In [1]:
```python
"""
ECGR 5105 - Intro to Machine Learning
Homework 2 - Part 1
Phillip Harmon
"""
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
```

In [2]:
```python
#Import Dataset
csvData = pd.read_csv('diabetes.csv')
csvData
```

Out[2]:

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunctic |
|---|---|---|---|---|---|---|---|
| 0 | 6 | 148 | 72 | 35 | 0 | 33.6 | 0.62 |
| 1 | 1 | 85 | 66 | 29 | 0 | 26.6 | 0.35 |
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | 0.67 |
| 3 | 1 | 89 | 66 | 23 | 94 | 28.1 | 0.16 |
| 4 | 0 | 137 | 40 | 35 | 168 | 43.1 | 2.28 |
| ... | ... | ... | ... | ... | ... | ... | |
| 763 | 10 | 101 | 76 | 48 | 180 | 32.9 | 0.17 |
| 764 | 2 | 122 | 70 | 27 | 0 | 36.8 | 0.34 |
| 765 | 5 | 121 | 72 | 23 | 112 | 26.2 | 0.24 |
| 766 | 1 | 126 | 60 | 0 | 0 | 30.1 | 0.34 |
| 767 | 1 | 93 | 70 | 31 | 0 | 30.4 | 0.31 |

768 rows × 9 columns

In [3]:
```python
#Sort Dataset
x = csvData.iloc[:,0:-1].values
y = csvData.iloc[:,-1].values
```

In [4]:
```python
#Train-Test Split
from sklearn.model_selection import train_test_split
xt, xv, yt, yv = train_test_split(x, y, train_size = 0.8, test_size = 0.2, ran
```

In [5]:
```python
#Clean the Dataset
from sklearn.preprocessing import MinMaxScaler, StandardScaler
 # scaler = StandardScaler() #MinMaxScaler gave better results here
scaler = MinMaxScaler()
xt = scaler.fit_transform(xt)
xv = scaler.fit_transform(xv)
```

In [6]:
```python
#Perform the Training
from sklearn.linear_model import LogisticRegression
training_montage = LogisticRegression(random_state=1337)
training_montage.fit(xt,yt);
```

In [7]:
```python
#Test the Model
p = training_montage.predict(xv)
```
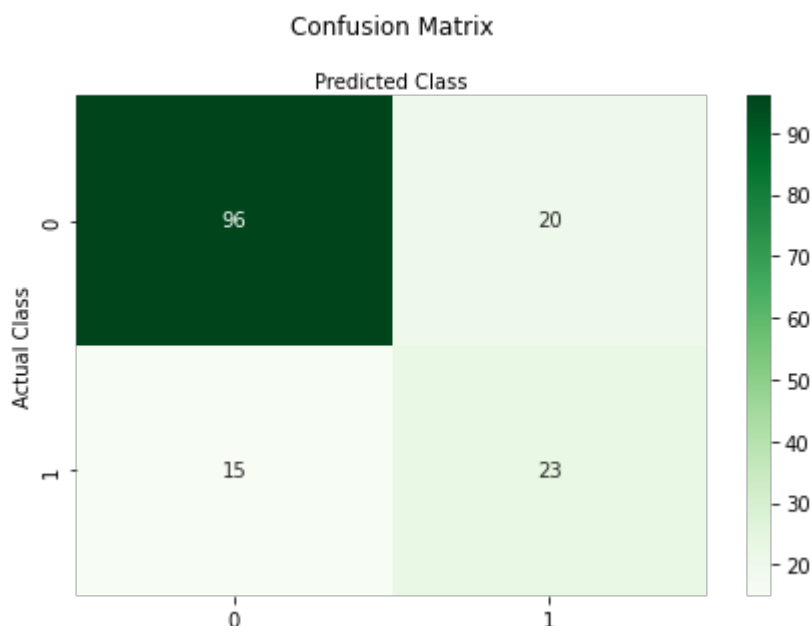
In [8]:
```python
#Evaluate the model metrics
from sklearn import metrics
print("Model Accuracy:  {:.3f}%".format(metrics.accuracy_score(yv,p)*100))
print("Model Precision: {:.3f}%".format(metrics.precision_score(yv,p)*100))
print("Model Recall:    {:.3f}%".format(metrics.recall_score(yv,p)*100))
```

```
Model Accuracy:  77.273%
Model Precision: 53.488%
Model Recall:    60.526%
```

In [9]:
```python
#Analyze using the Confusion Matrix
from sklearn.metrics import confusion_matrix
import seaborn as sns
classes = ['Not Diabetes','Diabetes']
figure, axis = plt.subplots()
ticks = np.arange(len(classes))
plt.xticks(ticks, classes)
plt.yticks(ticks, classes)
sns.heatmap(pd.DataFrame(confusion_matrix(yv, p)), annot=True, cmap="Greens",
axis.xaxis.set_label_position("top")
plt.tight_layout()
plt.title('Confusion Matrix', y=1.1)
plt.ylabel('Actual Class')
plt.xlabel('Predicted Class');
```



In [ ]:

In [1]:
```python
"""
ECGR 5105 - Intro to Machine Learning
Homework 2 - Part 2
Phillip Harmon
"""
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
```

In [2]:
```python
#Import Dataset
csvData = pd.read_csv('diabetes.csv')
csvData
```

Out[2]:

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunctio |
|---|---|---|---|---|---|---|---|
| 0 | 6 | 148 | 72 | 35 | 0 | 33.6 | 0.62 |
| 1 | 1 | 85 | 66 | 29 | 0 | 26.6 | 0.35 |
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | 0.67 |
| 3 | 1 | 89 | 66 | 23 | 94 | 28.1 | 0.16 |
| 4 | 0 | 137 | 40 | 35 | 168 | 43.1 | 2.28 |
| ... | ... | ... | ... | ... | ... | ... | |
| 763 | 10 | 101 | 76 | 48 | 180 | 32.9 | 0.17 |
| 764 | 2 | 122 | 70 | 27 | 0 | 36.8 | 0.34 |
| 765 | 5 | 121 | 72 | 23 | 112 | 26.2 | 0.24 |
| 766 | 1 | 126 | 60 | 0 | 0 | 30.1 | 0.34 |
| 767 | 1 | 93 | 70 | 31 | 0 | 30.4 | 0.31 |

768 rows × 9 columns

In [3]:
```python
#Sort Dataset
x = csvData.iloc[:,0:-1].values
y = csvData.iloc[:,-1].values
```

In [4]:
```python
#Clean the Dataset
from sklearn.preprocessing import MinMaxScaler, StandardScaler
scaler = MinMaxScaler()
x = scaler.fit_transform(x)
```

In [9]:
```python
#Perform the Training with K=5
from sklearn.model_selection import KFold
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import cross_val_score
kcup = KFold(n_splits=5, random_state=1337, shuffle=True)
model = LogisticRegression(random_state=1337)
results = cross_val_score(model,x,y,cv=kcup)
print("K=5 | Accuracy: {:.3f}% ({:.3f}%)".format(results.mean()*100, results.s
```

K=5 | Accuracy: 76.819% (2.656%)

In [11]:
```python
#Perform the Training with K=10
kcup = KFold(n_splits=10, random_state=1337, shuffle=True)
model = LogisticRegression(random_state=1337)
results = cross_val_score(model,x,y,cv=kcup)
print("K=10 | Accuracy: {:.3f}% ({:.3f}%)".format(results.mean()*100, results.
```

K=10 | Accuracy: 76.811% (3.108%)

In [ ]:

In [1]:
```python
"""
ECGR 5105 - Intro to Machine Learning
Homework 2 - Part 3
Phillip Harmon
"""
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
```

In [2]:
```python
#Load and Build the Dataset
from sklearn.datasets import load_breast_cancer
loaded  = load_breast_cancer()
labels  = np.reshape(loaded.target, (len(loaded.target),1))
inputs  = pd.DataFrame(loaded.data)
names   = np.append(loaded.feature_names, 'label')
dataset = pd.DataFrame(np.concatenate([inputs,labels],axis=1))
dataset.columns = names
dataset
```

