped to a metric that is designed to follow PD progression.
- A typical metric of such is the [Unified Parkinson's Disease Rating Scale \(UPDRS\)](#), which reflects the presence and severity of symptoms of PD.

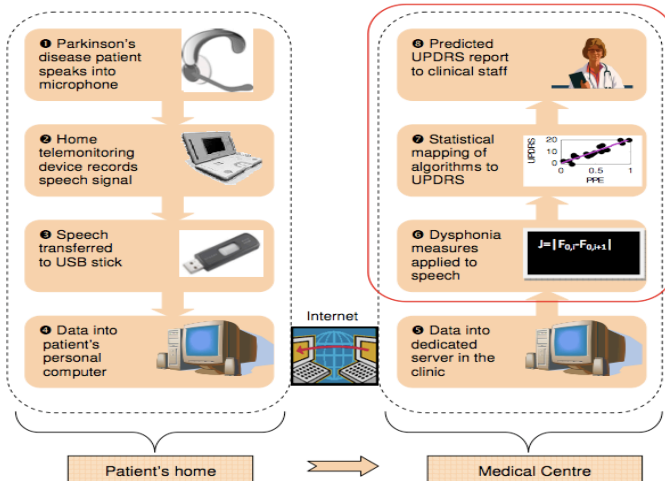


- However, **UPDRS** requires the patient's presence in the clinic and is also a time-consuming physical examination by trained clinical staff.
- Thus, symptom monitoring is **costly** and **logistically inconvenient** for both patients and clinical staff.
- An affordable, cost-effective and reliable alternative:

### Telemonitoring

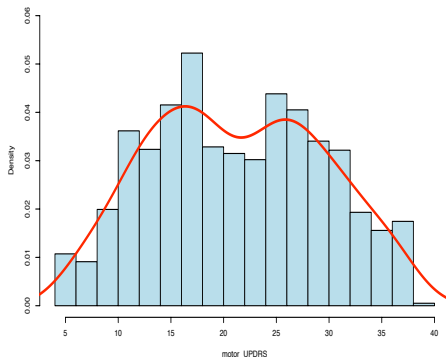
- Noninvasive telemonitoring is an emerging option in general medical care.
- Such options are also considered in PD research.
- The potential for telemonitoring of PD depends heavily on the design of simple tests that can be self-administrated quickly and remotely.
- Recording **speech or vocal properties** are good candidates in this regard.
- Research has shown that approximately 90% of people with PD exhibit some form of vocal impairment.

Tsanas et al. (2010)



- Whether dysphonic features extracted from speech signals recorded at home can be used as surrogate to study PD severity and progression?
- That is to find a statistical model between **speech properties** and **UPDRS**.
- The data under our consideration contains observations on  $p$  speech features.
- The goal is to select an optimally reduced subset of the  $p$  features,  $\mathbf{x} = (x_1, x_2, \dots, x_p)$ , that are predictive of the **UPDRS**, the response variable,  $Y$ , leading to a clinically useful model.

## Histogram of $Y$ (UPDRS)



Notice the bimodality of the histogram !

- Simple exploration indicates that a single regression model is unlikely to be able to accommodate the multimodality nature of the data.
- Some of the speech features are highly correlated.
- Not all the covariates  $x_1, x_2, \dots, x_p$  are significant in explaining the random behaviour of the UPDRS.
- In general, as sample size increases, more features will likely be extracted from additional algorithms in the hope that greater prediction accuracy can be achieved. That is,  $p$  increases as  $n$  increases.

- From statistical standpoint, the challenging issue in analyzing the PD data falls into the general category of *feature or variable selection* in a regression model.
- We will use an **FMR** model to analyze the data.

- We discuss the problem of feature selection in *Finite Mixture of Regression (FMR)* Models.



- Response variable:  $Y \in \mathcal{Y} \subset \mathbb{R}$

Covariates:  $\mathbf{x}^\top = (x_1, x_2, \dots, x_p) \in \mathcal{X} \subset \mathbb{R}^p$ .

- In a GLM model the  $Y$  and  $\mathbf{x}$  are related through a density function

$$h(y; \theta(\mathbf{x}), \phi)$$

with a known link function  $\theta(\mathbf{x}) = \mathcal{L}(\beta_{0k} + \mathbf{x}^\top \beta_k)$ .

- The model claims that the response  $y$  is a function of  $\mathbf{x}$  through a known link function  $\mathcal{L}(\cdot)$ .

