

```
In [1]: from IPython.display import Image  
Image(url= "https://media.nationalgeographic.org/assets/photos/000/273/27302.jpg")
```

Out[1]:

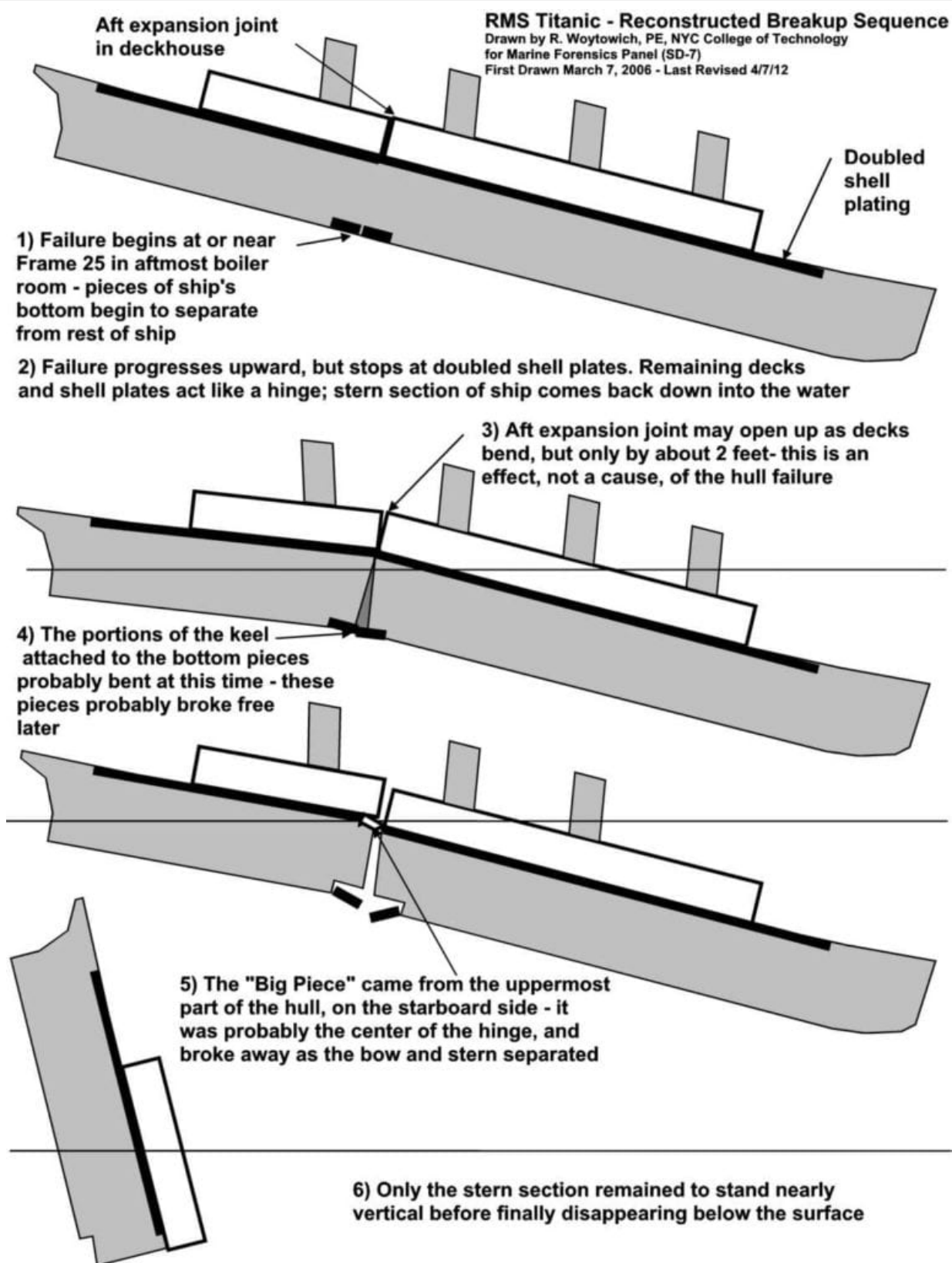


RMS Titanic was a British passenger liner operated by the White Star Line that sank in the North Atlantic Ocean on 15 April 1912, after striking an iceberg during her maiden voyage from Southampton to New York City. The RMS Titanic was the largest ship afloat at the time it entered service and was the second of three Olympic-class ocean liners operated by the White Star Line. The Titanic was built by the Harland and Wolff shipyard in Belfast. Thomas Andrews, her architect, died in the disaster.

In [2]:

<https://www.wufbv-wpengine.netdna-ssl.com/wp-content/uploads/2017/12/titanicBreakupRecons>

Out[2]:



```
In [3]: # linear algebra
import numpy as np

# data processing
import pandas as pd

# data visualization
import seaborn as sns
%matplotlib inline
from matplotlib import pyplot as plt
from matplotlib import style

# Algorithms
from sklearn import linear_model
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import Perceptron
from sklearn.linear_model import SGDClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.svm import SVC, LinearSVC
from sklearn.naive_bayes import GaussianNB
```

```
In [4]: test_df = pd.read_csv("test.csv")
train_df = pd.read_csv("train.csv")
```

In [5]: `train_df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
#   Column          Non-Null Count  Dtype
---  -
0   PassengerId     891 non-null    int64
1   Survived        891 non-null    int64
2   Pclass         891 non-null    int64
3   Name            891 non-null    object
4   Sex            891 non-null    object
5   Age            714 non-null    float64
6   SibSp          891 non-null    int64
7   Parch          891 non-null    int64
8   Ticket         891 non-null    object
9   Fare           891 non-null    float64
10  Cabin          204 non-null    object
11  Embarked       889 non-null    object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB
```

In [6]: `test_df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 418 entries, 0 to 417
Data columns (total 11 columns):
#   Column          Non-Null Count  Dtype
---  -
0   PassengerId     418 non-null    int64
1   Pclass         418 non-null    int64
2   Name            418 non-null    object
3   Sex            418 non-null    object
4   Age            332 non-null    float64
5   SibSp          418 non-null    int64
6   Parch          418 non-null    int64
7   Ticket         418 non-null    object
8   Fare           417 non-null    float64
9   Cabin          91 non-null     object
10  Embarked       418 non-null    object
dtypes: float64(2), int64(4), object(5)
memory usage: 36.0+ KB
```

Data Dictionary Survived: 0 = No, 1 = Yes pclass: Ticket class 1 = 1st, 2 = 2nd, 3 = 3rd sibsp: # of siblings / spouses aboard the Titanic parch: # of parents / children aboard the Titanic ticket: Ticket number cabin: Cabin number embarked: Port of Embarkation C = Cherbourg, Q = Queenstown, S = Southampton Total rows and columns

We can see that there are 891 rows and 12 columns in our training dataset.

```
In [7]: train_df.describe()
```

Out[7]:

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
count	891.000000	891.000000	891.000000	714.000000	891.000000	891.000000	891.000000
mean	446.000000	0.383838	2.308642	29.699118	0.523008	0.381594	32.204208
std	257.353842	0.486592	0.836071	14.526497	1.102743	0.806057	49.693429
min	1.000000	0.000000	1.000000	0.420000	0.000000	0.000000	0.000000
25%	223.500000	0.000000	2.000000	20.125000	0.000000	0.000000	7.910400
50%	446.000000	0.000000	3.000000	28.000000	0.000000	0.000000	14.454200
75%	668.500000	1.000000	3.000000	38.000000	1.000000	0.000000	31.000000
max	891.000000	1.000000	3.000000	80.000000	8.000000	6.000000	512.329200

Above we can see that 38% out of the training-set survived the Titanic. We can also see that the passenger ages range from 0.4 to 80. On top of that we can already detect some features, that contain missing values, like the 'Age' feature.

what data is actually missing:

```
In [9]: total = train_df.isnull().sum().sort_values(ascending=False)
percent_1 = train_df.isnull().sum()/train_df.isnull().count()*100
percent_2 = (round(percent_1, 1)).sort_values(ascending=False)
missing_data = pd.concat([total, percent_2], axis=1, keys=['Total', '%'])
missing_data.head(5)
```

Out[9]:

	Total	%
Cabin	687	77.1
Age	177	19.9
Embarked	2	0.2
Fare	0	0.0
Ticket	0	0.0

```
In [10]: train_df.columns.values
```

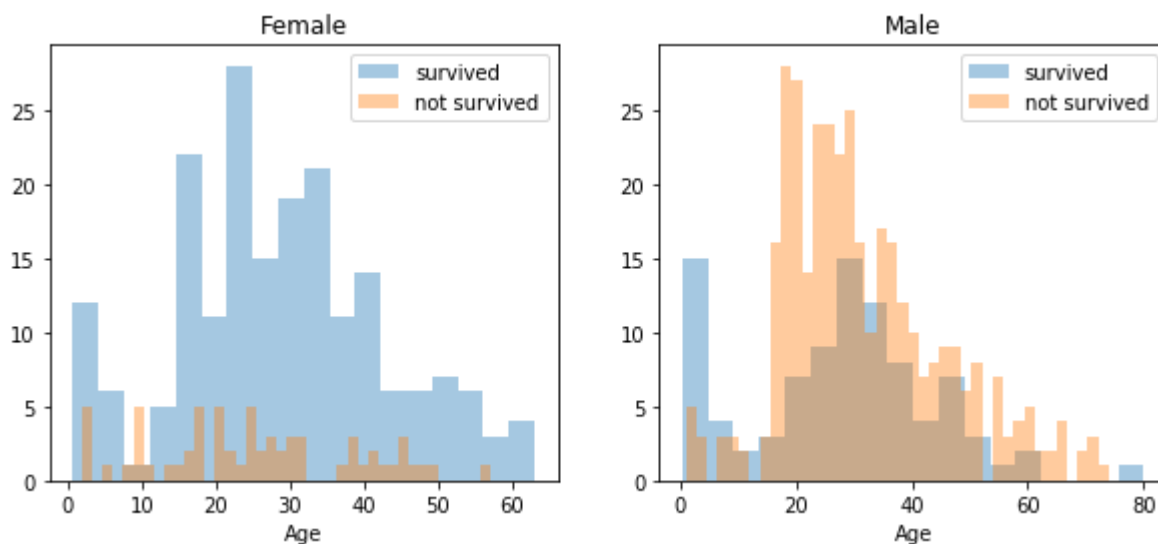
```
Out[10]: array(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',
                'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'], dtype=object)
```

