# DSP Kaggle Project II
## Titanic: Machine Learning from Disaster

### Group 2

Brian Lai
Kirill Evseev
Alejandra Bermudez

California State University, Northridge

June 20, 2019

# Hello Data

# A look at the Data

```
In [20]:  # Mean: survival rate
          # Count: total observations
          # Sum: people survived
          # Sex: 0 female, 1 male
          data[['Sex', 'Survived']].groupby(['Sex'], as_index=False).agg(['mean', 'count', 'sum'])
```

Out[20]:

| | Survived | | |
|---|---|---|---|
| | mean | count | sum |
| Sex | | | |
| 0 | 0.754789 | 261 | 197 |
| 1 | 0.205298 | 453 | 93 |

# Another Look at the Data

```python
# survival rate per age
group_by_age = pd.cut(train["Age"], np.arange(0, 90, 10))
age_grouping = train.groupby(group_by_age).mean()
age_grouping['Survived'].plot.bar()
plt.ylabel('Survival Rate')
plt.grid(c="lightblue")
```

# How we Cleaned the Data

## First Attempt

```
In [12]:   1  # dropping the data that does not contribute
           2  data = train.drop(columns=['Cabin','Embarked','Name','Ticket','Fare'],axis=0)
```

```
[Feature]
Cabin           687
Age             177
Embarked          2
Fare              0
Ticket            0
Parch             0
SibSp             0
Sex               0
Name              0
Pclass            0
Survived          0
PassengerId       0
dtype: int64
```

# Cleaning the Data

```
1   # Cleaning the data: first attempt
```

```
In [70]:   1   train3=pd.DataFrame(train)
```

```
In [71]:   1   train3.drop(['Name', 'SibSp', 'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'], axis=1, inplace=True)
```

```
In [72]:   1   train3=train3.dropna(axis=0)
```

```
In [73]:   1   train3['Age']=np.where(train3['Age'].between(0,10),0,train3['Age'])
           2   train3['Age']=np.where(train3['Age'].between(10.5,30),1,train3['Age'])
           3   train3['Age']=np.where(train3['Age'].between(30.5,60),2,train3['Age'])
           4   train3['Age']=np.where(train3['Age'].between(60.5,100),3,train3['Age'])
```

```
In [75]:   1   train3=pd.get_dummies(train3,columns=['Sex'])
           2   train3=pd.get_dummies(train3,columns=['Age'])
           3   train3=pd.get_dummies(train3,columns=['Pclass'])
```

# Cleaned Data

In [88]: `train3[:10]`

Out[88]:

|    | PassengerId | Survived | Sex_female | Sex_male | Age_0.0 | Age_1.0 | Age_2.0 | Age_3.0 | Pclass_1 | Pclass_2 | Pclass_3 |
|----|-------------|----------|------------|----------|---------|---------|---------|---------|----------|----------|----------|
| 0  | 1           | 0        | 0          | 1        | 0       | 1       | 0       | 0       | 0        | 0        | 1        |
| 1  | 2           | 1        | 1          | 0        | 0       | 0       | 1       | 0       | 1        | 0        | 0        |
| 2  | 3           | 1        | 1          | 0        | 0       | 1       | 0       | 0       | 0        | 0        | 1        |
| 3  | 4           | 1        | 1          | 0        | 0       | 0       | 1       | 0       | 1        | 0        | 0        |
| 4  | 5           | 0        | 0          | 1        | 0       | 0       | 1       | 0       | 0        | 0        | 1        |
| 6  | 7           | 0        | 0          | 1        | 0       | 0       | 1       | 0       | 1        | 0        | 0        |
| 7  | 8           | 0        | 0          | 1        | 1       | 0       | 0       | 0       | 0        | 0        | 1        |
| 8  | 9           | 1        | 1          | 0        | 0       | 1       | 0       | 0       | 0        | 0        | 1        |
| 9  | 10          | 1        | 1          | 0        | 0       | 1       | 0       | 0       | 0        | 1        | 0        |
| 10 | 11          | 1        | 1          | 0        | 1       | 0       | 0       | 0       | 0        | 0        | 1        |

# Second Cleaning Attempt (SCA)

```
In [30]:   1  tr4=pd.DataFrame(train) #going with tr4 instead of train4
```

```
In [31]:   1  tr4.drop(['Ticket', 'Fare', 'Cabin', 'Embarked'], axis=1, inplace=True)
```

```
In [32]:   1  data=[tr4]
           2  titles = {"Mr": 1, "Miss": 2, "Mrs": 3, "Master": 4, "Rare": 5}
           3  for dataset in data:
           4      dataset['Title'] = dataset.Name.str.extract(' ([A-Za-z]+)\.', expand=False)
           5      dataset['Title'] = dataset['Title'].replace(['Lady', 'Countess','Capt', 'Col','Don', 'Dr',\
           6                                                   'Major', 'Rev', 'Sir', 'Jonkheer', 'Dona'], 'Rare')
           7      dataset['Title'] = dataset['Title'].replace('Mlle', 'Miss')
           8      dataset['Title'] = dataset['Title'].replace('Ms', 'Miss')
           9      dataset['Title'] = dataset['Title'].replace('Mme', 'Mrs')
          10      dataset['Title'] = dataset['Title'].map(titles)
          11      dataset['Title'] = dataset['Title'].fillna(0)
          12  tr4 = tr4.drop(['Name'], axis=1)
          13  #taken from https://towardsdatascience.com/predicting-the-survival-of-titanic-passengers-30870ccc7e8
```

```
In [33]:   1  titleNull=tr4[tr4['Title']==0] #there is nobody with a null title
           2  len(titleNull)
```

```
Out[33]:  0
```

