

UNIVERSITAT POLITÈCNICA DE CATALUNYA (UPC)  
– BARCELONATECH

BACHELOR'S THESIS

---

# Using Random Fourier Features with Random Forest

---

*Author:*  
Albert RIBES

*Supervisor:*  
Lluís A. BELANCHE

Computer Science

April 11, 2019



UNIVERSITAT POLITÈCNICA DE CATALUNYA (UPC) – BARCELONATECH

# *Abstract*

Facultat d'Informàtica de Barcelona (FIB)  
Computer Science

Bachelor Degree in Computer Science

**Using Random Fourier Features with Random Forest**

by Albert RIBES

The Thesis Abstract is written here (and usually kept to just this page). The page is kept centered vertically so can expand into the blank space above the title too...

- En 3 idiomas



# Contents

<b>Abstract</b>	<b>iii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Problem to study . . . . .	1
1.2 Project proposal . . . . .	1
<b>2 Background Information and Theory</b>	<b>5</b>
2.1 Machine Learning . . . . .	5
2.2 Some currently used Machine Learning models . . . . .	6
2.2.1 Decision Tree . . . . .	6
2.2.2 Logistic Regression . . . . .	6
2.2.3 Support Vector Machines . . . . .	6
2.3 Ensemble Methods . . . . .	7
2.4 The kernel trick . . . . .	9
2.5 Random Fourier Features . . . . .	10
2.6 Nyström . . . . .	10
<b>3 Project Development</b>	<b>11</b>
3.1 General Idea . . . . .	11
3.1.1 State of the art con las RFF . . . . .	12
3.1.2 State of the art con las Nyström . . . . .	13
3.2 Hyper-parameters . . . . .	13
3.3 Hypothesis . . . . .	14
3.3.1 Experiments Proposal . . . . .	16
3.4 Datasets . . . . .	18
<b>4 Experimental Results</b>	<b>21</b>
4.1 Enfrentar resultados 2 a 2 . . . . .	21
4.2 2-8 . . . . .	22
4.3 Contrastar hipótesis con resultados . . . . .	23
<b>5 Conclusion and Future Directions</b>	<b>25</b>
<b>6 Sustainability Report</b>	<b>27</b>
6.1 Environmental . . . . .	27
6.2 Economic . . . . .	27
6.3 Social . . . . .	27
6.3.1 Impacto Personal . . . . .	27
6.3.2 Impacto Social . . . . .	27
6.3.3 Riesgos Sociales . . . . .	27
<b>A Results of experiment 1.1</b>	<b>29</b>
<b>B Results of experiment 2.1</b>	<b>31</b>

<b>C</b>	<b>Results of experiment 2.2</b>	<b>33</b>
<b>D</b>	<b>Results of experiment 2.3</b>	<b>37</b>
<b>E</b>	<b>Results of experiment 2.4</b>	<b>39</b>
<b>F</b>	<b>Results of experiment 3.1</b>	<b>43</b>
<b>G</b>	<b>Results of experiment 3.2</b>	<b>45</b>
<b>H</b>	<b>Results of experiment 3.3</b>	<b>47</b>
<b>I</b>	<b>Results of experiment 3.4</b>	<b>49</b>
<b>J</b>	<b>Results of experiment 4.1</b>	<b>51</b>
<b>K</b>	<b>Results of experiment 4.2</b>	<b>53</b>
	<b>References</b>	<b>57</b>
	<b>Datasets</b>	<b>59</b>

# List of Figures

3.1 Model Boxes Diagram . . . . .	13
-----------------------------------	----





## Chapter 1

# Introduction

### 1.1 Problem to study

Supervised Learning uses statistical and mathematical models to predict a response variable from a classification or regression problem. Usually it does so by trying to minimise an error function or to minimise an accuracy function. The better a model is, the higher the accuracy it will obtain on new, unseen data.

Increasing the accuracy of the models is one of the main topics in the field. But it is not easy to achieve it. Usually, it comes at the cost of increasing the computation time to produce the model. Therefore, a trade-off needs to be made between the accuracy obtained and the training time.

In this project we study some new approaches to improve the trade-off of currently used Machine Learning methods. In particular, we study how Kernel approximation techniques could make some procedures feasible in very large datasets and also how to increase the accuracy of some models at the expense of some more time.

### 1.2 Project proposal

- Existe una batería de técnicas que son buenas, pero que nadie las ha combinado. Son:
  - Modelos simples
  - Ensembles
  - kernel trick
  - Aproximaciones de kernel
- La propuesta es combinar todo esto para mejorar el trade-off
- Sostenemos las siguientes hipótesis:
  - Se podría hacer un ensemble con modelos distintos a DT
  - Se puede aproximar una RBF-SVM pero con el coste de una lineal
  - RFF + Bootstrap quizá es demasiado aleatorio
  - Los modelos que no se basan en productos escalares no se beneficiarán tanto de usar RFF
- Lo que se hará en cada capítulo del trabajo

The current development of Machine Learning has opened many fronts and techniques trying to solve many of the difficulties that the field has.

One known issue is the Bias–Variance dilemma. While solving a classification or regression problem, the expected generalization error of a model is the sum of three error terms: the bias, the variance and the irreducible error. While the last one, as the name suggests, cannot be reduced, since it is caused by the inherent random noise in the data, the other two terms seem to have an inverse proportion: trying to reduce one of them increases the other most of the times. Ensemble methods were developed in an attempt to reduce both of them (or at least their sum). Although they can outperform some models their usage is mostly restricted to a small subset. This is because they only increase significantly the accuracy on unstable methods.

Another advance in Machine Learning has been the usage of kernel methods. They are useful to implicitly transform the data into another feature space with better properties, such as a clear dividing margin between classes of data. They are very effective, but their high computational costs has caused their use to be limited to just some specific problems or with a small number of instances. There are some less expensive approaches to approximate these methods, but they are not widely used.

And then there is a collection of classical algorithms which have the advantage of being very simple and straightforward, although they don't usually get the highest scores.

There is a collection of techniques that have shown some good results by their own, but still they haven't been tested in combination with the others. If we could mix some of these methods, maybe we could find new Machine Learning methods, with better accuracy or trade-off.

In this project, we try some combinations of currently known techniques which could produce better results or show new useful model designs. For the one hand, we try to extend the usage of ensemble methods to new basic

models. Currently, it doesn't make much sense to train an ensemble of Support Vector Machines or of Logistic Regression models, because they are so stable and most of the estimators would predict the same answer. We propose the use of random kernel approximations such as Random Fourier Features (RFF) or the Nyström method to increase the unstability of these models and thus be able to successfully train and ensemble with them, hopefully increasing the accuracy of a single one.

On the other hand, the usage of these Random Kernel approximations could allow us to use some methods which right now are not accessible. Support Vector Machines cannot be used with non-linear Kernels such as the Radial Basis Function (RBF) kernel on big datasets, since the cost is very high. But if we transform the data to some space almost equivalent to the one of the RBF kernel, we can use a Linear Support Vector Machine, which is less expensive to train, and achieve a similar accuracy. In fact, with the usage of an ensemble, the results could be improved. Similar approaches have already been studied, and they have showed good results [9] [13] [8].

For this study, we have formulated some hypothesis and we try to confirm or refuse them based on the experimental results. The hypothesis are:

1. An SVM using RFF or Nyström could achieve an accuracy similar to using the RBF kernel but with much less training time.
2. Training an ensemble of SVMs or Logistic Regression models using RFF or Nyström could increase the accuracy of a single estimator.
3. Mixing a random mapping with the Bootstrap (from Bagging) could cause an excess of randomness and hence a bad accuracy compared to just using the random mapping.
4. Basic algorithms which are not based on the dot product of the input data such as the Decision Tree will not benefit so much of the usage of RFF and Nyström than those that do, like SVM or Logistic Regression.

In the following pages we show an experimental set up to check these hypothesis and the results of the experiments. In Chapter 2 we explain some of the Machine Learning concepts which are needed to understand the rest of the project. In 3 we explain with more detail how this study is developed and what experiments are executed. In 4 we show the results obtained with the experiments and discuss the hypothesis suggested based on them. In 5 we present the conclusions of this project and propose some future work related to this topic. Finally, in Chapter 6 we present a sustainability report of this project.



## Chapter 2

# Background Information and Theory

### 2.1 Machine Learning

- Una definición rápida
- Clasificación y regresión
- Cross-validation
- Qué son los datos de train y test, y por qué se hace esa partición
- Qué es el sobre-ajuste

Machine Learning uses statistical and mathematical models to give answers to problems when there is no known formula or procedure to compute the answer. In the subfield of Supervised Learning, the objective is to predict a numerical or categorical variable in response to some input data, and the way of doing it is to feed the model with lots of different examples for which we already know the correct answer, and we expect the models to be able to predict the correct answer to instances that it hasn't seen before. When it does, we say that the model is able to generalize.

When a model is trained with some data, there is always a risk of overfitting. For a model to overfit means that it adjusts very well to the data that it has seen, but can't predict the correct answer to new, unseen data. This happens because it has not only learned from the relevant information, but also from the random noise that the data had, and so it can only memorize, but not generalize. For this reason, when a Machine Learning algorithm is trained the data is split in two subsets, a *Training dataset* and a *Testing dataset*. The training dataset will be used to train the model, while the testing datasets will be used only to check it. If a model has generalized well, it will achieve a good accuracy score on both the training and the testing dataset, but if it has overfitted it will show good results in the training dataset and bad ones in the testing dataset. Many models need some parameters to tune the behaviour of the algorithm. For example, some of them are used to adjust how much a model will fit to the data. We usually call these "hyperparameters". The correct value for them is not straightforward, and it is normally chosen with a process called "cross-validation". This process consists of splitting the training dataset in many subsets and checking many possible values for the hyperparameters in order to

see which one gets a higher accuracy with unseen data.

## 2.2 Some currently used Machine Learning models

### 2.2.1 Decision Tree

- No se basa en productos escalares
- Es extremadamente rápido
- Es más fácil de interpretar que otros modelos
- Es extremadamente inestable
- Cuando se hace un Random Forest, se randommiza un poco, de modo que árboles distintos entrenados con los mismos datos pueden ser distintos
- Es un modelo no lineal

Decision Tree is a predictive model which uses the training data to build a tree where each node splits the data in two sets according to some feature, and the leaves contain the set of instances that belong to some class (in classification problems) or that has a similar numerical response variable (for regression problems).

To predict the answer to a new instance, it uses the features to “decide” the nodes to cross until it reaches a leaf. The response given is the most prevailing class in the leaf for classification problems, or the mean of the values of the rest of the instances in the leaf.

Decision Trees have the advantages that it is easy to interpret the tree produced and that it is very fast to build the tree. The way to avoid overfitting is to limit its growth.

These models are very unstable. This means that small differences in the training data can produce very different Decision Trees. This property is very useful to build an ensemble of estimators to produce better answers. Random Forest is an algorithm that trains many Decision Trees with some sort of randomization.

### 2.2.2 Logistic Regression

### 2.2.3 Support Vector Machines

- Inicialmente pensadas para clasificación en 2 clases
- Pero se puede más clases con *one-vs-rest* y también hay formas de hacer regresión
- Se basa únicamente en el producto escalar de sus entradas

- Intenta separar los datos con un hiper-plano
- Actualmente es poco eficiente usarlas porque su coste es cúbico con la cantidad de entradas.
- Las fórmulas que quiere optimizar

Support Vector Machine (SVM) is a model that finds in hyperplane that divides the data in two sets. In two-class classification problems, each side of the hyperplane contains the instances of each of the classes. It does so by converting the problem to an optimization one.

Given some data  $D = \{\chi, \mathbf{y}\}$ , where  $\chi = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $\mathbf{y} = \{-1, +1\}^n$ , the optimization problem consists on finding  $\alpha \in \mathbb{R}^n$  the maximises

$$L = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \quad (2.1)$$

subject to

$$0 \leq \alpha_i \leq C; \forall i \quad (2.2)$$

$$\sum_{i=1}^n \alpha_i y_i = 0 \quad (2.3)$$

$C$  is an hyperparameter to tune the ammount of penalization for miss-classified instances.

If we compute

$$w = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \quad (2.4)$$

and

$$b = y_i - w \mathbf{x}_i \quad (2.5)$$

for any  $i$  so that  $\alpha_i \neq 0$ , we can compute the class of  $\mathbf{x}_0$  with  $\text{sign}(w \mathbf{x}_0 + b)$ .

Note that this algorithm just uses the dot product of the input data, not the data itself. This property allows us to use the Kernel Trick with them. See [2.4](#)

## 2.3 Ensemble Methods

- Bagging
  - Inventado por Leo Breiman (referencia)
  - Pretende reducir el sesgo
  - Wikipedia dice que pretende reducir la varianza
  - Es el boosting el que pretende reducir el sesgo
  - Entrenamiento de los estimadores es independiente, se podría hacer en paralelo

- Actualmente casi solo se usa con DT, debido a su inestabilidad
- Bootstrap
  - Intenta solucionar el problema de que para bagging es bueno que los estimadores sean distintos
  - Idealmente usaríamos un dataset distinto para cada estimador
  - Consiste en hacer un resampling con repetición
  - Si la cantidad de instancias del original es la misma que la de cada uno de los subconjuntos, se espera que la proporción de elementos únicos sea de  $1 - \frac{1}{e} \approx 0.632$ .
  - Si el conjunto original tiene  $n$  elementos, y tu haces un subconjunto de tamaño  $r$ , puedes esperar que la proporción de elementos del original que sí tienen presencia en el nuevo sea de  $1 - e^{-\frac{r}{n}}$
- Random Forest

Ensemble methods are a technique used in Machine Learning to reduce the overall accuracy error of a basic classification or regression model. The idea is that a comitee of models is expected to learn better than a single one.

Some ensemble methods are focused on decreasing the error caused by the variance of the model. One example is *Bagging*. Others are focused on error decreasing the bias error, like *Boosting*.

In Bagging, every model in the ensemble vote with equal weight. Thus, it is important to promote the variance among each of the models, since not doing it would be equivalent to training just one model. Ideally, one would train each of the models with totally different datasets, with no correlation among them. But in practice this is not always possible, because of a limited number of instances to train. One alternative is to use a technique called *Bootstrap*. Bootstrap allows to generate many different instances of a dataset by performing a resampling.

Given a dataset  $D$  of size  $n$ , Bootstrap generates  $m$  new datasets  $D_i$  of size  $n$  by sampling instances from  $D$  uniformly and with replacement. This means that some of the instances in  $D$  may be repeated in  $D_i$ , and others may not appear at all. With a large  $n$ , it is expected that each dataset  $D_i$  will contain 63.2% of the instances in  $D$ .

Theoretically Bagging could be used with any kind of method. However, for most of them Bootstrap is not enough to decorrelate the estimators. In practice, Bagging is mostly used with Decision Tree, given that this method produces very different trees with a small variation in the data. Random Forest is an algorithm that trains many Decision Trees with a Bagging. Instead of building the tree in a deterministic way, in each split it choses a random subset of features on which to perform the separation. Besides, it lets the estimators overfit, since it has a positive impact in reducing the overall variance of the Forest.



## 2.4 The kernel trick

- Teorema de Bochner
- El kernel RBF
  - Su fórmula es ...
  - Equivalencia entre  $\gamma$  y  $\sigma$
  - La noción de similitud que tiene
  - $\mathcal{H}$  es de dimensionalidad infinita
  - Permite ajustarse infinitamente a los datos, tuneando el hiperparámetro
  - $\sigma$  más pequeño, más sobreajuste
  - $\gamma$  más grande, más sobreajuste

A Kernel is a function that equals to the inner product of inputs mapped into some Hilbert Space <sup>a</sup>, i.e:

$$\kappa(x, y) = \langle \phi(x) \phi(y) \rangle \quad (2.6)$$

They are interesting in Machine Learning because we don't need to know the explicit function  $\phi(\cdot)$ . In fact,  $\phi(\cdot)$  could map the data to a Hilbert Space with infinite dimensions, and we could still compute  $\phi(x)\phi(y)$  through the kernel  $\kappa$

Support Vector Machines can benefit a lot of Kernel Functions. SVMs solve an optimization problem to maximise

$$L = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \quad (2.7)$$

in order to find an hyperplane that separates the data points in two classes. But with some problems there may not exist such hyperplane, and so it would be needed to map the data to a different feature space. If we did that, then the function to maximise would be

$$L = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \quad (2.8)$$

As we said previously, SVMs don't work with the data points alone, but just with their inner products. Thus, a Kernel could be used to define the optimization problem as

$$L = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \quad (2.9)$$

This approach has one big advantage: as long as the learning technique relies only on the inner product of the input, the underlying mapping  $\phi(\cdot)$  does not need to be explicitly calculated and can, in fact, be unknown.

