

Contents

Introduction	iii
1 The Wishart Distribution	1
1.1 Definition	1
1.2 PDF Computation for Invertible Σ	1
1.2.1 Visualizing the Wishart Distribution	2
1.3 The Wishart Distribution in Bayesian Conjugate Prior Analysis	3
1.3.1 Bayesian Inference and Priors Distributions	3
1.3.2 The Wishart Conjugate Prior	5
2 The WISDoM Multiple Order Classifier	7
2.1 Wishart Sampling and Log-Likelihood Ratio Distance	7
2.1.1 Computing the Estimated Distribution	8
2.2 Log-Likelihood Ratio Distance	10
2.2.1 Complete Matrix Distance	11
2.2.2 Single Feature Distance and Multiple Order Reduction	11

Introduction

Chapter 1

The Wishart Distribution

1.1 Definition

The *wishart distribution* $W_p(n, \Sigma)$ is a probability distribution of random nonnegative-definite $p \times p$ matrices that is used to model random covariance matrices.

The parameter n is the number of degrees of freedom, and Σ is a nonnegative-definite symmetric $p \times p$ matrix, called the *scale matrix*.

Def. Let $X_1 \dots X_n$ be independent $N_p(0, \Sigma)$ distributed vectors, forming a data matrix $p \times n$, $X = [X_1 \dots X_n]$. The distribution of a $p \times p$, $M = XX' = \sum_{i=1}^n X_i X_i'$ *random matrix* is a Wishart distribution. [1]

We have then by definition:

$$M \sim W_p(n, \Sigma) \sim \sum_{i=1}^n X_i X_i' \quad X_i \sim N_p(0, \Sigma) \quad (1.1)$$

so that $M \sim W_p(n, \Sigma)$ is the distribution of a sum of n rank-one matrices defined by independent normal $X_i \in R^p$ with $E(X) = 0$ and $Cov(X) = \Sigma$.

In particular, it holds for the present case:

$$E(M) = nE(X_i X_i') = nCov(X_i) = n\Sigma \quad (1.2)$$

1.2 PDF Computation for Invertible Σ

In general, any $X \sim N(\mu, \Sigma)$ can be represented as

$$X = \mu + AZ, \quad Z \sim N(0, I_p) \quad (1.3)$$

so that

$$\Sigma = Cov(X) = ACov(Z)A' = AA' \quad (1.4)$$

The easiest way to find A in terms of Σ is the LU-decomposition, which finds a unique lower diagonal matrix A with $A_{ii} \geq 0$ such that $AA' = \Sigma$.

Then by 1.1 and 1.4, with $\mu = 0$ we have:

$$W_p(n, \Sigma) \sim \sum_{i=1}^n (AZ_i)(AZ_i)' \sim A \left(\sum_{i=1}^n Z_i Z_i' \right) A' \sim AW_p(n)A' \quad (1.5)$$

where $Z_i \sim N(0, I_p)$ and $W_p(n) = W_p(I_p, n)$.

Assuming that $n \geq p$ and Σ is invertible, the density of the random $p \times p$ matrix M in 1.1 can be written ¹ :

$$f(M, n, \Sigma) = \frac{1}{2^{\frac{np}{2}} \Gamma_p(\frac{n}{2}) \|\Sigma\|^{\frac{n}{2}}} \|M\|^{\frac{n-p-1}{2}} \exp\left[-\frac{1}{2} \text{trace}(\Sigma^{-1}M)\right] \quad (1.6)$$

so that $f(M, n, \Sigma) = 0$ unless M is *symmetric and positive-definite*. [2]

Note that in 1.6 we define $\Gamma_p(\alpha)$ as the *generalized gamma function*
 $\Gamma_p(\alpha) = \pi^{\frac{p(p-1)}{4}} \prod_{i=1}^p \Gamma\left(\frac{2\alpha+1-i}{2}\right)$

1.2.1 Visualizing the Wishart Distribution

The Wishart distribution is a generalization to multiple dimensions of the *chi-squared distribution*, or in the case of non-integer degrees of freedom, of the *gamma distribution*.

We show as a proof in fig.1.1 that for a 1-dimensional and equal to 1 Σ scale matrix, the Wishart distribution $W_1(n, 1)$ collapses to the $\chi^2(n)$ distribution.