Above you can see the 11 features + the target variable (survived). What features could contribute to a high survival rate ? To me it would make sense if everything except 'PassengerId', 'Ticket' and 'Name' would be correlated with a high survival rate.

```

In [11]: survived = 'survived'
not_survived = 'not survived'
fig, axes = plt.subplots(nrows=1, ncols=2, figsize=(10, 4))
women = train_df[train_df['Sex']=='female']
men = train_df[train_df['Sex']=='male']
ax = sns.distplot(women[women['Survived']==1].Age.dropna(), bins=18, label = survived)
ax = sns.distplot(women[women['Survived']==0].Age.dropna(), bins=40, label = not_survived)
ax.legend()
ax.set_title('Female')
ax = sns.distplot(men[men['Survived']==1].Age.dropna(), bins=18, label = survived)
ax = sns.distplot(men[men['Survived']==0].Age.dropna(), bins=40, label = not_survived)
ax.legend()
_ = ax.set_title('Male')

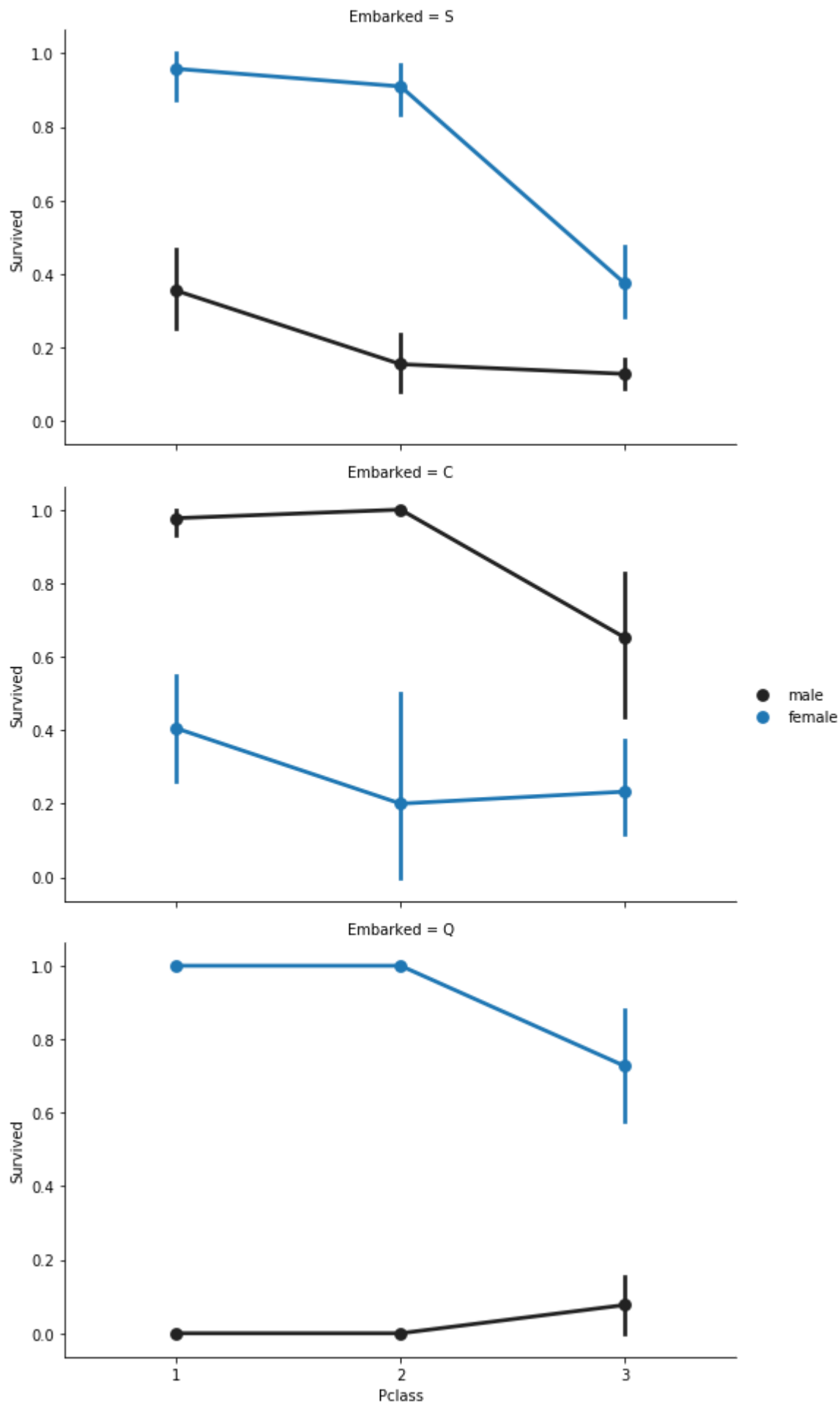
```



```
In [12]: FacetGrid = sns.FacetGrid(train_df, row='Embarked', size=4.5, aspect=1.6)
FacetGrid.map(sns.pointplot, 'Pclass', 'Survived', 'Sex', palette=None, order=None)
FacetGrid.add_legend()
```

```
C:\Users\work\anaconda3\lib\site-packages\seaborn\axisgrid.py:243: UserWarning:
The `size` parameter has been renamed to `height`; please update your code.
  warnings.warn(msg, UserWarning)
```

```
Out[12]: <seaborn.axisgrid.FacetGrid at 0x154400d5ac8>
```



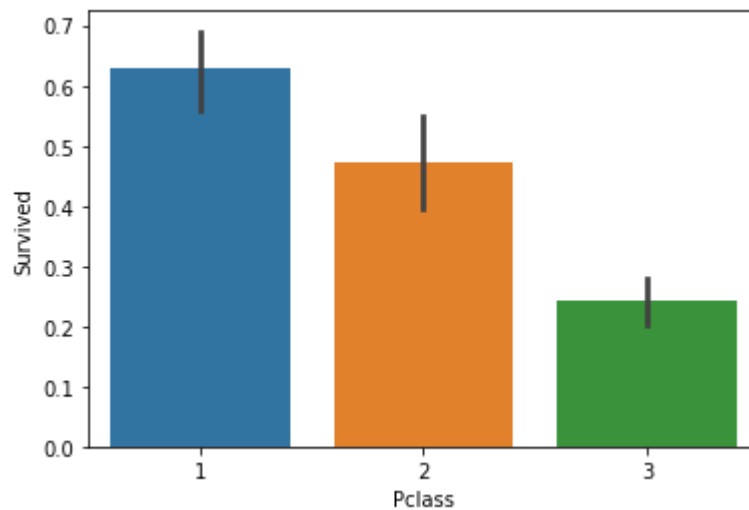
In [13]:

```
ve a high survival probability if they are on port C, but a low probability if th
```

File "<ipython-input-13-c746ddd5f596>", line 1

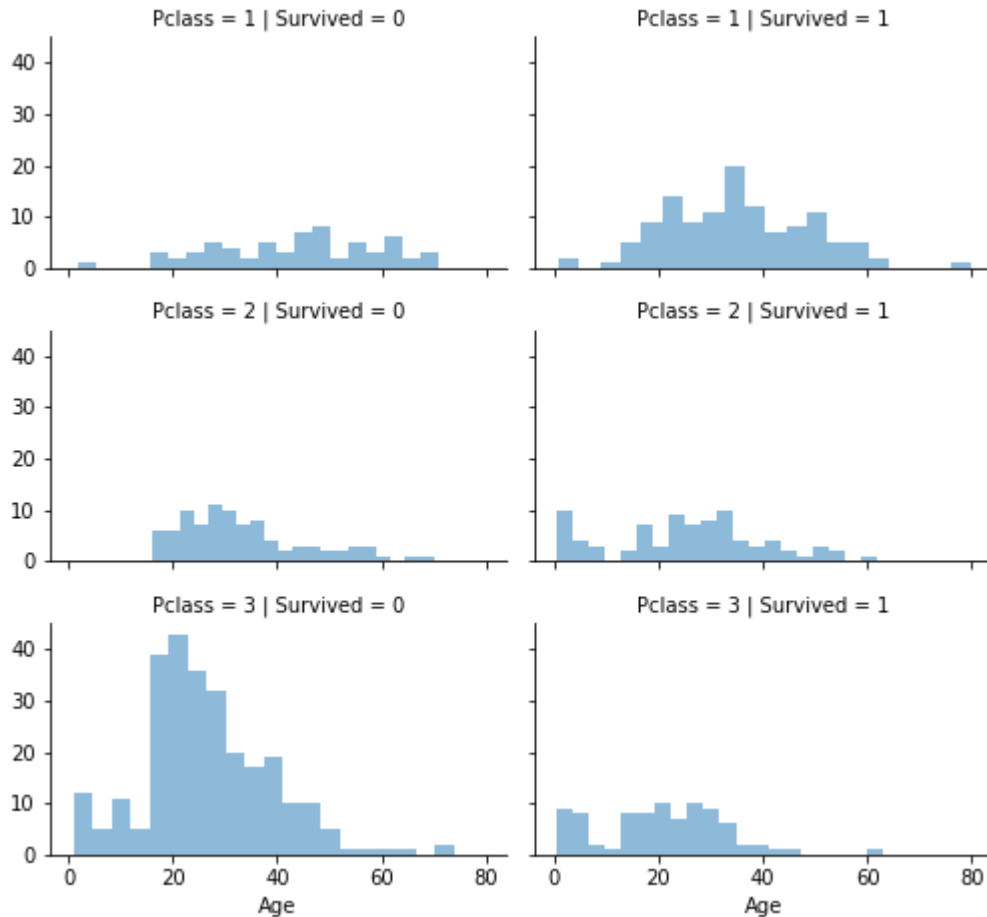
Embarked seems to be correlated with survival, depending on the gender.