Kernel functions can be characterised with the Mercer's condition [7]. It says

that given a function  $\kappa(\mathbf{x}, \mathbf{y})$ , there exists a mapping  $\phi(\cdot)$  so that  $\kappa(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})\phi(\mathbf{y})$  if and only if for any  $g(\mathbf{x})$  such that  $\int g(\mathbf{x})^2 d\mathbf{x}$  is finite then  $\kappa(\mathbf{x}, \mathbf{y})g(\mathbf{x})g(\mathbf{y}) \geq 0$ .

There are many known Kernels. One that is very popular is the Radial Basis Function Kernel, RBF. This kernel is defined as:

$$\kappa(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|^2} \quad (2.10)$$

where  $\gamma$  is a free parameter. The value of this Kernel decreases with the euclidean distance of the parameters, so it can be interpreted as a measure of similarity. The feature space of this kernel has infinite number of dimensions. When a kernel is used with an SVM, the answer can be computed with  $\text{sign}(\sum_{i=1}^n \alpha_i y_i \kappa(\mathbf{x}_i, \mathbf{x}))$ . SVMs using the RBF kernel have a huge ability to fit to the data, and is able to separate classes for very difficult problems. The problem is that the optimization of the function

<sup>a</sup>A Hilbert space is just a generalization of the Euclidean Space which contains the structure of an inner product that allows length and angle to be measured.

## 2.5 Random Fourier Features

## 2.6 Nyström

## Chapter 3

# Project Development

### 3.1 General Idea

- Hemos visto que se puede sacar una aproximación aleatoria de la función implícita de un shift invariant kernel. Esto tiene 2 ventajas
  - Podemos transformar los datos directamente
  - Podemos producir pequeñas variaciones de un mismo dataset, todas ellas válidas
- Las 4 tipos de modelos que he definido. Referencia a la foto
- ¿Por qué he cogido estos 4 modelos? ¿No podrían haber sido otros? ¿Que tienen estos de bueno? Me he inspirado en Random Forest
- Hay por ahí algún paper que compara RFF y Nyström

- Usar RFF añade un coste al procedimiento, pero es lineal con la cantidad de instancias
- Por lo tanto, el uso de las RFF probablemente solo será útil con datasets muy grandes. Si el dataset es pequeño, los costes que tiene son demasiado grandes
- De los 9 problemas con los que hemos experimentado, 7 de ellos son pequeños y únicamente sirven para ilustrar la mejoría de precisión, pero no los tiempos
- El MNIST y el Fashion MNIST se han usado para sacar conclusiones de tiempo

RFF and Nyström can be very useful for two main reasons:

- They allow us to perform an explicit mapping approximating the Kernel feature space.
- They can be used to generate many different equally valid datasets with the required number of features from a single one.

We will study how these advantages could be used to increase the accuracy of some models. In particular, we expect they will be useful in two different ways. First, some methods could achieve a higher accuracy if they were

trained with data in a Kernel feature space instead of the general one. Second, generating many datasets could make it possible to increase the accuracy of some models by training an ensemble.

Although both approaches are expected to increase the accuracy of some models, they will also increase a little bit the training time. Performing a random mapping has a linear cost with the number of instances, and training an ensemble clearly multiplies the amount of work to be done. Thus, the benefit of the methods presented will depend on the circumstances of each problem. It should be noted that both increases in time are linear with the number of instances, so they can scale well with very large datasets. This is in contrast to training SVM with non-linear kernels, whose cost is cubic with the number of instances.

For this project we have used 9 different datasets to check the hypothesis. 7 of them were quite small (5000 instances approx.) and the other 2 had a much bigger size (70000 instances). For this reason, the first ones will be useful to compare the accuracies obtained, but will show a very big increase in the training time. In contrast, the bigger datasets will show the real strength of these methods.

We propose two ways of using the random mappings. The first one is very straightforward: simply using a random mapping of the data to train a single model, and use the same mapping in the prediction step with the instances in the input. The second one involves mixing the mapping with some ensemble algorithm.

There are many ways to mix these two methods. The ones that we have studied are based on the Random Forest and are the most straightforward, but other ways could be studied in a future work.

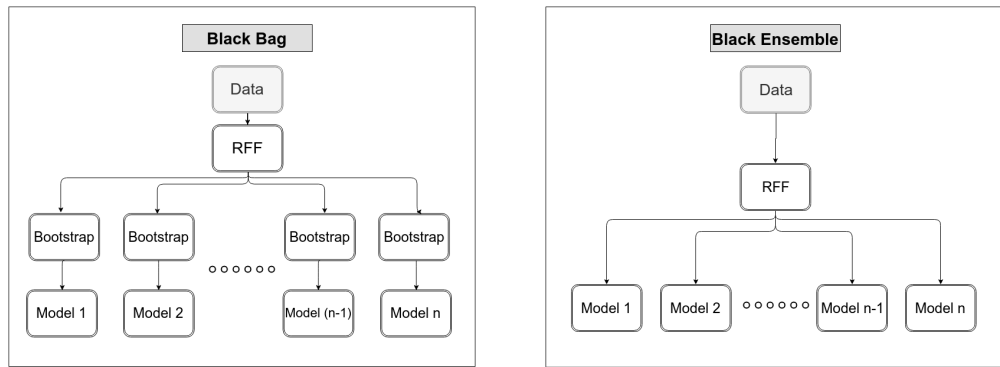
The first thing to decide is where to place the random mapping in the whole process. We have taken two different approaches: we can think of an ensemble as a Black Box, where we can only affect the inputs and the outputs, or we can differentiate the different parts it contain, like in a White Box.

Then, there's what kind of ensemble to use. We've chosen to use the Bagging technique, which is the one that Random Forest uses, but we've also defined a modified version which doesn't perform the Bootstrap. This is because we expected that using Bootstrap together with a random mapping would produce too much randomness in the data, affecting negatively the overall accuracy. We've called the first method a "Bag", and the second method an "Ensemble".

From the combination of these approaches we have defined four different ways of mixing: "Black Bag", "Black Ensemble", "White Bag" and "White ensemble". See 3.1.

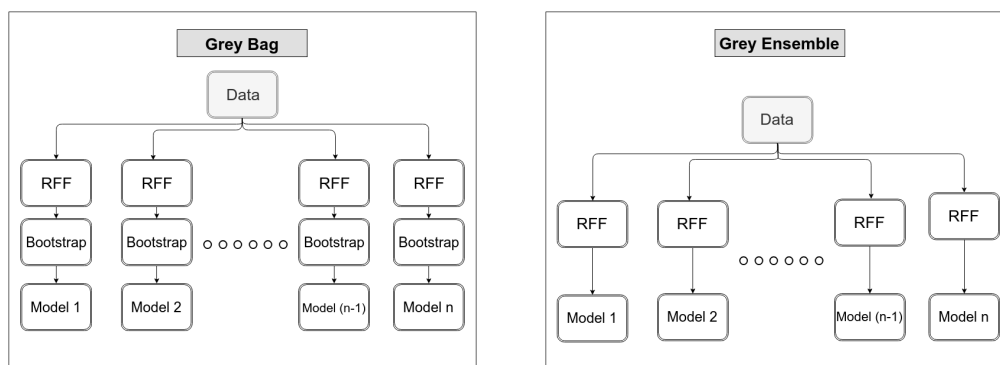
For this project we have used Random Fourier Features and the Nyström method to approximate the Radial Basis Function (RBF) Kernel with well known Machine Learning algorithms: Support Vector Machine, Logistic Regression and Decision Tree. We study how we can use them to increase the accuracy that we can achieve with them and the computational costs.

### 3.1.1 State of the art on las RFF



Black Bag. Each model receives data from the same Random Mapping, but they take a different resampling with replacement.

Black Ensemble. Each model receives exactly the same data, from just one Random Mapping. Only useful for randomized models.



White Bag. Each model receives data from a different Random Mapping, which has then had a resampling with replacement.

White Ensemble. Each model receives data from a different Random Mapping.

FIGURE 3.1:

Four different ways to mix Random Fourier Features with ensemble methods.

- Se ha trabajado poco con ellas. Solo he encontrado 2 usos:
  - Stacked kernel network (referencia): usarlas junto a una red neuronal para tener más niveles de aprendizaje no lineal
  - RFF with SVM (referencia): usar una SVM sin kernel con los datos mapeados usando RFF

### 3.1.2 State of the art con las Nyström

## 3.2 Hyper-parameters

- Existen los siguientes:
  - min-impurity-decrease para DT
  - C para SVM

- gamma para RFF y Nyström
- cantidad de features para RFF y Nyström
- cantidad de estimadores para ensembles
- Hemos usado los siguientes valores:
  - Cantidad de features a 500
  - Cantidad de estimadores a 50
  - En modelos simples, el parámetro por crossvalidation
  - En modelos simples con RFF, el parámetro por crossvalidation y una gamma que sobreajuste
  - En modelos con ensemble, parámetros que sobreajusten y la gamma por crossvalidation
  - En RBF-SVM, la gamma por gamest y el parámetro por crossvalidation

With the models defined in this project there are many hyper-parameters to tune the models. These are the hyper-parameters that have been used in the experiments:

**Number of features extracted from the kernel** The higher this value is, the better the approximation of the kernel function. We have fixed a value of 500, since this is not important as long as it is not too low.

**Ammount of estimators** Having a large number of estimator doesn't affect negatively the accuracy obtained, but increases the computation time, so the ideal number depends on the computational resources available. For this project we have picked 50 estimators for each Bag/Ensemble.

**Gamma parameter of the RBF Kernel** A higher value will generate a higher overfit. There is a fast method to find a suitable value for this parameter, explained in [4]. We have chosen this estimation.

**Parameters of the simple models** Decision Tree use *min\_impurity\_decrease* to tune the overfit, and SVM use a penalty *C* to do the same. When we train these models without any ensemble, we use Cross-Validation to find a suitable value. When we train an ensemble of these models, as we want them to overfit we set *min\_impurity\_decrease* to 0 and *C* to 1000, whih is enoght to achieve it.

### 3.3 Hypothesis

1. Podemos aproximar bien una RBF-SVM
2. Puede tener sentido hacer ensembles con otros modelos a DT
3. RFF + Bootstrap puede ser malo

4. Si el modelo no se basa en productos escalares no se beneficiará tanto

El paper que hace Linear SVM con RFF no está indexado en ningún sitio. Es solo un pdf que hay por la red

We had proposed these four hypothesis:

**1. It is possible to achieve an accuracy close to using the RBF Kernel but with a lower cost**

When the number of instances available is too big it is not possible to use an SVM with the RBF kernel, because the cost is cubic with the number of instances, and the optimization problem is too complex. A linear Kernel needs less time, but it may not be suitable for some problems, since data may not be easy to separate.

If we could first map the data to the new feature space, we could then feed a Linear SVM with it and have the same accuracy with less costs. But this can't be done with the RBF kernel, since the new feature space has infinite dimensions. However, with the use of RFF and Nyström, we can get an approximation of the feature space of the RBF. Using them with a Linear SVM could increase the accuracy on some datasets at almost the same cost.

**2. It could make sense to train ensembles of SVM and Logistic Regression algorithms**

Since these models are very stable, having an ensemble of them is useless: all of them will always predict the same answer. There are some methods to randomize a little bit the data, such as Bootstrap, but with these models it is not enough.

Since RFF and Nyström generate a random mapping of the data, we can achieve a higher level of randomization of the data, while still being a good representation of the real data. Random Mapping can allow us to build ensembles with these two models, increasing the overall accuracy at the expense of some computation time.

**3. Bootstrap together with a Random Mapping may be too much randomization**

With a simple mix of Bagging with RFF there are two different sources of randomness. For the one hand, Bootstrap generates a random sample of the data with replacement, and on the other hand, RFF and Nyström perform a Random Mapping of the data to a different feature space.

It is possible that for some models, this is too much randomization of the data, and it could have a bad effect on the learning process.

**4. Decision Tree does not benefit from RFF and Nyström as much as Logit and SVM do**

Kernels were originally used on Support Vector Machines because they were a fast way to implicitly compute the inner product of two vectors

in a feature space where data was separable by an hyper-plane. They were useful because SVM just needed the inner products of their input to work.

RFF and Nyström are ways to explicitly compute an approximation of that mapping, which doesn't necessarily fits the requirements of Decision Tree, which has nothing to do with the inner products. That's the reason why Decision Tree may not benefit so much of these Random Mappings.

### 3.3.1 Experiments Proposal

1. Hipótesis: Aproximar RBF-SVM
  - 1.1. Comparar una RBF-SVM con SVM normal que use RFF
2. Hipótesis: Ensembles con otros
  - 2.1. Logit normal vs. Logit con RFF
  - 2.2. Logit normal vs. Logit con RFF Black Bag
  - 2.3. Logit normal vs. Logit con RFF Grey Bag
  - 2.4. Logit normal vs. Logit con RFF Grey Ensemble

---

  - 2.5. Linear-SVM vs Linear-SVM con RFF
  - 2.6. Linear-SVM vs Linear-SVM con RFF Black Bag
  - 2.7. Linear-SVM vs Linear-SVM con RFF Grey Bag
  - 2.8. Linear-SVM vs Linear-SVM con RFF Grey Ensemble
3. Hipótesis: RFF + Bootstrap
  - 3.1. Logit con RFF Grey Bag vs Logit con RFF Grey Ensemble
  - 3.2. Logit con RFF Black Bag vs Logit con RFF Black Ensemble (los dos con un solo estimador)

---

  - 3.3. Linear-SVM con RFF Grey Bag vs Linear-SVM con RFF Grey Ensemble
  - 3.4. Linear-SVM con RFF Black Bag vs Linear-SVM con RFF Black Ensemble (los
4. Hipótesis: DT + RFF
  - DT vs DT con RFF
  - DT vs DT con RFF Black Bag
  - DT vs DT con RFF Black Ensemble
  - DT vs DT con RFF Grey Bag
  - DT vs DT con RFF Grey Ensemble



There are two factors we need to check in order to accept or refute the hypothesis: we want to know if these methods increase the accuracy of the models (and if so, how much) and also at what costs does it come, the increase of training time. The 7 first datasets (Covertypes, Digits, Fall Detection, Pen Digits, Satellite, Segment and Vowel) are suitable to measure the increase of the accuracy, but not to measure the time, since the overhead caused by these techniques is too much compared with training with a problem as small as that. The other two datasets (MNIST and Fashion MNIST) have a larger number of instances, and are expected to reflect better the benefits of using these methods. However, some of the models proposed in the experiments are very expensive to train on very large datasets, and it was not possible to run the whole training with them. Some simplifications have been made to extrapolate the real training time.

To check the hypothesis previously suggested we have proposed a set of experiments:

### Hypothesis 1

In this hypothesis we assumed that we could train an SVM with results comparable to the ones using the RBF Kernel, but with much less training time. To check that we have defined the experiment 1.1, where we compare four different models, an SVM without a kernel, an SVM with the RBF kernel and two SVM without kernel but using the random mappings, RFF and Nyström. In order to accept the hypothesis we need to see that the models using RFF and Nyström significantly increase the accuracy compared to the single SVM. For the method to be useful it would be needed for the training time to be much less using the random mappings compared to the RBF-SVM on large datasets. To check that we will need to focus on MNIST and Fashion-MNIST, since they are the only large datasets.

Theoretically, for RBF-SVM a cross-validation should have been performed to find a suitable value for  $C$ , but the training time needed to do that for the large datasets is too much. To avoid that, we have used the same value found for the SVM with RFF in MNIST and Fashion-MNIST, and have approximated the time that would have been needed if cross-validation had been performed.

Results of this experiment can be seen in [Appendix A](#)

### Hypothesis 2

For this hypothesis it was asserted that training an ensemble of SVM and Logistic Regression using a random mapping could increase significantly the accuracy obtained. Nevertheless, the training time is expected to be significantly greater since training an ensemble requires a lot more steps.

Experiments 2.1 through 2.4 were defined to check that. The first two are for Logistic Regression and the others for SVM. 2.1 and 2.3 will compare a single model alone against a single model using a random mapping. 2.2 and 2.4 will compare a single model against the model with Black Bag, White Bag and White Ensemble. The Black Ensemble model is not tested because these estimators are not randomized, and each one would predict the same answer. For the hypothesis to be accepted we need to see that we can achieve an increase in the accuracy using any of the ensemble methods proposed. Then we

will need to study the increase in the training time, and consider if it is worth it.

Results of these experiments can be seen in Appendices [B](#), [C](#), [D](#) and [E](#).

