

Assignment 2 - Group 85

February 10, 2023

Authors: *Niclas Lindmark, Anton Johansson, Noa Sjöstrand*

```
[5]: import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
from scipy import stats
import seaborn as sns
from sklearn.datasets import load_iris
from sklearn.cluster import KMeans
from sklearn import datasets
from sklearn import metrics
from sklearn import linear_model
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
```

```
[6]: d = pd.read_csv("./data_assignment2.csv")

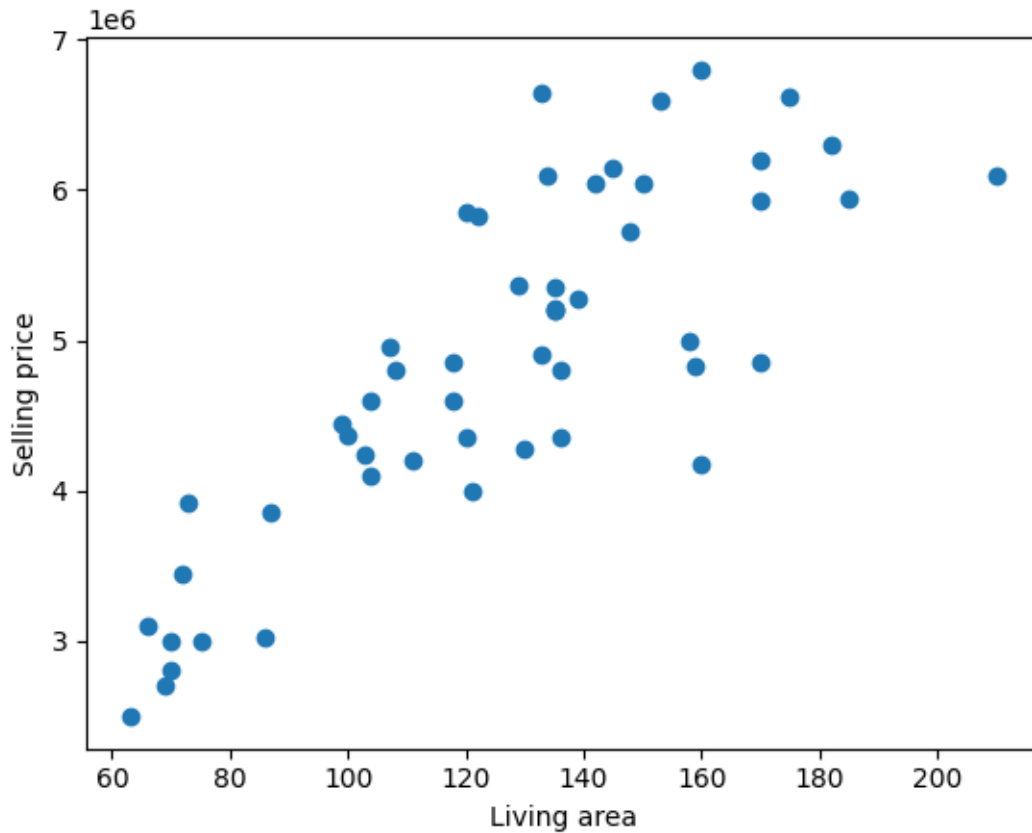
# remove outliers
old_data = d.copy()
d = d.drop(d.index[[45, 40, 24, 9]])
d.head(5)
```

```
[6]:
```

	ID	Living_area	Rooms	Land_size	Biarea	Age	Selling_price
0	1	104	5.0	271.0	25.0	33	4600000
1	2	99	5.0	1506.0	6.0	88	4450000
2	3	133	6.0	486.0	NaN	44	4900000
3	4	175	7.0	728.0	NaN	14	6625000
4	5	118	6.0	1506.0	NaN	29	4600000

```
[7]: x,y = d['Living_area'], d['Selling_price']
x_not_cleaned,y_not_cleaned = old_data['Living_area'], old_data['Selling_price']
plt.scatter(x,y)
plt.xlabel('Living area')
plt.ylabel('Selling price')
```

```
[7]: Text(0, 0.5, 'Selling price')
```



1 Task 1

1.1 Part a)

Find a linear regression model that relates the living area to the selling price. If you did any data cleaning step(s), describe what you did and explain why

```
[8]: x, y = list(x), list(y)
slope, intercept, r, p, std_err = stats.linregress(x, y)

# Data that hasn't been cleaned
x_not, y_not = list(x_not_cleaned), list(y_not_cleaned)
slope_not, intercept_not, r_not, p_not, std_not = stats.linregress(x_not, y_not)

# Regression function
def regression(x):
    return slope * x + intercept
```

```
print(f'Without cleaning the score was {round(r_not, 3)}, after manual cleaning,  
→the score was {round(r, 3)}')
```

Without cleaning the score was 0.562, after manual cleaning the score was 0.817

The correlation value was 0.56 in the original dataset. We plotted the data in a scatter plot, and removed 4 houses that we considered not representative for the rest of the houses. For example, we had a newly built house that were way more expensive than the rest. The high price most likely did not occur because of the living area. When we removed the outliers the correlation value increased to 0.82, greatly improving the model.

