

Prediction of Pilot's Reaction Time Based on EEG Signals - Appendix

Bartosz Binias^{1*}, Dariusz Myszor², Henryk Palus¹ and Krzysztof A. Cyran³

¹Department of Data Mining and Engineering, Faculty of Automatic Control, Electronics and Computer Science, Silesian University of Technology, Gliwice, Poland

²Department of Department of Algorithmics and Software, Faculty of Automatic Control, Electronics and Computer Science, Silesian University of Technology, Gliwice, Poland

³Department of Computer Vision Graphics and Digital Systems, Faculty of Automatic Control, Electronics and Computer Science, Silesian University of Technology, Gliwice, Poland

Correspondence*:

Bartosz Binias

bartbinias@gmail.com

REFERENCES

- 2 Binias, B., Myszor, D., and Cyran, K. A. (2018). A machine learning approach to the detection of pilot's
3 reaction to unexpected events based on EEG signals. *Computational Intelligence and Neuroscience*
4 2018
- 5 Binias, B., Palus, H., and Niezabitowski, M. (2016). Elimination of bioelectrical source overlapping effects
6 from the EEG measurements. In *Carpathian Control Conference (ICCC), 2016 17th International*
7 (IEEE), 70–75
- 8 Efron, B., Hastie, T., Johnstone, I., Tibshirani, R., et al. (2004). Least angle regression. *The Annals of*
9 *Statistics* 32, 407–499
- 10 Friedman, J., Hastie, T., and Tibshirani, R. (2010). Regularization paths for generalized linear models via
11 coordinate descent. *Journal of Statistical Software* 33, 1
- 12 Kraskov, A., Stögbauer, H., and Grassberger, P. (2004). Estimating mutual information. *Physical Review E*
13 69, 066138
- 14 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., et al. (2011). Scikit-learn:
15 Machine Learning in Python. *Journal of Machine Learning Research* 12, 2825–2830
- 16 Robert, C. (2014). Machine learning, a probabilistic perspective. *CHANCE* 27, 62–63. doi:10.1080/
17 09332480.2014.914768
- 18 Ross, B. C. (2014). Mutual information between discrete and continuous data sets. *PloS One* 9, e87357
- 19 Smola, A. J. and Schölkopf, B. (2004). A tutorial on Support Vector Regression. *Statistics and Computing*
20 14, 199–222

A DESCRIPTION OF STATISTICAL AND MACHINE LEARNING METHODS

21 A.1 Feature standardization and mean removal

22 Feature standardization and mean removal i.e., normalization, is a common practice in many machine
23 learning approaches. The procedure is performed by removing the mean of the feature vector and scaling

24 it to unit variance. This procedure is then applied to each feature vector independently by computing the
 25 relevant statistics on the basis of samples from the training set. Mean and standard deviation are then stored
 26 for use on an independent test data.

27 If we denote the mean value of training samples from feature f by $\mu(X_f)$ and their standard deviation as
 28 $\sigma(X_f)$, then for a sample x_{fi} ($i = 1, \dots, M_f$, M_f - is a number of samples of feature f), the described
 29 transformation can be calculated using the following equation:

$$z_{fi} = \frac{x_{fi} - \mu(X_f)}{\sigma(X_f)} \quad (1)$$

30 In this research study f denotes the bandpower features calculated from 10 frequency bands defined in
 31 Sec. 2.3.3.

32 A.2 Mutual Information

33 Mutual Information (MI) between two random variables is a non-negative value that describes the
 34 dependency between these variables. MI is equal to zero if and only if two random variables are independent.
 35 At the same time, higher values of the MI score indicate a higher dependency. If $X_f \in \mathbf{R}^{M_f}$ denotes the
 36 vector of M_f samples of feature f and $Y \in \mathbf{R}^{M_f}$ is a vector of corresponding targets, then the MI of these
 37 two discrete variables can be defined as in Eq. 2. In this study, MIs for discrete variables were obtained
 38 with nonparametric methods based on entropy estimation from k-nearest neighbors distances (Kraskov
 39 et al., 2004; Ross, 2014).

$$MI(X_f; Y) = \sum_{y \in Y} \sum_{x_f \in X_f} p(x_f, y) \log \left(\frac{p(x_f, y)}{p(x_f)p(y)} \right) \quad (2)$$

40 In Eq. 2, $p(x_f, y)$ is the joint probability function of X_f and Y , and $p(x_f)$ and $p(y)$ are the marginal
 41 probability distribution functions of X_f and Y respectively. In this study, x_f is a vector of 10 bandpower
 42 features obtained for a single event, while y is the corresponding time of delay in reaction to that event.

43 The MI criterion is known for being capable of capturing any kind of dependency between variables.
 44 Use of MI-based feature selection methods have been proven to yield highly satisfactory results in many
 45 approaches to EEG signal processing (Binias et al., 2016, 2018).

