

You can find my results on my github: <https://github.com/Barrel-Titor/homework-MLDL>

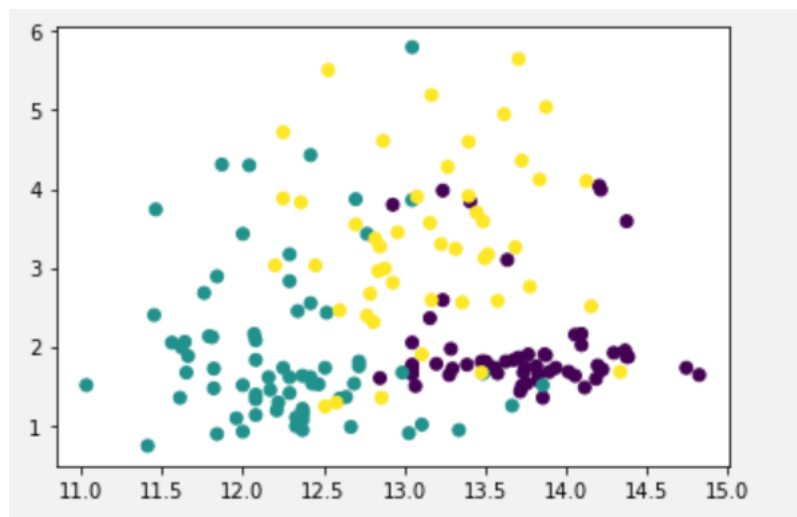
# KNN

## 1. Load Wine dataset

```
1 from sklearn.datasets import load_wine
2 data, target = load_wine(True)
```

## 2. Select two attributes

```
1 data = data[:, :2]
2 data.shape # (178, 2)
```



## 3. Split data into train, validation and test

Use `train_test_split` twice to randomly cut data into three sets, with proportion 5: 2: 3

```

1 from sklearn.model_selection import train_test_split
2
3 # Train V.S. Val + Test
4 X_train, X_test, y_train, y_test = train_test_split(
5     data, target, test_size=0.5, random_state=0
6 )
7
8 # Train V.S. Val V.S. Test
9 X_val, X_test, y_val, y_test = train_test_split(
10     X_test, y_test, test_size=0.6, random_state=0
11 )
12
13 len(y_train), len(y_val), len(y_test)    # (89, 35, 54)

```

## 4. Try K = [1, 3, 5, 7]

Define function `plot_boundary`, using `plt.contourf()` to help visualize the decision boundaries of KNN

```

1 lc1 = ListedColormap(['#FFAAAA', '#AAFFAA', '#AAAAFF'])
2 lc2 = ListedColormap(['#FF0000', '#00FF00', '#0000FF'])
3
4 def plot_boundary(func, points, labels):
5     x_min = X_train[:, 0].min() - .5
6     x_max = X_train[:, 0].max() + .5
7     y_min = X_train[:, 1].min() - .5
8     y_max = X_train[:, 1].max() + .5
9     h = 0.1
10
11     xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min,
12 y_max, h))
13     z = func.predict(np.c_[xx.ravel(), yy.ravel()])
14     z = z.reshape(xx.shape)
15
16     plt.contourf(xx, yy, z, cmap=lc1)
17     plt.scatter(points[:, 0], points[:, 1], c=labels, cmap=lc2)

```

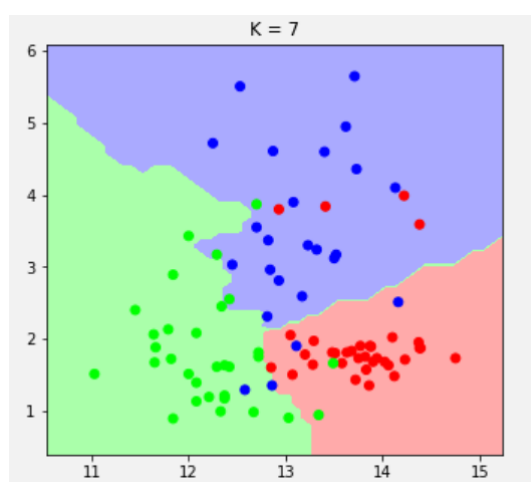
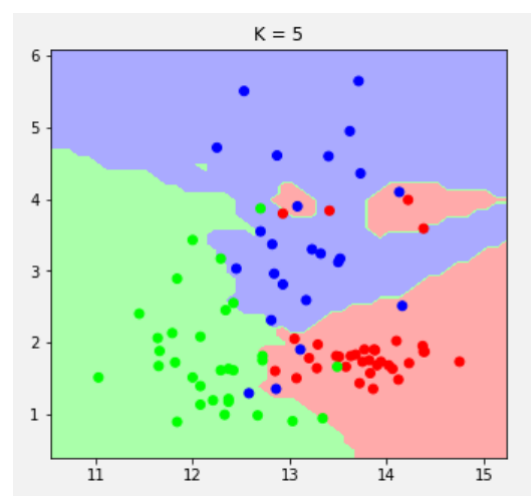
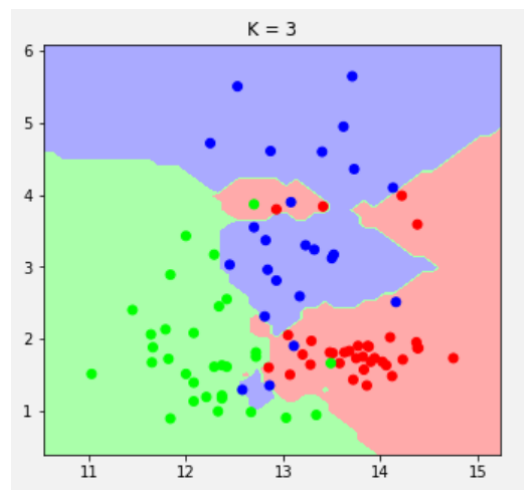
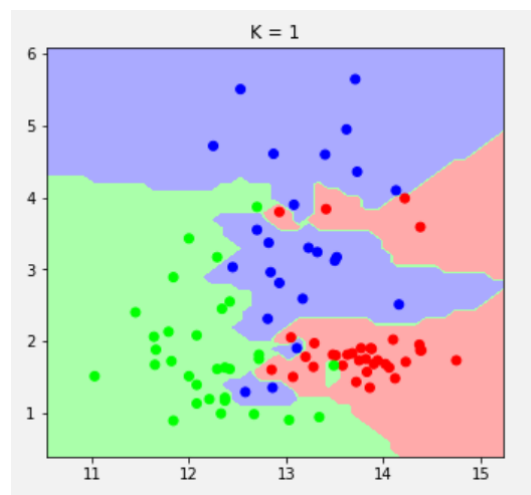
Define function `K_KNN` to fit data and validate

```

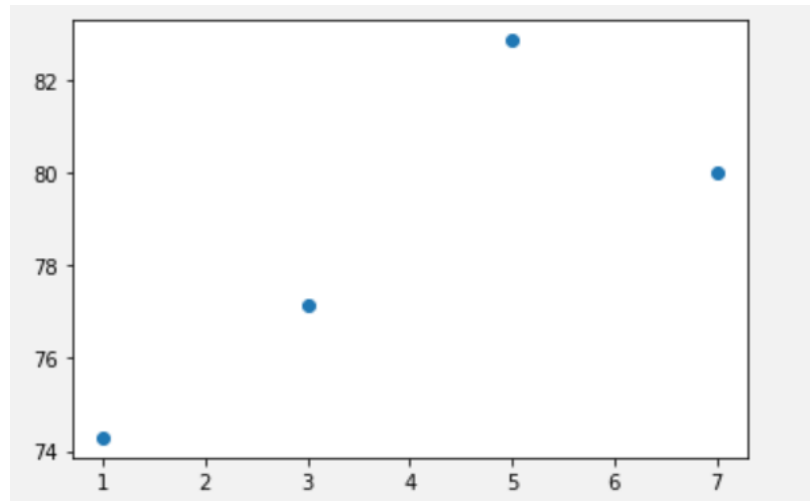
1 def K_KNN(K):
2     clf = KNeighborsClassifier(K)
3     clf.fit(X_train, y_train)
4     clf_list.append(clf)
5     plot_boundary(clf, X_train, y_train)
6     y_val_predicted = clf.predict(X_val)
7     accuracy = np.mean(y_val_predicted == y_val) * 100
8     accuracies.append(accuracy)

```

The results look as follows:



## 5. Show accuracy on validation set



It shows it's better when  $K = 5$

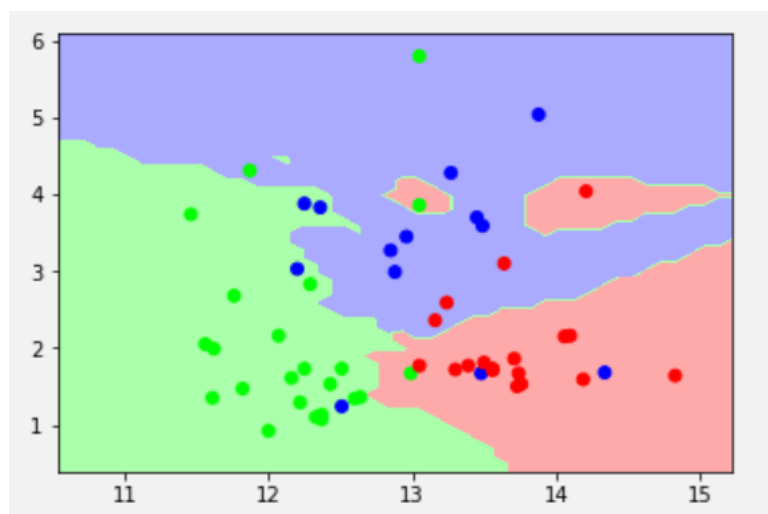
## 6. How the boundaries change?

When  $K$  is too small, neighborhood of each point is very small, which can lead to overfitting. Therefore, the boundary can be easily disturbed by noise and outliers.

As  $K$  increases, the boundary becomes simpler has more powerful generalization capabilities. However, if  $K$  is too large, points of different labels become its neighbors and the accuracy will become lower.

## 7. Use the best value of $K$ on the test set

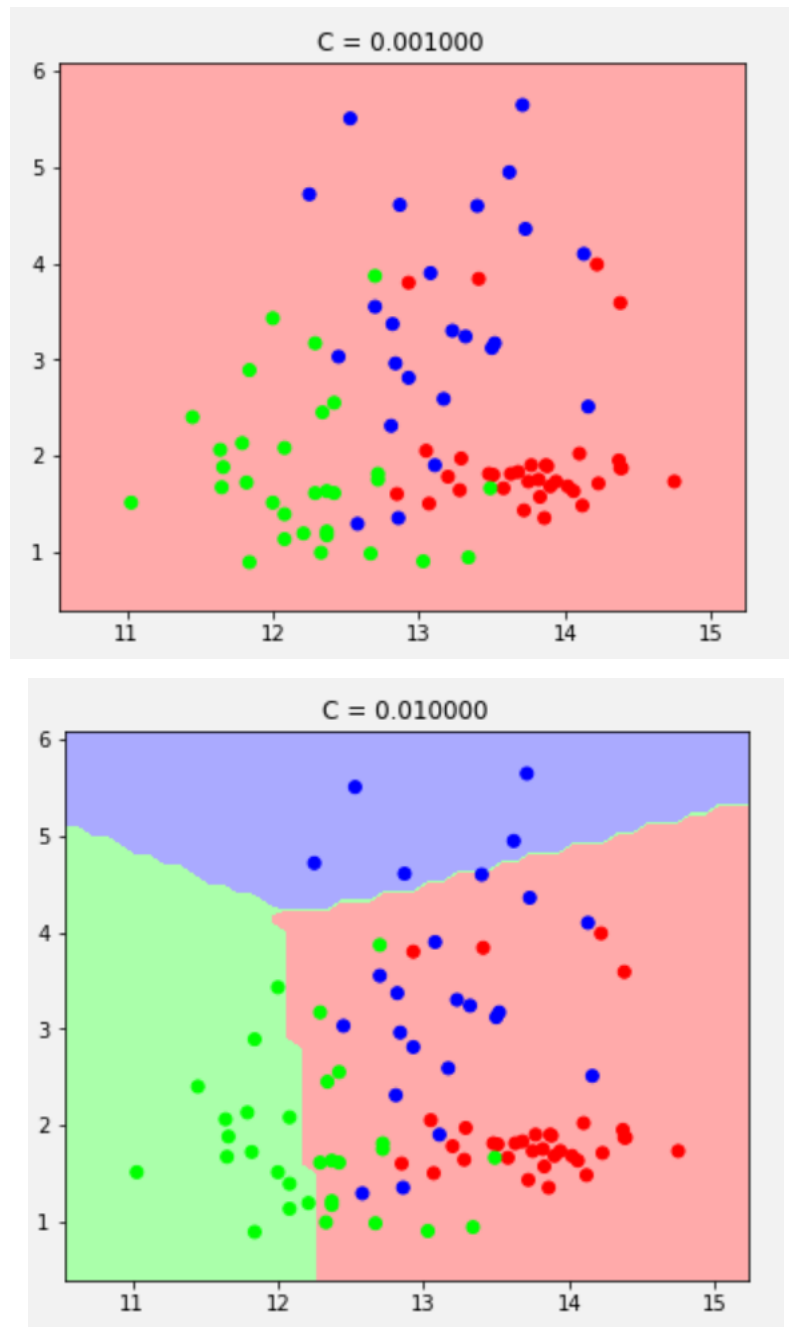
```
1 ! [6](6.png)y_test_predicted = clf_list[2].predict(X_test)
2 accuracy = np.mean(y_test_predicted == y_test) * 100
3 accuracy    # 75.92592592592592
```

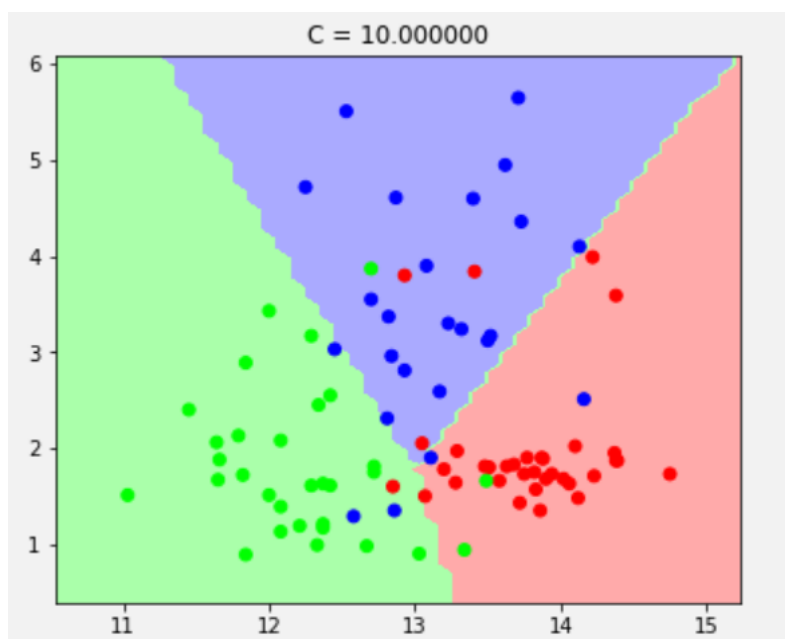
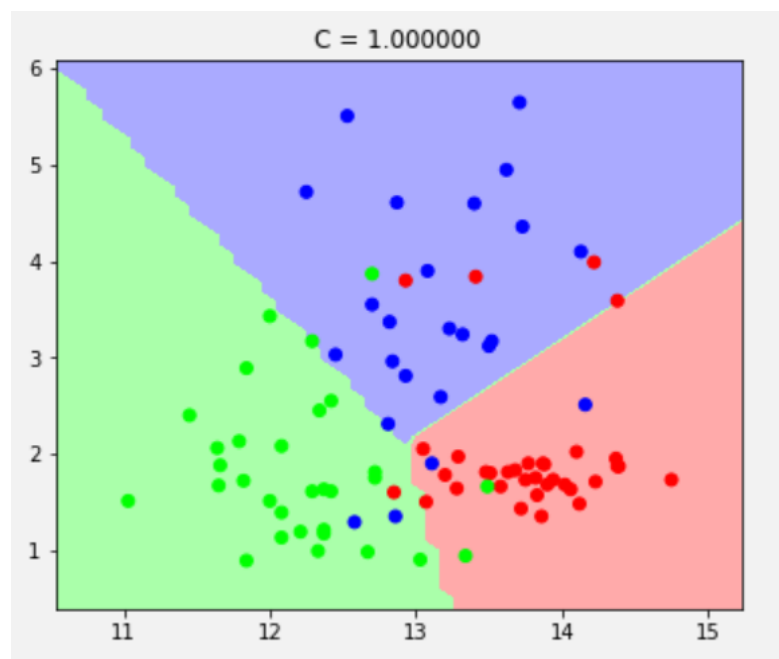
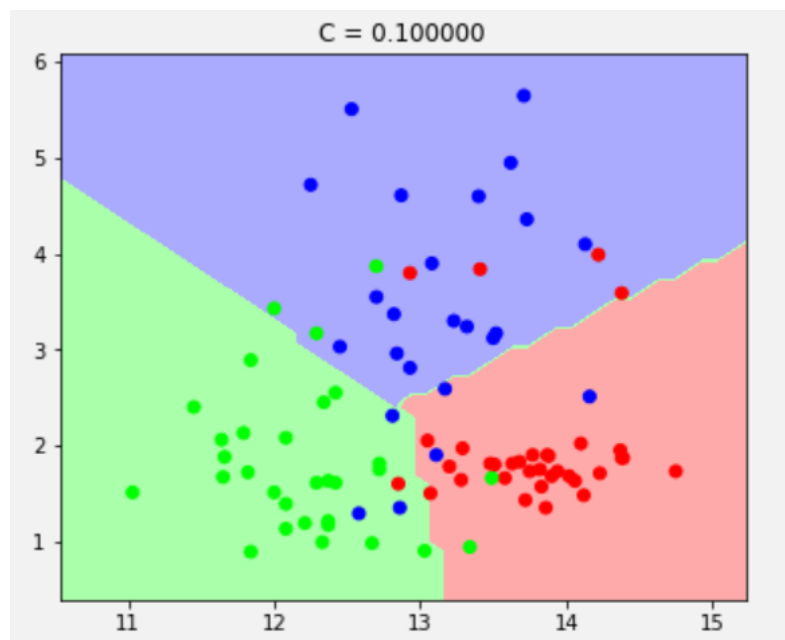


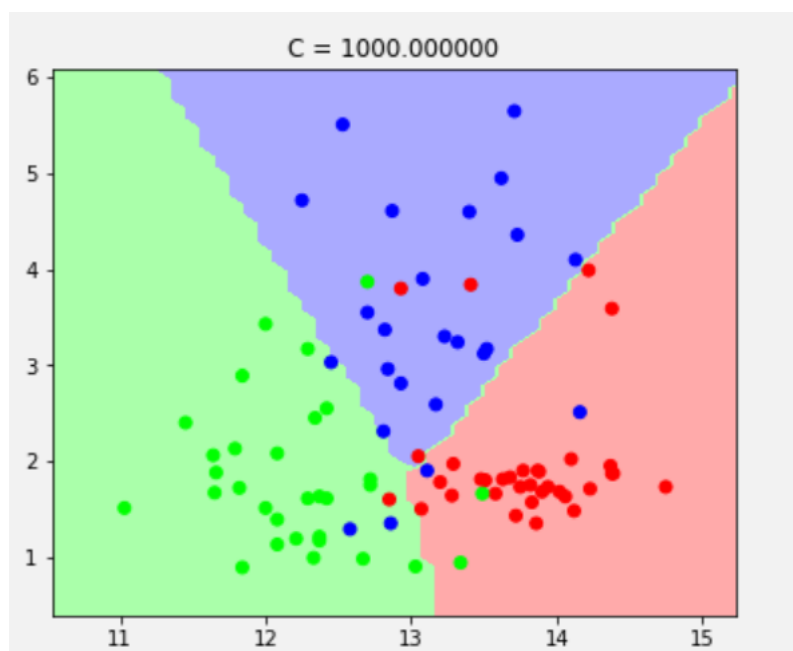
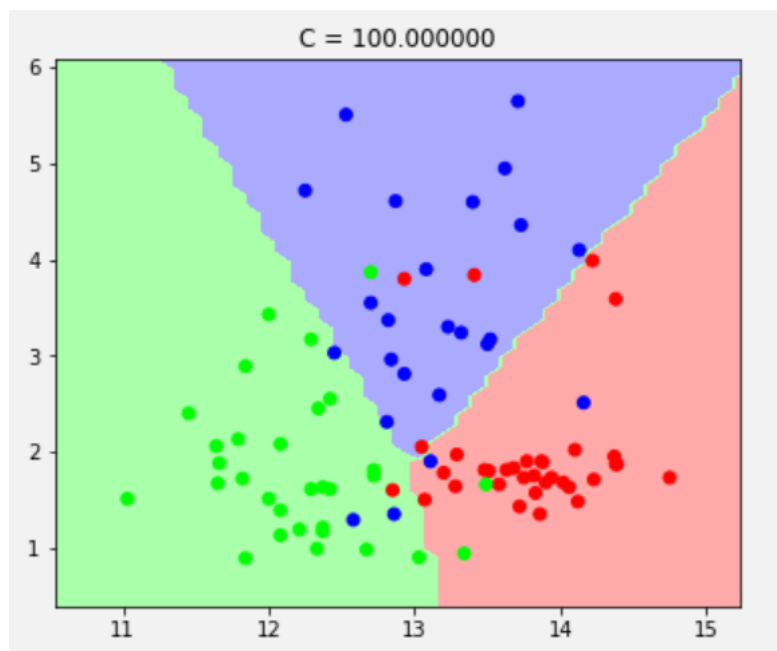
# Linear SVM

## 8. Try $C = [0.001, 0.01, 0.1, 1, 10, 100, 1000]$

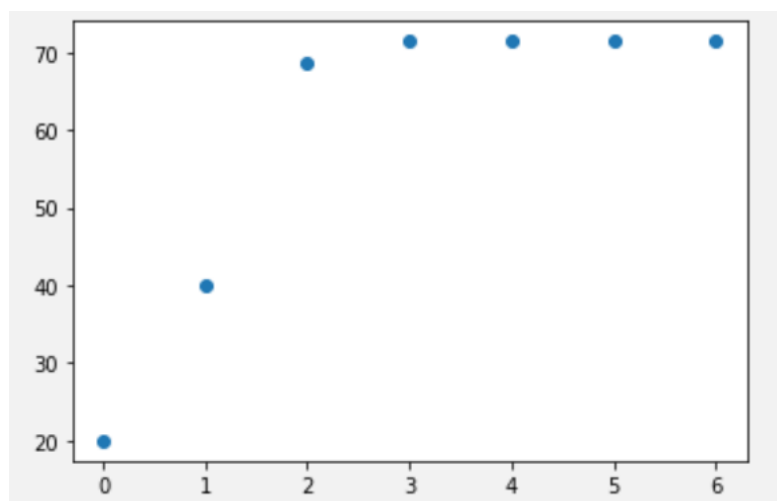
Use `plot_boundary` above to plot decision boundary of SVM







## 9. Show accuracy on validation set



It seems when  $C \geq 1$ , the accuracy doesn't vary much.

## 10. How the boundaries change?

---

The larger the value of  $C$ , the greater the penalty for outliers, and the less willing the classifier is to allow outliers.

## 11. Use the best value of $C$ on the test set

---

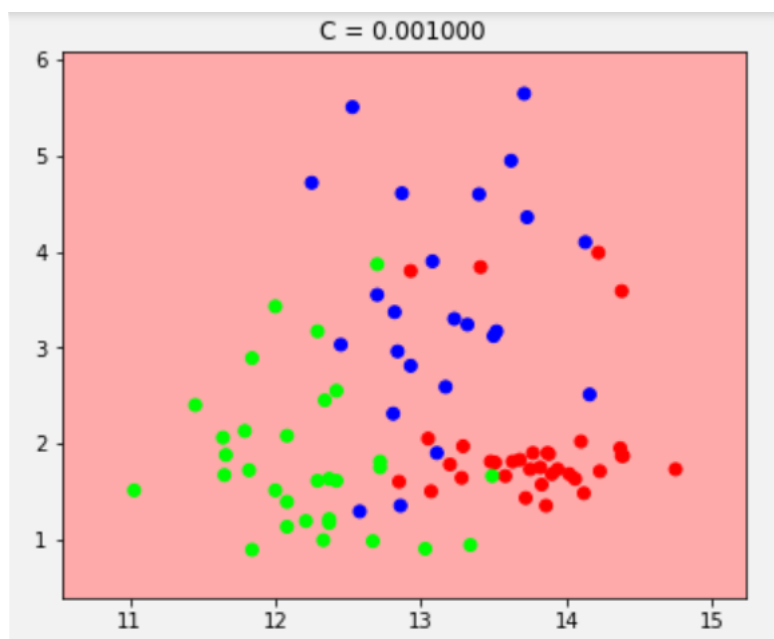
```
1 accuracy = np.mean(y_test_predicted == y_test) * 100
2 accuracy  # 79.62962962962963 when C = 1
3 accuracy  # 83.33333333333334 when C = 10
4 accuracy  # 83.33333333333334 when C = 100
5 accuracy  # 83.33333333333334 when C = 1000
```

## SVM with RBF Kernel

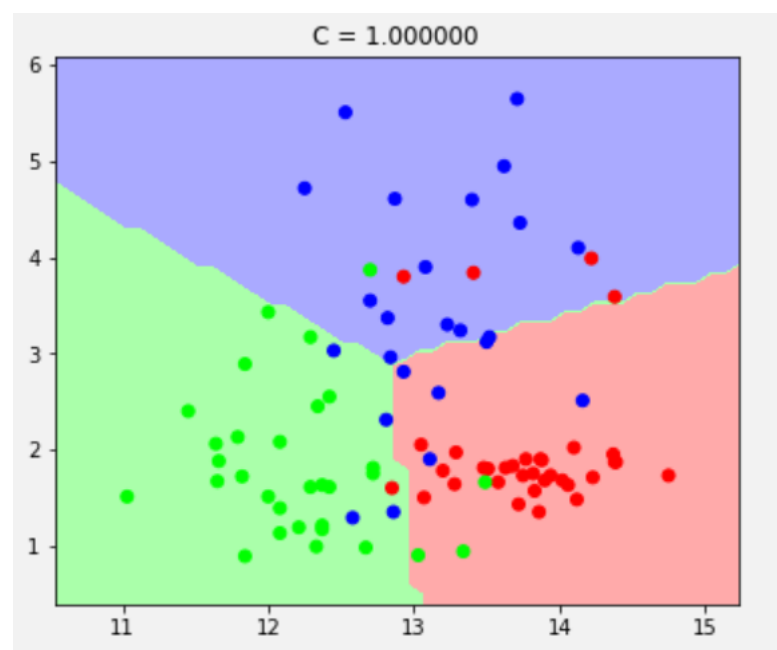
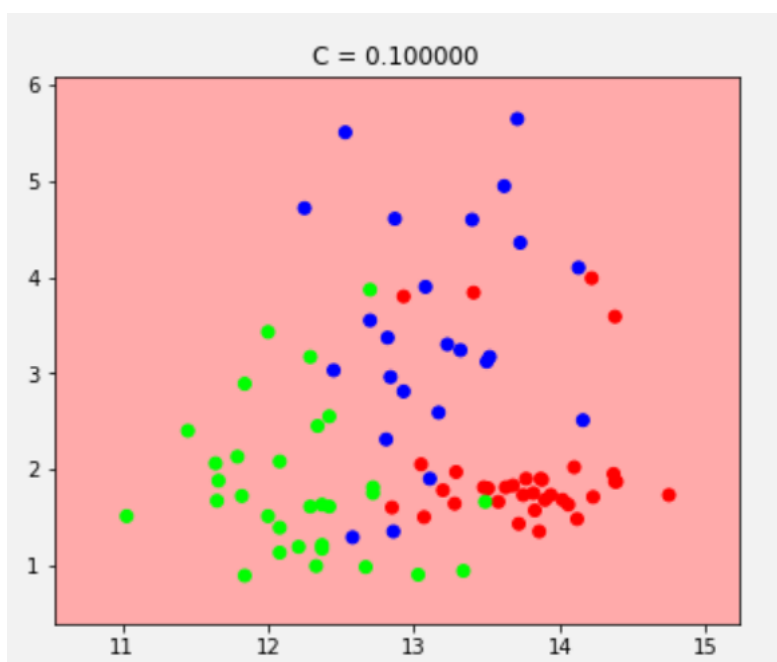
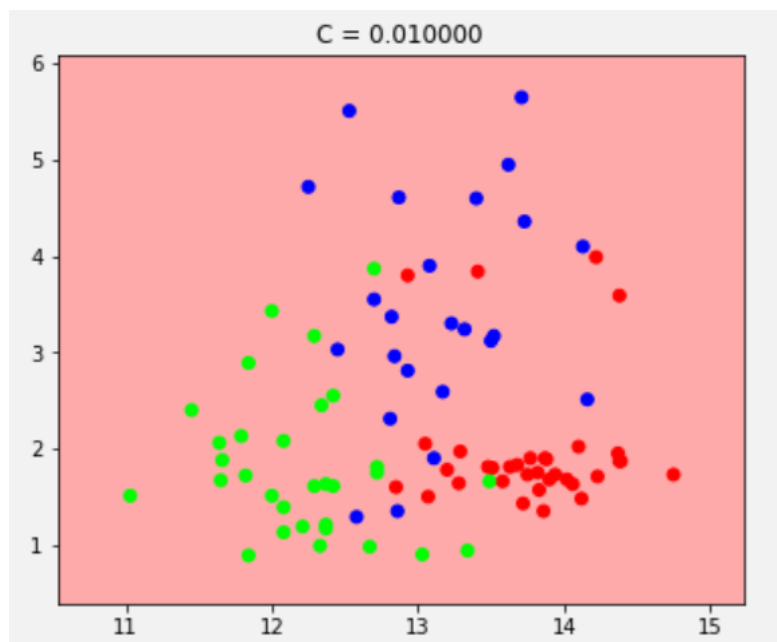
---

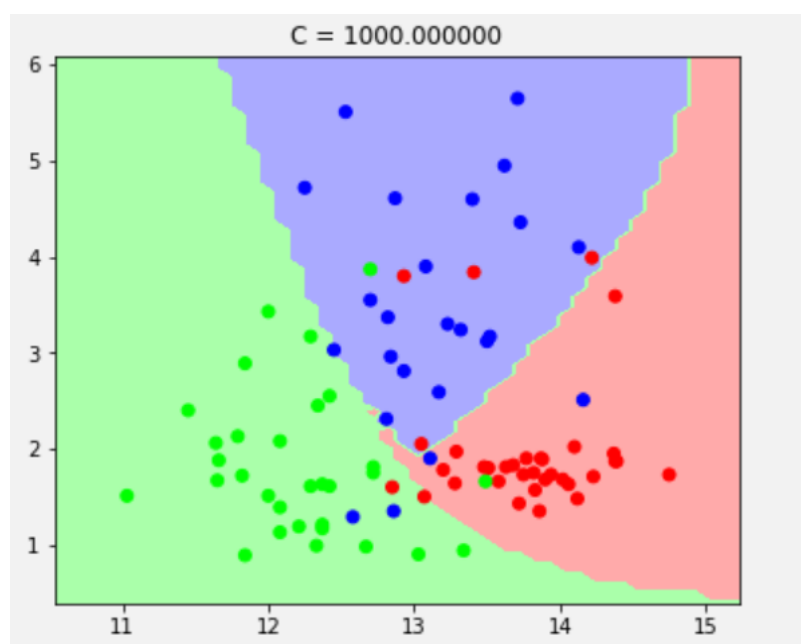
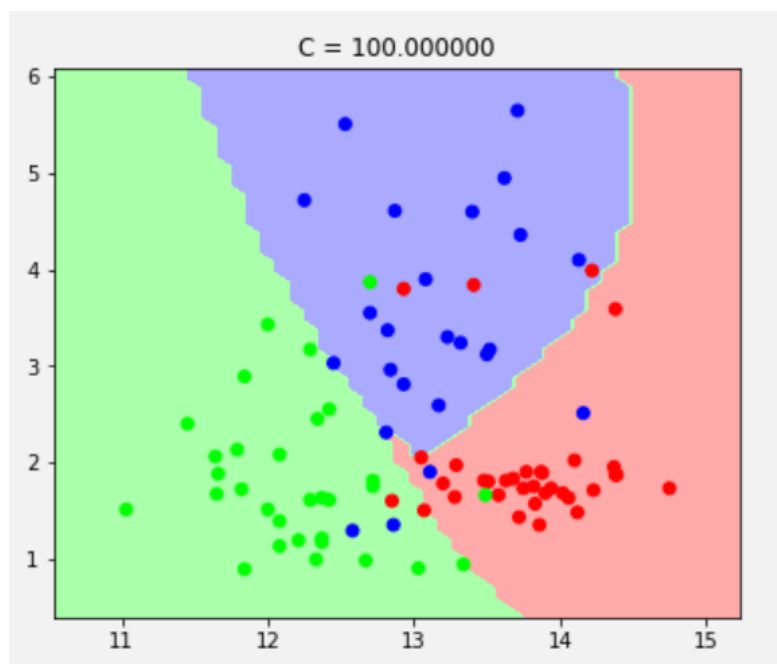
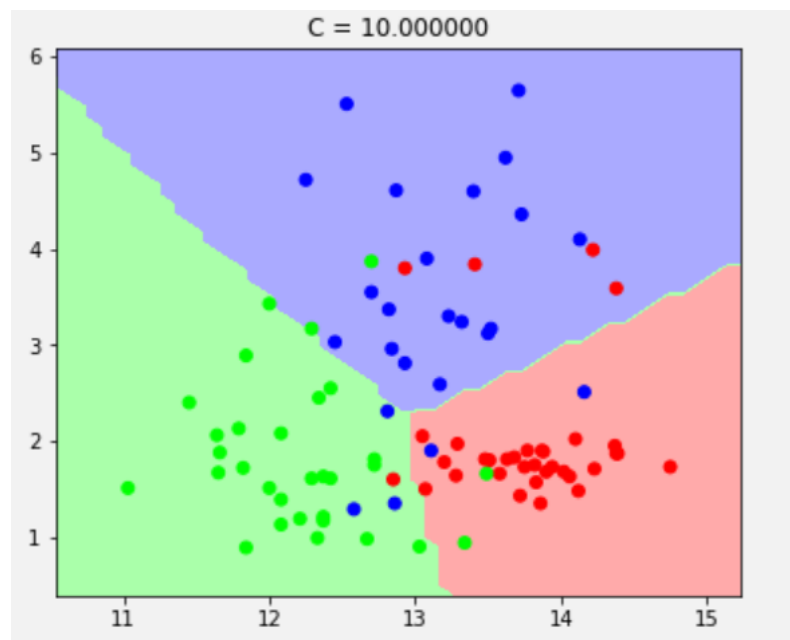
## 12. Repeat point 8. (train, plot) but using RBF kernel

---









**13. Evaluate the best  $C$  on the test set**

---

```

1 accuracies # [20.0, 20.0, 20.0, 60.0, 71.42857142857143, 68.57142857142857,
2   77.14285714285715]
3 y_test_predicted = clf_list[6].predict(X_test)
4 accuracy = np.mean(y_test_predicted == y_test) * 100
5 accuracy # 81.48148148148148

```

## 14. Differences compared to the linear kernel

The boundaries become non-linear because of the kernel.

## 15. Perform grid search for both gamma and C at the same time

By checking the source code of `svm.SVC()`, I calculate two default values of gamma

```

1 # value of gamma='scale' in SVC
2 1 / (X_train.shape[0] * X_train.var()) # 0.000377073287583712
3
4 # value of gamma='auto' in SVC
5 1 / (X_train.shape[0]) # 0.011235955056179775

```

Therefore I try gamma from 1e-5 to 1, and C from 0.1 to 1000

```

1 best_score = 0
2 for gamma in [1e-5, 1e-4, 1e-3, 1e-2, 0.1, 1]:
3     for C in [0.1, 1, 10, 100, 1000]:
4         clf = svm.SVC(gamma=gamma, C=C)
5         clf.fit(X_train, y_train)
6         score = clf.score(X_val, y_val)
7         if score > best_score:
8             best_score = score
9             best_parameters = {'gamma': gamma, 'C': C}
10            best_clf = deepcopy(clf)

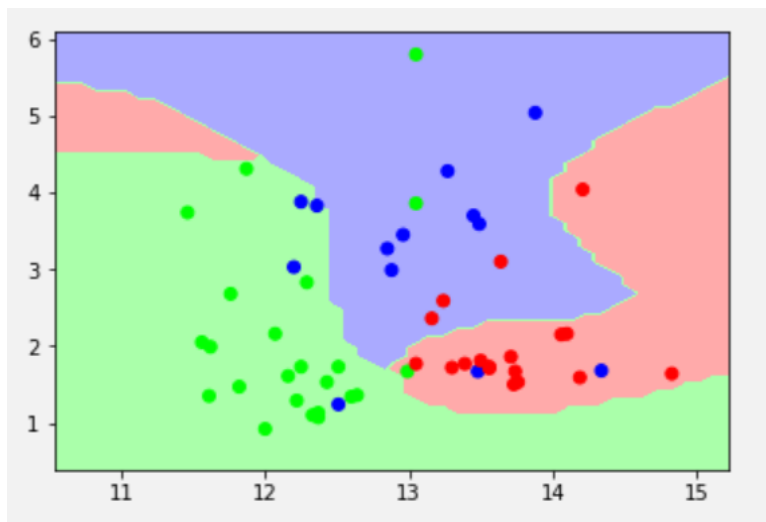
```

Best parameters are as follow:

```

1 best_score, best_parameters # (0.8571428571428571, {'C': 1000, 'gamma': 0.1})
2 best_clf.score(X_test, y_test) # 0.7777777777777778

```



## K-Fold

### 16. Merge training and validation split

```
1 X_train.shape, X_val.shape # ((89, 2), (35, 2))
2
3 X_train = np.concatenate((X_train, X_val), axis=0)
4 y_train = np.concatenate((y_train, y_val), axis=0)
5 X_train.shape, y_train.shape # ((124, 2), (124,))
```

### 17. Grid search with 5-fold cross validation for gamma and C

`sklearn.model_selection` provides `GridSearchCV` for grid search with cross validation

```
1 params = {
2     "gamma": [1e-5, 1e-4, 1e-3, 1e-2, 0.1, 1],
3     "C": [0.1, 1, 10, 100, 1000]
4 }
5 grid = GridSearchCV(svm.SVC(), params, cv=5)
6 grid.fit(X_train, y_train)
```

### 18. Evaluate on test set

```
1 grid.best_score_, grid.best_params_ # (0.8226666666666667, {'C': 1000,
2   'gamma': 0.1})
3 grid.score(X_test, y_test) # 0.7777777777777778
```

# Extra

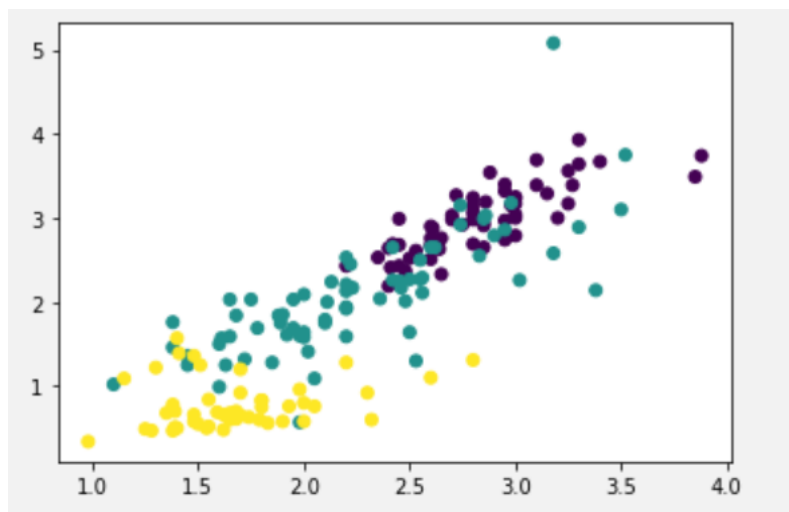
## 19. Discuss the difference between KNN and SVM

1. SVM is less affected by outliers than KNN
2. Once SVM is trained, we can quickly predict labels, while we don't have a training phase on KNN and we have to calculate distances with neighbors every time new data comes in

## 20. Try also with different pairs of attributes

Methodology is the same with above. Using `Gridsearchcv` is convenient and the result is reliable

I choose the 5th and 6th dimensions to perform SVM with RBF kernel



```
1 grid.best_score_, grid.best_params_ # (0.7993333333333333, {'C': 1000,
  'gamma': 1})
2 grid.score(X_test, y_test) # 0.7592592592592593
```

Decision boundary:

