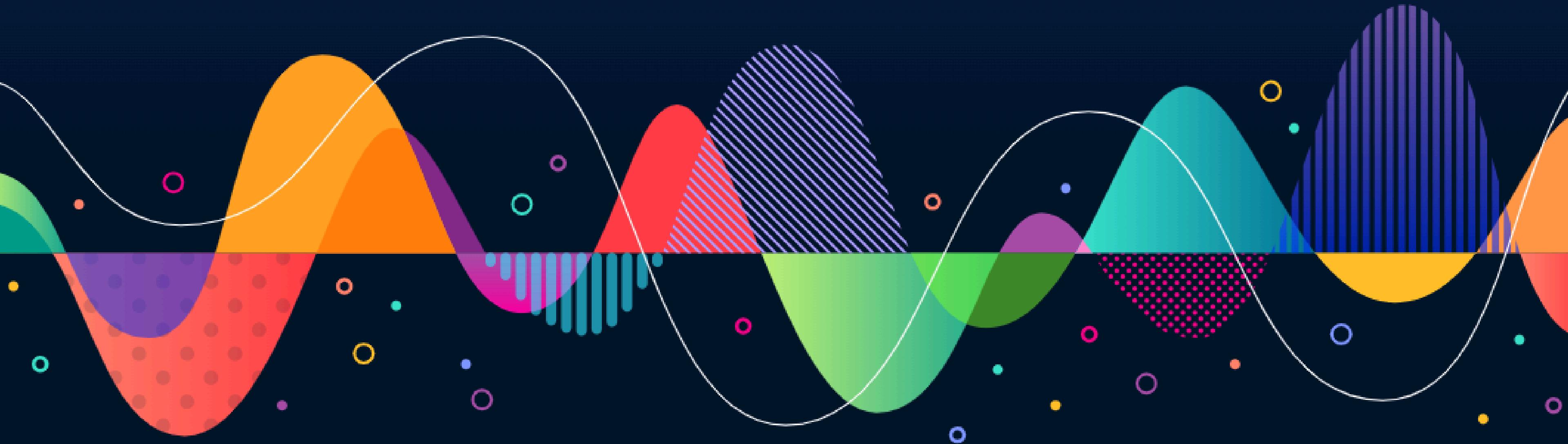


Team:

- Maricarmen Daniela Barillas Duarte a01369993@tec.mx
- Daniel Alejandro Olivares Ángeles a01754838@tec.mx
- Gael Eduardo Pérez Gómez a01753336@tec.mx
- Santiago Martínez Vallejo a00571878@tec.mx

# Kaggle – Titanic classification

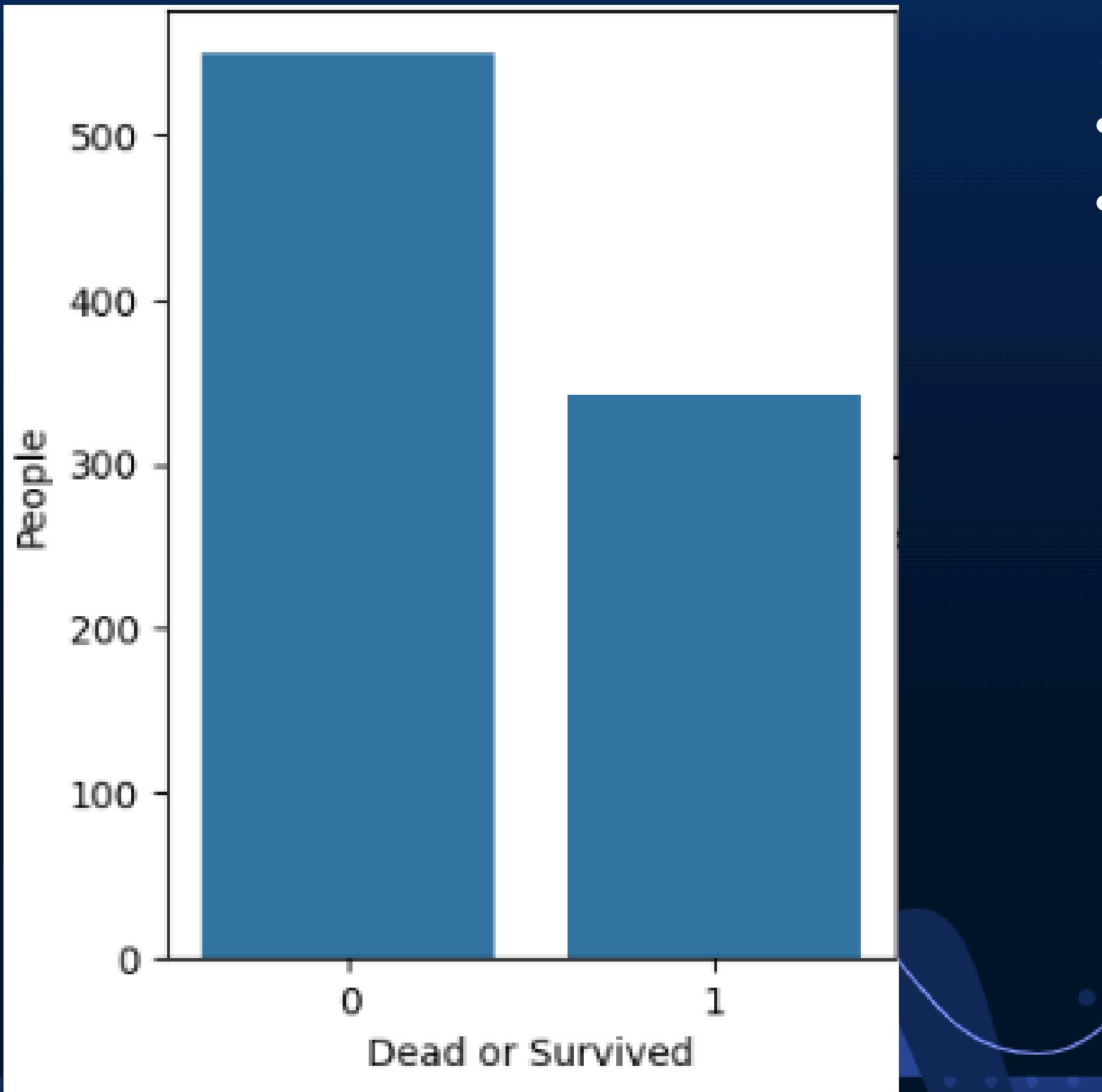


# Problem & Goals

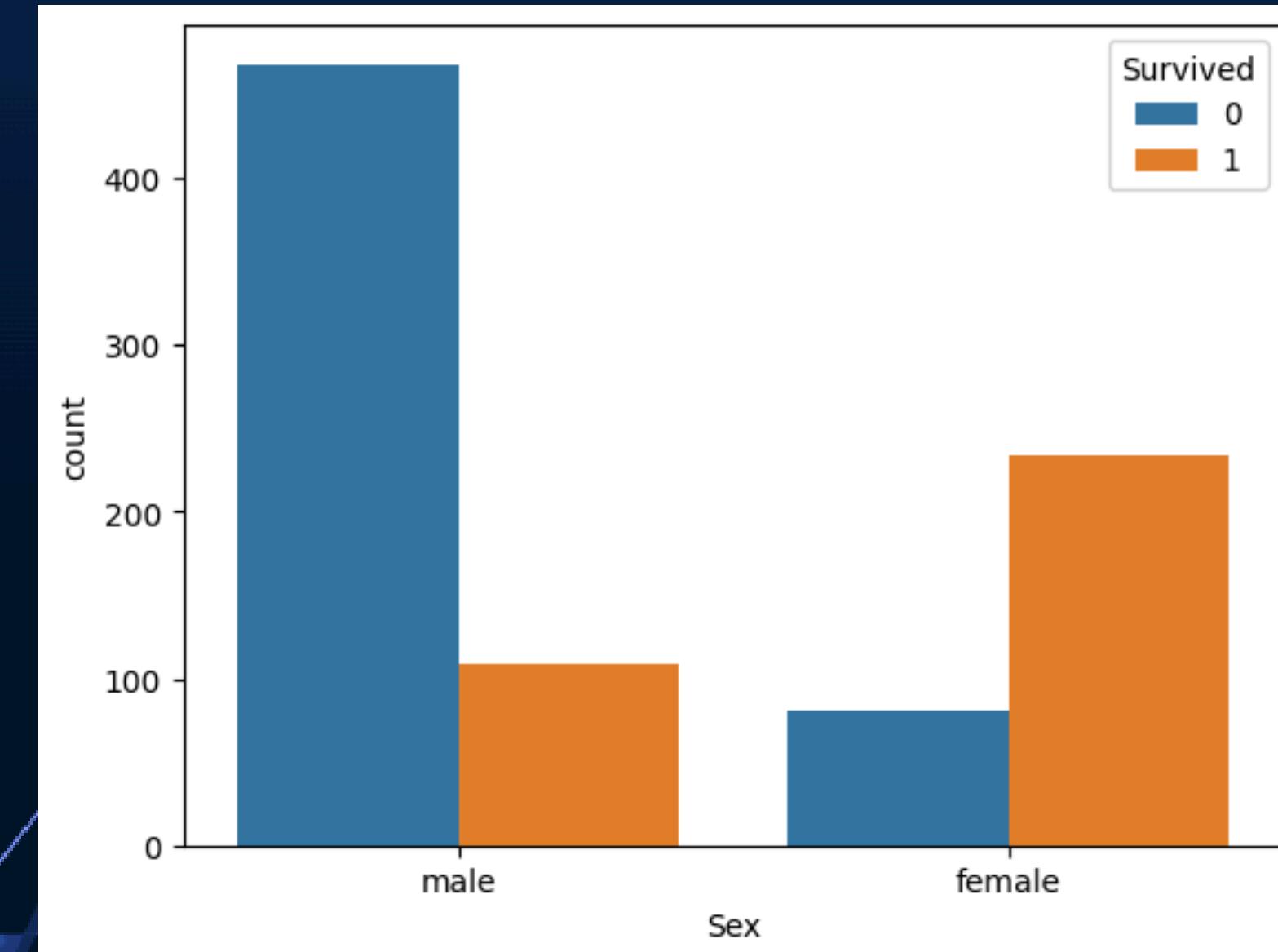


- Identify the most relevant features that determine a passengers survival probability.
- Select the best model for predicting the survival probability.
  - Identify the model with an accuracy of 85.
  - Consistent K-cross fold validation avarage.
- Hypothesis: Random Forest will be the best model.
  - Randomness in sampling of data points and selection of features. Less overfitting and improved generalization.

# Data Distribution



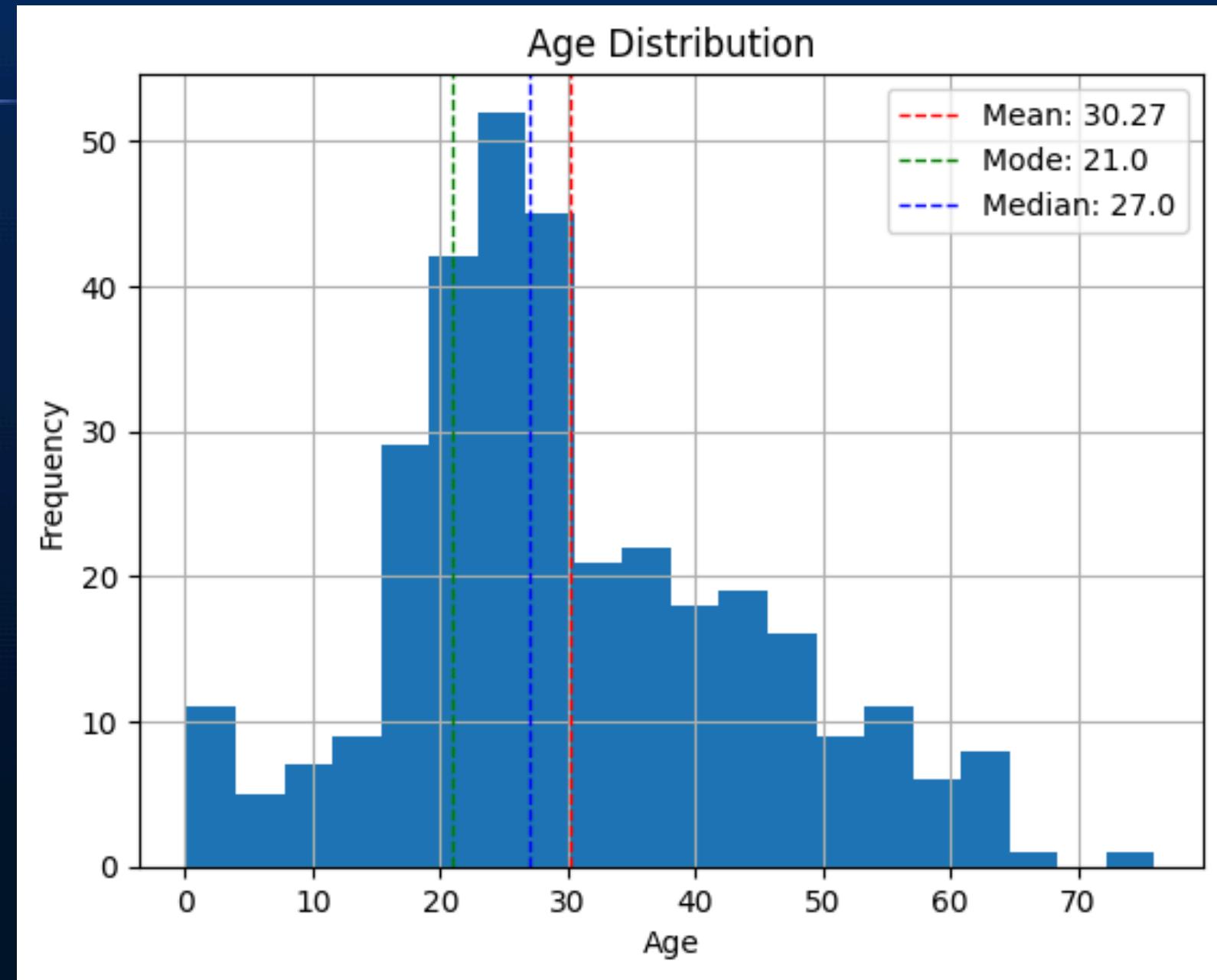
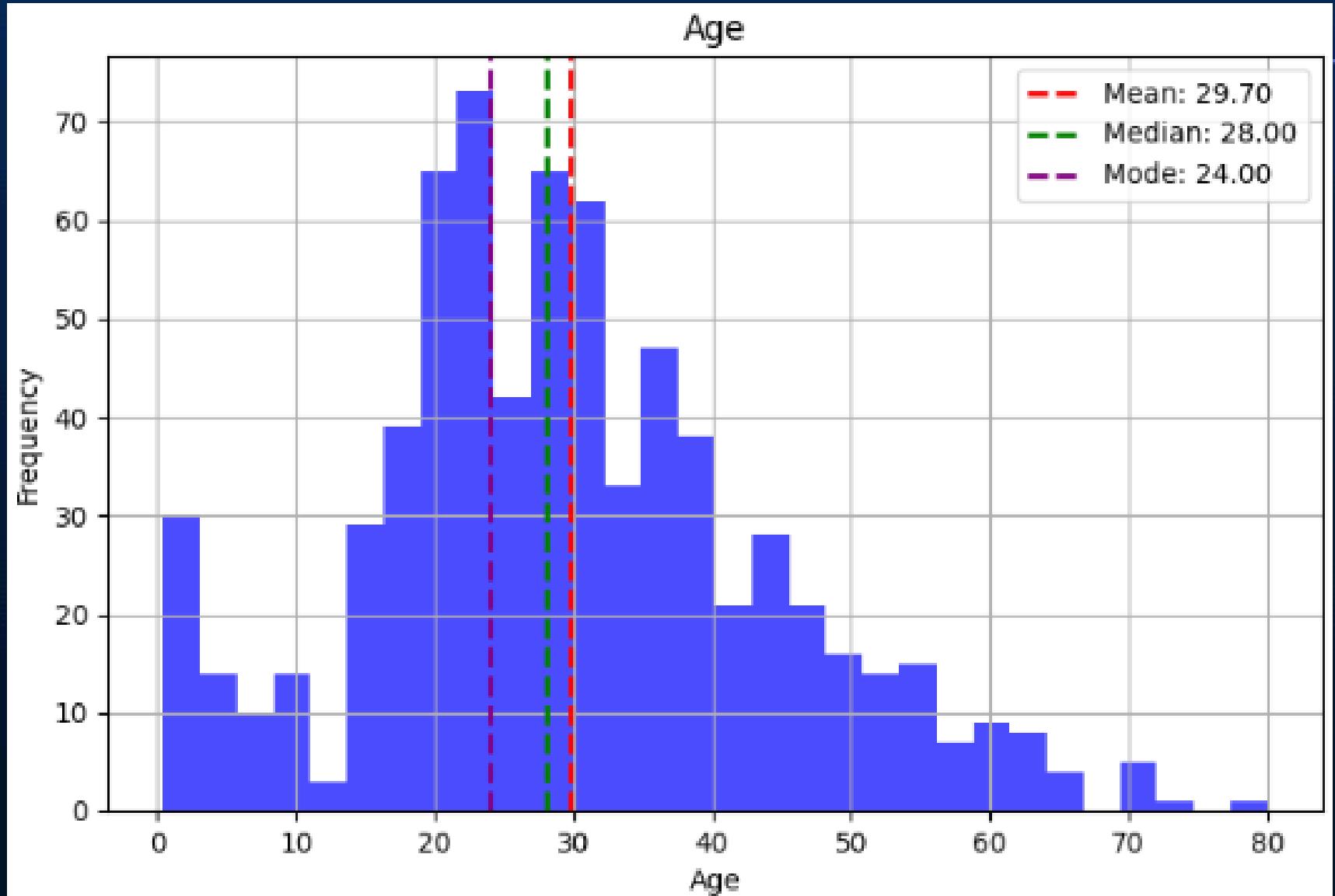
- The data leans on dead rather than not survived.
- Sex distribution with survived or dead.