- Let  $\mathcal{H} = \{h(y; \theta, \phi); (\theta, \phi) \in \Theta \times \Phi \subset \mathbb{R} \times \mathbb{R}^+\}$  be a parametric family of density functions for  $Y$  with respect to a  $\sigma$ -finite measure  $\mu$ , and  $\phi$  is a dispersion parameter.
- The conditional density function of  $Y$  given  $\mathbf{x}$  in a FMR model with  $K$  components is:

$$f(y; \mathbf{x}, \Psi) = \sum_{k=1}^K \pi_k h(y; \theta_k(\mathbf{x}), \phi_k)$$

with a known link function  $\theta_k(\mathbf{x}) = \mathcal{L}(\beta_{0k} + \mathbf{x}^\top \beta_k)$ , and

$$\sum_{k=1}^K \pi_k = 1, \pi_k > 0.$$

- In some of the recent feature selection problems the dimension of the model is large and may grow with the sample size.
- In this talk we allow the dimension  $p$  of the feature vector  $\mathbf{x}$  to increase with the sample size  $n$ , that is:

$$\mathbf{x}_n^\top = (x_1, x_2, \dots, x_{p_n}) \in \mathcal{X} \subset \mathbb{R}^{p_n}, p_n \rightarrow \infty \text{ as } n \rightarrow \infty.$$

- The component-wise regression vectors:  $\beta_{nk}$
- The diverging vector of parameters:

$$\Psi_n = \{(\beta_{0k}, \beta_{nk}, \pi_k, \phi_k); k = 1, 2, \dots, K\}$$

- Naturally, when there are a large number of covariates, there may be only a small subset of covariates that are significant within each subpopulation.
- That is, for each  $k$ , many of the  $\beta_{kj}$ 's are zero.
- The goal is to identify those  $x_j$ 's for which  $\beta_{kj} \neq 0$ .
- In such situations the FMR model of interest is called:

SPARSE FMR MODEL

- Let  $\mathbf{s}$  be a subset of  $\{1, 2, \dots, p_n\}$ . We denote  $\mathbf{x}_n[\mathbf{s}]$  as the sub-vector of  $\mathbf{x}_n$  with elements in  $\mathbf{s}$ .
- Let  $\beta_{nk}[\mathbf{s}]$  be the subvector of  $\beta_{nk}$  with  $\beta_{kj} = 0, j \notin \mathbf{s}$ .
- For any  $\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_K$ , we get an FMR submodel:

$$f(y; \mathbf{x}, \boldsymbol{\Psi}_n, \mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_K) = \sum_{k=1}^K \pi_k h(y; \theta_k(\mathbf{x}_n[\mathbf{s}_k]), \phi_k)$$

with  $\theta_k(\mathbf{x}_n[\mathbf{s}_k]) = \mathcal{L}(\beta_{0k} + \mathbf{x}_n^\top[\mathbf{s}_k]\beta_{nk}[\mathbf{s}_k])$ .

- A feature selection method aims at selecting  $\mathbf{s}_k$ 's such that the resulting FMR submodel best balances the model parsimony and goodness of fit of the data.

- All-subset selection methods such as AIC or BIC are not feasible even for moderate values of  $p$  and  $K$ .

e.g. If  $K = 3$  and  $p = 20$ , there are:  $2^{20} \times 2^{20} \times 2^{20}$  possible submodels to be test by AIC or BIC.

- Khalili and Chen (2007) proposed a computationally efficient penalized likelihood approach for feature selection in FMR models.
- Statistical properties of their proposed method:
  - (1)  $\sqrt{n}$ -CONSISTENCY,
  - (2) ASYMPTOTIC NORMALITY,
  - (3) SPARSITY.
- In this work the dimension  $p$  of the feature space is assumed fixed with respect to the sample size  $n$ .

- We consider the problem of feature selection in FMR models when  $p = p_n$ , and  $p_n \rightarrow \infty$  as  $n \rightarrow \infty$ .
- Huber (1973), Fan and Peng (2004).



- Let  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$  be a random sample of observations from the FMR model.
- The (conditional) log-likelihood function is given by

$$l_n(\boldsymbol{\Psi}_n) = \sum_{i=1}^n \log\{f(y_i; \mathbf{x}_i, \boldsymbol{\Psi}_n)\}.$$

- Maximum likelihood estimation (MLE) of  $\boldsymbol{\Psi}_n$ :

$$\tilde{\boldsymbol{\Psi}}_n = \operatorname{argmax}_{\boldsymbol{\Psi}_n} l_n(\boldsymbol{\Psi}_n).$$

- But if  $\beta_{kj}^0 = 0$ , the MLE  $\tilde{\beta}_{kj}$  is not necessarily **zero** !

