

Supervised Approach to Learning Multivariate Linear Systems

Luminița State, Iuliana Paraschiv-Munteanu

Abstract

We consider the problem of developing a learning from data scheme for the unknown input-output dependency of a system \mathbf{S} of linear type in the sense that the m -dimensional output of \mathbf{S} results by combining in a linear way the effects of n observable variables and the effects of several unobservable latent variables. The effects of the latent variables on the output is treated as additive noise, that is being given the observable vector x , the system computes the output $y = \beta^T \begin{pmatrix} 1 \\ x \end{pmatrix} + \varepsilon$, where β is a $(n+1) \times m$ matrix and ε is a m -dimensional Gaussian variable. In the paper the mathematical arguments for the estimation scheme based exclusively on a finite size set of observations is provided. We present an experimental evaluation of the quality of the resulted learning scheme in order to establish conclusions concerning their accuracy and generalization capacities, the evaluation being performed in terms of metric, probabilistic and informational criterion functions.

1 Introduction

The tremendous growth in practical applications of machine learning over the past decade has been accompanied by a wide variety of important developments in the underlying algorithms and techniques that make use of concepts and results coming from several areas as mathematical statistics, computer science and engineering.

Since the main aim of machine learning is to obtain computer programs that are able to extract information from samples of data and as well as knowledge from the past experience and include them in the process of solving problems of high complexity, the research methodology in the field of machine learning is essentially based on a large class of concepts and results coming from mathematical statistics, neural and evolutionary computation, brain models, adaptive control theory and so on.

The aim of the research was to develop a model free learning methodology in order to predict a system behavior, conventionally denoted by \mathbf{S} on the basis of finite size sequence of observations. In real life applications, the input data sequence are either obtained in a controlled way, that is the observer knows in advance the generating mechanism, or, in an uncontrolled way when the generating mechanism is ignored by the observer. We denote by \mathbf{G} (Generator) the component that supplies a series of samples from an n -dimensional space loaded as inputs to us. The learning environment is assumed to be of supervised type, that is the output of \mathbf{S} is available to the observer. The goal is to develop a learning component \mathbf{L} on the basis of a finite size set of input-output observations, that is to infer the unknown the input-output dependency of \mathbf{S} and use it for further predictions. The general scheme of our model is presented in Figure 1 ([3]).

The learning component \mathbf{L} is responsible for a class of possible models for the unknown dependency corresponding to \mathbf{S} . In other words, the learning component \mathbf{L} implements a class of hypothesis (models) Ω , such that to each particular hypothesis $\omega \in \Omega$ corresponds a function $\varphi_\omega : \mathcal{X} \rightarrow \mathcal{Y}$ defined on the space

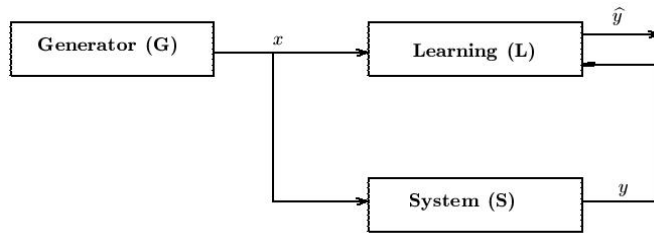


Figure 1: The learning environment.

of inputs \mathcal{X} and taking values in the space of outputs \mathcal{Y} . For each particular input x_0 , $\hat{y}_0 = \varphi_\omega(x_0)$ is the estimate of the \mathbf{S} 's output corresponding to x_0 in case of the model ω . Being given a criterion function \mathcal{C} that expresses numerically the fitness of each model with respect to the available evidence E , about \mathbf{S} , the best model $\omega_0(E)$ is a solution of the optimization problem

$$\arg(\text{optimize}_{\omega \in \Omega} \mathcal{C}(\omega, E)) . \quad (1)$$

In the case of supervised learning the available evidence E is represented by a finite set of pairs $\{(x_i, y_i), 1 \leq i \leq N\} \subset \mathcal{X} \times \mathcal{Y}$, where each y_i is the actual output of \mathbf{S} for the input x_i . If we assume that the unknown dependency is of deterministic type, that is the inputs and the outputs of \mathbf{S} are functionally related a reasonable choice of the criterion function \mathcal{C} is the arithmetic mean of the square errors, that is for each $\omega \in \Omega$,

$$\mathcal{C}(\omega, E) = \frac{1}{N} \sum_{i=1}^N \|y_i - \varphi_\omega(x_i)\|^2 . \quad (2)$$

The optimization problem (1) becomes

$$\arg\left(\min_{\omega \in \Omega} \mathcal{C}(\omega, E)\right) , \quad (3)$$

and its solutions are called the Minimum Square Errors (MSE) models computed on the basis of $\{(x_i, y_i), 1 \leq i \leq N\}$ ([14]).

A more realistic approach is to accept that besides the inputs, the outputs of \mathbf{S} are also influenced by a series of unknown number of unobservable factors referred as latent variables. In this case the cumulated effects of the latent variables can be only modelled in probabilistic terms, assuming some class of multivariate probability distributions, each hypothesis corresponding to a certain input-output dependency combined with a probabilistic model for the latent vector. For simplicity sake, we consider that the latent vector is a continuous random vector, that is to each hypothesis $\omega \in \Omega$ corresponds a conditional density function $f(\cdot|\cdot; \omega)$ ([9], [11]). Put in other words, for each $\omega \in \Omega$, $x \in \mathcal{X}$, $y \in \mathcal{Y}$, $f(y|x; \omega)$ expresses 'the chance' of getting the output y for the input x in case of the model ω . If the available evidence about \mathbf{S} is $\{(x_i, y_i), 1 \leq i \leq N\}$ then a reasonable choice of $\mathcal{C}(\omega, E)$ is the likelihood function. If we assume that the inputs x_1, \dots, x_N are independently generated by \mathbf{G} then

$$\mathcal{C}(\omega, E) = \prod_{i=1}^N f(y_i|x_i; \omega) , \quad (4)$$

and the optimization problem (1) becomes

$$\arg\left(\max_{\omega \in \Omega} \mathcal{C}(\omega, E)\right) . \quad (5)$$

The solutions of (5) are the MLE (Maximum Likelihood Estimation) models computed on the basis of $\{(x_i, y_i), 1 \leq i \leq N\}$.

2 A MLE-based approach in learning multivariate linear systems

In our work we assume that the inputs come from a n -dimensional real space and for each input x the output y of \mathbf{S} is a tuple of m real variables. It is also assumed that the output dependency on the input is of linear type, that is the influence of the input is given by $\beta^T x$ where β is an unknown $(n+1) \times m$ matrix. The effects of the latent variables can be taken as noise additively superimposed on the influence of the inputs, that is the output of \mathbf{S} is $y = \beta^T \begin{pmatrix} 1 \\ x \end{pmatrix} + \varepsilon$ where ε is a m -dimensional Gaussian distributed random vector, $\varepsilon \sim \mathcal{N}(\mu, \Sigma)$. The parameters $\mu \in \mathbb{R}^m$ and $\Sigma \in \mathcal{M}_m(\mathbb{R})$, symmetric positive definite matrix, are unknown. Therefore, the dimension of the space of hypotheses Ω is $m(m+n+2)$, and for each particular tuple of $m(m+n+2)$ parameters, $\omega = (\beta, \mu, \Sigma)$ defines a model of \mathbf{S} . We denote by $f(y|x; \omega)$ the conditional density function on the output space corresponding to the model ω .

Several intuitively justified criterion functions can be considered in order to quantitatively express the quality of each possible model with respect to available sequence of observations. For each criterion function, the identification of the "fittest" model reduces to solving an constrained/unconstrained optimization problem ([6]).

Let $\mathcal{S}_N = \{(x_i, y_i), 1 \leq i \leq N\}$ be the sequence of input-output observations taken on \mathbf{S} . Being given the model $\omega = (\beta, \mu, \Sigma)$, the estimate \tilde{y} of the actual output of \mathbf{S} being given the input x is a random vector of density function

$$f(\tilde{y}|x, \omega) = \frac{1}{\sqrt{(2\pi)^m |\Sigma|}} \exp \left\{ -\frac{1}{2} (\tilde{y} - \beta^T z - \mu)^T \Sigma^{-1} (\tilde{y} - \beta^T z - \mu) \right\},$$

where $z = \begin{pmatrix} 1 \\ x \end{pmatrix}$.

Consequently, the log-likelihood function is

$$l(\beta, \mu, \Sigma, \mathcal{S}_N) = -\frac{mN}{2} \ln(2\pi) - \frac{N}{2} \ln |\Sigma| - \frac{1}{2} \sum_{i=1}^N \left((y_i - \mu)^T - z_i^T \beta \right) \Sigma^{-1} \left((y_i - \mu) - \beta^T z_i \right),$$

and the best model from the point of view of maximum likelihood principle is a solution of the constrained optimization problem

$$\begin{cases} \min_{\beta, \mu, \Sigma} (\Phi(\beta, \mu, \Sigma)) \\ \Sigma \in \mathcal{M}_m(\mathbb{R}) \text{ symmetric and positive defined,} \end{cases}$$

where

$$\Phi(\beta, \mu, \Sigma) = N \ln |\Sigma| + \sum_{i=1}^N \left((y_i - \mu)^T - z_i^T \beta \right) \Sigma^{-1} \left((y_i - \mu) - \beta^T z_i \right).$$

Let us denote

$$\begin{aligned} Y &= (y_1, \dots, y_N) \in \mathcal{M}_{m \times N}(\mathbb{R}), \quad Z = (z_1, \dots, z_N) \in \mathcal{M}_{(n+1) \times N}(\mathbb{R}), \\ u &= (1, \dots, 1)^T \in \mathbb{R}^N, \quad A = \mathbf{I}_N - \frac{1}{N} uu^T. \end{aligned} \quad (6)$$

For any matrix B , we denote by B^+ the generalized inverse (Penrose pseudo-inverse) of B .

Theorem 1 *The objective function $\Phi(\beta, \mu, \Sigma)$ has an unique critical point $(\beta_0, \mu_0, \Sigma_0)$ where*

$$\beta_0 = \left(Y (ZA)^+ \right)^T, \quad \mu_0 = \frac{1}{N} \left(Y u - Y (ZA)^+ Z u \right), \quad \Sigma_0 = \frac{1}{N} Y \left(A - (ZA)^+ (ZA) \right) Y^T, \quad (7)$$

and Σ_0 is a symmetric and positive semi-defined matrix.