## Training Set

# Missing Data

## Test Set

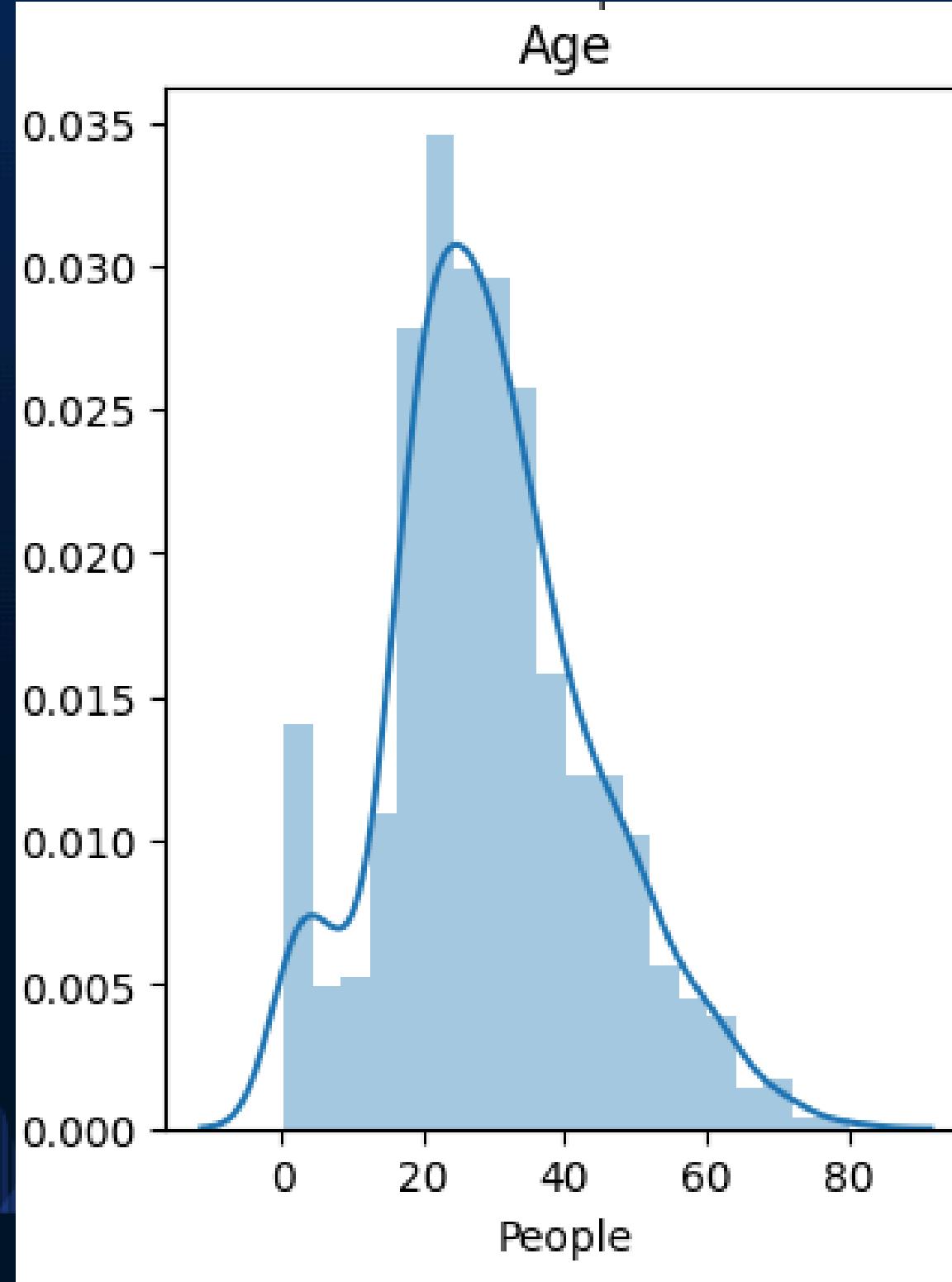


- Age: 19.9% missing
- Median for missing values

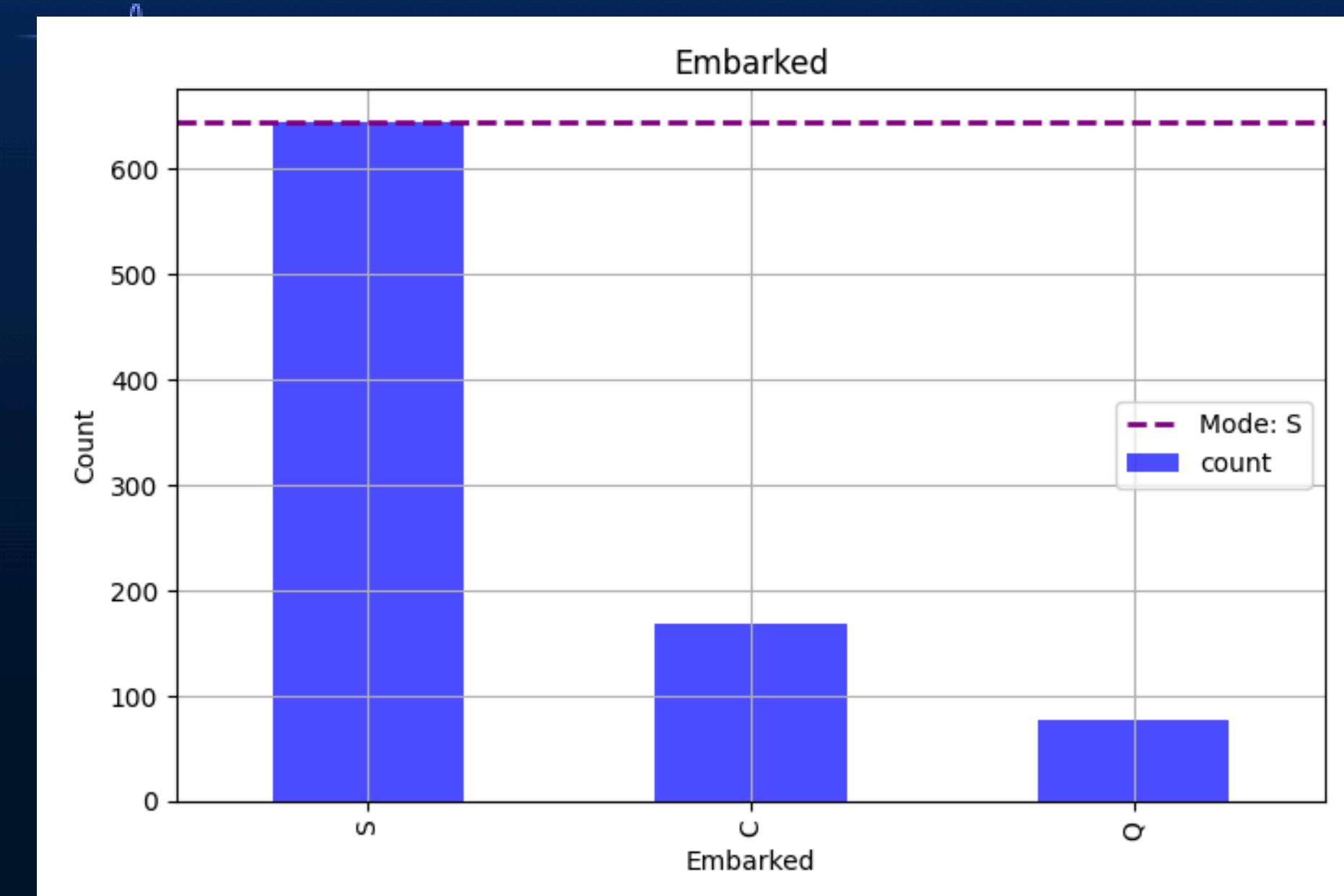
- Age: 20.6% missing
- Median for missing values

# Missing Data

## Age distribution

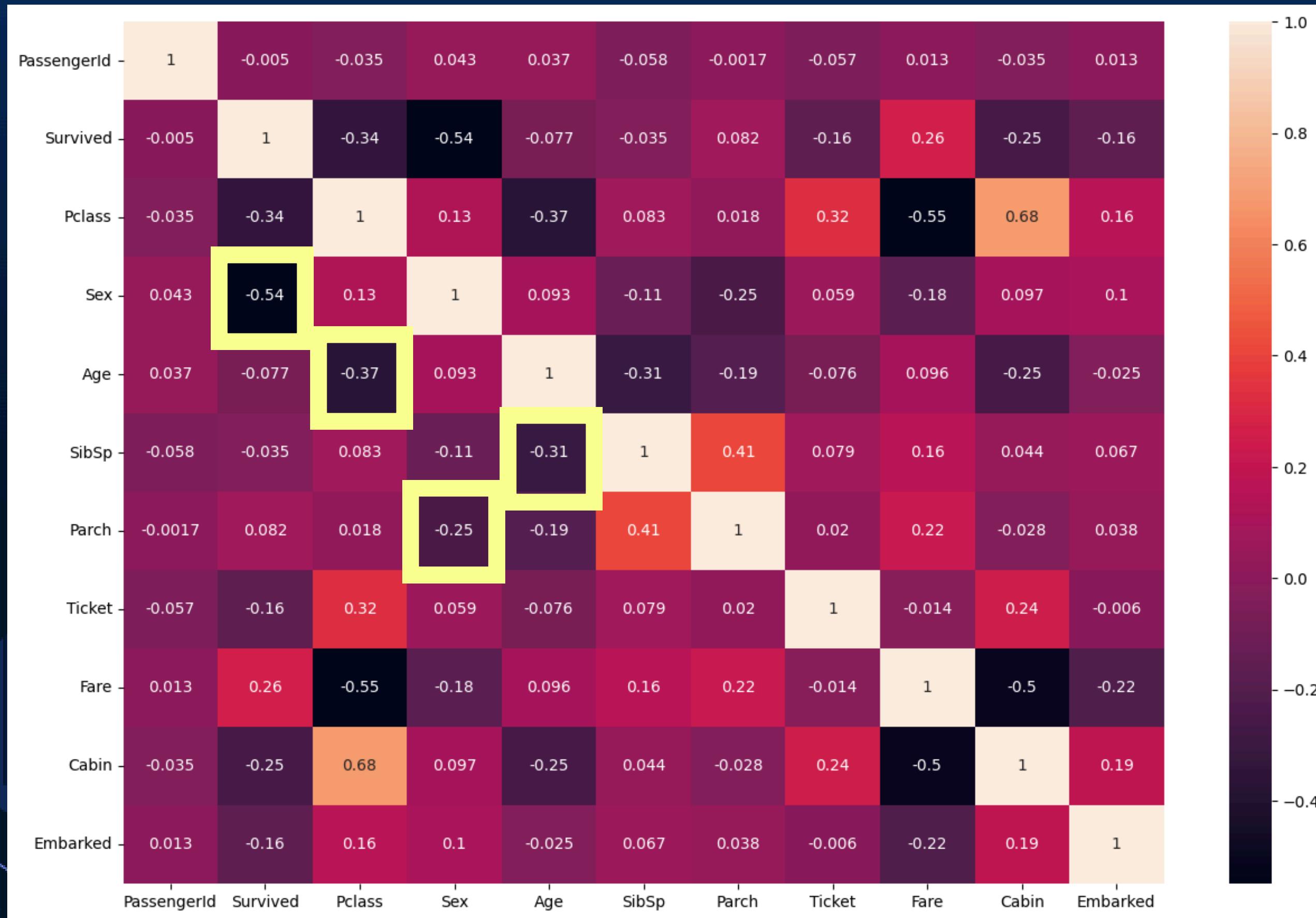


## Embarked



- Embarked: 2 values missing
- Mode for missing values

# Correlation Analysis

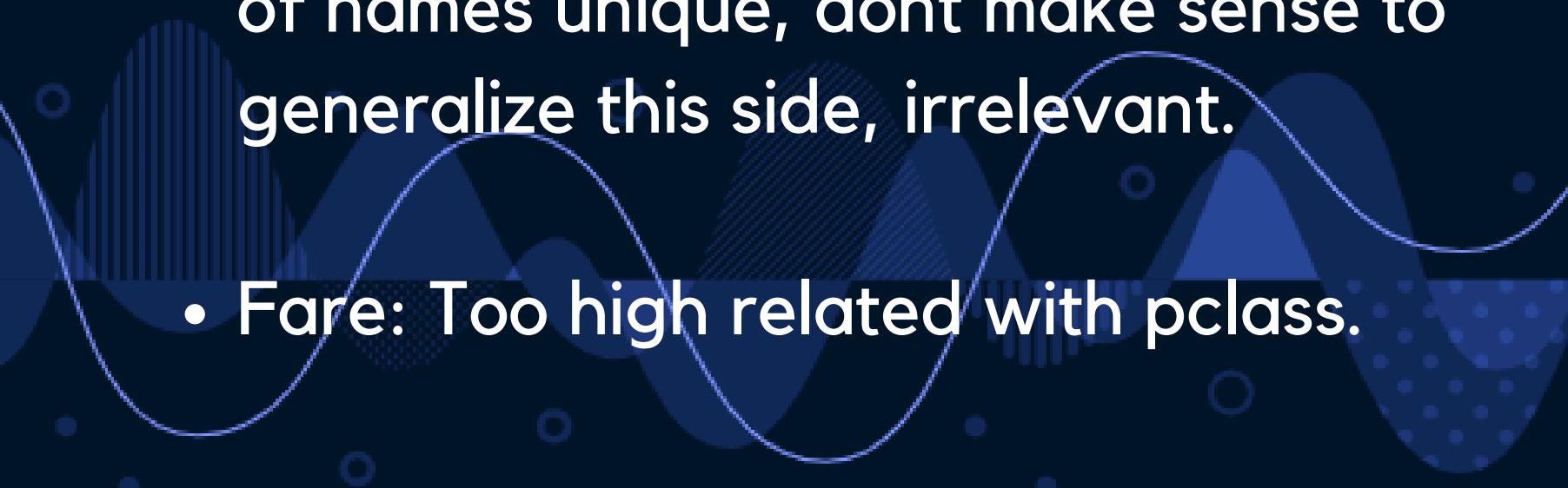


- Sex
- Age
- Pclass
- Embarked
- Parch + Sibsp \*

\* go to Data transformation

# Dropped Features

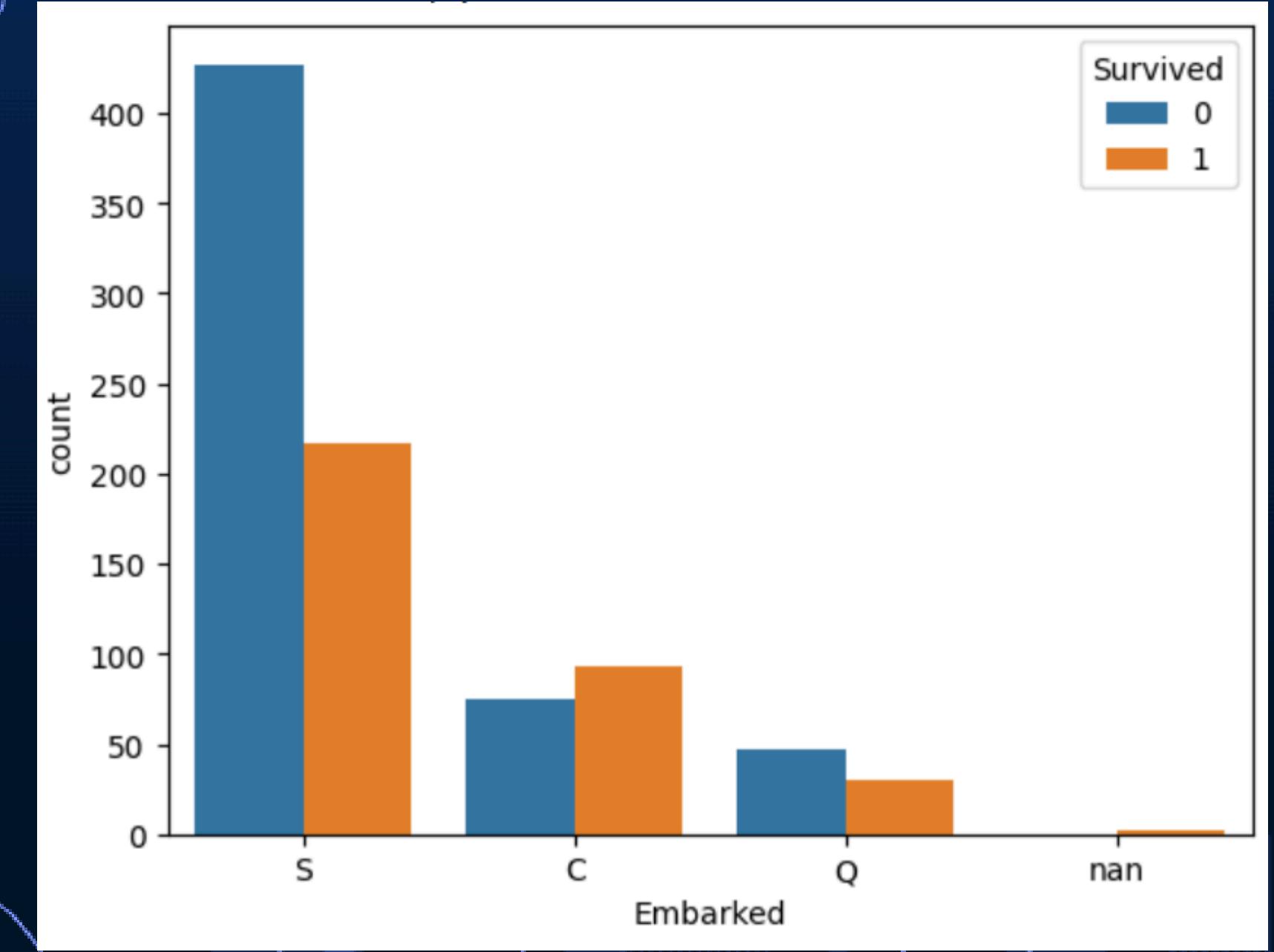
- Passenger Id: Irrelevant feature, just an assignment to the tickets.
- Ticket and Cabin: Both are related in the correlation matrix with Pclass, so we can prescend with both, also, there are a lot of missing values in cabin.
- Name: As we saw before, there are a lot of names unique, dont make sense to generalize this side, irrelevant.
- Fare: Too high related with pclass.



PassengerId	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
PassengerId	1	-0.005	-0.035	0.043	0.037	-0.058	-0.0017	-0.057	0.013	-0.035	0.013
Survived	-0.005	1	-0.34	-0.54	-0.077	-0.035	0.082	-0.16	0.26	-0.25	-0.16
Pclass	-0.035	-0.34	1	0.13	-0.37	0.083	0.018	0.32	-0.55	0.68	0.16
Sex	0.043	-0.54	0.13	1	0.093	-0.11	-0.25	0.059	-0.18	0.097	0.1
Age	0.037	-0.077	-0.37	0.093	1	-0.31	-0.19	-0.076	0.096	-0.25	-0.025
SibSp	-0.058	-0.035	0.083	-0.11	-0.31	1	0.41	0.079	0.16	0.044	0.067
Parch	-0.0017	0.082	0.018	-0.25	-0.19	0.41	1	0.02	0.22	-0.028	0.038
Ticket	-0.057	-0.16	0.32	0.059	-0.076	0.079	0.02	1	-0.014	0.24	-0.006
Fare	0.013	0.26	-0.55	-0.18	0.096	0.16	0.22	-0.014	1	-0.5	-0.22
Cabin	-0.035	-0.25	0.68	0.097	-0.25	0.044	-0.028	0.24	-0.5	1	0.19
Embarked	0.013	-0.16	0.16	0.1	-0.025	0.067	0.038	-0.006	-0.22	0.19	1

# Data Transformation

- **Sex:**
  - male -> 1
  - female -> 0
- **Embarked:**
  - C -> 1
  - Q -> 1
  - S -> 0
- **fit\_transform** from scikit learn to standardize the features.

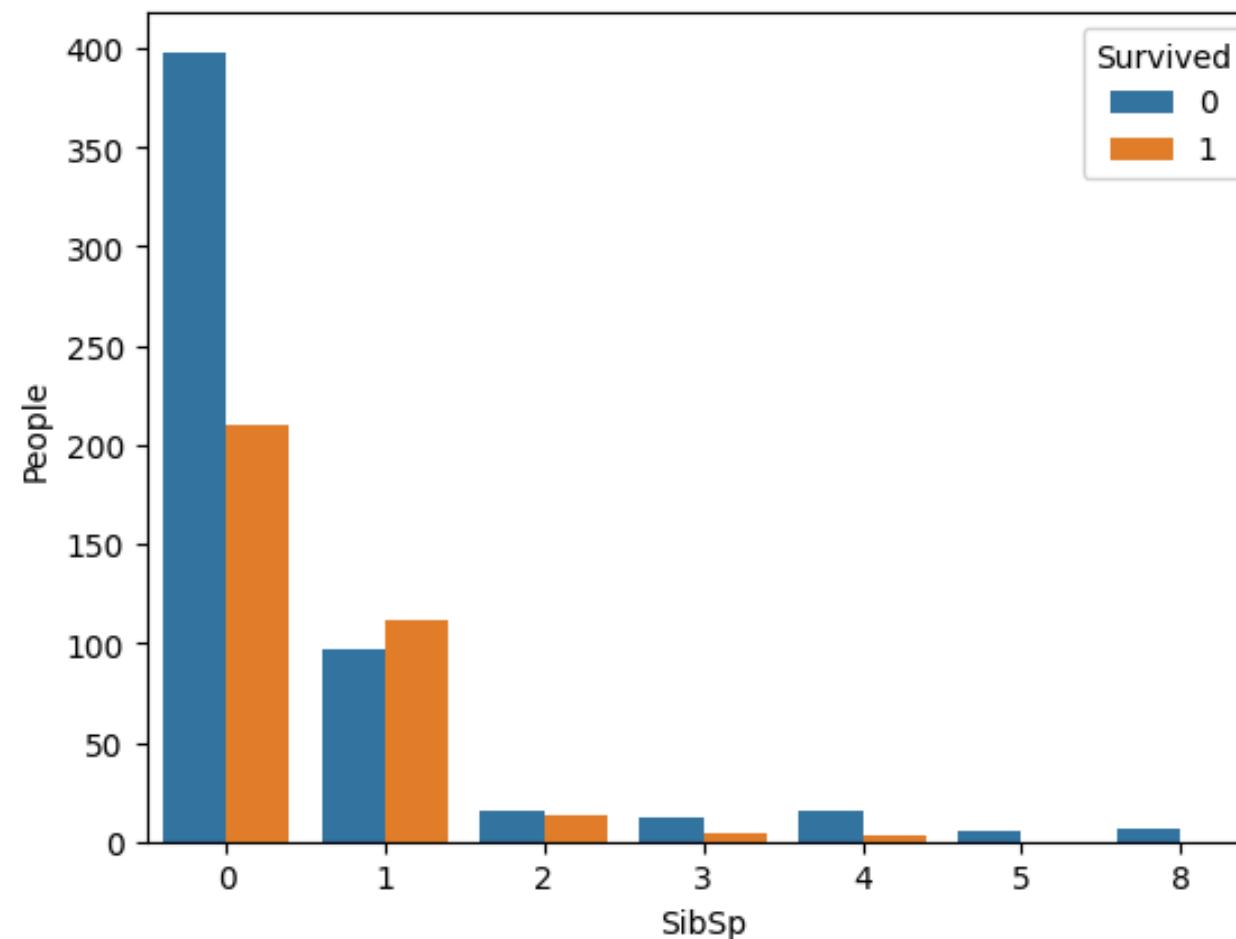


# Data Transformation

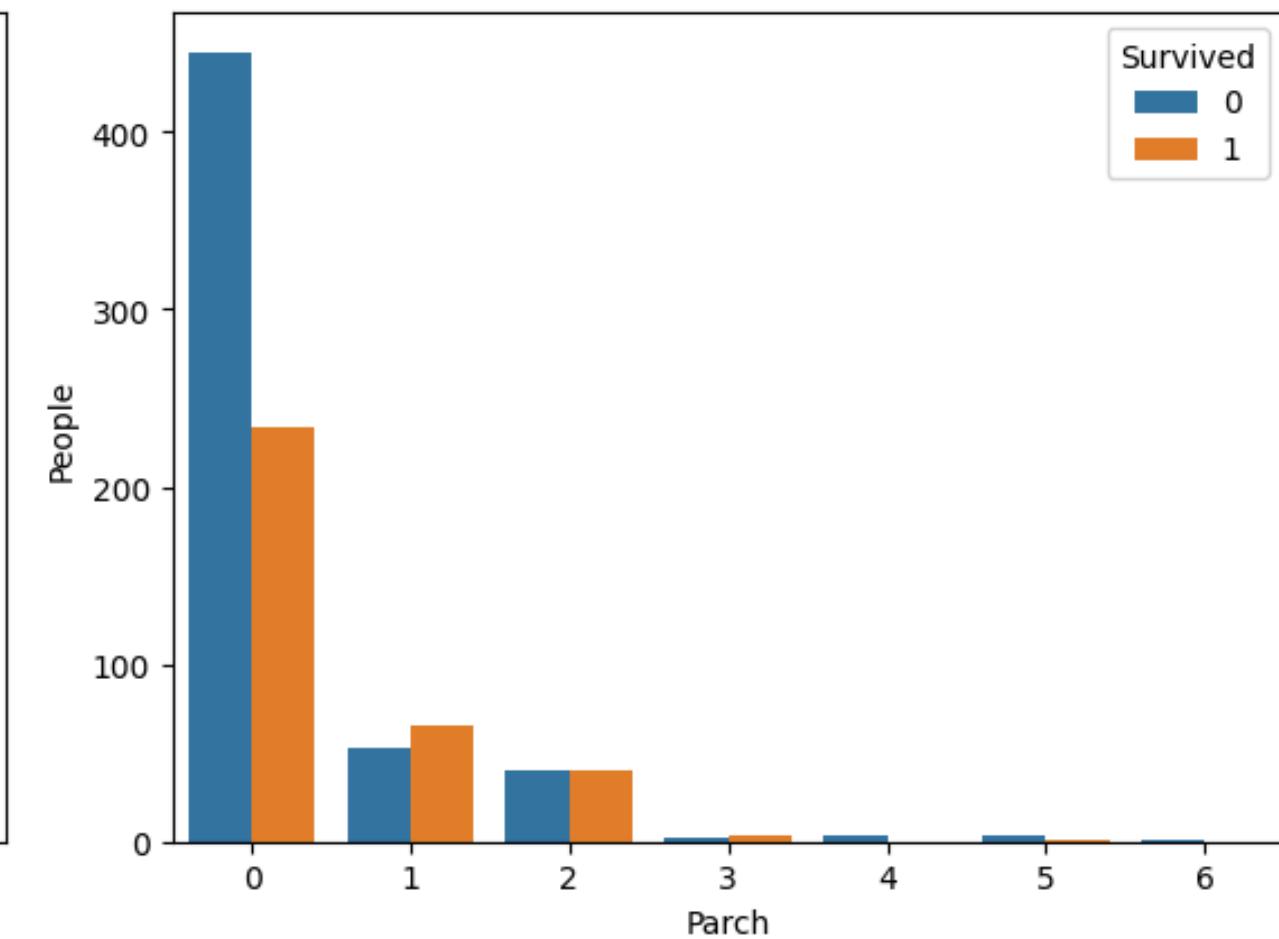


- We created “companians” Parch + Sibsp

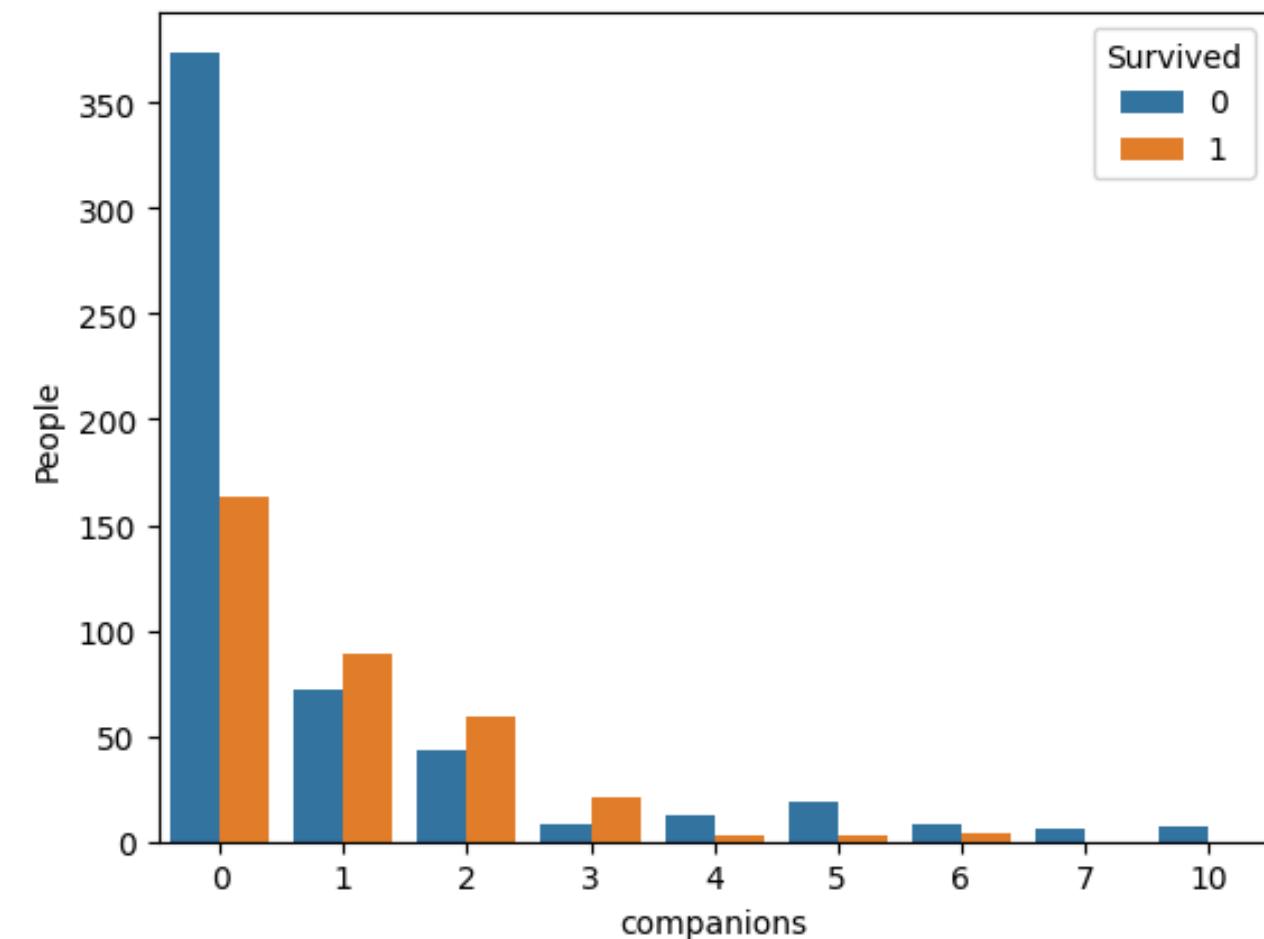
SibSp

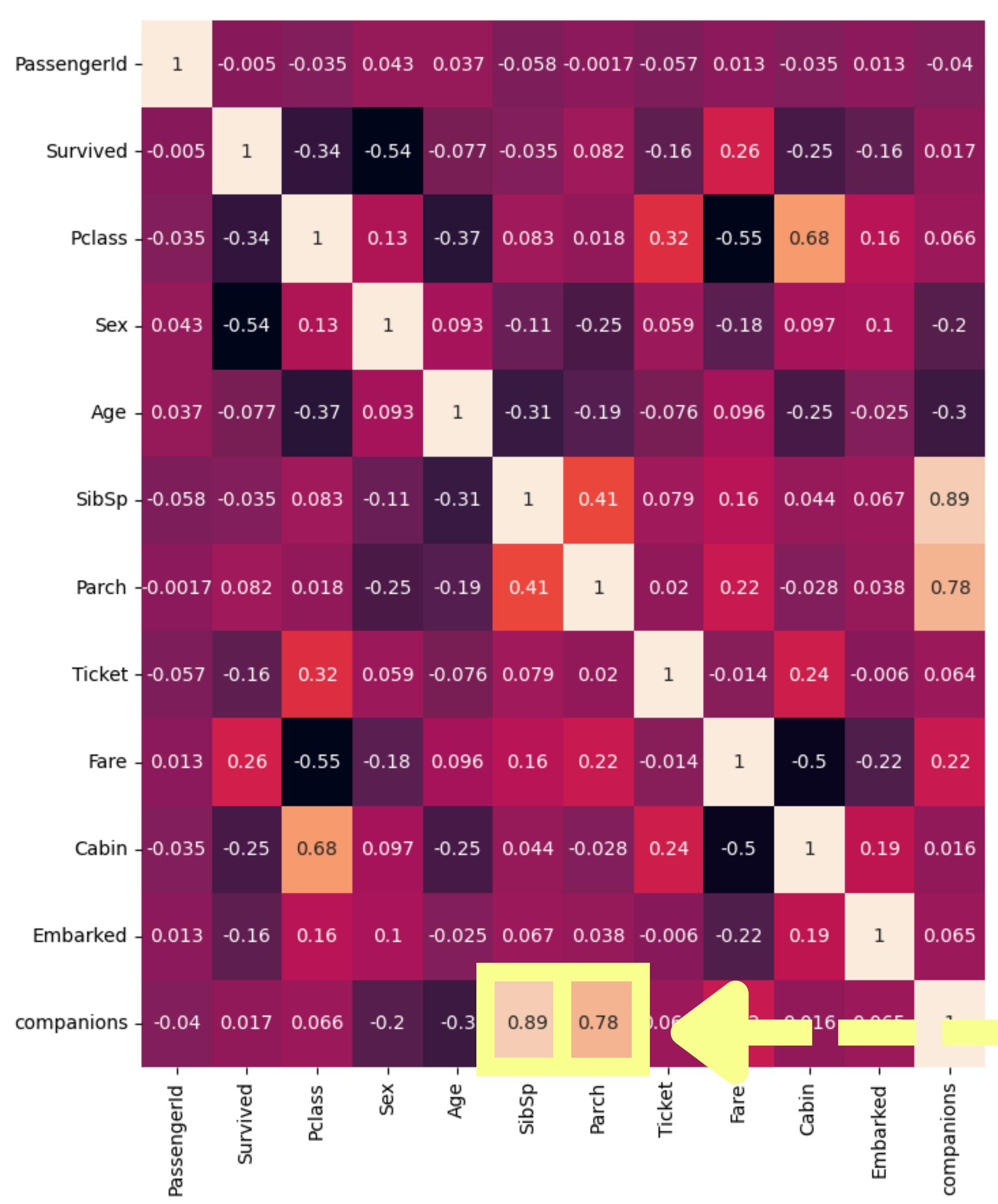


Parch



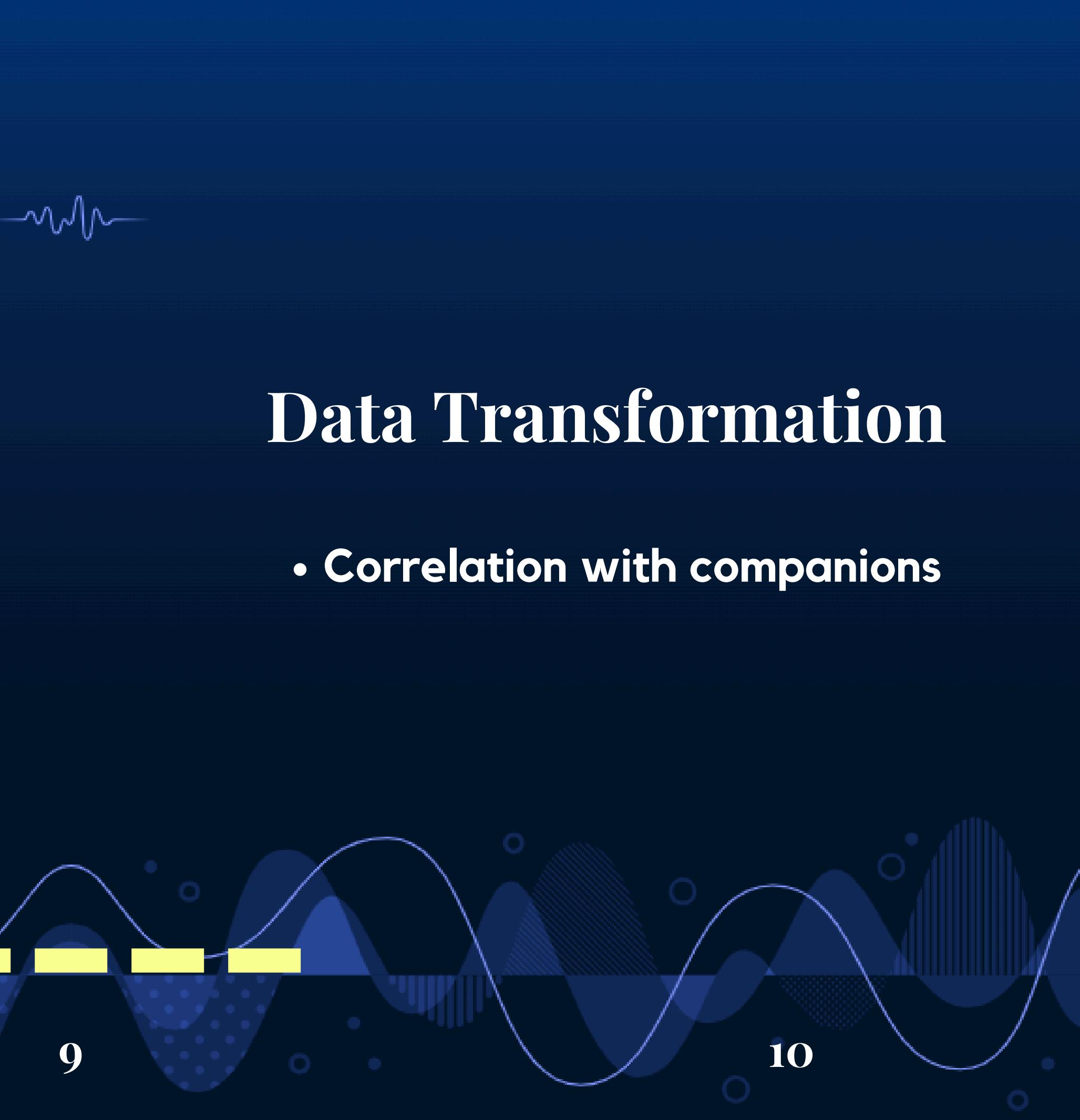
Companions



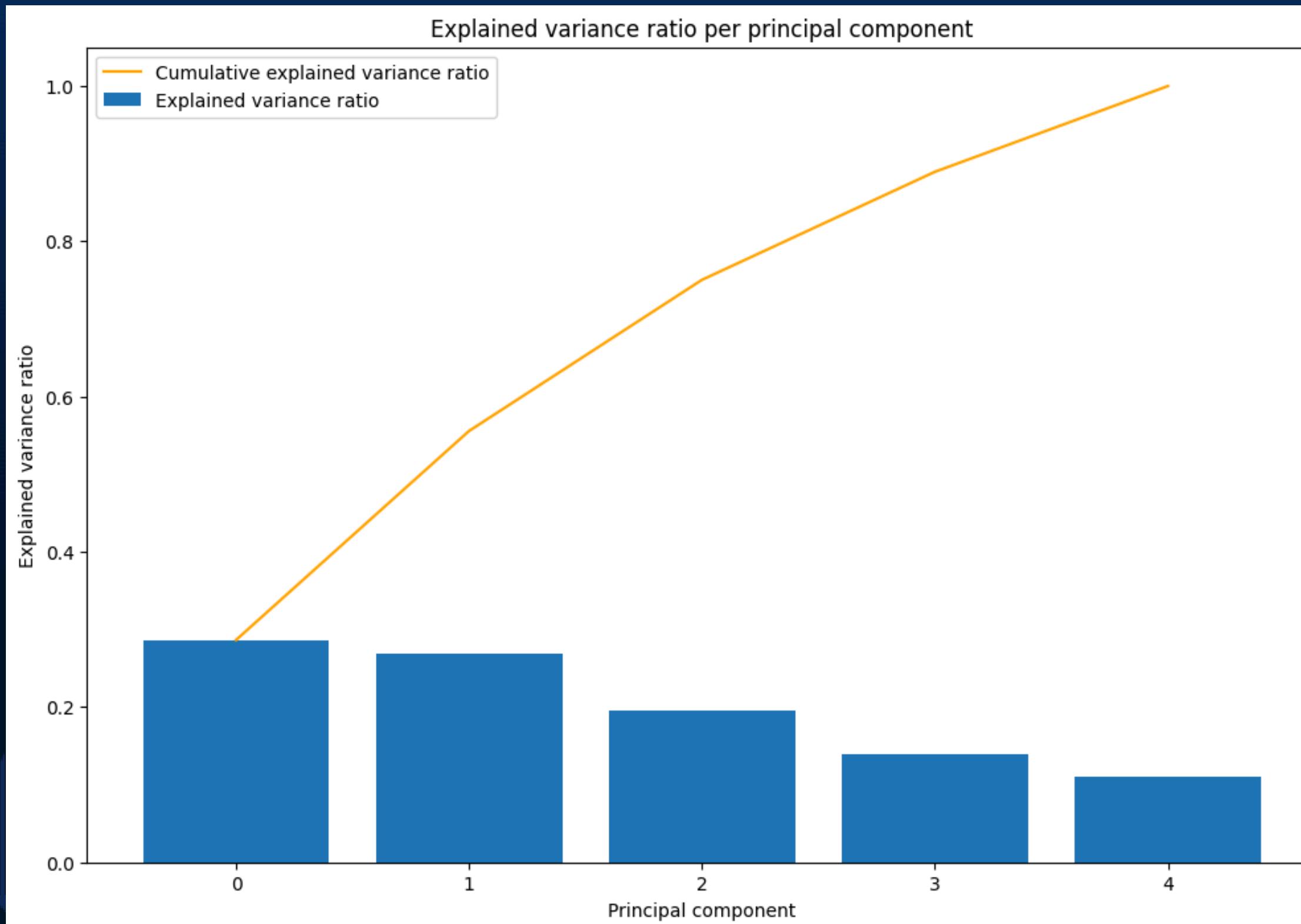


# Data Transformation

- Correlation with companions



# Principal Component Analysis



85% of variance

PC1

PC2

PC3

PC4

# Random Forest

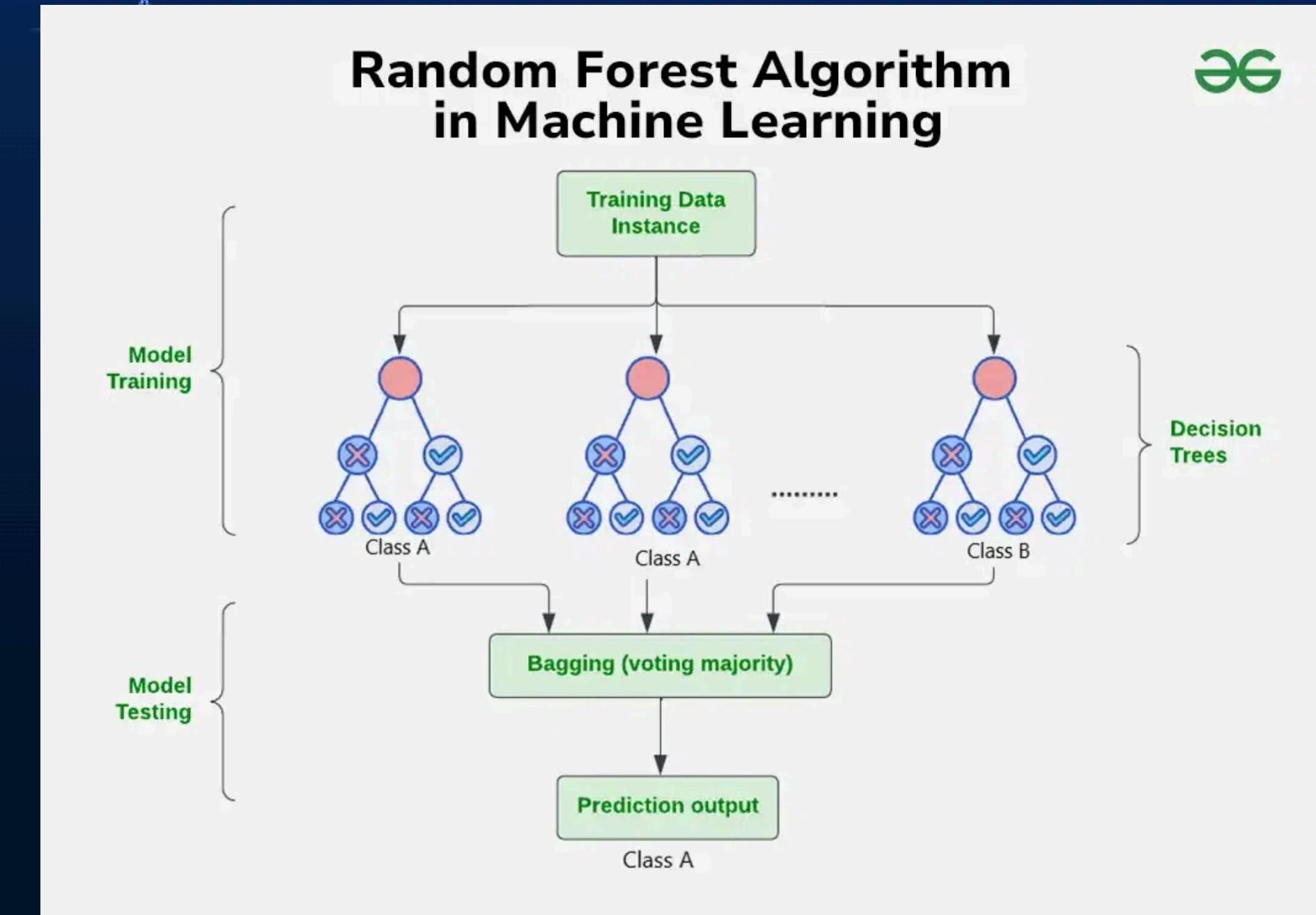


- Handles both categorical and numerical features well.
- Provides feature importance, which can help identify influential factors.
- Often achieves high accuracy
- Strengths:
  - Handles missing data and non-linear relationships well.
  - Robust to overfitting with enough trees.
- Weaknesses:
  - Harder to interpret.
  - May require tuning (e.g., number of trees, depth).

