

# Startup Succes Prediction Model II

---

Aprendizaje Automático II

Máster en Informática Industrial y Robótica

---

*Víctor Caínzos López & David Adrián Rodríguez García*

## Introducción

El proyecto consiste en una segunda versión del código desarrollado en la primera parte del trabajo:

- [Startup Succes Prediction Model](#)

Utilizando el mismo conjunto de datos, pero ampliando el contenido con nuevas técnicas de reducción de la dimensionalidad de los datos y aprendizaje no supervisado: clustering, detección de anomalías.

En este documento se incluye la descripción de las ampliaciones realizadas y resultados obtenidos

No obstante, se recomienda consultar la [documentación](#).

# Índice

---

1. [Técnicas de aprendizaje](#)
2. [Ejecución de tareas en paralelo](#)
3. [Medición de tiempos](#)
4. [Resultados obtenidos](#)
5. [Conclusiones](#)

# 1. Técnicas de aprendizaje

Se han implementado técnicas de reducción de la dimensionalidad de los datos (PCA), agrupamiento (clustering) y detección de anomalías.

- [Código](#)
- [Documentación](#)

## 1.1. Reducción de la dimensionalidad

Se define una nueva clase en la que se incluyen los métodos apropiados para cada técnica.

```
class Dimension:
    """Class with implemented techniques for reducing dimensions of features"""
```

### Matriz de correlación

Para estudiar la correlación entre las variables del conjunto de datos se incluye el siguiente método que calcula la matriz de correlación.

```
@staticmethod
def get_correlation_matrix(X: pd.DataFrame):
    """Calculates the correlation matrix of data features

    :param X: Dataset of features and samples
    :type X: pd.DataFrame
    :return: Dataframe with correlation between variables
    :rtype: pd.DataFrame
    """
    samples, nvar = X.shape
    corr_mat = np.corrcoef(np.c_[X].T)
    etiquetas = X.columns.values.tolist()
    # DataFrame with correlations between variables
    return pd.DataFrame(corr_mat, columns=etiquetas, index=etiquetas)
```

### PCA

También se desarrolla un estudio de los componentes principales del dataset empleado, analizando aquellos que expliquen la mayor varianza de los mismos, para conocer qué variables son de mayor relevancia de cara a reducir la dimensionalidad del problema.

```
@staticmethod
def get_PCA(X: pd.DataFrame, **kwargs):
    """Applies PCA technique to reduce dimensionality

    :param X: Dataset of features and samples
    :type X: pd.DataFrame
```

```

: return: PCA results
: rtype: obj
"""
samples, nvar = X.shape
pca = PCA(**kwargs)
pca.fit(X)
return pca

```

## 1.2. Clustering

Análogamente, se diseña una clase que tenga embebida las técnicas de agrupamiento de interés cuando se afronta el proyecto en base a este problema.

```

class Clustering:
    """Class with clustering techniques implemented for unsupervised learning"""

```

### K-means

Dentro de esta clase se incluyen una serie de métodos que utilizan como herramientas técnicas de aprendizaje no supervisado que son aplicables al conjunto de datos cuando no se tiene en cuenta las etiquetas. A continuación se presenta el principal método de agrupamiento por particiones, el K-means. El código fuente del método devuelve las inercias o distancias al cuadrado de cada punto con su centroide más cercano.

```

@staticmethod
def perform_kmeans(X: pd.DataFrame, **kwargs):
    """Static method wrapped in Clustering class which performs Kmeans
    algorithm

    :param X: Dataset with samples
    :type X: pd.DataFrame
    :return: Inertias from data samples to clusters
    :rtype: dict
    """
    k_list = np.arange(1, len(X.columns) + 1)
    inertia_dic = {}
    for k in k_list:
        kmeans = KMeans(n_clusters=k, random_state=0).fit(X)
        inertia_dic[k] = kmeans.inertia_
    return inertia_dic

```

### DBSCAN

Además del K-means, también se define un método referido al agrupamiento por densidad de los datos que utiliza la técnica DBSCAN. En el siguiente bloque de código se calculan las distancias de los k-vecinos más cercanos, siendo k el doble del número de variables, y haciendo uso de este algoritmo para estudiar el valor óptimo del hiperparámetro epsilon que define la densidad de cada punto.

```

@staticmethod
def perform_DBSCAN(X: pd.DataFrame, **kwargs):
    """Static method wrapped in Clustering class which performs

    :param X: Dataset with samples
    :type X: pd.DataFrame
    :return: Distances for all data to its k-nearest-neighbor
    :rtype: list
    """
    neighbors = NearestNeighbors(n_neighbors=2 * len(X.columns))
    neighbors_fit = neighbors.fit(X)
    distances, _ = neighbors_fit.kneighbors(X)
    distances = np.sort(distances, axis=0)
    distances = distances[:, -1]
    return distances

```

### 1.3 Detección de anomalías

Para finalizar con el apartado de nuevas técnicas de aprendizaje, se crea una nueva clase específica para los métodos dedicados a la detección de anomalías.

```

class Anomalies:
    """Class with clustering techniques implemented for unsupervised learning."""

```

En este caso, dado que el problema principal no fue concebido para afrontarse según este criterio, es necesario manipular el conjunto de datos considerando una clase mayoritaria como normal y otra minoritaria como anómala. Esta labor la lleva a cabo el constructor de la clase, adaptando y creando un escenario en la detección de anomalías donde Los datos de entrenamiento solo contienen datos de tipo normal, mientras que los datos de test pueden contener anomalías.

### Isolation Forest

El primer método que se presenta ejecuta el procedimiento conocido como Isolation Forest, utilizando el algoritmo empleado en los bosques aleatorios para aislar los datos con distintas características del resto. El código devuelve un informe con los resultados de la clasificación.

```

def perform_IsolationForest(
    self, n_estimators=100, target_names=None, contamination=0, **kwargs
):
    """Isolation Forest algorithm.

    :param n_estimators: number of estimators used, defaults to 100
    :type n_estimators: int, optional
    :param contamination: percentage os samples included as anomalies in
    trining set, defaults to 0
    :type contamination: int, optional
    :return: classification report with results

```

```

        :rtype: str
        """
        model = IsolationForest(
            n_estimators=n_estimators, contamination=contamination, **kwargs
        )
        model.fit(self.X_train)
        y_test = model.predict(self.X_test)
        return classification_report(
            self.t_test, y_test, target_names=target_names
        )

```

## LOF

También se define un método basado que utiliza el algoritmo Local Outlier Factor (LOF) y la densidad de los datos como factor para detectar valores atípicos, comparándola con la de sus vecinos. De igual manera, el método definido devuelve un informe de clasificación.

```

def perform_LOF(self, n_neighbors=10, target_names=None, novelty=True):
    """LOF algorithm.

    :param n_neighbors: number of data neighbours used, defaults to 10
    :type n_neighbors: int, optional
    :param novelty: param necessary to detect anomalies, defaults to True
    :type novelty: bool, optional
    :return: classification report with results
    :rtype: str
    """
    model = LocalOutlierFactor(n_neighbors=n_neighbors, novelty=novelty)
    model.fit(self.X_train)
    y_test = model.predict(self.X_test)
    return classification_report(
        self.t_test, y_test, target_names=target_names
    )

```

## Red de neuronas autoencoder

Terminando con los métodos definidos en la clase específica para la detección de anomalías, se incluyen una serie de funciones para la creación y utilización de una red neuronas autoencoder, que comprime los datos a una dimensión inferior, y luego a la salida del modelo, intenta reconstruir los datos de entrada, extrayendo la información esencial.

```

def perform_autoencoding(self):
    """Create an autoencoder neural network."""
    _, variables = self.X_train.shape
    autoencoder = models.Sequential()
    autoencoder.add(layers.Dense(1, input_dim=variables, activation="relu"))
    autoencoder.add(layers.Dense(variables, activation="relu"))
    opt = keras.optimizers.Adam(learning_rate=0.001)

```

```
autoencoder.compile(loss="mean_squared_error", optimizer=opt)
self.autoencoder = autoencoder
```

Durante el entrenamiento, la capa oculta aprenderá la representación de los datos de entrada normales y en la fase de test, una anomalía será diferente a un dato normal por lo que la red tendrá problemas para reconstruir el dato anómalo, provocando que el error a la salida sea alto.

```
def train_autoencoding(self, epochs=200, batch_size=100, **kwargs):
    """Fits autoencoder neural network.

    :param epochs: number of times all data is passed to network, defaults to
200
    :type epochs: int, optional
    :param batch_size: number of splits of data in batches, defaults to 100
    :type batch_size: int, optional
    """
    self.history = self.autoencoder.fit(
        self.X_train,
        self.X_train,
        validation_data=(self.X_test, self.X_test),
        **kwargs
    )
```

Se define una función encargada de graficar la curva de validación con los resultados de entrenamiento y test, utilizada posteriormente a la hora de generar el report donde se recogen todos los resultados de los experimentos realizados.

```
def plot_autoencoder_validation(
    self,
    xlabel: str = "Mean Square Error (MSE)",
    ylabel: str = "Iteration (epoch)",
    legend: tuple = ("Entrenamiento", "Test"),
    figsize: tuple = (12, 4),
):
    """Plot autoencoder validation curve.

    :param xlabel: plot xlabel, defaults to "Mean Square Error (MSE)"
    :type xlabel: str, optional
    :param ylabel: plot ylabel, defaults to "Iteration (epoch)"
    :type ylabel: str, optional
    :param legend: plot legend, defaults to ("Entrenamiento", "Test")
    :type legend: tuple, optional
    :param figsize: size of plot, defaults to (12, 4)
    :type figsize: tuple, optional
    :return: plot figure
    :rtype: obj
    """
    fig, ax = plt.subplots(figsize=figsize)
    ax.plot(self.history.history["loss"])
```

```

ax.plot(self.history.history["val_loss"])
ax.set_ylabel(xlabel)
ax.set_xlabel(ylabel)
ax.legend(legend, loc="upper right")
return fig

```

Para etiquetar un nuevo dato como anómalo se emplea un umbral para el error de reconstrucción, determinado en base a los errores de entrenamiento.

```

def plot_autoencoder_threshold(
    self,
    xlabel: str = "Reconstruction error (training)",
    ylabel: str = "Number of data",
    legend: tuple = ("Threshold"),
    figsize: tuple = (12, 4),
):
    """Predicts anomalies by reconstructing samples passed to encoder.

    :param xlabel: xlabel plot, defaults to "Reconstruction error (training)"
    :type xlabel: str, optional
    :param ylabel: ylabel plot, defaults to "Number of data"
    :type ylabel: str, optional
    :param legend: legend plot, defaults to ("Threshold")
    :type legend: tuple, optional
    :param figsize: size of plot, defaults to (12, 4)
    :type figsize: tuple, optional
    :return: figure
    :rtype: obj
    """
    y_train = self.autoencoder.predict(self.X_train)
    self.mse_train = np.mean(np.power(self.X_train - y_train, 2), axis=1)
    threshold = np.max(self.mse_train)

    fig, ax = plt.subplots(figsize=figsize)
    ax.hist(self.mse_train, bins=50)
    ax.set_xlabel(xlabel)
    ax.set_ylabel(ylabel)
    ax.axvline(threshold, color="r", linestyle="--")
    ax.legend(legend, loc="upper center")
    return fig

```

Adicionalmente, se define una función gráfica para la detección en función del umbral seleccionado.

```

def plot_autoencoder_error(
    self,
    xlabel: str = "Data Index",
    ylabel: str = "Reconstruction Error",
    legend: tuple = ("Training", "Test", "Threshold"),
    figsize: tuple = (12, 4),

```



```

):
    """Plot autoencoder error.

    :param xlabel: xlabel plot, defaults to "Data Index"
    :type xlabel: str, optional
    :param ylabel: ylabel plot, defaults to "Reconstruction Error"
    :type ylabel: str, optional
    :param legend: legend plot, defaults to ("Training", "Test", "Threshold")
    :type legend: tuple, optional
    :param figsize: size of figure, defaults to (12, 4)
    :type figsize: tuple, optional
    :return: figure
    :rtype: obj
    """
    e_test = self.autoencoder.predict(self.X_test)
    self.mse_test = np.mean(np.power(self.X_test - e_test, 2), axis=1)
    threshold = np.max(self.mse_train)
    fig, ax = plt.subplots(figsize=(12, 4))
    ax.plot(range(1, self.X_train.shape[0] + 1), self.mse_train, "b.")
    ax.plot(
        range(
            self.X_train.shape[0] + 1,
            self.X_train.shape[0] + self.X_test.shape[0] + 1,
        ),
        self.mse_test,
        "r.",
    )
    ax.axhline(threshold, color="r", linestyle="--")
    ax.set_xlabel(xlabel)
    ax.set_ylabel(ylabel)
    ax.legend(legend, loc="upper left")
    return fig

```

Y para generar un informe con los resultados de clasificación.

```

def get_autoencoder_clreport(self, target_names=None):
    """Generate classification report with autoencoders results.

    :return: classification report with results
    :rtype: str
    """
    y_test = np.ones((self.t_test.shape))
    threshold = np.max(self.mse_train)
    y_test[self.mse_test > threshold] = -1

    return classification_report(
        self.t_test, y_test, target_names=target_names
    )

```

## 2. Ejecución de tareas en paralelo

En este apartado se explican las partes del programa en las que se ha implementado la posibilidad de realizar las declaraciones en paralelo. Concretamente utilizando multithreading o programación en hilos de ejecución.