Out[2]:

| | mean radius | mean texture | mean perimeter | mean area | mean smoothness | mean compactness | mean concavity | mean concave points | mea symmetr |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 17.99 | 10.38 | 122.80 | 1001.0 | 0.11840 | 0.27760 | 0.30010 | 0.14710 | 0.241 |
| 1 | 20.57 | 17.77 | 132.90 | 1326.0 | 0.08474 | 0.07864 | 0.08690 | 0.07017 | 0.181 |
| 2 | 19.69 | 21.25 | 130.00 | 1203.0 | 0.10960 | 0.15990 | 0.19740 | 0.12790 | 0.206 |
| 3 | 11.42 | 20.38 | 77.58 | 386.1 | 0.14250 | 0.28390 | 0.24140 | 0.10520 | 0.259 |
| 4 | 20.29 | 14.34 | 135.10 | 1297.0 | 0.10030 | 0.13280 | 0.19800 | 0.10430 | 0.180 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | . |
| 564 | 21.56 | 22.39 | 142.00 | 1479.0 | 0.11100 | 0.11590 | 0.24390 | 0.13890 | 0.172 |
| 565 | 20.13 | 28.25 | 131.20 | 1261.0 | 0.09780 | 0.10340 | 0.14400 | 0.09791 | 0.175 |
| 566 | 16.60 | 28.08 | 108.30 | 858.1 | 0.08455 | 0.10230 | 0.09251 | 0.05302 | 0.159 |
| 567 | 20.60 | 29.33 | 140.10 | 1265.0 | 0.11780 | 0.27700 | 0.35140 | 0.15200 | 0.239 |
| 568 | 7.76 | 24.54 | 47.92 | 181.0 | 0.05263 | 0.04362 | 0.00000 | 0.00000 | 0.158 |

569 rows × 31 columns

In [3]:
```python
#Sort Dataset
x = dataset.iloc[:,0:-1].values
y = dataset.iloc[:,-1].values
```

In [4]:
```python
#Train-Test Split
from sklearn.model_selection import train_test_split
xt, xv, yt, yv = train_test_split(x, y, train_size = 0.8, test_size = 0.2, ran
```

In [5]:
```python
#Clean the Dataset
from sklearn.preprocessing import MinMaxScaler, StandardScaler
scaler = StandardScaler()
 # scaler = MinMaxScaler() #StandardScaler gave better results here
xt = scaler.fit_transform(xt)
xv = scaler.fit_transform(xv)
```

In [6]:
```python
#Perform the Training
from sklearn.linear_model import LogisticRegression
training_montage = LogisticRegression(random_state=1337)
training_montage.fit(xt,yt);
```

