In [18]:

```python
import pandas as pd
```

In [20]:

```python
train=pd.read_csv('train.csv')
```

In [21]:

```python
train
```

Out[21]:

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S |
| 1 | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C |
| 2 | 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S |
| 3 | 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | S |
| 4 | 5 | 0 | 3 | Allen, Mr. William Henry | male | 35.0 | 0 | 0 | 373450 | 8.0500 | NaN | S |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 886 | 887 | 0 | 2 | Montvila, Rev. Juozas | male | 27.0 | 0 | 0 | 211536 | 13.0000 | NaN | S |
| 887 | 888 | 1 | 1 | Graham, Miss. Margaret Edith | female | 19.0 | 0 | 0 | 112053 | 30.0000 | B42 | S |
| 888 | 889 | 0 | 3 | Johnston, Miss. Catherine Helen "Carrie" | female | NaN | 1 | 2 | W./C. 6607 | 23.4500 | NaN | S |
| 889 | 890 | 1 | 1 | Behr, Mr. Karl Howell | male | 26.0 | 0 | 0 | 111369 | 30.0000 | C148 | C |
| 890 | 891 | 0 | 3 | Dooley, Mr. Patrick | male | 32.0 | 0 | 0 | 370376 | 7.7500 | NaN | Q |

891 rows × 12 columns

In [22]:

```python
test=pd.read_csv('test.csv')
```

In [23]:

```python
test
```

Out[23]:

| | PassengerId | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 892 | 3 | Kelly, Mr. James | male | 34.5 | 0 | 0 | 330911 | 7.8292 | NaN | Q |
| 1 | 893 | 3 | Wilkes, Mrs. James (Ellen Needs) | female | 47.0 | 1 | 0 | 363272 | 7.0000 | NaN | S |
| 2 | 894 | 2 | Myles, Mr. Thomas Francis | male | 62.0 | 0 | 0 | 240276 | 9.6875 | NaN | Q |
| 3 | 895 | 3 | Wirz, Mr. Albert | male | 27.0 | 0 | 0 | 315154 | 8.6625 | NaN | S |
| 4 | 896 | 3 | Hirvonen, Mrs. Alexander (Helga E Lindqvist) | female | 22.0 | 1 | 1 | 3101298 | 12.2875 | NaN | S |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 413 | 1305 | 3 | Spector, Mr. Woolf | male | NaN | 0 | 0 | A.5. 3236 | 8.0500 | NaN | S |
| 414 | 1306 | 1 | Oliva y Ocana, Dona. Fermina | female | 39.0 | 0 | 0 | PC 17758 | 108.9000 | C105 | C |
| 415 | 1307 | 3 | Saether, Mr. Simon Sivertsen | male | 38.5 | 0 | 0 | SOTON/O.Q. 3101262 | 7.2500 | NaN | S |
| 416 | 1308 | 3 | Ware, Mr. Frederick | male | NaN | 0 | 0 | 359309 | 8.0500 | NaN | S |
| 417 | 1309 | 3 | Peter, Master. Michael J | male | NaN | 1 | 1 | 2668 | 22.3583 | NaN | C |

418 rows × 11 columns

In [24]:

```python
train.head(5)
```

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S |
| **1** | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C |
| **2** | 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S |
| **3** | 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | S |
| **4** | 5 | 0 | 3 | Allen, Mr. William Henry | male | 35.0 | 0 | 0 | 373450 | 8.0500 | NaN | S |

In [25]:

```
train.shape
```

Out[25]:

```
(891, 12)
```

In [26]:

```
test.shape
```

Out[26]:

```
(418, 11)
```

In [28]:

```
train.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 [29]:

```
test.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

In [30]:

```
train.isnull().sum()
```

Out[30]:

```
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
```
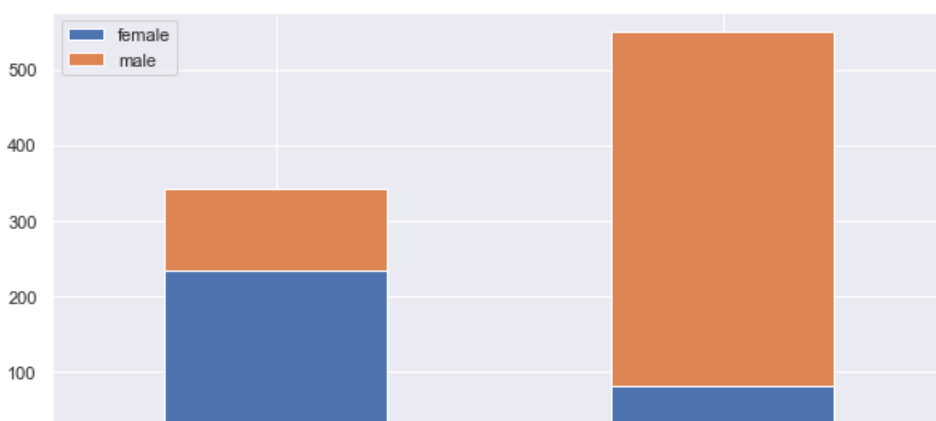
In [31]:

```
test.isnull().sum()
```

Out[31]:

```
PassengerId     0
Pclass          0
Name            0
Sex             0
Age            86
SibSp           0
Parch           0
Ticket          0
Fare            1
Cabin         327
Embarked        0
dtype: int64
```

In [32]:

```
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
sns.set()
```

In [33]:

```
def bar_chart(feature):
    survived = train[train['Survived']==1][feature].value_counts()
    dead = train[train['Survived']==0][feature].value_counts()
    df = pd.DataFrame([survived,dead])
    df.index = ['Survived','Dead']
    df.plot(kind='bar',stacked=True, figsize=(10,5))
```
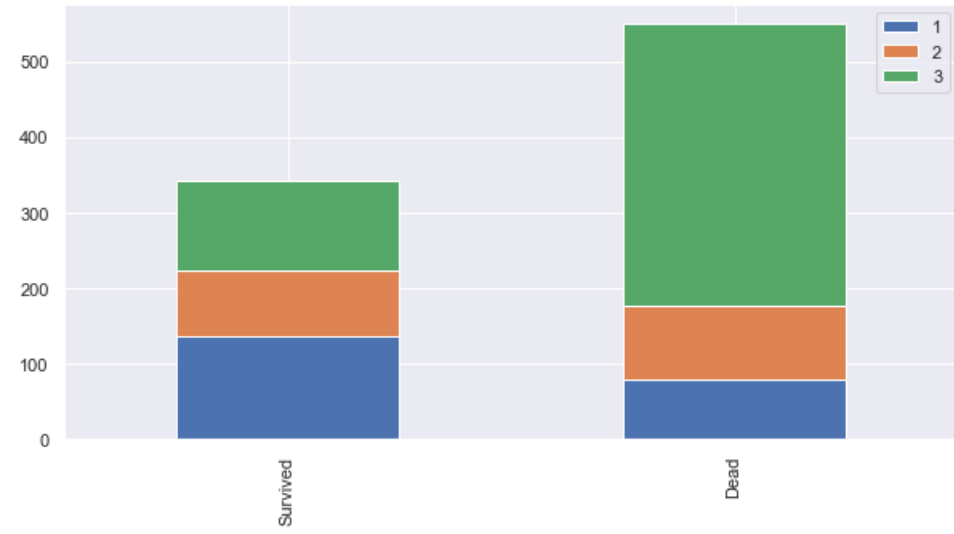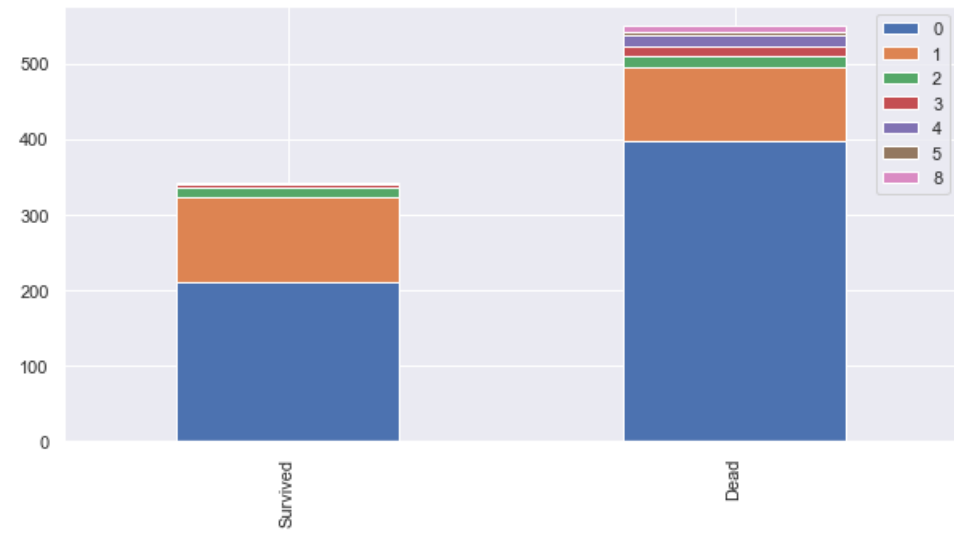
In [35]:

```
bar_chart('Sex')
```
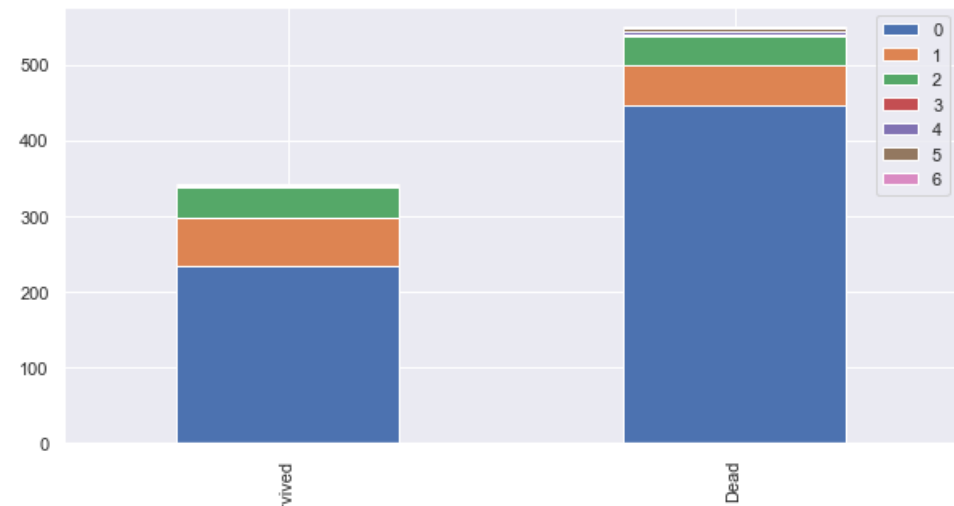
In [37]:

```
bar_chart('Pclass')
```



In [38]:

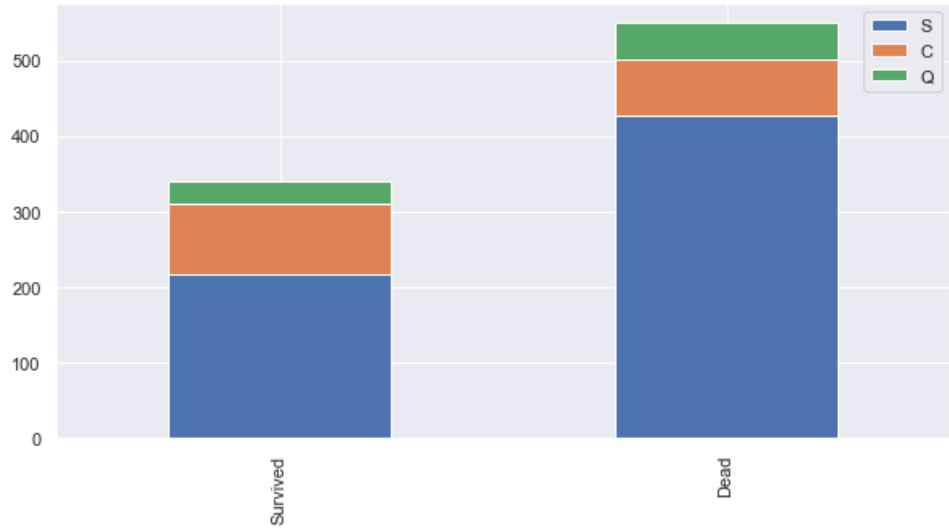```
bar_chart('SibSp')
```



In [39]:

```
bar_chart('Parch')
```

In [40]:

```
bar_chart('Embarked')
```



In [41]:

```
train_test_data = [train, test]

for dataset in train_test_data:
    dataset['Title'] = dataset['Name'].str.extract(' ([A-Za-z]+)\.', expand=False)
```

In [42]:

```
train['Title'].value_counts()
```

Out[42]:

```
Mr          517
Miss        182
Mrs         125
Master       40
Dr            7
Rev           6
Col           2
Mlle          2
Major         2
Capt          1
Ms            1
Countess      1
Lady          1
Don           1
Jonkheer      1
Sir           1
Mme           1
Name: Title, dtype: int64
```

In [43]:

```
test['Title'].value_counts()
```

Out[43]:

```
Mr        240
Miss       78
Mrs        72
Master     21
Col         2
Rev         2
Ms          1
Dr          1
Dona        1
Name: Title, dtype: int64
```

```
title_mapping = {"Mr": 0, "Miss": 1, "Mrs": 2,
        "Master": 3, "Dr": 3, "Rev": 3, "Col": 3, "Major": 3, "Mlle": 3,"Countess": 3,
        "Ms": 3, "Lady": 3, "Jonkheer": 3, "Don": 3, "Dona" : 3, "Mme": 3,"Capt": 3,"Sir": 3 }
for dataset in train_test_data:
    dataset['Title'] = dataset['Title'].map(title_mapping)
```

In [45]:

```
train.head()
```

Out[45]:

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | Title |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S | 0 |
| 1 | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C | 2 |
| 2 | 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S | 1 |
| 3 | 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | S | 2 |
| 4 | 5 | 0 | 3 | Allen, Mr. William Henry | male | 35.0 | 0 | 0 | 373450 | 8.0500 | NaN | S | 0 |

In [46]:

```
test.head()
```

Out[46]:

| | PassengerId | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | Title |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 892 | 3 | Kelly, Mr. James | male | 34.5 | 0 | 0 | 330911 | 7.8292 | NaN | Q | 0 |
| 1 | 893 | 3 | Wilkes, Mrs. James (Ellen Needs) | female | 47.0 | 1 | 0 | 363272 | 7.0000 | NaN | S | 2 |
| 2 | 894 | 2 | Myles, Mr. Thomas Francis | male | 62.0 | 0 | 0 | 240276 | 9.6875 | NaN | Q | 0 |
| 3 | 895 | 3 | Wirz, Mr. Albert | male | 27.0 | 0 | 0 | 315154 | 8.6625 | NaN | S | 0 |
| 4 | 896 | 3 | Hirvonen, Mrs. Alexander (Helga E Lindqvist) | female | 22.0 | 1 | 1 | 3101298 | 12.2875 | NaN | S | 2 |

In [47]:

```
bar_chart('Title')
```



In [48]:

```
# delete unnecessary feature from dataset
train.drop('Name', axis=1, inplace=True)
test.drop('Name', axis=1, inplace=True)
```

In [49]:

```
train.head()
```

Out[49]:

| | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | Title |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S | 0 |
| 1 | 2 | 1 | 1 | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C | 2 |
| 2 | 3 | 1 | 3 | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S | 1 |
| 3 | 4 | 1 | 1 | female | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | S | 2 |
| 4 | 5 | 0 | 3 | male | 35.0 | 0 | 0 | 373450 | 8.0500 | NaN | S | 0 |

In [50]:

```
test.head()
```

Out[50]:

| | PassengerId | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | Title |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 892 | 3 | male | 34.5 | 0 | 0 | 330911 | 7.8292 | NaN | Q | 0 |
| 1 | 893 | 3 | female | 47.0 | 1 | 0 | 363272 | 7.0000 | NaN | S | 2 |
| 2 | 894 | 2 | male | 62.0 | 0 | 0 | 240276 | 9.6875 | NaN | Q | 0 |
| 3 | 895 | 3 | male | 27.0 | 0 | 0 | 315154 | 8.6625 | NaN | S | 0 |
| 4 | 896 | 3 | female | 22.0 | 1 | 1 | 3101298 | 12.2875 | NaN | S | 2 |

In [51]:

```
sex_mapping = {"male": 0, "female": 1}
for dataset in train_test_data:
    dataset['Sex'] = dataset['Sex'].map(sex_mapping)
```

In [52]:

```
bar_chart('Sex')
```



In [53]:

```
train.head()
```

Out[53]:

| | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | Title |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | 0 | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S | 0 |
| 1 | 2 | 1 | 1 | 1 | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C | 2 |
| 2 | 3 | 1 | 3 | 1 | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S | 1 |
| 3 | 4 | 1 | 1 | 1 | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | S | 2 |

In [54]:

```
# fill missing age with median age for each title (Mr, Mrs, Miss, Others)
train["Age"].fillna(train.groupby("Title")["Age"].transform("median"), inplace=True)
test["Age"].fillna(test.groupby("Title")["Age"].transform("median"), inplace=True)
```

In [55]:

```
train.groupby("Title")["Age"].transform("median")
train.head()
```

Out[55]:

| | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | Title |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | 0 | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S | 0 |
| 1 | 2 | 1 | 1 | 1 | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C | 2 |
| 2 | 3 | 1 | 3 | 1 | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S | 1 |
| 3 | 4 | 1 | 1 | 1 | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | S | 2 |
| 4 | 5 | 0 | 3 | 0 | 35.0 | 0 | 0 | 373450 | 8.0500 | NaN | S | 0 |

In [56]:

```
test.groupby("Title")["Age"].transform("median")
test.head()
```

Out[56]:

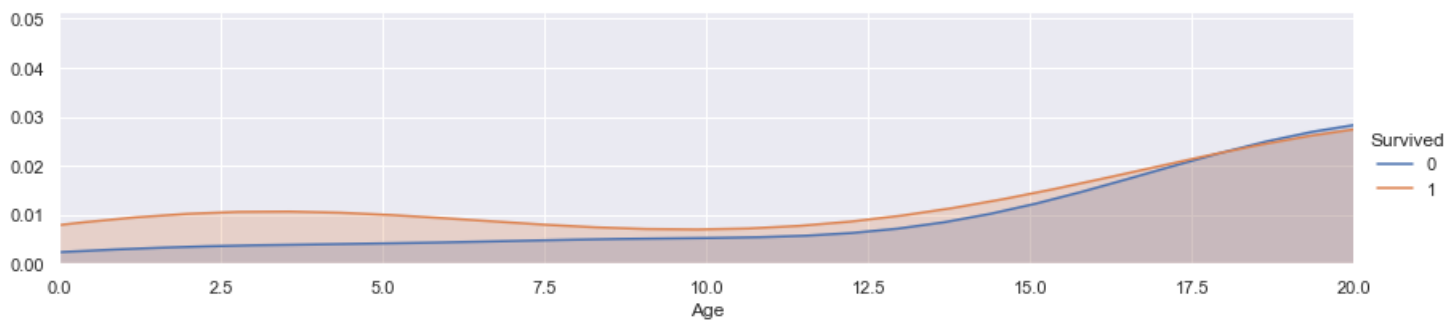| | PassengerId | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | Title |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 892 | 3 | 0 | 34.5 | 0 | 0 | 330911 | 7.8292 | NaN | Q | 0 |
| 1 | 893 | 3 | 1 | 47.0 | 1 | 0 | 363272 | 7.0000 | NaN | S | 2 |
| 2 | 894 | 2 | 0 | 62.0 | 0 | 0 | 240276 | 9.6875 | NaN | Q | 0 |
| 3 | 895 | 3 | 0 | 27.0 | 0 | 0 | 315154 | 8.6625 | NaN | S | 0 |
| 4 | 896 | 3 | 1 | 22.0 | 1 | 1 | 3101298 | 12.2875 | NaN | S | 2 |

In [57]:

```
facet = sns.FacetGrid(train, hue="Survived",aspect=4)
facet.map(sns.kdeplot,'Age',shade= True)
facet.set(xlim=(0, train['Age'].max()))
facet.add_legend()

plt.show()
```



In [58]:

```
facet = sns.FacetGrid(train, hue="Survived",aspect=4)
facet.map(sns.kdeplot,'Age',shade= True)
facet.set(xlim=(0, train['Age'].max()))
facet.add_legend()
plt.xlim(0, 20)
```
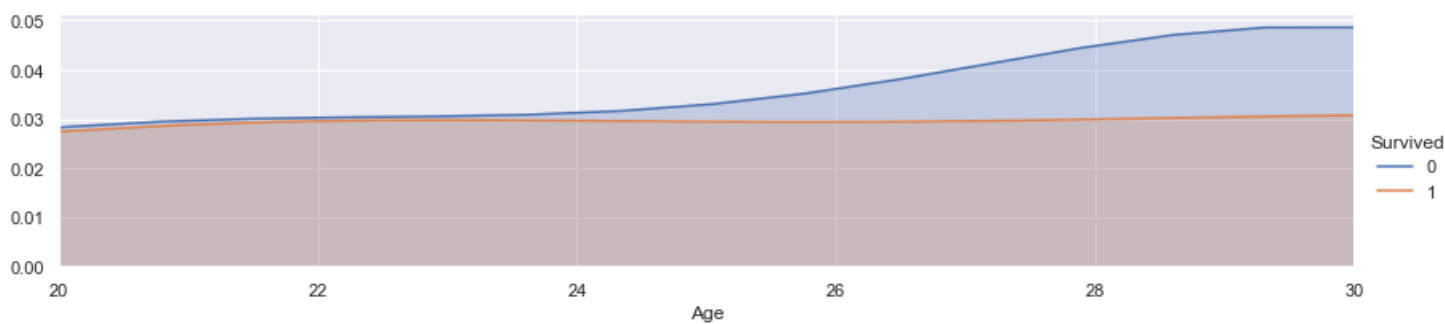
Out[58]:

(0.0, 20.0)

```
facet = sns.FacetGrid(train, hue="Survived",aspect=4)
facet.map(sns.kdeplot,'Age',shade= True)
facet.set(xlim=(0, train['Age'].max()))
facet.add_legend()
plt.xlim(20, 30)
```
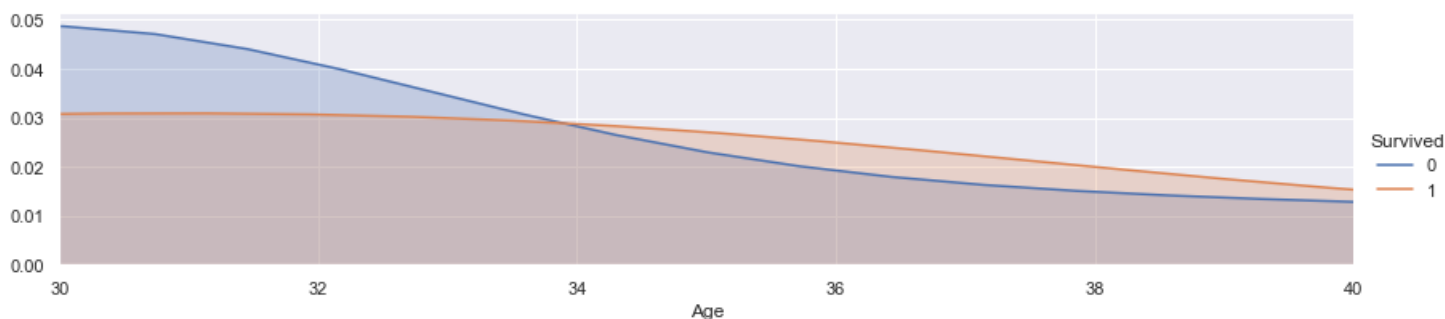
Out[59]:

(20.0, 30.0)



In [60]:

```
facet = sns.FacetGrid(train, hue="Survived",aspect=4)
facet.map(sns.kdeplot,'Age',shade= True)
facet.set(xlim=(0, train['Age'].max()))
facet.add_legend()
plt.xlim(30, 40)
```
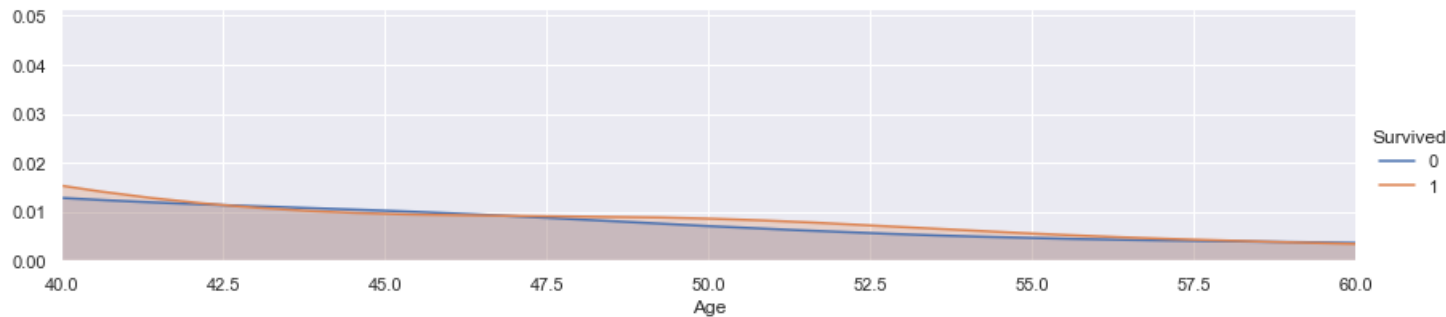
Out[60]:

(30.0, 40.0)



In [61]:

```
facet = sns.FacetGrid(train, hue="Survived",aspect=4)
facet.map(sns.kdeplot,'Age',shade= True)
facet.set(xlim=(0, train['Age'].max()))
facet.add_legend()
plt.xlim(40, 60)
```

Out[61]:

(40.0, 60.0)

```
train.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   Sex          891 non-null    int64
 4   Age          891 non-null    float64
 5   SibSp        891 non-null    int64
 6   Parch        891 non-null    int64
 7   Ticket       891 non-null    object
 8   Fare         891 non-null    float64
 9   Cabin        204 non-null    object
 10  Embarked     889 non-null    object
 11  Title        891 non-null    int64
dtypes: float64(2), int64(7), object(3)
memory usage: 83.7+ KB
```

```
test.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   Sex          418 non-null    int64
 3   Age          418 non-null    float64
 4   SibSp        418 non-null    int64
 5   Parch        418 non-null    int64
 6   Ticket       418 non-null    object
 7   Fare         417 non-null    float64
 8   Cabin        91 non-null     object
 9   Embarked     418 non-null    object
 10  Title        418 non-null    int64
dtypes: float64(2), int64(6), object(3)
memory usage: 36.0+ KB
```