- [Código](#)
- [Documentación](#)

Este procedimiento se implementa en aquellas zonas en las que se llevan a cabo tareas recurrentes como el entrenamiento de los modelos o la red neuronal utilizando validación cruzada. De esta manera se adapta y amplía el programa para aceptar la ejecución de esta parte del código fuente de forma secuencial o en varios hilos, según especifique el usuario por teclado junto con el número de intentos por modelo y particiones o k-folds en la validación cruzada. Es importante anotar en este caso el interés de introducir unos datos sustanciales de forma que se obtengan soluciones representativas para los modelos evaluados.

Para el caso de los modelos de Scikit-Learn, se define una versión de la función original de optimización donde se utilizan tantos hilos como intentos y modelos se tengan, es decir, un hilo por cada versión del modelo a optimizar.

```
def perform_optimizing_models_multithread(
    self,
    tag,
    X: np.ndarray,
    t: np.ndarray,
    train_size: float = 0.85,
    scoring: dict = {
        "accuracy": "accuracy",
        "recall": "recall",
        "specificity": make_scorer(recall_score, pos_label=0),
        "precision": "precision",
        "f1": "f1",
        "roc_auc": "roc_auc",
    },
    cv: int = 10,
    trials: int = 25,
) -> dict:
    """Optimize hyperparameters of the given model and returns the best of
    them using cross validation and multithreading.

    :param tag: model to be evaluated.
    :type tag: str
    :param X: Characteristic matrix numpy array of the dataset which will be
    evaluated
    :type X: np.ndarray
    :param t: Vector labels numpy array of the dataset which will be evaluated
    :type t: np.ndarray
    :param train_size: % of the data to be splitted into train and test
    values, defaults to 0.85
    :type train_size: float, optional
    :param scoring: A dictionary with the wanted metrics to compare and
    evaluate the different models, defaults to { "accuracy": "accuracy", "recall":
```

```

"recall", "specificity": make_scorer(recall_score, pos_label=0), "precision":
"precision", "f1": "f1", "roc_auc": "roc_auc", }
        :type scoring: dict, optional
        :param cv: Number of folds for the cross validation algorithm, defaults to
10
        :type cv: int, optional
        :param trials: Number of trials used to generate random models with
different hyperparameters, defaults to 25
        :type trials: int, optional
        :return: A dictionary with the scores and models trained using cross
validation
        :rtype: dict
        """
        if "accuracy" not in scoring:
            scoring["accuracy"] = "accuracy"

X_train, _, t_train, _ = train_test_split(X, t, train_size=train_size)

print(f"\n***Optimizing {tag} hyperparameters***")
self.threads_dict[tag] = []
for _ in range(trials):
    current_model = HpModels.random_model(tag)
    macro_model = HpModels.build_macro_model(current_model)
    model_thread = mth.ModelThread(
        X_train, np.ravel(t_train), macro_model, cv, scoring
    )
    model_thread.start()
    self.threads_dict[tag].append(model_thread)

last_accuracy = 0
print(f"\n***Cross validation results for {tag}***")
for thread in self.threads_dict[tag]:
    thread.join()
    try:
        mean_accuracy = np.mean(thread.scores["test_accuracy"])
        if mean_accuracy > last_accuracy:
            self.best_models[tag] = (thread.scores, thread.model)
            last_accuracy = mean_accuracy
            print("\nBest accuracy so far: ", last_accuracy)
    except:
        raise Exception(
            f"No solution valid for {tag}. Increase number of trials or
kfolds."
        )
    try:
        print(
            "\nScore:",
            round(np.mean(last_accuracy), 4),
            "-",
            self.best_models[tag][1].steps[1][1],
        )
    except:
        raise Exception(
            f"No solution valid for {tag}. Increase number of trials or

```

```
kfolds."
    )
```

Análogamente para la red de neuronas profunda de Keras, se crea una versión de la optimización específica para multithreading en la que se utilizan tantos hilos como intentos y kfolds se definan, en otras palabras, para cada versión de la red de neuronas, con sus correspondientes hiperparámetros, se ejecutan en paralelo las kfolds de la validación cruzada.

```
def perform_optimize_DNN_multithread(
    self,
    X: np.ndarray,
    t: np.ndarray,
    kfolds: int = 10,
    train_size: float = 0.8,
    trials: int = 5,
    epochs: int = 50,
    batch_size: int = 40,
    metrics: tuple = (
        "accuracy",
        "Recall",
        cm.specificty,
        "Precision",
        cm.f1_score,
        "AUC",
    ),
) -> tuple:
    """Train the current model using cross validation and multithreading and
    register its score comparing with new ones in each trial.

    :param X: Characteristic matrix numpy array of the dataset which will be
    evaluated
    :type X: np.ndarray
    :param t: Vector labels numpy array of the dataset which will be evaluated
    :type t: np.ndarray
    :param kfolds: Number of folds for the cross validation algorithm,
    defaults to 10
    :type kfolds: int, optional
    :param train_size: % of the data to be splitted into train and test
    values, defaults to 0.8
    :type train_size: float, optional
    :param trials: Number of trials used to generate random neural networks
    with different hyperparameters, defaults to 5
    :type trials: int, optional
    :param epochs: Number of maximum iterations allow to converge the
    algorithm, defaults to 50
    :type epochs: int, optional
    :param batch_size: Size of the batch used to calculate lossing function,
    defaults to 40
    :type batch_size: int, optional
    :param metrics: Tuple with the metrics to compare and evaluate the
    differene neural networks, defaults to
```

```

("accuracy", "Recall", cm specificity, "Precision", cm.f1_score, "AUC")
    :type metrics: tuple, optional
    :return: A tuple containing a tuple with the model and the history of the
training, and another tuple with the X_test and t_test arrays to evaluate the
model

    :rtype: tuple
    """
    # Its needed the accuracy metric, if it is not passed it will be auto-
included
    if "accuracy" not in metrics:
        metrics.append("accuracy")

    _, m = X.shape
    n_classes = len(np.unique(t))
    cv = KFold(kfolds)
    last_mean = 0

    # Split the data into train and test sets. Note that test set will be
reserved for evaluate the best model later with data unseen previously for it
    X_train_set, X_test, t_train_set, t_test = train_test_split(
        X, t, train_size=train_size
    )

    # Loop for different trials or models to train in order to find the best
threads_dict = {}
    for row in range(trials):

        model_aux, optimizer, loss = HpDNN.create_random_network(
            m, n_classes, metrics=metrics
        )
        params, comp = HpDNN.get_hyperparams(model_aux)

        print(f"\n***Trial {row+1} hyperparameters***", end="\n\n")
        print(params, "\n", comp)
        # Lists to store threads to manage them
        threads_dict[row] = []
        # Loop that manage cross validation using training set
        for train_index, test_index in cv.split(X_train_set):
            # This sentence carefully clones the untrained model in each fold
in order to avoid unwanted learning weights between them
            model = keras.models.clone_model(model_aux)
            model.compile(optimizer=optimizer, loss=loss, metrics=metrics)

            X_train, X_val = (
                X_train_set[train_index],
                X_train_set[test_index],
            )
            t_train, t_val = (
                t_train_set[train_index],
                t_train_set[test_index],
            )

            t_train = to_categorical(t_train, num_classes=n_classes)
            t_val = to_categorical(t_val, num_classes=n_classes)

```

```

        # Training of the model using multithreading
        thread = mth.DnnThread(
            X_train, t_train, X_val, t_val, model, epochs, batch_size
        )
        thread.start()
        threads_dict[row].append(thread)

# Save the train and test results, of the current model
for row in range(trials):
    # Lists to store historical data and means per fold
    historial = []
    mean_folds = []
    print(f"\n***Trial {row+1} results***", end="\n\n")
    for thread in threads_dict[row]:
        thread.join()
        historial.append(thread.history.history)
        mean = np.mean(thread.history.history["val_accuracy"])
        mean_folds.append(mean)
        print(
            f"\nKFold{threads_dict[row].index(thread)+1} --- Current
score: {mean}"
        )

    # Criterion to choose the best model with the highest accuracy score
    in validation
    means = np.mean(mean_folds)
    history_metrics = HpDNN.split_history_metrics(historial)
    print(
        f"\n\tMin. Train score:
{min(np.concatenate(history_metrics['accuracy']))} | Max. Train score:
{max(np.concatenate(history_metrics['accuracy']))}"
    )
    if means > last_mean:
        # Save the model and its historial results
        self.best_model = (thread.model, history_metrics)
        self.test_set = (X_test, t_test)
        last_mean = means

```

Cabe recordar que las versiones secuenciales de las funciones que contienen el código fuente que definen las instrucciones de los procedimientos de optimización tanto para los modelos Scikit-Learn como para la red neuronal de Keras siguen disponibles tras la ampliación de esta parte del programa conservando la estructura principal de la primera parte del trabajo:

- [Startup Succes Prediction](#)

A excepción de la implementación en clases, que es de utilidad para la medición de los tiempo de ejecución tal y como se describe en el siguiente apartado.

### 3. Medición de tiempos

Se lleva a cabo una ampliación que permite realizar mediciones de tiempos de entrenamiento medios de los modelos de aprendizaje desarrollados para la clasificación supervisada. Para ello se utiliza una función envoltorio que decora aquellas funciones en las que resulta de interés medir el tiempo de ejecución, de forma que una vez decoradas, con solo llamarlas ya se calcularía el tiempo de ejecución.

- [Código](#)
- [Documentación](#)

```
import logging
import time

from functools import wraps

logger = logging.getLogger(__name__)

# Misc logger setup so a debug log statement gets printed on stdout.
logger.setLevel("DEBUG")
handler = logging.StreamHandler()
log_format = "%(asctime)s %(levelname)s -- %(message)s"
formatter = logging.Formatter(log_format)
handler.setFormatter(formatter)
logger.addHandler(handler)
```

Como se venía diciendo al final del apartado anterior, se ha optado por reestructurar los módulos de optimización de los modelos en clases con el fin de hacer uso de los atributos de instancia para referenciar los resultados con los mejores modelos y facilitando así la devolución de las mediciones de los tiempos de ejecución utilizando la función decoradora.

```
def timed(func):
    """This decorator prints the execution time for the decorated function."""

    @wraps(func)
    def wrapper(*args, **kwargs):
        start = time.time()
        result = func(*args, **kwargs)
        end = time.time()
        logger.debug(
            "{} ran in {}s".format(func.__name__, round(end - start, 2))
        )
        return round(end - start, 2)

    return wrapper
```

De la misma forma, para poder medir los tiempos de cada modelo por separado se generaliza la función de optimización original empleada para los modelos de Scikit-Learn, que se encarga de realizar el proceso de

optimización y se define una nueva función dedicada a la llamada y recogida de los tiempos de ejecución de la anterior para cada modelo. Esta modificación no es necesaria para la red de neuronas puesto que es exclusiva el módulo específico.

Se decoran las funciones de optimización tanto secuencial como multihilo utilizando la función envoltorio anterior.

```
@timed
def perform_optimizing_model
```

```
@timed
def perform_optimizing_models_multithread
```

Y se llaman a las funciones de nivel superior que gestionan la medición de los tiempos.

```
def optimizing_models(self, models: list, *args, **kwargs) -> dict:

    time_models = dict()
    for tag in models:
        time_models[tag] = self.perform_optimizing_model(
            tag, *args, **kwargs
        )
    return self.best_models, pd.DataFrame(
        time_models, index=["Time Models"]
    )
```

```
def optimizing_models_multithread(self, models: list, *args, **kwargs):

    time_models = dict()
    for tag in models:
        time_models[tag] = self.perform_optimizing_models_multithread(
            tag, *args, **kwargs
        )

    return self.best_models, pd.DataFrame(
        time_models, index=["Time Models"]
    )
```

Para la red de neuronas el modo de operación sería el mismo. Decorar las funciones de optimización con la función envoltorio de la medición de tiempos:

```
@timed
def perform_optimizing_DNN(
```



```
@timed
def perform_optimize_DNN_multithread
```

Y ejecutar las funciones que llaman a la optimización y guardan el tiempo de ejecución:

```
def optimize_DNN(self, *args, **kwargs):
    timeDNN_seq = self.perform_optimizing_DNN(*args, **kwargs)
    return (self.best_model, self.test_set), timeDNN_seq
```

```
def optimize_DNN_multithread(self, *args, **kwargs):
    timeDNN_thr = self.perform_optimize_DNN_multithread(*args, **kwargs)
    return (self.best_model, self.test_set), timeDNN_thr
```

## 4. Resultados obtenidos

Todos los resultados extraídos de los métodos explicados anteriormente se recogen en un report de manera automática, se presenta el informe completo para la primera parte del trabajo más las aportaciones del actual.

- [Código](#)
- [Documentación](#)

### Preprocesado: Preparación de los datos

En primer lugar se realizará un análisis del dataset **startup\_data.csv** para el cual se realizará la limpieza de los datos espurios o nulos y se procederá al filtrado de las columnas representativas y la recodificación de variables cualitativas a cuantitativas.