Trying to find out the mean age of those who are possibly children or adults with the thinking spouse/sib<1 would mean siblings and 0<parents/children<=2 would probably mean parents if the first was satisfied

```
In [34]:   1  tr4pc=tr4[tr4['Parch']<=2] #tr4pc is possible child
           2  tr4pc=tr4pc[tr4pc['Parch']>0] #trying no parents when # is in front
           3  tr4pc=tr4pc[tr4pc['SibSp']>1]
           4  #last might mess things up since the assumption is that spouse with no kids but it could be someone with only one sibling
           5  len(tr4pc)
```

```
Out[34]:  55
```

```
In [35]:   1  tr4pa=pd.concat([tr4,tr4pc,tr4pc]).drop_duplicates(keep=False)#tr4pa is possible adult
           2  len(tr4pa)
```

```
 3  meanAgeMr=mr["Age"].mean()
 4  medAgeMr=mr["Age"].median()
 5  print(meanAgeMr,medAgeMr,len(mr))
 6
 7  miss=tr4LA[tr4LA['Title']==2]
 8  #miss=miss[miss['Parch']==0] #likely older
 9  meanAgeMs=miss["Age"].mean()
10  medAgeMs=miss["Age"].median()
11  print(meanAgeMs,medAgeMs,len(miss))
12
13  mrs=tr4LA[tr4LA['Title']==3]
14  meanAgeMrs=mrs["Age"].mean()
15  medAgeMrs=mrs["Age"].median()
16  print(meanAgeMrs,medAgeMrs,len(mrs))
17
18  rare=tr4LA[tr4LA['Title']==5]
19  meanAgeRare=rare["Age"].mean()
20  medAgeRare=rare["Age"].median()
21  print(meanAgeRare,medAgeRare,len(rare))
```

```
32.56106870229008 30.0 509
27.68617021276596 25.5 121
35.898148148148145 35.0 125
45.54545454545455 48.5 23
```
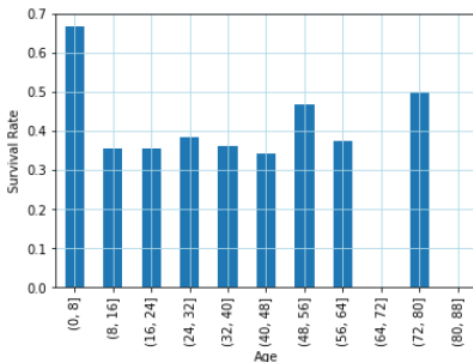
# Making an Informed Guess (MIG)

In [45]:
```
1  #Replace age with mean of each larger group rounded to the near
2  tr4LCh=tr4LCh.fillna(9.5)
3  tr4LA=tr4LA.fillna(33)
```

# Resulting Age Distribution (RAD)

```
# Survival rate per age of combined
group_by_age = pd.cut(train_combine["Age"], np.arange(0, 90, 8))
age_grouping = train_combine.groupby(group_by_age).mean()
age_grouping['Survived'].plot.bar()
plt.grid(c="lightblue")
plt.ylabel('Survival Rate')
```
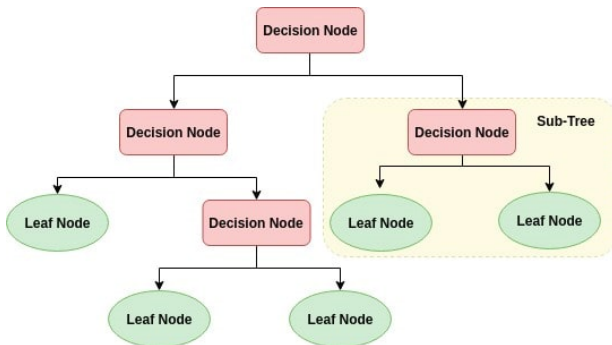
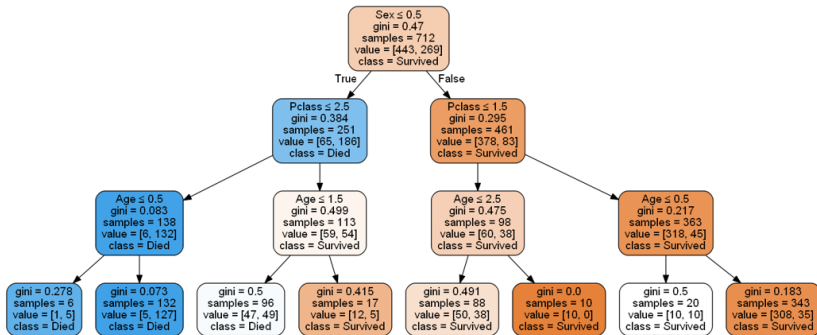Text(0, 0.5, 'Survival Rate')

# What is our Model?

What exactly is a decision tree?

- A Decision Tree is a support tool that uses a flowchart-like model of decisions and potential consequences
- A supervised learning algorithm that works with both continuous and categorical variables
- It is a tool that only contains conditional control statements
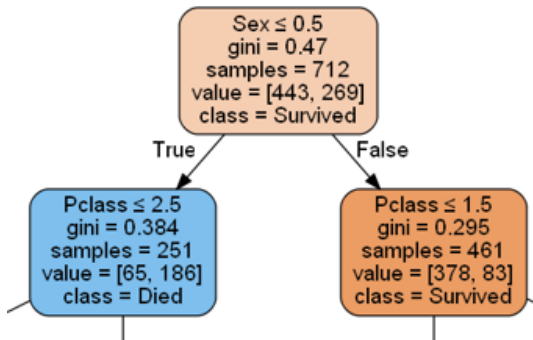- Each leaf node is a class label which is the outcome

# Model: Decision Tree

- Each branch is the outcome of the test
- It's rarely ever balanced
- A depth of 3:
  Sex, Pclass, and Age

# Model: Decision Tree

- Most important factor is Sex, it splits off to determine survival/death outcome
- Blue boxes = most likely to die
- Orange boxes = most likely to survive

# Fitting the Model

```
In [78]:   1  from sklearn import tree
           2  from sklearn.model_selection import train_test_split
           3  from sklearn.metrics import accuracy_score
```

```
In [79]:   1  X=train3.values
           2  Y=train3['Survived'].values
```

```
In [80]:   1  X=np.delete(X,1,axis=1)
```

```
In [81]:   1  X=np.delete(X,0,axis=1)
```

```
In [85]:   1  XTRAIN, XTEST, YTRAIN, YTEST=train_test_split(X,Y)
           2  nsplit=1000
           3  depth=1
```

```
In [87]:   1  while (depth < 7):
           2      errs=[]
           3      for j in range(nsplit):
           4          XTRAIN, XTEST, YTRAIN, YTEST=train_test_split(X,Y)
           5          DT=tree.DecisionTreeClassifier(max_depth=depth)
           6          DT.fit(XTRAIN,YTRAIN)
           7          YP=DT.predict(XTEST)
           8          errs.append(1-accuracy_score(YTEST,YP))
           9      print("Decision Tree Depth = %d mean error = %7.6f SD=%7.6f"\
          10          %(depth,np.mean(errs),np.std(errs)))
          11      depth = depth + 1
```
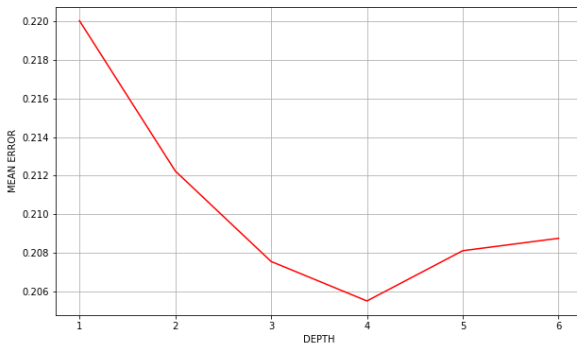
```
Decision Tree Depth = 1 mean error = 0.220544 SD=0.026738
Decision Tree Depth = 2 mean error = 0.213799 SD=0.026424
Decision Tree Depth = 3 mean error = 0.209737 SD=0.026049
Decision Tree Depth = 4 mean error = 0.204296 SD=0.025150
Decision Tree Depth = 5 mean error = 0.209514 SD=0.024734
Decision Tree Depth = 6 mean error = 0.209944 SD=0.024799
```

# Determining the Optimal Depth

```
Decision Tree Depth = 1 mean error = 0.220028 SD=0.026047
Decision Tree Depth = 2 mean error = 0.212240 SD=0.027499
Decision Tree Depth = 3 mean error = 0.207553 SD=0.026333
Decision Tree Depth = 4 mean error = 0.205508 SD=0.024900
Decision Tree Depth = 5 mean error = 0.208112 SD=0.025645
Decision Tree Depth = 6 mean error = 0.208754 SD=0.025076
```

In [99]:
```python
plt.plot((range(1,7)), errors, color='red')
plt.xlabel("DEPTH")
plt.ylabel("MEAN ERROR")
plt.grid()
plt.gcf().set_size_inches(10,6)
```

# Results: Initial Attempt

From our first attempt:
Score $= 0.55980$

# Results: Second Attempt

From our second attempt:
Score = 0.78468

An increase of 0.22488!

# References

https://towardsdatascience.com/predicting-the-survival-of-titanic-passengers-30870ccc7e8?gi=339b274bc177

https://www.datacamp.com/community/tutorials/decision-tree-classification-python

https://blog.patricktriest.com/titanic-machine-learning-in-python/