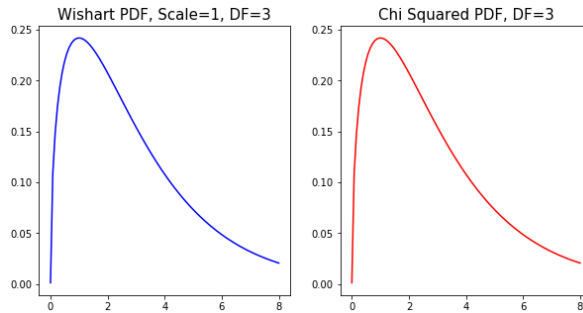


Figure 1.1: *Monodimensional Wishart Distribution and $\chi^2(n)$ distribution comparison*

¹Note: $\|\Sigma, N\| = \det(\Sigma, M)$

1.3. THE WISHART DISTRIBUTION IN BAYESIAN CONJUGATE PRIOR ANALYSIS 3

Save for this simple case, being the Wishart a distribution over matrices, it is a generally hard task to visualize it as a density function.

We can however sample from it and use the eigenvectors and eigenvalues of the resulting sampled matrix to define an ellipse.

An example of this technique is shown in fig.1.2. A set of five sampled matrices is drawn for each plot. While the parameter $n = 2$ (*degrees of freedom*) is the same for both the samplings shown, a different *scale matrix* Σ is used for each plot.

Note that for $\Sigma = I_2$ (left plot in fig.1.2) the sample would look *on average* like circles.

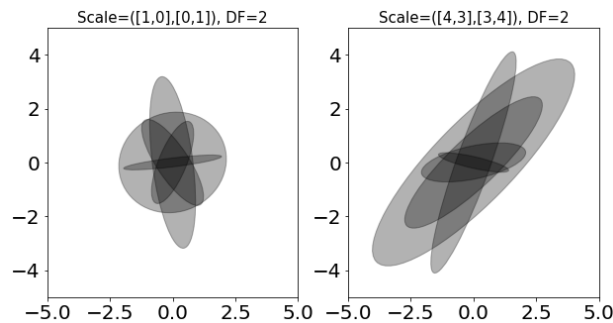


Figure 1.2: Plot of eigenvalue and eigenvectors defined ellipses, drawn from different scale matrix defined Wishart-sampled distribution.

1.3 The Wishart Distribution in Bayesian Conjugate Prior Analysis

An important use of the Wishart distribution is as a conjugate prior for *multivariate normal sampling*. We now recall some basics concepts about Bayesian inference and prediction in order to show the application of the Wishart in those fields.

1.3.1 Bayesian Inference and Priors Distributions

The distinctive feature of the Bayesian approach underlies in its way of defining probability.

Probability is treated as *belief* and not as *frequency*, thus introducing a fundamental difference between the Bayesian and the *frequentist* approach and shifting the goal toward the analysis and statement of a *belief* [3].

We can sum up the process of Bayesian inference as follows:

- A probability density called *prior distribution* $\pi(\theta)$ is chosen, expressing the *beliefs* about a parameter θ before any data are seen.
- A statistical model $p(x | \theta)$ is chosen, which must reflect the beliefs about x given θ .
- After observing the data $D_n = [X_1 \dots X_n]$, the beliefs is updated and the *posterior distribution* $p(\theta | D_n)$ is computed.

By Bayes' theorem the posterior distribution can be written as

$$p(\theta | X_1 \dots X_n) = \frac{p(X_1 \dots X_n | \theta) \pi(\theta)}{p(X_1 \dots X_n)} = \frac{L_n(\theta) \pi(\theta)}{c_n} \propto L_n(\theta) \pi(\theta) \quad (1.7)$$

where $L_n(\theta) = \prod_{i=1}^n p(X_i | \theta)$ is the likelihood function and the *normalizing constant* c_n is defined as follows:

$$c_n = p(X_1 \dots X_n) = \int p(X_1 \dots X_n | \theta) \pi(\theta) d\theta = \int L_n(\theta) \pi(\theta) d\theta \quad (1.8)$$

the normalizing constant is also called the *evidence*.

We now define the general properties of a *conjugate prior*.

If, for a given problem, the posterior distribution $p(\theta | X_1 \dots X_n)$ and the prior $\pi(\theta)$ belong to the same family of distribution, they're called *conjugated distributions* and the prior is said to be a *conjugate prior* for the given *likelihood function* $L_n = p(X_1 \dots X_n | \theta)$.