# Random Forest

Hyperparameters used

- max\_depth: 5 <
- max\_features :log2 -
- min\_samples\_leaf: 2 >
- min\_samples\_split: 10 >
- n\_estimators: 100 >



# Logistic regression



- Nature of the Problem
- Inherent binary classification algorithm
- Interpretable.
- Linear Relationship
- Provides probability estimates.
- Weaknesses:
  - May underperform with non-linear relationships.
  - Can be sensitive to outliers and multicollinearity.
- Regularization: Elasticnet (L1 and L2)

# Support Vector Machine (SVM)



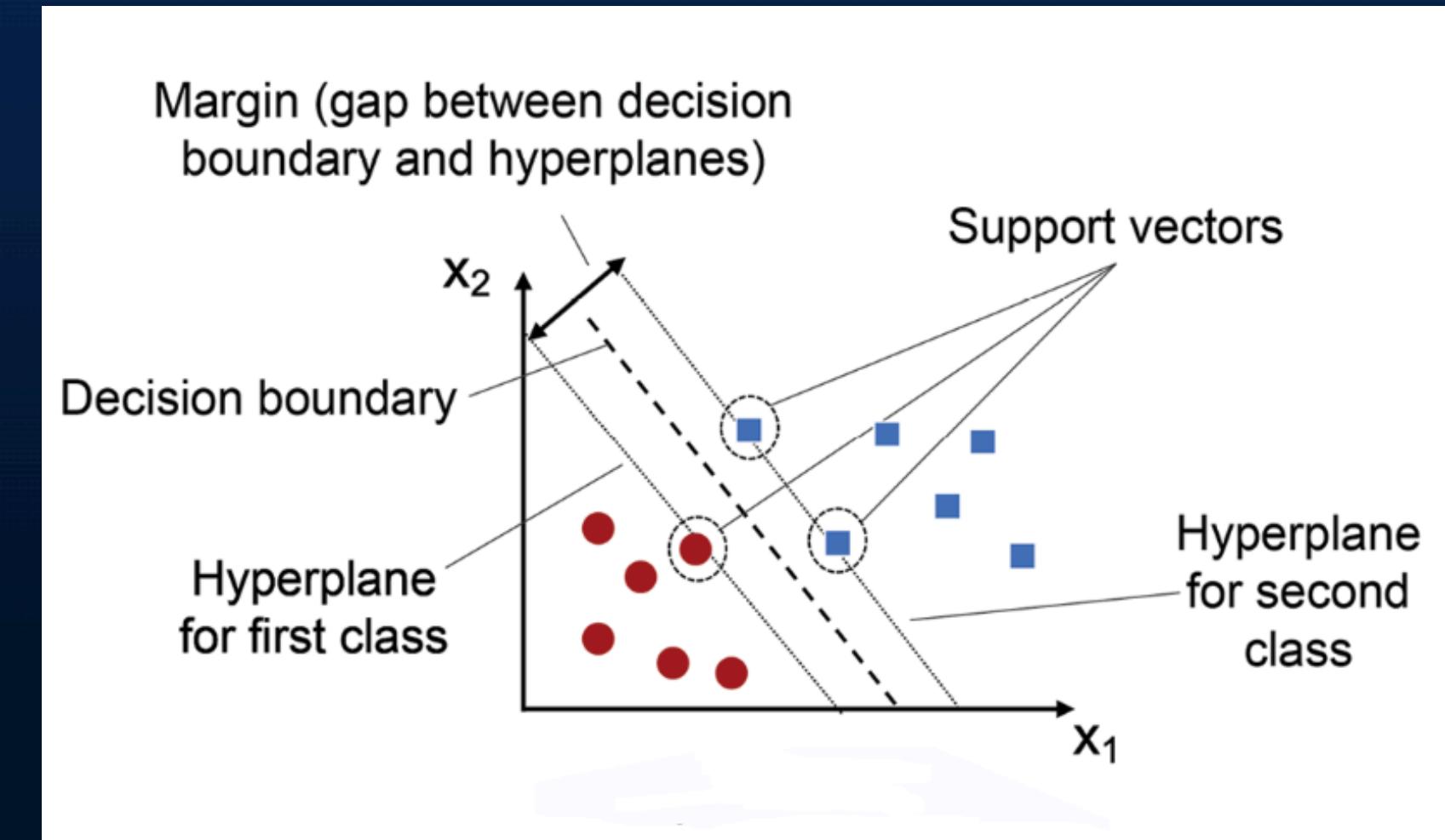
- Effective for many dimensions.
- Can handle non-linear relationships.
- Requires hyperparameter tuning.

Hyperparameters:

- Regularization (C): controls the trade-off between maximizing the margin between classes and minimizing the classification error on the training data.
- Kernel: defines the mathematical function used to transform the input data into a higher-dimensional space.
- Gamma: defines how far the influence of a single training example reaches.

# Support Vector Machine (SVM)

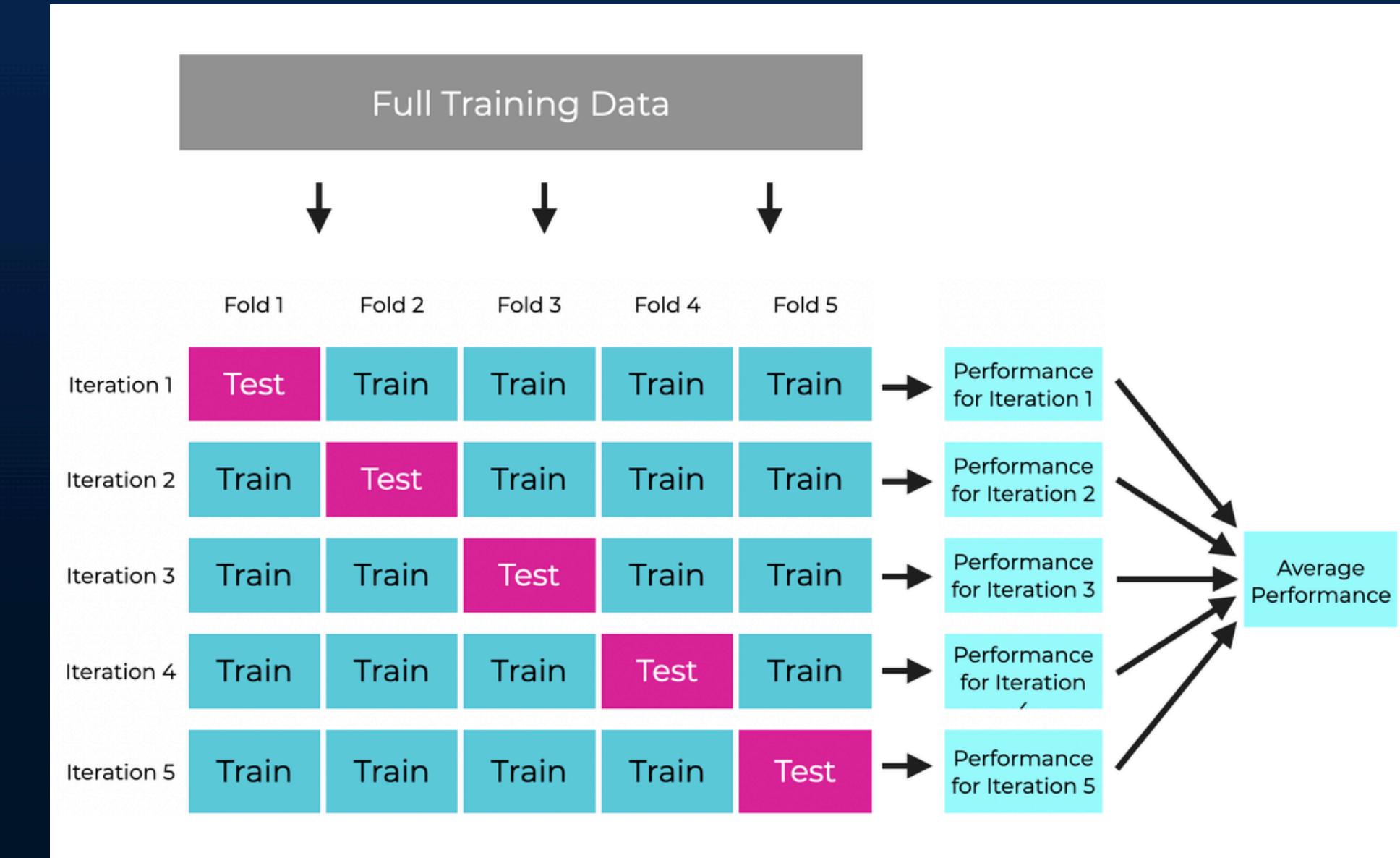
- 'C': 50: relatively high, because we prefer the correct classification of all points.
- 'gamma': 0.1: the influence of each training point on the decision boundary is moderate.
- 'kernel': 'rbf': the data was best separated using the Radial Basis Function (RBF) kernel.



<https://vitalflux.com/classification-model-svm-classifier-python-example/>

# K-cross fold validation

- Iteration 100 times
  - Time
- With k-cross
  - Avg Confusion Matrix
  - Avg Auc