1.2 Part b)

What are the values of the slope and intercept of the regression line?

```
[9]: print(f'The slope is = {round(slope, 2)} kr/m2')  
print(f'The intercept is = {round(intercept, 2)} kr')
```

The slope is = 26686.12 kr/m2

The intercept is = 1504030.05 kr

1.3 Part c)

Use this model to predict the selling prices of houses which have living area 100 m2, 150 m2 and 200 m2

```
[10]: #100m2 selling price  
a = regression(100)  
print(f'100m2 predicted selling price: {int(a)} kr')  
#150m2 selling price  
a = regression(150)  
print(f'150m2 predicted selling price: {int(a)} kr')  
#200m2 selling price  
a = regression(200)  
print(f'200m2 predicted selling price: {int(a)} kr')
```

100m2 predicted selling price: 4172642 kr

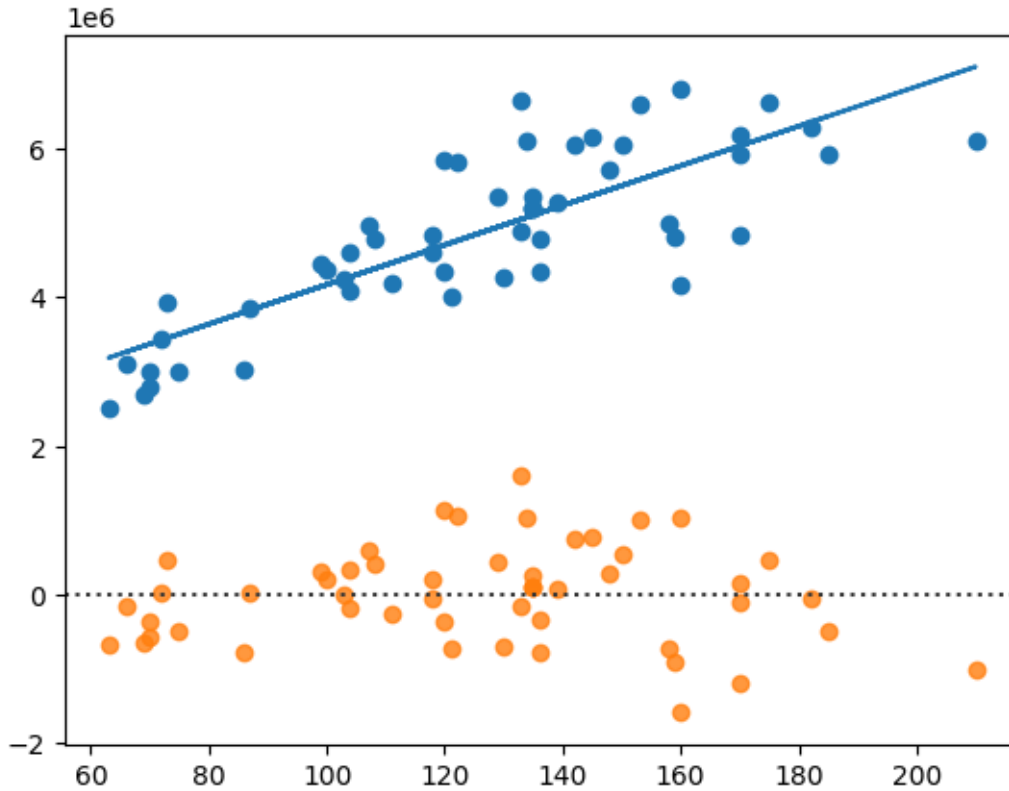
150m2 predicted selling price: 5506948 kr

200m2 predicted selling price: 6841254 kr

1.4 Part d)

Draw a residual plot.

```
[11]: model = list(map(regression, x))  
plt.scatter(x, y)  
plt.plot(x, model)  
sns.residplot(x=x, y=y, data=d)  
plt.show()
```



2 Task 2

2.1 Part a)

Use a confusion matrix to evaluate the use of logistic regression to classify the iris data set

```
[12]: data = load_iris()
      # data
```

