

MACHINE LEARNING

Project Presentation
DS105
Clement Leo

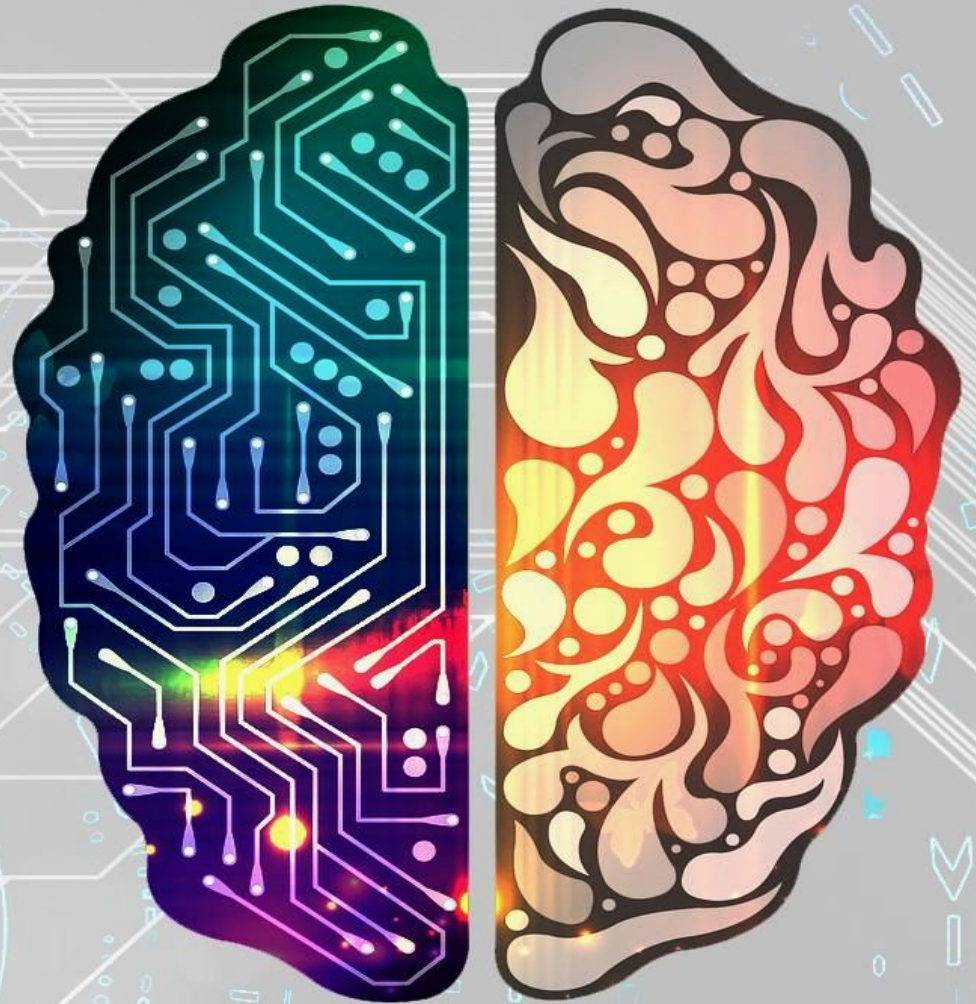


Table of Contents

- **Dataset Intro**
 - Brief glance at the set
- **General workflow**
 - Breakdown on the thought process
- **Preprocessing & Cleaning**
 - Setting up the dataset
- **Analysis & Insights**
 - Some observations
- **Machine Learning Models**
 - Models used and results

Dataset

This dataset is on the unsinkable ship, RMS Titanic.