<https://www.sharpsightlabs.com/blog/cross-validation-explained/>

# Comparison



## Random Forest

Accuracy: 0.8170

Precision: 0.8199

Recall: 0.8170

F1-Score: 0.8118

## Logistic Regression

Accuracy: 0.7789

Precision: 0.7772

Recall: 0.7789

F1-Score: 0.7776

## Support Vector Machine

Accuracy: 0.8193

Precision: 0.8013

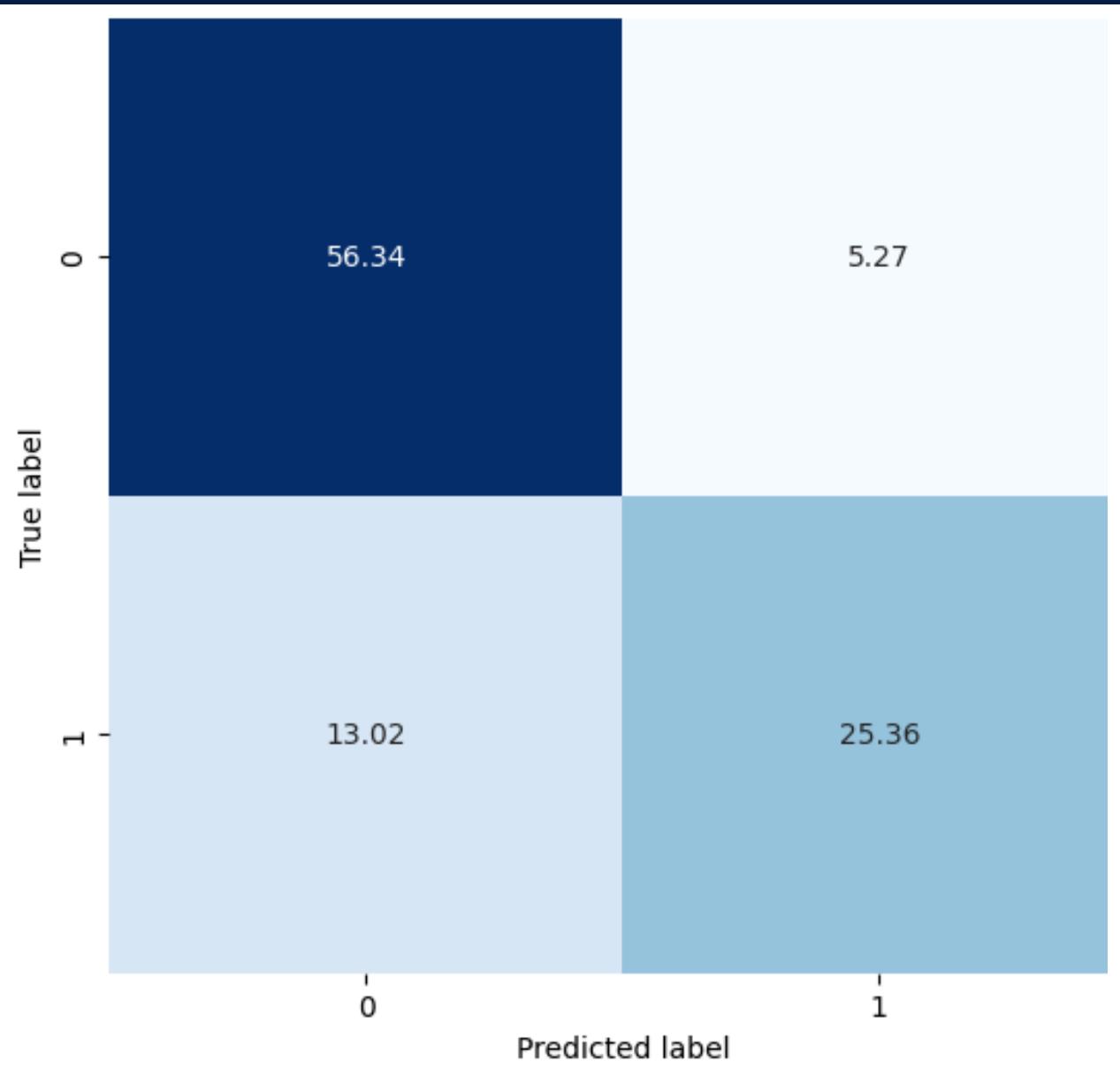
Recall: 0.7991

F1-Score: 0.7947

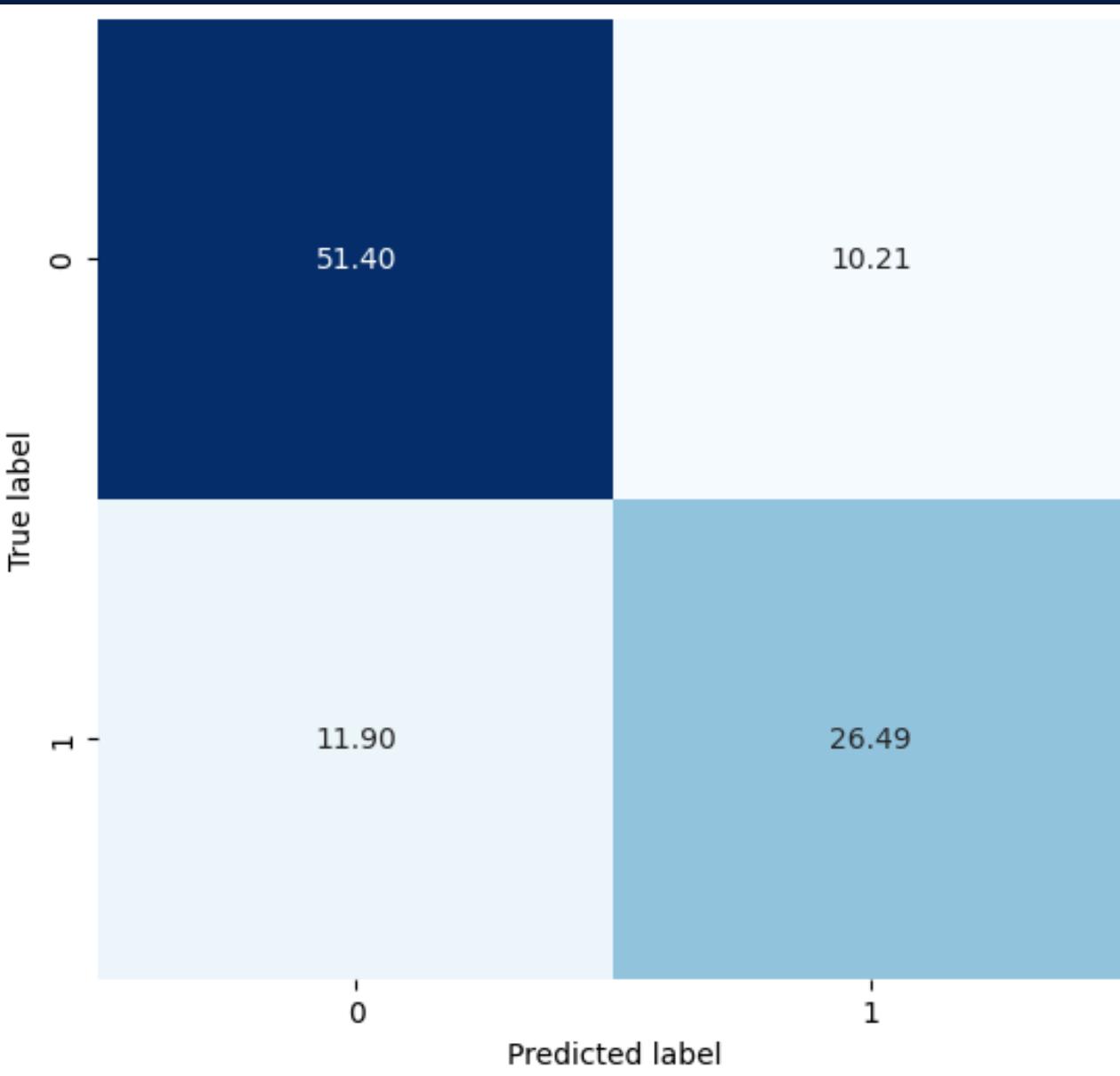
# Confusion matrix comparison



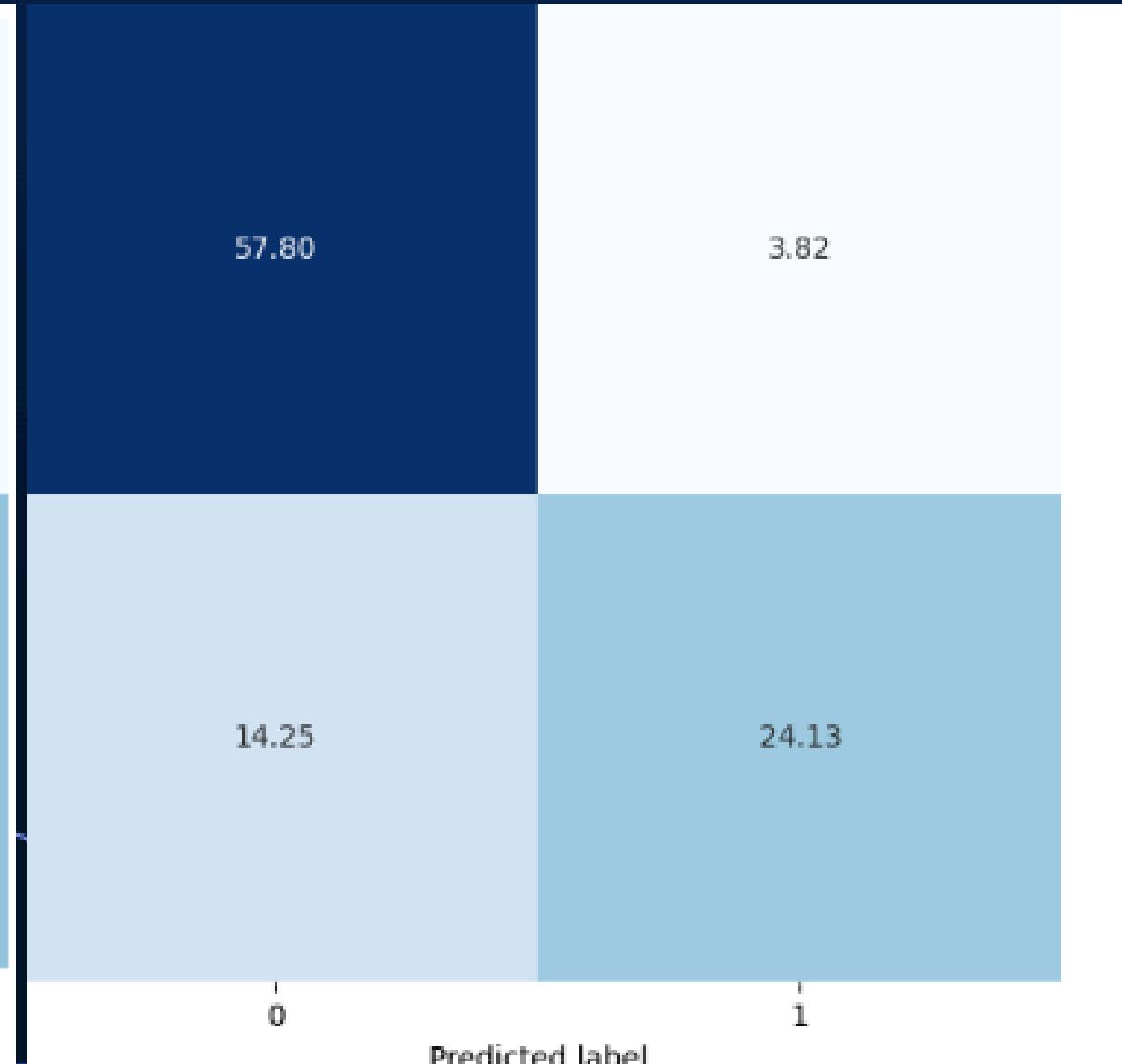
Random Forest



Logistic Regression



Support Vector Machine



# Discussion

Through our analysis, we identified that the most significant variables contributing to the model's performance are five key variables.

We applied PCA using four components to reduce dimensionality and improve model efficiency.

Despite the improvements, the model's accuracy did not reach the 85% threshold.

# Best Model

SVM performed better because various hyperparameters were adjusted and different combinations of these parameters were evaluated to determine which configuration yielded the best performance.

# Future Improvement



- One-Hot Encoding: apply one-hot encoding to categorical features to ensure they are correctly represented in the model.
- Model Diagnostics: conduct residual analysis and analyze ROC curves to understand model performance better.
- Apply transformation to variables to linearize not linear relations.
- Hyperparameter Tuning: perform grid search or randomized search for hyperparameter optimization to enhance model performance.
- Polynomial Kernels: experiment with polynomial kernels to capture interactions between features.
- Compare models with others (Naive-Bayes, knn, etc).

# References

- "A Practical Guide to Support Vector Classification" por Chih-Wei Hsu, Chih-Chung Chang y Chih-Jen Lin: Este artículo ampliamente citado proporciona una guía práctica para la clasificación con SVM, incluyendo una discusión detallada sobre la selección de gamma.
- "An Introduction to Support Vector Machines and Other Kernel-based Learning Methods" por Nello Cristianini y John Shawe-Taylor: Este libro ofrece una introducción completa a los SVM y otros métodos de aprendizaje basados en kernels, incluyendo una explicación detallada de diferentes tipos de kernels y cómo elegir el kernel adecuado para un problema específico.
- "Support Vector Networks" por Corinna Cortes y Vladimir Vapnik: Este artículo fundamental presenta el concepto de SVM y discute el papel del parámetro C en el control de la complejidad del modelo y la prevención del sobreajuste.
- Breiman, L. (2001). Random Forests. *Machine Learning*, 45, 5–32. © 2001 Kluwer Academic Publishers.
- Géron, A. (2019). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*. O'Reilly

Media.

Bishop, C. M. (2006). *Pattern recognition and machine learning*. Springer

# References

- Ekinci, E., İlhan Omurca, S., & Acun, N. (n.d.). A comparative study on machine learning techniques using Titanic dataset. Kocaeli University, Faculty of Engineering, Computer Engineering Department. Retrieved from [https://www.researchgate.net/profile/Neytullah-Acun/publication/324909545\\_A\\_Comparative\\_Study\\_on\\_Machine\\_Learning\\_Techniques\\_Using\\_Titanic\\_Dataset/links/607533bc299bf1f56d51db20/A-Comparative-Study-on-Machine-Learning-Techniques-Using-Titanic-Dataset.pdf](https://www.researchgate.net/profile/Neytullah-Acun/publication/324909545_A_Comparative_Study_on_Machine_Learning_Techniques_Using_Titanic_Dataset/links/607533bc299bf1f56d51db20/A-Comparative-Study-on-Machine-Learning-Techniques-Using-Titanic-Dataset.pdf)
- Titanic Machine Learning Study from Disaster. (n.d.). UDSpace. University of Delaware. Retrieved September 11, 2024, from <https://udspace.udel.edu/items/61cb7b16-2590-4cd9-9d5b-fa23284bf2a8>
- Khan, M., & Qamar, U. (2018). Predicting survival on Titanic by applying exploratory data analytics and machine learning techniques. International Journal of Computer Applications, 179(44), 18-25.  
<https://doi.org/10.5120/ijca2018916972>
- Rigatti, S. J. (2017). Random forest. Journal of Insurance Medicine, 47(1), 31–39.  
<https://doi.org/10.17849/insm-47-01-31-39.1>

# Thank you

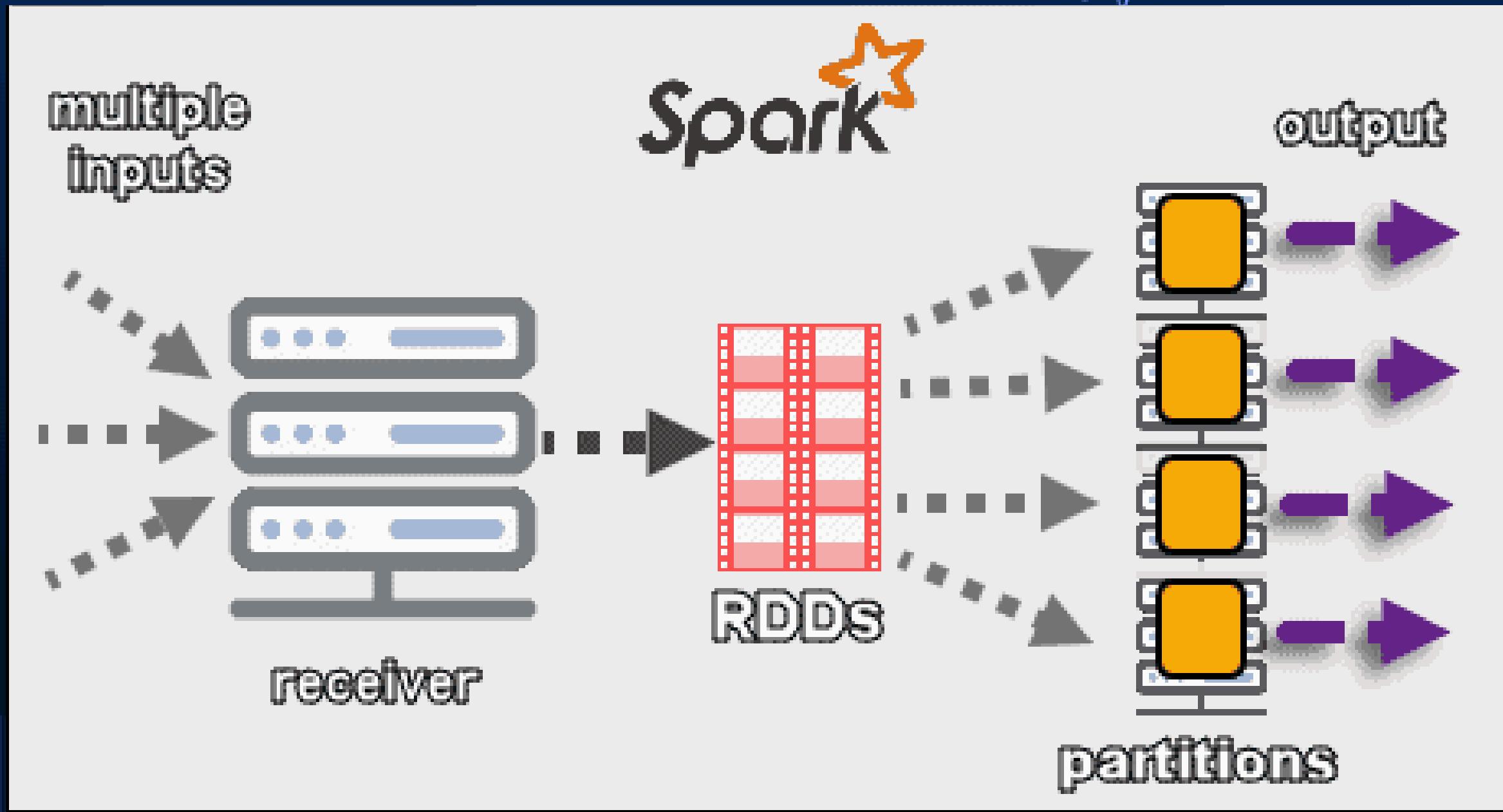


Video: [https://youtu.be/aYYxbzGhE\\_Q](https://youtu.be/aYYxbzGhE_Q)

Code:

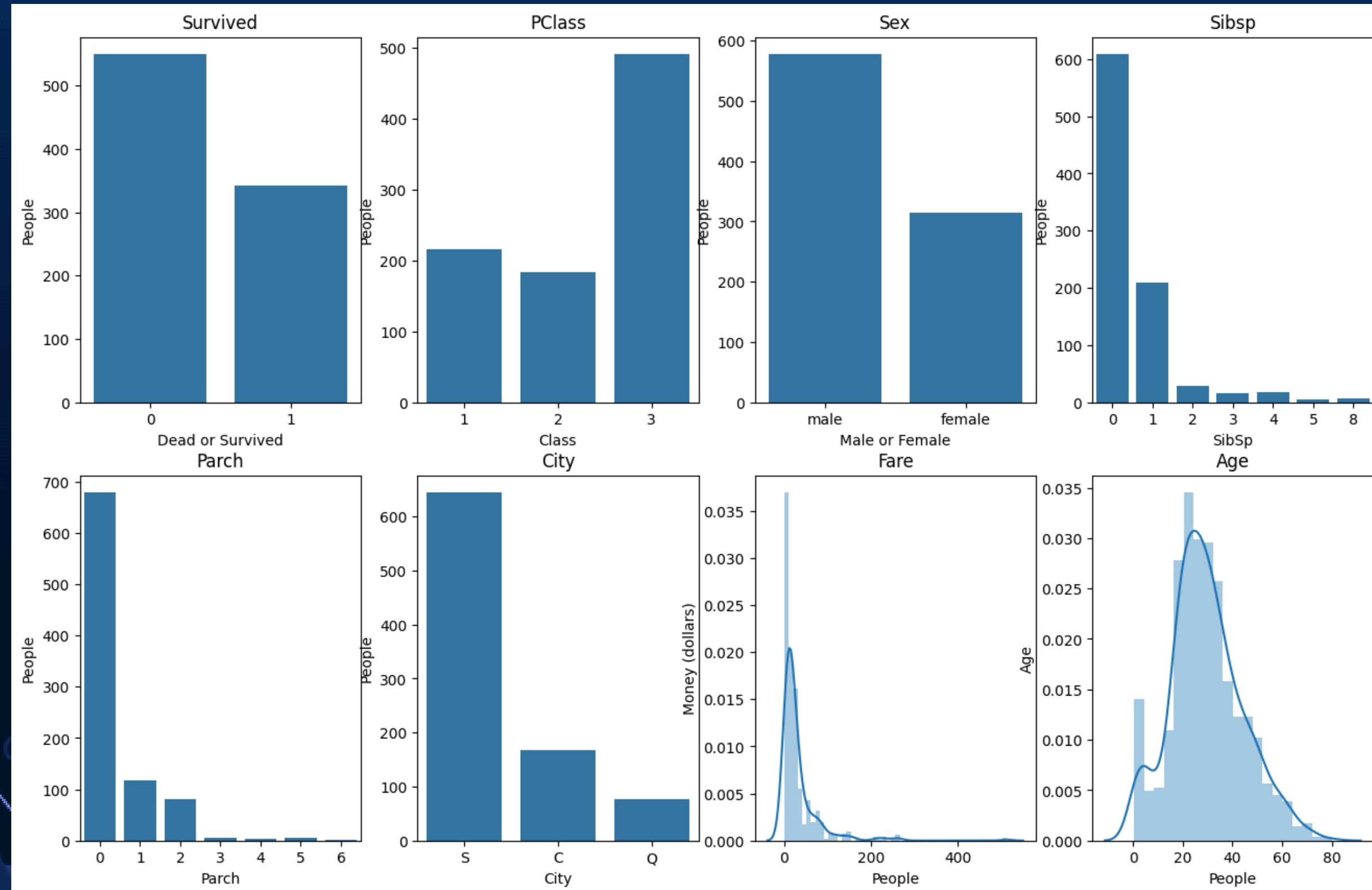
Big Data Aproach:

# Big Data Approach

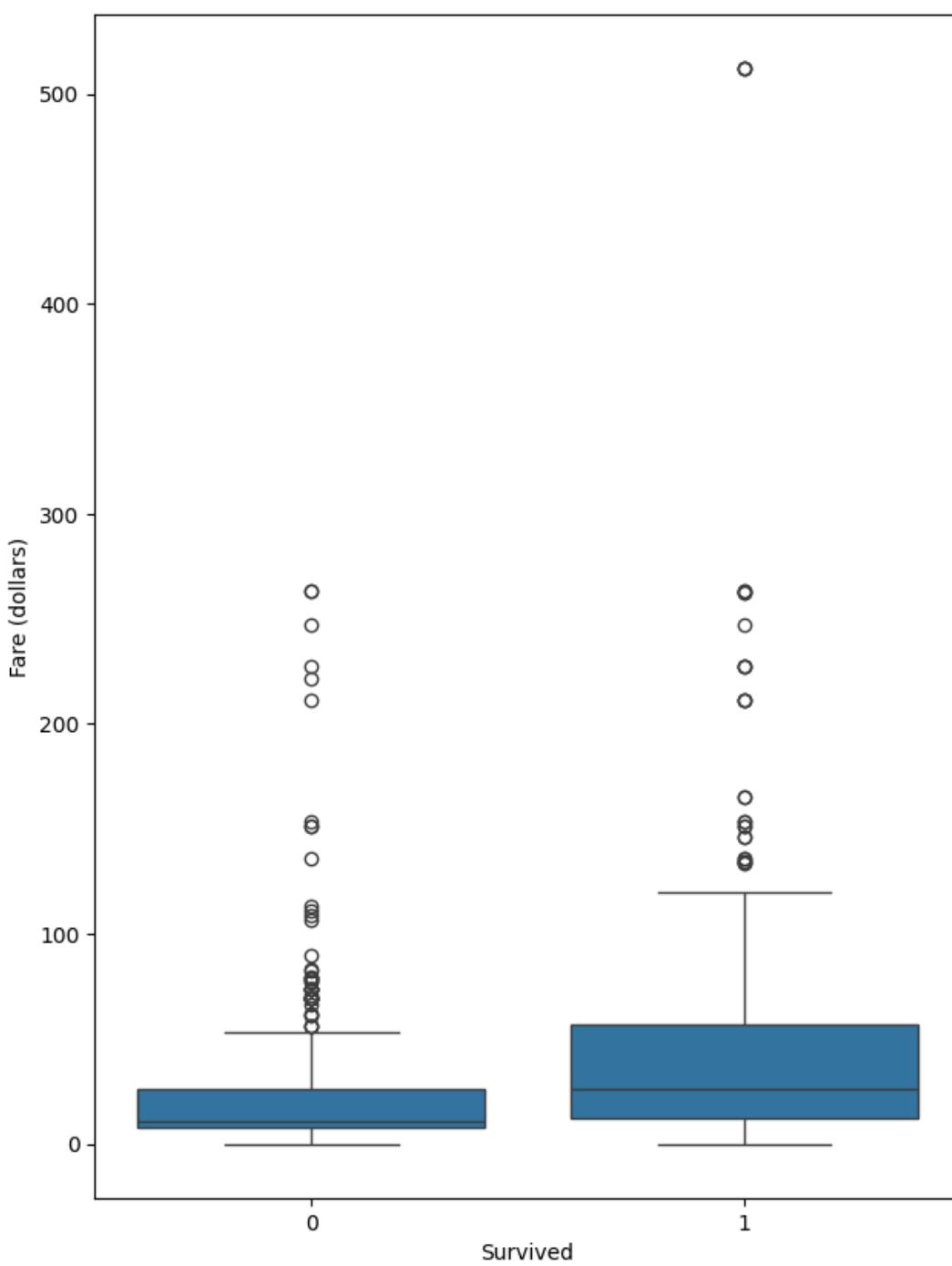
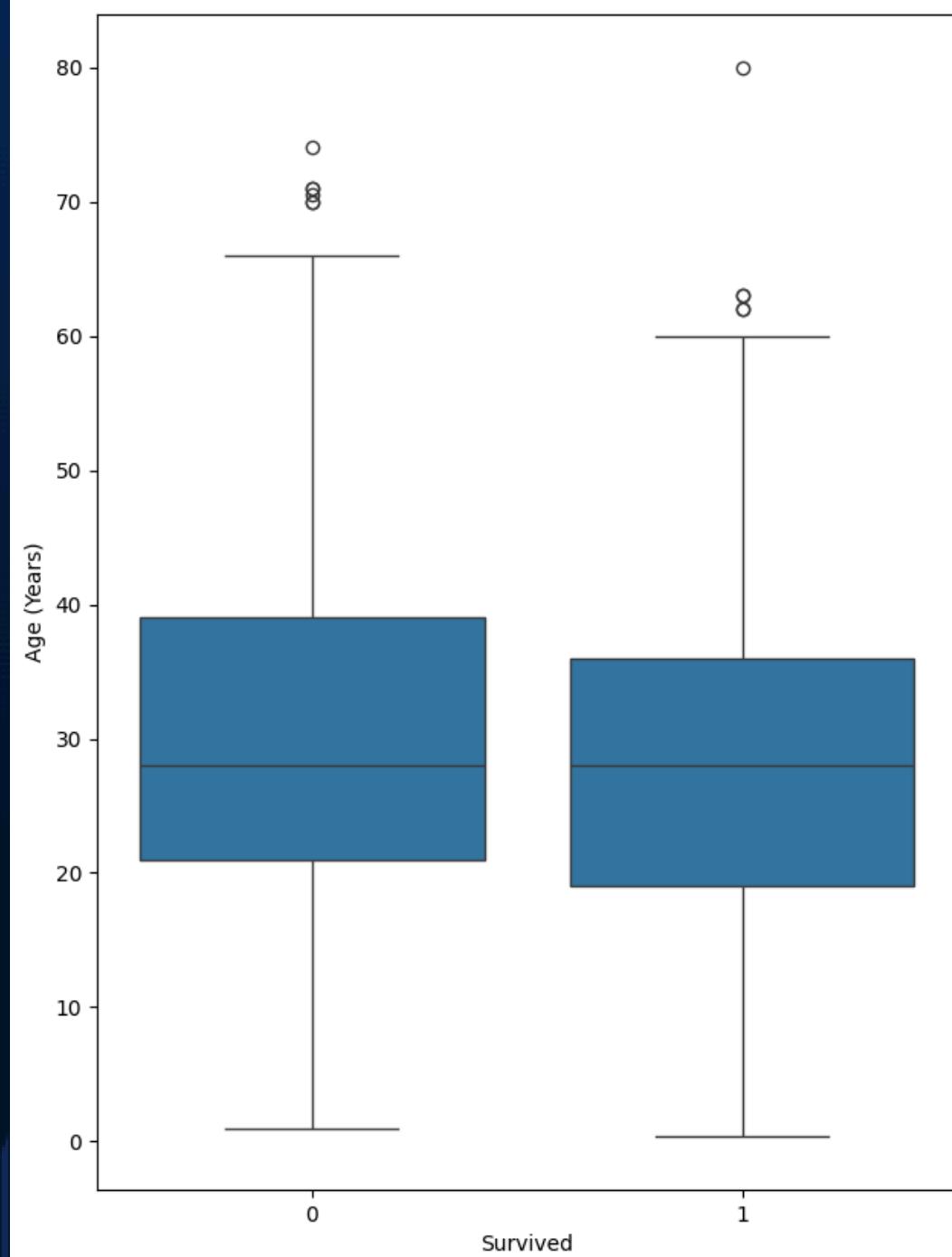


PySpark

# Data Distribution



# Data Distribution Age & Fare



We can conclude from here:  
Just a few outliers from older  
people.