In [7]:
```python
#Test the Model
p = training_montage.predict(xv)
```

In [8]:
```python
#Evaluate the model metrics
from sklearn import metrics
print("Model Accuracy:  {:.3f}%".format(metrics.accuracy_score(yv,p)*100))
print("Model Precision: {:.3f}%".format(metrics.precision_score(yv,p)*100))
print("Model Recall:    {:.3f}%".format(metrics.recall_score(yv,p)*100))
```
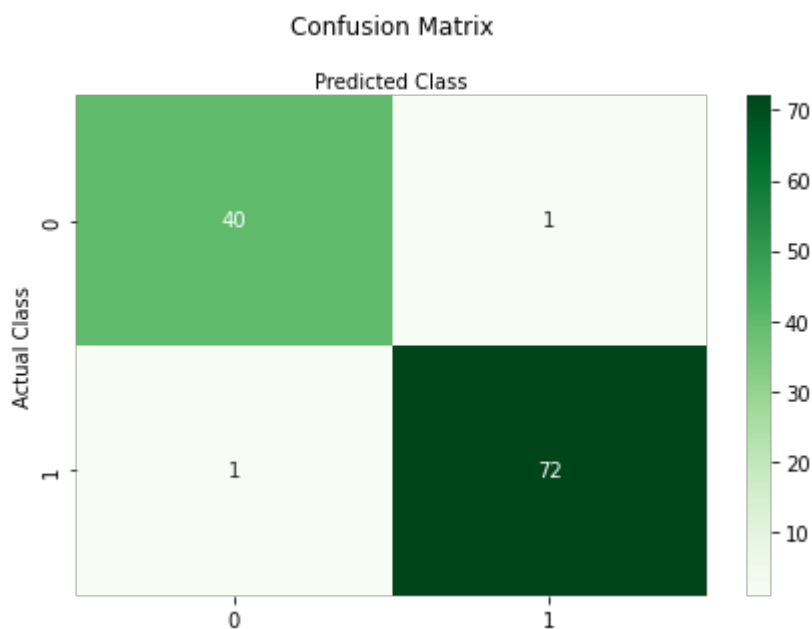
```
Model Accuracy:  98.246%
Model Precision: 98.630%
Model Recall:    98.630%
```
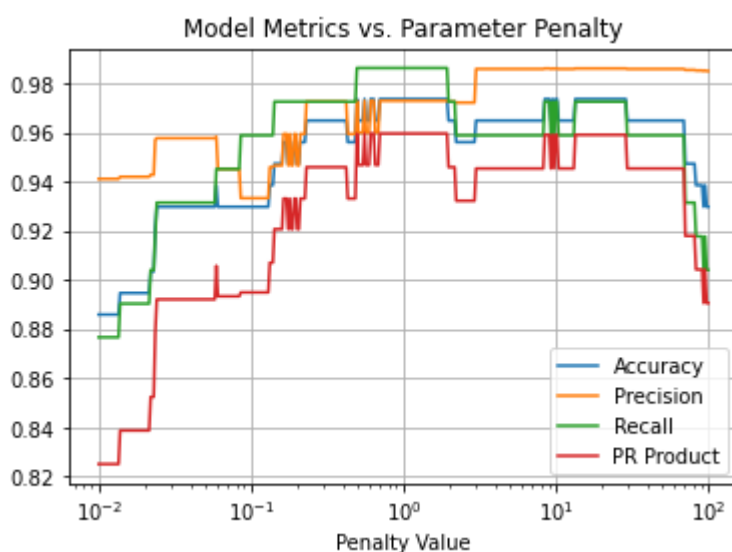
```python
In [9]: #Analyze using the Confusion Matrix
        from sklearn.metrics import confusion_matrix
        import seaborn as sns
        classes = ['Not Diabetes','Diabetes']
        figure, axis = plt.subplots()
        ticks = np.arange(len(classes))
        plt.xticks(ticks, classes)
        plt.yticks(ticks, classes)
        sns.heatmap(pd.DataFrame(confusion_matrix(yv, p)), annot=True, cmap="Greens",
        axis.xaxis.set_label_position("top")
        plt.tight_layout()
        plt.title('Confusion Matrix', y=1.1)
        plt.ylabel('Actual Class')
        plt.xlabel('Predicted Class');
```



```python
In [10]: #Reevaluate using a variety of weight penalties
         lambdas = np.logspace(-2,2,num=400)
         acc_log = []
         prc_log = []
         rec_log = []
         for lam in lambdas:
             model = LogisticRegression(penalty='l1',C=lam,solver='liblinear',random_st
             model.fit(xt,yt)
             p = model.predict(xv)
             rec_log.append(metrics.recall_score(yv,p))
             prc_log.append(metrics.precision_score(yv,p))
             acc_log.append(metrics.accuracy_score(yv,p))
```

In [11]:
```python
#Plot the results
plt.semilogx(lambdas,acc_log,label='Accuracy')
plt.semilogx(lambdas,prc_log,label='Precision')
plt.semilogx(lambdas,rec_log,label='Recall')
PRval = np.multiply(prc_log,rec_log)
plt.semilogx(lambdas,PRval,label='PR Product')
plt.grid()
plt.xlabel('Penalty Value')
plt.title('Model Metrics vs. Parameter Penalty')
plt.legend();
```
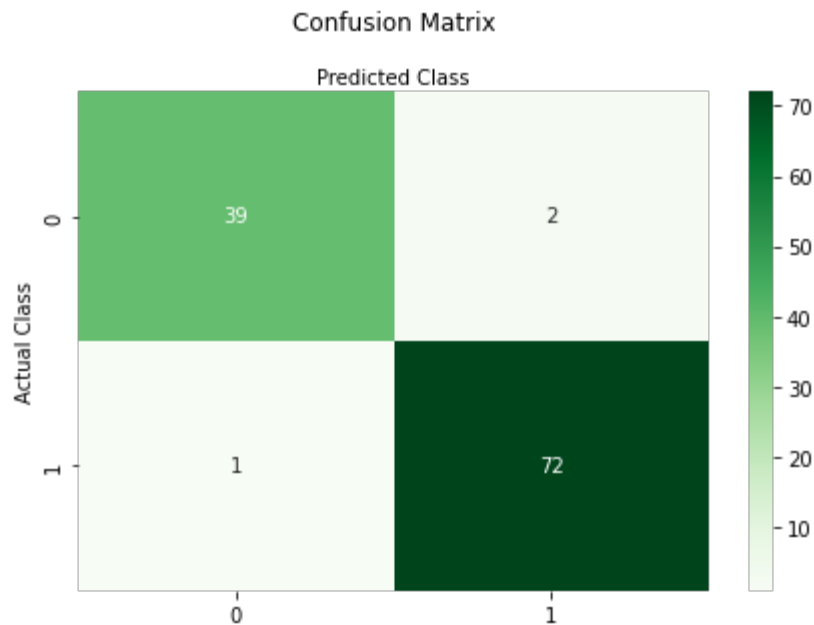
In [12]:
```python
#According to the plot, lambda = 1 is about the best it gets
model = LogisticRegression(penalty='l1',C=1,solver='liblinear',random_state=13
model.fit(xt,yt)
p = model.predict(xv)
print("Model Accuracy:  {:.3f}%".format(metrics.accuracy_score(yv,p)*100))
print("Model Precision: {:.3f}%".format(metrics.precision_score(yv,p)*100))
print("Model Recall:    {:.3f}%".format(metrics.recall_score(yv,p)*100))
classes = ['Not Diabetes','Diabetes']
figure, axis = plt.subplots()
ticks = np.arange(len(classes))
plt.xticks(ticks, classes)
plt.yticks(ticks, classes)
sns.heatmap(pd.DataFrame(confusion_matrix(yv, p)), annot=True, cmap="Greens",
axis.xaxis.set_label_position("top")
plt.tight_layout()
plt.title('Confusion Matrix', y=1.1)
plt.ylabel('Actual Class')
plt.xlabel('Predicted Class');
```

```
Model Accuracy:  97.368%
Model Precision: 97.297%
Model Recall:    98.630%
```