- We propose to estimate  $\Psi_n$  through the regularized log-likelihood function

$$pl_n(\Psi_n) = l_n(\Psi_n) - p_n(\Psi_n)$$

with the regularization function

$$p_n(\Psi_n) = \sum_{k=1}^K \pi_k \sum_{j=1}^{p_n} r_n(\beta_{kj}; \lambda_{nk}) + \frac{\tau_n}{2} \sum_{k=1}^K \sum_{j=1}^{p_n} \beta_{kj}^2.$$

## Why the penalty function $r_n(\beta; \lambda_{nk})$ ?

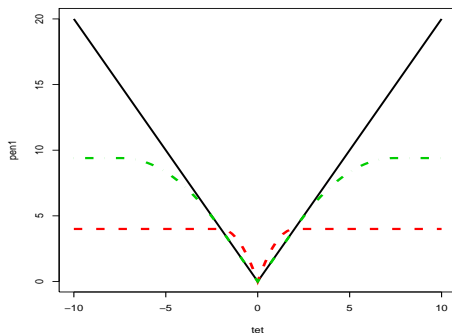
- For each  $k$ , let  $\beta_{nk}^0 = (\beta_{1k,n}^0, \beta_{2k,n}^0)^\tau$  such that  $\beta_{2k,n}^0$  contains the 0 effects.
- *Consistency in feature selection (sparsity):*

The penalty function  $r_n(\beta_{kj}; \lambda_{nk})$  is designed such that

$$P(\hat{\beta}_{2k,n} = 0) \rightarrow 1, \quad k = 1, 2, \dots, K, \quad \text{as } n \rightarrow \infty.$$

- By having an estimator with the above property, we are in fact performing *feature selection* !

- Most common choices are: LASSO, SCAD and HARD.



- *Controlling the variance of  $\hat{\beta}_{1k,n}$ :*

Similar to the ridge regression of [Hoerl and Kennard \(1970\)](#), the quadratic terms  $\beta_{kj}^2$  are to prevent wild estimates of the true nonzero  $\beta_{kj}^0$ 's corresponding to the highly correlated features  $x_j$ 's.

- See also [Hastie, Tibhsirani and Friedman \(2009\)](#).

- The performance of the new method is studied:
  - Theoretically:  
consistency, asymptotic normality, and sparsity.
  - An extensive simulation study.
  - The PD data is also analyzed to demonstrate the use of the new method in real applications.

- Let:  $\hat{\Psi}_n = \operatorname{argmax}_{\Psi_n} pl_n(\Psi_n)$
- Consider the joint density function  $f(\mathbf{z}_i; \Psi_n)$  of the random variable  $\mathbf{Z}_i = (\mathbf{x}_i, Y_i)$ .
- Consider the quantities:
  - $q_{1n} = \max_{k,j} \{r_n(\beta_{kj}^0; \lambda_{nk}) / \sqrt{n} : \beta_{kj}^0 \neq 0\}$
  - $q_{2n} = \max_{k,j} \{|r'_n(\beta_{kj}^0; \lambda_{nk})| / \sqrt{n} : \beta_{kj}^0 \neq 0\}$
  - $q_{3n} = \max_{k,j} \{|r''_n(\beta_{kj}^0; \lambda_{nk})| / n : \beta_{kj}^0 \neq 0\}$

( $\mathcal{P}_0$ ). For all  $n$  and  $\lambda_{nk}$ ,  $r_n(0; \lambda_{nk}) = 0$ , and  $r_n(\beta; \lambda_{nk})$  is symmetric and non-negative. It is non-decreasing and twice differentiable for all  $\beta$  in  $(0, \infty)$  with at most a few exceptions.