**Data Missing dataframe: Contiene los datos eliminados del dataset original.**

	feature	missing	(%) of total
0	closed_at	587	63.67
1	age_first_milestone_year	152	16.49
2	age_last_milestone_year	152	16.49
3	state_code.1	1	0.11

**Data Spurious dataframe: Contiene los datos sin sentido del dataset original.**

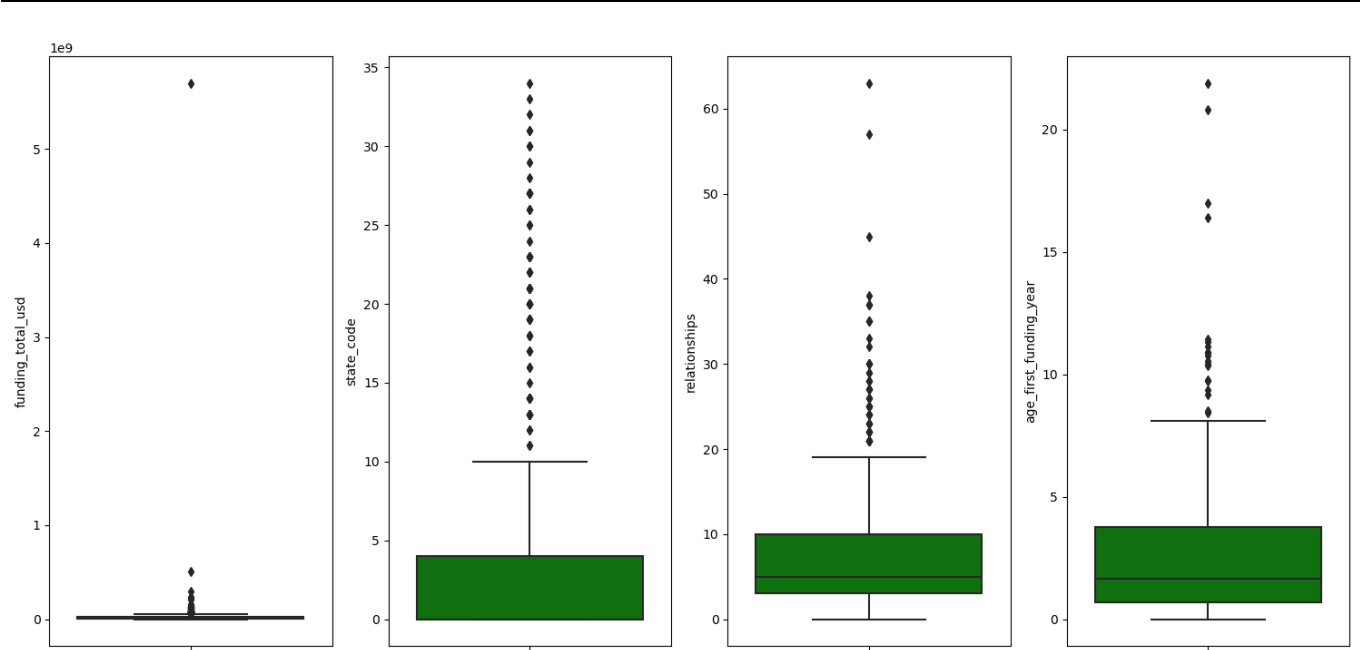
	age_first_funding_year	age_last_funding_year	age_first_milestone_year	age_last_milestone_year	age
88	0.8822	0.8822	0	0	-8
558	-9.0466	-9.0466	-6.0466	-3.8822	-4
73	1.6685	9.337	7.3808	10.474	-2
350	0.3288	0.3288	-0.4192	-0.4192	0
690	0	0.6904	0	0.6904	0

**Data Skewness dataframe: Contiene los datos con alta dispersión del dataset original.**

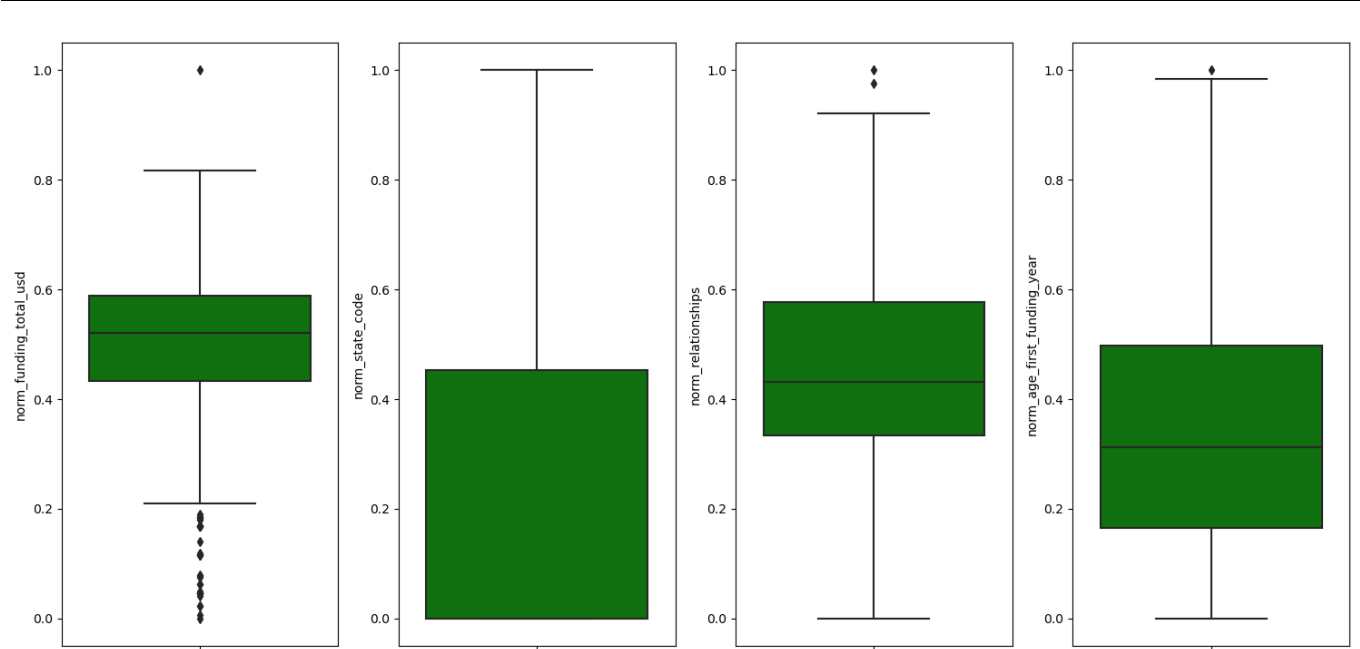
	feature	skewness
0	funding_total_usd	27.7868
1	state_code	2.54287
2	relationships	2.33271
3	age_first_funding_year	2.32396

	feature	skewness
4	avg_participants	1.75802

**Boxplot Feature Skewness > 2:** Muestra la dispersión de los datos para las características con asimetría mayor que 2.



**Boxplot Norm Features:** Muestra la dispersión de los datos para las características normalizadas.



**X dataframe:** Contiene la matriz de características.

	state_code	age_last_funding_year	age_first_milestone_year	age_last_milestone_year	funding_rounds
--	------------	-----------------------	--------------------------	-------------------------	----------------

	state_code	age_last_funding_year	age_first_milestone_year	age_last_milestone_year	funding_rounds
0	0	3.0027	4.6685	6.7041	3
1	0	9.9973	7.0055	7.0055	4
2	0	1.0329	1.4575	2.2055	1
3	0	5.3151	6.0027	6.0027	3
4	0	1.6685	0.0384	0.0384	2
	milestones	is_CA	is_NY	is_MA	is_TX
0	3	1	0	0	0
1	1	1	0	0	0
2	2	1	0	0	0
3	1	1	0	0	0
4	1	1	0	0	0
	is_otherstate	is_software	is_web	is_mobile	is_enterprise
0	0	0	0	0	0
1	0	0	0	0	1
2	0	0	1	0	0
3	0	1	0	0	0
4	0	0	0	0	0
	is_advertising	is_gamesvideo	is_ecommerce	is_biotech	is_consulting
0	0	0	0	0	0
1	0	0	0	0	0
2	0	0	0	0	0
3	0	0	0	0	0
4	0	1	0	0	0
	is_othercategory	has_VC	has_angel	has_roundA	has_roundB
0	1	0	1	0	0
1	0	1	0	0	1
2	0	0	0	1	0
3	0	0	0	0	1
4	0	1	1	0	0

	has_roundC	has_roundD	avg_participants	is_top500	age
0	0	0	1	0	7
1	1	1	4.75	1	14
2	0	0	4	1	5
3	1	1	3.3333	1	12
4	0	0	1	1	2
		norm_funding_total_usd	norm_age_first_funding_year	norm_relationships	
0		0.268198	0.376383	0.333333	
1		0.623283	0.57891	0.553655	
2		0.415358	0.226596	0.430827	
3		0.623093	0.453101	0.430827	
4		0.36268	0	0.26416	

**t dataframe:** Contiene el vector de etiquetas.

	labels
0	1
1	1
2	1
3	1
4	0

Entrenamiento: Comparativa de modelos de aprendizaje automático

Se procederá a comparar los resultados obtenidos de diferentes modelos de aprendizaje automático variando tanto el tipo de modelo como los hiperparámetros de los que depende con el objetivo de obtener el mejor modelo que prediga el éxito o fracaso de las diferentes startups

**Results dataframe:** Muestra los resultados de los mejores modelos obtenidos

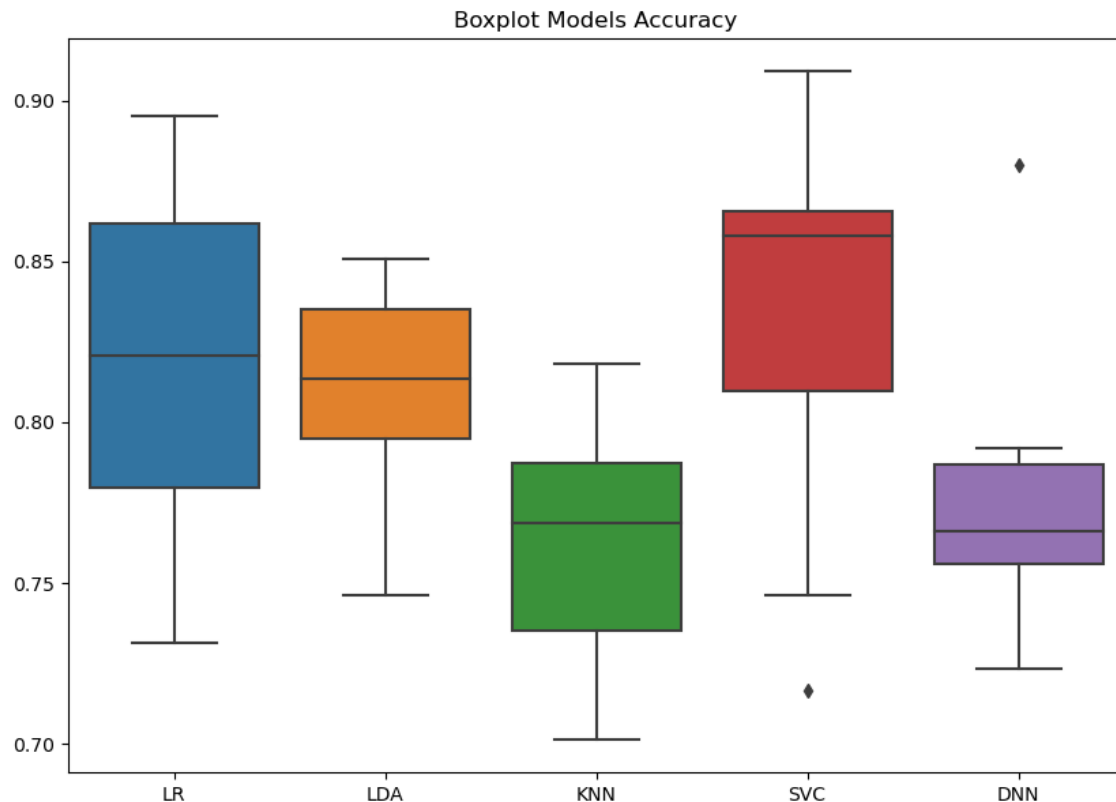
Folds	LR_train_accuracy	LR_val_accuracy	LDA_train_accuracy	LDA_val_accuracy	KNN_train_accuracy
1	0.847176	0.80597	0.818937	0.835821	0.770764
2	0.833887	0.865672	0.830565	0.80597	0.77907
3	0.853821	0.776119	0.820598	0.820896	0.759136
4	0.842193	0.895522	0.82392	0.80597	0.784053

Folds	LR_train_accuracy	LR_val_accuracy	LDA_train_accuracy	LDA_val_accuracy	KNN_train_accuracy
5	0.850498	0.850746	0.825581	0.850746	0.775748
Folds	KNN_val_accuracy	SVC_train_accuracy	SVC_val_accuracy	DNN_train_accuracy	DNN_val_accuracy
1	0.791045	0.863787	0.865672	0.80495	0.769478
2	0.716418	0.863787	0.80597	0.814311	0.757537
3	0.746269	0.853821	0.820896	0.766329	0.76306
4	0.701493	0.872093	0.746269	0.777932	0.77306
5	0.776119	0.873754	0.716418	0.789693	0.791791
Folds	LR_train_recall	LR_val_recall	LDA_train_recall	LDA_val_recall	KNN_train_recall
1	0.914573	0.911111	0.906329	0.909091	0.965
2	0.904523	0.977778	0.911392	0.886364	0.97
3	0.91206	0.933333	0.906329	0.909091	0.9625
4	0.914787	0.931818	0.916456	0.863636	0.965
5	0.927318	0.863636	0.908861	0.931818	0.9625
Folds	KNN_val_recall	SVC_train_recall	SVC_val_recall	DNN_train_recall	DNN_val_recall
1	1	0.918575	0.953488	0.80495	0.769478
2	0.933333	0.916031	0.906977	0.814311	0.757537
3	0.933333	0.918575	0.860465	0.766329	0.76306
4	0.933333	0.928571	0.863636	0.777932	0.77306
5	0.955556	0.928571	0.818182	0.789693	0.791791
Folds	LR_train_specificity	LR_val_specificity	LDA_train_specificity	LDA_val_specificity	KNN_train_specificity
1	0.715686	0.590909	0.652174	0.695652	0.386139
2	0.696078	0.636364	0.676329	0.652174	0.40099
3	0.740196	0.454545	0.657005	0.652174	0.356436
4	0.699507	0.826087	0.647343	0.695652	0.425743
5	0.699507	0.826087	0.666667	0.695652	0.405941
Folds	KNN_val_specificity	SVC_train_specificity	SVC_val_specificity	DNN_train_specificity	DNN_val_specificity
1	0.363636	0.760766	0.708333	0.805547	0.77534
2	0.272727	0.76555	0.625	0.813461	0.75875
3	0.363636	0.732057	0.75	0.765211	0.772583
4	0.227273	0.766667	0.521739	0.777461	0.773525