A classical example concerns the Gaussian Distribution: the Gaussian family is conjugate to itself (or self-conjugate) with respect to a Gaussian likelihood function: if the likelihood function is Gaussian, choosing a Gaussian prior over the mean will ensure that the posterior distribution is also Gaussian [4].

Considering the general problem of inferring a distribution for a parameter θ given some observations $D_n = [X_1 \dots X_n]$ and referring to theorem 1.7, by which we let the likelihood function be considered fixed as it is usually well-determined from a statement of the data-generating process, it is clear that different choices of the prior distribution $\pi(\theta)$ may make the integral in 1.8 more or less difficult to compute. The product $L_n(\theta) \pi(\theta)$ will also be influenced, gaining the possibility to take one algebraic form or another.

If for certain choices of the prior the posterior has the same algebraic form as the prior, those choices are said to yield a *conjugate prior*.

It is then possible to state that a conjugate prior is an algebraic convenience giving a closed-under-sampling-form expression for the posterior.

1.3.2 The Wishart Conjugate Prior

We now show how the Wishart Distribution is correlated to the *Inverse Gamma Distribution* in a multidimensional setting, by considering a Gaussian model with known mean μ , so that the free parameter is the variance σ^2 , as in [3].

The likelihood function is defined as follows:

$$p(X_1 \dots X_n \mid \sigma^2) \propto (\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} n \overline{(X - \mu^2)}\right), \quad \overline{(X - \mu^2)} = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 \quad (1.9)$$

The conjugate prior is an inverse Gamma distribution. Recall that θ has an inverse Gamma distribution with parameters (α, β) when $\frac{1}{\theta} \sim \text{Gamma}(\alpha, \beta)$.

The density is then bound to take the form

$$\pi_{\alpha, \beta}(\theta) \propto \theta^{-(\alpha+1)} e^{-\frac{\beta}{\theta}} \quad (1.10)$$

Using this prior, the posterior distribution of σ^2 is given by

$$p(\sigma^2 \mid X_1 \dots X_n) \sim \text{InvGamma}\left(\alpha + \frac{n}{2}, \beta + \frac{n}{2} \overline{(X - \mu^2)}\right) \quad (1.11)$$

An alternative way of parameterization of the prior is given by the *Inverse Scaled χ^2 Distribution*, whose density is defined as

$$\pi_{\nu_0, \sigma_0^2} \propto \theta^{-(1+\frac{\nu_0}{2})} \exp\left(-\frac{\nu_0 \sigma_0^2}{2\theta}\right) \quad (1.12)$$

Under this kind of parameterization of the prior, the posterior takes the form

$$p(\sigma^2 \mid X_1 \dots X_n) \sim \text{ScaledInv}\chi^2\left(\nu_0 + n, \frac{\nu_0 \sigma_0^2}{\nu_0 + n} + \frac{n \overline{(X - \mu^2)}}{\nu_0 + n}\right) \quad (1.13)$$

In the multidimensional setting, the inverse Wishart takes the place of the inverse Gamma.

It has already been stated that the Wishart distribution is a distribution over *symmetric positive semi-definite* $d \times d$ matrices W . A more compact form of the density is given by

$$\pi_{\nu_0, S_0}(W) \propto |W|^{\frac{(\nu_0 - d - 1)}{2}} \exp\left(-\frac{1}{2} \text{trace}(S_0^{-1} W)\right), \quad |W| = \det(W) \quad (1.14)$$

where the parameters are the degrees of freedom ν_0 and the positive-definite *scale matrix* S_0 .

If $W^{-1} \sim \text{Wishart}(\nu_0, S_0)$ we can then state that W has an *Inverse Wishart Distribution*, whose density has the form

$$\pi_{\nu_0, S_0}(W) \propto |W|^{-\frac{(\nu_0+d+1)}{2}} \exp\left(-\frac{1}{2}\text{trace}(S_0 W^{-1})\right), \quad |W| = \det(W) \quad (1.15)$$