Size of dataset: 12 columns, 891 entries

Some of the features:

- Survived, Passenger Class, Sex, Age

Source:

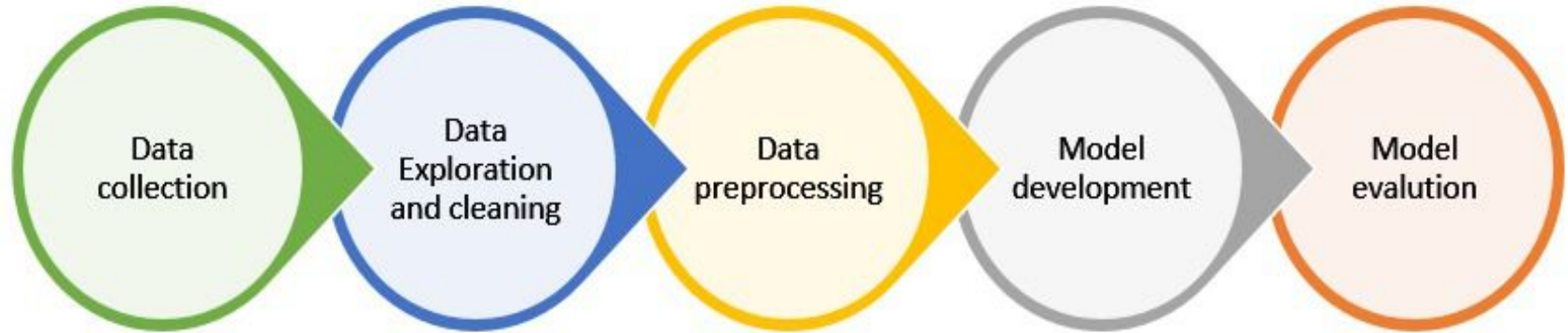
<https://www.kaggle.com/competitions/titanic>

Goal

To predict the survivability of a person on the ship using ML models and to see which would be the best model to use.

Passenger Id	Survived	Pclass	Name	Sex	Age	Sib Sp	Parch	Ticket	Fare	Cabin	Embarked
1	0	3	Braund, Mr. Owen	male	22	1	0	A/5 21171	7.25		S
2	1	1	Cummings, Mrs. John	female	38	1	0	PC 17599	71.28	C85	S
3	1	3	Heikkinen, Miss.	female	26	0	0	STON/O2 31012	7.93		S
4	1	1	Futrelle, Mrs. Jack	female	35	1	0	113803	53.1	C123	S
5	0	3	Allen, Mr. William	male	35	0	0	373450	8.05		S
6	0	3	Moran, Mr. James	male		0	0	330877	8.46		Q
7	0	1	McCarthy, Mr. Thomas	male	54	0	0	17463	51.86	E46	S
8	0	3	Palsson, Master.	male	2	3	1	349909	21.08		S
9	1	3	Johnson, Mrs. Oscar	female	27	0	2	347742	11.13		S
10	1	2	Nasser, Mrs. Nicholas	female	14	1	0	237736	30.07		C
11	1	3	Sandstrom, Miss.	female	4	1	1	PP 9549	16.7	G6	S
12	1	1	Bonnell, Miss. Elizabeth	female	58	0	0	113783	26.55	C103	S
13	0	3	Saunders, Mr.	male	20	0	0	A/5. 2151	8.05		S
14	0	3	Andersson, Mr. Anders	male	39	1	5	347082	31.28		S
15	0	3	Vestrom, Miss. Helene	female	14	0	0	350406	7.85		S
16	1	2	Hewlett, Mrs. (Mrs.)	female	55	0	0	248706	16		S
17	0	3	Rice, Master. Eugene	male	2	4	1	382652	29.13		Q
18	1	2	Williams, Mr. Charles	male		0	0	244373	13		S
19	0	3	Vander Planke, Mr.	male	31	1	0	345763	18		S
20	1	3	Massei, Mrs.	female		0	0	2649	7.23		C
21	0	2	Fynney, Mr. Joseph	male	35	0	0	239865	26		S
22	1	2	Beesley, Mr. Lawrence	male	34	0	0	248698	13	D56	S
23	1	3	McGowan, Miss.	female	15	0	0	330923	8.03		Q
24	1	1	Sloper, Mr. William	male	28	0	0	113788	35.5	A6	S
25	0	3	Palsson, Miss. Torvald	female	8	3	1	349909	21.08		S
26	1	3	Asplund, Mrs. Carl	female	38	1	5	347077	31.39		S
27	0	3	Emir, Mr. Farred	male		0	0	2631	7.23		C
28	0	1	Fortune, Mr. Charles	male	19	3	2	19950	263	C23 C25 C27	S
29	1	3	O'Dwyer, Miss. Elsie	female		0	0	330959	7.88		Q
30	0	3	Todoroff, Mr. Laila	male		0	0	349216	7.9		S
31	0	1	Urchurtu, Don.	male	40	0	0	PC 17601	27.72		C
32	1	1	Spencer, Mrs. William	female		1	0	PC 17569	146.52	B78	C
33	1	3	Glynn, Miss. Mary	female		0	0	335677	7.75		Q
34	0	2	Wheardon, Mr. Edgar	male	66	0	0	C.A. 24579	10.5		S
35	0	1	Meyer, Mr. Edgar	male	28	1	0	PC 17604	82.17		C
36	0	1	Holmerson, Mr. Anders	male	42	1	0	113789	52		S
37	1	3	Mamee, Mr. Han	male		0	0	2677	7.23		C
38	0	3	Cann, Mr. Ernest	male	21	0	0	A./5. 2152	8.05		S
39	0	3	Vander Planke, Mr.	male	18	2	0	345764	18		S
40	1	3	Nicola-Yarred, Mr.	male	14	1	0	2651	11.24		C
41	0	3	Ahlin, Mrs. Johan	female	40	1	0	7546	9.48		S
42	0	2	Turpin, Mrs. William	female	27	1	0	11668	21		S