46 A.3 F-regression

47 F-test statistics can be used as a criterion for ranking features. This approach utilizes univariate linear
 48 regression for testing the individual effect of the regression variables. To extract this information, the first
 49 step requires that the correlation between the vector of regressors $X_f \in \mathbf{R}^{M_f}$ and the vector of targets
 50 $Y \in \mathbf{R}^{M_f}$ is computed, according to the following equation:

$$R_f^2 = \frac{(X_f - \mu(X_f))^T (Y - \mu(Y))}{\sigma(X_f)\sigma(Y)} \quad (3)$$

51 The R_f^2 is then converted to an F-score to obtain the final result. If we denote the number of observations
 52 as M_f and the degrees of freedom as p_f , then the relation between the F-score F_f and R_f^2 is expressed as
 53 in Eq.4.

$$R_f^2 = 1 - \left(1 + F_f \frac{p_f - 1}{M_f - p_f}\right)^{-1} \quad (4)$$

54 It must be noted that the F-test expresses only a linear dependency between variables. In this study, x_f is
 55 a vector of 10 bandpower features obtained for a single event, while y is the corresponding time of delay in
 56 reaction to that event.

57 **A.4 Least Absolute Shrinkage and Selection Operator**

58 Least Absolute Shrinkage and Selection Operator (LASSO) is a linear model that estimates sparse
 59 coefficients. Mathematically, the optimization objective for trained linear model is the $L1$ norm regularizer
 60 defined by the following equation (Friedman et al., 2010):

$$L1 = \min_w \frac{1}{2N} \|Xw - Y\|_2^2 + \alpha \|w\|_1 \quad (5)$$

61 where:

- 62 • $X \in \mathbf{R}^{M \times N}$ - input data (bandpower features),
- 63 • $Y \in \mathbf{R}^M$ - target (vector of reaction times),
- 64 • $\|w\|_1$ - L1-norm of the parameter vector,
- 65 • α - constant,
- 66 • M - number of samples,
- 67 • N - number of features (10 bandpower features were being used in this study).

68 The implementation of the LASSO used in this work was taken from the Python library *scikit-learn* and
 69 uses the coordinate descent as the algorithm to fit the coefficients (Pedregosa et al., 2011).

70 **A.4.1 LASSO with Least-Angle Regression**

71 Least Absolute Shrinkage and Selection Operator with Least-Angle Regression (LASSO-LARS) is a
 72 LASSO model implemented using the LARS algorithm rather than the coordinate descent *scikit-learn*.
 73 LARS is a regression algorithm that is similar to the forward stepwise regression (Efron et al., 2004).
 74 Although its detailed description is beyond the scope of this article, some most important features of LARS
 75 will be listed in this section. The algorithm has numerous advantages over the classical implementation of
 76 LASSO. One of the most important advantages is the numeric efficiency for high-dimensional data with a
 77 relatively small sample size. Additionally, LARS is fast in terms of computation time and has proven to be
 78 more stable. On the other hand, the LARS algorithm may be particularly sensitive to noise. Since EEG data
 79 can be considered noisy by nature, this might have a crucial impact on the effectiveness of LASSO-LARS
 80 in this study.

81 **A.5 Ridge Regression with Radial Kernel**

82 Ridge Regression with Radial Kernel (KernelRidge) is a combination of a linear least squares with $L2$
 83 norm regularization and kernel transformation (Robert, 2014). The $L2$ can be defined as presented in Eq.6.

$$L2 = \min_w \|Xw - Y\|_2^2 + \alpha \|w\|_2^2 \quad (6)$$

84 where:

- 85 • $X \in \mathbf{R}^{M \times N}$ - input data (bandpower features),
- 86 • $Y \in \mathbf{R}^M$ - target (vector of reaction times),
- 87 • $\|w\|_1$ - L2-norm of the parameter vector,
- 88 • α - complexity parameter that controls the amount of shrinkage,
- 89 • M - number of samples,
- 90 • N - number of features (10 bandpower features were being used in this study).

91 In this study, a Radial Basis Function (RBF) was used for kernel transformation. The RBF for a feature
 92 vector $X_f \in \mathbf{R}^M$ is defined as presented in Eq.7.

$$RBF = \exp(-\gamma \|X_f - X'_f\|^2) \quad (7)$$

93 A.6 Support Vector Machine with Radial Basis Function

94 Support Vector Machine (SVM) is a supervised learning method that can be used for classification and
 95 regression problems. The mathematical formulation of SVM for regression problems can be found below
 96 (Smola and Schölkopf, 2004).