Folds	KNN_val_specificity	SVC_train_specificity	SVC_val_specificity	DNN_train_specificity	DNN_val_specificity
5	0.409091	0.771429	0.521739	0.79075	0.78819
Folds	LR_train_precision	LR_val_precision	LDA_train_precision	LDA_val_precision	KNN_train_precision
1	0.862559	0.82	0.832558	0.851064	0.756863
2	0.853081	0.846154	0.843091	0.829787	0.762279
3	0.872596	0.777778	0.834499	0.833333	0.747573
4	0.856808	0.911111	0.832184	0.844444	0.768924
5	0.858469	0.904762	0.838785	0.854167	0.762376
Folds	KNN_val_precision	SVC_train_precision	SVC_val_precision	DNN_train_precision	DNN_val_precision
1	0.762712	0.878345	0.854167	0.80495	0.769478
2	0.724138	0.880196	0.8125	0.814311	0.757537
3	0.75	0.865707	0.860465	0.766329	0.76306
4	0.711864	0.881356	0.77551	0.777932	0.77306
5	0.767857	0.883495	0.765957	0.789693	0.791791
Folds					
1					
2					
3					
4					
5					

Boxplot models: Muestra los valores de exactitud de los diferentes modelos

---



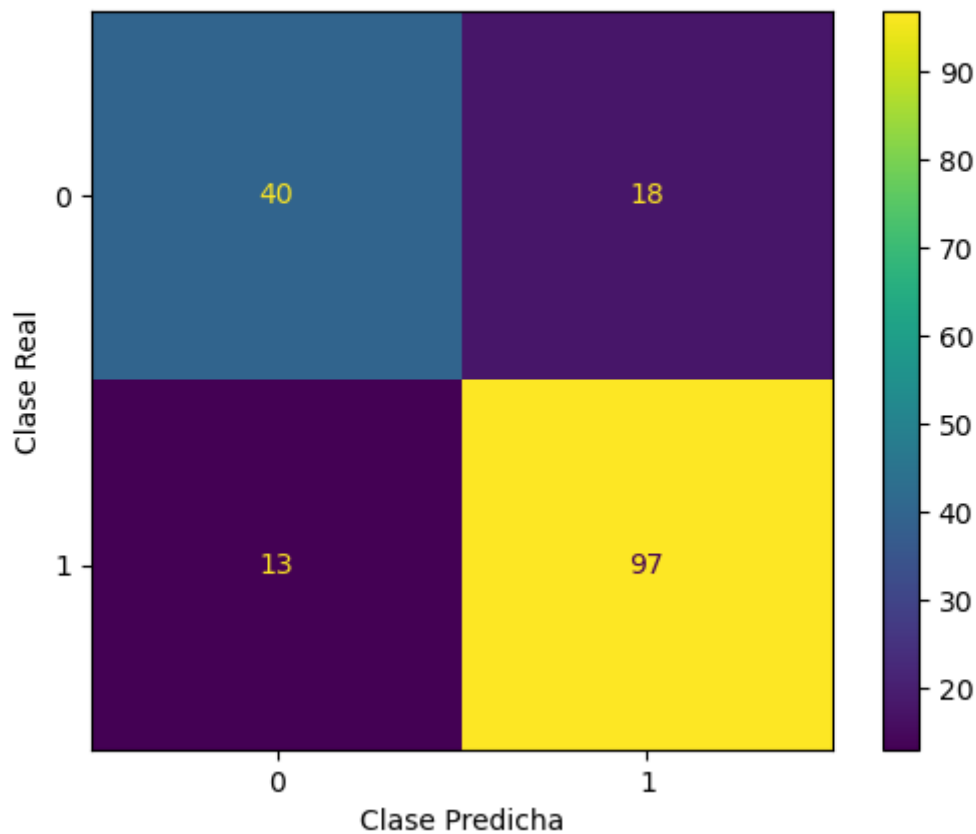
### Contraste de hipótesis: Comparación de modelos mediante el test de Kruskal-Wallis

```
p-valor KrusW:0.01136941395786562
Hypotheses are being rejected: the models are different
Multiple Comparison of Means - Tukey HSD, FWER=0.05
=====
group1  group2  meandiff p-adj  lower  upper  reject
-----
modelDNN modelKNN  -0.0151    0.9 -0.0763 0.0461  False
modelDNN modelLDA   0.0357  0.4718 -0.0255 0.0969  False
modelDNN modelLR    0.0417  0.3131 -0.0195 0.1029  False
modelDNN modelSVC   0.0582  0.0698 -0.003 0.1194  False
modelKNN modelLDA   0.0508  0.1462 -0.0104 0.112   False
modelKNN modelLR    0.0568  0.0804 -0.0044 0.118   False
modelKNN modelSVC   0.0733  0.0117  0.0121 0.1345  True
modelLDA modelLR    0.0061    0.9 -0.0551 0.0673  False
modelLDA modelSVC   0.0225  0.8122 -0.0387 0.0837  False
modelLR  modelSVC   0.0164    0.9 -0.0448 0.0776  False
-----
```

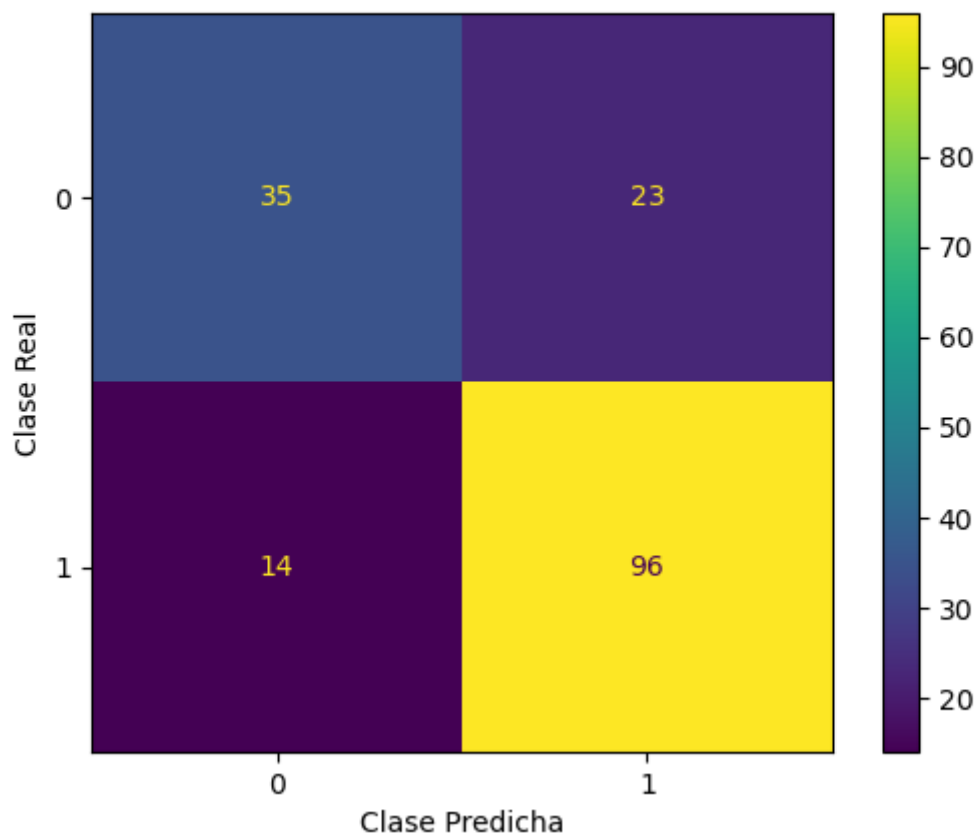
### Matrices de confusión: Compara los valores reales con los valores predichos para cada modelo