- Solo por aclarar alguna idea, incluir también una gráfica que compare un solo rff con el resto de ensembles, a ver si la mejoría viene por eso

### Hypothesis 3

With this hypothesis we wanted to know if using Bootstrap together with a random mapping will produce worse results than just using the random mapping. To check that we will compare the “Bags” against the “Ensembles”. Experiments 3.1 and 3.3 will compare “White Bag” and “White ensemble” for Logistic Regression and SVM respectively, and 3.2 and 3.4 will compare “Black Bag” and “Black Ensemble”. To make it fairer only one estimator has been trained, since “Black Ensemble” doesn’t benefit from using more.

In order to accept this hypothesis we would need to see a meaningful differences in the accuracy of these models.

Results of these experiments can be seen in Appendices [F](#), [G](#), [H](#) and [I](#).

### Hypothesis 4

In this hypothesis we asserted that Decision Tree, which doesn’t work with the inner product of the inputs, would not experiment a meaningful increase in the accuracy with the usage of a random mapping approximating the RBF feature space. Decision Tree can normally increase their accuracy by using a normal Bagging. That algorithm is called Random Forest. For this reason the ensemble methods using a random mapping should be compared against the Random Forest instead of a single Decision Tree.

Experiment 4.1 compares a single Decision Tree with one that uses a random mapping, and 4.2 compares the Random Forest with “Black Bag”, “Black Ensemble”, “White Bag” and “White Ensemble”

Results of these experiments can be seen in Appendices [J](#) and [K](#)

### Additionally

- Quizá, para los de hipótesis 2 puedo poner gráficas para contrastar RFF con Nystrom

## 3.4 Datasets

- 8 Datasets
- Normalizados
- Únicamente tienen variables numéricas, no categóricas
- Únicamente problemas de clasificación

TABLE 3.1: Information on the datasets used in this project

Dataset	N. Instances	N. Features	N. Classes
Coverttype[3]	4900	12	7
Digits[2]	5000	10	10
Fall Detection[14]	5000	6	6
MNIST[6]	5000	717	10
Pen Digits[1]	5000	16	10
Satellite[10]	5000	36	6
Segment[11]	2310	19	7
Vowel[5]	990	11	10
Fashion MNIST[12]	70000	784	10

- Algunas cosas particulares que he hecho:
  - Mezclar datos de train y de test para luego hacer mi propia separación
  - Cuando había poca presencia de una clase, hacer un resampling para igualar las cantidades
  - No trabajar cosas como el skiwness o los outliers
  - Eliminar columnas en las que todo eran 0
  - Reducir el conjunto de instancias

We have chosen 9 different datasets to run the experiments proposed. All of them are from classification problems since that was the scope of the project. All the data has been normalized to have a mean of 0 and a variance of 1. Since the kernel RBF can only work with numerical data, all categorical variables have been converted to integers. Here follows some of the preprocessing that has been done on the datasets:

**Coverttype** In the original dataset the classes were not balanced. So, we performed a resampling of the data, taking 700 random instances of each of the clases. There were 40 binary columns that were converted to a single numerical variable.

**MNIST** We took a subset of 5000 instances from the original dataset. Since in the new sample some of the columns contained only zeros, they were removed.



## Chapter 4

# Experimental Results

hi

### 4.1 Enfrentar resultados 2 a 2

#### Exp 1-1

- Nunca jamás conseguimos sacar mejores resultados que una RBF-SVM
- Pero los resultados que conseguimos son bastante parecidos en algunos casos
- En el experimento que hemos hecho, los tiempos de los métodos lineales son mucho mayores que los de la RBF-SVM
- Pero hay que atribuirlo al pequeño tamaño del dataset que hemos usado. En un dataset más grande sí se notaría
- Efectivamente, he realizado el experimento en fashion-mnist, y no hay comparación

#### Exp 2-1

- En general sí conseguimos mejorar sustancialmente un logit
- Pero los tiempos son un poco mayores
- Con datasets mayores, es de esperar que los tiempos no se disparen demasiado

#### Exp 2-2

- En los mismos datasets en los que un solo logit no había conseguido hacer nada, aquí tampoco han hecho nada
- Entre logit normal y logit con ensemble, la precisión es más o menos la misma, pero el tiempo es mucho mayor
- En el caso logit, realmente puede deberse a que no hemos hecho nada para sobreajustar más.

**2-3**

- Exactamente igual que el anterior

**2-4**

- Exactamente igual que el anterior

**2-5**

- Conseguimos mejorar sustancialmente la precisión de la mayoría de datasets
- Los tiempos son mayores, pero todavía son aceptables
- Suponemos que en datasets más grandes los tiempo serían más parecidos

**2-6**

- El dataset que se resistía se sigue resistiendo cuando hacemos un ensemble
- Ahora los tiempo sí que són muchísimo más grandes, quizá no sale a cuenta hacer ensemble
- Para ver si sale a cuenta hacer ensemble, habrá que verlo más adelante

**2-7**

- Exactamente igual que el anterior

**4.2 2-8**

- Exactamente igual que el anterior

**3-1**

- No podemos decir que haya una ensemble que claramente sea mejor que los demás
- Da la impresión que en algunos casos el grey es mejor que los otros, pero no es nada concluyente

**3-2**

- Da la impresión que los black ensembles son un poco mejor que los grey, pero no es nada concluyente

**3-3**

- Los resultados no son nada concluyentes

**3-4**

- Exactamente igual que el anterior

**4-1**

- Usar rff con un solo árbol no es nada beneficioso
- Tanto los errores como los tiempos son más altos

**4-2**

- En algunos casos ha mejorado sustancialmente, y en otros no
- Pero estamos comparando un solo árbol con 50: claramente los 50 tendrían que ser siempre mejores, y ese no es el caso
- Los tiempos no tienen ninguna comparación
- En algunos casos, incluso empeora el usar un ensemble

**4-3**

- Exactamente igual que el anterior

**4-4**

- Exactamente igual que el anterior

**4-5**

- Exactamente igual que el anterior

**4.3 Contrastar hipótesis con resultados**





## Chapter 5

# Conclusion and Future Directions

- Problemas de regresión
- Aproximar otros kernels a RBF
- Ver el comportamiento con problemas que no sean tan bonitos (con missings, clases desbalanceadas, etc)
- Otros tipos de ensembles, como el boosting

- Todo lo que hemos hecho solo sirve con datasets muy grandes. Para pequeños, en general salimos perdiendo
- Hemos conseguido hacerle un boost a logit
- Si usamos RFF, en general no sale nada a cuenta hacer un ensemble. No se mejora demasiado
- Sí que podemos aproximar una RBF-SVM con una lineal a coste lineal
- No hemos observado ninguna diferencia significativa entre usar RFF o Nyström

En general los únicos éxitos de este trabajo son:

- Ahora podemos hacer un ensemble de logit y de svm, que antes no se podría
- también hemos conseguido mejorar un poco un solo logit y svm
- Hemos aprendido que da igual el tipo de ensemble que cojamos