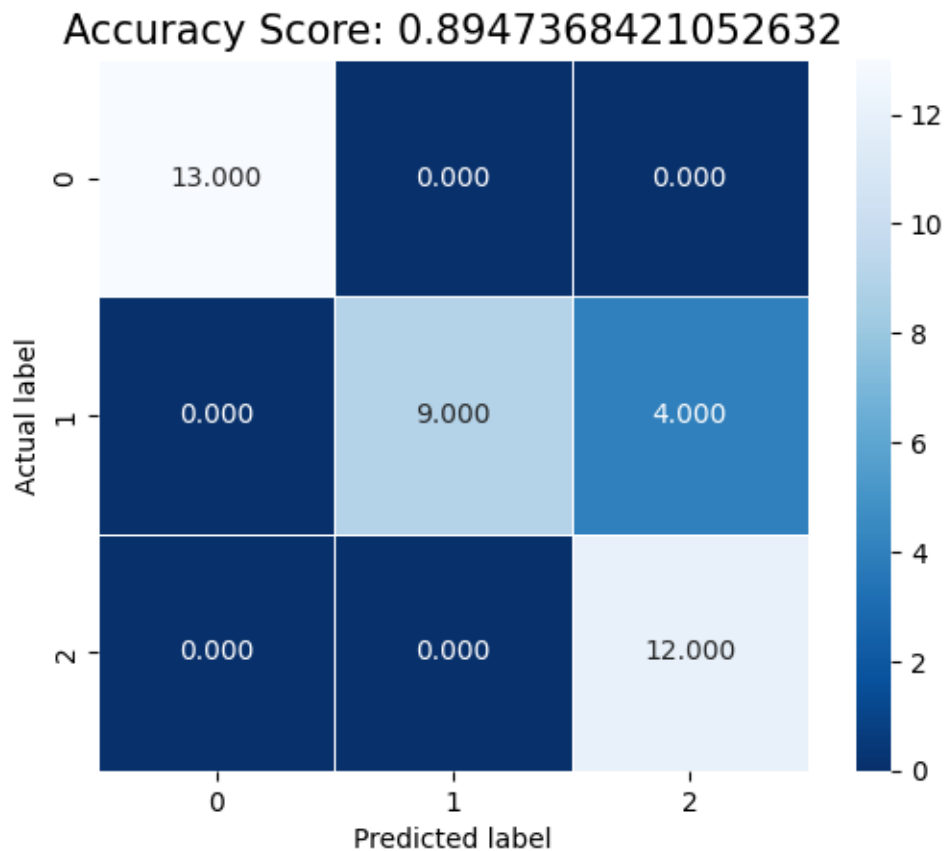
```
[13]: # make a logistic regressor
x_train, x_test, y_train, y_test = train_test_split(data.data, data.target)
clf = linear_model.LogisticRegression(multi_class='ovr', solver='liblinear').
    ↪ fit(x_train, y_train)
clf.predict(x_test[0].reshape(1,-1))
score = clf.score(x_test, y_test)
```

```
[14]: # make a confusion matrix
predictions = clf.predict(x_test)
cm = metrics.confusion_matrix(y_test, predictions)
```

```

sns.heatmap(cm, annot=True, fmt=".3f", linewidths=.5, square = True, cmap = 'Blues_r')
plt.ylabel('Actual label');
plt.xlabel('Predicted label');
all_sample_title = 'Accuracy Score: {0}'.format(score)
plt.title(all_sample_title, size = 15);
plt.savefig('toy_Digits_ConfusionSeabornCodementor.png')
plt.show()

```



2.2 Part b)

Use k-nearest neighbours to classify the iris data set with some different values for k, and with uniform and distance-based weights. What will happen when k grows larger for the different cases? Why?

[28]: *# making a for loop to iterate between different values of k-neighbours, to find the optimal number of neighbours.*

```

ret_list_dis = []
ret_list_uni = []