General Workflow



Titanic Dataset

Cleaning up, Preprocess & Insights

Use different algorithms

Compare the models

Preprocessing & Cleaning

```
1 🔦 Null Checking  
2 df_train.isnull().sum()
```

✓ 0.4s

PassengerId	0
Survived	0
Pclass	0
Name	0
Sex	0
Age	177
SibSp	0
Parch	0
Ticket	0
Fare	0
Cabin	687
Embarked	2

dtype: int64

There are a total of 891 rows, 12 Columns.

However there are missing values Age - 177 (177 Missing) Cabin - 687 (687 Missing) Embarked - 2 (2 Missing)

Missing/NaN Values

3 features with null values.

Age ~20%

Cabin ~77%

Embarked ~0%

What I did

Age - Imputed median

Cabin - Totally dropped

Embarked - Input the most used value, 'S'

Preprocessing & Cleaning

```
1 # Combine Sibsp & Parch together as Companions, since these are there together as a group
2 ⚡
3 df_train['Companion'] = df_train['SibSp'] + df_train['Parch']
✓ 0.3s

1 # Drop the SibSp and Parch columns
2 ⚡
3 cols = ['SibSp', 'Parch']
4 df_train = df_train.drop(cols, axis=1)
✓ 0.3s
```

```
1 cols = ['Name', 'Ticket', 'Cabin']
2 df_train = df_train.drop(cols, axis=1)
✓ 0.4s
```

Combined similar features

'SibSp' & 'Parch'