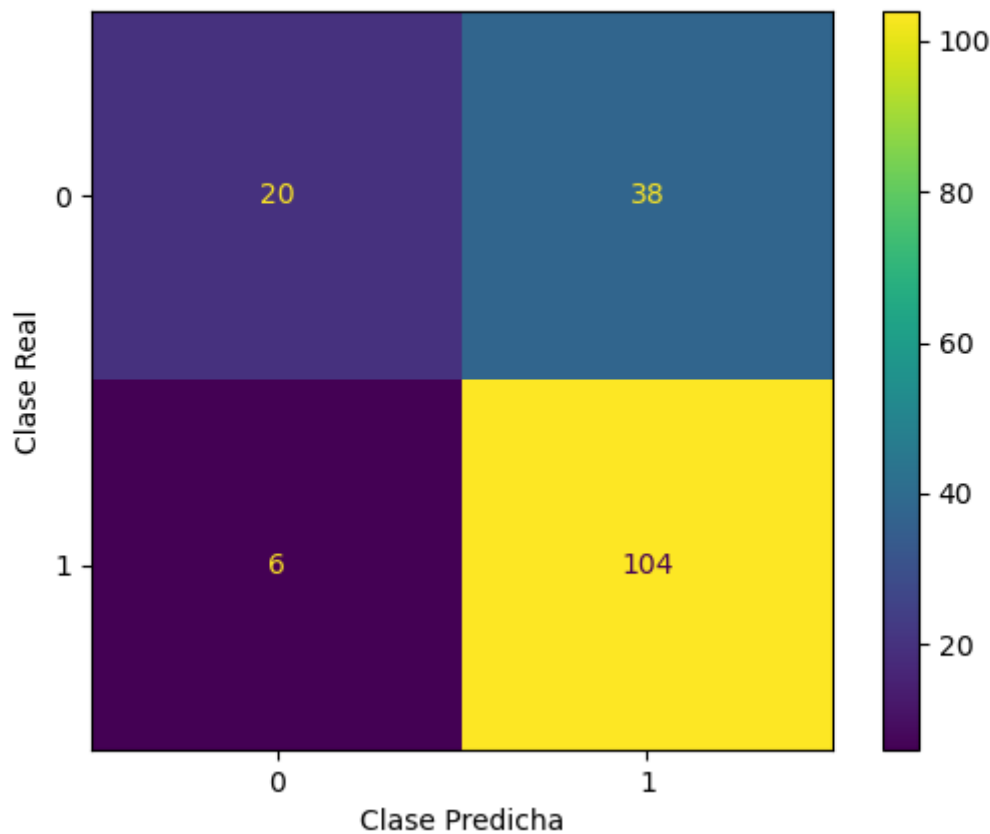
Matriz de confusión LR



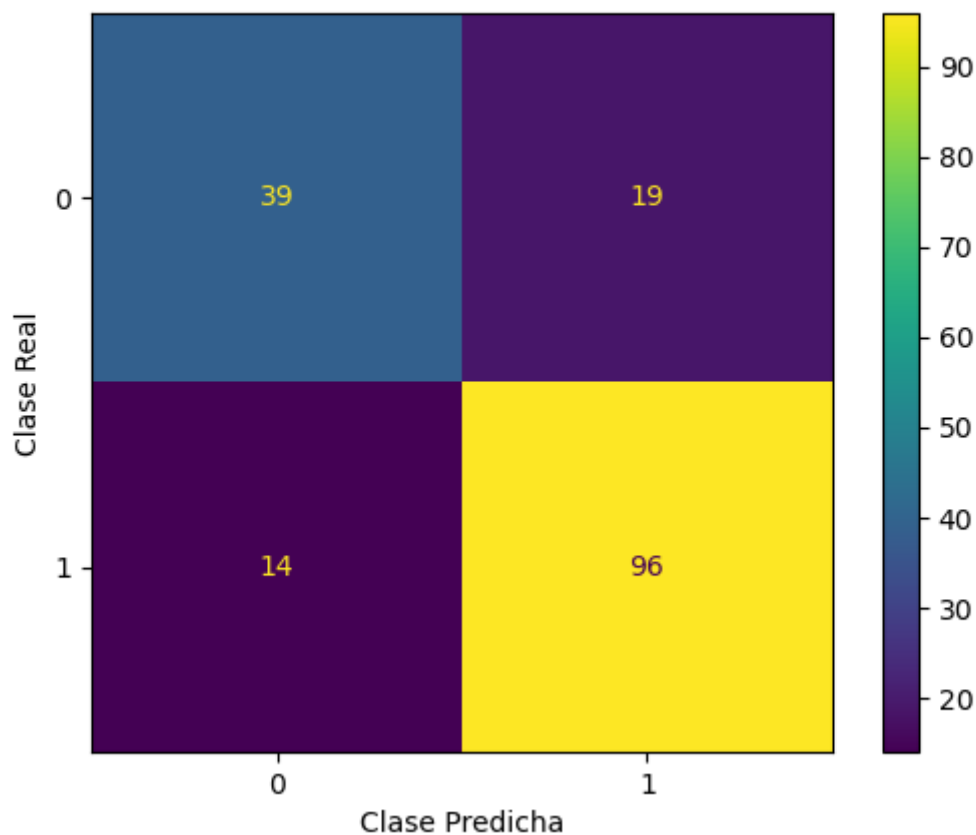
Matriz de confusión LDA



Matriz de confusión KNN



Matriz de confusión SVC

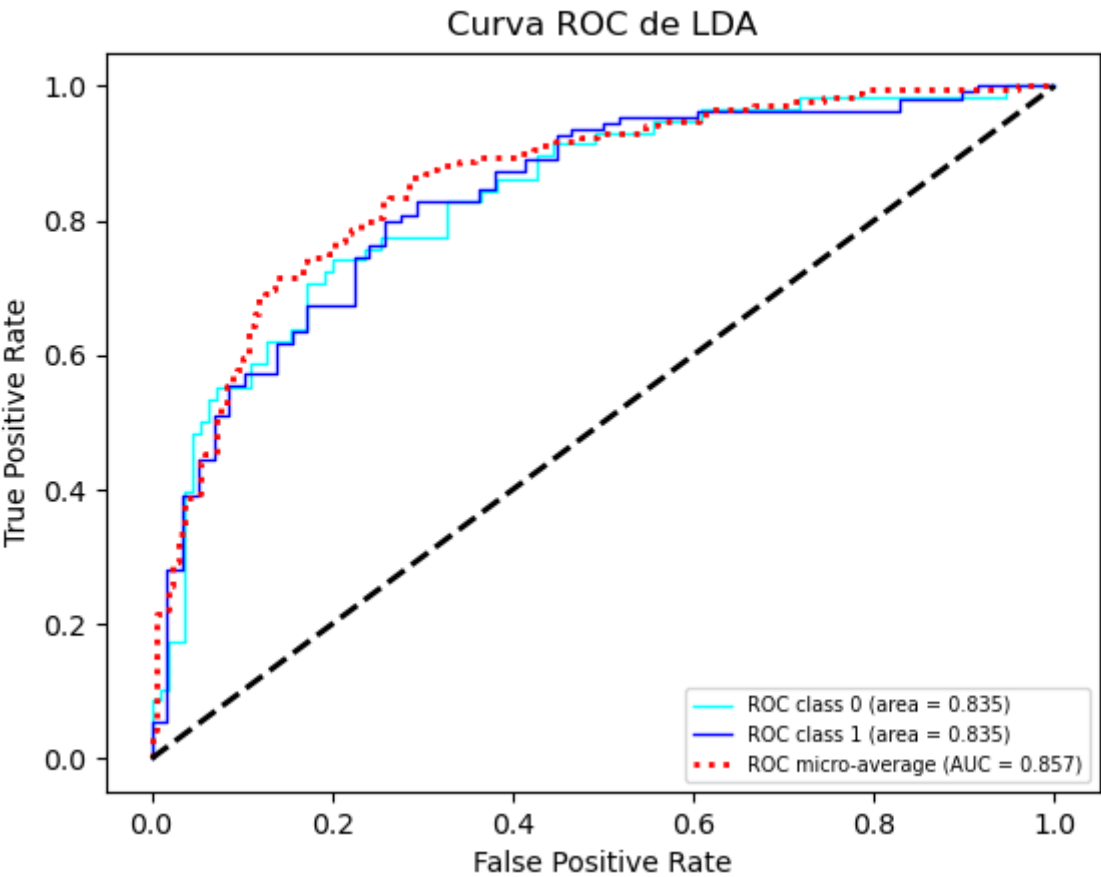
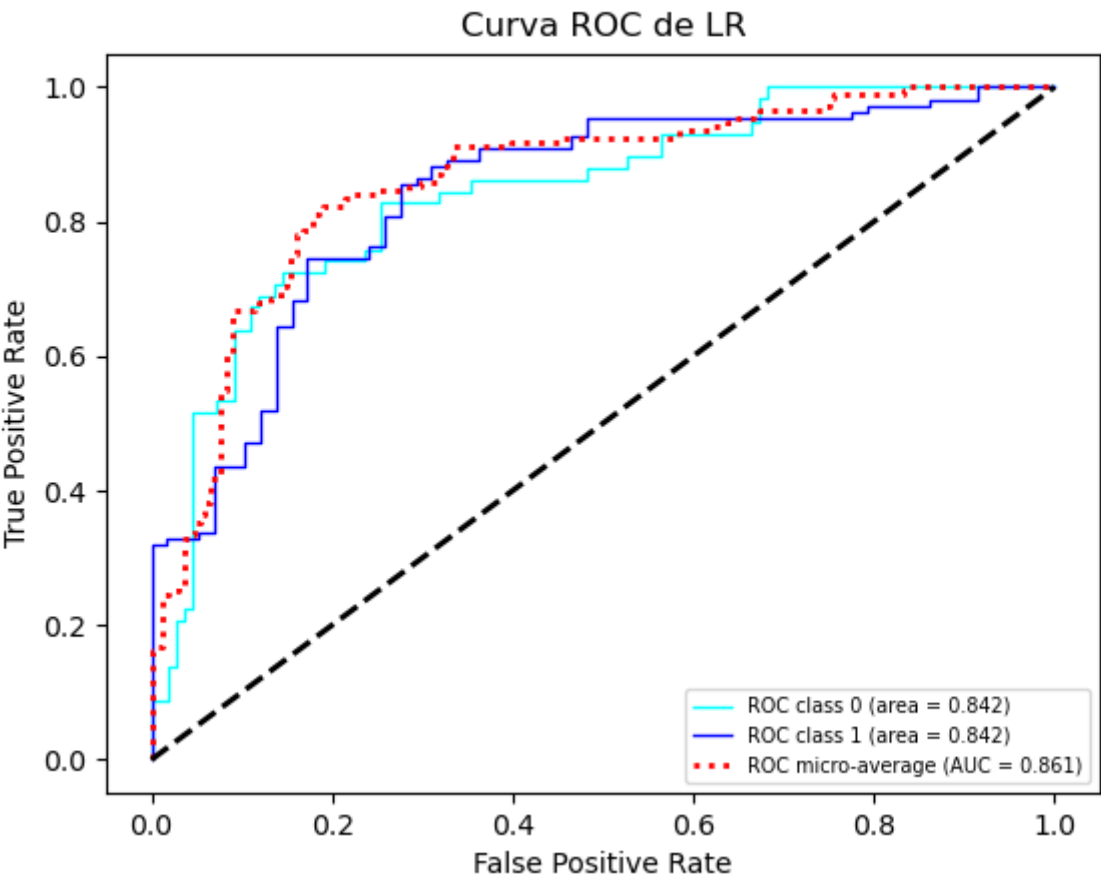


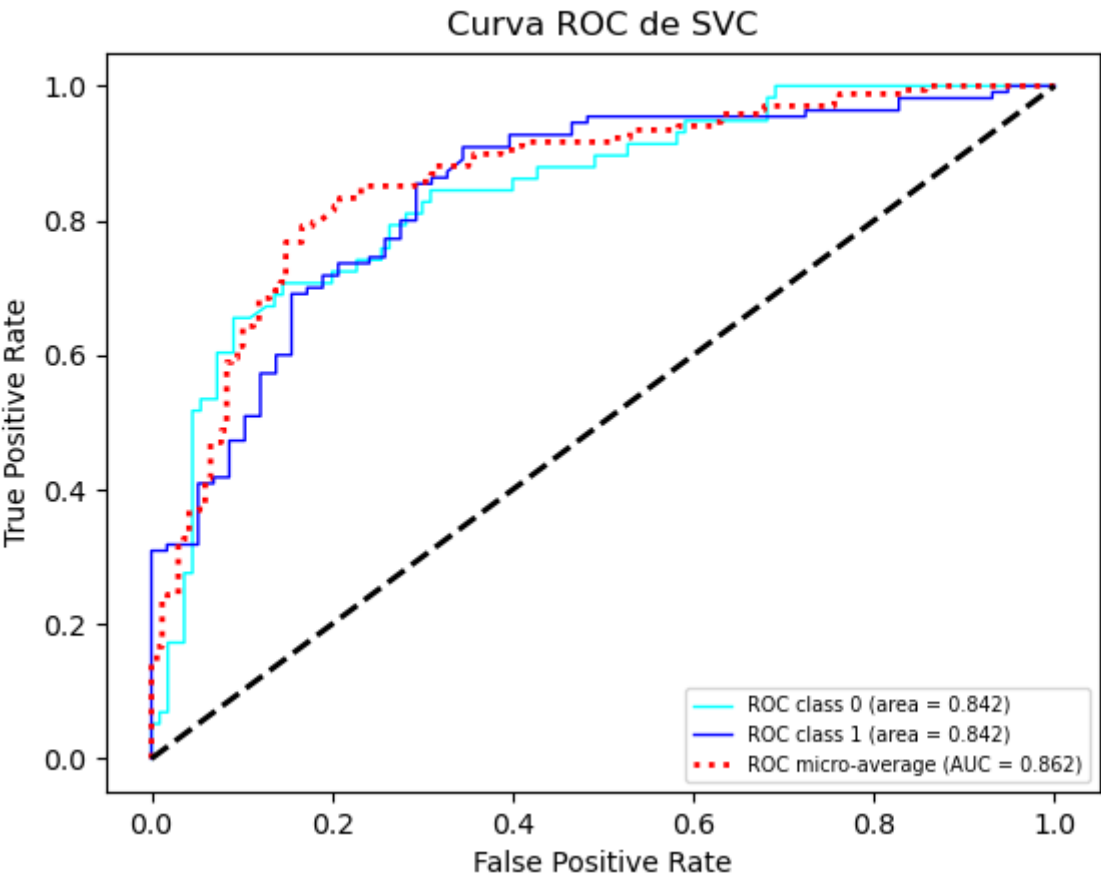
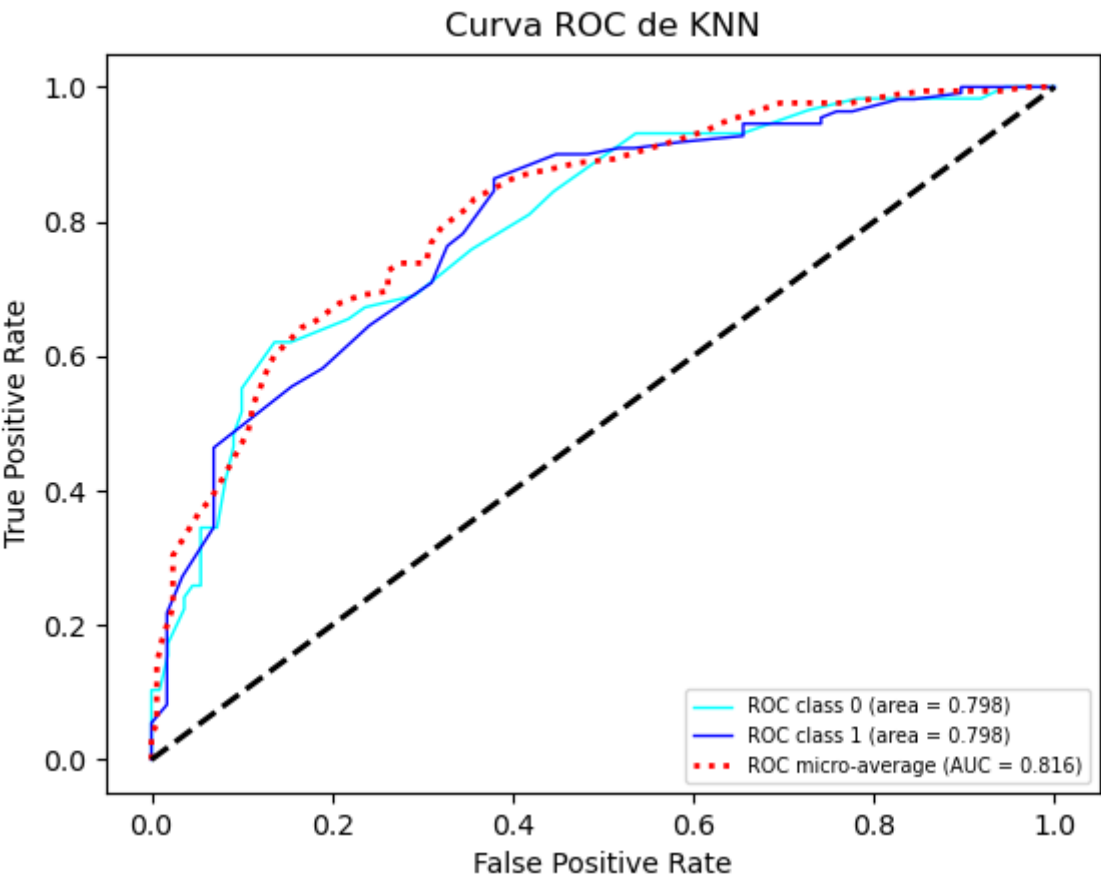
Matriz de confusión DNN