In addition, there exists constants  $A_1$  and  $A_2$  such that when  $\beta_1 > A_1 \lambda_{nk}$ ,  $\beta_2 > A_1 \lambda_{nk}$ , then:

$$\frac{1}{n} |r_n''(\beta_1; \lambda_{nk}) - r_n''(\beta_2; \lambda_{nk})| \leq A_2 |\beta_1 - \beta_2|.$$

( $\mathcal{P}_1$ ). As  $n \rightarrow \infty$ ,

$$\sqrt{p_n} q_{1n} = o(1 + q_{2n}), q_{2n} = o(p_n^{-1/2}), q_{3n} = o(1).$$



( $\mathcal{P}_2$ ). As  $n \rightarrow \infty$ ,

$$\frac{\tau_n}{\sqrt{n}} \max_{k,j} |\beta_{kj}^0| = o(1 + q_{2n}); \tau_n = o(n).$$

( $\mathcal{P}_3$ ). For  $T_n = \{\beta; 0 < \beta \leq \sqrt{\frac{p_n}{n}} \log n\}$ ,

$$\lim_{n \rightarrow \infty} \inf_{\beta \in T_n} \frac{r'_n(\beta; \lambda_{nk})}{\sqrt{np_n}} = \infty.$$

- Let  $\mathbf{Z}_i = (\mathbf{x}_i, Y_i)$ ,  $i = 1, 2, \dots, n$ , be a random sample from a density function  $f(\mathbf{z}; \Psi_n)$  that satisfies certain regularity Conditions.
- Assume that the function  $r_n(\beta; \lambda_{nk})$  satisfies Conditions  $\mathcal{P}_0$  and  $\mathcal{P}_1$ , and the tuning parameter  $\tau_n$  in ridge penalty satisfies Condition  $\mathcal{P}_2$ .

If  $\frac{p_n^2}{\sqrt{n}} \rightarrow 0$ , then there exists a local maximizer  $\hat{\Psi}_n$  of the function  $pl_n(\Psi_n)$  for which

$$\|\hat{\Psi}_n - \Psi_n^0\| = O_p\left\{\sqrt{\frac{p_n}{n}}(1 + q_{2n})\right\}.$$

## Theorem 2: Sparsity and asymptotic normality

- Assume conditions in [Theorem 1](#) are fulfilled, and also assume that the regularization function satisfy Conditions  $\mathcal{P}_0$ - $\mathcal{P}_3$ .

If  $\frac{p_n^{2.5}}{\sqrt{n}} \rightarrow 0$ , then for any  $\sqrt{n/p_n}$ -consistent maximum regularized likelihood estimator  $\hat{\Psi}_n$ , as  $n \rightarrow \infty$ :

- SPARSITY:**  $P(\hat{\beta}_{2k,n} = \mathbf{0}) \rightarrow 1, k = 1, 2, \dots, K$ .
- ASYMPTOTIC NORMALITY:**

$$\sqrt{n} \mathbf{B}_n \mathbf{I}_1^{-1/2}(\Psi_{01}) \left\{ \left[ \mathbf{I}_1(\Psi_{01}) - \frac{\mathbf{p}_n''(\Psi_{01})}{n} \right] (\hat{\Psi}_{n1} - \Psi_{01}) + \frac{\mathbf{p}_n'(\Psi_{01})}{n} \right\} \\ \longrightarrow^d N(\mathbf{0}, \mathbf{G})$$

where  $\mathbf{B}_n \mathbf{B}_n^\tau \rightarrow \mathbf{G}$ .

1. The estimator  $\hat{\Psi}_n$  is  $\sqrt{n/p_n}$ -consistent if  $q_{2n} = O(1)$ .
  - Proper choices of the tuning parameter  $\lambda_{nk}$  in **LASSO** and **SCAD** will lead to the desired consistency rate.
  - Regarding the ridge penalty, if for example  $\tau_n = \log n$ , then Condition  $\mathcal{P}_2$  is also guaranteed.
2. To have consistency in feature selection, the  $\lambda_{nk}$  needs to be chosen such that  $\sqrt{\frac{n}{p_n}} \lambda_{nk} \rightarrow \infty$ .
  - The properties 1 and 2 **cannot** be achieved simultaneously by **LASSO**. Achievable by **SCAD**!

- The data are the results of a clinical trial.
- Information on the speech signals can be extracted using both linear and nonlinear algorithms to characterize clinically relevant properties.
- Some features may address the ability of the voice folds to sustain simple vibration whereas some others may be able to characterize the extent of turbulent noise in the speech signals.
- $n = 5875$  observations on 16 speech characteristics, leading to a  $5875 \times 16$  design matrix.

