

1 Pattern Recognition

Pattern recognition is a technique whose goal is classification and the principle by which it is performed is learned independently from the data, i.e., training set. There are two main types of pattern recognition: supervised and unsupervised. Supervised pattern recognition implies that the classes of the training set are known and are used to obtain the model. New inputs are identified as one of the predetermined classes. On the other hand, unsupervised pattern recognition is used when no labels are available and samples are assigned to unknown classes. This technique is more appropriate for the clustering problem because the classes are determined automatically by the system, whereas supervised approach is more appropriate for the classification because classes are defined by the system designer.

In statistical pattern recognition, each sample is composed of m measures that form the pattern, i.e., features $(x_0, x_1, \dots, x_{m-1})$. The goal of the algorithm is to obtain a decision rule, i.e., the decision boundary which separates well samples of different classes. There are many *state-of-the-art* classifiers that use various principles to construct these boundaries. However, many researchers agree that the fidelity of the classification in EMG applications depends mostly on selection of features. In other words, with appropriate selection of features, all classifiers will give similar classification result. Since pattern recognition is not the topic of this thesis, only a short introduction is provided on the two methods used in the thesis: Linear Discriminant Analysis and Support Vector Machine.

1.1 Linear Discriminant Analysis

All models are wrong; some models are useful.
George E. P. Box

Linear Discriminant Analysis is a computationally simple and efficient classifier with linear decision boundary and it is based on the Bayesian equation.

In a classical problem with n samples in training set, which consist of m features, the dataset of available samples is a matrix $\mathbf{x}_{[n \times m]}$, whereas the vector of labels that describe the belonging of each sample to one of the classes is \mathbf{y} , where $y \in (0, 1, 2, \dots, K - 1)$.

According to Bayesian equation, the probability that a sample \mathbf{x}_0 belongs to a class k is equivalent to the:

$$P(y = k | \mathbf{x} = \mathbf{x}_0) = \frac{P(\mathbf{x} = \mathbf{x}_0 | y = k) P(y = k)}{P(\mathbf{x} = \mathbf{x}_0)} \quad (1)$$

,where k represents the class. Term $P(\mathbf{x} = \mathbf{x}_0 | y = k)$ is called the *class-conditional* probability and describes the probability that the sample with exact features \mathbf{x}_0 is encountered within the group of samples belonging only to the class k . Term $P(y = k)$ is called the *a priori* probability and describes the probability that the sample with class k is found within the group of all samples, regardless of the features. Finally, the term $P(\mathbf{x} = \mathbf{x}_0)$ is called *marginal* probability and describes the probability of finding the exact set of features in the dataset, regardless of the class. Marginal probability can be written as a sum of class-conditional probabilities multiplied by the a priori probabilities for each class:

$$\begin{aligned} P(\mathbf{x} = \mathbf{x}_0) &= P(\mathbf{x} = \mathbf{x}_0 | y = 1) P(y = 1) + \\ &P(\mathbf{x} = \mathbf{x}_0 | y = 2) P(y = 2) + \dots + \\ &P(\mathbf{x} = \mathbf{x}_0 | y = K) P(y = K) \end{aligned} \quad (2)$$

Following Bayesian theory, the hypothesis, i.e., the predicted class of a sample \mathbf{x}_0 is chosen as the class which has the highest probability $P(y = k | \mathbf{x} = \mathbf{x}_0)$:

$$h(\mathbf{x}_0) = \underset{k}{\operatorname{argmax}} P(y = k | \mathbf{x} = \mathbf{x}_0) \quad (3)$$

Statistically speaking, this is the best possible classifier. The problem arises in implementation. The exact probability density functions are unknown and have to be estimate from the available data, which is the source of error. Estimated version of the stated probabilities will be marked with a different symbols to stress out the fact they are just an estimates:

$$p_k(\mathbf{x}) := P(y = k \mid \mathbf{x}) \quad (4)$$

$$g_k(\mathbf{x}) := P(\mathbf{x} \mid y = k) \quad (5)$$

$$\pi_k := P(y = k) \quad (6)$$

Linear Discriminant Analysis estimates marginal probability term (π_k) as ratio of number of samples belonging to class k and the total number of samples, whereas class-conditional probability term in the Bayesian equation is estimated as a Gaussian function:

$$g_k(\mathbf{x}) = \frac{1}{(2\pi)^{1/2} |\Sigma_k|^{1/2}} e^{-1/2(\mathbf{x}-\mu_k)^T \Sigma_k^{-1} (\mathbf{x}-\mu_k)} \quad (7)$$