## Chapter 6

# Sustainability Report

### 6.1 Environmental

### 6.2 Economic

### 6.3 Social

#### 6.3.1 Impacto Personal

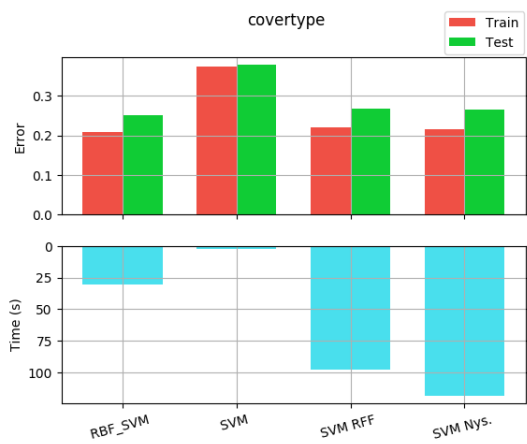
#### 6.3.2 Impacto Social

#### 6.3.3 Riesgos Sociales

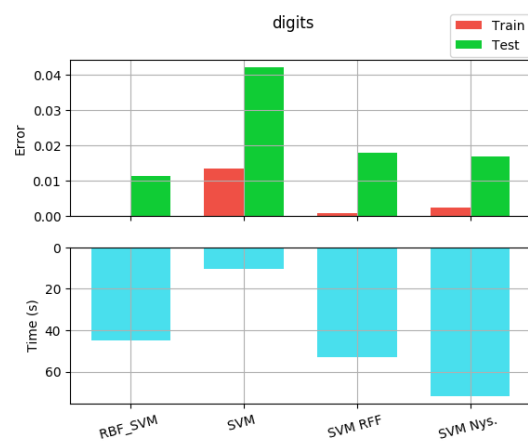


## Appendix A

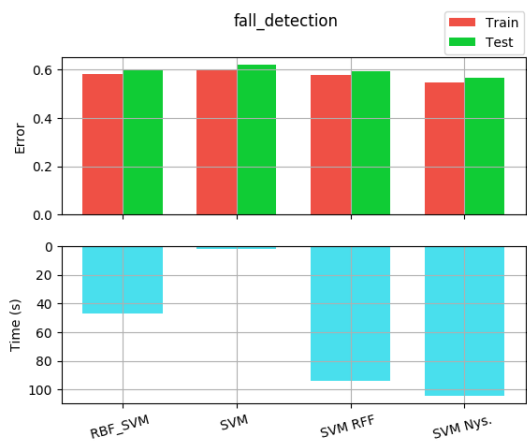
# Results of experiment 1.1



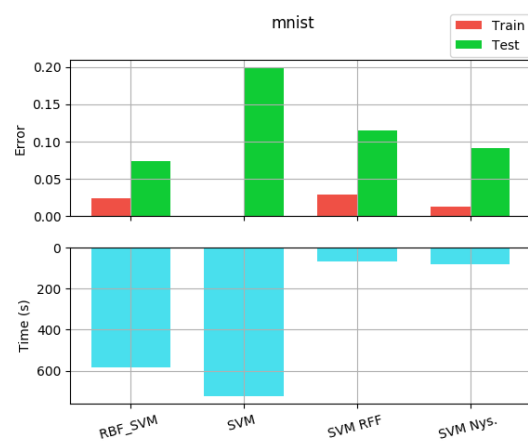
prueba covertype



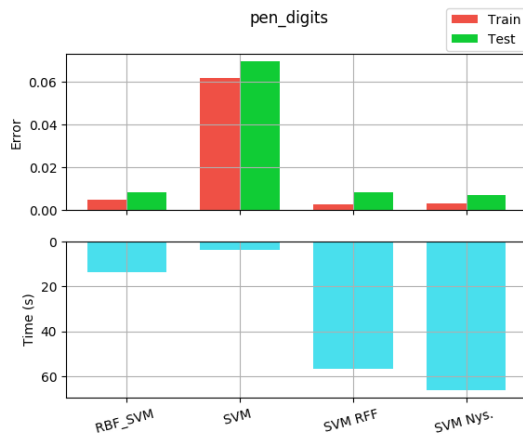
prueba digits



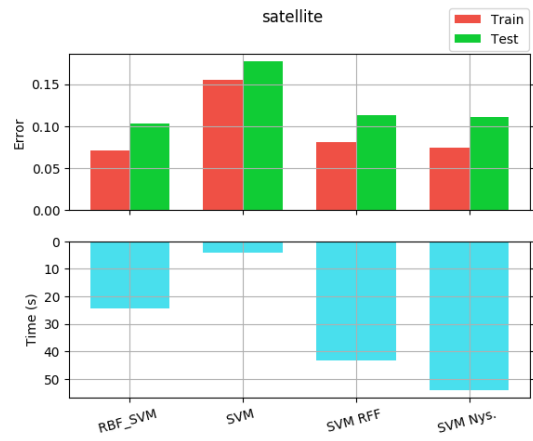
prueba fall-detection



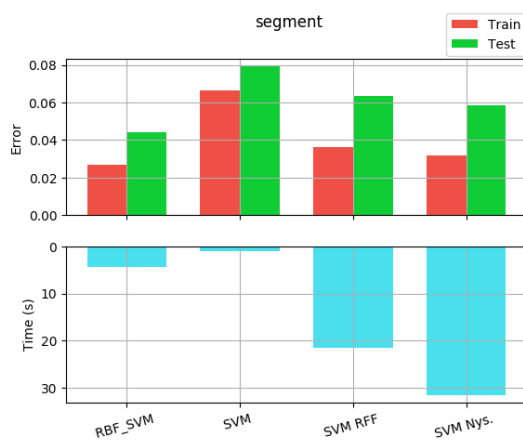
prueba mnist



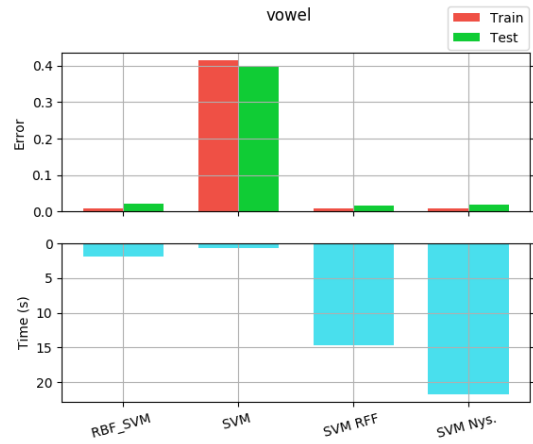
prueba pen-digits



prueba satellite



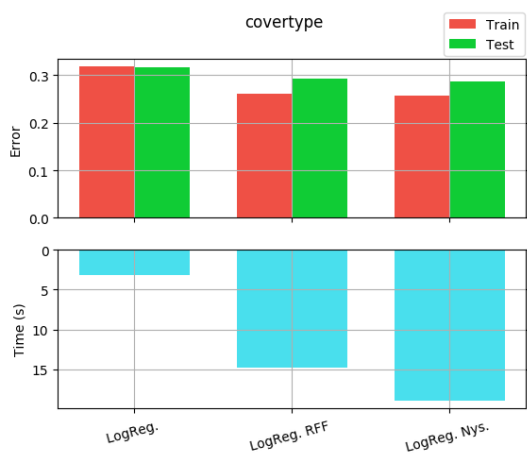
prueba segment



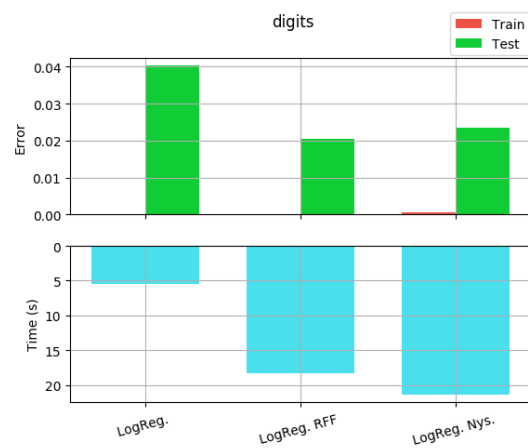
prueba vowel

## Appendix B

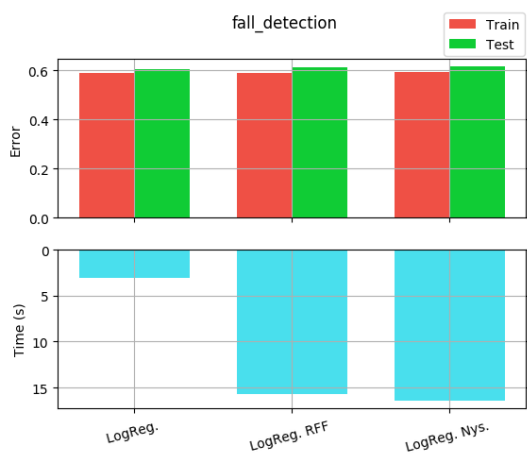
# Results of experiment 2.1



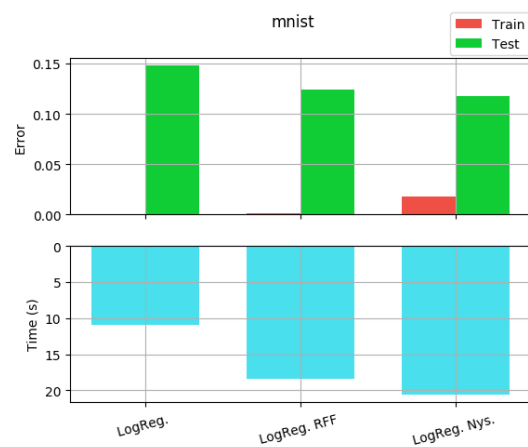
prueba covertypes



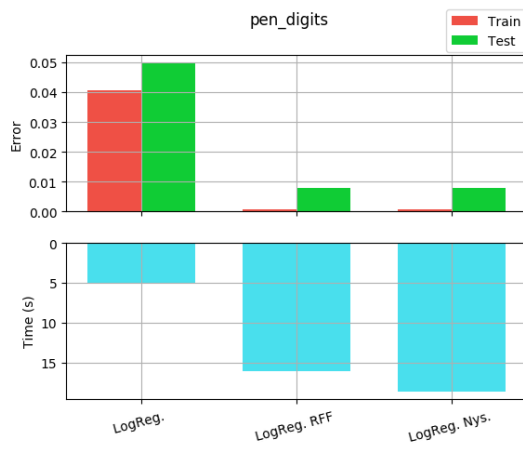
prueba digits



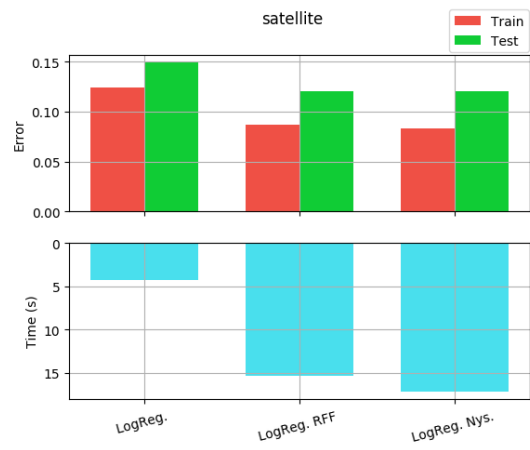
prueba fall-detection



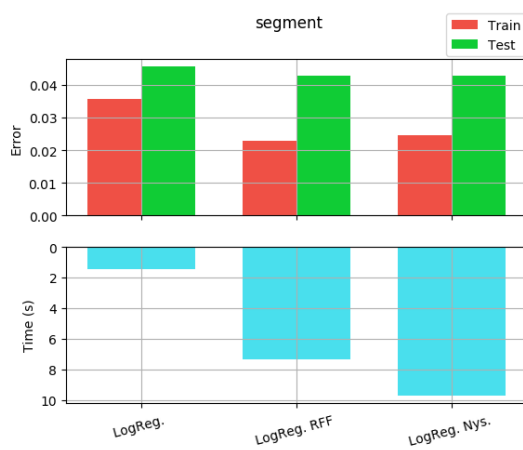
prueba mnist



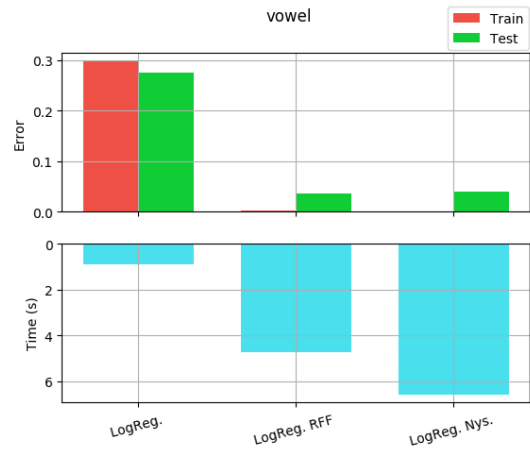
prueba pen-digits



prueba satellite



prueba segment

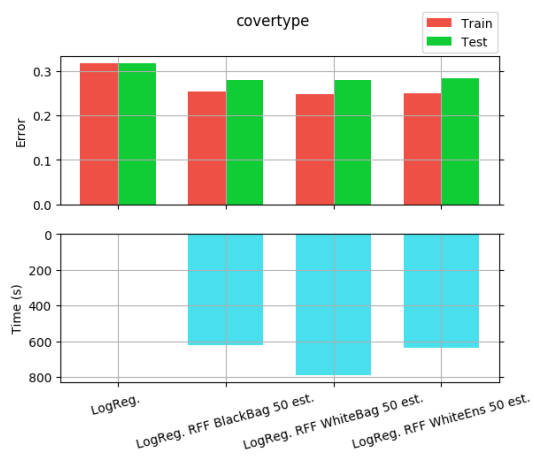


prueba vowel

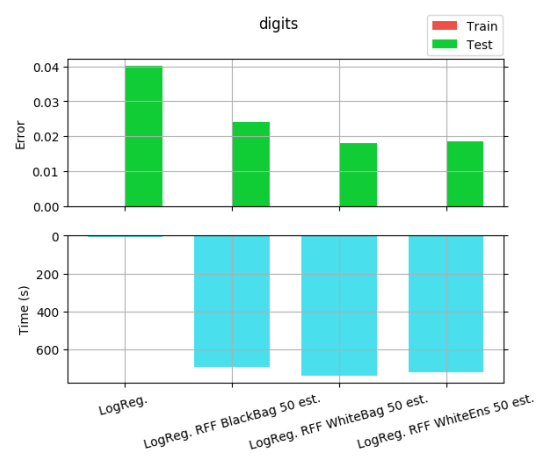


## Appendix C

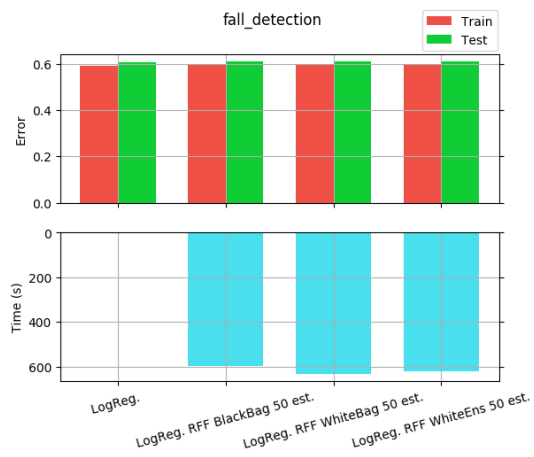
# Results of experiment 2.2



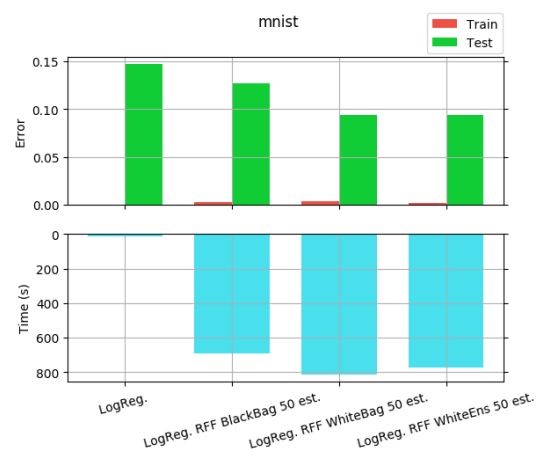
prueba covertypes



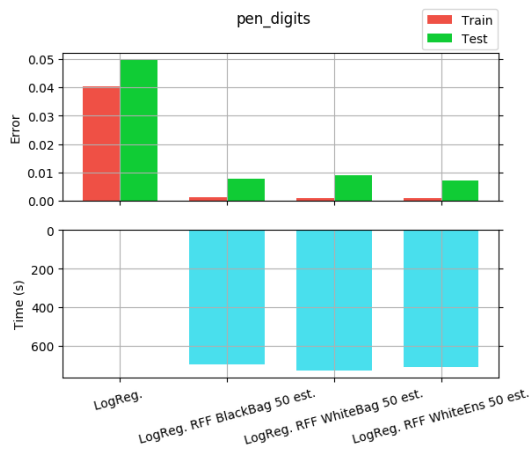
prueba digits



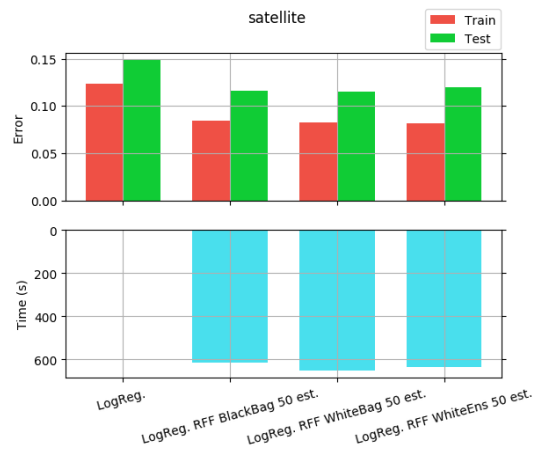
prueba fall-detection



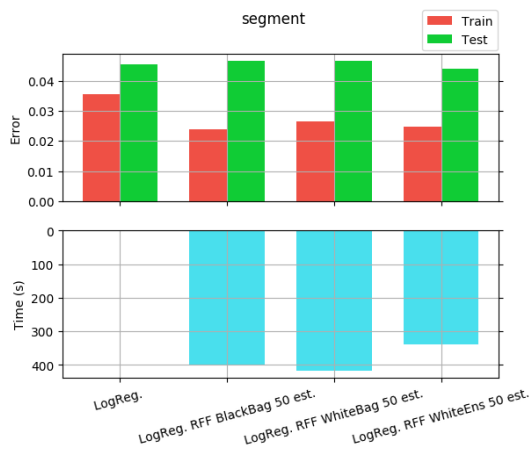
prueba mnist



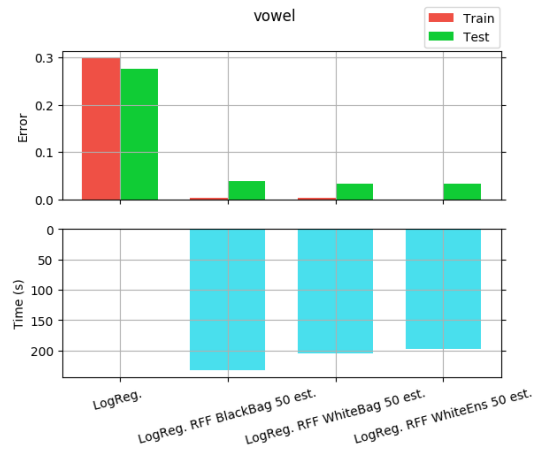
prueba pen-digits



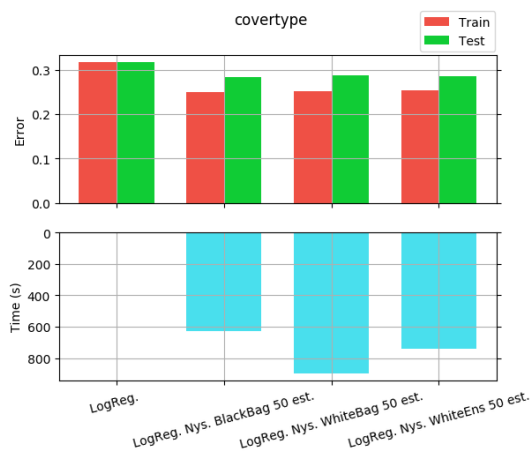
prueba satellite



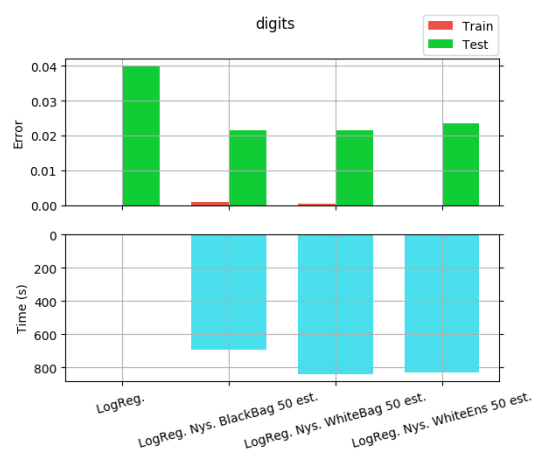
prueba segment



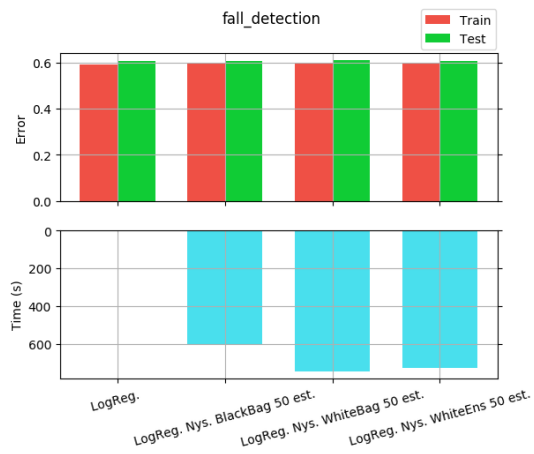
prueba vowel



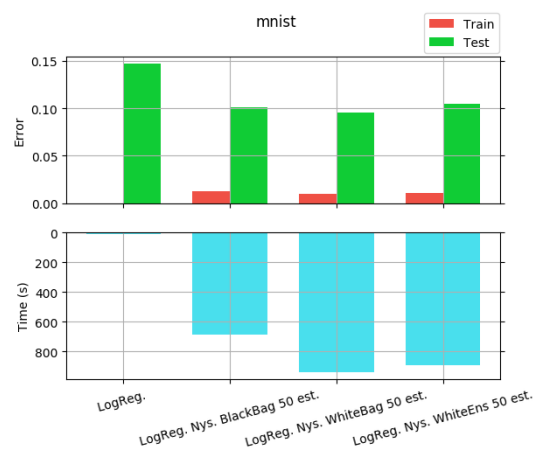
prueba coverttype