- We fitted an FMR model with  $K = 2$ .
- The final selected model is

$$0.57 \phi(y; \hat{\mu}_1(\mathbf{x}), 4.36^2) + 0.43 \phi(y; \hat{\mu}_2(\mathbf{x}), 4.36^2)$$

$$\begin{aligned}\hat{\mu}_1(\mathbf{x}) &= 15.87 + \mathbf{x}^\tau[s_1]\hat{\beta}_1[s_1] \\ \hat{\mu}_2(\mathbf{x}) &= 28.70 + \mathbf{x}^\tau[s_2]\hat{\beta}_1[s_2].\end{aligned}$$

$$\begin{aligned}s_1 &= \{9, 10, 13, 15, 16\} \\ s_2 &= \{2, 3, 4, 16\}.\end{aligned}$$

## Fitted FMR model to the PD data

Voice features	Mixture components	
	Com <sub>1</sub>	Com <sub>2</sub>
Intercept	15.87 <sub>(.10)</sub>	28.70 <sub>(.12)</sub>
MDVP:Jitter (%)	—	—
MDVP:Jitter(Abs)	—	-3.65 <sub>(.13)</sub>
MDVP: RAP	—	3.77 <sub>(.12)</sub>
MDVP: PPQ	—	-1.85 <sub>(.13)</sub>
Jitter:DDP	—	—
MDVP:Shimmer	—	—
MDVP:Shimmer (dB)	—	—
Shimmer:APQ3	—	—
Shimmer:APQ5	-1.92 <sub>(.12)</sub>	—
MDVP: APQ	1.73 <sub>(.12)</sub>	—
Shimmer:DDA	—	—
NHR	—	—
HNR	-1.16 <sub>(.12)</sub>	—
RPDE	—	—
DFA	-1.89 <sub>(.10)</sub>	—
PPE	1.50 <sub>(.12)</sub>	2.16 <sub>(.12)</sub>

- For those patients whose PD symptoms are mild (smaller UPDRS scores), dysphonic measures such as **Shimmer:APQ5**, **MDVP:APQ**, **HNR**, **DFA**, and **PPE** are important features for monitoring disease progression.
- For those with more severe symptoms (larger UPDRS scores), a different set of features, **MDVP:Jitter (Abs)**, **MDVP:RAP**, **MDVP:PPQ**, appear to be more relevant in monitoring the progression of PD.



- Although all 16 features represent dysphonic measures, and multiple measures may characterize some similar aspect of the speech, still subtle difference in different  $x_j$ 's may have a profound implication in PD research.
- For instance, fundamental frequency variations (**MDVP:Jitter (Abs)**) and variations in signal amplitude (**Shimmer:APQ5**) both capture symptoms manifested in vocal fold vibration and lung efficiency, but the former is predictive of severity in advanced PD patients whereas the latter is predictive of severity in mild PD patients.
- Also our results suggest that UPDRS is affected by harmonics to noise ratio (**HNR**), but only in mild PD patients, with decreasing **HNR** leading to progression of the disease. However, as PD progresses to an advance stage, **HNR** may reach such a low level that it may no longer be useful for further severity prediction.

# References

- Huber, P. J. (Ann. Stat., 1973). Robust regression: Asymptotic, conjectures and Monte Carlo.
- Fan and Peng (Ann. Stat., 2004). Nonconcave penalized likelihood with a diverging number of parameters.
- Little, M.A., McSharry, P.E., Hunter, E.J., Spielman, J., Ramig, L.O. (2009). Suitability of dysphonia measurements for telemonitoring of Parkinson's disease. *IEEE Trans. Biomedical Engineering*, **56**, 1015-1022.
- Little, M.A., Costello, D.A.E., Harries, M.L. (2011). Objective dysphonia quantification in vocal fold paralysis: comparing nonlinear with classical measures. *Journal of Voice*, (in press).
- Khalili and Chen (JASA, 2007). Variable selection in finite mixture of regression models.
- Khalili, A. (Can. J. Statist, 2010). New estimation and feature selection methods in mixture-of-experts models.
- Khalili, Chen and Lin (Biostatistics, 2011). Feature selection in mixture of sparse normal linear models in high dimensional feature space.
- Khalili, A., Lin, S. (2011). Regularization in finite mixture of regression models with diverging number of parameters (submitted).
- Tsanas, A., Little, M. A., McSharry, P. E., Ramig, L.O. (2010). Accurate telemonitoring of Parkinson's disease progression by non-invasive speech tests. *IEEE Transactions on Biomedical Engineering*, **57**, 884-893.

Thank you for your attention!