Proof. The generalized gradients of $\Phi(\beta, \mu, \Sigma)$ with respect to β , μ , and Σ respectively, are

$$\begin{cases} \nabla_{\beta} \Phi(\beta, \mu, \Sigma) = -(ZY^T - Zu\mu^T - ZZ^T\beta)\Sigma^{-1}, \\ \nabla_{\mu} \Phi(\beta, \mu, \Sigma) = -\Sigma^{-1}(Yu - N\mu - \beta^T Zu), \\ \nabla_{\Sigma} \Phi(\beta, \mu, \Sigma) = \Sigma^{-1}D\Sigma^{-1} - \frac{1}{2} \text{diag}(\Sigma^{-1}D\Sigma^{-1}), \end{cases}$$

where

$$D = N\Sigma - YY^T + (\mu(Yu)^T + (Yu)\mu^T) + ((ZY^T)^T\beta + \beta^T ZY^T) - N(\mu\mu^T) - (\mu(Zu)^T\beta + \beta^T(Zu)\mu^T) - \beta^T ZZ^T\beta.$$

From the system

$$\begin{cases} \nabla_{\mu} \Phi(\beta, \mu, \Sigma) = \mathbf{0}_m, \\ \nabla_{\beta} \Phi(\beta, \mu, \Sigma) = \mathbf{0}_{n+1, m}, \end{cases}$$

since $|\Sigma| \neq 0$, we get

$$\begin{cases} Yu - N\mu - \beta^T Zu = \mathbf{0}_m \\ ZY^T - Zu\mu^T - ZZ^T\beta = \mathbf{0}_{n+1, m}, \end{cases}$$

that is

$$\beta = (Y(ZA)^+)^T = \beta_0, \quad \mu = \frac{1}{N} (Yu - Y(ZA)^+Zu) = \mu_0.$$

Replacing μ_0, β_0 in the system

$$\nabla_{\Sigma} \Phi(\beta, \mu, \Sigma) = \mathbf{0}_m$$

we obtain

$$\Sigma^{-1}D\Sigma^{-1} - \frac{1}{2} \text{diag}(\Sigma^{-1}D\Sigma^{-1}) = \mathbf{0}_m,$$

where $\text{diag}(\Sigma^{-1}D\Sigma^{-1}) \in \mathcal{M}_m(\mathbf{R})$ is the diagonal matrix that retains only the entries placed on the main diagonal of $\Sigma^{-1}D\Sigma^{-1}$. Since Σ is a positive definite matrix, we get $D = \mathbf{0}_m$ and consequently,

$$\begin{aligned} \Sigma &= \frac{1}{N} (YY^T + \beta_0^T ZZ^T \beta_0 - YZ^T \beta_0 - \beta_0^T ZY^T + \\ &\quad \frac{1}{N} (\beta_0^T Zuu^T Y^T - \beta_0^T Zuu^T Z^T \beta_0 + Yuu^T Z^T \beta_0 - Yuu^T Y^T)) = \\ &\quad \frac{1}{N} (YAY^T + \beta_0^T ZAZ^T \beta_0 - YAZ^T \beta_0 - \beta_0^T ZAY^T). \end{aligned}$$

Using the well-known properties of the Penrose pseudo-inverse, the expression of Σ becomes

$$\begin{aligned} \Sigma &= \frac{1}{N} \left(YAY^T + Y(ZA)^+(ZA)AZ^T(Y(ZA)^+)^T - YAZ^T(Y(ZA)^+)^T - Y(ZA)^+ZAY^T \right) = \\ &\quad \frac{1}{N} \left(YAY^T + Y((ZA)^+(ZA))^T(ZA)^T((ZA)^+)^T Y^T - Y((ZA)^+(ZA))^T Y^T - Y(ZA)^+(ZA)Y^T \right) = \\ &\quad \frac{1}{N} Y \left(A + ((ZA)^+(ZA)(ZA)^+(ZA))^T - (ZA)^+(ZA) - (ZA)^+(ZA) \right) Y^T = \\ &\quad \frac{1}{N} Y \left(A - (ZA)^+(ZA) \right) Y^T \stackrel{\text{not}}{=} \Sigma_0. \end{aligned}$$

A further simplification can be obtained by noting that $B = A - (ZA)^+ZA$ is a symmetric matrix and $B^2 = B$. Which finally yields to

$$\Sigma_0 = \frac{1}{N} YBB^T Y^T.$$

Obviously, the estimate Σ_0 is a symmetric and positive semi-defined matrix. \square

Theorem 2 Let β_0 and μ_0 be given by Theorem 1. Then for any (β, μ, Σ) in the parameter space,

$$l(\beta, \mu, \Sigma, \mathcal{S}_N) \leq l(\beta_0, \mu_0, \Sigma, \mathcal{S}_N). \quad (8)$$

Proof. Using (6) and $v^T v = \text{tr}(vv^T)$ the expression of the log-likelihood function can be written as

$$l(\beta, \mu, \Sigma, \mathcal{S}_N) = -\frac{mN}{2} \ln(2\pi) - \frac{N}{2} \ln |\Sigma| - \frac{1}{2} \text{tr} \left(\Sigma^{-1} (Y - \mu u^T - \beta^T Z) (Y - \mu u^T - \beta^T Z)^T \right).$$

Therefore

$$l(\beta_0, \mu_0, \Sigma, \mathcal{S}_N) = -\frac{mN}{2} \ln(2\pi) - \frac{N}{2} \ln |\Sigma| - \frac{1}{2} \text{tr} \left(\Sigma^{-1} Y (A - (ZA)^+ ZA) (A - (ZA)^+ ZA)^T Y^T \right).$$