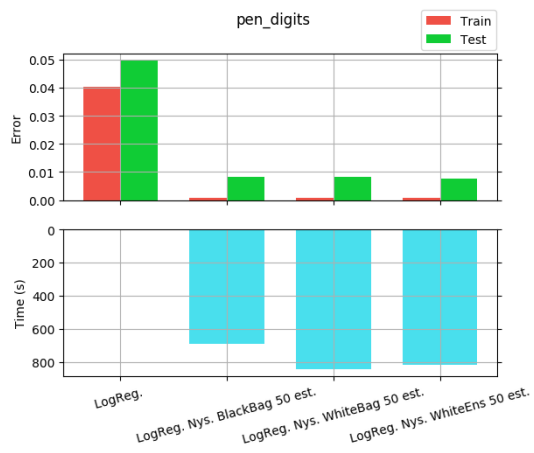
prueba digits



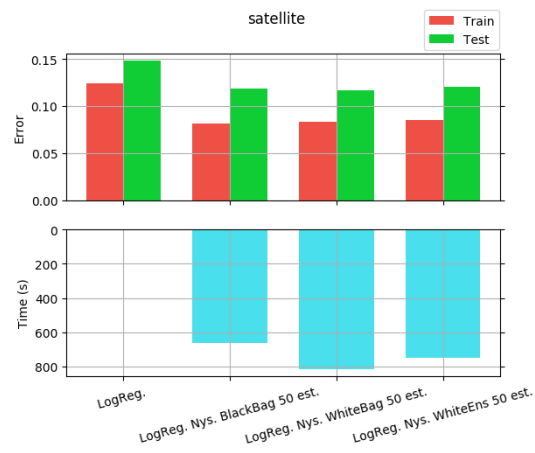
prueba fall-detection



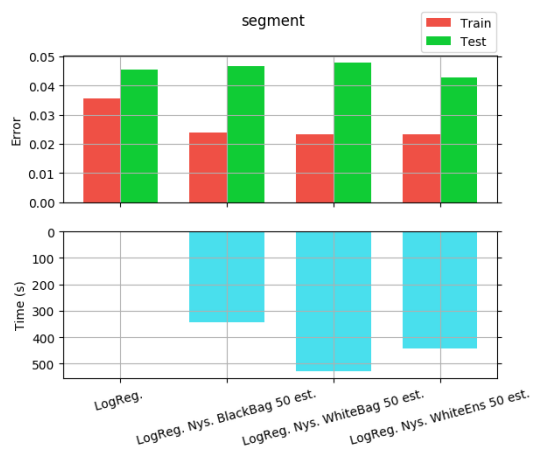
prueba mnist



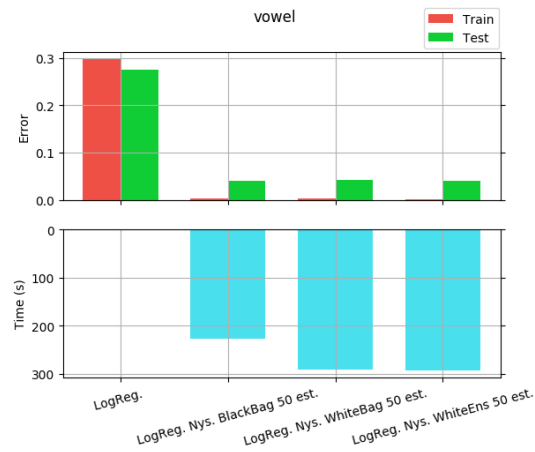
prueba pen-digits



prueba satellite



prueba segment

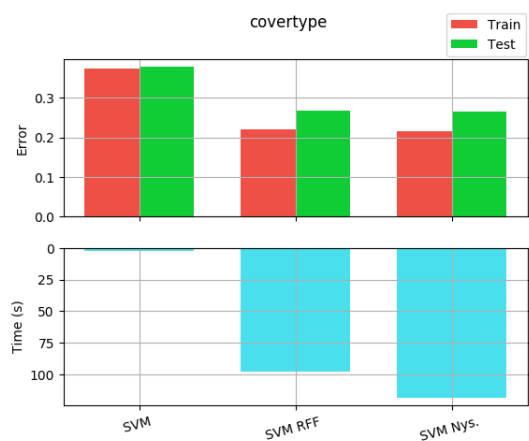


prueba vowel

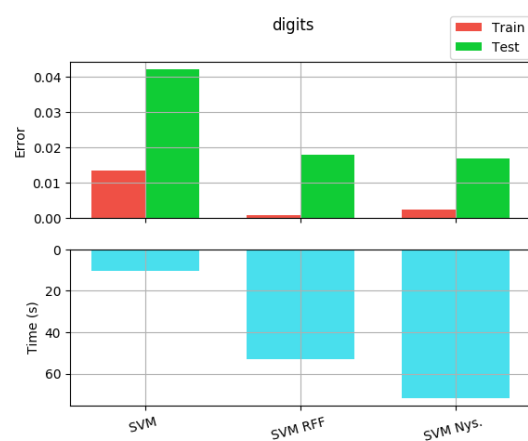


## Appendix D

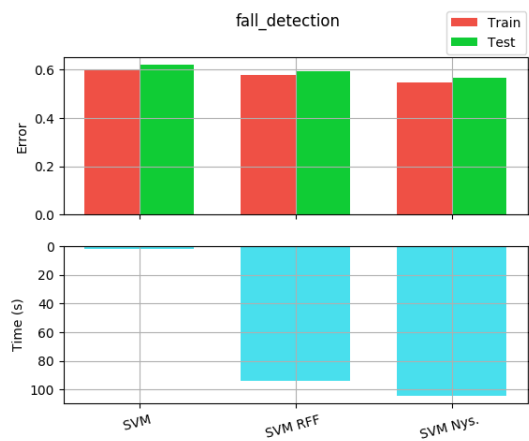
# Results of experiment 2.3



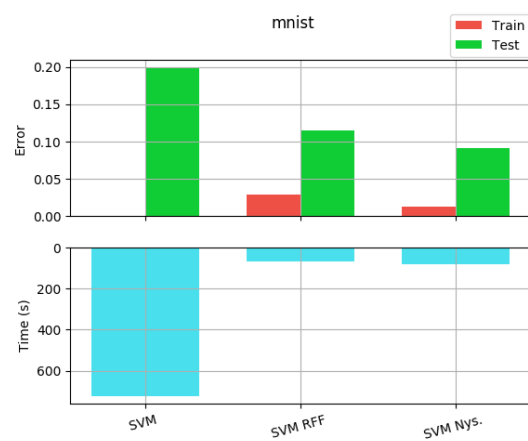
prueba covertype



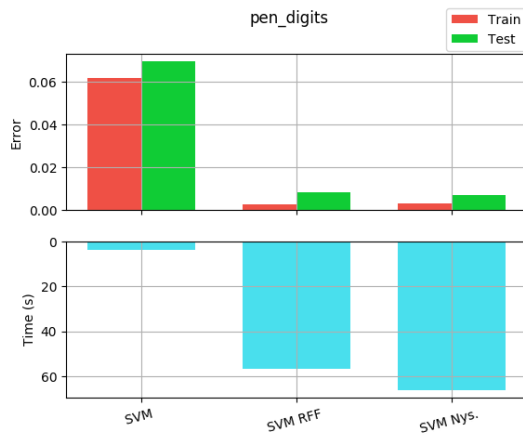
prueba digits



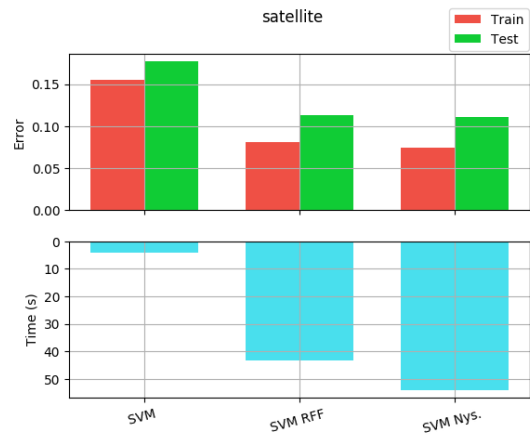
prueba fall-detection



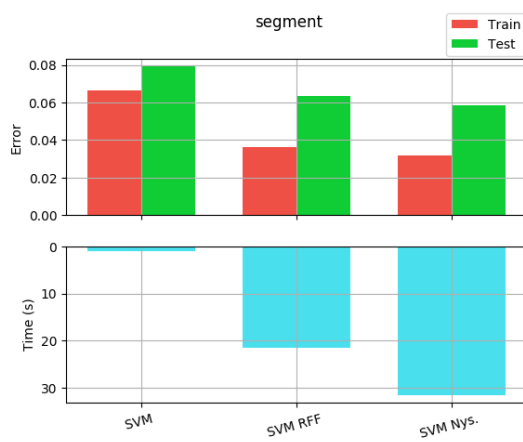
prueba mnist



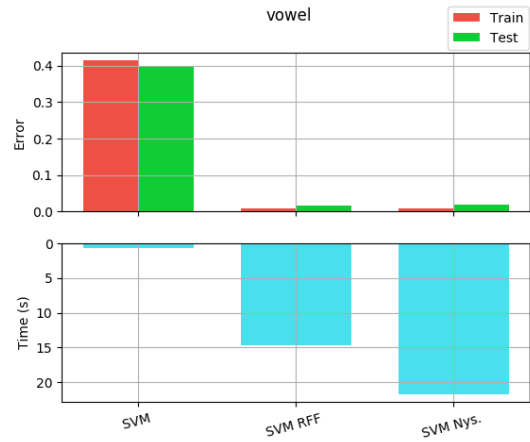
prueba pen-digits



prueba satellite



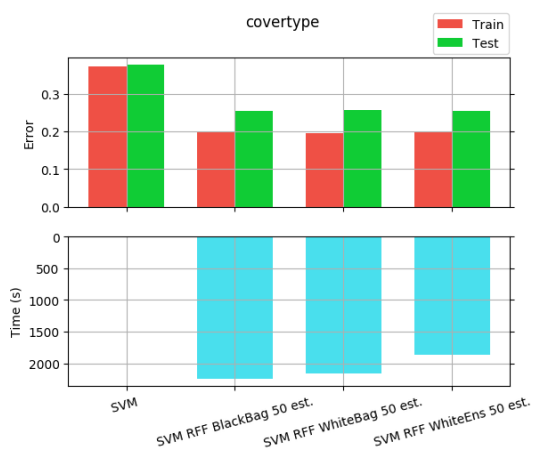
prueba segment



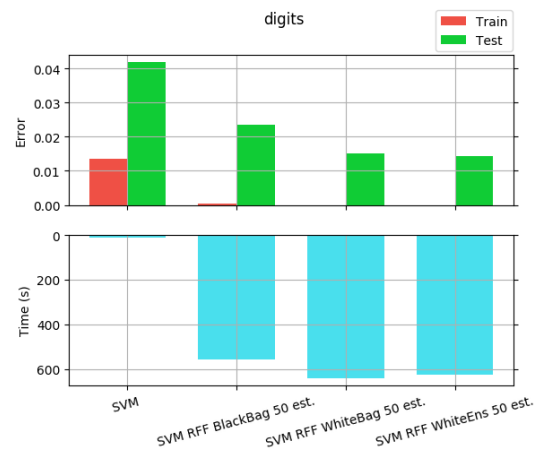
prueba vowel

## Appendix E

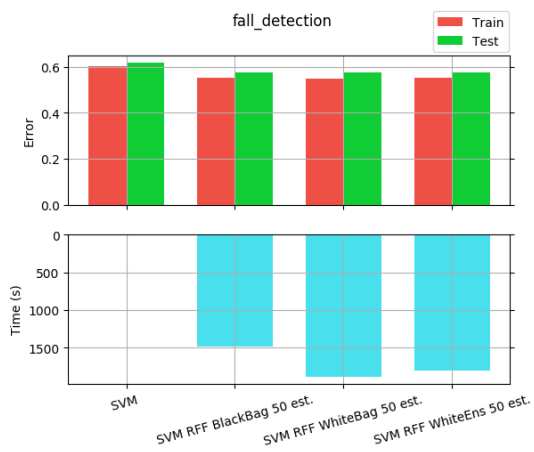
# Results of experiment 2.4



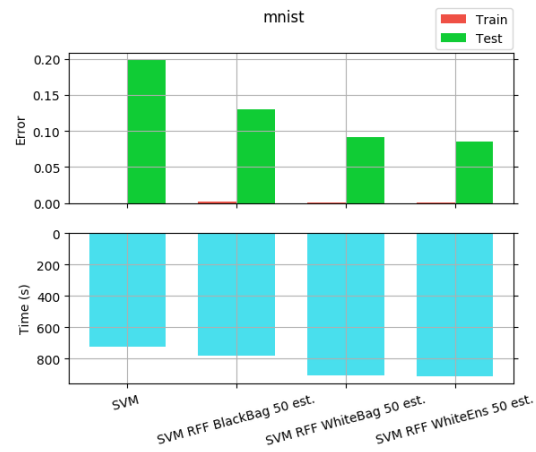
prueba covertypes



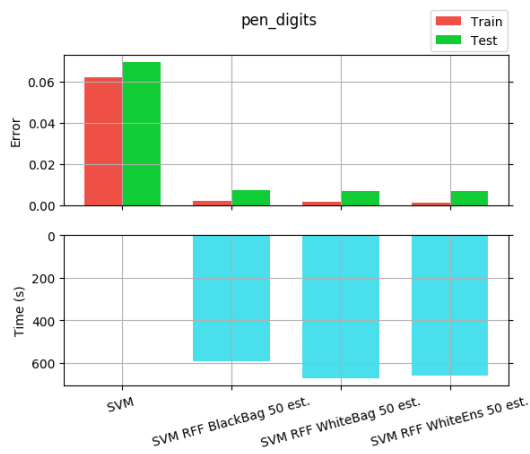
prueba digits



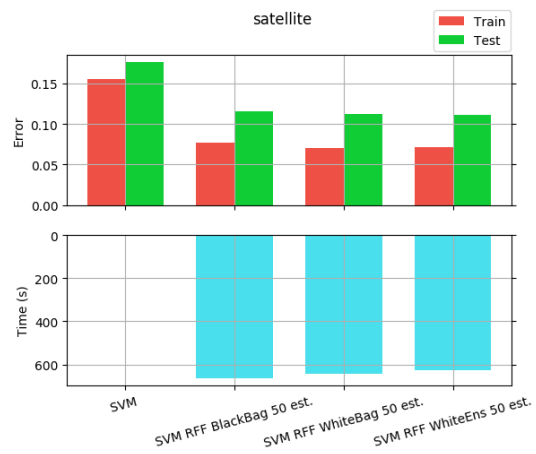
prueba fall-detection



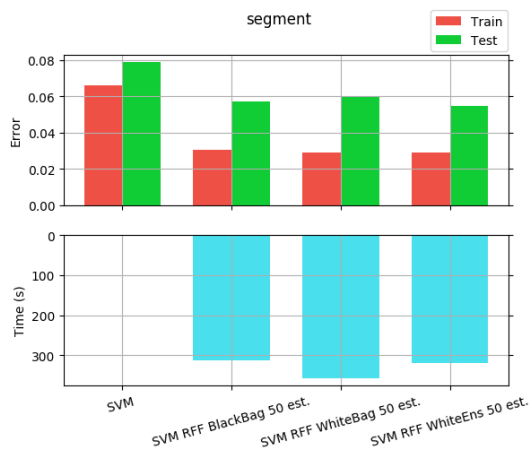
prueba mnist



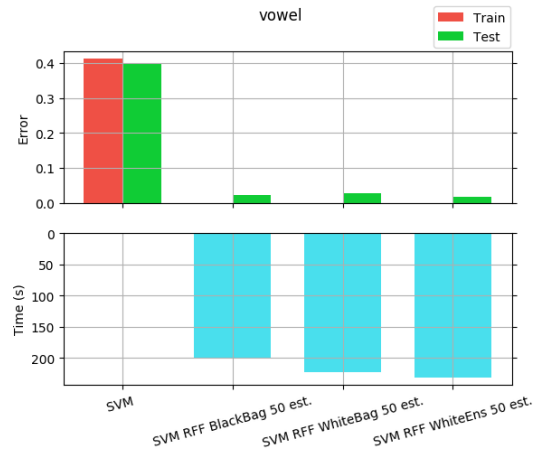
prueba pen-digits



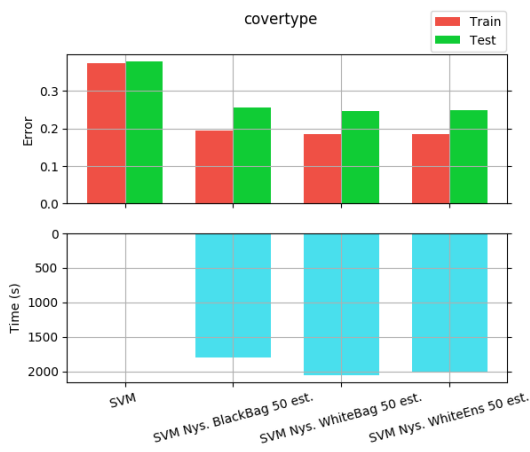
prueba satellite



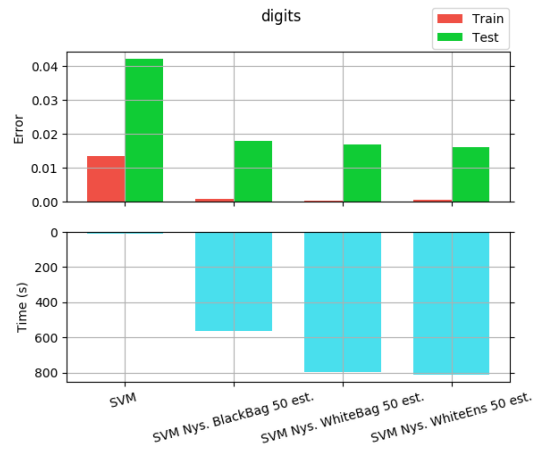
prueba segment



prueba vowel

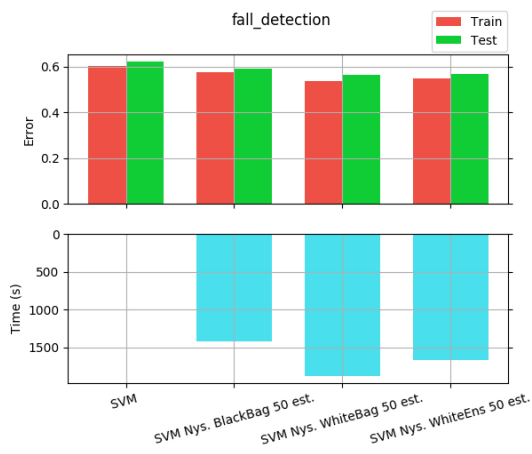


prueba covertype

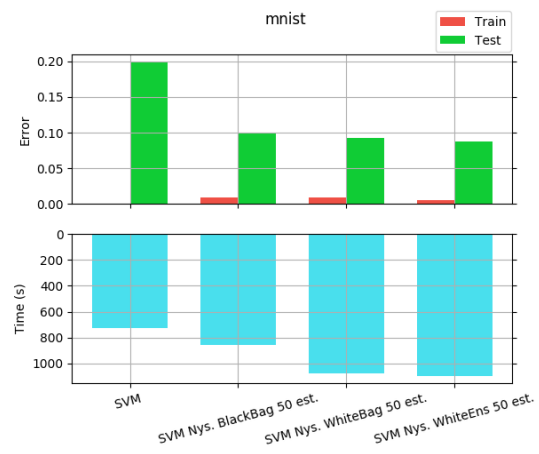


prueba digits

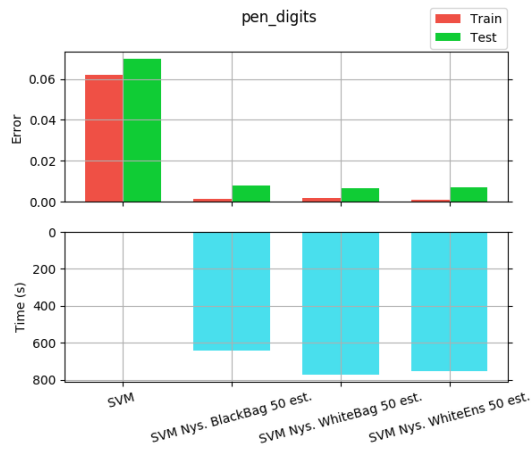




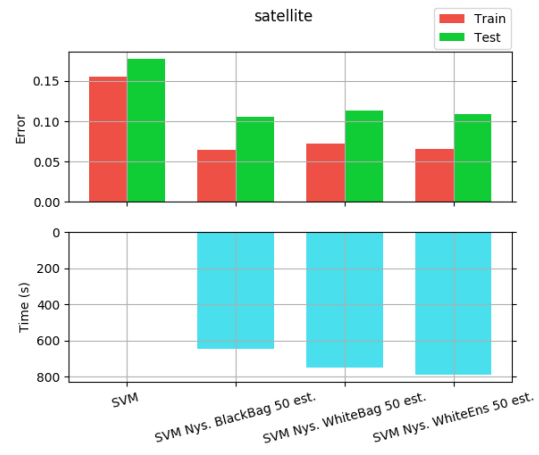
prueba fall-detection



prueba mnist



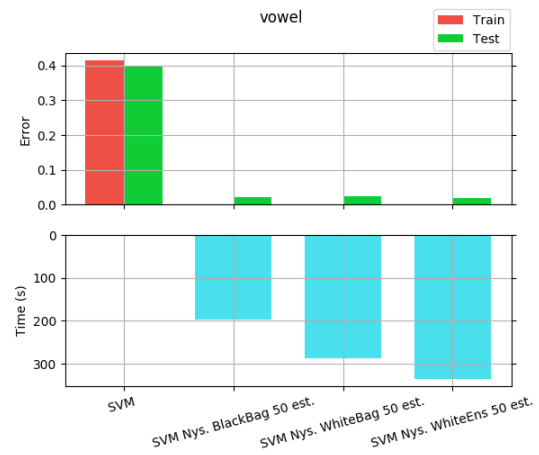
prueba pen-digits



prueba satellite



prueba segment

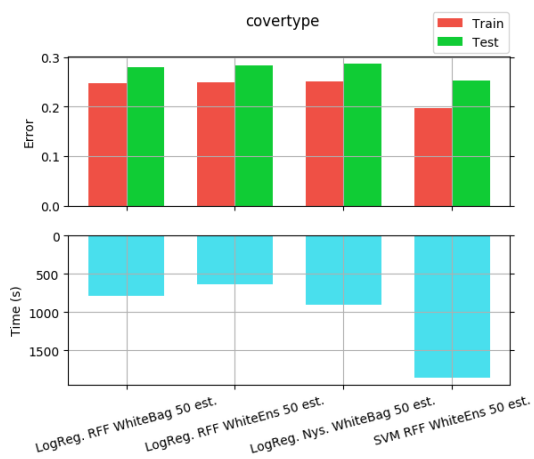


prueba vowel

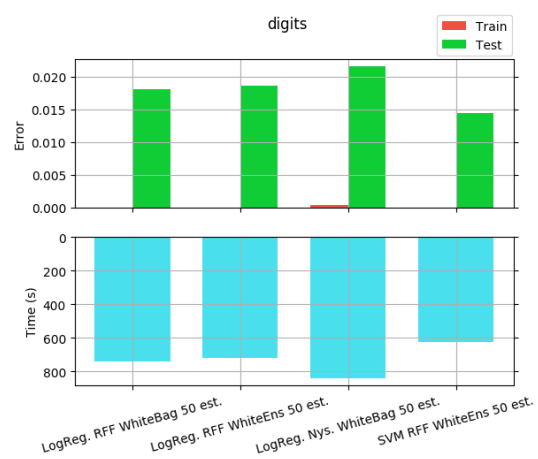


## Appendix F

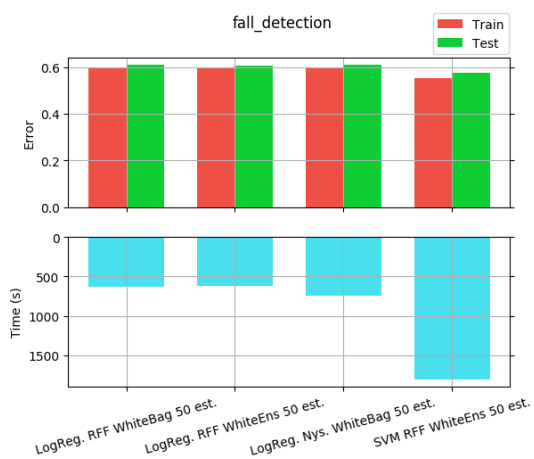
# Results of experiment 3.1



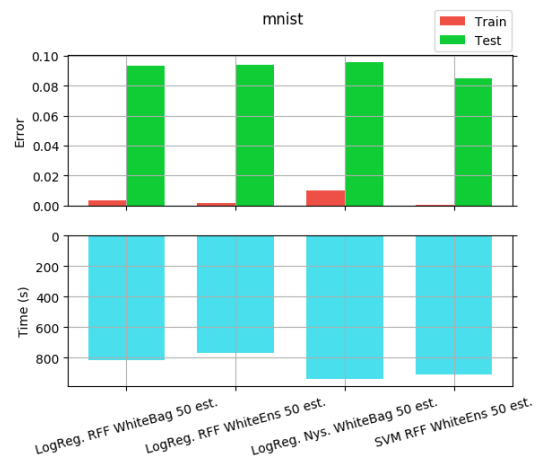