In [ ]:

In [1]:
```python
"""
ECGR 5105 - Intro to Machine Learning
Homework 2 - Part 4
Phillip Harmon
"""
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
```

In [2]:
```python
#Load and Build the Dataset
from sklearn.datasets import load_breast_cancer
loaded  = load_breast_cancer()
labels  = np.reshape(loaded.target, (len(loaded.target),1))
inputs  = pd.DataFrame(loaded.data)
names   = np.append(loaded.feature_names, 'label')
dataset = pd.DataFrame(np.concatenate([inputs,labels],axis=1))
dataset.columns = names
dataset
```

Out[2]:

| | mean radius | mean texture | mean perimeter | mean area | mean smoothness | mean compactness | mean concavity | mean concave points | mean symmetr |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 17.99 | 10.38 | 122.80 | 1001.0 | 0.11840 | 0.27760 | 0.30010 | 0.14710 | 0.241 |
| 1 | 20.57 | 17.77 | 132.90 | 1326.0 | 0.08474 | 0.07864 | 0.08690 | 0.07017 | 0.181 |
| 2 | 19.69 | 21.25 | 130.00 | 1203.0 | 0.10960 | 0.15990 | 0.19740 | 0.12790 | 0.206 |
| 3 | 11.42 | 20.38 | 77.58 | 386.1 | 0.14250 | 0.28390 | 0.24140 | 0.10520 | 0.259 |
| 4 | 20.29 | 14.34 | 135.10 | 1297.0 | 0.10030 | 0.13280 | 0.19800 | 0.10430 | 0.180 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 564 | 21.56 | 22.39 | 142.00 | 1479.0 | 0.11100 | 0.11590 | 0.24390 | 0.13890 | 0.172 |
| 565 | 20.13 | 28.25 | 131.20 | 1261.0 | 0.09780 | 0.10340 | 0.14400 | 0.09791 | 0.175 |
| 566 | 16.60 | 28.08 | 108.30 | 858.1 | 0.08455 | 0.10230 | 0.09251 | 0.05302 | 0.159 |
| 567 | 20.60 | 29.33 | 140.10 | 1265.0 | 0.11780 | 0.27700 | 0.35140 | 0.15200 | 0.239 |
| 568 | 7.76 | 24.54 | 47.92 | 181.0 | 0.05263 | 0.04362 | 0.00000 | 0.00000 | 0.158 |

569 rows × 31 columns

In [3]:
```python
#Sort Dataset
x = dataset.iloc[:,0:-1].values
y = dataset.iloc[:,-1].values
```

In [4]:
```python
#Clean the Dataset
from sklearn.preprocessing import MinMaxScaler, StandardScaler
scaler = StandardScaler()
 # scaler = MinMaxScaler() #StandardScaler gave better results here
x = scaler.fit_transform(x)
```

