# Data Mining & Machine Learning
## Computer Exercise 8 - Support Vector Machines

Steinarr Hrafn Höskuldsson

October 2022

## Section 1.1
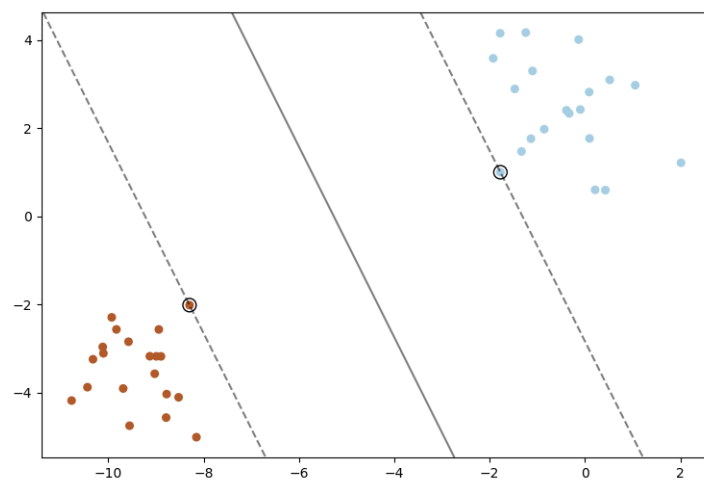


Figure 1: 1_1_1.png The decision boundary and Support Vector margins for generated data.

## Section 1.2

In the support vector machine plotted in Figure 1 there is one support vector for each of the two classes. Thus the decision boundary is a straight line.
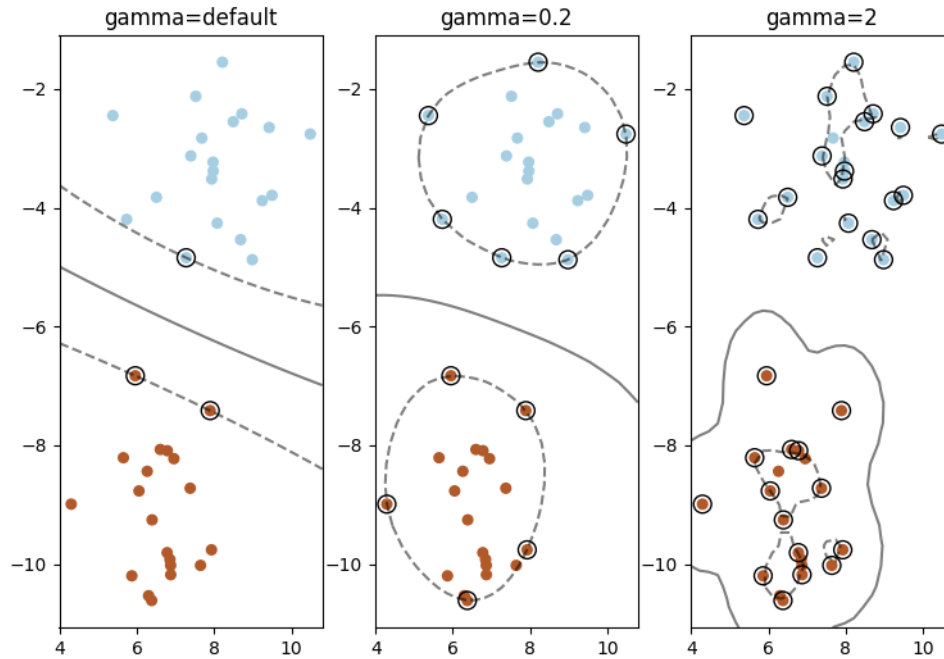
# Section 1.3



Figure 2: 1_3_1.png The decision boundary and Support Vector margins for three Support Vector machines that differ only in the value of *gamma*.

# Section 1.4

When plotting the decision boundary for different values of gamma as can be seen in Figure 2 the amount of support vectors were printed:

```
For default gamma, number of support vectors is: [1 2]
For gamma=0.2, number of support vectors is: [6 5]
For gamma=2, number of support vectors is: [18 15]
```

As can be seen on the plots the shape of the decision boundary gets more and more complicated, With default gamma it is a simple curve with a constant curvature, with gamma=0.2 it appears to be a curve with changing curvature and with gamma=2 it is a very complicated shape.

# Section 1.5
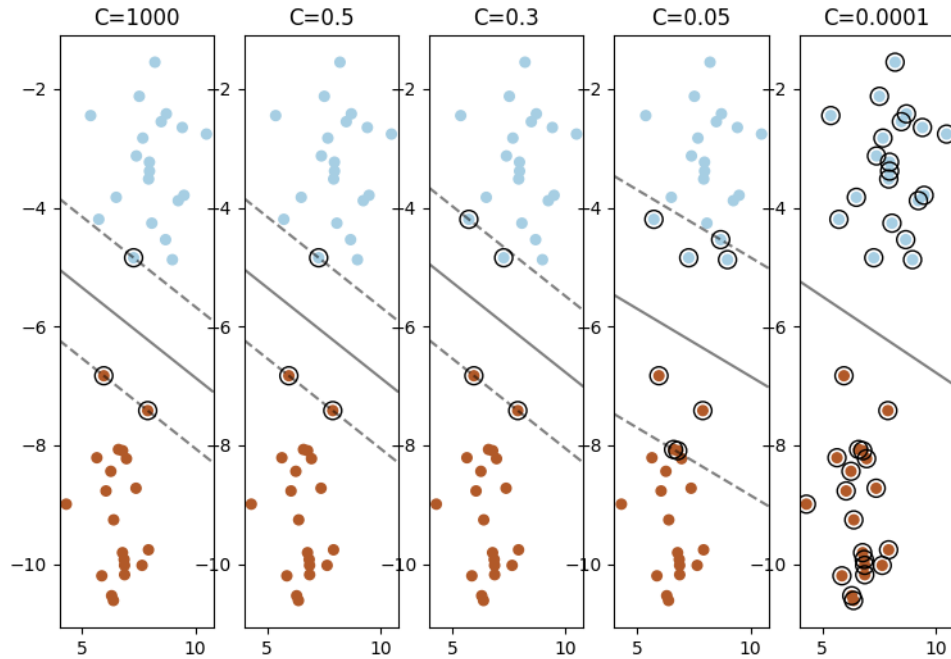


Figure 3: 1_5_1.png The decision boundary and Support Vector margins for five Support Vector machines that differ only in the value of C

# Section 1.6

When plotting the decision boundary for different values of C as can be seen in Figure 3 the amount of support vectors were printed:

```
For C=1000, number of support vectors is: [1 2]
For C=0.5, number of support vectors is: [1 2]
For C=0.3, number of support vectors is: [2 2]
For C=0.05, number of support vectors is: [4 4]
For C=0.0001, number of support vectors is: [20 20]
```

As C decreases there are more support vectors inside the margins until at $C = 0.0001$ all the datapoints are support vectors inside the margins. No support vectors are misclassified although the soft margin classifier being used does allow for that. But since the dataset is linearly seperable it does not occur.

## Section 2.2

The different kernels were tested on the breast cancer dataset. The results can be found in Table 1

|              | Accuracy | Test Score | Recall |
|--------------|----------|------------|--------|
| Linear       | 95.9%    | 95.4%      | 97.3%  |
| Radial Basis | 94.7%    | 94.7%      | 97.3%  |
| Polynomial   | 93.6%    | 93.8%      | 96.4%  |

Table 1: The results of Support Vector machine with different kernels tested on the breast cancer dataset.

Looking at the results it looks like the linear kernel gave slightly better results.

## Independent

A collection of monthly average weather observations in Reykjavík from 1949 to 2022 was collected from the Icelandic Meteorological Office's website. The objective being to train a classifier to guess the month based on Average weather observations. A short script was written to parse the values for Average Temperature, Average Pressure, Average Sun hours, Average Wind Speed and The number for the corresponding month. These were then parsed into a train test splitter.

A Support vector machine was trained on the training data and then made to predict the month from the test data. The Accuracy of the model was around 70%. Here the confusion matrix can be seen:

|           | Jan. | Feb. | March | April | May | June | July | Aug. | Sept. | Oct. | Nov. | Dec. |
|-----------|------|------|-------|-------|-----|------|------|------|-------|------|------|------|
| January   | 6    | 0    | 3     | 3     | 1   | 0    | 0    | 0    | 0     | 0    | 0    | 0    |
| February  | 0    | 11   | 0     | 0     | 0   | 2    | 0    | 0    | 0     | 0    | 0    | 0    |
| March     | 0    | 0    | 6     | 0     | 4   | 0    | 0    | 0    | 0     | 0    | 0    | 0    |
| April     | 1    | 0    | 0     | 15    | 0   | 0    | 0    | 0    | 0     | 0    | 0    | 0    |
| May       | 0    | 0    | 2     | 0     | 13  | 1    | 0    | 0    | 0     | 0    | 0    | 0    |
| June      | 0    | 2    | 0     | 0     | 4   | 11   | 0    | 0    | 0     | 0    | 0    | 0    |
| July      | 0    | 1    | 0     | 0     | 0   | 4    | 10   | 0    | 0     | 0    | 0    | 1    |
| August    | 0    | 0    | 0     | 0     | 0   | 0    | 2    | 12   | 0     | 0    | 0    | 0    |
| September | 0    | 0    | 0     | 0     | 0   | 0    | 0    | 1    | 8     | 0    | 0    | 2    |
| October   | 0    | 0    | 0     | 0     | 0   | 0    | 0    | 0    | 0     | 13   | 5    | 1    |
| November  | 0    | 0    | 0     | 0     | 0   | 0    | 0    | 0    | 2     | 10   | 5    | 0    |
| December  | 0    | 1    | 0     | 0     | 0   | 0    | 0    | 1    | 0     | 0    | 0    | 12   |

More than half inaccuracies are only off by one month and, interestingly, the model seems to have had particular issues with classifying between October and November.

# Appendix

## A Code

Listing 1: The code used

```python
1  # Author:
2  # Date:
3  # Project:
4  # Acknowledgements:
5  #
6
7  # NOTE: Your code should NOT contain any main functions or code that is executed
8  # automatically.  We ONLY want the functions as stated in the README.md.
9  # Make sure to comment out or remove all unnecessary code before submitting.
10
11
12 from tools import plot_svm_margin, load_cancer
13 from sklearn import svm
14 from sklearn.datasets import make_blobs
15 from sklearn.metrics import accuracy_score, precision_score, recall_score,
      confusion_matrix
16
17 import numpy as np
18 import matplotlib.pyplot as plt
19
20
21 def _plot_linear_kernel(show=True):
22     X, t = make_blobs(40, centers=2)
23
24     svc = svm.SVC(C=1000, kernel="linear")
25     svc.fit(X, t)
26     print(" number of support vectors for each class:", svc.n_support_)
27     print("thus the shape of decision boundary is a line")
28
29     plot_svm_margin(svc, X, t)
30     if show:
31         plt.show()
32
33 if __name__ == "_Q_main__":
34     _plot_linear_kernel(show=False)
35     plt.gcf().set_figheight(6)
36     plt.gcf().set_figwidth(9)
37     plt.savefig("08_SVM/1_1_1.png")
38     plt.show()
39
40
41 def _subplot_svm_margin(
42     svc,
43     X: np.ndarray,
44     t: np.ndarray,
45     num_plots: int,
46     index: int,
47     title=None
48 ):
49     '''
50     Plots the decision boundary and decision margins
51     for a dataset of features X and labels t and a support
52     vector machine svc.
```

5

```
53
54         Input arguments:
55         * svc: An instance of sklearn.svm.SVC: a C-support Vector
56         classification model
57         * X: [N x f] array of features
58         * t: [N] array of target labels
59         '''
60         # similar to tools.plot_svm_margin but added num_plots and
61         # index where num_plots should be the total number of plots
62         # and index is the index of the current plot being generated
63         plt.subplot(1, num_plots, index)
64         if title is not None:
65             plt.gca().title.set_text(title)
66         plot_svm_margin(svc, X, t)
67
68
69  def _compare_gamma(show=True):
70         X, t = make_blobs(n_samples=40, centers=2, random_state=6)
71
72         clf = svm.SVC(C=1000, kernel="rbf")
73         clf.fit(X,t)
74         print(f"For default gamma, number of support vectors is: {clf.n_support_}")
75         _subplot_svm_margin(clf, X, t, 3, 1, title="gamma=default")
76
77         clf = svm.SVC(C=1000, kernel="rbf", gamma=0.2)
78         clf.fit(X, t)
79         print(f"For gamma=0.2, number of support vectors is: {clf.n_support_}")
80         _subplot_svm_margin(clf, X, t, 3, 2, title="gamma=0.2")
81
82
83         clf = svm.SVC(C=1000, kernel="rbf", gamma=2, )
84         clf.fit(X, t)
85         print(f"For gamma=2, number of support vectors is: {clf.n_support_}")
86         _subplot_svm_margin(clf, X, t, 3, 3, title="gamma=2")
87
88         if show:
89             plt.show()
90
91  if __name__ == "_Q_main__":
92         _compare_gamma(show=False)
93         plt.gcf().set_figheight(6)
94         plt.gcf().set_figwidth(9)
95         plt.savefig("08_SVM/1_3_1.png")
96         plt.show()
97
98
99  def _compare_C(show=True):
100        X, t = make_blobs(n_samples=40, centers=2, random_state=6)
101
102        Cs = [1000, 0.5, 0.3, 0.05, 0.0001]
103        for i, C in enumerate(Cs):
104            clf = svm.SVC(C=C, kernel="linear")
105            clf.fit(X, t)
106            print(f"For C={C}, number of support vectors is: {clf.n_support_}")
107            _subplot_svm_margin(clf, X, t, len(Cs), i+1, title=f"{C=}")
108        if show:
109            plt.show()
110
111  if __name__ == "__main__":
112        _compare_C(False)
```

```
113        plt.gcf().set_figheight(6)
114        plt.gcf().set_figwidth(9)
115        plt.savefig("08_SVM/1_5_1.png")
116        plt.show()
117
118
119    def train_test_SVM(
120        svc,
121        X_train: np.ndarray,
122        t_train: np.ndarray,
123        X_test: np.ndarray,
124        t_test: np.ndarray,
125    ):
126        '''
127        Train a configured SVM on <X_train> and <t_train>
128        and then measure accuracy, precision and recall on
129        the test set
130
131        This function should return (accuracy, precision, recall)
132        '''
133        svc.fit(X_train, t_train)
134        y = svc.predict(X_test)
135        return accuracy_score(t_test, y), precision_score(t_test, y), recall_score(t_test
           , y)
136
137    if __name__ == "_Q_main__":
138        (X_train, t_train), (X_test, t_test) = load_cancer()
139
140        kernels = ["linear", "rbf","poly"]
141        for kernel in kernels:
142            svc = svm.SVC(C=1000, kernel=kernel)
143            print(f"Using {kernel=}, gives: {train_test_SVM(svc, X_train, t_train, X_test
           , t_test)}")
144
145    """ Independent"""
146    if __name__ == "__main__":
147
148
149        def split_train_test(
150            features: np.ndarray,
151            targets: np.ndarray,
152            train_ratio: float = 0.8
153        ):
154            '''
155            Shuffle the features and targets in unison and return
156            two tuples of datasets, first being the training set,
157            where the number of items in the training set is according
158            to the given train_ratio
159            '''
160            p = np.random.permutation(features.shape[0])
161            features = features[p]
162            targets = targets[p]
163
164            split_index = int(features.shape[0] * train_ratio)
165
166            train_features, train_targets = features[0:split_index, :],\
167                targets[0:split_index]
168            test_features, test_targets = features[split_index:-1, :],\
169                targets[split_index: -1]
170
```

```python
171            return (train_features, train_targets), (test_features, test_targets)
172
173
174     def get_weather_data():
175         temps = []
176         pressures = []
177         suns = []
178         winds = []
179         month = []
180         with open("08_SVM/medalvedur_rvk.txt", 'r') as f:
181             for line in f.readlines():
182                 stuff = line.split("\t")
183                 temps.append(stuff[3])
184                 pressures.append(stuff[14])
185                 suns.append(stuff[16])
186                 winds.append(stuff[17])
187                 month.append(stuff[2])
188         return split_train_test(np.vstack([temps, pressures, suns, winds]).T, np.
        array(month).T)
189
190     (X_train, t_train), (X_test, t_test) = get_weather_data()
191     # kernels = ["linear", "rbf", "poly"]
192     # for kernel in kernels:
193     #     svc = svm.SVC(C=1000, kernel=kernel,)
194     #     print(f"Using {kernel=}, gives: {train_test_SVM(svc, X_train, t_train,
        X_test, t_test)}")
195
196     svc = svm.SVC(C=1000, kernel='linear')
197     svc.fit(X_train, t_train),
198     y = svc.predict(X_test)
199     confuse = confusion_matrix(t_test, y)
200
201
202     def format_matrix(matrix, environment="bmatrix", formatter=str):
203         """Format a matrix using LaTeX syntax"""
204
205         if not isinstance(matrix, np.ndarray):
206             try:
207                 matrix = np.array(matrix)
208             except Exception:
209                 raise TypeError("Could not convert to Numpy array")
210
211         if len(shape := matrix.shape) == 1:
212             matrix = matrix.reshape(1, shape[0])
213         elif len(shape) > 2:
214             raise ValueError("Array must be 2 dimensional")
215
216         classes = ("January",
217                    "February",
218                    "March",
219                    "April",
220                    "May",
221                    "June",
222                    "July",
223                    "August",
224                    "September",
225                    "October",
226                    "November",
227                    "December")
228         body_lines = [classes[i] + " & " +
```

```
229                        " & ".join(map(formatter, row)) for i, row in enumerate(matrix)]
230
231         body = "\\\\\n".join(body_lines)
232         return f"""\\begin{{environment}}
233         {"(empty) &" + " & ".join(classes)}\\\\
234         {body}
235         \\end{{environment}}"""
236
237     print(format_matrix(confuse))
238     print(t_test.shape)
```