, where g_k is estimated class-conditional probability of class k , and μ_k and Σ_k are the mean and co-variance matrix for class k , respectively, and they are estimated from the available data as:

$$\mu_k = \frac{1}{n_k} \sum_i \mathbf{x}_i \Big|_{\forall \mathbf{x} \in k} \quad (8)$$

$$\Sigma_k = \frac{1}{n_k - K} \sum_i (\mathbf{x}_i - \mu_k) (\mathbf{x}_i - \mu_k)^T \Big|_{\forall \mathbf{x} \in k} \quad (9)$$

where n_k represents the number of samples belonging to a class k . To simplify the model, LDA assumes that the co-variance matrices Σ_k are the same for all classes:

$$\Sigma_0 = \Sigma_1 = \dots = \Sigma_{K-1} = \Sigma \quad (10)$$

and they are usually calculated using weighted average:

$$\Sigma = \frac{\sum_{k=1}^K n_k \Sigma_k}{\sum_{k=1}^K n_k} \quad (11)$$

The consequence of this assumption is the linearity of the decision boundary.

In a two class example ($y \in \{0, 1\}$), all samples on the decision boundary ($D.B.$) will have the same probability of belonging to class 0 or 1:

$$D.B. = \left\{ \mathbf{x} \mid P(y = 0 \mid \mathbf{x} = \mathbf{x}_0) = P(y = 1 \mid \mathbf{x} = \mathbf{x}_0) \right\} \quad (12)$$

Following this idea, the decision boundary can be estimated by solving the equation:

$$\frac{g_0(\mathbf{x}) \pi_0}{\sum_{k=1}^K g_k \pi_k} = \frac{g_1(\mathbf{x}) \pi_1}{\sum_{k=1}^K g_k \pi_k} \quad (13)$$

$$\frac{1}{(2\pi)^{d/2} |\Sigma_0|^{1/2}} e^{-1/2(\mathbf{x}-\mu_0)^T \Sigma_0^{-1} (\mathbf{x}-\mu_0)} \pi_0 = \frac{1}{(2\pi)^{d/2} |\Sigma_1|^{1/2}} e^{-1/2(\mathbf{x}-\mu_1)^T \Sigma_1^{-1} (\mathbf{x}-\mu_1)} \pi_1 \quad (14)$$

, where d is the dimensionality of the feature space, i.e., number of features representing each sample. If making the assumption on the equal co-variance matrices for both classes:

$$\Sigma_0 = \Sigma_1 = \Sigma \quad (15)$$

and taking the logarithm, the equation takes the form:

$$\begin{aligned}
& -\frac{1}{2}(\mathbf{x} - \mu_0)^T \Sigma^{-1}(\mathbf{x} - \mu_0) + \log(\pi_0) = \\
& -\frac{1}{2}(\mathbf{x} - \mu_1)^T \Sigma^{-1}(\mathbf{x} - \mu_1) + \log(\pi_1)
\end{aligned} \tag{16}$$

This equation can be written as the linear function $\mathbf{x}^T \beta + \alpha = 0$ as:

$$\mathbf{x}^T (\Sigma^{-1} \mu_0 - \Sigma^{-1} \mu_1) + \frac{1}{2} (\mu_1^T \Sigma^{-1} \mu_1 - \mu_0^T \Sigma^{-1} \mu_0) + \log\left(\frac{\pi_0}{\pi_1}\right) = 0 \tag{17}$$

The equation represents the decision boundary between classes, i.e., all samples lying on this line will have equal probability of belonging to class 0 and class 1. It is interesting to note that the slope of the line depends only on the class means and co-variance matrix, whereas a priori probabilities (which are the result of number of samples belonging to class 0 or 1) have effect only on the y -intercept term, i.e., the offset of the function. This is an interesting point that demands caution. If groups are unbalanced, that is, number of samples of one group is higher than in the other group, y -intercept of the decision boundary will be affected and the classifier will be biased by this disproportion.

When considering multiclass classification problem, probability of a sample belonging to each class is firstly estimated by the equation:

$$p_k = -\frac{1}{2} \log |\Sigma| - \frac{1}{2} (\mathbf{x} - \mu_k)^T \Sigma^{-1} (\mathbf{x} - \mu_k) + \log(\pi_k) \tag{18}$$

and then the class is estimated as the one with the highest probability as:

$$h(\mathbf{x}) = \operatorname{argmax}_k p_k(\mathbf{x}) \tag{19}$$

1.2 Support Vector Machine

Try to solve the problem directly and never solve a more general problem as an intermediate step.

Vladimir Vapnik

Support vector machine is nowadays known as a very powerful classifier with a lot of different applications. The big advantage over LDA is the fact that it is a *non-parametric* classifier. The model is not obtained using assumption of the form of the class density function and estimation of its parameters, which is inevitably erroneous. Instead, SVM forms the decision boundary using the samples (not their density estimates) by maximizing the distance between samples and the boundary. This was the idea Vladimir Vapnik, the inventor of this method stood for. It is better to try to solve the problem directly and simply, without many intermediate steps that can be complicated and inaccurate.

This classification method essentially has a very simple principle, but using several mathematical tricks, it became a very powerful tool. The optimization problem does not depend on \mathbf{x} , but on $\mathbf{x}^T \mathbf{x}$. This enables the use of *kernel trick* and allows nonlinear transform of the feature space at no additional cost.

In pattern recognition, the decision rule (h) is usually obtained by multiplying the sample (\mathbf{x}) by predefined weights (Θ):

$$\Theta^T \mathbf{x} + \Theta_0 \tag{20}$$

, where Θ_0 is a constant.

If samples \mathbf{x}_0 and \mathbf{x}_1 lay on the decision boundary $D.B.$, following statements are true:

$$\Theta^T \mathbf{x}_0 + \Theta_0 = \Theta^T \mathbf{x}_1 + \Theta_0 \tag{21}$$

$$\Theta^T (\mathbf{x}_0 - \mathbf{x}_1) = 0 \tag{22}$$

This result implies that Θ is perpendicular to $D.B.$.

The goal of the SVM is to find the decision boundary so that the distance between the $D.B.$ and samples, i.e., the margin is maximized. The distance (d) from a sample to the decision boundary can be defined as the distance between the sample \mathbf{x} and any point on the boundary, \mathbf{x}_0 , projected onto the vector Θ .

$$d = \frac{\Theta^T (\mathbf{x} - \mathbf{x}_0)}{|\Theta|} \quad (23)$$

Term $|\Theta|$ is introduced to normalize the vector Θ

Since \mathbf{x}_0 is on the decision boundary, the expression $\Theta^T \mathbf{x}_0 + \Theta_0 = 0$ is valid, and, therefore, the expression for the distance can be written as:

$$d = \frac{\Theta^T \mathbf{x} + \Theta_0}{|\Theta|} \quad (24)$$

Margin (M) is the distance from the boundary to the closest samples:

$$M = \min_i d_i \quad (25)$$

Depending on which side of the boundary the sample is located, distance can be positive or negative. In order to keep the strictly positive, term y is introduced, where $y \in -1, 1$:

$$M = \min_i \{y_i d_i\} \quad (26)$$

$$M = \min_i \left\{ \frac{y_i (\Theta^T \mathbf{x} + \Theta_0)}{|\Theta|} \right\} \quad (27)$$

The goal is to maximize the margin. Since Θ can be rescaled, a certain Θ exists so that $y_i (\Theta^T \mathbf{x} + \Theta_0) = 0$, which implies

$$\exists \Theta, y_i (\Theta^T \mathbf{x} + \Theta_0) = 1 \Rightarrow M = \min_i \left\{ \frac{1}{|\Theta|} \right\} \quad (28)$$

Therefore, to maximize the margin, a hyperplane should be found such that a norm of vector orthogonal to the hyperplane (Θ) is minimal. L_2 norm is preferred because it has continuous derivative

For every point not on the boundary the following term is valid:

$$y_i (\Theta^T \mathbf{x} + \Theta_0) > 0 \quad (29)$$

Value C can be selected such that:

$$y_i (\Theta^T \mathbf{x} + \Theta_0) > C \quad (30)$$

$$y_i \left(\frac{\Theta^T \mathbf{x}}{C} + \frac{\Theta_0}{C} \right) > 1 \quad (31)$$

Since Θ and Θ_0 can be rescaled, it can be written:

$$\Theta := \frac{\Theta}{C}, \quad \Theta_0 := \frac{\Theta_0}{C} \quad (32)$$

$$y_i (\Theta^T \mathbf{x} + \Theta_0) > 1 \quad (33)$$

Finally the optimization problem states:

$$\min \frac{1}{2} |\Theta|^2, \quad s.t. \quad y_i (\Theta^T \mathbf{x} + \Theta_0) > 1 \quad (34)$$

and it can be solved using the Lagrange polynomial as:

$$L(\Theta, \Theta_0, \alpha) = \frac{1}{2} |\Theta|^2 - \sum_{i=1}^n [y_i (\Theta^T \mathbf{x} + \Theta_0) - 1] \quad (35)$$

$$\frac{\partial L}{\partial \Theta} = \Theta - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i = 0 \Rightarrow \Theta = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \quad (36)$$

$$\frac{\partial L}{\partial \Theta} = \sum_{i=1}^n \alpha_i y_i = 0 \quad (37)$$

By enlisting this terms into the equation HABA,

$$L(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_j \sum_i \alpha_j \alpha_i y_j y_i \mathbf{x}_i^T \mathbf{x}_j \quad (38)$$