```
for dataset in train_test_data:
    dataset.loc[ dataset['Age'] <= 16, 'Age'] = 0,
    dataset.loc[(dataset['Age'] > 16) & (dataset['Age'] <= 26), 'Age'] = 1,
    dataset.loc[(dataset['Age'] > 26) & (dataset['Age'] <= 36), 'Age'] = 2,
    dataset.loc[(dataset['Age'] > 36) & (dataset['Age'] <= 62), 'Age'] = 3,
    dataset.loc[ dataset['Age'] > 62, 'Age'] = 4
```
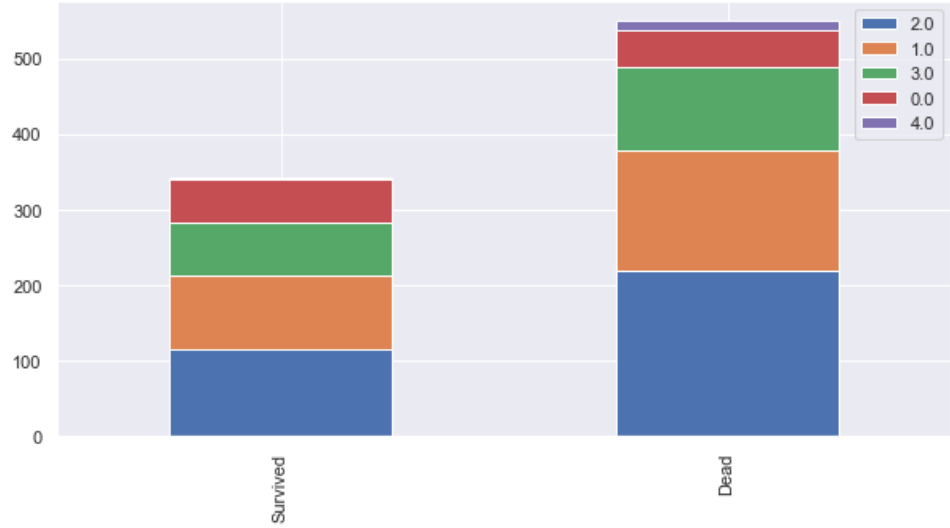
```
train.head()
```

| PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | Title |
|---|---|---|---|---|---|---|---|---|---|---|---|

| | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | Title |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | 0 | 1.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S | 0 |
| 1 | 2 | 1 | 1 | 1 | 3.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C | 2 |
| 2 | 3 | 1 | 3 | 1 | 1.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S | 1 |
| 3 | 4 | 1 | 1 | 1 | 2.0 | 1 | 0 | 113803 | 53.1000 | C123 | S | 2 |
| 4 | 5 | 0 | 3 | 0 | 2.0 | 0 | 0 | 373450 | 8.0500 | NaN | S | 0 |

In [66]:
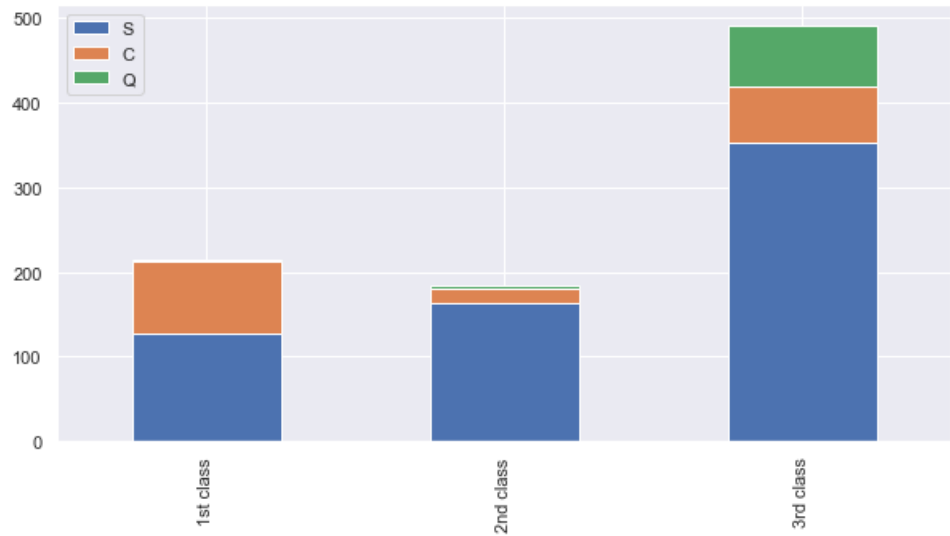
```
bar_chart('Age')
```



In [67]:

```
Pclass1 = train[train['Pclass']==1]['Embarked'].value_counts()
Pclass2 = train[train['Pclass']==2]['Embarked'].value_counts()
Pclass3 = train[train['Pclass']==3]['Embarked'].value_counts()
df = pd.DataFrame([Pclass1, Pclass2, Pclass3])
df.index = ['1st class','2nd class', '3rd class']
df.plot(kind='bar',stacked=True, figsize=(10,5))
```

Out[67]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x29dbc7cfeb0>
```



In [68]:

```
for dataset in train_test_data:
    dataset['Embarked'] = dataset['Embarked'].fillna('S')
```

In [69]:

```
train.head()
```

Out[69]:

| | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | Title |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | 0 | 1.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S | 0 |
| 1 | 2 | 1 | 1 | 1 | 3.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C | 2 |
| 2 | 3 | 1 | 3 | 1 | 1.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S | 1 |
| 3 | 4 | 1 | 1 | 1 | 2.0 | 1 | 0 | 113803 | 53.1000 | C123 | S | 2 |
| 4 | 5 | 0 | 3 | 0 | 2.0 | 0 | 0 | 373450 | 8.0500 | NaN | S | 0 |

In [70]:

```
embarked_mapping = {"S": 0, "C": 1, "Q": 2}
for dataset in train_test_data:
    dataset['Embarked'] = dataset['Embarked'].map(embarked_mapping)
```

In [71]:

```
train.head()
```

Out[71]:

| | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | Title |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | 0 | 1.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | 0 | 0 |
| 1 | 2 | 1 | 1 | 1 | 3.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | 1 | 2 |
| 2 | 3 | 1 | 3 | 1 | 1.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | 0 | 1 |
| 3 | 4 | 1 | 1 | 1 | 2.0 | 1 | 0 | 113803 | 53.1000 | C123 | 0 | 2 |
| 4 | 5 | 0 | 3 | 0 | 2.0 | 0 | 0 | 373450 | 8.0500 | NaN | 0 | 0 |

In [72]:

```
# fill missing Fare with median fare for each Pclass
train["Fare"].fillna(train.groupby("Pclass")["Fare"].transform("median"), inplace=True)
test["Fare"].fillna(test.groupby("Pclass")["Fare"].transform("median"), inplace=True)
train.head(5)
```
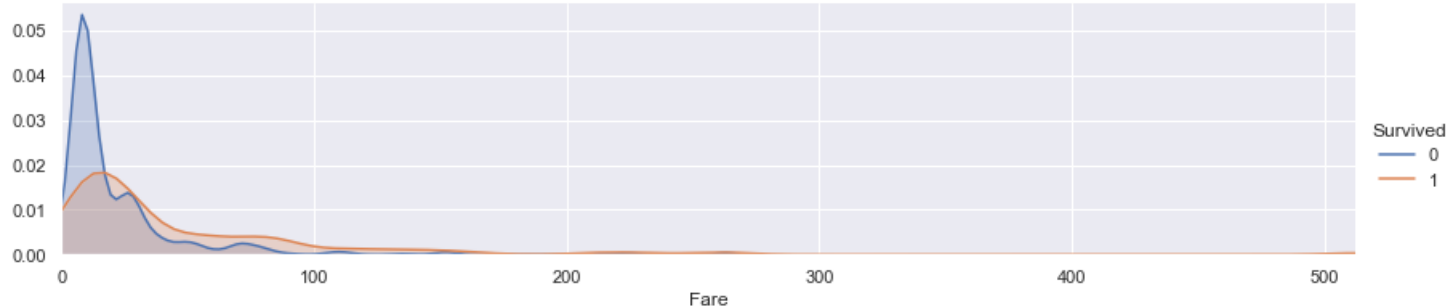
Out[72]:

| | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | Title |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | 0 | 1.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | 0 | 0 |
| 1 | 2 | 1 | 1 | 1 | 3.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | 1 | 2 |
| 2 | 3 | 1 | 3 | 1 | 1.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | 0 | 1 |
| 3 | 4 | 1 | 1 | 1 | 2.0 | 1 | 0 | 113803 | 53.1000 | C123 | 0 | 2 |
| 4 | 5 | 0 | 3 | 0 | 2.0 | 0 | 0 | 373450 | 8.0500 | NaN | 0 | 0 |

In [73]:

```
facet = sns.FacetGrid(train, hue="Survived",aspect=4)
facet.map(sns.kdeplot,'Fare',shade= True)
facet.set(xlim=(0, train['Fare'].max()))
facet.add_legend()

plt.show()
```

```
facet = sns.FacetGrid(train, hue="Survived",aspect=4)
facet.map(sns.kdeplot,'Fare',shade= True)
facet.set(xlim=(0, train['Fare'].max()))
facet.add_legend()
plt.xlim(0, 20)
```
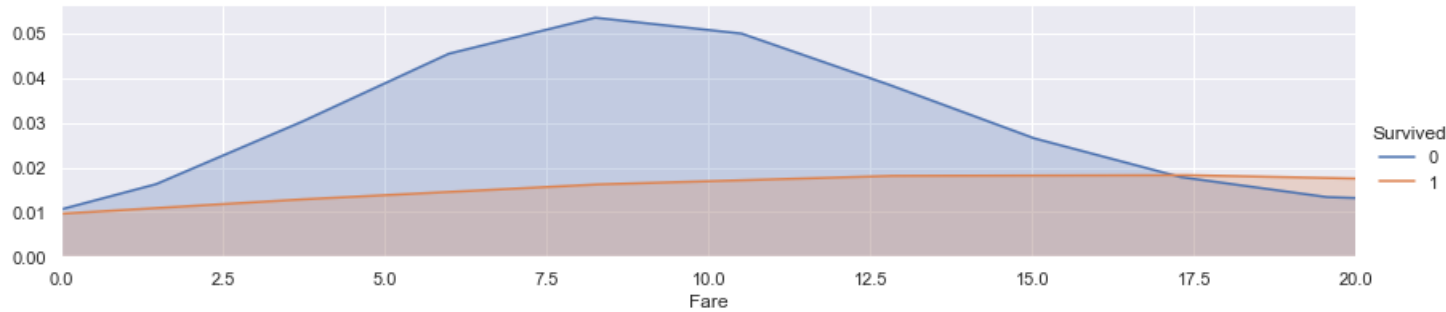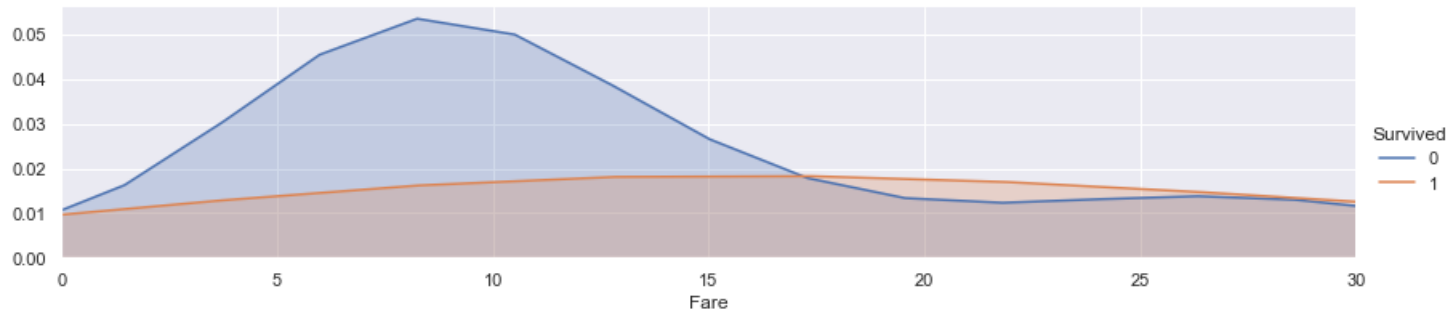
Out[74]:

(0.0, 20.0)



In [75]:

```
facet = sns.FacetGrid(train, hue="Survived",aspect=4)
facet.map(sns.kdeplot,'Fare',shade= True)
facet.set(xlim=(0, train['Fare'].max()))
facet.add_legend()
plt.xlim(0, 30)
```

Out[75]:

(0.0, 30.0)



In [76]:

```
for dataset in train_test_data:
    dataset.loc[ dataset['Fare'] <= 17, 'Fare'] = 0,
    dataset.loc[(dataset['Fare'] > 17) & (dataset['Fare'] <= 30), 'Fare'] = 1,
    dataset.loc[(dataset['Fare'] > 30) & (dataset['Fare'] <= 100), 'Fare'] = 2,
    dataset.loc[ dataset['Fare'] > 100, 'Fare'] = 3
```

In [77]:

```
train.head()
```

Out[77]:

| | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | Title |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | 0 | 1.0 | 1 | 0 | A/5 21171 | 0.0 | NaN | 0 | 0 |
| 1 | 2 | 1 | 1 | 1 | 3.0 | 1 | 0 | PC 17599 | 2.0 | C85 | 1 | 2 |
| 2 | 3 | 1 | 3 | 1 | 1.0 | 0 | 0 | STON/O2. 3101282 | 0.0 | NaN | 0 | 1 |
| 3 | 4 | 1 | 1 | 1 | 2.0 | 1 | 0 | 113803 | 2.0 | C123 | 0 | 2 |
| 4 | 5 | 0 | 3 | 0 | 2.0 | 0 | 0 | 373450 | 0.0 | NaN | 0 | 0 |