SyntaxError: invalid syntax

In [14]: `sns.barplot(x='Pclass', y='Survived', data=train_df)`Out[14]: `<matplotlib.axes._subplots.AxesSubplot at 0x15440444bc8>`

```
In [15]: grid = sns.FacetGrid(train_df, col='Survived', row='Pclass', size=2.2, aspect=1.6)
grid.map(plt.hist, 'Age', alpha=.5, bins=20)
grid.add_legend();
```

C:\Users\work\anaconda3\lib\site-packages\seaborn\axisgrid.py:243: UserWarning:
The `size` parameter has been renamed to `height`; please update your code.
warnings.warn(msg, UserWarning)



```
In [16]: data = [train_df, test_df]
for dataset in data:
    dataset['relatives'] = dataset['SibSp'] + dataset['Parch']
    dataset.loc[dataset['relatives'] > 0, 'not_alone'] = 0
    dataset.loc[dataset['relatives'] == 0, 'not_alone'] = 1
    dataset['not_alone'] = dataset['not_alone'].astype(int)
train_df['not_alone'].value_counts()
```

```
Out[16]: 1    537
         0    354
         Name: not_alone, dtype: int64
```

Missing Data:

Cabin:

As a reminder, we have to deal with Cabin (687), Embarked (2) and Age (177). First I thought, we have to delete the 'Cabin' variable but then I found something interesting. A cabin number looks like 'C123' and the letter refers to the deck. Therefore we're going to extract these and create a new feature, that contains a persons deck. Afterwards we will convert the feature into a numeric variable. The missing values will be converted to zero. In the picture below you can see the actual decks of the titanic, ranging from A to G.

```
In [17]: import re
deck = {"A": 1, "B": 2, "C": 3, "D": 4, "E": 5, "F": 6, "G": 7, "U": 8}
data = [train_df, test_df]

for dataset in data:
    dataset['Cabin'] = dataset['Cabin'].fillna("U0")
    dataset['Deck'] = dataset['Cabin'].map(lambda x: re.compile("([a-zA-Z]+)").search(x).group(1))
    dataset['Deck'] = dataset['Deck'].map(deck)
    dataset['Deck'] = dataset['Deck'].fillna(0)
    dataset['Deck'] = dataset['Deck'].astype(int)
# we can now drop the cabin feature
train_df = train_df.drop(['Cabin'], axis=1)
test_df = test_df.drop(['Cabin'], axis=1)
```

```
In [18]: data = [train_df, test_df]

for dataset in data:
    mean = train_df["Age"].mean()
    std = test_df["Age"].std()
    is_null = dataset["Age"].isnull().sum()
    # compute random numbers between the mean, std and is_null
    rand_age = np.random.randint(mean - std, mean + std, size = is_null)
    # fill NaN values in Age column with random values generated
    age_slice = dataset["Age"].copy()
    age_slice[np.isnan(age_slice)] = rand_age
    dataset["Age"] = age_slice
    dataset["Age"] = train_df["Age"].astype(int)
train_df["Age"].isnull().sum()
```