prueba covertype



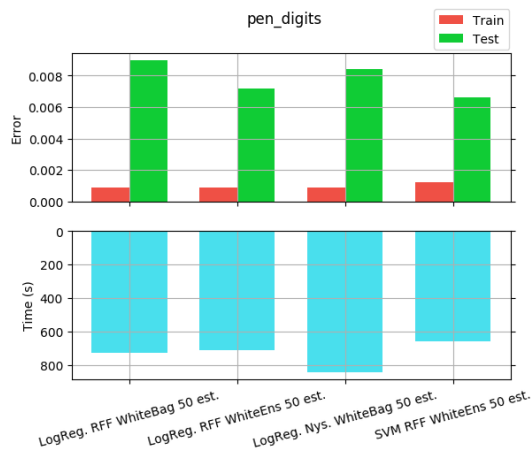
prueba digits



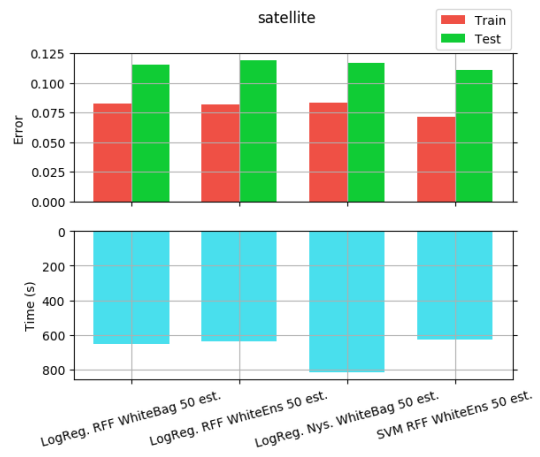
prueba fall-detection



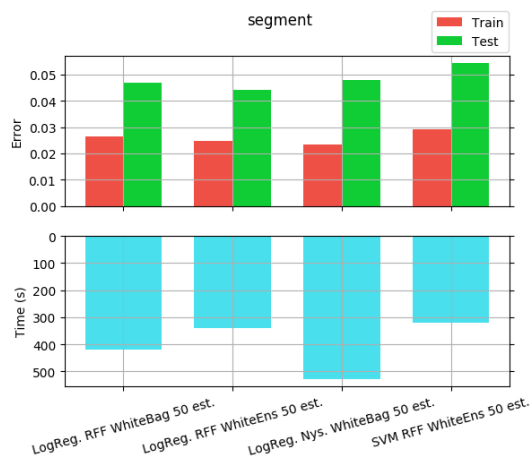
prueba mnist



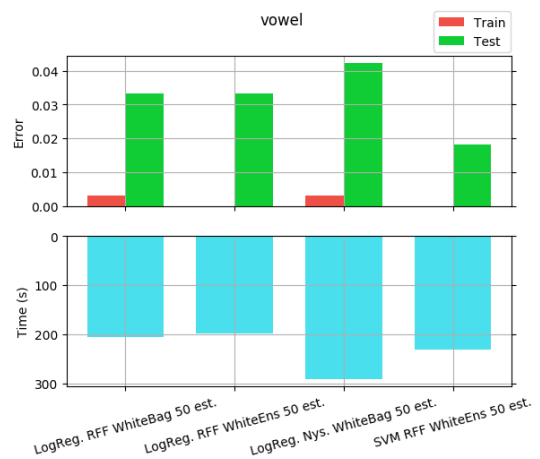
prueba pen-digits



prueba satellite



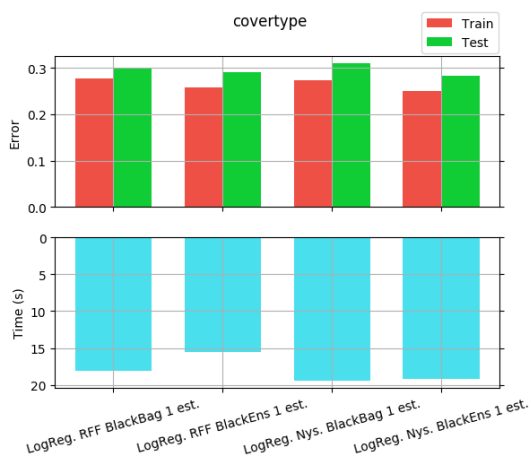
prueba segment



prueba vowel

## Appendix G

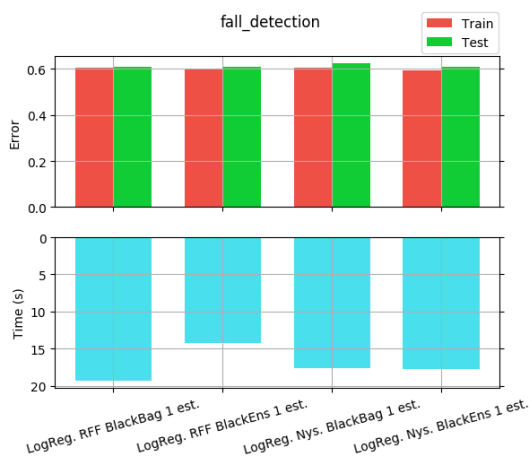
# Results of experiment 3.2



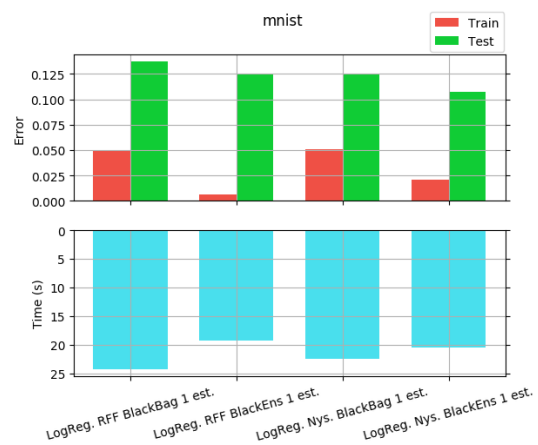
prueba covertypes



prueba digits



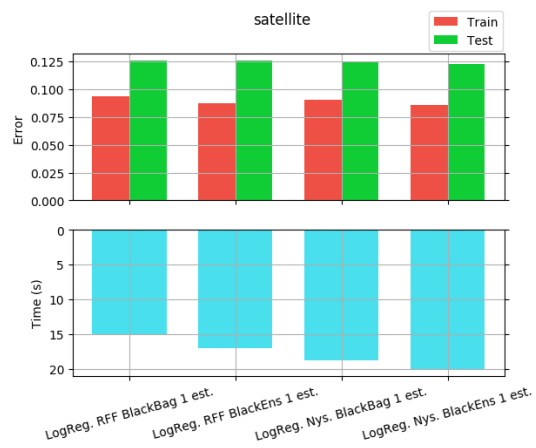
prueba fall-detection



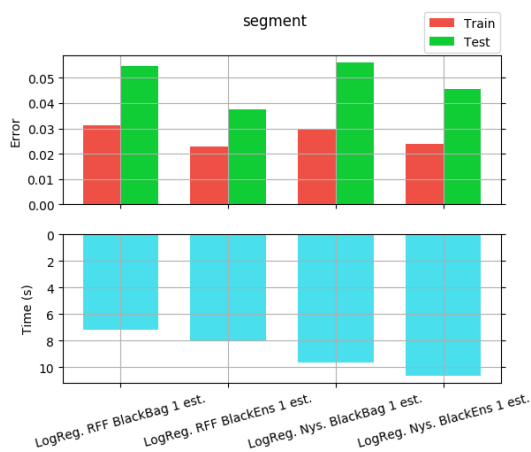
prueba mnist



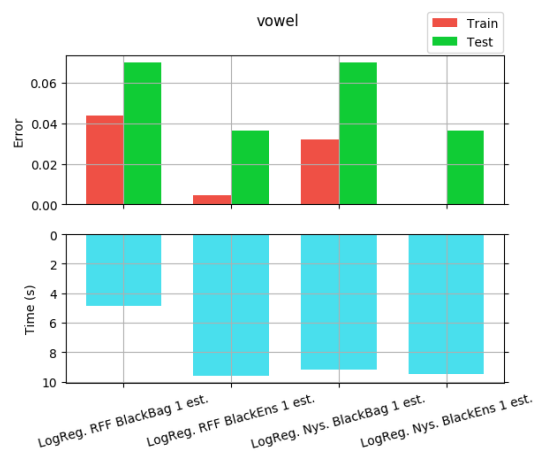
prueba pen-digits



prueba satellite



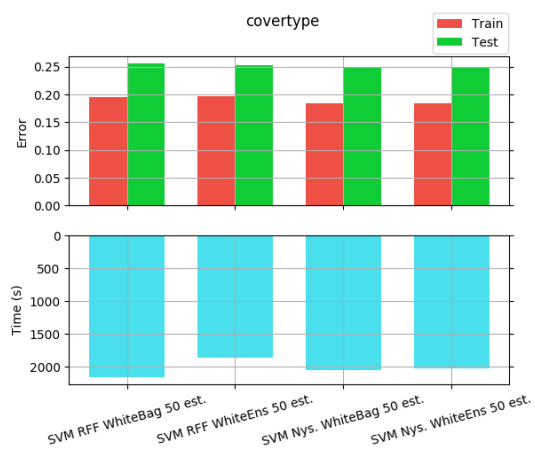
prueba segment



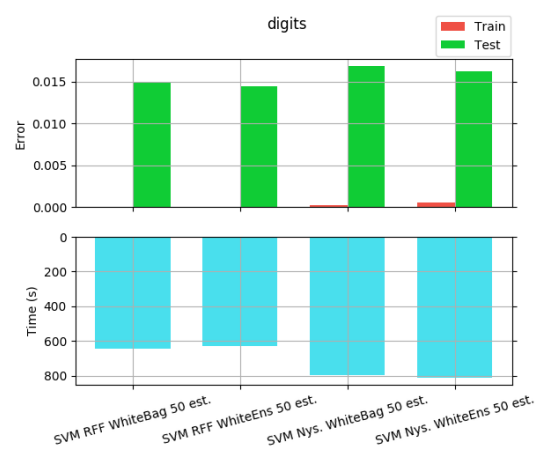
prueba vowel

## Appendix H

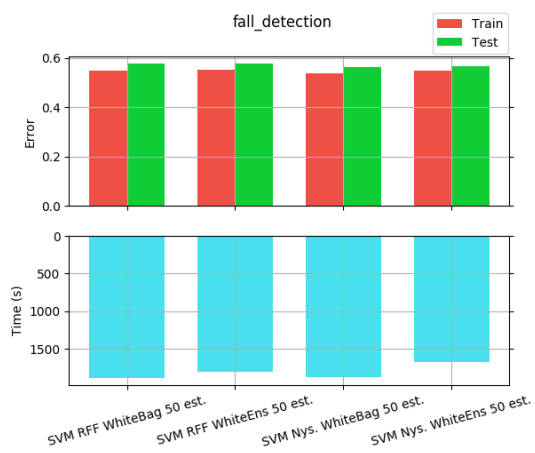
# Results of experiment 3.3



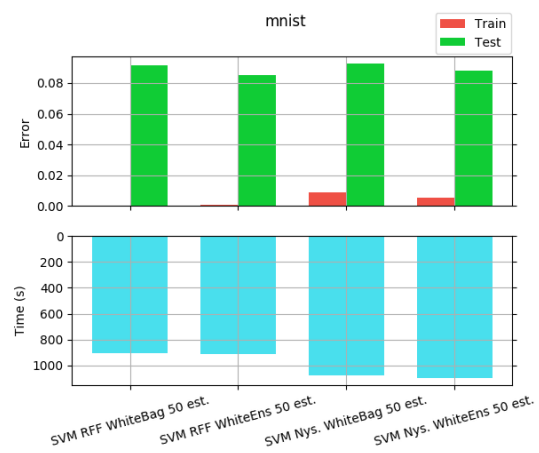
prueba covertypes



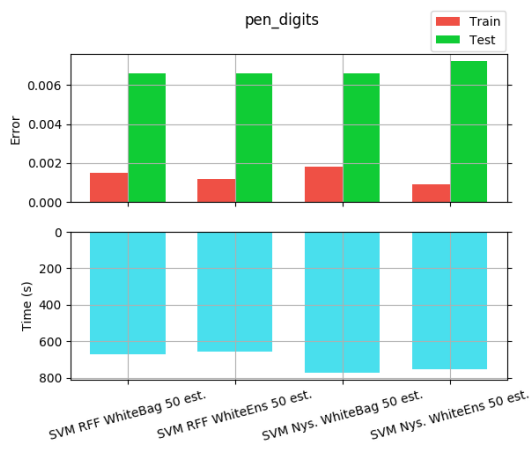
prueba digits



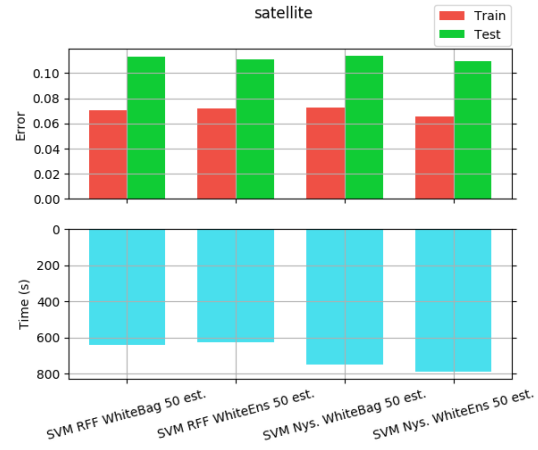
prueba fall-detection



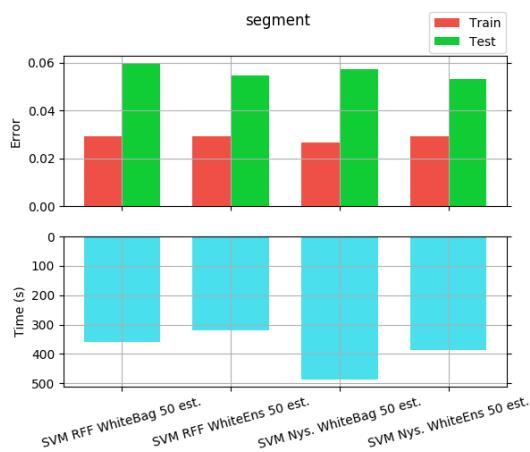
prueba mnist



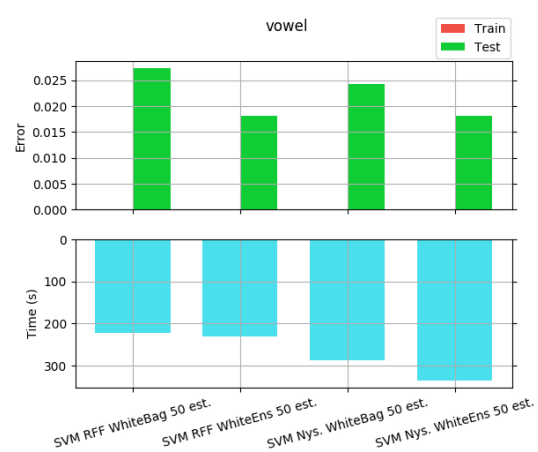
prueba pen-digits



prueba satellite



prueba segment

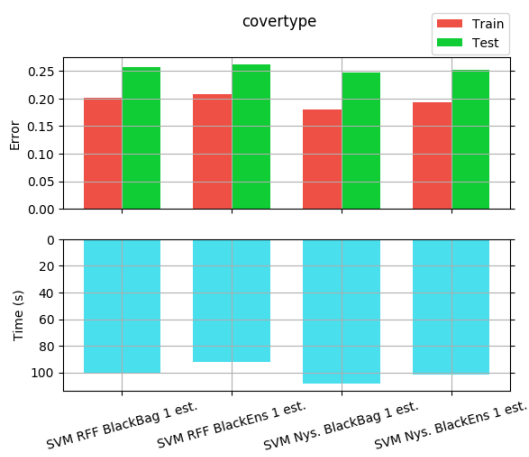


prueba vowel

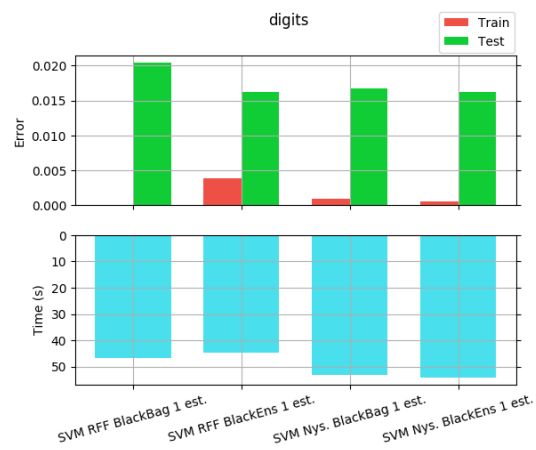


## Appendix I

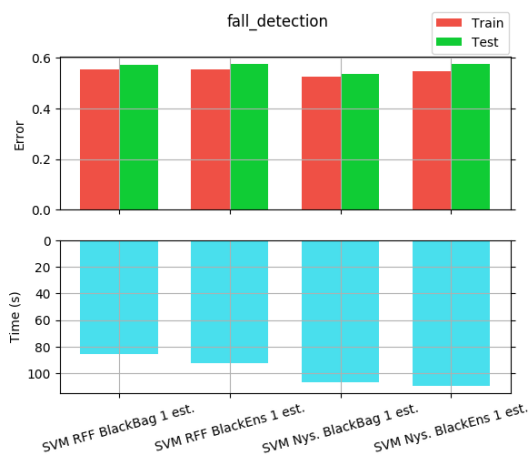
# Results of experiment 3.4



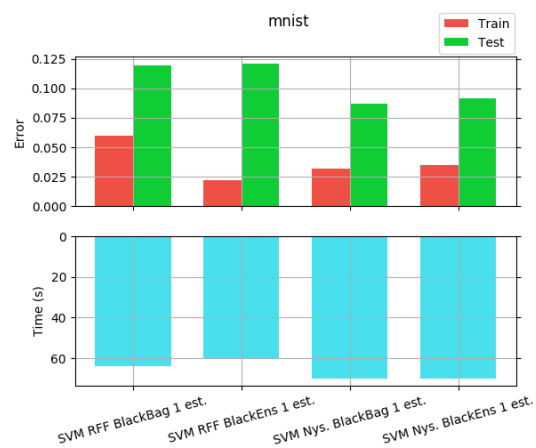
prueba covertypes



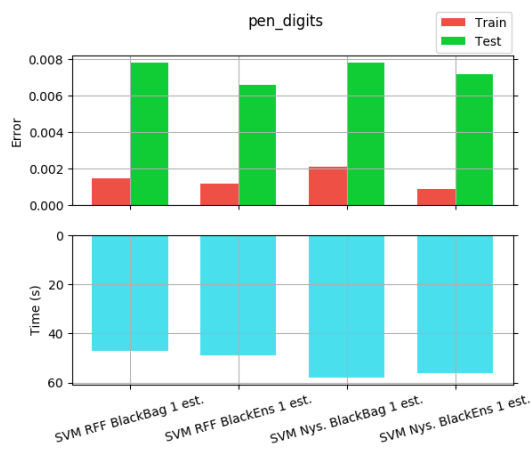
prueba digits



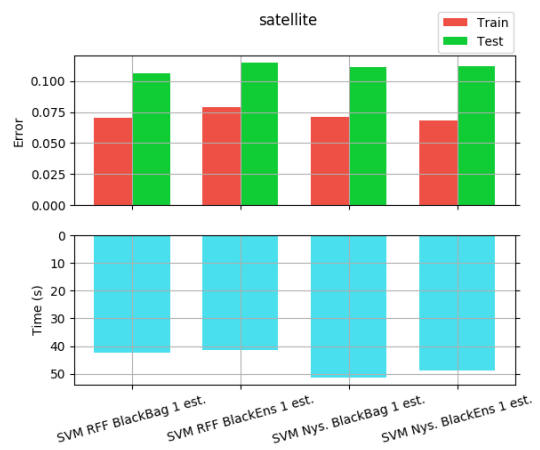
prueba fall-detection



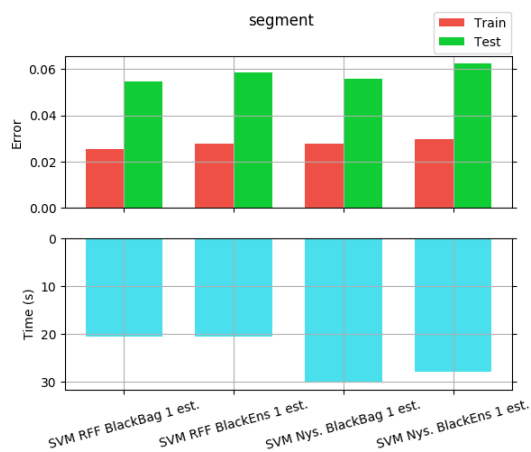
prueba mnist



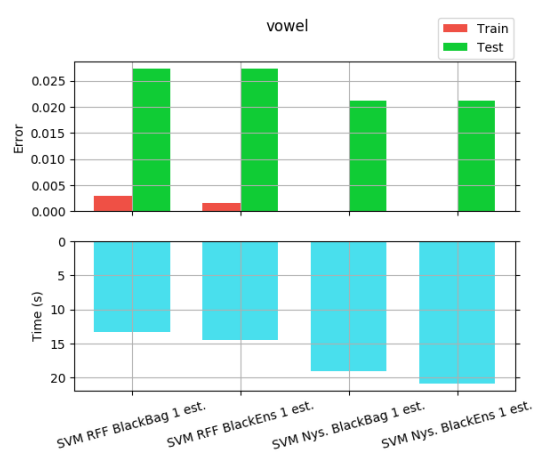
prueba pen-digits



prueba satellite



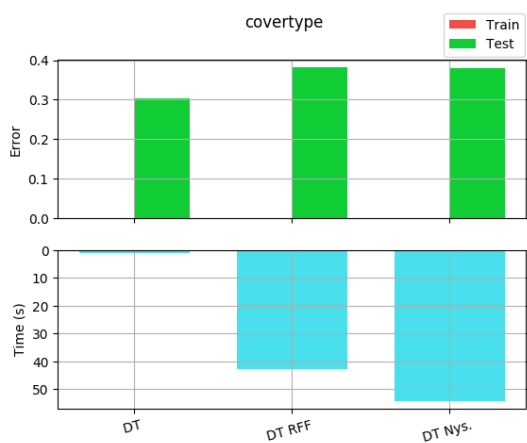
prueba segment



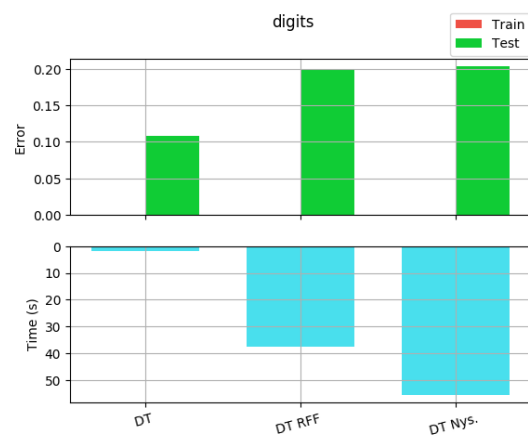
prueba vowel

## Appendix J

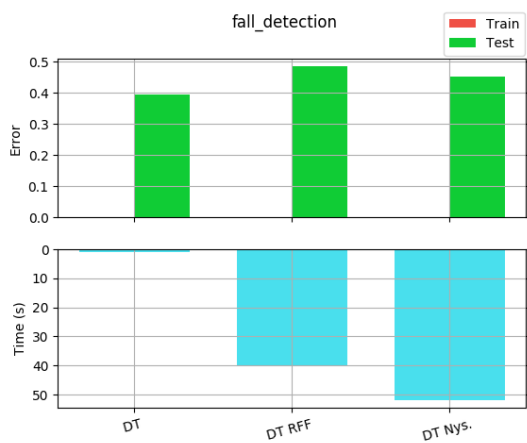