Combined into Companion

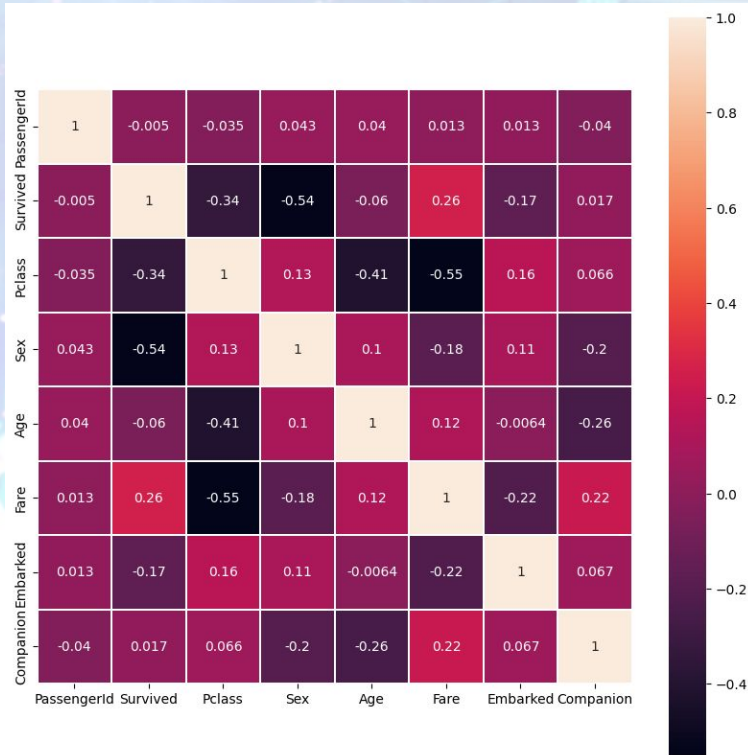
Dropped unhelpful features

Cabin

Name

Ticket

Analysis & Insights



Correlationship Heatmap

The darker the color, the worse the relationship

-Age & Pclass,

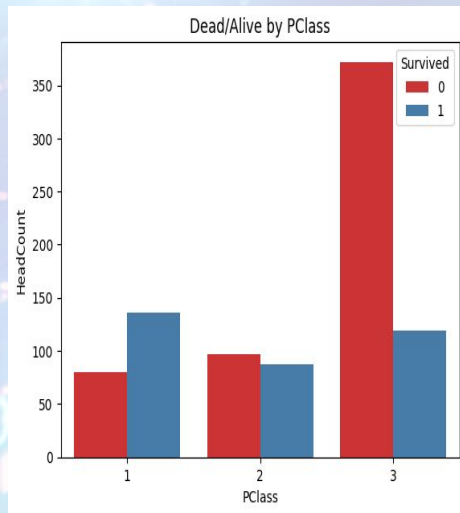
-Fare & PClass,

-Pclass & Survived

-Sex & Survived,

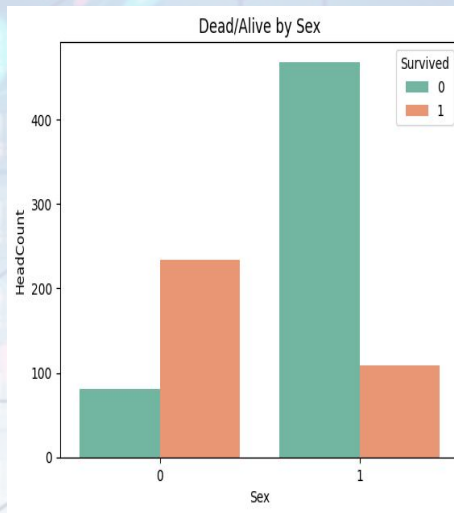
Have terrible relations.

Analysis & Insights



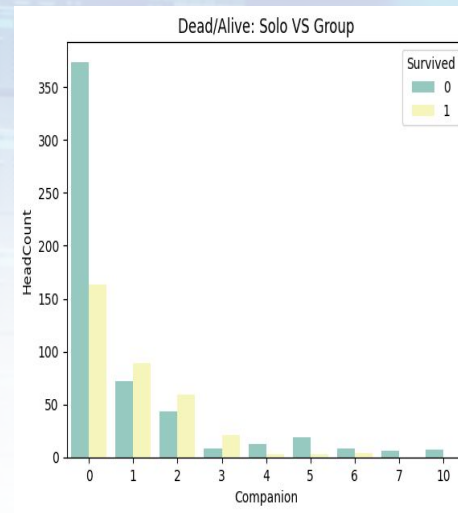
Dead/Alive by Pclass

1st class has the highest probability of surviving



Dead/Alive by Sex

Females are more likely to survive



Dead/Alive wrt Companion

Solo passengers are much more likely to die

Groups between 1-3 have better chances of surviving

Machine Learning Models

Score	Model
92.14	Random Forest
80.25	Logistic Regression
79.35	Naive Bayes
78.34	Support Vector Machines
73.51	KNN

Machine Learning Models (RandomForest)

```
1 # Random Forest
2
3 random_forest = RandomForestClassifier(criterion='gini',
4   n_estimators=700,
5   min_samples_split=10,
6   min_samples_leaf=1,
7   max_features='auto',
8   oob_score=True,
9   random_state=1,
10  n_jobs=-1)
11
12 random_forest.fit(X_train, Y_train)
13
14 Y_pred = random_forest.predict(X_test)
15
16 acc_random_forest = round(random_forest.score(X_train, Y_train) * 100, 2)
17 acc_random_forest
✓ 1.7s

c:\Users\user\anaconda3\envs\name-py10\lib\site-packages\sklearn\ensemble\_forest.py:425: FutureWarning: The attribute `max_features` will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove
RandomForestClassifiers and ExtraTreesClassifiers.
  warn(
92.14

1 rf = RandomForestClassifier(n_estimators=700)
2 scores = cross_val_score(rf, X_train, Y_train, cv=10, scoring = "accuracy")
3 print("Scores:", scores)
4 print("Mean:", scores.mean())
5 print("Standard Deviation:", scores.std())
✓ 13.9s

Scores: [0.77777778 0.78651685 0.75280899 0.79775281 0.88764045 0.85393258
0.85393258 0.80898876 0.82022472 0.82022472]
Mean: 0.8159800249687889
Standard Deviation: 0.03840732207728619

1 # Initial Run - 92.14
2 ⚡ X-Validation shows that it has an accuracy of 81.5% with deviation of 3.8%
```

Referenced a tuned RandomForest code of a user

Simpler Randomforest code gave a much higher score

Model ran a score of 92.14

Cross-Validation gave accuracy of 81.5%, deviation of 3.8%

Machine Learning Models (Log Regression)

```
1 # Log Regression (variable) Y_train: Series
2 (variable) Y_train: Any
3 logreg = LogisticReg
4 logreg.fit(X_train, Y_train)
5
6 Y_pred = logreg.predict(X_test)
7
8 acc_log = round(logreg.score(X_train, Y_train) * 100, 2)
9 acc_log
```

✓ 0.6s

c:\Users\user\anaconda3\envs\name-py10\lib\site-packages\sklearn\linear_model_logistic.py:450: UserWarning: SolverWarning: Solver did not converge. STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
<https://scikit-learn.org/stable/modules/preprocessing.html>
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
n_iter_i = _check_optimize_result(

80.25

```
1 rf = LogisticRegression()
2 scores = cross_val_score(rf, X_train, Y_train, cv=10, scoring = "accuracy")
3 print("Scores:", scores)
4 print("Mean:", scores.mean())
5 print("Standard Deviation:", scores.std())
```

✓ 0.2s

Scores: [0.78888889 0.80898876 0.76404494 0.84269663 0.80898876 0.7752809
0.80898876 0.79775281 0.83146067 0.79775281]
Mean: 0.8024843945068664
Standard Deviation: 0.02242945427644189

Model ran a score of 80.25

Cross-Validation gave accuracy of 80.24%,
deviation of 2.24%

However, would require more iterations as
per the advice.

Could potentially be a good model to use

Machine Learning Models (Naive Bayes)

```
1 # Naive Bayes
2
3 gaussian = GaussianNB()
4 gaussian.fit(X_train, Y_train)
5
6 Y_pred = gaussian.predict(X_test)
7
8 acc_gaussian = round(gaussian.score(X_train, Y_train) * 100, 2)
9 acc_gaussian
```

✓ 0.1s

79.35

```
1 rf = GaussianNB()
2 scores = cross_val_score(rf, X_train, Y_train, cv=10, scoring = "accuracy")
3 print("Scores:", scores)
4 print("Mean:", scores.mean())
5 print("Standard Deviation:", scores.std())
```

✓ 0.1s

Scores: [0.71111111 0.76404494 0.7752809 0.79775281 0.79775281 0.7752809
0.82022472 0.83146067 0.7752809 0.82022472]

Mean: 0.7868414481897628

Standard Deviation: 0.0333370363637702

```
1 # Initial Run - 79.35
2 ⚡ X-Validation shows that it has an accuracy of 78.7% with deviation of 3.33%
```