```

```

for w in ['distance', 'uniform']:
    for i in range(1, 20):
        neigh = KNeighborsClassifier(n_neighbors=i, weights= w)
        avg_list = []
        for j in range(1000):
            x_train, x_test, y_train, y_test = train_test_split(data.data, data.
→target)
            neigh.fit(x_train, y_train)
            y_pred = neigh.predict(x_test)
            avg_list.append(metrics.accuracy_score(y_test, y_pred))
        if w == 'distance':
            ret_list_dis.append((i, np.mean(avg_list)))
        else:
            ret_list_uni.append((i, np.mean(avg_list)))

# finding best accuracy for uniform and distance based
best_dis = [0,0]
for (n, s) in ret_list_dis:
    if s > best_dis[1]:
        best_dis[0] = n
        best_dis[1] = s

best_uni = [0,0]
for (n, s) in ret_list_uni:
    if s > best_uni[1]:
        best_uni[0] = n
        best_uni[1] = s

# plotting the data

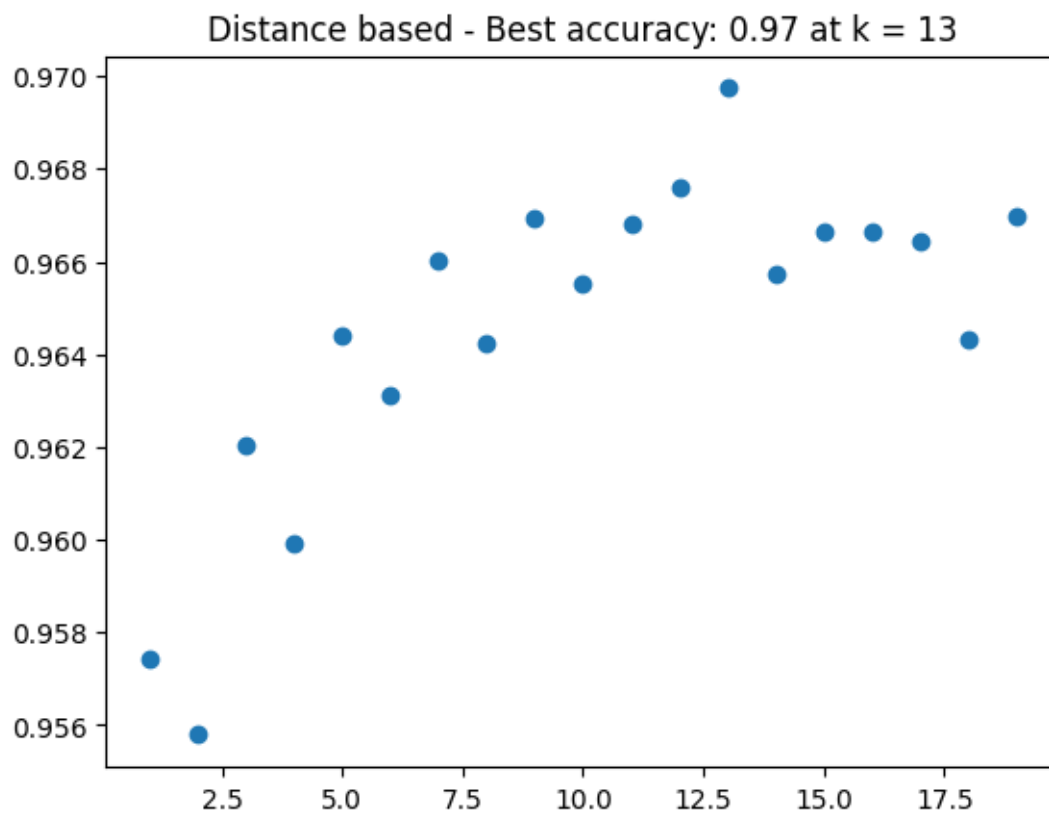
x_vals_dis = [x for (x, y) in ret_list_dis]
y_vals_dis = [y for (x, y) in ret_list_dis]

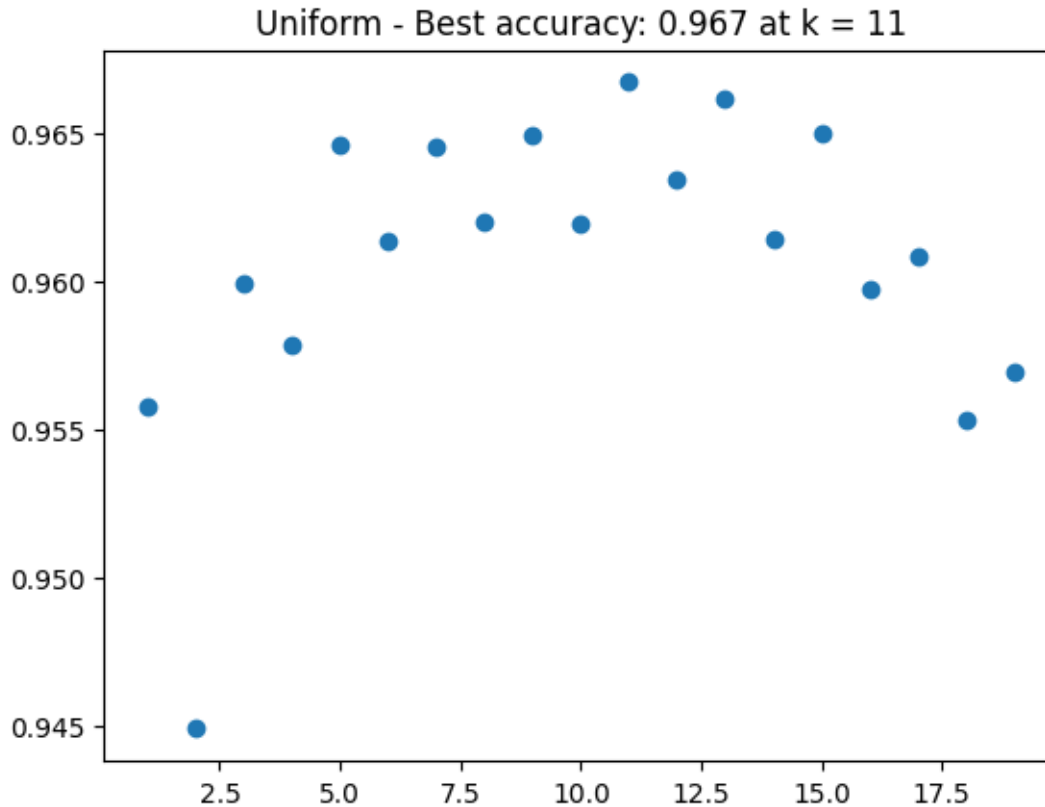
x_vals_uni = [x for (x, y) in ret_list_uni]
y_vals_uni = [y for (x, y) in ret_list_uni]

plt.scatter(x_vals_dis, y_vals_dis)
plt.title(f'Distance based - Best accuracy: {round(best_dis[1], 3)} at k =_
→{best_dis[0]}')
plt.show()

plt.scatter(x_vals_uni, y_vals_uni)
plt.title(f'Uniform - Best accuracy: {round(best_uni[1], 3)} at k =_
→{best_uni[0]}')
plt.show()

```





Answer: We have iterated through different values for k . We started from 1 to 100 and saw that the accuracy dropped dramatically at around $k=60$. We saw that the most accurate values were in the interval 1-20 and chose to run iterations through this interval with many repetitions (1000 repetitions) to get a reliable average value. We did this for both 'distance based' and 'uniform' KNN.

Uniform: We plotted the values and saw that the odd numbers performed noticeably better, and think that is because there always will be more of one type of neighbour. And because the number of neighbours is the only thing that matters in the uniform based model, this poses a problem when the number of neighbours is even. If you have an even number it might be a case where you have 4 neighbours of one type and 4 neighbours of another type, and then the model has to choose randomly, thus even numbers reduce the accuracy.

Distance based: We plotted the values and found that there wasn't as big of a difference between even and odd numbers. That is as expected since this model also uses euclidean distance to determine the type of neighbour.

From this data we chose following k -values: - 13 for distance based - 11 for uniform

The reason why the accuracy drops as k -grows is because we are underfitting the model.