# Results of experiment 4.1



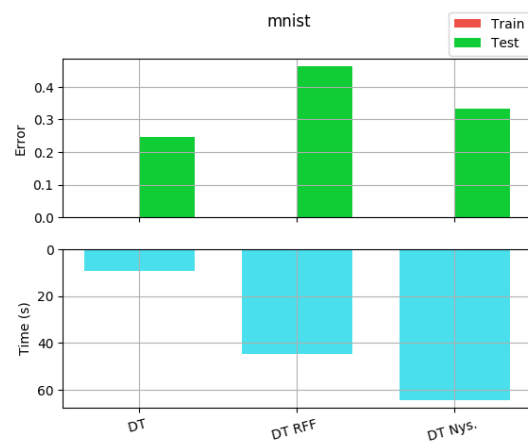
prueba covertypes



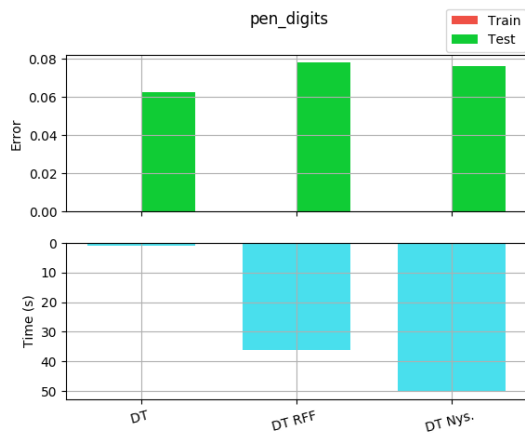
prueba digits



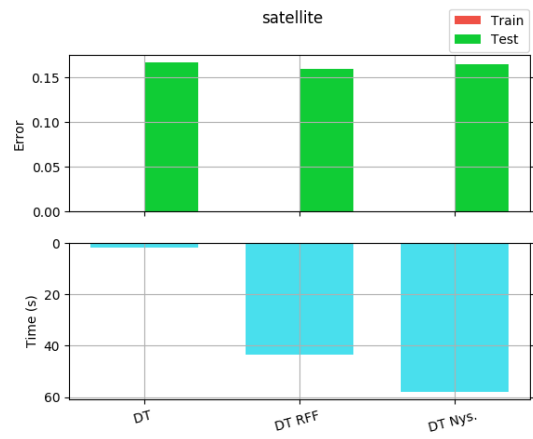
prueba fall-detection



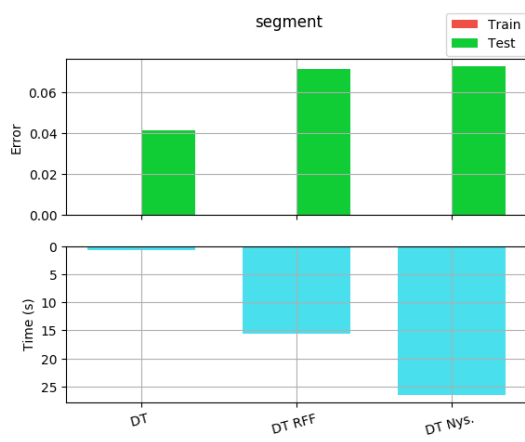
prueba mnist



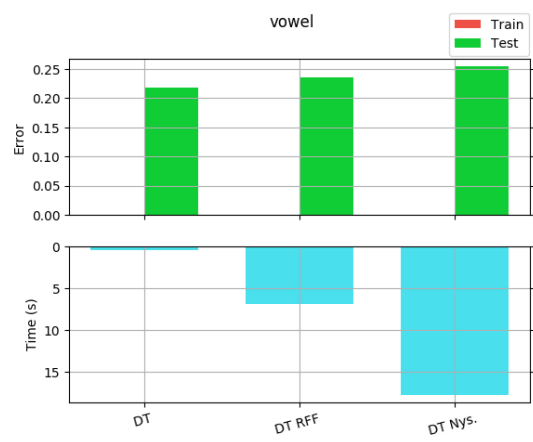
prueba pen-digits



prueba satellite



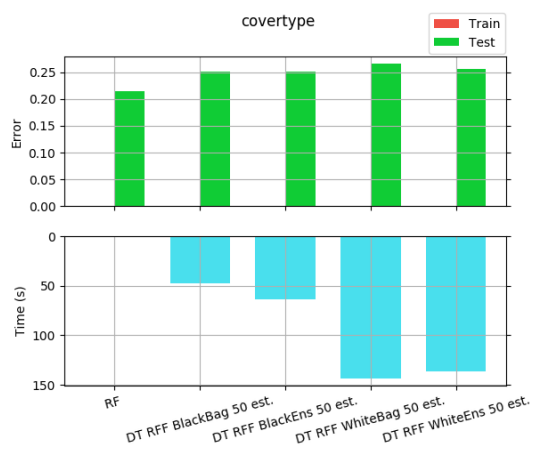
prueba segment



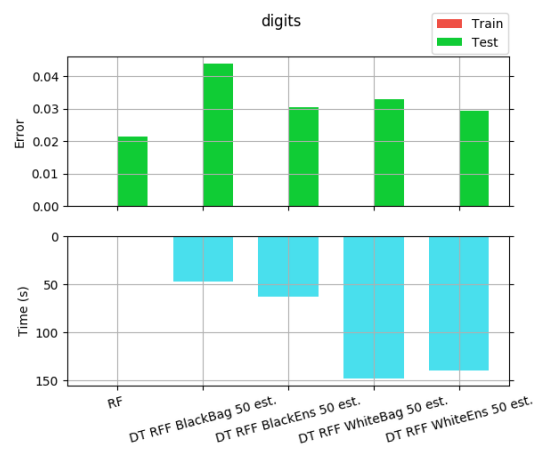
prueba vowel

## Appendix K

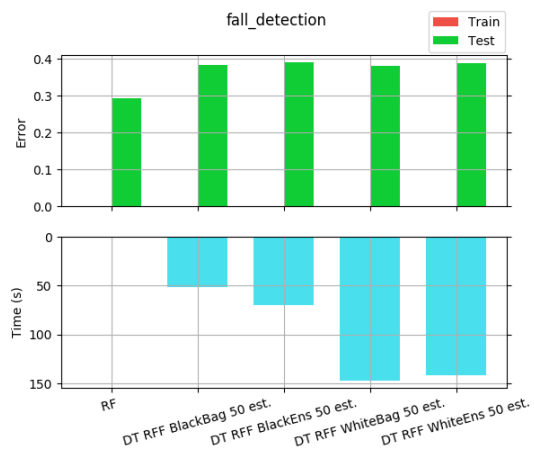
# Results of experiment 4.2



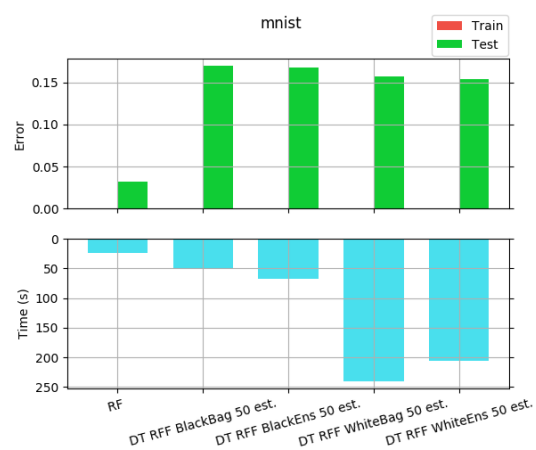
prueba covertype



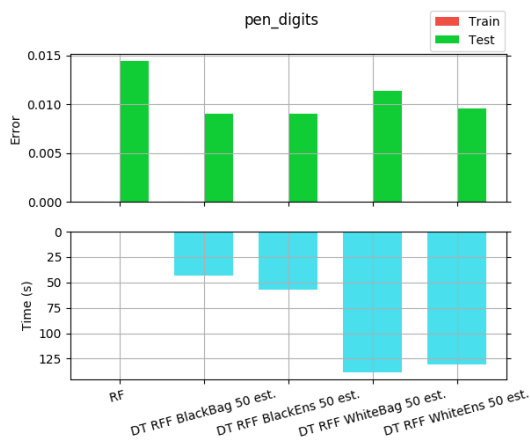
prueba digits



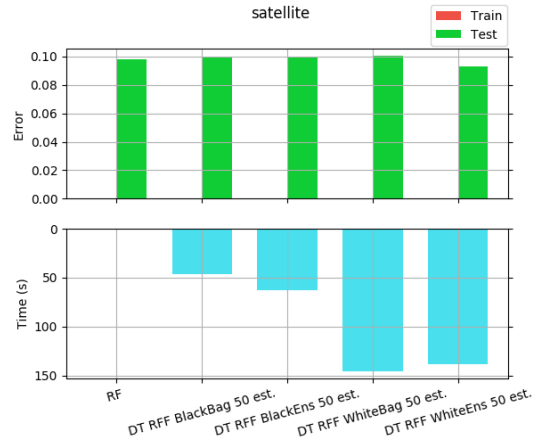
prueba fall-detection



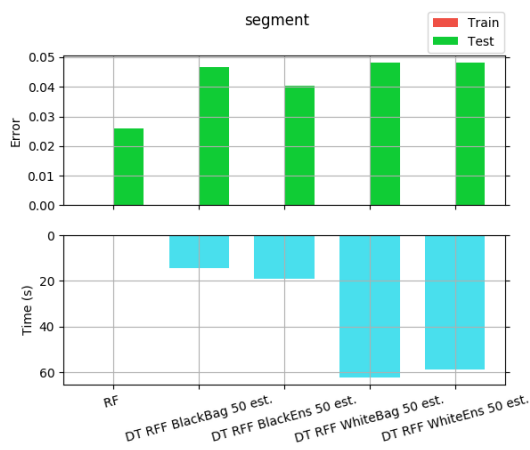
prueba mnist



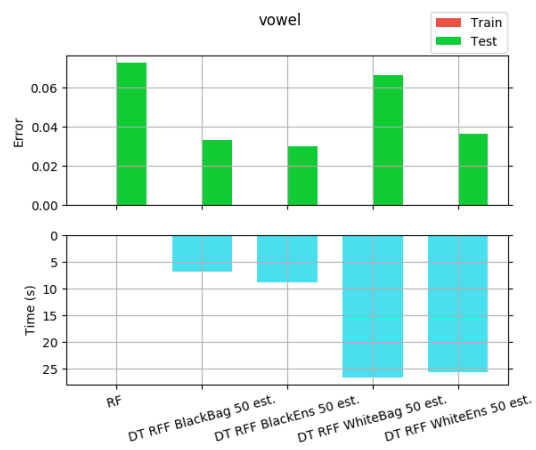
prueba pen-digits



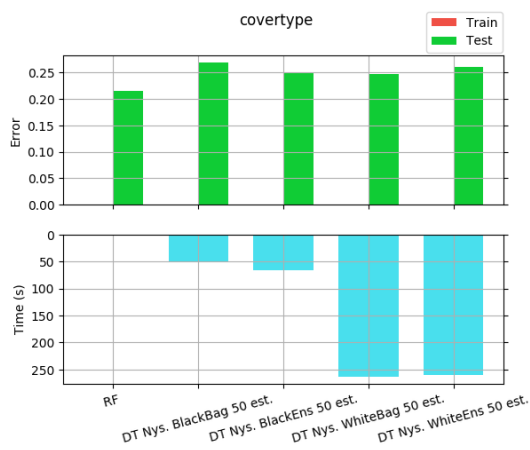
prueba satellite



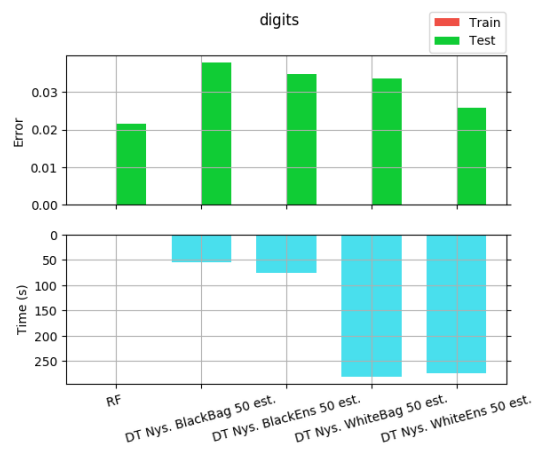
prueba segment



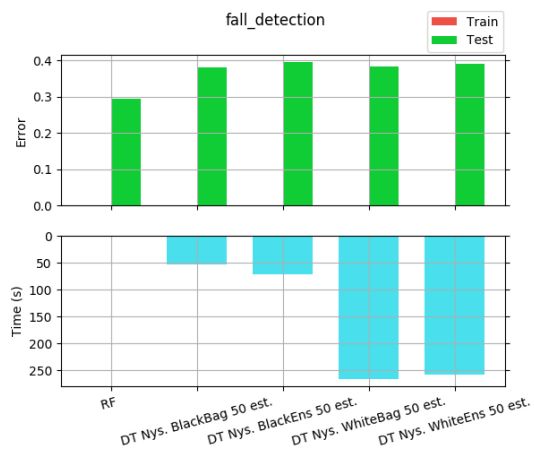
prueba vowel



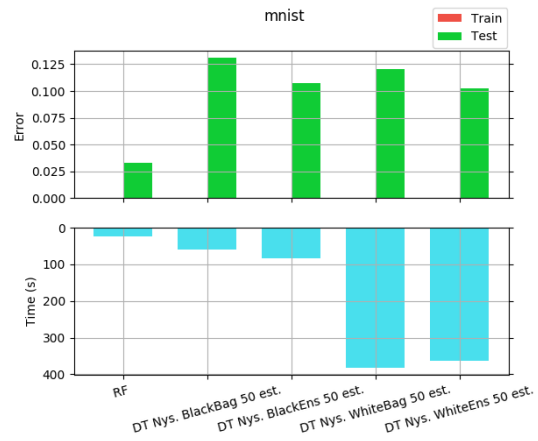
prueba coverytype



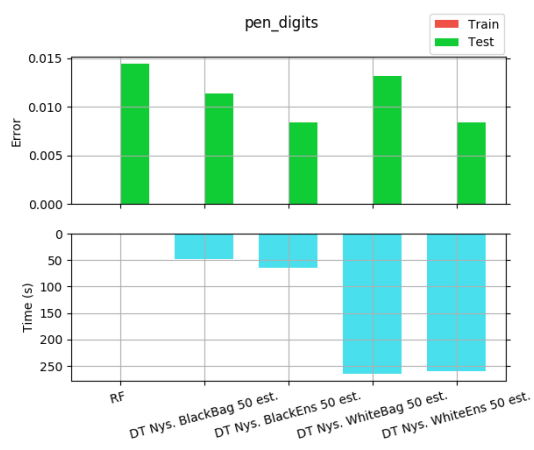
prueba digits



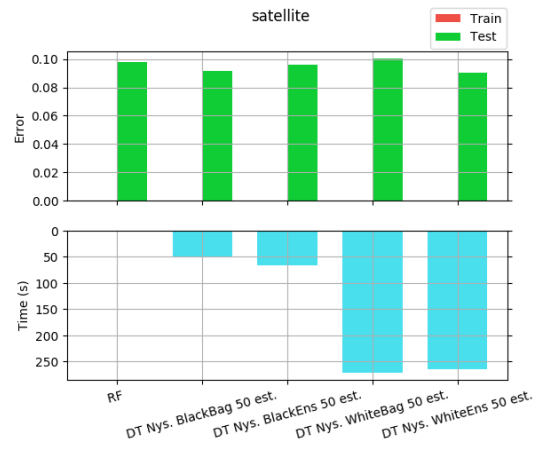
prueba fall-detection



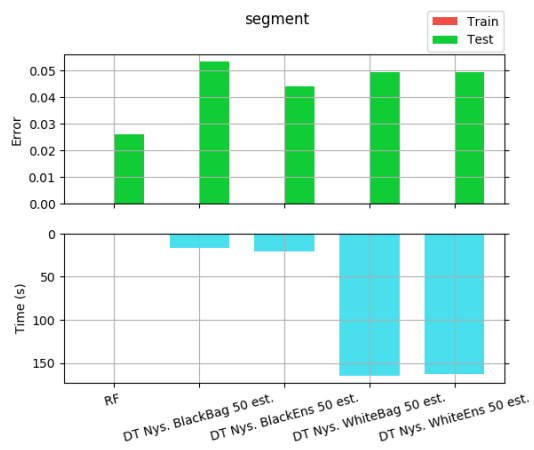
prueba mnist



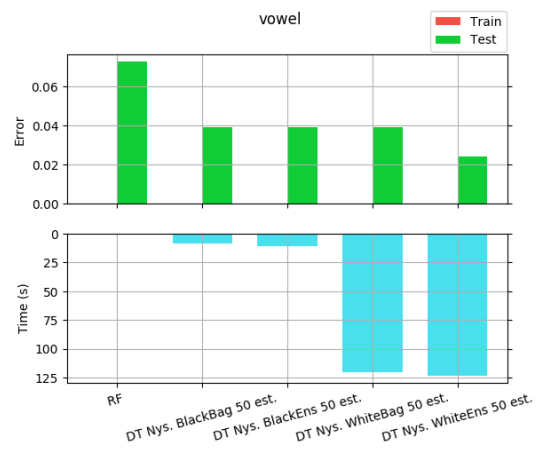
prueba pen-digits



prueba satellite



prueba segment



prueba vowel





# References

- [4] Barbara Caputo et al. "Appearance-based object recognition using SVMs: which kernel should I use?" In: *Proc of NIPS workshop on Statistical methods for computational experiments in visual processing and computer vision*, Whistler. Vol. 2002. 2002.
- [7] J Mercer. *Functions of positive and negative type and their connection with the theory of integral equations*, *Philosophical Transsaction of the Royal Society of London, Ser.* 1909.
- [8] Ali Rahimi and Benjamin Recht. "Random features for large-scale kernel machines". In: *Advances in neural information processing systems*. 2008, pp. 1177–1184.