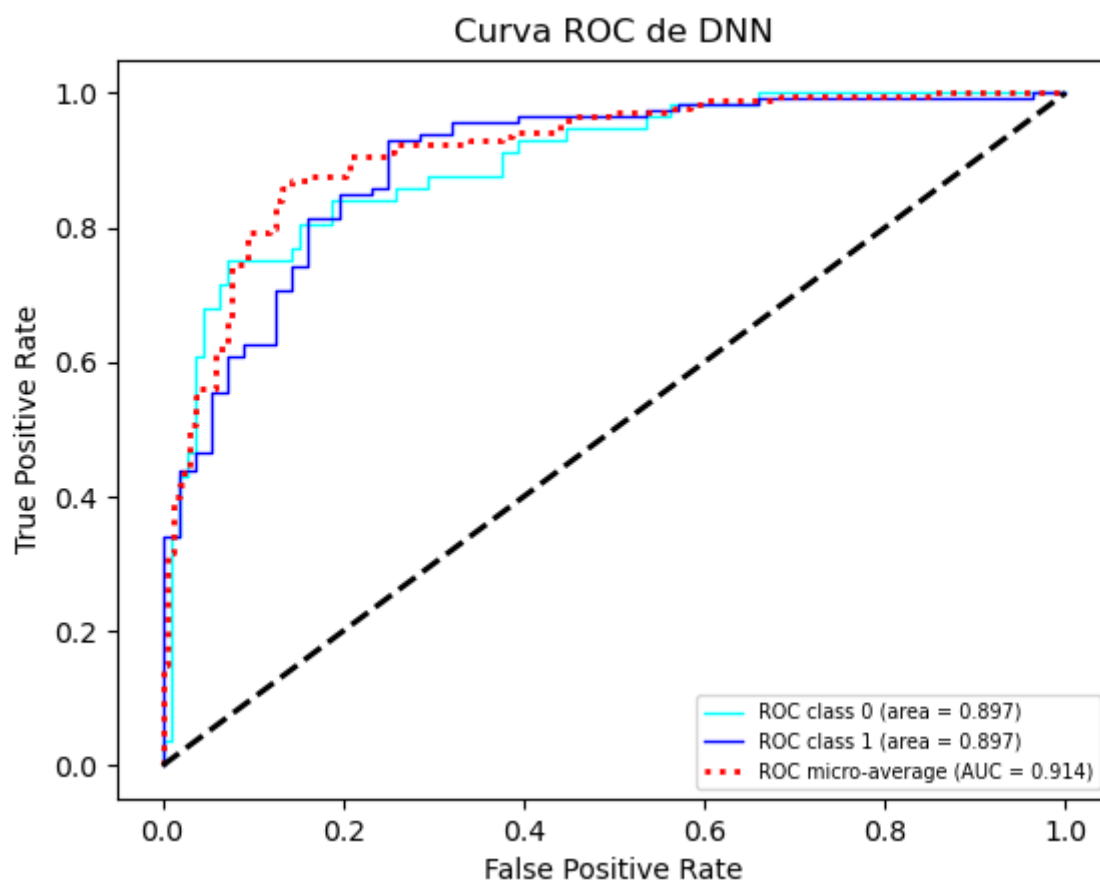
	Clase Real 0	Clase Real 1
Clase Predicha 0	42	14
Clase Predicha 1	10	102

**Curva ROC: Compara el ajuste entre la especificidad y la sensibilidad para cada modelo**

---







### Informe de clasificación: Compara los resultados de cada modelo de clasificación

Classification report for model LR:

	precision	recall	f1-score	support
0	0.75	0.69	0.72	58
1	0.84	0.88	0.86	110
accuracy			0.82	168
macro avg	0.80	0.79	0.79	168
weighted avg	0.81	0.82	0.81	168

Classification report for model LDA:

	precision	recall	f1-score	support
0	0.71	0.60	0.65	58
1	0.81	0.87	0.84	110
accuracy			0.78	168
macro avg	0.76	0.74	0.75	168

weighted avg	0.77	0.78	0.77	168
--------------	------	------	------	-----

Classification report for model KNN:

	precision	recall	f1-score	support
0	0.77	0.34	0.48	58
1	0.73	0.95	0.83	110
accuracy			0.74	168
macro avg	0.75	0.65	0.65	168
weighted avg	0.75	0.74	0.70	168

Classification report for model SVC:

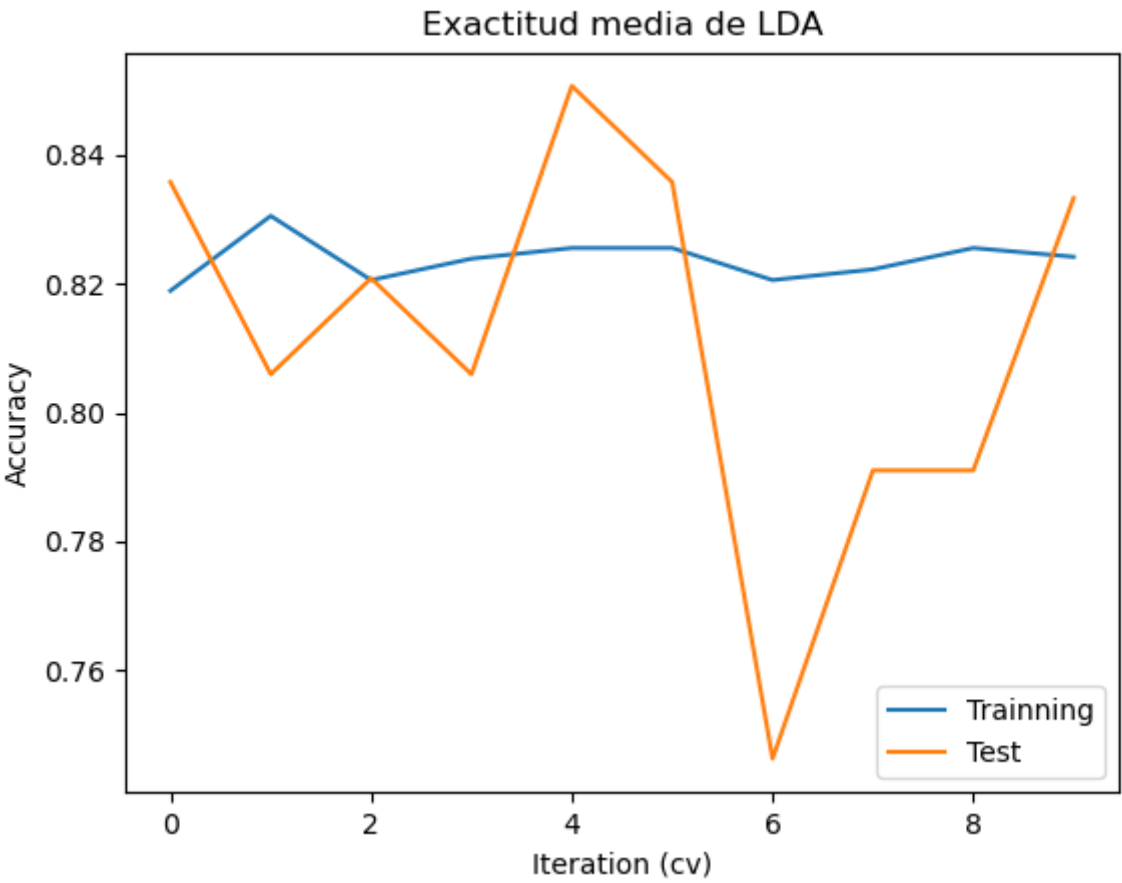
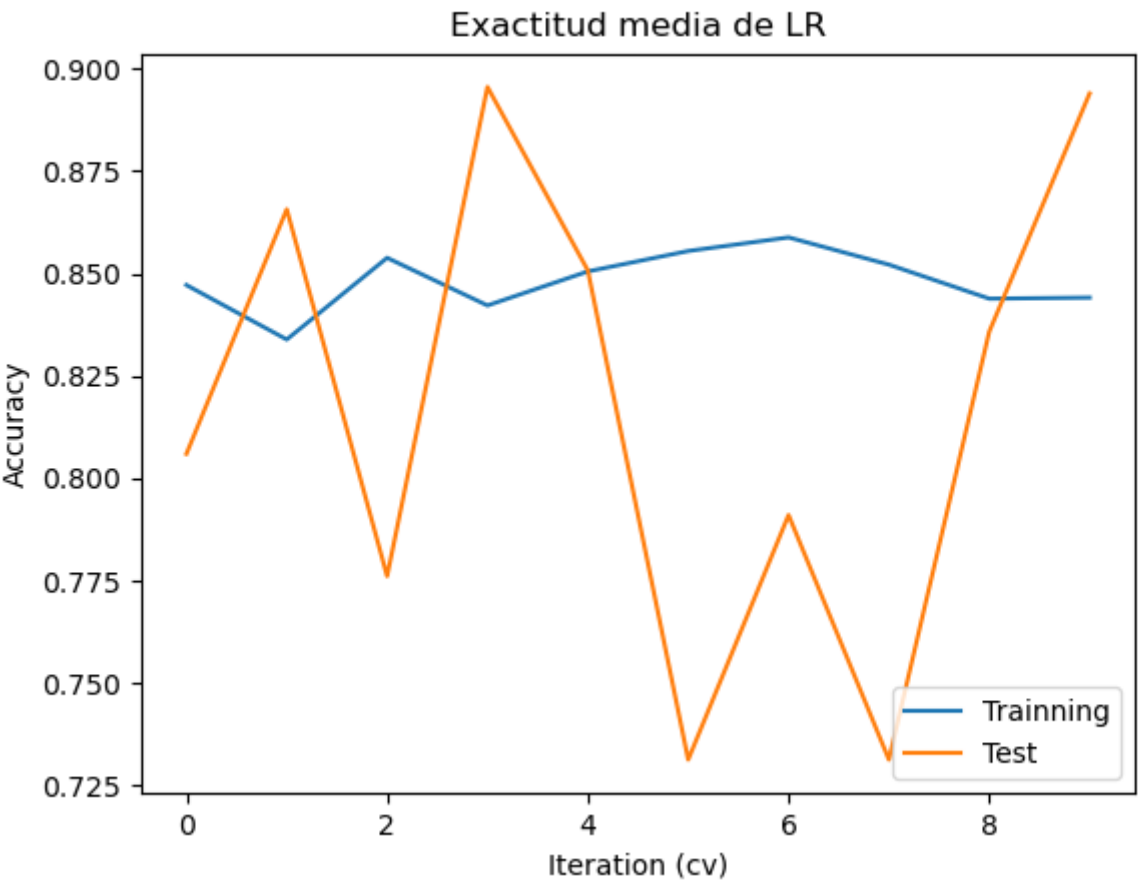
	precision	recall	f1-score	support
0	0.74	0.67	0.70	58
1	0.83	0.87	0.85	110
accuracy			0.80	168
macro avg	0.79	0.77	0.78	168
weighted avg	0.80	0.80	0.80	168

Classification report for model DNN:

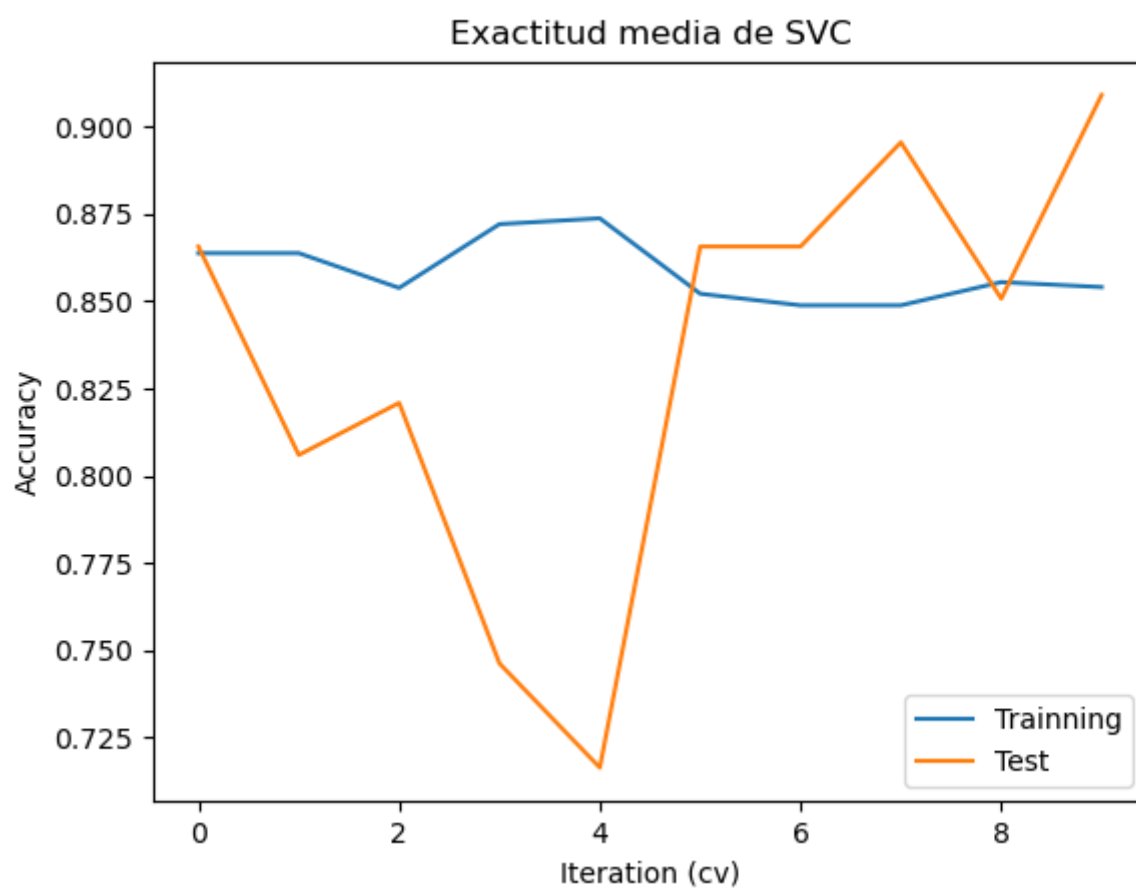
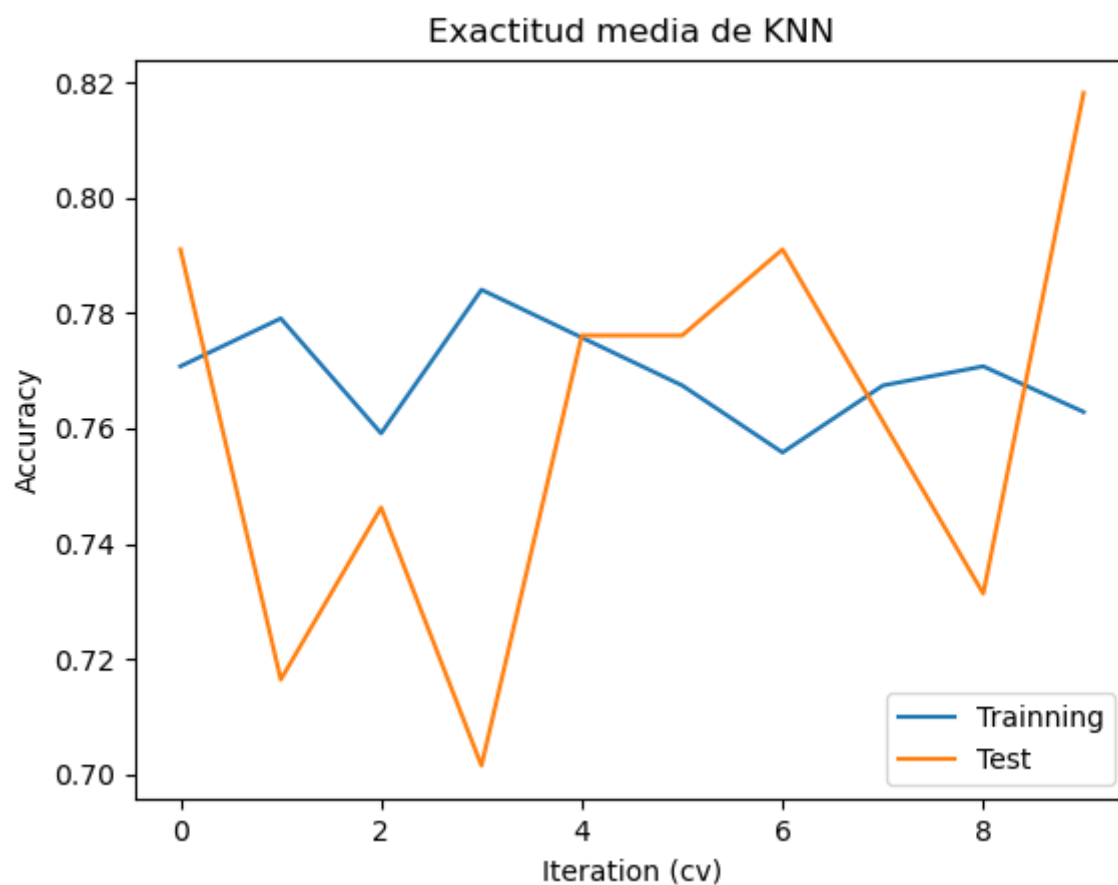
	precision	recall	f1-score	support
0	0.81	0.75	0.78	56
1	0.88	0.91	0.89	112
accuracy			0.86	168
macro avg	0.84	0.83	0.84	168
weighted avg	0.86	0.86	0.86	168

**Exactitud media: Compara la exactitud de cada modelo en función de sus hiperparámetros**

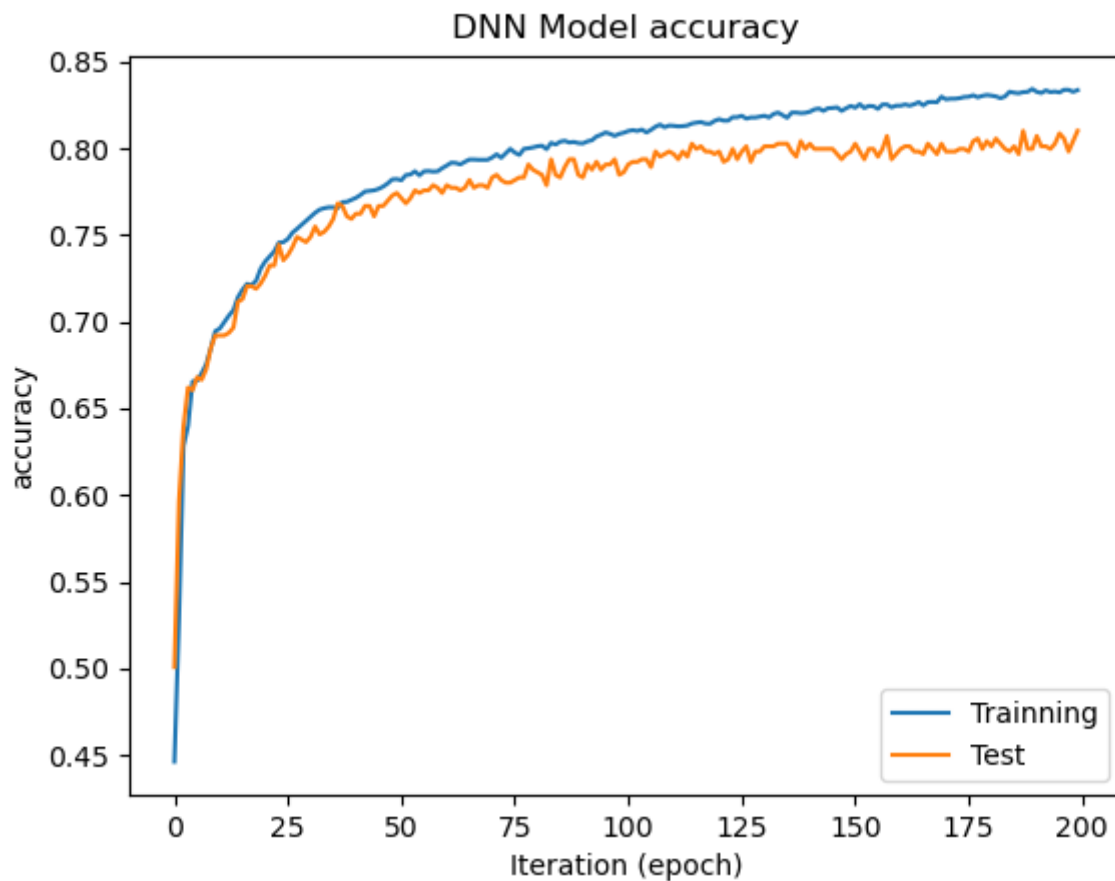
---







**Curva de validación: Compara los resultados del modelo en función de sus hiperparámetros**



**Hiperparámetros: Muestra el dataframe con los hiperparámetros usados en el entrenamiento**

#### Hiperparámetros del modelo LR

	C	max_iter	penalty	solver
hyperparams	7.12741	575	l2	lbfgs

#### Hiperparámetros del modelo LDA

	shrinkage	solver	tol
hyperparams	auto	lsqr	0.00099

#### Hiperparámetros del modelo KNN

	n_neighbors	weights
hyperparams	64	uniform

#### Hiperparámetros del modelo SVC

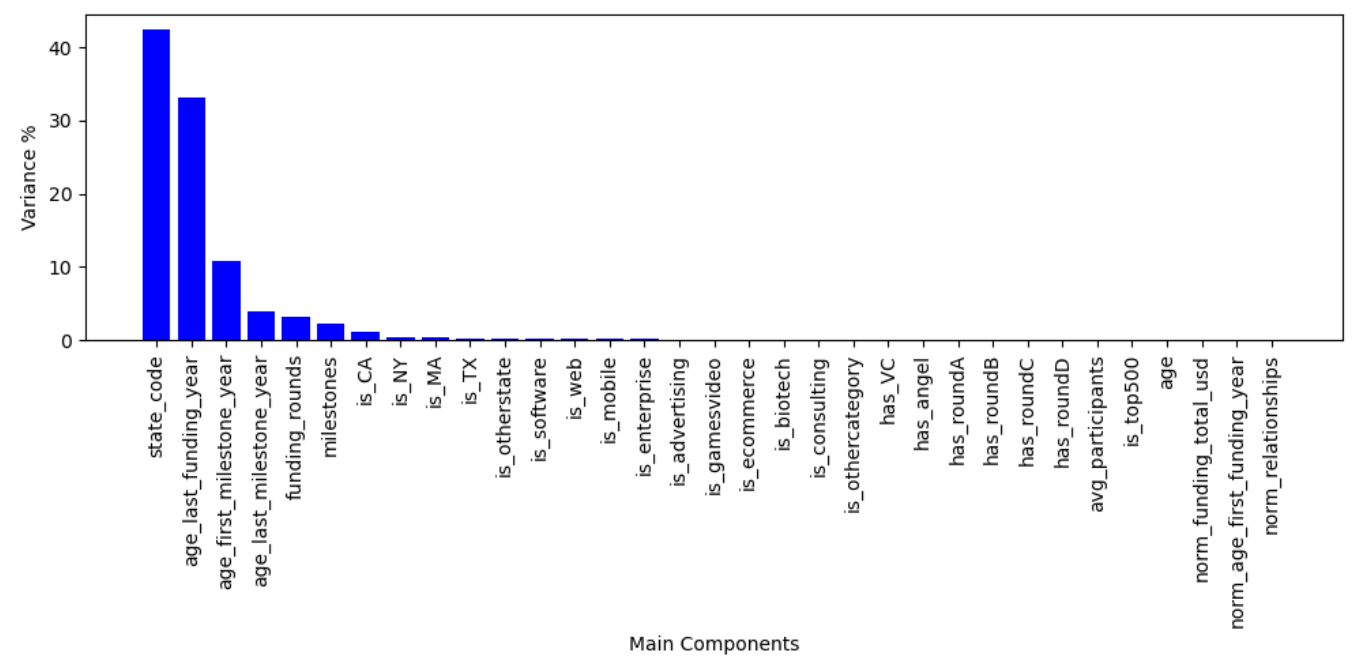
	C	decision_function_shape	gamma	kernel	probability
hyperparams	8.51396	ovr	auto	linear	True

Hiperparámetros de la red neuronal

	neurons	activation
layer 0	11	relu
layer 1	17	sigmoid
layer 2	14	relu
layer 3	11	relu
layer 4	2	softmax
	optimizer	lr
compiler	Adam	0.00017

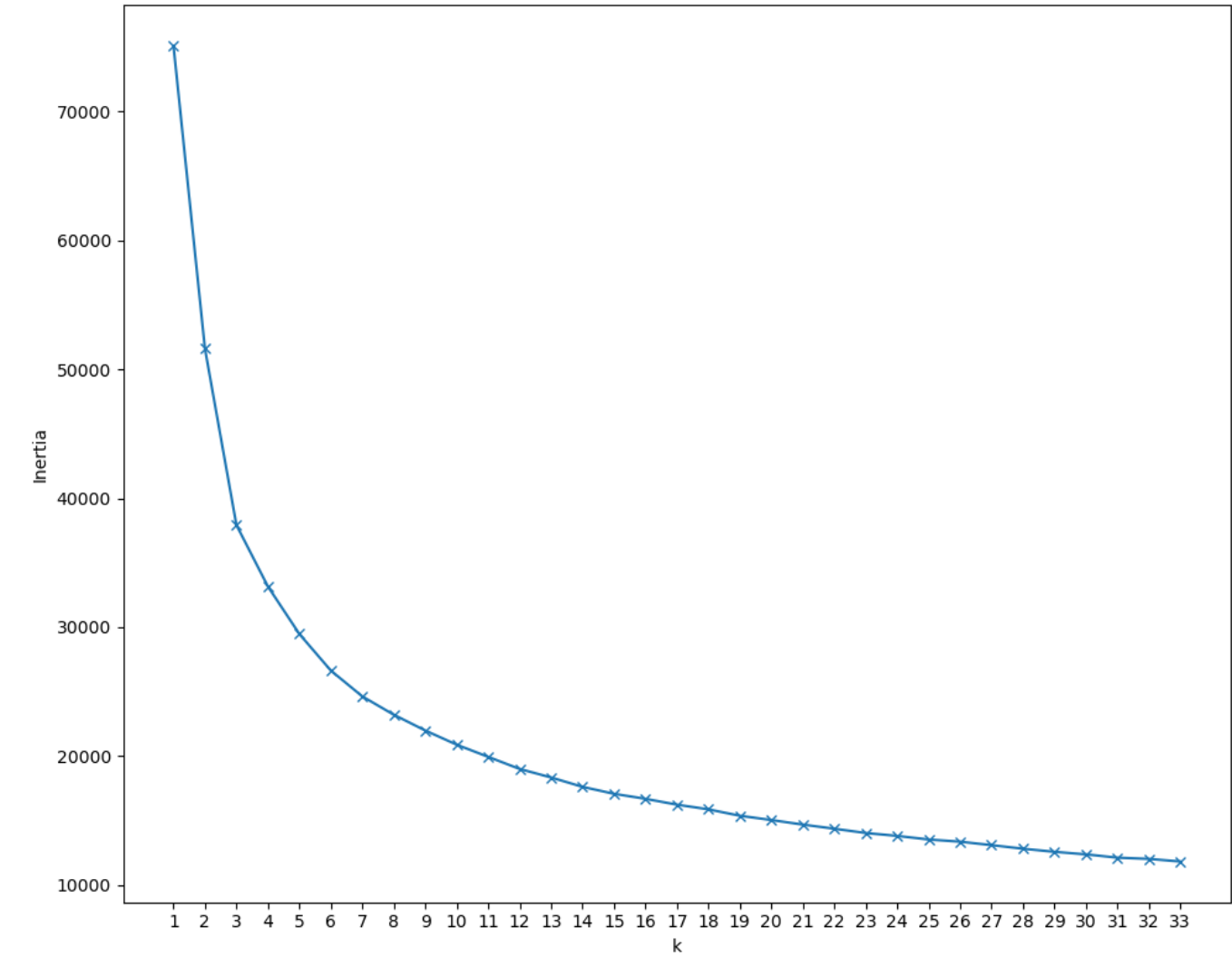
Clasificación No Supervisada

Reducción de la dimensionalidad

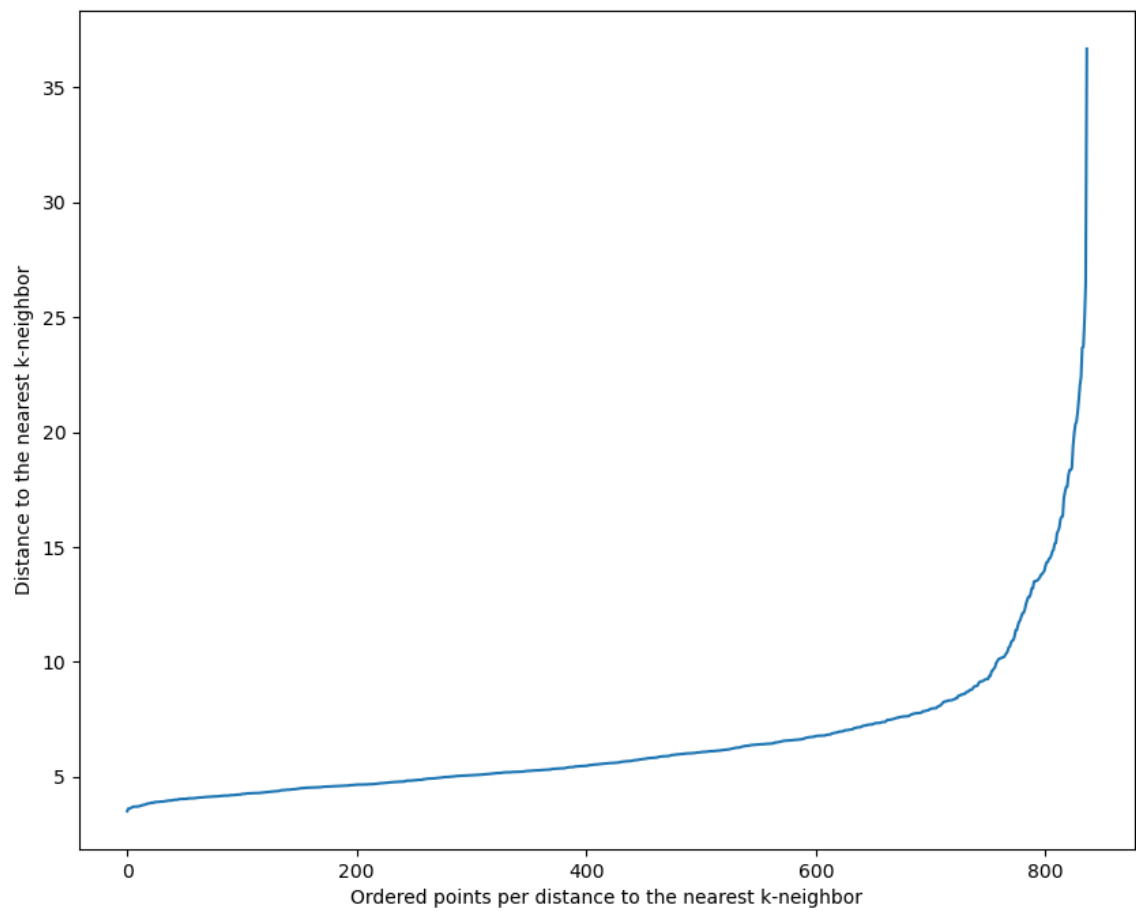


Clustering

K-means clustering



DBSCAN clustering



Detección de anomalías

Isolation Forest - Classification Report

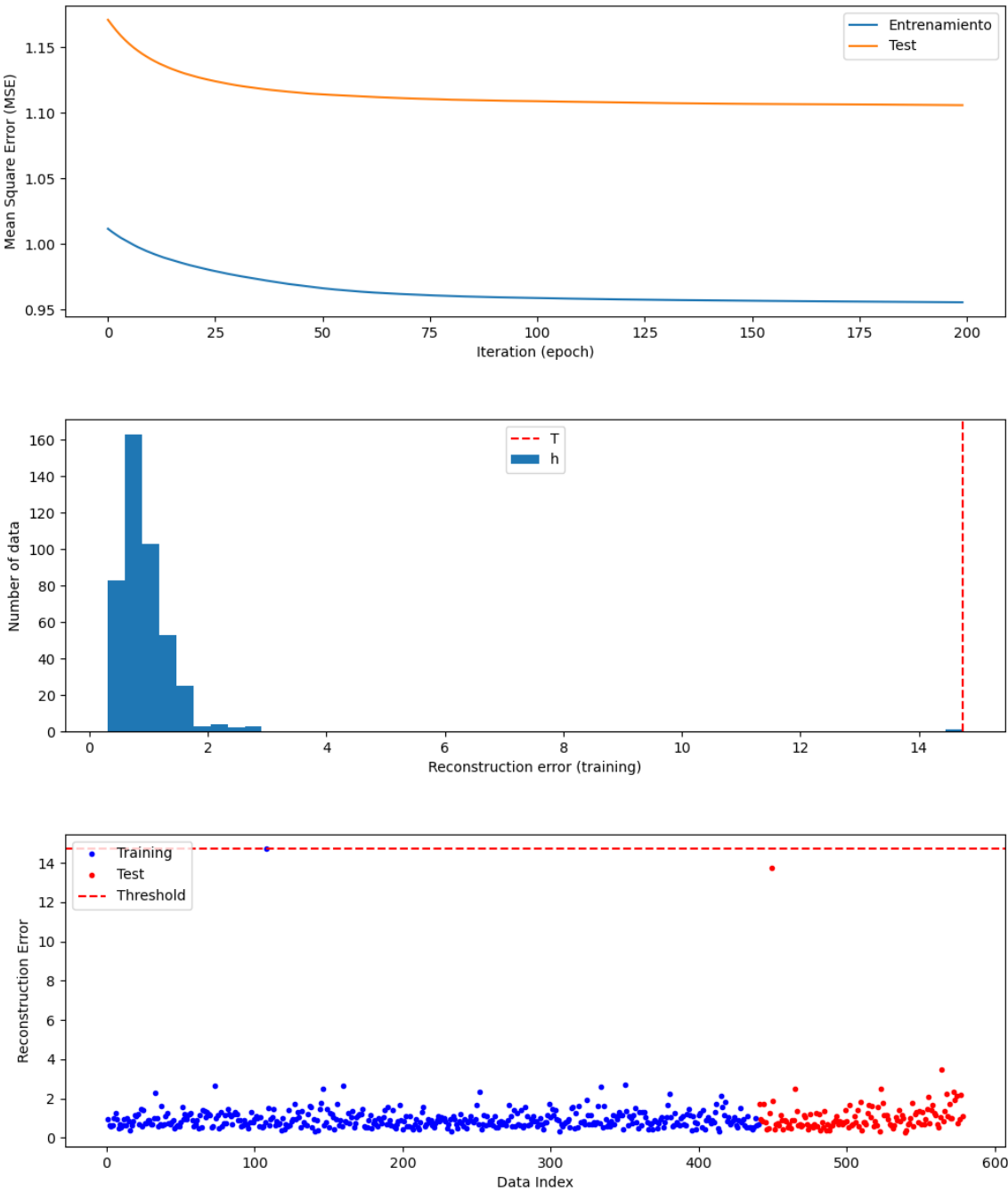
	precision	recall	f1-score	support
-1	0.00	0.00	0.00	0
0	0.00	0.00	0.00	28
1	0.80	0.99	0.88	110
accuracy			0.79	138
macro avg	0.27	0.33	0.29	138
weighted avg	0.63	0.79	0.70	138

Local Oulier Factor (LOF) - Classification Report

	precision	recall	f1-score	support
-1	0.00	0.00	0.00	0

0	0.00	0.00	0.00	28
1	0.80	0.99	0.88	110
accuracy			0.79	138
macro avg	0.27	0.33	0.29	138
weighted avg	0.63	0.79	0.70	138

Autoencoding



Autoencoding Classification Report

--	--	--	--	--	--

	precision	recall	f1-score	support
0	0.00	0.00	0.00	28
1	0.80	1.00	0.89	110
accuracy			0.80	138
macro avg	0.40	0.50	0.44	138
weighted avg	0.64	0.80	0.71	138

## Registro de tiempos de entrenamiento

---

El entrenamiento se ha realizado en computación multihilo

### Tiempos del entrenamiento de los modelos de clasificación supervisada

	LR	LDA	KNN	SVC	DNN
Time Models	2.32	2.04	8.99	9.46	740.15

## 5. Conclusiones

De acuerdo con las novedades incluidas en el presente documento, extensible de la primera parte:

- [Startup Success Prediction](#)

Se pueden añadir un grupo de conclusiones.

En cuanto a las técnicas de reducción de dimensionalidad, se puede estudiar el grado de influencia entre las variables del conjunto de datos mediante la matriz de correlación y analizar los componentes principales que explican la mayor varianza de los mismos desarrollando una PCA.

Tanto las técnicas de agrupamiento por particiones como por densidad, explicadas anteriormente: K-means y DBSCAN, permiten tratar el estudio como un modelo de aprendizaje no supervisado donde no se utilizan etiquetas y el objetivo consiste en encontrar agrupaciones óptimas para los datos. Para determinar el número de clusters necesarios, el K-means se basa en el conocido como método elbow para calcular las inercias o distancias al cuadrado de los datos con sus respectivos centroides, indentificando el número para el cual se produce la mayor disminución consecutiva. Análogamente, el DBSCAN utiliza el algoritmo de K-Nearest Neighbours para extraer el valor del hiperparámetro epsilon del modelo, que hace referencia a la densidad de punto empleada.

Por lo que a las funciones declaradas en la la clase dedicada a la detección de anomalías respecta, merece la pena destacar la necesidad de adaptar el conjunto de datos para poder tratarlo como un problema de este tipo. Esto pasa por escoger una clase mayoritaria como normal, en este caso que las empresas tengan éxito, aprovechando que se tienen más muestras globales con esta etiqueta, y otra minoritaria que se considera la clase anómala y hace referencia al fracaso de las startup. Se trata de un modelo de aprendizaje semisupervisado, ya que se utiliza una partición de los datos para entrenamiento que no contiene muestras anómalas y una para test que presenta una pequeña porción de las mismas. Dentro de este apartado se definen una serie de funciones que proponen la operación de algunos de los principales métodos empleados en este tipo de problemas: Isolation Forest, basado en los algoritmos utilizados en los bosques aleatorios para aislar las muestras con características diferentes al resto; LOF, utiliza la densidad de los datos para identificar las muestras atípicas comparándola con la de sus vecinos y una red de neuronas autoencoder, que durante el proceso de entrenamiento extrae las características esenciales de los datos, como consecuencia de la definición de la capa oculta de menor dimensión, y las utiliza para reconstruirlos a la salida, de esta manera las muestras anómalas presentes en el conjunto de test producirán un error de reconstrucción alto pudiendo ser identificadas en base a un umbral.

Si se analizan las partes del programa multihilo en comparativa con sus versiones secuenciales, se puede concluir una aceleración del proceso, como es el caso de la ejecución paralela en varios cores de los k experimentos de la validación cruzada para la red neuronal, o el proceso de optimización por modelo, teniendo en cuenta el número de intentos para cada uno. Esta ampliación favorece notablemente los aspectos del programa que abordan tareas recursivas, mejorando la organización de la ejecución de las instrucciones, que es donde se implantan. Este análisis tiene en cuenta los hilos virtuales configurados en el el código fuente, sin embargo y como cabría esperar, los resultados pueden diferir de una CPU a otra según el hardware.

Para terminar con el apartado de conclusiones, no se puede dejar de lado el interés y la importancia que presenta el hecho de realizar mediciones de tiempos de ejecución, especialmente allí donde el proceso demanda más recursos a la CPU, como es el caso del proceso de optimización donde se desarrollan varios



entrenamientos utilizando validación cruzada. Esta mejora, proporciona mayor trazabilidad al sistema y explica una comparativa con los modelos más difícilmente escalables de cara a encontrar una solución representativa al problema de optimización, pues el tiempo de ejecución sería superior.