In [78]:

```
train.Cabin.value_counts()
```

Out[78]:

```
C23 C25 C27    4
G6             4
B96 B98        4
C22 C26        3
F33            3
              ..
B86            1
E77            1
C47            1
C86            1
A19            1
Name: Cabin, Length: 147, dtype: int64
```
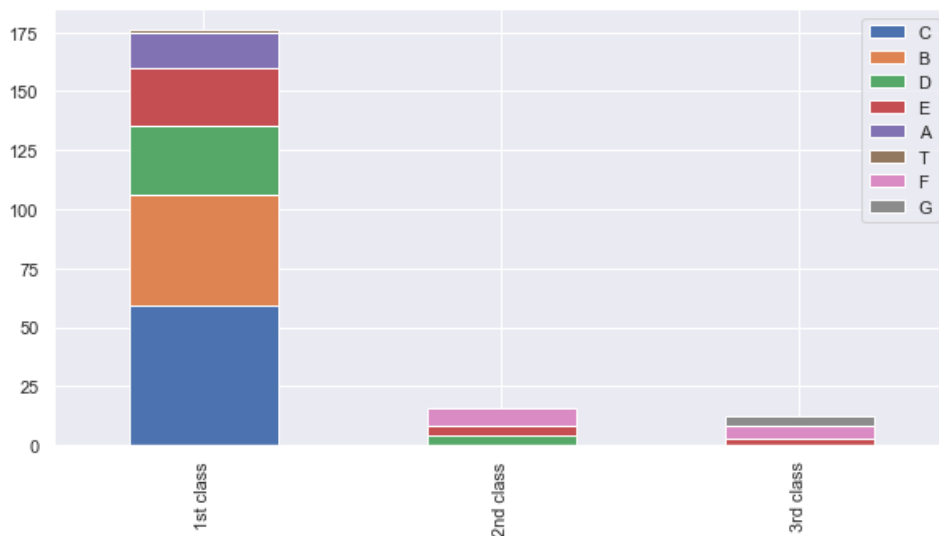
In [79]:

```python
for dataset in train_test_data:
    dataset['Cabin'] = dataset['Cabin'].str[:1]
```

In [80]:

```python
Pclass1 = train[train['Pclass']==1]['Cabin'].value_counts()
Pclass2 = train[train['Pclass']==2]['Cabin'].value_counts()
Pclass3 = train[train['Pclass']==3]['Cabin'].value_counts()
df = pd.DataFrame([Pclass1, Pclass2, Pclass3])
df.index = ['1st class','2nd class', '3rd class']
df.plot(kind='bar',stacked=True, figsize=(10,5))
```

Out[80]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x29dbc9e2a60>
```



In [81]:

```python
cabin_mapping = {"A": 0, "B": 0.4, "C": 0.8, "D": 1.2, "E": 1.6, "F": 2, "G": 2.4, "T": 2.8}
for dataset in train_test_data:
    dataset['Cabin'] = dataset['Cabin'].map(cabin_mapping)
```

In [85]:

```python
#fill missing Fare with median fare for each Pclass
train["Cabin"].fillna(train.groupby("Pclass")["Cabin"].transform("median"), inplace=True)
test["Cabin"].fillna(test.groupby("Pclass")["Cabin"].transform("median"), inplace=True)
```

In [86]:

```python
train.head()
```

Out[86]:

| | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | Title |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | 0 | 1.0 | 1 | 0 | A/5 21171 | 0.0 | 2.0 | 0 | 0 |
| 1 | 2 | 1 | 1 | 1 | 3.0 | 1 | 0 | PC 17599 | 2.0 | 0.8 | 1 | 2 |
| 2 | 3 | 1 | 3 | 1 | 1.0 | 0 | 0 | STON/O2. 3101282 | 0.0 | 2.0 | 0 | 1 |
| 3 | 4 | 1 | 1 | 1 | 2.0 | 1 | 0 | 113803 | 2.0 | 0.8 | 0 | 2 |
| 4 | 5 | 0 | 3 | 0 | 2.0 | 0 | 0 | 373450 | 0.0 | 2.0 | 0 | 0 |

In [87]:

```
test.head()
```

Out[87]:

| | PassengerId | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | Title |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 892 | 3 | 0 | 2.0 | 0 | 0 | 330911 | 0.0 | 2.0 | 2 | 0 |
| 1 | 893 | 3 | 1 | 3.0 | 1 | 0 | 363272 | 0.0 | 2.0 | 0 | 2 |
| 2 | 894 | 2 | 0 | 3.0 | 0 | 0 | 240276 | 0.0 | 2.0 | 2 | 0 |
| 3 | 895 | 3 | 0 | 2.0 | 0 | 0 | 315154 | 0.0 | 2.0 | 0 | 0 |
| 4 | 896 | 3 | 1 | 1.0 | 1 | 1 | 3101298 | 0.0 | 2.0 | 0 | 2 |

In [89]:

```
train["FamilySize"] = train["SibSp"] + train["Parch"] + 1
test["FamilySize"] = test["SibSp"] + test["Parch"] + 1
```
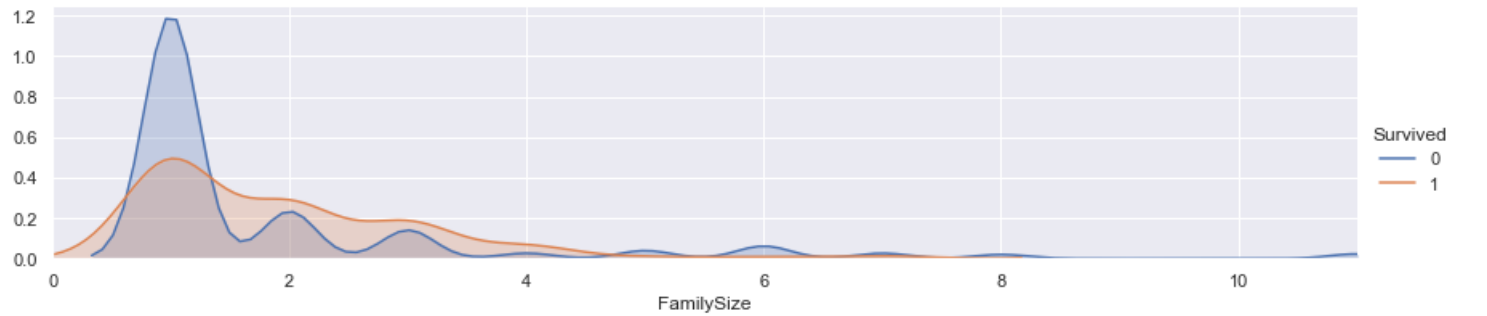
In [90]:

```
train.head()
```

Out[90]:

| | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | Title | FamilySize |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | 0 | 1.0 | 1 | 0 | A/5 21171 | 0.0 | 2.0 | 0 | 0 | 2 |
| 1 | 2 | 1 | 1 | 1 | 3.0 | 1 | 0 | PC 17599 | 2.0 | 0.8 | 1 | 2 | 2 |
| 2 | 3 | 1 | 3 | 1 | 1.0 | 0 | 0 | STON/O2. 3101282 | 0.0 | 2.0 | 0 | 1 | 1 |
| 3 | 4 | 1 | 1 | 1 | 2.0 | 1 | 0 | 113803 | 2.0 | 0.8 | 0 | 2 | 2 |
| 4 | 5 | 0 | 3 | 0 | 2.0 | 0 | 0 | 373450 | 0.0 | 2.0 | 0 | 0 | 1 |

In [91]:

```
facet = sns.FacetGrid(train, hue="Survived",aspect=4)
facet.map(sns.kdeplot,'FamilySize',shade= True)
facet.set(xlim=(0, train['FamilySize'].max()))
facet.add_legend()
plt.xlim(0)
```

Out[91]:

(0.0, 11.0)



In [92]:

```
family_mapping = {1: 0, 2: 0.4, 3: 0.8, 4: 1.2, 5: 1.6, 6: 2, 7: 2.4, 8: 2.8, 9: 3.2, 10: 3.6, 11: 4}
```

```
for dataset in train_test_data:
    dataset['FamilySize'] = dataset['FamilySize'].map(family_mapping)
```

In [93]:

```
train.head()
```

Out[93]:

| | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | Title | FamilySize |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | 0 | 1.0 | 1 | 0 | A/5 21171 | 0.0 | 2.0 | 0 | 0 | 0.4 |
| 1 | 2 | 1 | 1 | 1 | 3.0 | 1 | 0 | PC 17599 | 2.0 | 0.8 | 1 | 2 | 0.4 |
| 2 | 3 | 1 | 3 | 1 | 1.0 | 0 | 0 | STON/O2. 3101282 | 0.0 | 2.0 | 0 | 1 | 0.0 |
| 3 | 4 | 1 | 1 | 1 | 2.0 | 1 | 0 | 113803 | 2.0 | 0.8 | 0 | 2 | 0.4 |
| 4 | 5 | 0 | 3 | 0 | 2.0 | 0 | 0 | 373450 | 0.0 | 2.0 | 0 | 0 | 0.0 |

In [94]:

```
test.head()
```

Out[94]:

| | PassengerId | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | Title | FamilySize |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 892 | 3 | 0 | 2.0 | 0 | 0 | 330911 | 0.0 | 2.0 | 2 | 0 | 0.0 |
| 1 | 893 | 3 | 1 | 3.0 | 1 | 0 | 363272 | 0.0 | 2.0 | 0 | 2 | 0.4 |
| 2 | 894 | 2 | 0 | 3.0 | 0 | 0 | 240276 | 0.0 | 2.0 | 2 | 0 | 0.0 |
| 3 | 895 | 3 | 0 | 2.0 | 0 | 0 | 315154 | 0.0 | 2.0 | 0 | 0 | 0.0 |
| 4 | 896 | 3 | 1 | 1.0 | 1 | 1 | 3101298 | 0.0 | 2.0 | 0 | 2 | 0.8 |

In [95]:

```
features_drop = ['Ticket', 'SibSp', 'Parch']
train = train.drop(features_drop, axis=1)
test = test.drop(features_drop, axis=1)
train = train.drop(['PassengerId'], axis=1)
```

In [96]:

```
train_data = train.drop('Survived', axis=1)
target = train['Survived']

train_data.shape, target.shape
```

Out[96]:

```
((891, 8), (891,))
```

In [97]:

```
train_data.head()
```

Out[97]:

| | Pclass | Sex | Age | Fare | Cabin | Embarked | Title | FamilySize |
|---|---|---|---|---|---|---|---|---|
| 0 | 3 | 0 | 1.0 | 0.0 | 2.0 | 0 | 0 | 0.4 |
| 1 | 1 | 1 | 3.0 | 2.0 | 0.8 | 1 | 2 | 0.4 |
| 2 | 3 | 1 | 1.0 | 0.0 | 2.0 | 0 | 1 | 0.0 |
| 3 | 1 | 1 | 2.0 | 2.0 | 0.8 | 0 | 2 | 0.4 |
| 4 | 3 | 0 | 2.0 | 0.0 | 2.0 | 0 | 0 | 0.0 |

In [98]:

```
# Importing Classifier Modules
from sklearn.tree import DecisionTreeClassifier
```

```
from sklearn.ensemble import RandomForestClassifier

import numpy as np
```

In [99]:

```
train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 9 columns):
 #   Column      Non-Null Count  Dtype
---  ------      --------------  -----
 0   Survived    891 non-null    int64
 1   Pclass      891 non-null    int64
 2   Sex         891 non-null    int64
 3   Age         891 non-null    float64
 4   Fare        891 non-null    float64
 5   Cabin       891 non-null    float64
 6   Embarked    891 non-null    int64
 7   Title       891 non-null    int64
 8   FamilySize  891 non-null    float64
dtypes: float64(4), int64(5)
memory usage: 62.8 KB
```

In [100]:

```
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 [101]:

```
clf = DecisionTreeClassifier()
scoring = 'accuracy'
score = cross_val_score(clf, train_data, target, cv=k_fold, n_jobs=1, scoring=scoring)
print(score)
```

```
[0.76666667 0.83146067 0.7752809  0.76404494 0.88764045 0.76404494
 0.83146067 0.82022472 0.74157303 0.78651685]
```

In [102]:

```
# decision tree Score
round(np.mean(score)*100, 2)
```

Out[102]:

79.69

In [103]:

```
clf = RandomForestClassifier(n_estimators=13)
scoring = 'accuracy'
score = cross_val_score(clf, train_data, target, cv=k_fold, n_jobs=1, scoring=scoring)
print(score)
```

```
[0.77777778 0.83146067 0.82022472 0.80898876 0.93258427 0.80898876
 0.82022472 0.82022472 0.78651685 0.80898876]
```

In [104]:

```
# Random Forest Score
round(np.mean(score)*100, 2)
```

Out[104]:

82.16

In [105]:

```
clf = RandomForestClassifier(n_estimators=13)
clf.fit(train_data, target)
```

```
test_data = test.drop("PassengerId", axis=1).copy()
prediction = clf.predict(test_data)
```

```
submission = pd.DataFrame({
    "PassengerId": test["PassengerId"],
    "Survived": prediction
  })

submission.to_csv('submission.csv', index=False)
```

```
submission = pd.read_csv('submission.csv')
submission.head()
```

|   | PassengerId | Survived |
|---|-------------|----------|
| 0 | 892 | 0 |
| 1 | 893 | 0 |
| 2 | 894 | 0 |
| 3 | 895 | 0 |
| 4 | 896 | 1 |