- [9] Yauheni Selivonchyk. "Speeding up Support Vector Machines with Random Fourier Features". In: (2016).
- [13] Shuai Zhang et al. "Stacked Kernel Network". In: *CoRR* abs/1711.09219 (2017).



# Datasets

- [1] E. Alpaydin and Fevzi. Alimoglu. *Pen-Based Recognition of Handwritten Digits*. July 1998. URL: <https://archive.ics.uci.edu/ml/datasets/Pen-Based+Recognition+of+Handwritten+Digits>.
- [2] E. Alpaydin and C. Kaynak. *Optical Recognition of Handwritten Digits*. July 1998. URL: <http://archive.ics.uci.edu/ml/datasets/optical+recognition+of+handwritten+digits>.
- [3] Jock A. Blackard, Dr. Denis J. Dean, and Dr. Charles W. Anderson. *Forest Cover-type data*. Aug. 1998. URL: <https://archive.ics.uci.edu/ml/datasets/covertime>.
- [5] David Deterding, Mahesan Niranjan, and Tony Robinson. *Vowel Recognition (Deterding data)*. URL: [https://archive.ics.uci.edu/ml/datasets/Connectionist+Bench+\(Vowel+Recognition+-+Deterding+Data\)](https://archive.ics.uci.edu/ml/datasets/Connectionist+Bench+(Vowel+Recognition+-+Deterding+Data)).
- [6] Yann LeCun, Corinna Cortes, and Christopher J.C. Burges. *MNIST database*. URL: <http://yann.lecun.com/exdb/mnist/>.
- [10] Ashwin Srinivasan. *Statlog (Landsat Satellite)*. Feb. 1993. URL: [https://archive.ics.uci.edu/ml/datasets/Statlog+\(Landsat+Satellite\)](https://archive.ics.uci.edu/ml/datasets/Statlog+(Landsat+Satellite)).
- [11] University of Massachusetts Vision Group. *Image Segmentation data*. Nov. 1990. URL: <http://archive.ics.uci.edu/ml/datasets/image+segmentation>.
- [12] Han Xiao, Kashif Rasul, and Roland Vollgraf. *Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms*. Aug. 28, 2017. arXiv: [cs.LG/1708.07747](https://arxiv.org/abs/1708.07747) [cs.LG].
- [14] Özdemir, Ahmet Turan, and Billur Barshan. *Detecting Falls with Wearable Sensors Using Machine Learning Techniques*. 2017. URL: <https://www.kaggle.com/pitasr/falldata>.