Let $X_1 \dots X_n$ be $N(0, \Sigma)$ distributed observed data. Then an inverse Wishart prior multiplying the likelihood $p(X_1 \dots X_n | \Sigma)$ yields

$$\begin{aligned} p(X_1 \dots X_n | \Sigma) \pi_{\nu_0, S_0}(\Sigma) &\propto \\ | \Sigma |^{-\frac{n}{2}} \exp\left(-\frac{n}{2}\text{trace}(\bar{S} \Sigma^{-1})\right) &| \Sigma |^{-\frac{(\nu_0+d+1)}{2}} \exp\left(-\frac{1}{2}\text{trace}(S_0 \Sigma^{-1})\right) \quad (1.16) \\ = | \Sigma |^{-\frac{(\nu_0+d+n+1)}{2}} \exp\left(-\frac{1}{2}\text{trace}((n\bar{S} + S_0) \Sigma^{-1})\right) \end{aligned}$$

where \bar{S} is the *empirical covariance* $\bar{S} = \frac{1}{n} \sum_{i=1}^n X_i X_i^T$.

Thus, a posterior with the form

$$p(\Sigma | X_1 \dots X_n) \sim \text{InvWishart}(\nu_0 + n, n\bar{S} + S_0) \quad (1.17)$$

is obtained.

Analogously, it can be stated that for the inverse covariance (*precision*) matrix Σ^{-1} the conjugate prior is a Wishart distribution.

Chapter 2

The WISDoM Multiple Order Classifier

In this section, the classification method implemented and used on the ADNI2 and ABIDE databases is described, both in an analytical and technical way.

The "*distance*" used to train the classifier is defined as well as the *feature transformation* undergone by the each of the subject analyzed.

The general pipeline and the validation pipeline are then discussed while also introducing an example of possible parallelization for performance enhancing.

2.1 Wishart Sampling and Log-Likelihood Ratio Distance

Considering what has been said in the last section, using the Wishart distribution to model and sample the elements of a wide range of problems follows naturally.

As a matter of fact, every classification problem whose elements take the form of *symmetric positive-definite* matrices can be approached with the method we are about to discuss.

The main idea for the *WISDoM Classifier* is to use the *free parameters* of the Wishart distribution (the scale matrix S_0 and the number n of the degree of freedom, as shown in 1.6) to compute an estimation of the distribution for a certain class of elements, and then assign a single element to a given class by computing some sort of distance between the element being analyzed and the classes.

Furhermore, if we assume that the matrices are somehow representative of the *features* of the system studied (i.e. covariance matrices might be

taken into account), a score can be assigned to each feature by estimating the weight of said feature in terms of *Log Likelihood Ratio*.

In other words, a score can be assigned to each feature by analyzing the variation in terms of *LogLikelihood* caused by the deletion of it.

If the deletion of a feature causes significant increase (or decrease) in the *LogLikelihood* computed with respect to the *estimated distributions* for the classes, it can be stated that said feature is highly representative of the system analyzed.

It is now clear that the simplest usable objects to estimate the distribution for a class and to represent its elements is the *covariance matrix*. Further proofs for this statement will be given later on.

Thus, the aim of the WISDoM classifier is not only to assign a given element to the optimal class, but also to identify the features with the highest "*weights*" in the decision process.

2.1.1 Computing the Estimated Distribution

Let us briefly recall the parametrization of the Wishart Distribution in order to clearly define the application conditions for classification problems.

Let $X_1...X_n$ be independent $N_p(0, \Sigma)$ distributed vectors, forming a data matrix $p \times n$, $X = [X_1...X_n]$. The distribution of a $p \times p$, $M = XX' = \sum_{i=1}^n X_i X_i'$ *random matrix* is a Wishart distribution with parameters $W_p(n, S_0)$. In the previous chapter (1.6) it has been proved that for normal distributed data, for $S_0 = \Sigma$, a distribution of *random covariance matrices* is obtained.

In a similar fashion, if a good choice for the scale matrix S_0 is made for a given class, a representative distribution for the class can be estimated and samples can be drawn from it.

Covariance matrices are a good choice, although not limiting as long as the matrices are symmetric and positive-definite, both for the way they represent a system and for the property that *the mean of a set of covariance matrices is a covariance matrix*.

If each element of a given class C is represented by a covariance matrix Σ of its features, this property allows us to estimate a distribution for the class by choosing $S_0 = \hat{\Sigma}_C = \frac{1}{N} \sum_{i=1}^N \Sigma_i$.

The other necessary parameter for the estimation is the *degrees of freedom* n .

Assume that an $X_i = (x_1, ..., x_p)$ vector of p features is associated to each element i of a given class, while having n observation of said vector. The covariance matrix Σ_i computed over the n observations will represent the "interactions" between the features of element i .

2.1. WISHART SAMPLING AND LOG-LIKELIHOOD RATIO DISTANCE9

The degrees of freedom n of the Wishart distribution are then given by the number of times X_i is observed.

Let us introduce an example tied to *functional MR brain imaging* in order to further clarify the concepts being introduced.

An image of patient i 's brain is acquired; as usual these images are divided in a certain number p of zones (voxel, pixel etc.), each zone being sampled n times over a given time interval in order to observe a certain type of brain activity and functionality.

It is now clear that the features contained in vector $X_i = (x_1, \dots, x_p)$ associated to patient i are indeed the zones chosen to divide i 's brain image, each zone having been sampled n times during an acquisition interval.

The correlation $p \times p$ matrix Σ_i computed for i 's observation is then representative of the functional correlation between the p zones of i 's brain.

Repeating this procedure for N patients belonging to a known class C (i.e. a diagnostic group) and computing the $\hat{\Sigma}_C$ scale matrix for the class as stated before, will allow us to estimate a wishart distribution for that class correlation matrices and draw samples from it.

The module used for Wishart generation and sampling by the WISDoM calssifier is the *SciPy.Stats.Wishart* module of the *SciPy* Python3.6 library.

Further details on the generation and sampling algorithm used by the module can be found in [5].

Some samples drawn from Wishart distributions computed with different 5×5 scale matrices and degrees of freedom are shown in fig.2.1.

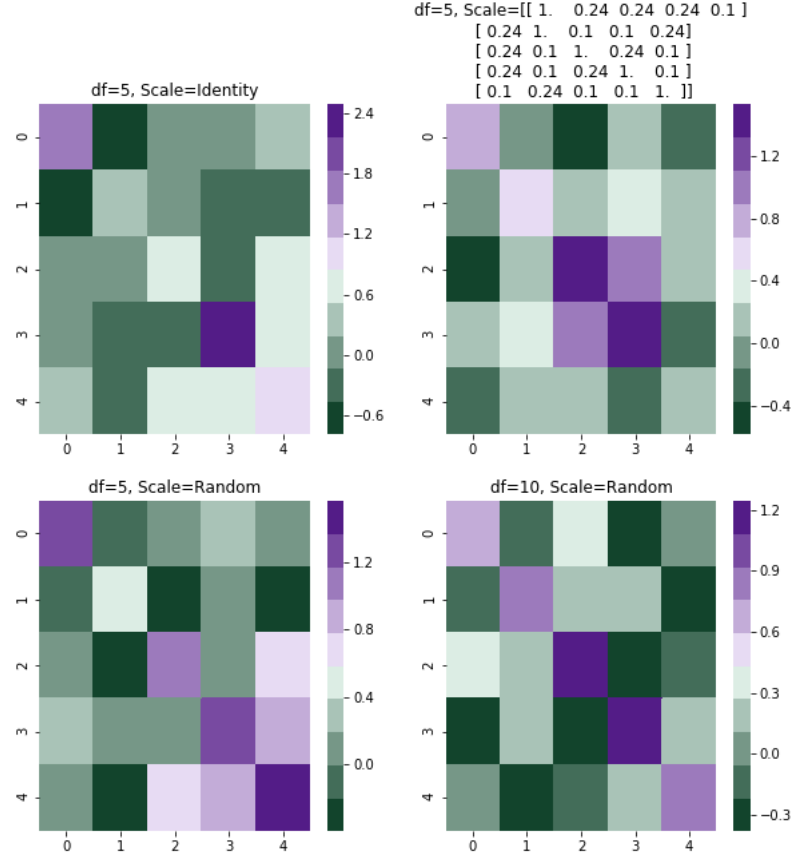


Figure 2.1: Various sampling from different wishart distribution. A diverging heatmap has been chosen to visualize the values of each sample's elements.

2.2 Log-Likelihood Ratio Distance

After the definition of the role of the Wishart distribution in symmetric positive definite matrices' modeling, it is necessary to define some sort of distance between the estimated distribution for a class C and its hypothetical elements.

As stated before, this will be done in terms of both entire matrices and *single features*, in order to achieve optimal classification and extract information about a system's most meaningful components.

2.2.1 Complete Matrix Distance

The scoring system used by the WISDoM Classifier relies on the *logpdf* function from the *SciPy.Stats.Wishart* module in order to compute the LogLikelihood of a matrix Σ_i with respect to the Wishart distribution estimated for a class C , using $\hat{\Sigma}_C$ as the scale matrix.

If a problem concerning two given classes C_A and C_B is taken into account, the score assigned to each Σ_i upon which the classification decision is based, can be defined as follows:

$$score_i = \log P_W(\Sigma_i | n, \hat{\Sigma}_A) - \log P_W(\Sigma_i | n, \hat{\Sigma}_B) \quad (2.1)$$

Where $\hat{\Sigma}_{A,B}$ are the scale matrix computed for the classes A, B and $\log P_W(\Sigma_i | n, \hat{\Sigma}_{A,B})$ can be seen as the logarithm of the probability of Σ_i belonging to the Wishart distribution estimated for one of the two classes A, B .

2.2.2 Single Feature Distance and Multiple Order Reduction

The aim of the WISDoM classifier is to further increase the informations obtained about the system's features during the classification.

To do this it is then necessary to introduce some mathematical properties of the symmetric positive definite matrices, upon which the method relies.

It will be shown that it is indeed possible to access different orders of information by scaling a matrix A to its *principal submatrices*.

Def. Let A be an $n \times n$ matrix. A $k \times k$ submatrix of A formed by deleting $n - k$ rows of A , and the same $n - k$ columns of A , is called *principal submatrix* of A . The determinant of a principal submatrix of A is called a *principal minor* of A .

Note that the definition does not specify which $n - k$ rows and columns to delete, only that their indices must be the same.

Let us introduce a 3×3 example.

For a general matrix $A_{3 \times 3}$

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \quad (2.2)$$

there are three *first order principal minors*:

- $|a_{11}|$ formed by deleting the last two rows and columns
- $|a_{22}|$ formed by deleting the first and third rows and columns
- $|a_{33}|$ formed by deleting the first two rows and columns

There are three *second order principal minors*:

- $\begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix}$ formed by deleting column 3 and row 3
- $\begin{vmatrix} a_{11} & a_{13} \\ a_{31} & a_{33} \end{vmatrix}$ formed by deleting column 2 and row 2
- $\begin{vmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{vmatrix}$ formed by deleting column 1 and row 1

There's one *third order principal minor*, namely $|A|$.

For the sake of completion, we also recall the following definition.

Def. Let A be an $n \times n$ matrix. The k^{th} order principal sub-matrix of A obtained by deleting the *last* $n - k$ rows and columns of A is called the k^{th} order **leading principal submatrix** of A , and its determinant is called the k^{th} **order leading principal minor** of A .

An important property for the principal submatrices of a symmetric positive definite matrix is that *any $(n - k) \times (n - k)$ partition is also symmetric and positive definite*.

It is now clear that such properties can be used to reduce both a class scale matrix $\hat{\Sigma}_C$ and any Σ_i matrix, in order to study its deviation from a class's estimated Wishart distribution derived from the deletion of one of its components (the features contained in vector X_p from which the matrix $\Sigma_{i,p \times p}$ is computed).

Iterating this process over all the features, or in other terms analyzing all of the $(p - 1) \times (p - 1)$ principal submatrices of Σ_i and $\hat{\Sigma}_C$, will allow us to assign a score to each feature, representing its weight in the decision for Σ_i to be assigned to one class or another.

Note that for such an order of principal submatrices, the process will reduce the $\Sigma_{i,p \times p}$ matrix into a *score vector* of length p for each element i undergoing the classification.

Bibliography

- [1] Hardle, Wolfgang and Leopold Simar
Applied Multivariate Statistical Analysis
Heidelberg: Springer Berlin Heidelberg, 2012
- [2] Anderson, T. W.
An Introduction to Multivariate Statistical Analysis
New York: John Wiley and Sons, 2003
- [3] Han Liu and Larry Wasserman
Statistical Machine Learning
Pittsburgh: CMU University, 2014
- [4] Murphy, Kevin P.
Conjugate Bayesian Analysis of the Gaussian Distribution
Vancouver: University of British Columbia, 2007
- [5] W.B. Smith and R.R. Hocking
Algorithm AS 53: Wishart Variate Generator
Applied Statistics, vol. 21, pp. 341-345, 1972.