
Local wavelet-vaguelette-based functional logistic regression for classification of gene expression data

M.M Rincón¹, M.D. Ruiz-Medina^{2,*}

¹ *mmrh@ugr.correo.es*

² *mrui@ugr.es*

**Corresponding author*

Abstract. *This paper focuses on the problem of functional statistical classification of gene expression curves. A local wavelet-vaguelette-based functional logistic regression approach is presented. This approach offers an alternative to the Functional-Principal-Component-Analysis-based logistic regression (see [4]). The performance of the methodology proposed is illustrated by implementing it for classification of yeast cell-cycle temporal gene expression data from [5] data set, where leave-one-out cross-validation error shows high accuracy of the model.*

Keywords. *Functional data; Functional logistic regression; Gene expression profile; Local wavelet-vaguelette decomposition; Yeast cell cycle gene expression data.*

1. INTRODUCTION

Functional wavelet bases have been widely used in the analysis of fractal biological signals, since their provide a localized multiscale decomposition of such signals. The wavelet transform of a random biological signal $\{X(t), t \in \mathbb{R}\}$ leads to a sequence of correlated random wavelet coefficients. To avoid redundancy in such coefficients a local version of the wavelet-vaguelette decomposition of a random signal is considered (see [1]), to obtain suitable response variables for a functional logistic regression, providing low-error rate classification for the yeast cell-cycle gene expression profiles analyzed.

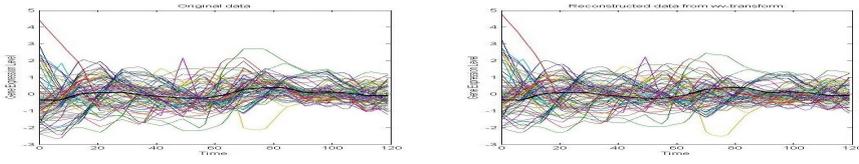


Figure 1: Left panel: Temporal gene expression profiles of yeast cell cycle. Right panel: Reconstruction of the temporal gene expression profiles in left panel from wavelet-vaguelette transform.

2. MODELS AND METHODS

The sample curves are assumed to be independent realizations of a mean-square integrable stochastic process $X(t)$ on $[0, S]$. Let $X_i(t_h)$ be the observation of the i th sample function at time t_h , for $h = 1, \dots, n$, and $i = 1, \dots, M$. Non-parametric kernel-based estimators, $\hat{\mu}(t)$ and $\hat{C}_X(s, t)$, are computed from a grid with $N = 2^p$, $p \in \mathbb{N}$, equally spaced points in $[0, S]$. Diagonal elements $\sigma^2(s) = \hat{C}_X(s, s)$, $s \in [0, S]$, are approximated by interpolated values $\hat{\sigma}^2(s)$, $s \in [0, S]$.

Multiresolution-like Analysis

The empirical eigenvalues $\hat{\lambda}_l$, $l = 1, \dots, N$, and the corresponding empirical eigenvectors $(\hat{\rho}_l(t_1), \dots, \hat{\rho}_l(t_N))$, $l = 1, \dots, N$, of the covariance estimate $\hat{C} = \hat{C}(t_l, t_m)$, $l, m = 1, \dots, N$, allow us to define the empirical kernel \hat{t}_X , factorizing the covariance function $\hat{C}_X(s, t)$, and the empirical kernel \hat{l}_X , approximating the inverse $L_X = \mathcal{T}_X^{-1}$ of operator \mathcal{T}_X , respectively as follows:

$$\hat{t}_X(t_h, t_m) = \sum_{l=1}^N \hat{\lambda}_l^{1/2} \hat{\rho}_l(t_h) \hat{\rho}_l(t_m), \quad \hat{l}_X(t_h, t_m) = \sum_{l=1}^N \hat{\lambda}_l^{-1/2} \hat{\rho}_l(t_h) \hat{\rho}_l(t_m). \quad (1)$$

for $h, m = 1, \dots, N$. The construction of the empirical wavelet-vaguelette functions is given in terms of kernels \hat{t}_X and \hat{l}_X , and a given orthonormal wavelet basis. We have chosen Haar system, with *the father wavelet*, $\phi(x) = I_{[0,1)}(x)$, and *the mother wavelet*, $\psi(x) = I_{[0,1/2)}(x) - I_{[1/2,1)}(x)$. Thus, (see [1],[6]), for $h = 1, \dots, N$,

$$\hat{\Phi}_0(t_h) = \sum_{m=1}^N \hat{t}_X(t_h, t_m) \phi(t_m), \quad \hat{\gamma}_{j,k}(t_h) = \sum_{m=1}^N \hat{t}_X(t_h, t_m) \psi_{j,k}(t_m), \quad k = 0, \dots, 2^j - 1, \quad j = 0, \dots, p-1.$$

In matrix form, we denote by $\Phi_0 = \{a_h\}$ the vector with entries $a_h = \hat{\Phi}_0(t_h)$, for $h = 1, \dots, N$, given by the product of the matrix $\hat{T} = \{b_{h,m}\}$, with $b_{h,m} = \hat{t}_X(t_h, t_m)$, for $h, m = 1, \dots, N$, and the vector $\Phi = \{c_m\}$, with $c_m = \phi(t_m)$, for $m = 1, \dots, N$. Similarly, for $j = 0, \dots, p-1$, the matrix $\Gamma_j = \{d_{h,k+1}\}$, with entries $d_{h,k+1} = \hat{\gamma}_{j,k}(t_h)$, for $h = 1, \dots, N$, and $k = 0, \dots, 2^j - 1$, is the product of matrices \hat{T} and Ψ_j , where $\Psi_j = \{l_{m,k+1}\}$ has entries $l_{m,k+1} = \psi_{j,k}(t_m)$, for $m = 1, \dots, N$, and $k = 0, \dots, 2^j - 1$. Additionally, we have $\Phi^0 = [\hat{T}^{-1}]^T \times \Phi$ and $\Gamma^j = [\hat{T}^{-1}]^T \times \Psi_j$.

For each sample curve X_i , evaluated at time $t \in [0, S]$, the following local empirical coefficients are computed:

$$\hat{X}_i^{\hat{\sigma}(t); \hat{\Phi}^0} = \sum_{m=1}^N (X_i(t_m) - \hat{\mu}(t_m)) \hat{\Phi}^{\hat{\sigma}(t); 0}(t_m) \quad \hat{X}_i^{\hat{\sigma}(t); j, k, \hat{\gamma}} = \sum_{m=1}^N (X_i(t_m) - \hat{\mu}(t_m)) \hat{\gamma}^{\hat{\sigma}(t); j, k}(t_m)$$

where $\{\hat{\Phi}^{\hat{\sigma}(t); 0}, \hat{\gamma}^{\hat{\sigma}(t); j, k}, k = 0, \dots, 2^j - 1, j = 0, \dots, p-1\} = \{\hat{\sigma}(t) \hat{\Phi}^0, \hat{\sigma}(t) \hat{\gamma}^{j, k}, k = 0, \dots, 2^j - 1, j = 0, \dots, p-1\}$ denotes the locally re-scaled empirical dual Riesz basis of $\{\hat{\Phi}_{\hat{\sigma}(t); 0}, \hat{\gamma}_{\hat{\sigma}(t); j, k}, k = 0, \dots, 2^j - 1, j = 0, \dots, p-1\} = \{(1/\hat{\sigma}(t)) \hat{\Phi}_0, (1/\hat{\sigma}(t)) \hat{\gamma}_{j, k}, k = 0, \dots, 2^j - 1, j = 0, \dots, p-1\}$. Note that

$$\langle \hat{\Phi}_{\hat{\sigma}(t); 0}, \hat{\Phi}^{\hat{\sigma}(t); 0} \rangle = 1, \quad \langle \hat{\gamma}_{\hat{\sigma}(t); j_1, k_1}, \hat{\gamma}^{\hat{\sigma}(t); j_2, k_2} \rangle = \delta_{j_1, j_2} \delta_{k_1, k_2}, \quad \langle \hat{\Phi}_{\hat{\sigma}(t); 0}, \hat{\gamma}^{\hat{\sigma}(t); j, k} \rangle = 0, \quad \langle \hat{\Phi}^{\hat{\sigma}(t); 0}, \hat{\gamma}_{\hat{\sigma}(t); j, k} \rangle = 0. \quad (2)$$

The M sample curves can then be approximated in terms of the following empirical local wavelet-vaguelette decomposition: For $i = 1, \dots, M$, and for each $t \in [0, S]$,

$$X_i(t) \simeq \hat{\mu} + \hat{X}_{\hat{\sigma}(t);i}^{\hat{\varphi}^0} \hat{\varphi}_{\hat{\sigma}(t);0}(t) + \sum_{j=0}^{p-1} \sum_{k=0}^{2^j-1} \hat{X}_{\hat{\sigma}(t);i}^{\hat{\gamma}^{j,k}} \hat{\gamma}_{\hat{\sigma}(t);j,k}(t), \quad t \in [0, S]. \quad (3)$$

This decomposition will be considered in the implementation of functional logistic regression to classify the data into two groups, G_0 and G_1 .

Functional Logistic Regression

Consider a response variable Y with Bernoulli distribution, having mean μ and variance $\sigma^2 = \mu(1 - \mu)$. The response variable Y takes the value $Y = 1$ if the sample curve is in group G_1 , or $Y = 0$ if it isn't. We define $\eta_i = g(\mu_i) = \alpha + \int \beta(t) Z_i(t) dt$, for α a constant, g the *logit* function, $g^{-1}(x) = e^x / (1 + e^x)$, and $Z_i(t) = X_i(t) - \hat{\mu}(t)$. Thus, $Y_i = g^{-1}(\eta_i) + e_i$, with errors $e_i, i = 1, \dots, M$, considered as independent random variables with zero-mean and finite variance.

Due to the square integrability of β , the functional parameter β admits the local decomposition:

$$\beta(t) = \beta_{\sigma(t); \varphi_0} \varphi^{\sigma(t); 0}(t) + \sum_{j=0}^{p-1} \sum_{k=0}^{2^j-1} \beta_{\sigma(t); j, k, \gamma} \gamma^{\sigma(t); j, k}(t), \quad t \in [0, S], \quad (4)$$

in terms of the dual local Riesz bases $\{\varphi_{\sigma(t);0}, \gamma_{\sigma(t);j,k}, k = 0, \dots, 2^j - 1, j = 0, \dots, p - 1\}$ and $\{\varphi^{\sigma(t);0,k}, \gamma^{\sigma(t);j,k}, k = 0, \dots, 2^j - 1, j = 0, \dots, p - 1\}$. In the development below, the local Fourier coefficients of parameter function β , with respect to the empirical scaled basis $\{(\varphi_{\hat{\sigma}(t);0}, \gamma_{\hat{\sigma}(t);j,k}, k = 0, \dots, 2^j - 1, j = 0, \dots, p - 1)\}$, will be denoted as $\hat{\beta}_{\varphi_0}^t = \beta_{\hat{\sigma}(t); \varphi_0}$, $\hat{\beta}_{j,k,\gamma}^t = \beta_{\hat{\sigma}(t); j, k, \gamma}$, $k = 0, \dots, 2^j - 1, j = 0, \dots, p - 1$, for each $t \in [0, S]$. The above approximations of $Z_i(t)$ from (3), and of $\beta(t)$ from (4), considering (2), lead to the following estimation of $\eta_i(t)$:

$$\eta_i(t) \simeq \hat{\alpha}^t + \sum_k \hat{Z}_i^{\hat{\sigma}(t); \varphi_0} \hat{\beta}_{\varphi_0}^t + \sum_{j=0}^{p-1} \sum_{k=0}^{2^j-1} \hat{Z}_i^{\hat{\sigma}(t); j, k, \gamma} \hat{\beta}_{j, k, \gamma}^t$$

The functional model is then reduced to a generalized linear model (see [2]), for each $t \in [0, S]$, where iterated weighted least square estimation is usually applied to compute $\hat{\beta}^t$, from the following equations: For $i = 1, \dots, M$,

$$\sum_i (Y_i - \mu_i(t)) = 0 \quad \sum_i (Y_i - \mu_i(t)) (\hat{Z}_i^t)^T = 0,$$

where $(\hat{Z}_i^t)^T$ is the vector of Fourier coefficients of $Z_i(t)$ on the empirical locally scaled wavelet-vaguelette basis $\{\hat{\varphi}_{\hat{\sigma}(t);0}, \hat{\gamma}_{\hat{\sigma}(t);j,k}, k = 0, \dots, 2^j - 1, j = 0, \dots, p - 1\}$.

The mean $\hat{\beta}$ over t of the obtained $\hat{\beta}^t$, for each $t = t_1, \dots, t_M$, is computed. A prior probability p_0 is considered for G_0 memberships, and similarly, a prior probability p_1 is considered for G_1 memberships. Thus, if $\hat{p}r(Y_i = 1 | X_i(t)) = \hat{\eta}_i = g^{-1}((\hat{Z}_i^t)^T, 1) * \hat{\beta} \geq p_1$, the i th curve is a member of G_1 . Otherwise, it belongs to G_0 .

3. RESULTS

Application to the analysis of yeast cell cycle gene expression profiles

We use the temporal gene expression data (α factor synchronized) for $M = 90$ genes involved in the yeast cell cycle obtained by [5] as sample curves. The gene expression is measured every 7 minutes between 0 and $S = 119$ minutes (both time instants included), thus, $n = 18$ observations for each gene. It is known that 44 of these genes are related to G_1 phase regulation and 46 to the $S, S/G_2, G_2/M$ and M/G_1 phases.

Figure (1) displays the original data with their approximation in terms of the local wavelet-vaguelette decomposition, considering a grid with $N = 64 = 2^6$ equally spaced time points. Convergence of the iterated weighted least squares algorithm is achieved for every point t in the grid after 100 iterations or less controlled by the deviance. $\hat{\beta}^t$ for a grid with $N = 64$ equally spaced time points, are displayed with the mean vector $\hat{\beta}$ in Figure 2. In order to measure the accuracy of the model, the cross-validation

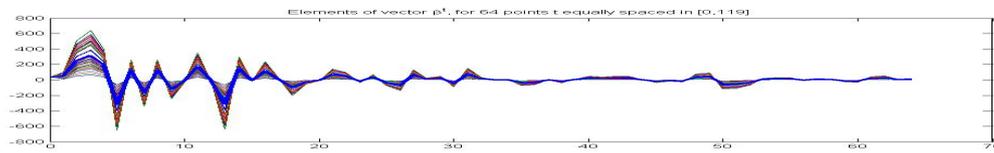


Figure 2: Components of $\hat{\beta}^t$ for 64 equally spaced points $t \in [0, 119]$, the coarsest line is the mean vector $\hat{\beta}$, note bigger influence of the first coefficients in the response variable.

classification error rate (*CVE*) is obtained. Suppose the i -th gene is missing, the mean and the covariance function estimates, based on the other 89 genes, and parameters $\hat{\beta}$, from the reduced functional sample, are then computed. These parameters are tested to obtain the approximation $\hat{\eta}^{-i}$ of $\hat{\eta}$, based on the sample information provided by the 89 gene expression curves, removing the i -th gene. This procedure is repeated with every gene, if $g^{-1}(\hat{\eta}^{-i}) \geq p_1$ the i -th gene is member of G_1 , otherwise is from G_0 . The *CVE* is defined as the quotient between the total number of genes misclassified under cross-validation, and the total number of genes. High accuracy of the model is assured by a $CVE = 0.13$.

4. CONCLUSIONS

In this paper, a local wavelet-vaguelette decomposition is considered for the non-redundant representation of gene expression profiles, since it holds for a large class of stochastic process, including processes with fractal and heavy-tailed covariance functions (see [3]).

Acknowledgments. This work has been supported in part by COLFUTURO, Colombia; projects MTM2009-13393 of the DGI, MEC, and P09-FQM-5052 of the Andalusian CICE, Spain.

References

- [1] Angulo, J.M. & Ruiz-Medina, M.D. (1999). Multiresolution approximation to the stochastic inverse problem. *Advanced Applied Probability* **31**, 1039–1057.
- [2] McCullagh, P. & Nelder J.A. (1989). *Generalized Linear Models*. Chapman & Hall
- [3] Kelbert, M., Leonenko, N.N. & Ruiz-Medina, M.D. (2005). Fractional random fields associated with stochastic fractional heat equations. *Advanced Applied Probability* **37**, 108–133.
- [4] Leng, X. & Müller, H.G. (2006). Classification using functional data analysis for temporal gene expression data. *Bioinformatics* **22** (1), 68–76. *Annals of Statistics* **33**, 774–805.
- [5] Spellman, P.T., Sherlock, G., Zhang, M.Q., Iyer, V.R., Anders, K., Eisen, M.B., Brown, P.O., Botstein, D. & Futcher, B. (1998). Comprehensive identification of cell-cycle regulated genes of the yeast *Saccharomyces cerevisiae* by microarray hybridization. *Molecular Biology of the Cell* **9**, 3273–3297.
- [6] Vidakovic, B. (2006) *Statistical Modelling by Wavelets*. John Wiley and sons. New York.