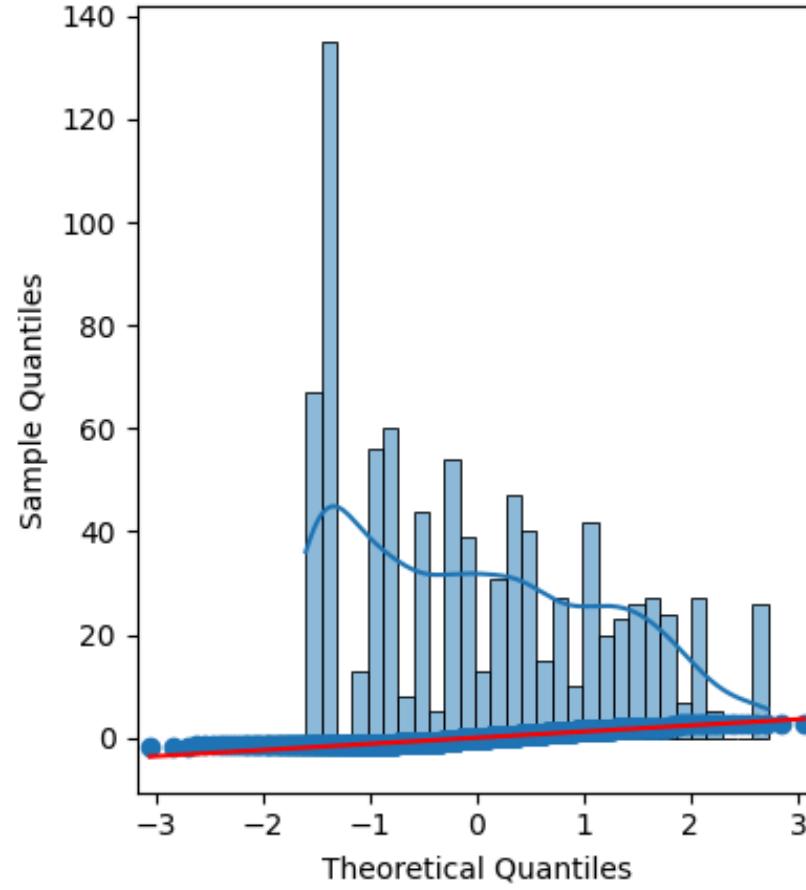
Cause of the IQR is similar,  
that told us that the age is not  
a determinant data for the  
survivors.

Passengers for first class, had  
more survivance probability

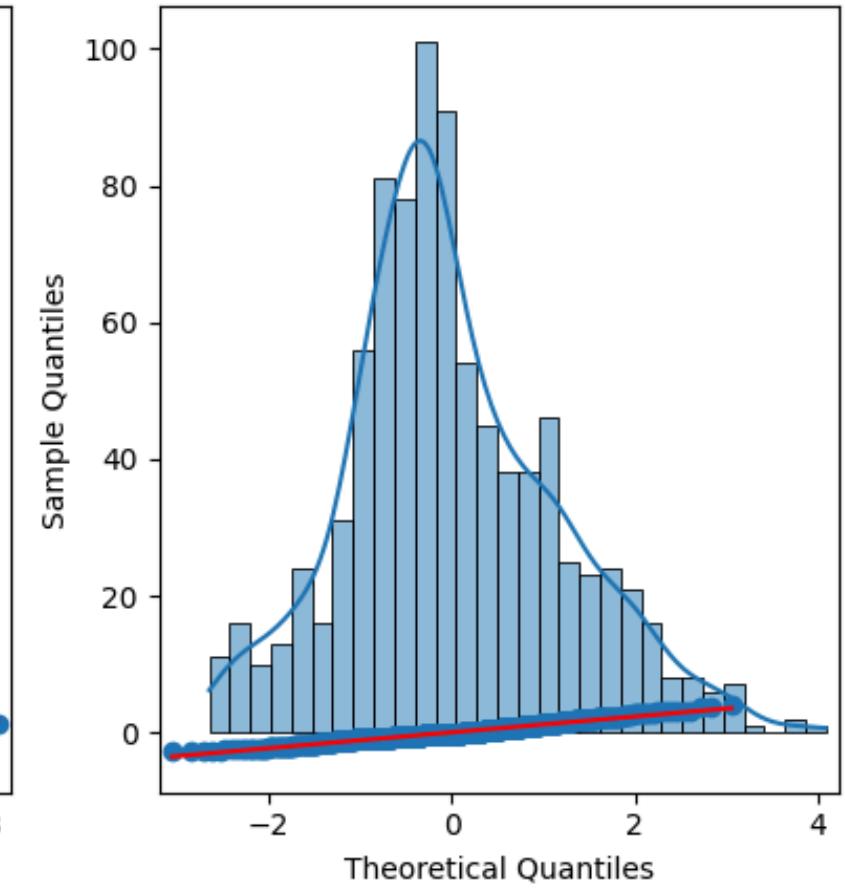
Younger people survived  
more.

### Histograms and QQ Plots for Numeric Features

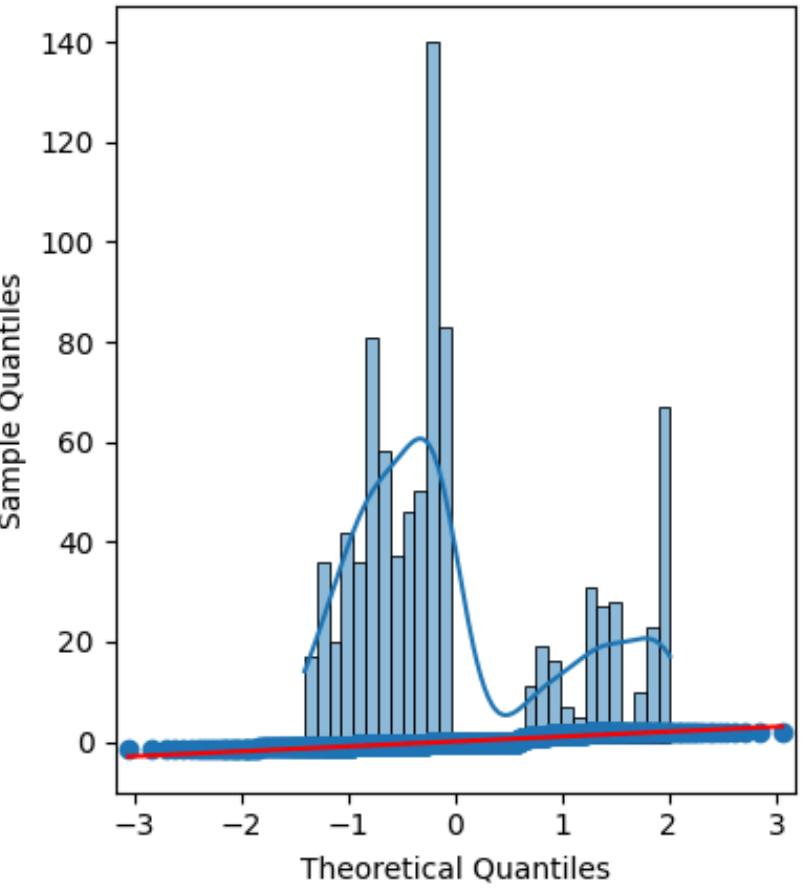
Histogram & KDE  
PC1



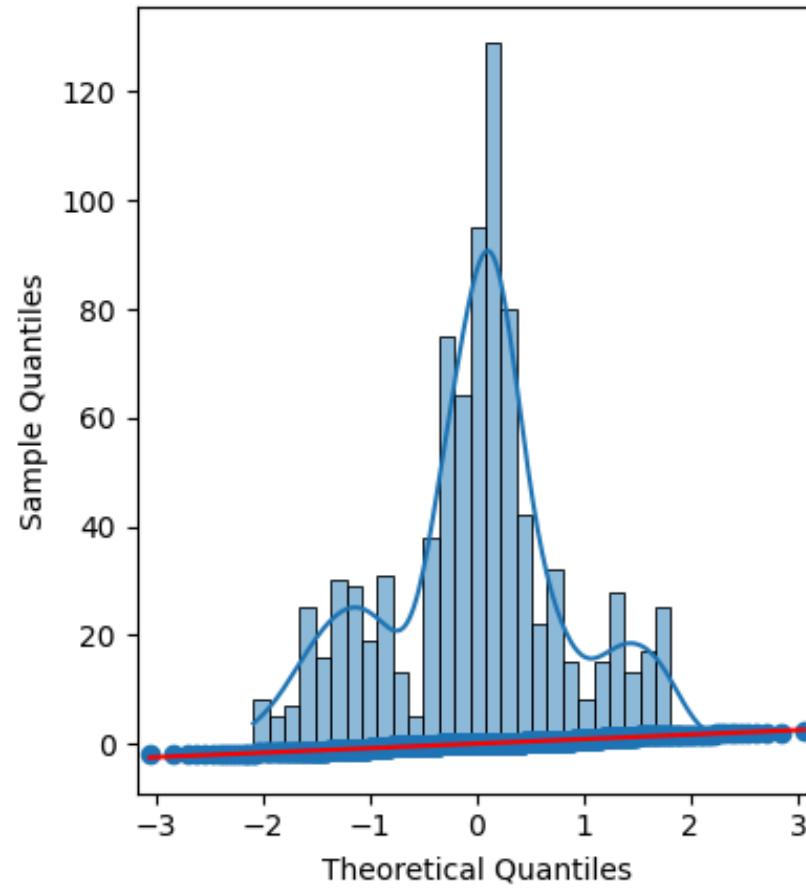
Histogram & KDE  
PC2



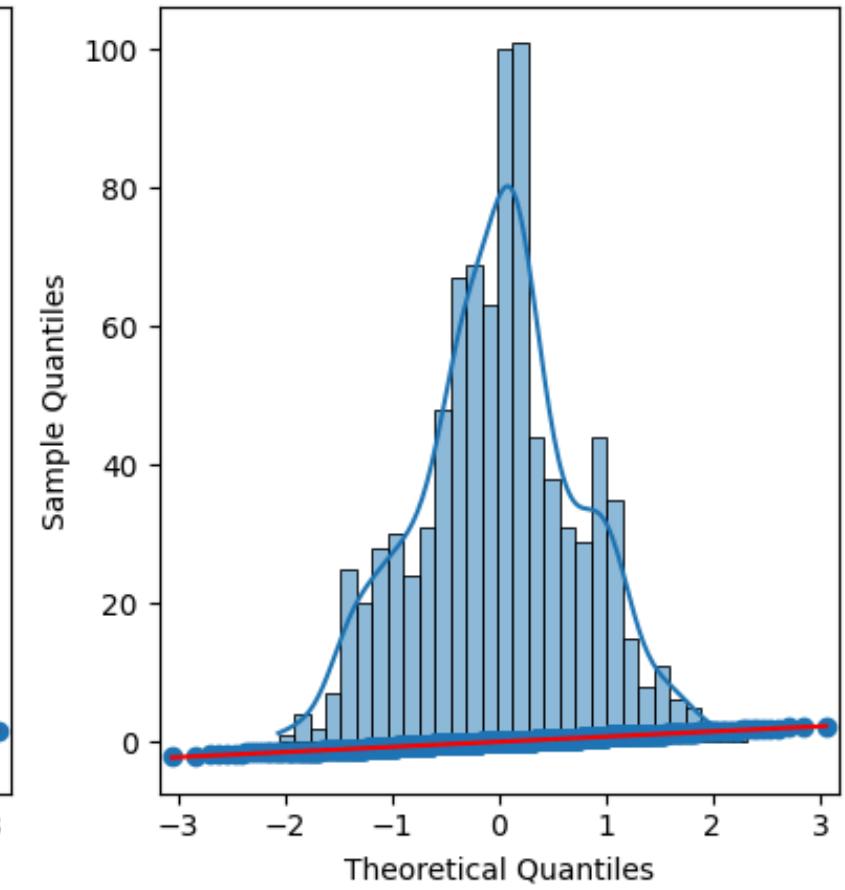
Histogram & KDE  
PC3



Histogram & KDE  
PC4



Histogram & KDE  
PC5

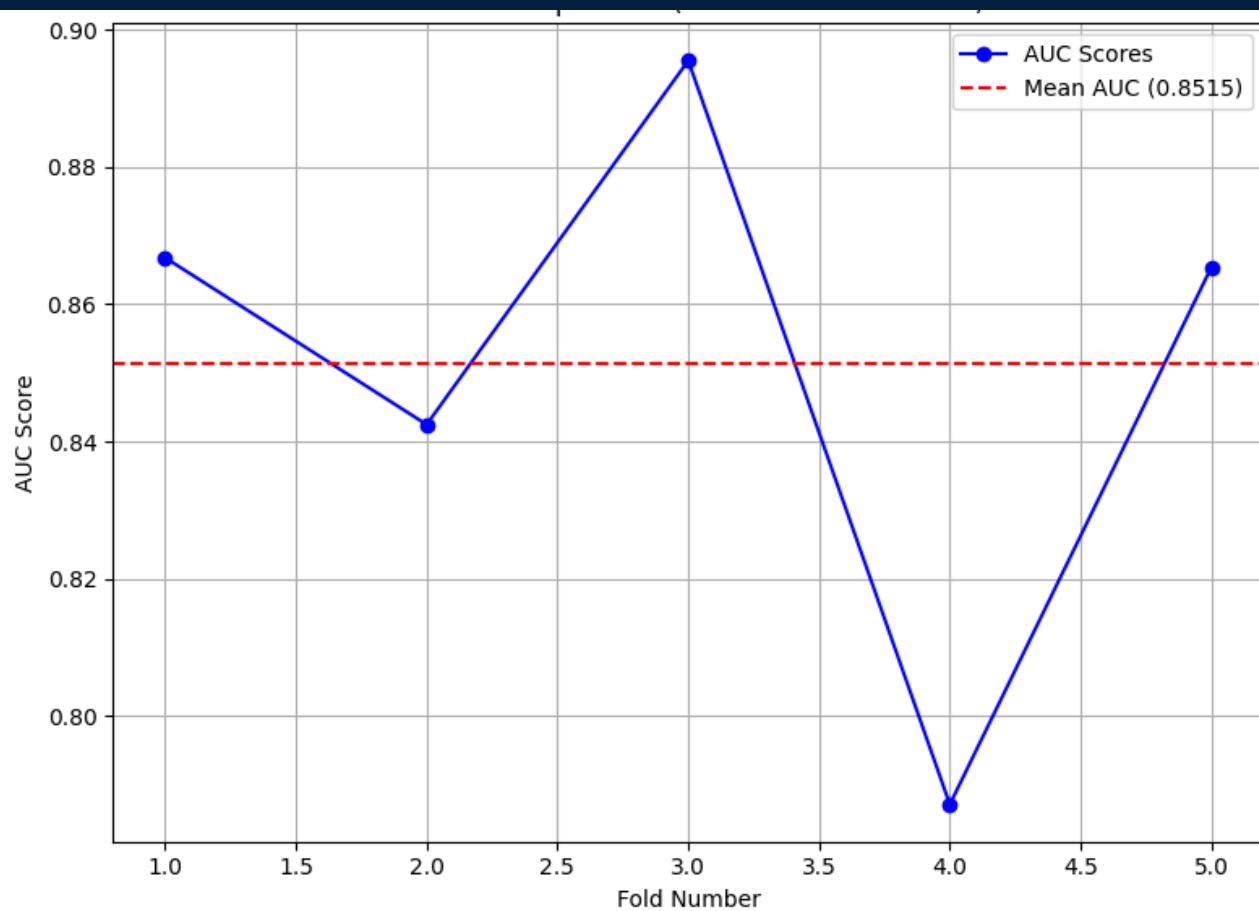


# Histograms and QQ Plots for Numeric Features

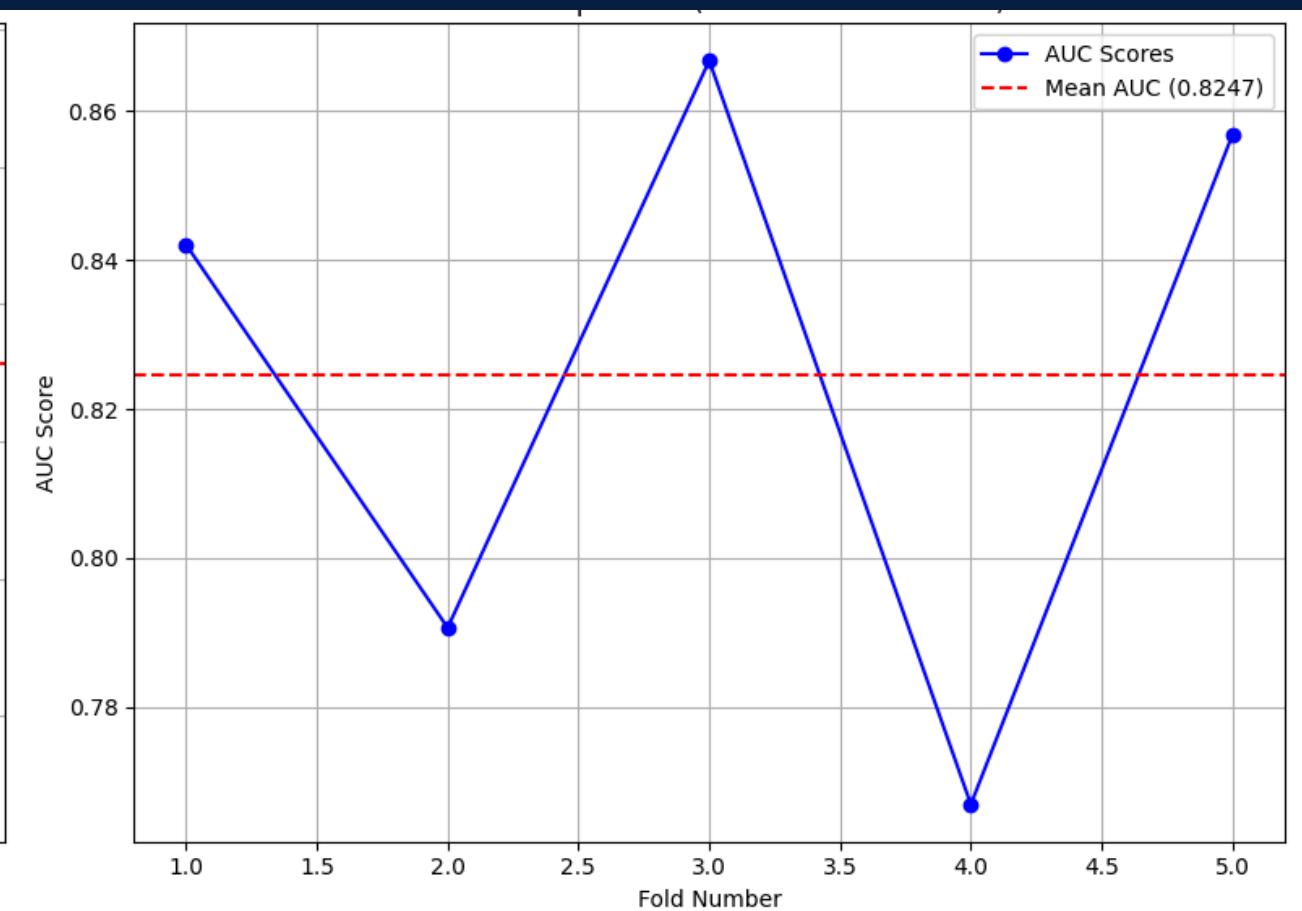
# Comparison AUC Score Per Fold



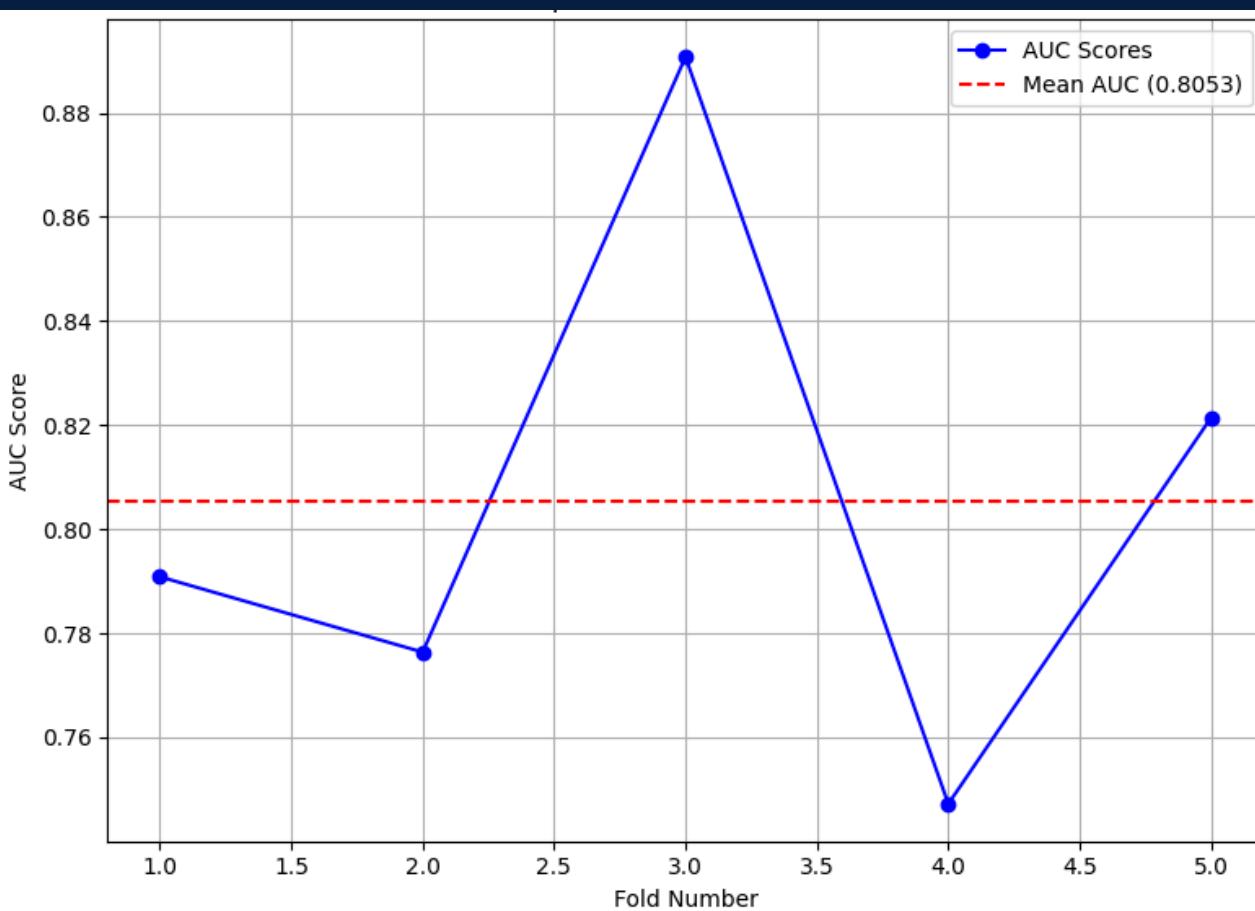
## Random Forest



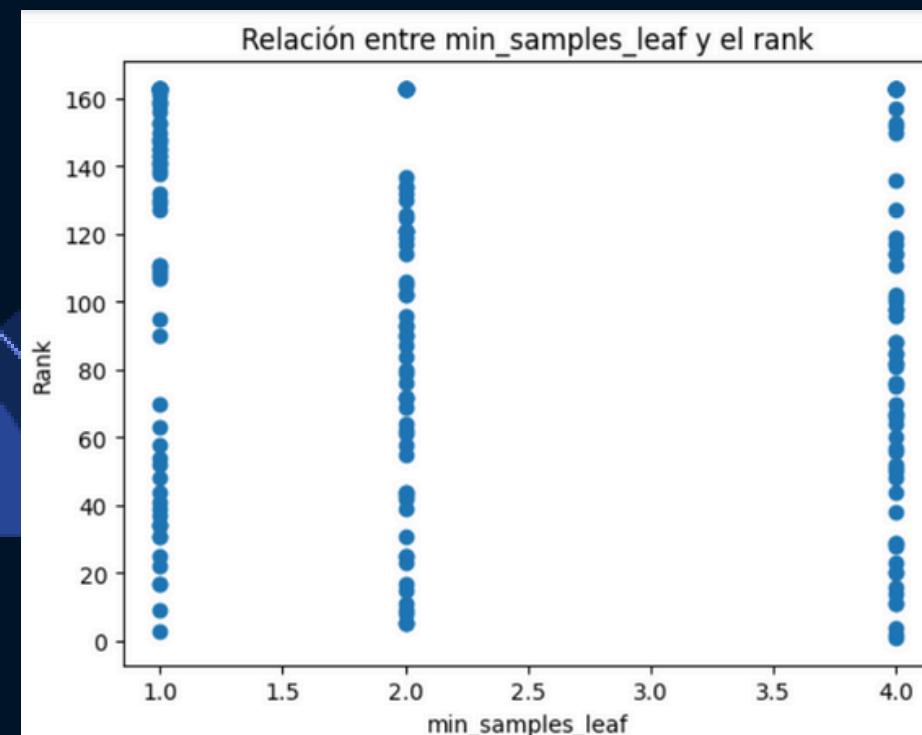
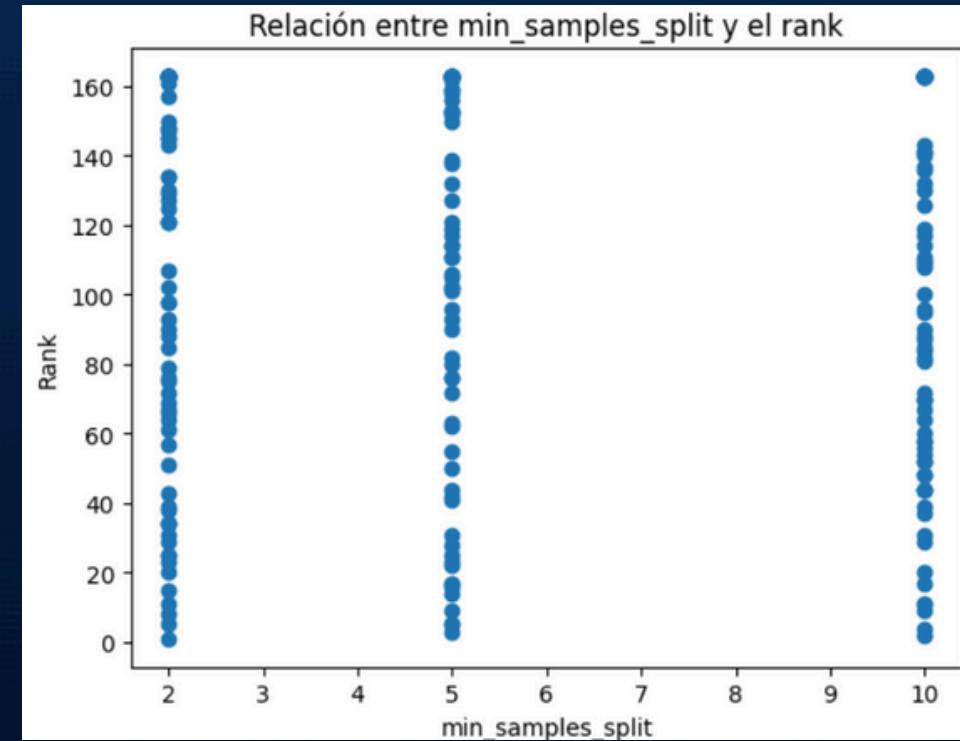
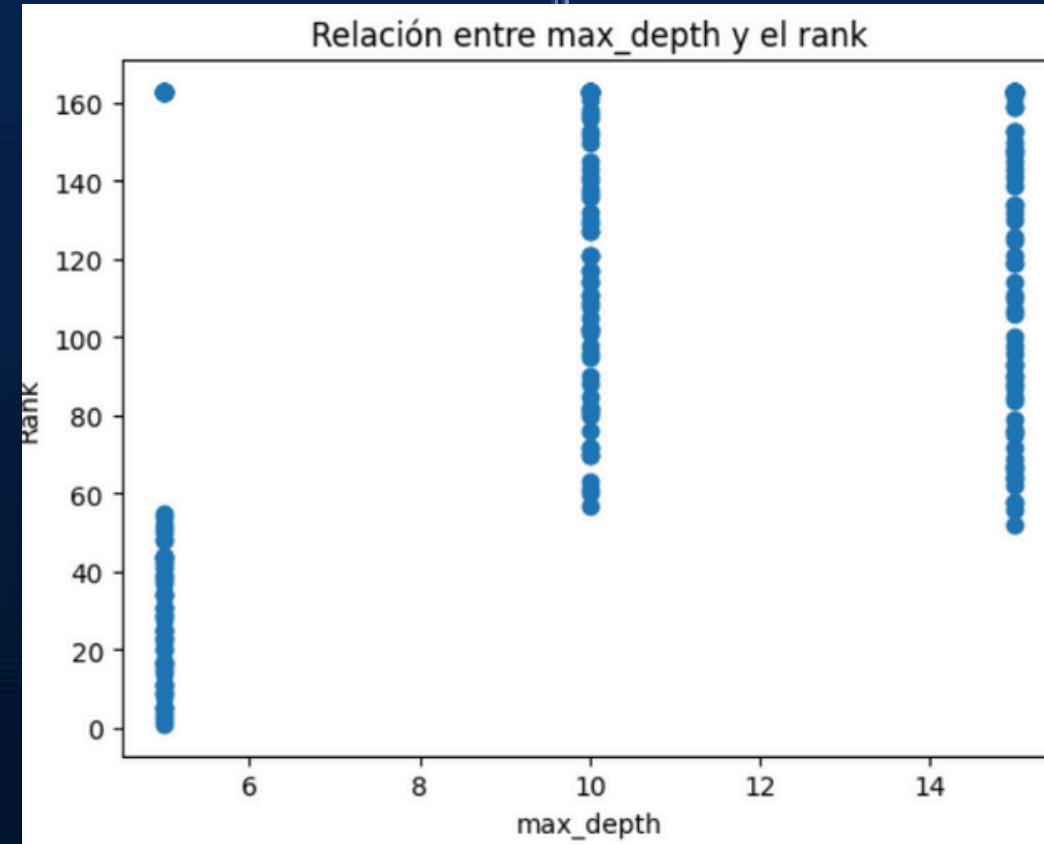
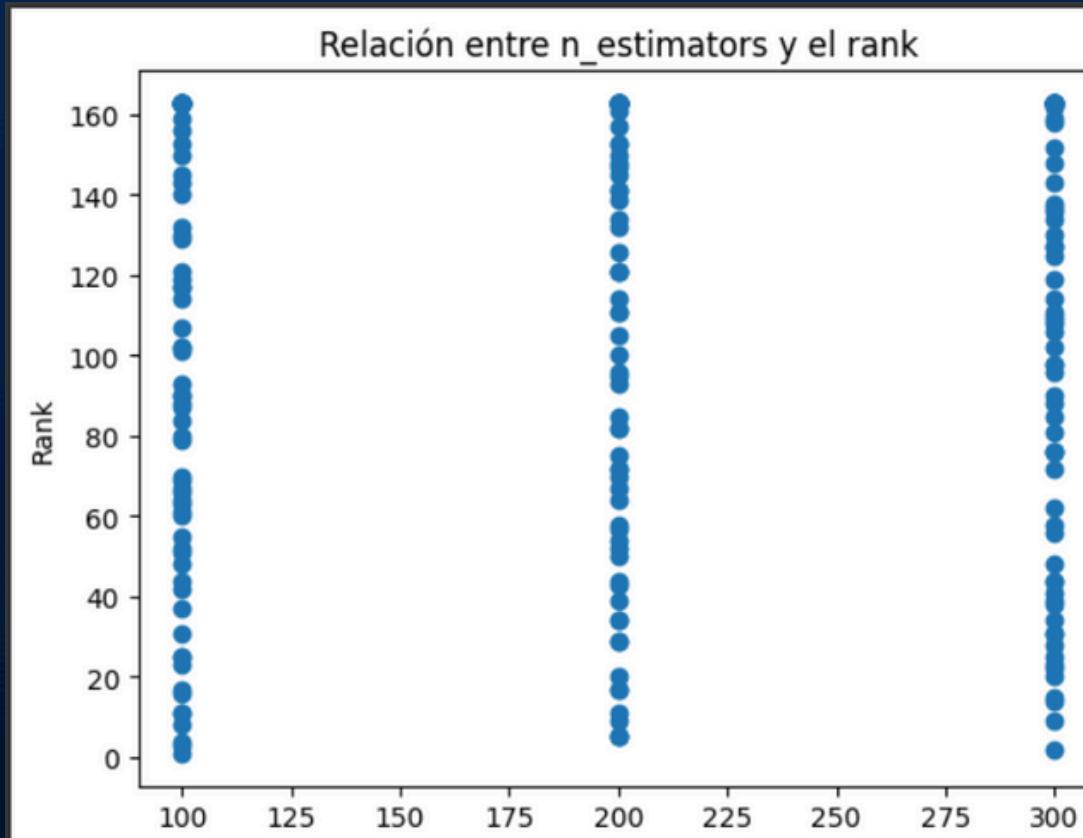
## Logistic Regression



## Support Vector Machine



# Anexos

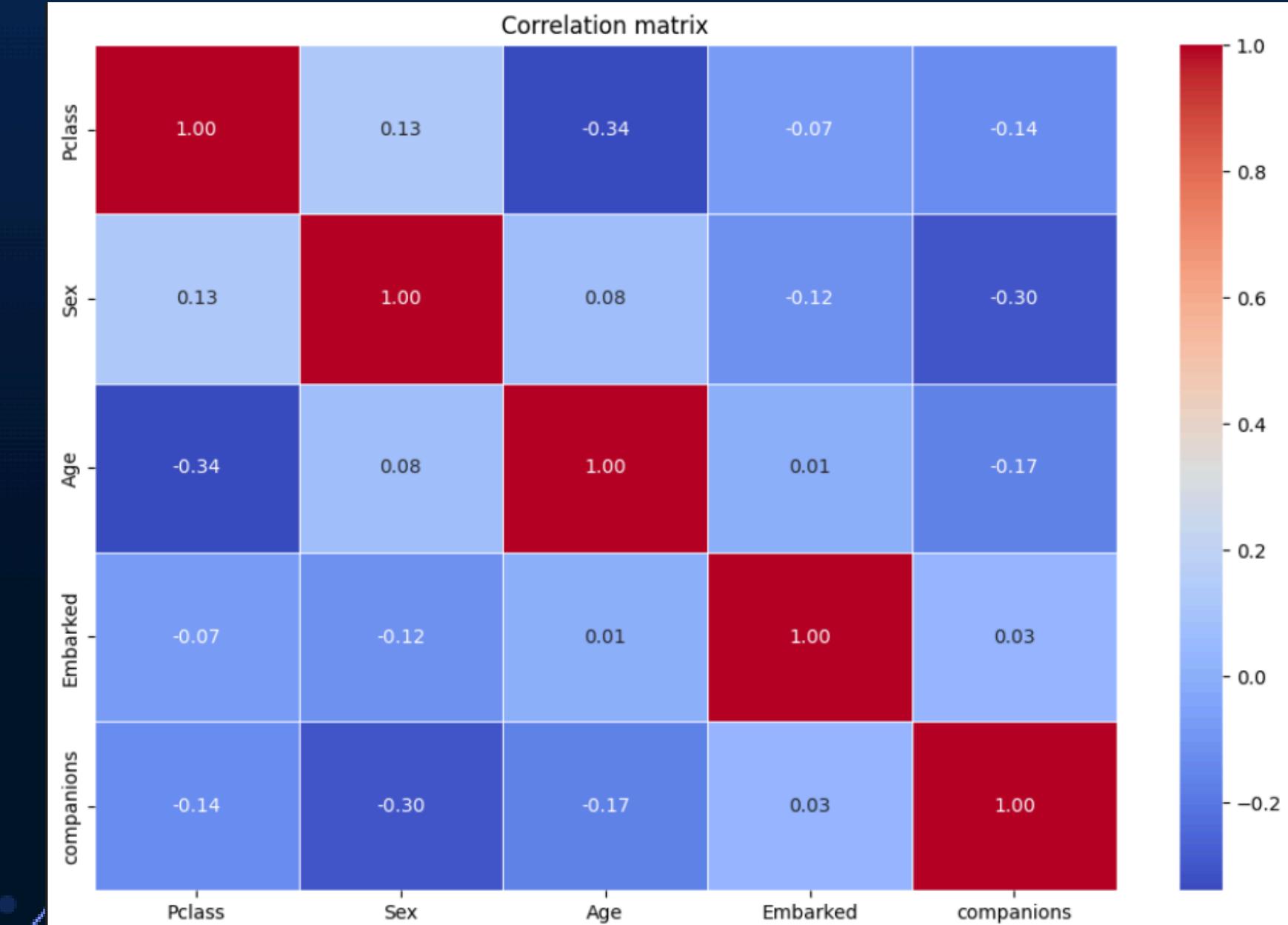


# Anexos

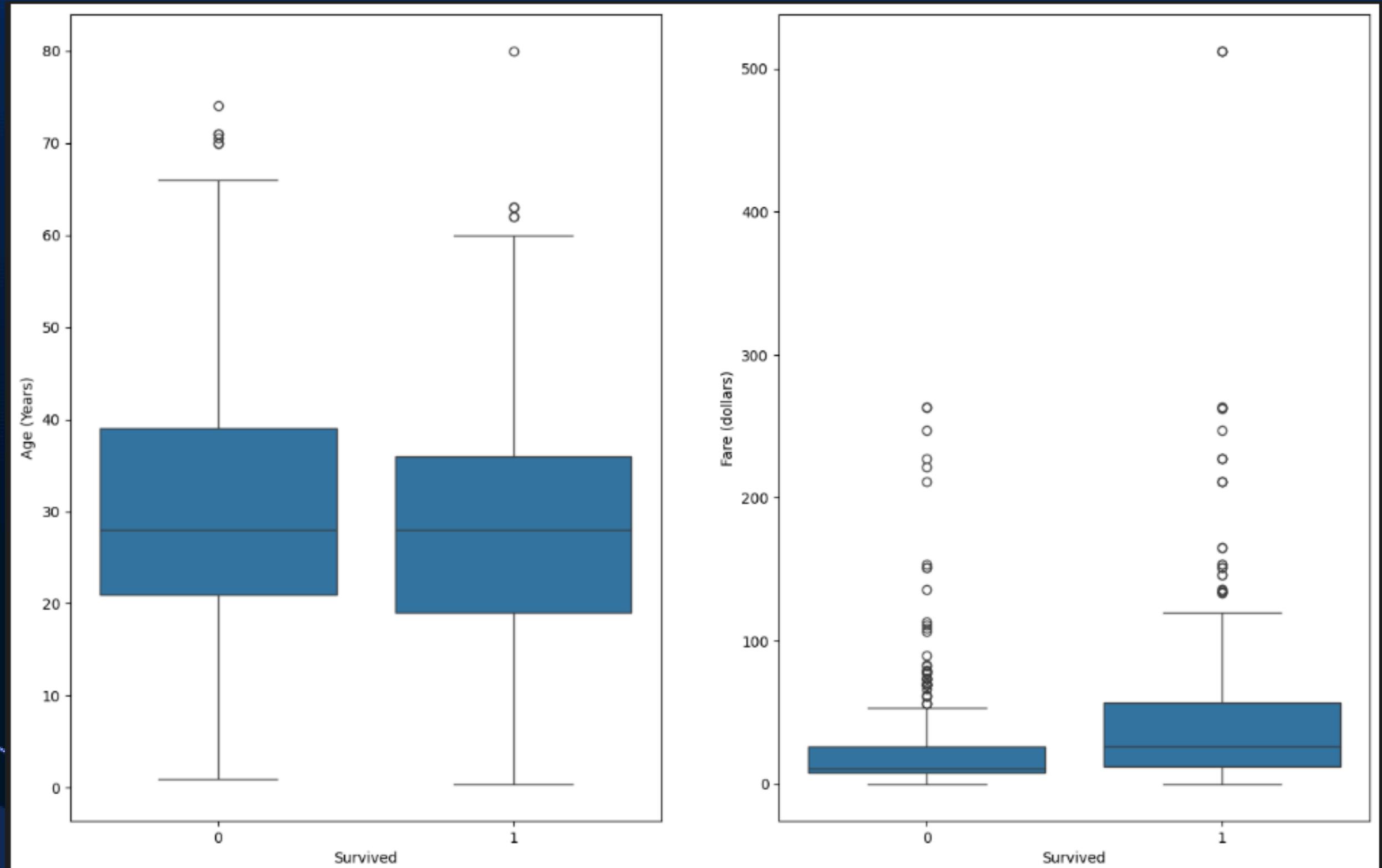


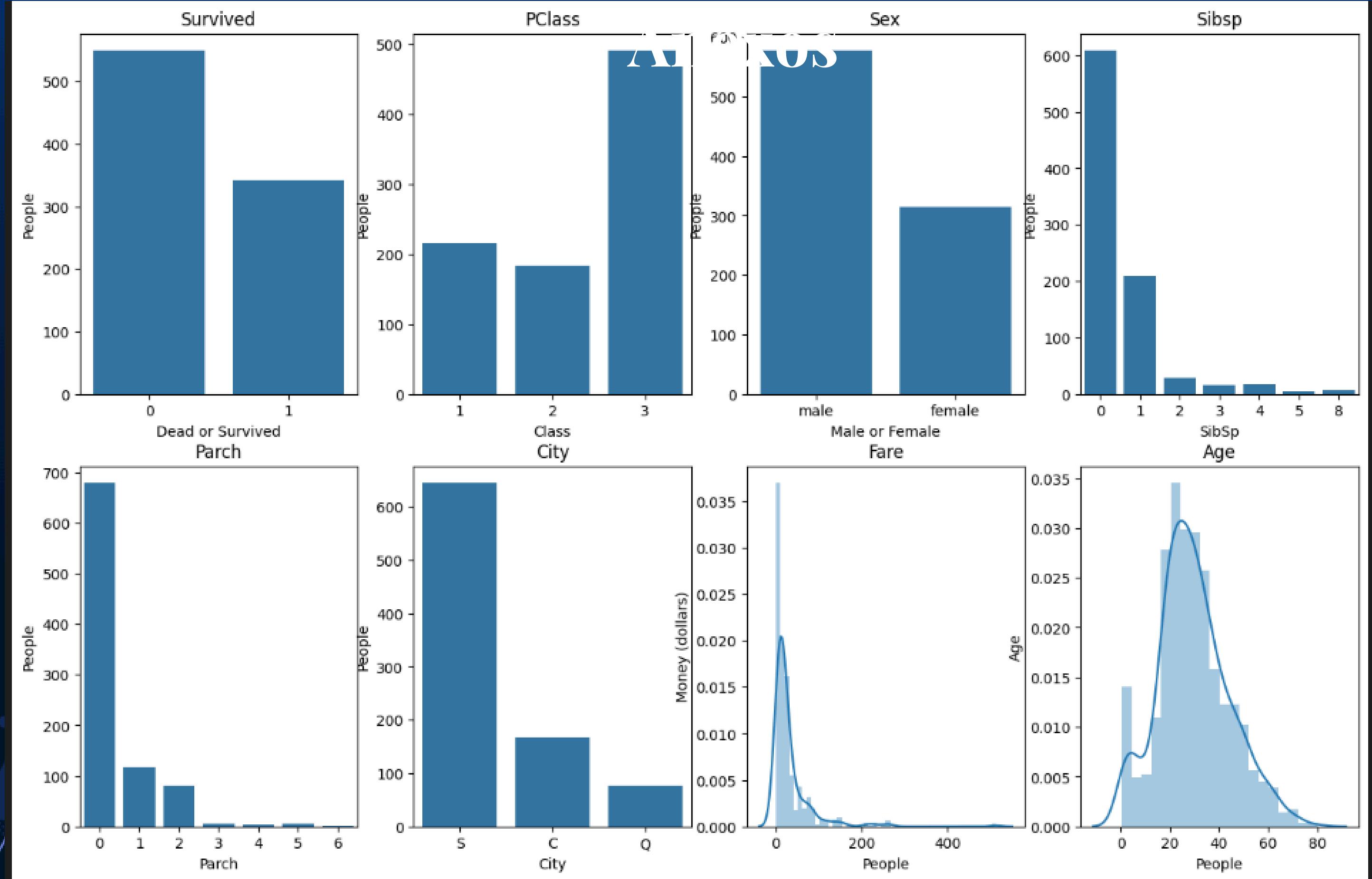
```
[1] data_frame
[✓] 0.0s
```

	Survived	Pclass	Sex	Age	Embarked	companions
0	0	3	1	22.0	0	1
1	1	1	0	38.0	1	1
2	1	3	0	26.0	0	0
3	1	1	0	35.0	0	1
4	0	3	1	35.0	0	0
...	...	...	...	...	...	...
886	0	2	1	27.0	0	0
887	1	1	0	19.0	0	0
888	0	3	0	28.0	0	1
889	1	1	1	26.0	1	0
890	0	3	1	32.0	1	0



# Anexos

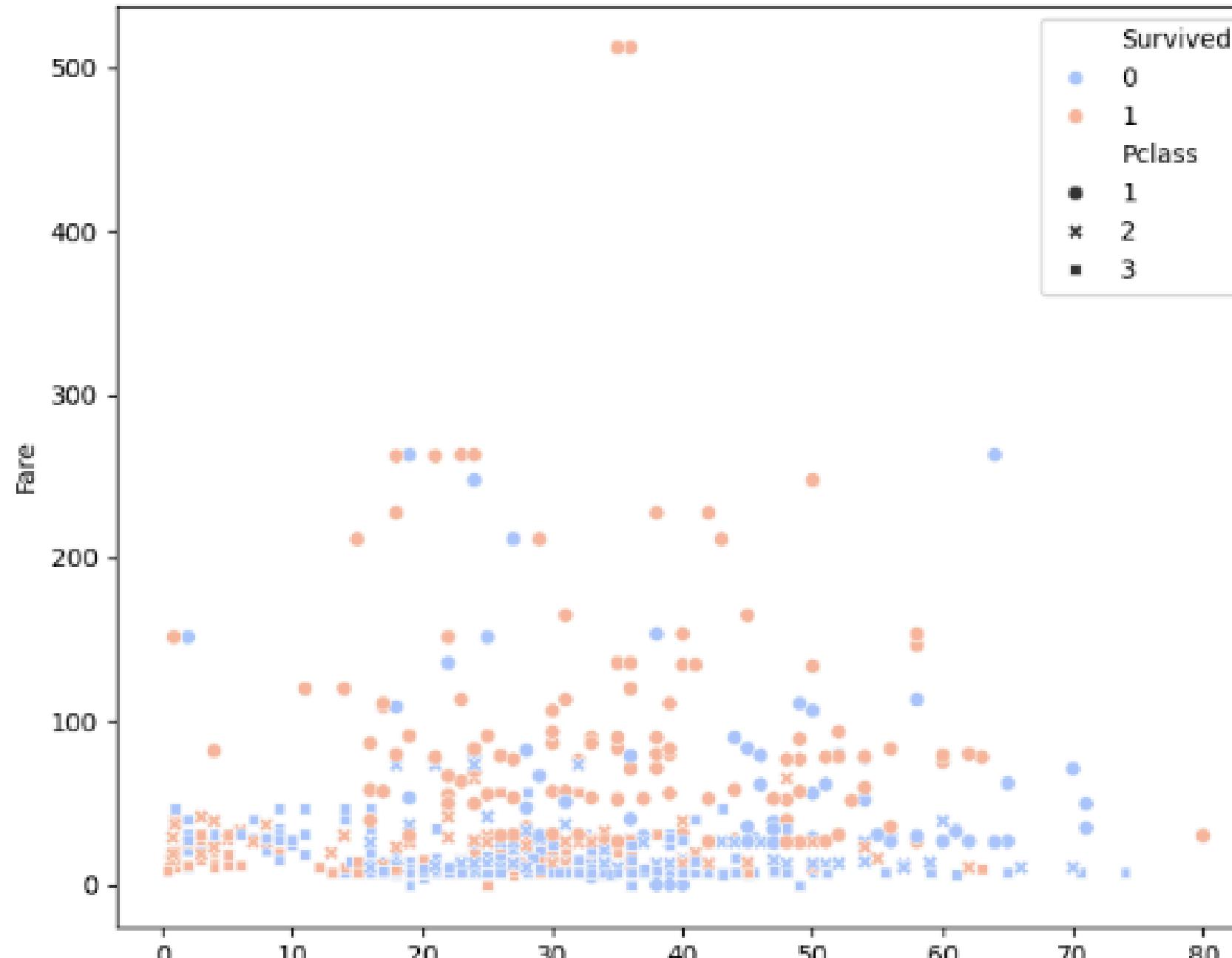




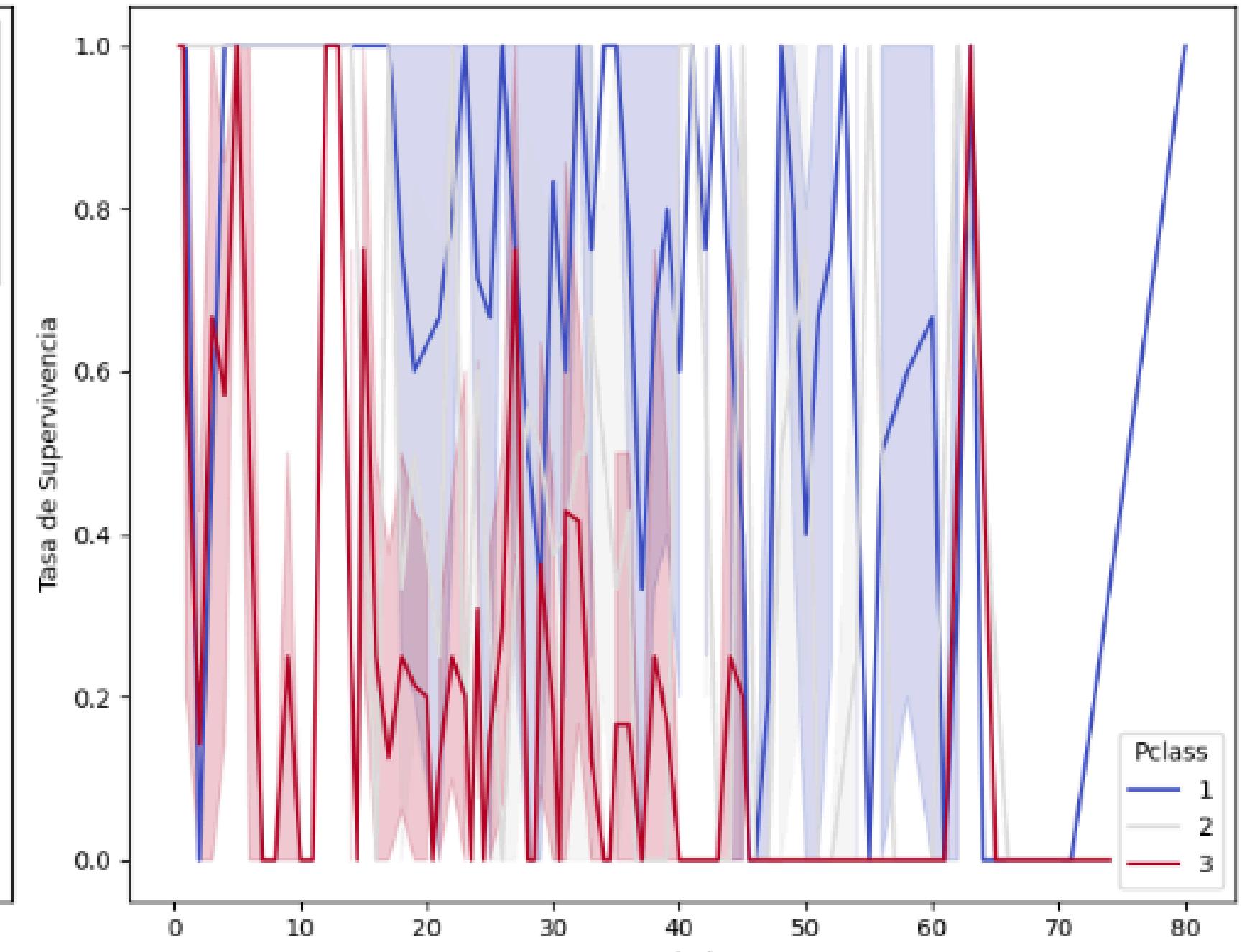
# Anexos



Interacción entre Edad, Fare y Supervivencia por Pclass



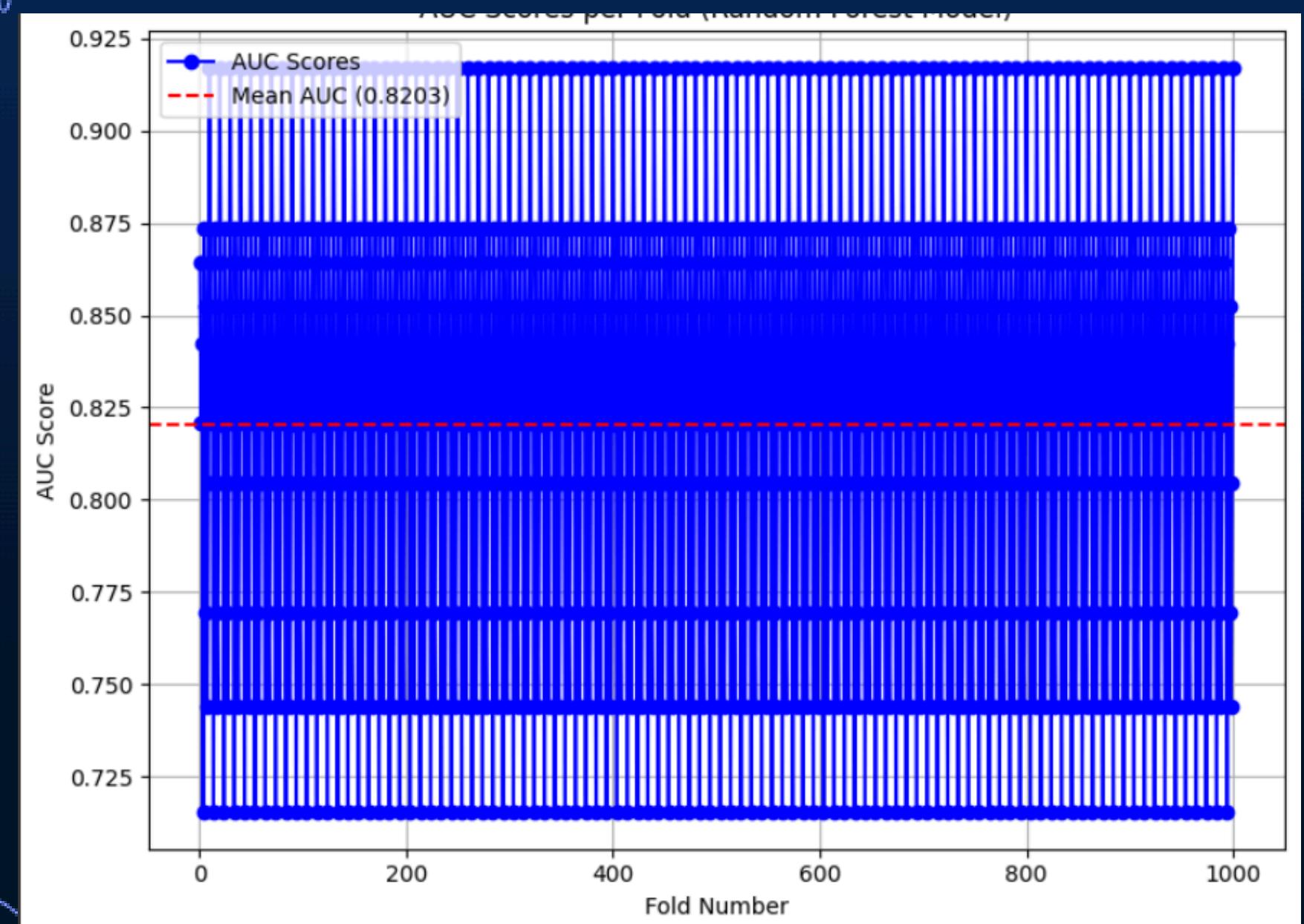
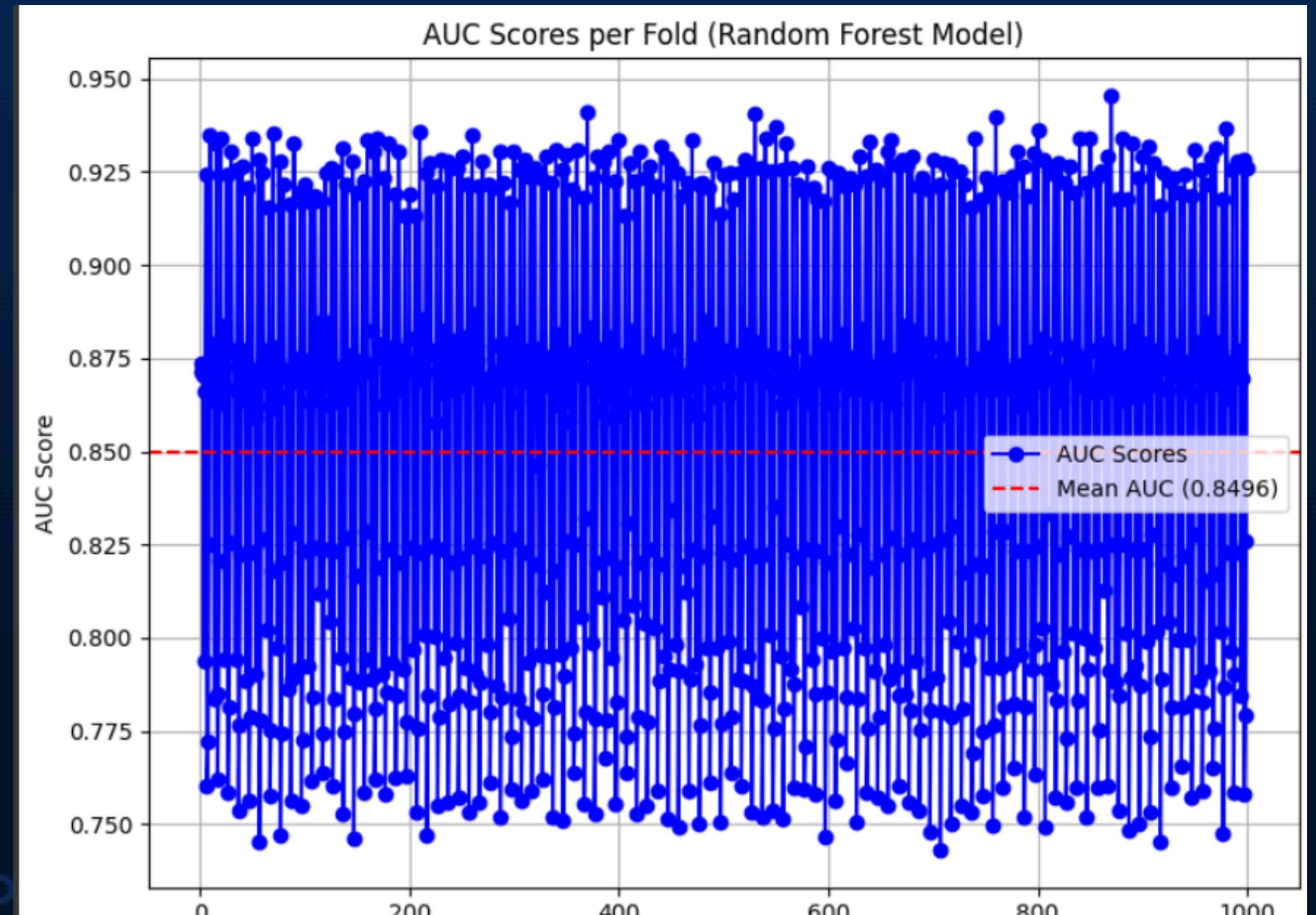
Tasa de Supervivencia por Edad y Pclass



# Anexos

RF

LR



# Anexos

SVM