Model ran a score of 79.35

Cross-Validation gave accuracy of 78.7%,
deviation of 3.33%

Not bad, but not the best

Machine Learning Models (SVM)

```

1 # SVM
2
3 linear_svc = LinearSVC()
4 linear_svc.fit(X_train, Y_train)
5
6 Y_pred = linear_svc.predict(X_test)
7
8 acc_linear_svc = round(linear_svc.score(X_train, Y_train) * 100, 2)
9 acc_linear_svc
✓ 0.5s

c:\Users\user\anaconda3\envs\name-py10\lib\site-packages\sklearn\svm\_base.py:1244: ConvergenceWarning: Liblinear failed to converge,
iterations.
  warnings.warn(

78.34

1 rf = LinearSVC()
2 scores = cross_val_score(rf, X_train, Y_train, cv=10, scoring = "accuracy")
3 print("Scores:", scores)
4 print("Mean:", scores.mean())
5 print("Standard Deviation:", scores.std())
✓ 0.3s

c:\Users\user\anaconda3\envs\name-py10\lib\site-packages\sklearn\svm\_base.py:1244: ConvergenceWarning: Liblinear failed to converge,
iterations.
  warnings.warn(

c:\Users\user\anaconda3\envs\name-py10\lib\site-packages\sklearn\svm\_base.py:1244: ConvergenceWarning: Liblinear failed to converge,
iterations.
  warnings.warn(

c:\Users\user\anaconda3\envs\name-py10\lib\site-packages\sklearn\svm\_base.py:1244: ConvergenceWarning: Liblinear failed to converge,
iterations.
  warnings.warn(

c:\Users\user\anaconda3\envs\name-py10\lib\site-packages\sklearn\svm\_base.py:1244: ConvergenceWarning: Liblinear failed to converge,
iterations.
  warnings.warn(

c:\Users\user\anaconda3\envs\name-py10\lib\site-packages\sklearn\svm\_base.py:1244: ConvergenceWarning: Liblinear failed to converge,
iterations.
  warnings.warn(

c:\Users\user\anaconda3\envs\name-py10\lib\site-packages\sklearn\svm\_base.py:1244: ConvergenceWarning: Liblinear failed to converge,
iterations.
  warnings.warn(

Scores: [0.71111111 0.80098876 0.65168539 0.84269663 0.7752809 0.60674157
0.80098876 0.70651685 0.70786517 0.79775281]
Mean: 0.7497627965043694

```

Model ran a score of 78.34

Cross-Validation gave accuracy of 75%, deviation of 7.3%

Repeated warnings of 'LibLinear failed to converge'

Current parameters are not good enough,
hence predictions from this model might be
useless

Machine Learning Models (KNN)

```
1 # KNN
2
3 knn = KNeighborsClass (variable) X_test: DataFrame
4 knn.fit(X_train, Y_train) (variable) X_test: Any
5
6 Y_pred = knn.predict(X_test)
7
8 acc_knn = round(knn.score(X_train, Y_train) * 100, 2)
9 acc_knn
```

✓ 0.8s

73.51

```
1 rf = KNeighborsClassifier(n_neighbors = 6)
2 scores = cross_val_score(rf, X_train, Y_train, cv=10, scoring = "accuracy")
3 print("Scores:", scores)
4 print("Mean:", scores.mean())
5 print("Standard Deviation:", scores.std())
```

✓ 0.1s

Scores: [0.48888889 0.62921348 0.49438202 0.46067416 0.4494382 0.52808989
0.50561798 0.50561798 0.68539326 0.65168539]

Mean: 0.5399001248439451

Standard Deviation: 0.07953945202686367

```
1 # Initial Run - 73.51
2 ⚡ X-Validation shows that it has an accuracy of 54% with deviation of 8%
3 # Hence this model isnt very good
```

Model ran a score of 73.51

Cross-Validation gave accuracy of 54%,
deviation of 8%

Not a good predictive model



<https://github.com/Clemykeileo/ML-Project.git>

Thank you!