2.3 Part c)

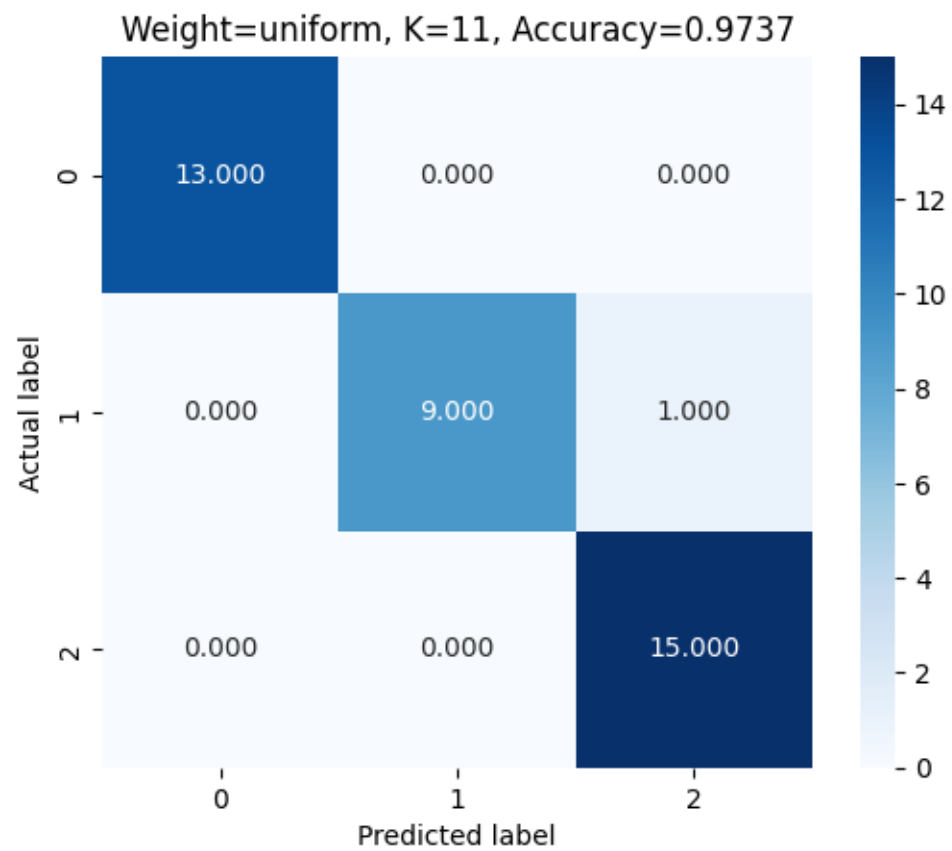
Compare the classification models for the iris data set that are generated by k-nearest neighbours (for the different settings from question b) and by logistic regression. Calculate confusion matrices for these models and discuss the performance of the various models

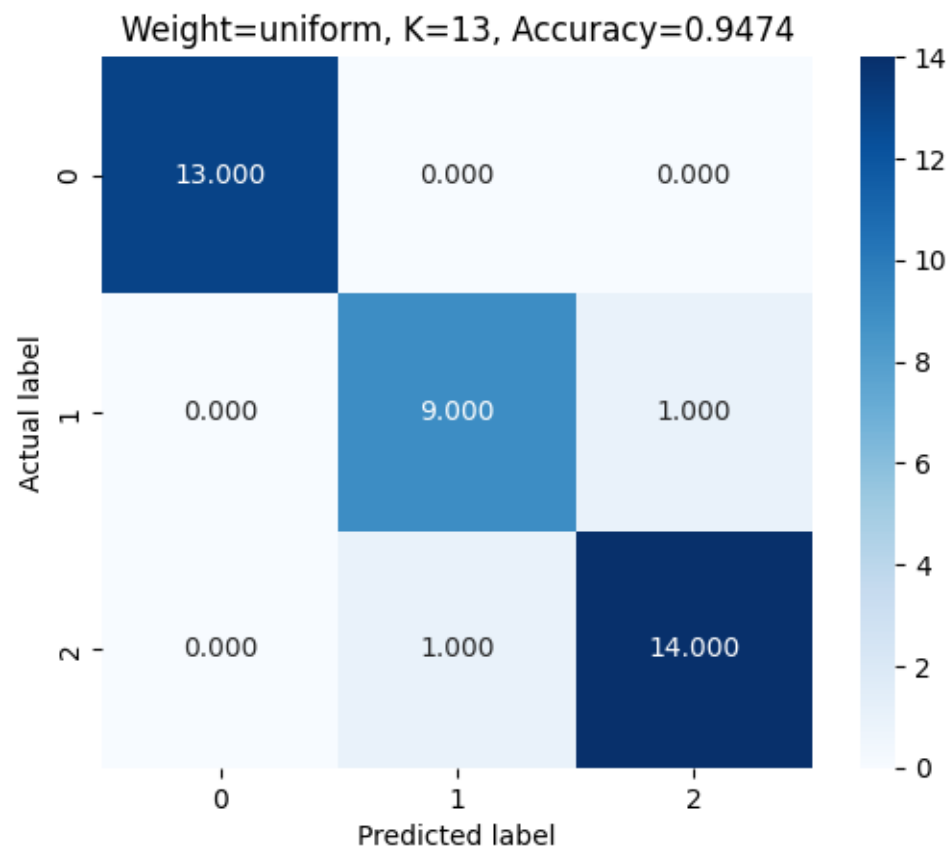
```
[66]: x_train, x_test, y_train, y_test = train_test_split(data.data, data.target)

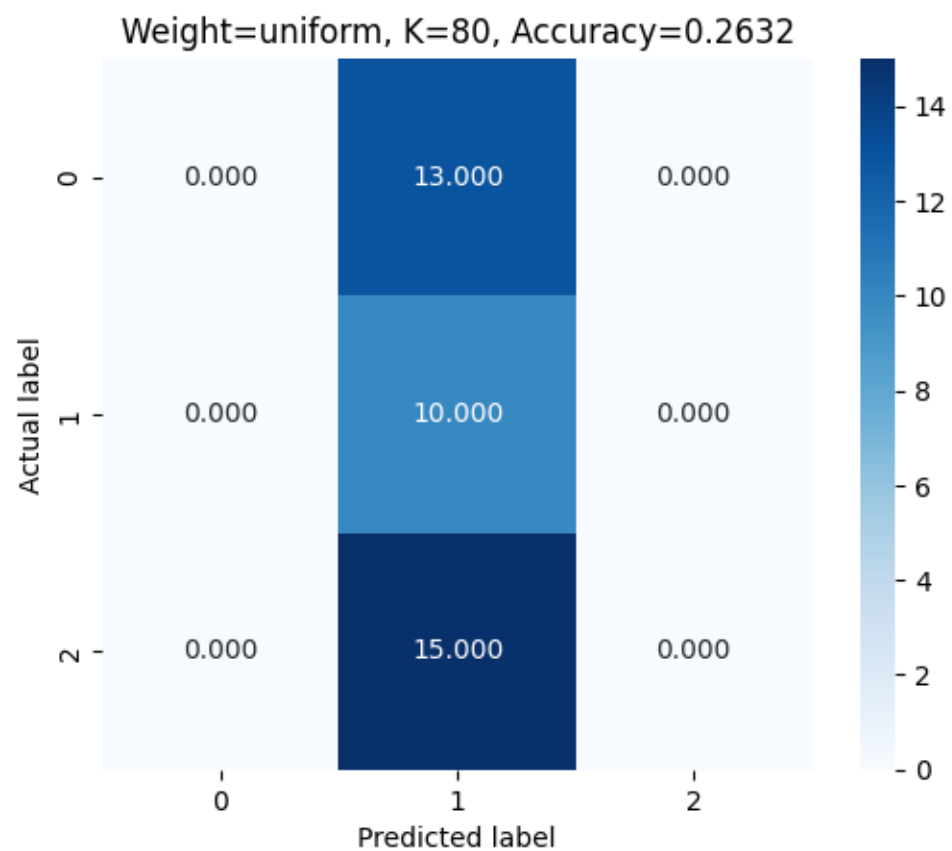
# k-neighbours
weights = ['uniform', 'distance']
values = [11, 13, 80]
for weight in weights:
    for value in values:
        neigh = KNeighborsClassifier(n_neighbors=value, weights=weight)
        neigh.fit(x_train, y_train)

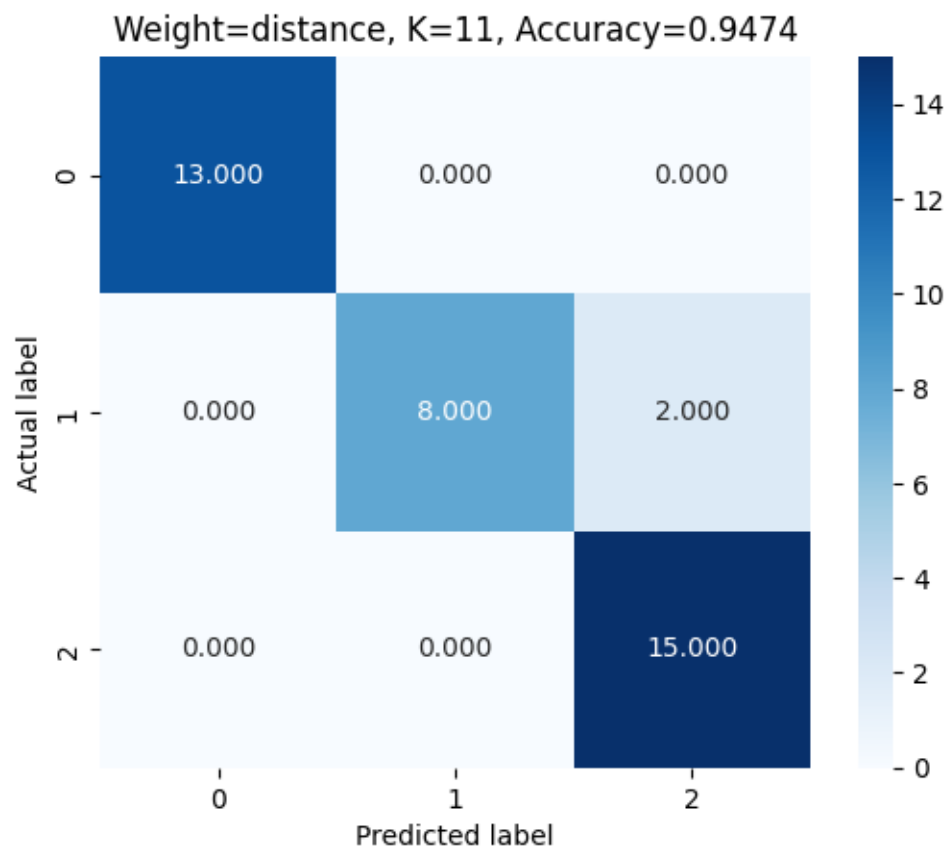
        pred_Kn = neigh.predict(x_test)
        cm_Kn = metrics.confusion_matrix(y_test, pred_Kn)
        score_Kn = metrics.accuracy_score(y_test, pred_Kn)

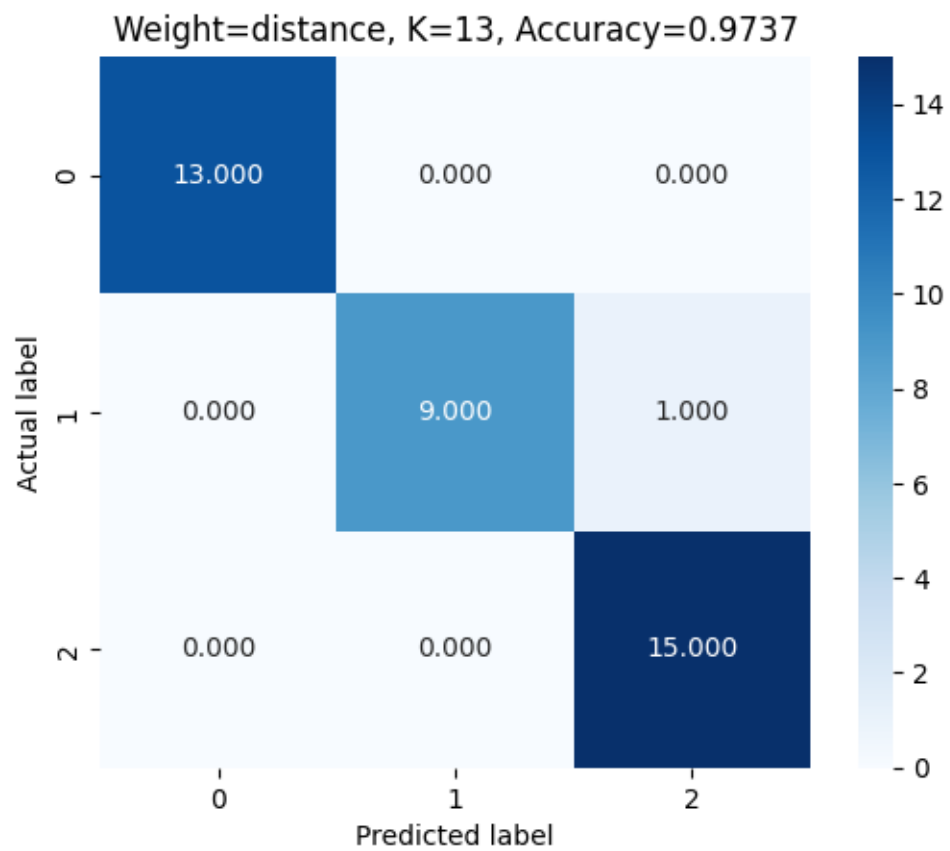
        sns.heatmap(cm_Kn, annot=True, fmt=".3f", linewidths=0, square = True,
→ cmap = 'Blues')
        plt.ylabel('Actual label')
        plt.xlabel('Predicted label')
        all_sample_title = (f'Weight={weight}, K={value},
→ Accuracy={round(score_Kn, 4)}')
        plt.title(all_sample_title, size = 12)
        plt.show()
```

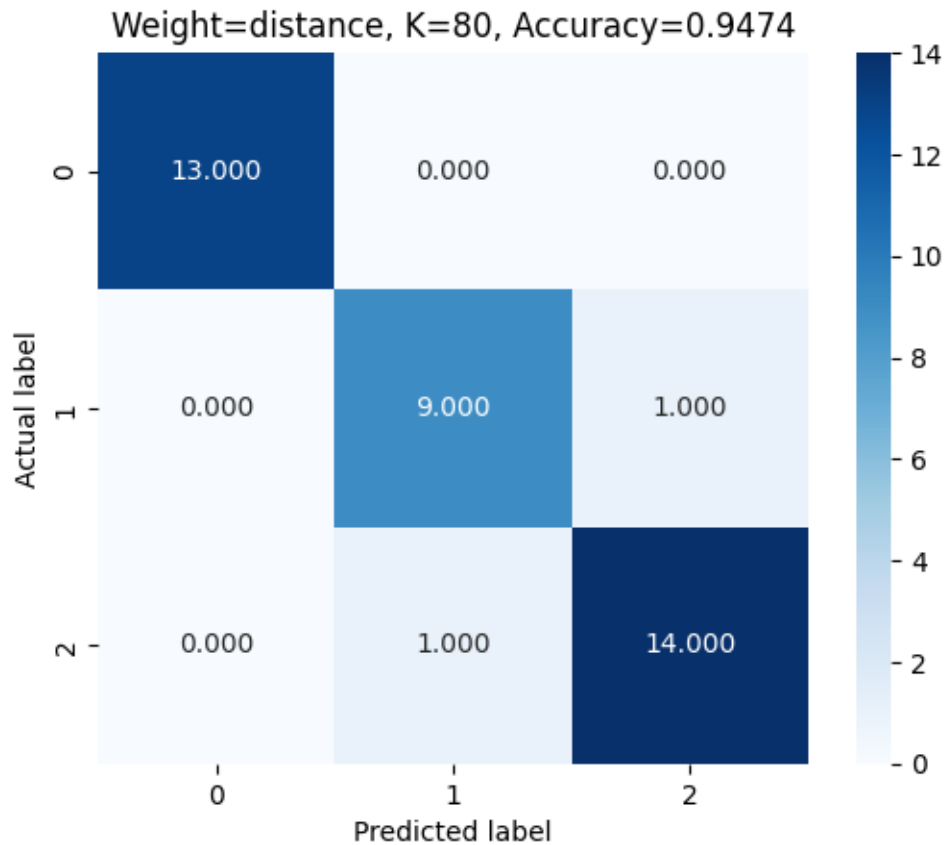








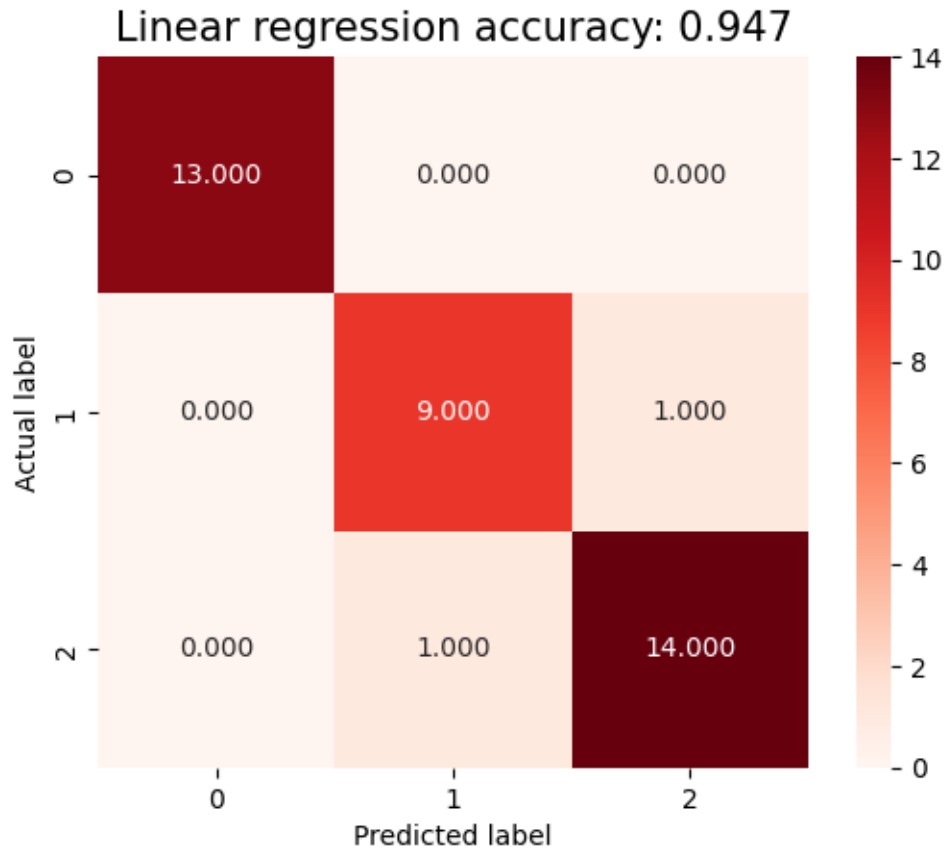




```
[72]: #logistic regression
clf = linear_model.LogisticRegression(multi_class='ovr', solver='liblinear').
    ↪ fit(x_train, y_train)

pred_lr = clf.predict(x_test)
cm_lr = metrics.confusion_matrix(y_test, pred_lr)
score_lr = clf.score(x_test, y_test)

sns.heatmap(cm_lr, annot=True, fmt=".3f", linewidths=0, square = True, cmap =
    ↪ 'Reds')
plt.ylabel('Actual label')
plt.xlabel('Predicted label')
all_sample_title = 'Linear regression accuracy: {0}'.format(round(score_lr, 3))
plt.title(all_sample_title, size = 15)
plt.show()
```



Answer: The performance is similar for KNN and linear regression. We tested with different data and the best performing model varied.

We can see a noticable difference between the distance based and uniform KNN in the confusion matrix. Both perform best at numbers around 11 however as you increase k the uniform based KNN deteriorates much quicker. The accuracy score for the uniform falls to around 25% when $k = 80$ while the distance based has an accuracy score of 95% for $k = 80$. A conclusion from this is that you should choose distance based KNN if you don't have time to find out what the best k-value is, because distance based is not as susceptible to 'bad' k-values.

With a 'good' value of k (11 for uniform, 13 for distance based) we found that the KNN had a slight upper hand compared to the linear regression, being a bit more consistent with the predictions. However the KNN requires more time as you have to find the best k-value, and that is especially true for the uniform KNN. If you are not careful when choosing your k value in the uniform KNN you will get a bad model.

Our conclusion is that the k-neighbors model is slightly more accurate but the linear regression model is easier to use and less likely to be improperly implemented.