Out[18]: 0

```
In [19]: train_df['Embarked'].describe()
```

```
Out[19]: count      889
         unique        3
         top          S
         freq        644
         Name: Embarked, dtype: object
```

```
In [20]: train_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 14 columns):
 #   Column          Non-Null Count  Dtype
---  -
 0   PassengerId     891 non-null   int64
 1   Survived        891 non-null   int64
 2   Pclass          891 non-null   int64
 3   Name            891 non-null   object
 4   Sex             891 non-null   object
 5   Age             891 non-null   int32
 6   SibSp           891 non-null   int64
 7   Parch           891 non-null   int64
 8   Ticket          891 non-null   object
 9   Fare            891 non-null   float64
10   Embarked        889 non-null   object
11   relatives       891 non-null   int64
12   not_alone       891 non-null   int32
13   Deck            891 non-null   int32
dtypes: float64(1), int32(3), int64(6), object(4)
memory usage: 87.1+ KB
```

```
In [21]: data = [train_df, test_df]

for dataset in data:
    dataset['Fare'] = dataset['Fare'].fillna(0)
    dataset['Fare'] = dataset['Fare'].astype(int)
```

```
In [22]: data = [train_df, test_df]
titles = {"Mr": 1, "Miss": 2, "Mrs": 3, "Master": 4, "Rare": 5}

for dataset in data:
    # extract titles
    dataset['Title'] = dataset.Name.str.extract('([A-Za-z+])\.', expand=False)
    # replace titles with a more common title or as Rare
    dataset['Title'] = dataset['Title'].replace(['Lady', 'Countess', 'Capt', 'Col',
                                                'Major', 'Rev', 'Sir', 'Jonkheer', 'Duke'], 'Rare')
    dataset['Title'] = dataset['Title'].replace('Mlle', 'Miss')
    dataset['Title'] = dataset['Title'].replace('Ms', 'Miss')
    dataset['Title'] = dataset['Title'].replace('Mme', 'Mrs')
    # convert titles into numbers
    dataset['Title'] = dataset['Title'].map(titles)
    # filling NaN with 0, to get safe
    dataset['Title'] = dataset['Title'].fillna(0)
train_df = train_df.drop(['Name'], axis=1)
test_df = test_df.drop(['Name'], axis=1)
```

Sex: Convert 'Sex' feature into numeric.

```
In [24]: genders = {"male": 0, "female": 1}
data = [train_df, test_df]

for dataset in data:
    dataset['Sex'] = dataset['Sex'].map(genders)
```

```
In [25]: ports = {"S": 0, "C": 1, "Q": 2}
data = [train_df, test_df]

for dataset in data:
    dataset['Embarked'] = dataset['Embarked'].map(ports)
```

```
In [26]: train_df = train_df.drop(['Ticket'], axis=1)
test_df = test_df.drop(['Ticket'], axis=1)
```

```
In [27]: data = [train_df, test_df]
for dataset in data:
    dataset['Age'] = dataset['Age'].astype(int)
    dataset.loc[ dataset['Age'] <= 11, 'Age'] = 0
    dataset.loc[(dataset['Age'] > 11) & (dataset['Age'] <= 18), 'Age'] = 1
    dataset.loc[(dataset['Age'] > 18) & (dataset['Age'] <= 22), 'Age'] = 2
    dataset.loc[(dataset['Age'] > 22) & (dataset['Age'] <= 27), 'Age'] = 3
    dataset.loc[(dataset['Age'] > 27) & (dataset['Age'] <= 33), 'Age'] = 4
    dataset.loc[(dataset['Age'] > 33) & (dataset['Age'] <= 40), 'Age'] = 5
    dataset.loc[(dataset['Age'] > 40) & (dataset['Age'] <= 66), 'Age'] = 6
    dataset.loc[ dataset['Age'] > 66, 'Age'] = 6

# Let's see how it's distributed train_df['Age'].value_counts()
```

```
In [28]: train_df.head(10)
```

Out[28]:

	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked	relatives	not_alone
0	1	0	3	0	2	1	0	7	0.0	1	0
1	2	1	1	1	5	1	0	71	1.0	1	0
2	3	1	3	1	3	0	0	7	0.0	0	1
3	4	1	1	1	5	1	0	53	0.0	1	0
4	5	0	3	0	5	0	0	8	0.0	0	1
5	6	0	3	0	4	0	0	8	2.0	0	1
6	7	0	1	0	6	0	0	51	0.0	0	1
7	8	0	3	0	0	3	1	21	0.0	4	0
8	9	1	3	1	3	0	2	11	0.0	2	0
9	10	1	2	1	1	1	0	30	1.0	1	0

```
In [29]: data = [train_df, test_df]
for dataset in data:
    dataset.loc[ dataset['Fare'] <= 7.91, 'Fare'] = 0
    dataset.loc[(dataset['Fare'] > 7.91) & (dataset['Fare'] <= 14.454), 'Fare'] = 1
    dataset.loc[(dataset['Fare'] > 14.454) & (dataset['Fare'] <= 31), 'Fare'] = 2
    dataset.loc[(dataset['Fare'] > 31) & (dataset['Fare'] <= 99), 'Fare'] = 3
    dataset.loc[(dataset['Fare'] > 99) & (dataset['Fare'] <= 250), 'Fare'] = 4
    dataset.loc[ dataset['Fare'] > 250, 'Fare'] = 5
    dataset['Fare'] = dataset['Fare'].astype(int)
```

```
In [30]: data = [train_df, test_df]
for dataset in data:
    dataset['Age_Class'] = dataset['Age'] * dataset['Pclass']
```

```
In [31]: for dataset in data:
        dataset['Fare_Per_Person'] = dataset['Fare']/(dataset['relatives']+1)
        dataset['Fare_Per_Person'] = dataset['Fare_Per_Person'].astype(int)
        # Let's take a last look at the training set, before we start training the models
train_df.head(10)
```

Out[31]:

	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked	relatives	not_alone
0	1	0	3	0	2	1	0	0	0.0	1	0
1	2	1	1	1	5	1	0	3	1.0	1	0
2	3	1	3	1	3	0	0	0	0.0	0	1
3	4	1	1	1	5	1	0	3	0.0	1	0
4	5	0	3	0	5	0	0	1	0.0	0	1
5	6	0	3	0	4	0	0	1	2.0	0	1
6	7	0	1	0	6	0	0	3	0.0	0	1
7	8	0	3	0	0	3	1	2	0.0	4	0
8	9	1	3	1	3	0	2	1	0.0	2	0
9	10	1	2	1	1	1	0	2	1.0	1	0

Machine Learning Models

```
In [38]: X_train = train_df.drop("Survived", axis=1)
        Y_train = train_df["Survived"]
        X_test  = test_df.drop("PassengerId", axis=1)
```

```
In [40]: from sklearn.model_selection import KFold
        from sklearn.model_selection import cross_val_score
        k_fold = KFold(n_splits=10, shuffle=True, random_state=0)
```

```
In [ ]: clf = KNeighborsClassifier(n_neighbors = 13)
        scoring = 'accuracy'
        score = cross_val_score(clf, train_data, target, cv=k_fold, n_jobs=1, scoring=scoring)
        print(score)
```

```
In [ ]: #Learning_rates = [0.05, 0.1, 0.25, 0.5, 0.75, 1]
clf = [KNeighborsClassifier(n_neighbors = 13),DecisionTreeClassifier(),
       RandomForestClassifier(n_estimators=13),GaussianNB(),SVC(),
       ]
def model_fit():
    scoring = 'accuracy'
    for i in range(len(clf)):
        score = cross_val_score(clf[i], X_train, cv=k_fold, n_jobs=1, scoring=scoring)
        print("Score of Model",i,":",round(np.mean(score)))
#     round(np.mean(score)*100,2)
#     print("Score of :\n",score)
model_fit()
```

```
In [ ]: clf1 = SVC()
clf1.fit(train_data, target)
test
test_data = test.drop(['Survived','PassengerId'], axis=1)
prediction = clf1.predict(test_data)
# test_data
```

```
In [ ]: in) Y_pred = gaussian.predict(X_test) acc_gaussian = round(gaussian.score(X_train, Y_train))
```

```
In [ ]: decision_tree.predict(X_test) acc_decision_tree = round(decision_tree.score(X_train, Y_train))
```

```
In [ ]: results = pd.DataFrame({
    'Model': ['Support Vector Machines', 'KNN', 'Logistic Regression',
             'Random Forest', 'Naive Bayes', 'Perceptron',
             'Stochastic Gradient Decent',
             'Decision Tree'],
    'Score': [acc_linear_svc, acc_knn, acc_log,
             acc_random_forest, acc_gaussian, acc_perceptron,
             acc_sgd, acc_decision_tree]})
result_df = results.sort_values(by='Score', ascending=False)
result_df = result_df.set_index('Score')
result_df.head(9)
```

```
In [ ]:
```