97 Let's denote the total number of features by N and a number of observations by M . Given training vectors
 98 $X_i \in \mathbf{R}^N$, $i = 1, \dots, M$ and a target vector $Y \in \mathbf{R}^M$, SVM solves the following regression problem:

$$\begin{aligned} \min_{w, b, \zeta, \zeta^*} & \frac{1}{2} w^T w + C \sum_{i=1}^M (\zeta_i + \zeta_i^*) \\ & Y_i - w^T \phi(X_i) - b \leq \varepsilon + \zeta_i, \\ & w^T \phi(X_i) + b - Y_i \leq \varepsilon + \zeta_i^*, \\ & \zeta_i, \zeta_i^* \geq 0, i = 1, \dots, M \end{aligned} \quad (8)$$

99 which is dual to:

$$\min_{\alpha, \alpha^*} \frac{1}{2} (\alpha - \alpha^*)^T Q (\alpha - \alpha^*) + \varepsilon e^T (\alpha + \alpha^*) - Y^T (\alpha - \alpha^*) \quad (9)$$

100 subject to

$$\begin{aligned} e^T (\alpha - \alpha^*) &= 0 \\ 0 &\leq \alpha_i, \alpha_i^* \leq C, i = 1, \dots, M \end{aligned} \quad (10)$$

101 where:

- 102 • e is the vector of all ones,
- 103 • $C > 0$ is the upper bound,
- 104 • $Q \in \mathbf{R}^{M \times M}$,

105 • $Q_{ij} \equiv K(X_i, X_j) = \phi(X_i)^T \phi(X_j)$ is the kernel function.

106 The decision function, with independent term ρ is presented in following equation:

$$\sum_{i=1}^M (\alpha_i - \alpha_i^*) K(X_i, X) + \rho \quad (11)$$

107 SVM algorithms support multiple kernel functions for input data transformation. These functions are
 108 particularly useful when dealing with complex problems that have many more features than observations.
 109 Since this is the case for the problem targeted in this study, a SVM with a RBF kernel (SVM-RBF) was
 110 used instead of a linear SVM. The RBF for a feature vector X_f is defined in Eq. 7.

B HYPERPARAMETERS OF REGRESSION METHODS

111 This section presents the range of hyperparameters that were used to optimize performance of the selected
 112 machine learning algorithms. For a detailed description of each of the presented hyperparameters, please
 113 refer to the documentation of the *scikit-learn* library (Pedregosa et al., 2011).

114 B.1 LASSO:

- 115 • $\epsilon = 0.001$ - length of the regularization path defined as $\frac{\alpha_{min}}{\alpha_{max}}$.
- 116 • α - the amount of penalization chosen based on minimizing cross-validated generalization error
 117 (method built-in to *scikit-learn* implementation).
- 118 • $tol = 0.0001$ - the tolerance for the optimization.
- 119 • Maximum number of iterations to perform was $1e6$.
- 120 • Coefficients were selected cyclically for the update every iteration.
- 121 • The interception point for the model was being calculated for the computations (i.e. data was not
 122 expected to be centered).

123 B.2 LASSO-LARS

- 124 • $\epsilon = 2e - 16$ - The machine-precision regularization in the computation of the Cholesky diagonal
 125 factors.
- 126 • α - the amount of penalization chosen based on minimizing cross-validated generalization error
 127 (method built-in to *scikit-learn* implementation).
- 128 • $tol = 0.0001$ - the tolerance for the optimization.
- 129 • Maximum number of iteration to perform was $1e5$.
- 130 • The maximum number of points (α) on the path used to compute the residuals in the cross-validation
 131 was 1000.
- 132 • The interception point for the model was being calculated for the computations (i.e. data was not
 133 expected to be centered).

134 B.3 KernelRidge

- 135 • Before training a subset of best features was selected.
- 136 • The criteria for feature selection was either the F -score or MI . The criteria that best suited each
 137 dataset was treated as a tuned hyperparameter.
- 138 • The number of best features that would be used was selected from the set $\{1, 2, \dots, 30\}$.

- 139 • $\alpha \in \left\{ \alpha_{min} + x \frac{\alpha_{max} - \alpha_{min}}{N_\alpha - 1} \mid x \in \{0, 1, \dots, N_\alpha - 1\}, \alpha_{min} = -1, \alpha_{max} = 10, N_\alpha = 100 \right\}$ -
 140 regularization strength term in L2 norm.
- 141 • $\gamma \in \left\{ \gamma_{min} + x \frac{\gamma_{max} - \gamma_{min}}{N_\gamma - 1} \mid x \in \{0, 1, \dots, N_\gamma - 1\}, \gamma_{min} = 10^{-3}, \gamma_{max} = 1, N_\gamma = 100 \right\}$ - gamma
 142 parameter for the RBF.

143 B.4 SVMRBF

- 144 • The same feature selection procedure as presented in B.3 was utilized.
- 145 • Shrinking was always enabled during the computations.
- 146 • $C \in \left\{ C_{min} + x \frac{C_{max} - C_{min}}{N_C - 1} \mid x \in \{0, 1, \dots, N_C - 1\}, C_{min} = 10^{-3}, C_{max} = 10^3, N_C = 100 \right\}$ -
 147 penalty parameter of the error term.
- 148 • $\gamma = 1/N_f$ where N_f denotes the number of features - kernel coefficient for RBF.