Using the relations $A = A^2 = A^T$ we get

$$\begin{aligned} & (A - (ZA)^+ ZA) (A - (ZA)^+ ZA)^T = \\ & A^2 - A(ZA)^+(ZA) - (ZA)^+(ZA)A + (ZA)^+(ZA)(ZA)^+(ZA) = \\ & A - A(ZA)^+(ZA) - (ZA)^+(ZA) + (ZA)^+(ZA) = A - A(ZA)^+(ZA). \end{aligned}$$

therefore the term $Y - \mu u^T - \beta^T Z$ becomes

$$\begin{aligned} Y - \mu u^T - \beta^T Z &= Y (A - (ZA)^+ (ZA)) + \frac{1}{N} (Y u - Y (ZA)^+ Z u - N \mu) u^T + (Y (ZA)^+ - \beta^T) Z = \\ & Y (A - (ZA)^+ (ZA)) + (\mu_0 - \mu) u^T + (\beta_0 - \beta)^T Z. \end{aligned}$$

Consequently,

$$\begin{aligned} l(\beta, \mu, \Sigma, \mathcal{S}_N) &= -\frac{mN}{2} \ln(2\pi) - \frac{N}{2} \ln |\Sigma| - \frac{1}{2} \text{tr} \left(\Sigma^{-1} Y (A - (ZA)^+ ZA) (A - (ZA)^+ ZA)^T Y^T \right) - \\ & \text{tr} \left(\Sigma^{-1} Y (A - (ZA)^+ ZA) ((\mu_0 - \mu) u^T)^T \right) - \text{tr} \left(\Sigma^{-1} Y (A - (ZA)^+ ZA) ((\beta_0 - \beta)^T Z)^T \right) - \\ & \frac{1}{2} \text{tr} \left(\Sigma^{-1} (\mu_0 - \mu) u^T u (\mu_0 - \mu)^T \right) - \frac{1}{2} \text{tr} \left(\Sigma^{-1} (\beta_0 - \beta)^T Z Z^T (\beta_0 - \beta) \right) = \end{aligned} \quad (9)$$

$$\begin{aligned} & l(\beta_0, \mu_0, \Sigma, \mathcal{S}_N) - \frac{1}{2} \text{tr} \left(\Sigma^{-1} (\mu_0 - \mu) u^T u (\mu_0 - \mu)^T \right) - \frac{1}{2} \text{tr} \left((\Sigma^{-1} \beta_0 - \beta)^T Z Z^T (\beta_0 - \beta) \right) - \\ & \text{tr} \left(\Sigma^{-1} Y (\mathbf{I}_N - (ZA)^+ Z) A u (\mu_0 - \mu)^T \right) - \text{tr} \left(\Sigma^{-1} Y (\mathbf{I}_N - (ZA)^+ Z) A Z^T (\beta_0 - \beta) \right). \end{aligned}$$

Since

$$A u = \left(\mathbf{I}_N - \frac{1}{N} u u^T \right) u = u - \frac{1}{N} u (u^T u) = \mathbf{0}_N$$

we get

$$Y (\mathbf{I}_N - (ZA)^+ Z) A u (\mu_0 - \mu)^T = \mathbf{0}_m.$$

Also, using the properties of the Penrose pseudo-inverse, we get

$$\begin{aligned} (\mathbf{I}_N - (ZA)^+ Z) A Z^T &= (ZA)^T - (ZA)^+ Z A A^T Z^T = (ZA)^T - ((ZA)^+ (ZA))^T (ZA)^T = \\ & (ZA)^T - ((ZA)(ZA)^+(ZA))^T = (ZA)^T - (ZA)^T = \mathbf{0}_{N, n+1}, \end{aligned}$$

that is $Y (\mathbf{I}_N - (ZA)^+ Z) A Z^T (\beta_0 - \beta) = \mathbf{0}_m$.

Taking into account these arguments we finally obtain,

$$l(\beta, \mu, \Sigma, \mathcal{S}_N) = l(\beta_0, \mu_0, \Sigma, \mathcal{S}_N) - \frac{1}{2} \text{tr} \left(\Sigma^{-1} (\mu_0 - \mu) u^T u (\mu_0 - \mu)^T \right) - \frac{1}{2} \text{tr} \left(\Sigma^{-1} (\beta_0 - \beta)^T Z Z^T (\beta_0 - \beta) \right).$$

Obviously, since Σ a positive definite matrix,

$$\text{tr} \left(\Sigma^{-1} (\mu_0 - \mu) u^T u (\mu_0 - \mu)^T \right) = N \text{tr} \left((\mu_0 - \mu)^T \Sigma^{-1} (\mu_0 - \mu) \right) = N (\mu_0 - \mu)^T \Sigma^{-1} (\mu_0 - \mu) \geq 0$$

and

$$\text{tr} \left(\Sigma^{-1} (\beta_0 - \beta)^T Z Z^T (\beta_0 - \beta) \right) = \text{tr} \left(Z^T (\beta_0 - \beta) \Sigma^{-1} (\beta_0 - \beta)^T Z \right) = \text{tr} \left(\left((\beta_0 - \beta)^T Z \right)^T \Sigma^{-1} (\beta_0 - \beta)^T Z \right) \geq 0$$

that is

$$l(\beta, \mu, \Sigma, \mathcal{S}_N) \leq l(\beta_0, \mu_0, \Sigma, \mathcal{S}_N). \quad \square$$

Remark. Although, a long series of tests pointed out that the estimate Σ_0 given by Theorem 1 is a positive matrix the mathematical proof is still an open problem. Also, it is not known whether the unique critical point $(\beta_0, \mu_0, \Sigma_0)$ corresponds to the best model in the sense of the maximum likelihood principle.

An adaptive learning procedure can be obtained using the gradient ascent method applied to the log-likelihood criterion function. The search developed by the adaptive procedure in a $m(m+n+2)$ -dimensional space aims to adjust the model parameters β, μ, Σ in order to maximize the log-likelihood function or, equivalently, to minimize $\Phi(\beta, \mu, \Sigma)$. The procedure should be implemented using a control parameter $\delta > 0$ and a stopping condition $\mathcal{C}(\delta)$ usually expressed in terms of the magnitude of the displacement in the parameter space due to the current iteration. In our tests $\mathcal{C}(\delta) = \text{true}$ if

$$\|\beta^{\text{new}} - \beta^{\text{old}}\| < \delta, \quad \|\mu^{\text{new}} - \mu^{\text{old}}\| < \delta, \quad \|\Sigma^{\text{new}} - \Sigma^{\text{old}}\| < \delta$$

where $\|\cdot\|$ is a conventionally norm, for instance Euclidian norm.

Also, the implementation of the procedure can be done using either a constant learning rate $\rho > 0$ or a decreasing sequence of positive learning rates (ρ_k) that refines the search while the search advances.

Since during the search process the estimates of Σ are not guaranteed to be invertible, the implementation procedure `MLE_gradient_ascent` uses approximations of the actual generalized gradients where the generalized inverse is used instead.

`procedure MLE_gradient_ascent`

Input: $\{(x_1, y_1), \dots, (x_N, y_N)\}$

Initializations: $\delta > 0, \rho > 0, \tilde{\beta}, \tilde{\mu}, \tilde{\Sigma},$

$$Z = \begin{pmatrix} 1 & & 1 \\ x_1 & \dots & x_N \end{pmatrix}, \quad Y = (y_1, \dots, y_N), \quad u = (1, \dots, 1)^T$$

$$\beta^{\text{old}} \leftarrow \tilde{\beta}, \quad \mu^{\text{old}} \leftarrow \tilde{\mu}, \quad \Sigma^{\text{old}} \leftarrow \tilde{\Sigma}$$

Compute $S = ZZ^T, Q = ZY^T, P = YY^T$

$$Z_1 = Zu, \quad Y_1 = Yu$$

repeat

$$\Sigma_1 \leftarrow (\Sigma^{\text{old}})^+$$

$$\beta^{\text{new}} \leftarrow \beta^{\text{old}} + \rho \left(Q - Z_1 (\mu^{\text{old}})^T - S \beta^{\text{old}} \right) \Sigma_1$$

$$\mu^{\text{new}} \leftarrow \mu^{\text{old}} + \rho \Sigma_1 \left(Y_1 - (\beta^{\text{old}})^T Z_1 - N \mu^{\text{old}} \right)$$

$$D = N \Sigma^{\text{old}} - P + Y_1 (\mu^{\text{old}})^T + \left(Y_1 (\mu^{\text{old}})^T \right)^T + Q^T \beta^{\text{old}} + (Q^T \beta^{\text{old}})^T -$$

$$N \mu^{\text{old}} (\mu^{\text{old}})^T - \mu^{\text{old}} (Z_1)^T \beta^{\text{old}} - \left(\mu^{\text{old}} (Z_1)^T \beta^{\text{old}} \right)^T - (\beta^{\text{old}})^T S \beta^{\text{old}}$$

$$\Sigma^{\text{new}} \leftarrow \Sigma^{\text{old}} + \rho \left(-\Sigma_1 D \Sigma_1 + \frac{1}{2} \text{diag}(\Sigma_1 D \Sigma_1) \right)$$

evaluate $\mathcal{C}(\delta)$

$$\beta^{\text{old}} \leftarrow \beta^{\text{new}}, \quad \mu^{\text{old}} \leftarrow \mu^{\text{new}}, \quad \Sigma^{\text{old}} \leftarrow \Sigma^{\text{new}}$$

until $\mathcal{C}(\delta)$

Output: $\beta^{\text{new}}, \mu^{\text{new}}, \Sigma^{\text{new}}.$

3 Experimental analysis

Because the model $(\beta_0, \mu_0, \Sigma_0)$ given by (7) is not theoretical guaranteed as the best model from the point of view of maximum likelihood principle, we have performed a long series of tests aiming to derive conclusions concerning the performance of the proposed method on experimental way. The test examples x_i 's were random generated from n -dimensional Gaussian repartition $\mathcal{N}(\mu_1, \Sigma_1)$. The target responses y_i 's were computed as $y_i = \tilde{\beta}^T x_i + \varepsilon$ for given β and ε randomly generated from known Gaussian repartition $\mathcal{N}(\tilde{\mu}, \tilde{\Sigma})$.

According to the previous arguments the expression of conditional density function on the output space corresponding to each example x_i being given the model $\omega = (\beta, \mu, \Sigma)$ is

$$f(y|x_i, \omega) = \frac{1}{\sqrt{(2\pi)^m |\Sigma|}} \exp \left\{ -\frac{1}{2} (y - \beta^T z_i - \mu)^T \Sigma^{-1} (y - \beta^T z_i - \mu) \right\},$$

where $z_i = \begin{pmatrix} 1 \\ x_i \end{pmatrix}$, therefore the most likely output predicted value is $y'_i = \beta^T z_i + \mu$.

In order to evaluate the quality of the resulted model we considered three indicators to evaluate the overall error.

The first indicator evaluates the overall mean error of miss-prediction for the given set of example $\{(x_i, y_i) \mid 1 \leq i \leq N\}$ corresponding to each possible model ω

$$error_1 = \frac{1}{N} \sum_{i=1}^N (1 - f(y_i|x_i, \omega)). \quad (10)$$

The second indicator is a mean error computed in terms of the actual responses and the most likely predicted values,

$$error_2 = \frac{1}{N} \sum_{i=1}^N \|y_i - y'_i\|^2 = \frac{1}{N} \sum_{i=1}^N \|y_i - \beta^T z_i - \mu\|^2. \quad (11)$$

The third measure, of informational type, aims to evaluate the informational correlation between the input and the computed output corresponding to each model, and it is expressed in terms of the empirical mutual information. Since x_1, \dots, x_N are randomly generated $\mathcal{N}(\mu_1, \Sigma_1)$, the probability distribution $\tilde{p} = (\tilde{p}(x_1), \dots, \tilde{p}(x_N))$ characterizes the collection of examples,

$$\tilde{p}(x_j) = \frac{p(x_j)}{\sum_{i=1}^N p(x_i)},$$

where $p(x_j) = \frac{1}{\sqrt{(2\pi)^n |\Sigma_1|}} \exp \left\{ -\frac{1}{2} (x_j - \mu_1)^T \Sigma_1^{-1} (x_j - \mu_1) \right\}$, $1 \leq j \leq N$.

The empirical entropy of the input samples x_1, \dots, x_N is given by the Shannon entropy corresponding to \tilde{p} ,

$$H(\tilde{p}) = - \sum_{i=1}^N \tilde{p}(x_i) \ln \tilde{p}(x_i). \quad (12)$$

Using the transition probabilities

$$\tilde{p}(y_j|x_i, \omega) = \frac{f(y_j|x_i, \omega)}{\sum_{k=1}^N f(y_k|x_i, \omega)}, \quad 1 \leq i, j \leq N,$$

we define the probability distribution $\tilde{q} = (\tilde{q}(y_1), \dots, \tilde{q}(y_N))$ on the set of target responses by

$\tilde{q}(y_j) = \sum_{i=1}^N \tilde{p}(x_i) \tilde{p}(y_j|x_i, \omega)$, $1 \leq j \leq N$, and let

$$H(\tilde{q}) = - \sum_{i=1}^N \tilde{q}(y_i) \ln \tilde{q}(y_i), \quad (13)$$

be the empirical Shannon entropy of the set of \mathbf{S} 's outputs.

Using the well-known expression of the mutual information ([4]), the empirical mutual information can be defined as

$$\mathcal{I}(\mathcal{S}_N) = H(\tilde{q}) - \sum_{i=1}^N \sum_{j=1}^N \tilde{p}(x_i) \tilde{p}(y_j | x_i, \omega) \ln \tilde{p}(y_j | x_i, \omega). \quad (14)$$

We evaluated the performance of the model $\omega_0 = (\beta_0, \mu_0, \Sigma_0)$ given by (7) for different sizes for learning samples \mathcal{S}_N .

Test 1. The settings are $n = 2$, $m = 1$, $\mu_1 = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$, $\Sigma_1 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$, $\tilde{\beta} = \begin{pmatrix} 0 \\ 2 \\ 4 \end{pmatrix}$, $\tilde{\mu} = 0.25$ and $\tilde{\Sigma} = 1$. Some of the results for data of different sizes N , corresponding to the quasi-optimal model ω_0 computed for each data set are summarized in Table 1a.

Table 1: Model evaluation in case $n = 2$, $m = 1$ and for different volumes of: a) learning data; b) test samples.

a						b					
N	$error_1$	$error_2$	$H(\tilde{p})$	$H(\tilde{q})$	$\mathcal{I}(\mathcal{S}_N)$	N_{test}	$error_1$	$error_2$	$H(\tilde{p})$	$H(\tilde{q})$	$\mathcal{I}(\mathcal{S}_N)$
10	0.603	0.533	2.214	2.205	1.335	15	0.706	1.149	2.458	2.582	1.346
20	0.702	0.802	2.767	2.796	1.304	30	0.710	1.257	3.243	3.276	1.367
50	0.711	0.892	3.768	3.831	1.226	50	0.683	1.039	3.752	3.829	1.573
100	0.712	0.950	4.477	4.524	1.233	70	0.650	1.136	4.146	4.172	1.573
200	0.716	0.959	5.098	5.188	1.221	90	0.707	1.292	4.309	4.400	1.637
300	0.721	1.019	5.510	5.590	1.106	100	0.709	1.276	4.376	4.483	1.662
400	0.713	0.953	5.822	5.897	1.155	200	0.715	1.361	5.124	5.235	1.616
500	0.725	1.025	6.023	6.113	1.150	300	0.719	1.375	5.547	5.633	1.579

Several tests aimed to establish conclusions concerning the generalization capacity of the quasi-optimal model. For instance, in case of a training sequence of volume $N = 100$, the computed quasi-optimal model is $\omega_0 = (\beta_0, \mu_0, \Sigma_0)$ is

$$\beta_0 = \begin{pmatrix} 0.000 \\ 2.614 \\ 4.077 \end{pmatrix} \quad \mu_0 = -0.218, \quad \Sigma_0 = 0.457,$$

the results are summarized in table 1b and the variations of the input, output empirical entropies and the empirical mutual information as functions of the test sample sizes are depicted in Figure 2b.

Test 2. The settings are $n = 2$, $m = 2$, $\mu_1 = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$, $\Sigma_1 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$, $\tilde{\beta} = \begin{pmatrix} 0 & 0 \\ 2 & 1 \\ 4 & 5 \end{pmatrix}$, $\tilde{\mu} = \begin{pmatrix} 0.25 \\ 0.25 \end{pmatrix}$ and $\tilde{\Sigma} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$. Some of the results for data of different sizes N , corresponding to the quasi-optimal model ω_0 computed for each data set are summarized in Table 2a.

Several tests aimed to establish conclusions concerning the generalization capacity of the quasi-optimal model. For instance, in case of a training sequence of volume $N = 100$, the computed quasi-optimal model is $\omega_0 = (\beta_0, \mu_0, \Sigma_0)$ is

$$\beta_0 = \begin{pmatrix} -0.000 & -0.000 \\ 1.872 & 1.012 \\ 4.134 & 5.047 \end{pmatrix} \quad \mu_0 = \begin{pmatrix} 0.274 \\ 0.104 \end{pmatrix}, \quad \Sigma_0 = \begin{pmatrix} 0.677 & -0.007 \\ -0.007 & 0.899 \end{pmatrix},$$

the results are summarized in table 2b and the variations of the input, output empirical entropies and the empirical mutual information as functions of the test sample sizes are depicted in Figure 3b.

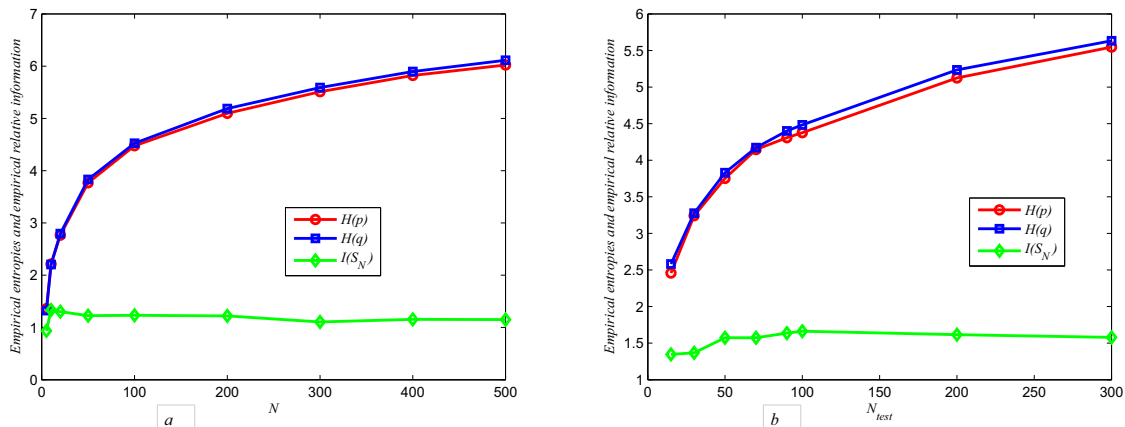


Figure 2: The variation of the empirical input/output entropies and empirical mutual information: a) as function of N , volume of learning data; b) for $N = 100$ and different volumes of test samples.

Table 2: Model evaluation in case $n = 2$, $m = 2$ and for different volumes of: a) learning data; b) test samples.

a						b					
N	$error_1$	$error_2$	$H(\hat{p})$	$H(\hat{q})$	$\mathcal{I}(S_N)$	N_{test}	$error_1$	$error_2$	$H(\hat{p})$	$H(\hat{q})$	$\mathcal{I}(S_N)$
15	0.887	1.436	2.500	2.483	1.610	15	0.917	2.154	2.542	2.311	1.277
20	0.877	1.282	2.816	2.720	1.521	30	0.912	1.874	3.190	3.116	1.763
50	0.892	1.473	3.748	3.675	1.810	50	0.890	1.438	3.764	3.699	1.627
100	0.926	2.030	4.441	4.369	1.583	70	0.910	2.216	4.042	3.889	1.772
200	0.908	1.774	5.090	5.065	1.704	90	0.907	1.881	4.349	4.277	1.769
300	0.923	2.072	5.513	5.495	1.638	100	0.912	2.014	4.432	4.292	1.783
400	0.919	1.982	5.810	5.768	1.661	200	0.906	1.938	5.098	5.062	1.698
500	0.918	1.954	6.014	5.987	1.763	300	0.902	1.860	5.512	5.457	1.822

4 Conclusions

The aim of the report research was to develop a model free learning methodology in order to predict a system behavior, conventionally denoted by \mathbf{S} on the basis of finite size sequence of observations. The learning environment is assumed to be of supervised type, that is a finite size sample of pairs of input-output values is available to the observer. The goal was to develop a learning component on the basis on the input-output sample, that is to infer the unknown the input-output dependency of \mathbf{S} and use it for further predictions. In Section 2, the input-output dependency of \mathbf{S} is modeled in terms of a Gaussian repartition of the output space. The model parameters are adjusted on the basis of the available sequence of observations using the Maximum Likelihood Principle. We managed to compute explicitly the quasi-optimal model ω_0 (Theorem 1) and establish its 'almost' optimality (Theorem 2).

The performance evaluation of the quasi-optimal model is evaluated in Section 3 aiming to establish conclusions concerning the quality of the predictions about \mathbf{S} computed by the learning component. The performance is evaluated in terms of three indicators, on the basis of the available sequence of observations and new test samples. The tests confirm that the proposed methodology assures fast learning of good quality, the predicted values being 'close' to the actual values in case of learning samples as well as in case of new test data. A series of further developments, still in progress, aiming to extend and refine the proposed methodology by taking into account non-linear models are going to be published in the near future.

Acknowledgement: The paper reports a series of results of the research performed in the framework of the Doctoral School in Computer Science at the University of Pitești, Romania.

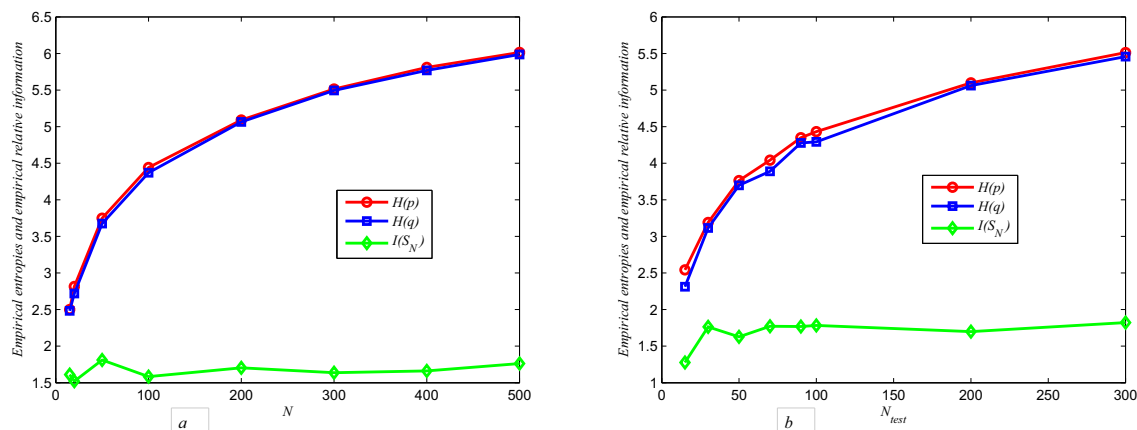


Figure 3: The variation of the empirical input/output entropies and empirical mutual information: a) as function of N , volume of learning data; b) for $N = 100$ and different volumes of test samples.

References

- [1] E. Alpaydin, *Introduction to Machine Learning*, The MIT Press, Massachusetts, 2010.
- [2] C. Bishop, *Pattern Recognition and Machine Learning*, Springer Verlag, 2006.
- [3] V. Cherkassky, F. Mulier, *Learning from Data Concepts, Theory, and Methods*, Second Edition, John Wiley & Sons, Inc., 2007.
- [4] T.M. Cover, J.A. Thomas, *Elements of information theory (Second Edition)*, John Wiley & Sons, Inc., 2006.
- [5] T. Hastie, R. Tibshirani, J. Friedman, *The Elements of Statistical Learning - Data mining, Inference and Prediction (Second Edition)*, Springer-Verlag, 2009.
- [6] A.J. Izenman, *Modern Multivariate Statistical Techniques: Regression, Classification and Manifold Learning*, Springer, 2008.
- [7] A.K. Jain, R. Duin, J. Mao, Statistical Pattern Recognition: A Review, *IEEE Transactions On Pattern Analysis and Machine Intelligence* 22(1): 4–37, 2000.
- [8] S. Marsland, *Machine Learning: An Algorithmic Perspective*, CRC Press, Taylor & Francis Group, Boca Raton - London - New York, 2009.
- [9] K.E. Muller, P.W. Stewart, *Linear Model Theory - Univariate, Multivariate and Mixed Models*, Wiley-Interscience, 2006.
- [10] I. Paraschiv-Munteanu, Theoretical approach in performance evaluation of a classification system, *University of Pitesti - Scientific Bulletin, Serie Mathematics and Computer Science*, 14: 139–150, 2008.
- [11] A.C. Rencher, *Methods of Multivariate Analysis*, Wiley-Interscience, 2002.
- [12] L. State, I. Paraschiv-Munteanu, *Introducere în teoria statistică a recunoașterii formelor*, Editura Universității din Pitești, Romania, 2009.
- [13] S. Sharma, *Applied Multivariate Techniques*, John Wiley & Sons, Inc., 1996.
- [14] K. Takeaki, K. Hiroshi, *Generalized Least Squares*, John Wiley & Sons, Inc., 2004.

Luminița State
 University of Pitești
 Faculty of Mathematics and Computer Science
 1 Târgu din Vale St., Pitești 110040
 ROMANIA
 E-mail: lstate@clicknet.ro

Iuliana Paraschiv-Munteanu
 University of Bucharest
 Faculty of Mathematics and Computer Science
 14 Academiei St., Bucharest 010014
 ROMANIA
 E-mail: pmiulia@fmi.unibuc.ro