In [5]:
```python
#Perform the Training with K=5
from sklearn.model_selection import KFold
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import cross_val_score
kcup = KFold(n_splits=5, random_state=1337, shuffle=True)
model = LogisticRegression(random_state=1337)
results = cross_val_score(model,x,y,cv=kcup)
print("K=5 | Accuracy: {:.3f}% ({:.3f}%)".format(results.mean()*100, results.s
```
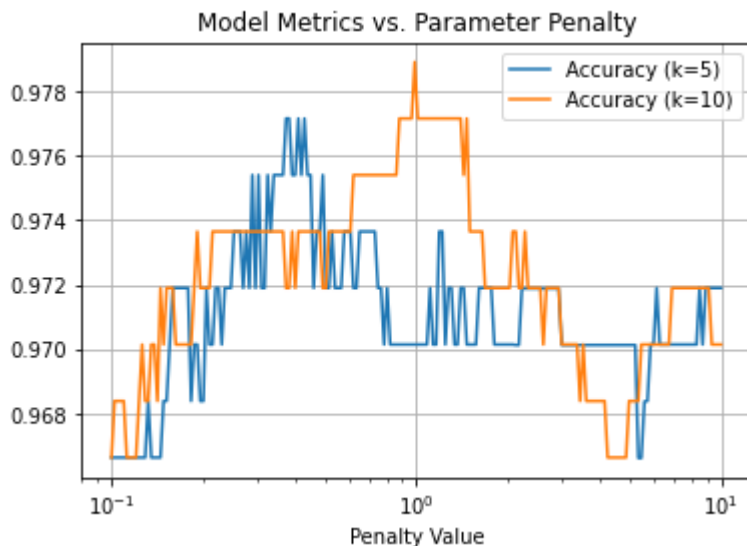
```
K=5 | Accuracy: 97.539% (1.290%)
```

In [6]:
```python
#Perform the Training with K=10
kcup = KFold(n_splits=10, random_state=1337, shuffle=True)
model = LogisticRegression(random_state=1337)
results = cross_val_score(model,x,y,cv=kcup)
print("K=10 | Accuracy: {:.3f}% ({:.3f}%)".format(results.mean()*100, results.
```

```
K=10 | Accuracy: 97.716% (1.122%)
```

In [7]:
```python
#Reevaluate using a variety of weight penalties
lambdas = np.logspace(-1,1,num=200)
k5_acc_log = []
kcup = KFold(n_splits=5, random_state=1337, shuffle=True)
for lam in lambdas:
    model = LogisticRegression(penalty='l1',C=lam,solver='liblinear',random_st
    results = cross_val_score(model,x,y,cv=kcup)
    k5_acc_log.append(results.mean())
k10_acc_log = []
kcup = KFold(n_splits=10, random_state=1337, shuffle=True)
for lam in lambdas:
    model = LogisticRegression(penalty='l1',C=lam,solver='liblinear',random_st
    results = cross_val_score(model,x,y,cv=kcup)
    k10_acc_log.append(results.mean())
```

In [8]:
```python
#Plot the results
plt.semilogx(lambdas,k5_acc_log,label='Accuracy (k=5)')
plt.semilogx(lambdas,k10_acc_log,label='Accuracy (k=10)')
plt.grid()
plt.xlabel('Penalty Value')
plt.title('Model Metrics vs. Parameter Penalty')
plt.legend();
```



In [9]:
```python
#Best k=5 weight is about 0.35
kcup = KFold(n_splits=5, random_state=1337, shuffle=True)
model = LogisticRegression(penalty='l1',C=0.35,solver='liblinear',random_state
results = cross_val_score(model,x,y,cv=kcup)
print("K=5 | Accuracy: {:.3f}% ({:.3f}%)".format(results.mean()*100, results.s
```

K=5 | Accuracy: 97.541% (0.653%)

In [10]:
```python
#Best k=10 weight is about 1
kcup = KFold(n_splits=10, random_state=1337, shuffle=True)
model = LogisticRegression(penalty='l1',C=1,solver='liblinear',random_state=13
results = cross_val_score(model,x,y,cv=kcup)
print("K=10 | Accuracy: {:.3f}% ({:.3f}%)".format(results.mean()*100, results.
```

K=10 | Accuracy: 97.892% (1.312%)

In [ ]: