

## Solution 1

---

### Step 1

The general form for a linear predictor is the following:

$$f(x) = w^T x + b$$

We are given three linear functions for each class:

$$\text{Class 1: } w_1 = (1, 1), b_1 = 0 \rightarrow f_1(x) = x_1 \cdot (w_1^1) + x_2 \cdot (w_2^1) + b_1 = x_1 \cdot (1) + x_2 \cdot (1) + 0 = x_1 + x_2$$

$$\text{Class 2: } w_2 = (1, 0), b_2 = 1 \rightarrow f_2(x) = x_1 \cdot (w_1^2) + x_2 \cdot (w_2^2) + b_2 = x_1 \cdot (1) + x_2 \cdot (0) + 1 = x_1 + 1$$

$$\text{Class 3: } w_3 = (0, 1), b_3 = -1 \rightarrow f_3(x) = x_1 \cdot (w_1^3) + x_2 \cdot (w_2^3) + b_3 = x_2 \cdot (0) + x_1 \cdot (1) - 1 = x_2 - 1$$

Hence, for each class, we have the following:

$$f_1(x) = x_1 + x_2,$$

$$f_2(x) = x_1 + 1,$$

$$f_3(x) = x_2 - 1.$$

### Step 2

Compare  $f_1$ ,  $f_2$ , and  $f_3$  at the intersection points of the decision boundaries:

$$(1 \text{ vs } 2) : f_1 = f_2 \rightarrow x_1 + x_2 = x_1 + 1 \rightarrow x_2 = 1 \rightarrow (0, 1)$$

$$(1 \text{ vs } 3) : f_1 = f_3 \rightarrow x_1 + x_2 = x_2 - 1 \rightarrow x_1 = -1 \rightarrow (-1, 0)$$

$$(2 \text{ vs } 3) : f_2 = f_3 \rightarrow x_1 + 1 = x_2 - 1 \rightarrow x_2 = x_1 + 2 \rightarrow (0, 2) \text{ and } (-2, 0)$$

Hence, we have the following decision boundaries:

$$\text{Decision Boundary between Class 1 and Class 2: } x_2 = 1$$

$$\text{Decision Boundary between Class 1 and Class 3: } x_1 = -1$$

$$\text{Decision Boundary between Class 2 and Class 3: } x_2 = x_1 + 2$$

**Step 3**

Note:  $x_2 = 1$ ,  $x_1 = -1$ , and  $x_2 = x_1 + 2$  all pass through the common point  $(-1, 1)$ , and slice the plane into 6 wedges.

Now, select six test points for each wedge, and evaluate  $f_1$ ,  $f_2$ , and  $f_3$  to determine the classification.

Wedge label	Test point	$f_1(x)$	$f_2(x)$	$f_3(x)$	Classification
$A$	$(1, 4)$	5	2	3	Class 1
$B$	$(1, 2)$	3	2	1	Class 1
$C$	$(1, -1)$	0	2	-2	Class 2
$D$	$(-2, -1)$	-4	-1	-3	Class 2
$E$	$(-3, 0)$	-3	-2	-1	Class 3
$F$	$(-3, 3)$	0	-2	2	Class 3

**Step 4**

Visualize points and classify each wedge.

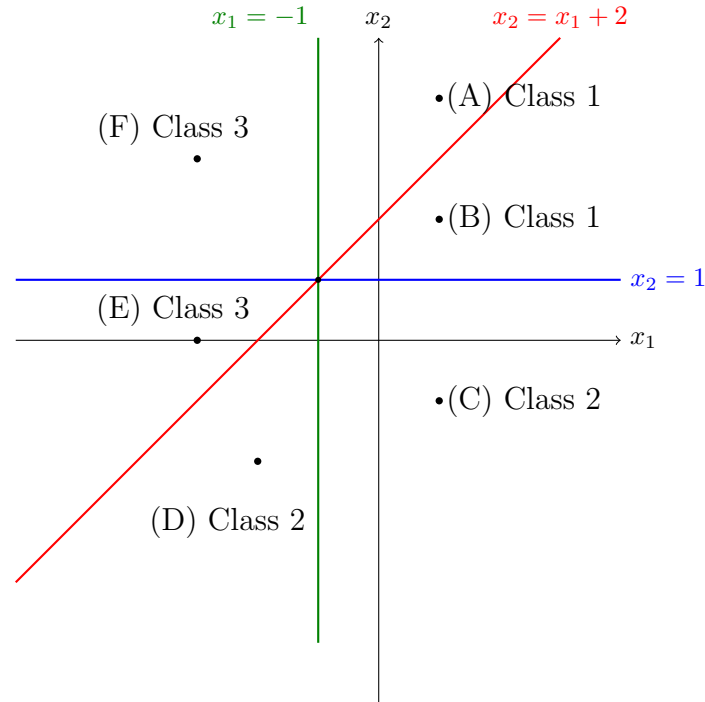


Figure 1: All six wedges produced by the three pairwise boundaries, with the winning class shown in each.

$\therefore$  the decision boundaries and the regions that belong to each class can be seen in *Figure 1*.

## Solution 2

---

### Step 1

In order to determine which option will best satisfy the statistical learning framework, we will evaluate each option.

- Option A: A web search on *American plants* will yeild subset of plant images outside of California.
- Option B: Plant photos from a similar region to that of California, will yeild similar flora but not identienal to Californian flora.
- Option C: Plant photos from the favorite Californian city will yeild photos of plants will have a higher bias towards urban Californian flora.

∴ Option B will best satisfy the statistical learning framework, since it will introduce the least amount of bias towards Californian plant identification.

---

## Solution 3

---

### Step 1

We will explore a few potential reasons why the SVM classifier performs well despite the large number of features.

- If the classes the SVM is predicting are sufficiently different, the SVM will be able to find a good decision boundary that separates the data. As a consequence of the difference between the classes, the SVM will have a large margin. This means that the SVM can find a good decision boundary even with a large number of features and minimal training points.
- The generalization error can be expressed as:

$$\text{generalisation error} \approx \sqrt{\frac{R^2/\gamma^2}{n}}$$

So we can see that as the margin grows larger, the generalization error decreases. This means that the SVM can find a good decision boundary even with a large number of features and minimal training points.

$\therefore$  a wide margin for the SVM is a key reason the 1000 training points are sufficient regardless of the one million features in the data.

---

**Solution 4 (a)**

---

The class priors changed: fewer *sports* articles, more *politics*. Within each topic (class), the way the text is written is assumed unchanged, so only the frequency of each label moves.

$\therefore$  this would be an example of a label shift

---

**Solution 4 (b)**

---

The vocabulary has drifted: new proper nouns appear and old ones vanish. This alters the distribution of the features within each topic, while the overall mix of topics (sports, politics, business, etc.) stays the same.

$\therefore$  this is an example of a covariate shift

---

## Solution 5

---

### Python Code

```
1  ## import libraries
2  import numpy as np
3  import matplotlib.pyplot as plt
4  from sklearn.utils import shuffle
5
6  ## load data0.txt
7  data = np.loadtxt('data0.txt', dtype=int)
8  x = data[:, :2]
9  y = data[:, 2]
10
11  ## get number of unique classes
12  k = len(np.unique(y))
13
14  ## define multi_perceptron function
15  def multi_perceptron(w, b, x):
16      predictions = np.dot(w, x) + b
17      best_prediction = predictions.argmax()
18      return(best_prediction)
19
20  ## define fit_multi_perceptron function
21  def fit_multi_perceptron(x, y, k, track_updates, set_seed):
22      if set_seed:
23          x, y = shuffle(x, y, random_state=42)
24      else:
25          x, y = shuffle(x, y)
26      w = np.zeros((k, x.shape[1]))
27      b = np.zeros(k)
28      updates = 0
29      max_updates = 1000
30
31      def make_prediction(w, b, x, y, updates):
32          error = False
33          for xi, yi in zip(x, y):
34              predicted_y = multi_perceptron(w, b, xi)
35              if predicted_y != yi:
36                  w[yi] += xi
37                  b[yi] += 1
38                  w[predicted_y] -= xi
39                  b[predicted_y] -= 1
40                  updates += 1
41                  error = True
42          if updates <= max_updates:
43              if error:
44                  return make_prediction(w, b, x, y, updates)
45              else:
46                  return (w, b, updates)
47          else:
48              print("Did not converge after {max_iterations} iterations")
49      if track_updates:
50          w, b, updates = make_prediction(w, b, x, y, updates)
51          return (w, b, updates)
52      else:
53          w, b, updates = make_prediction(w, b, x, y, updates)
54          return (w, b)
```

## Plots

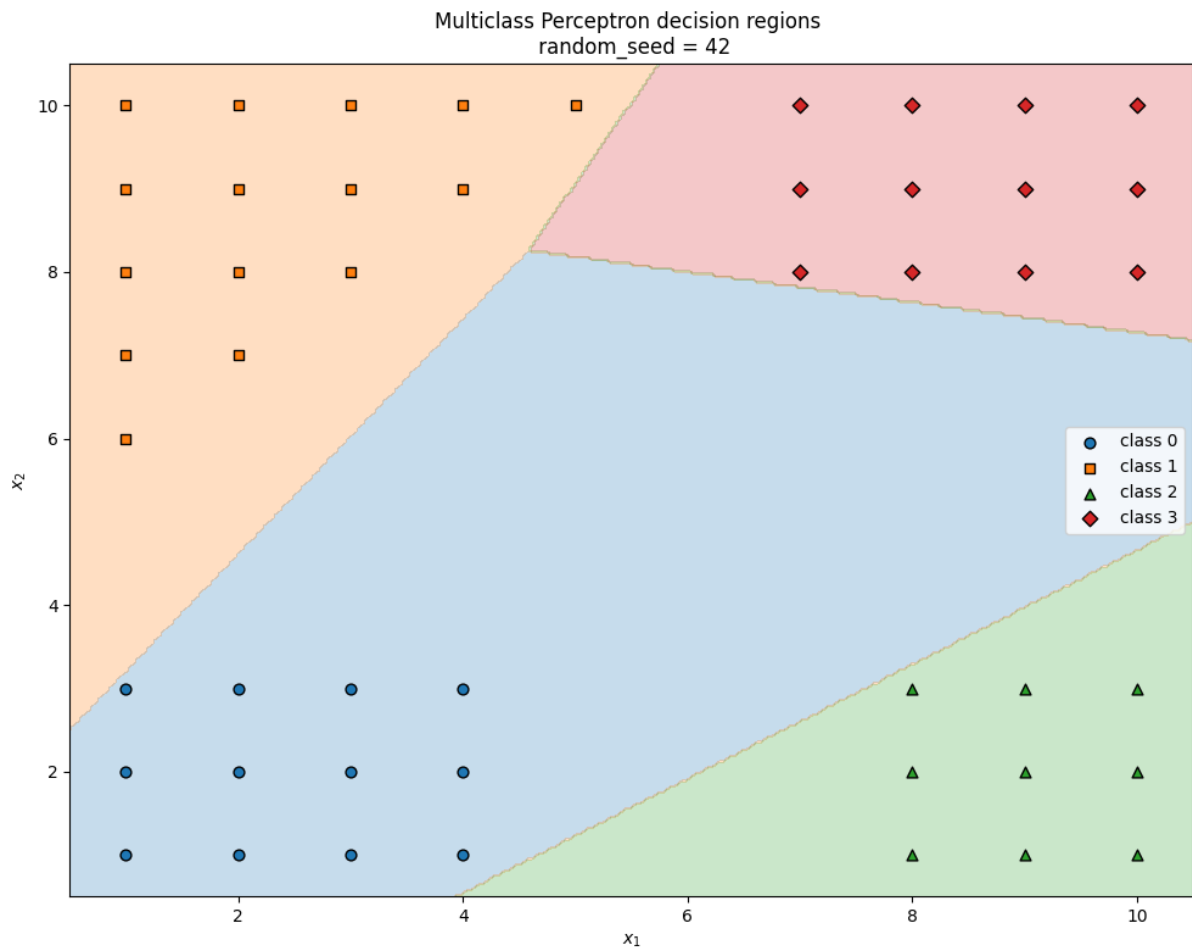


Figure 2: Multiclass Perceptron decision boundaries.

## Solution 6

### Plots

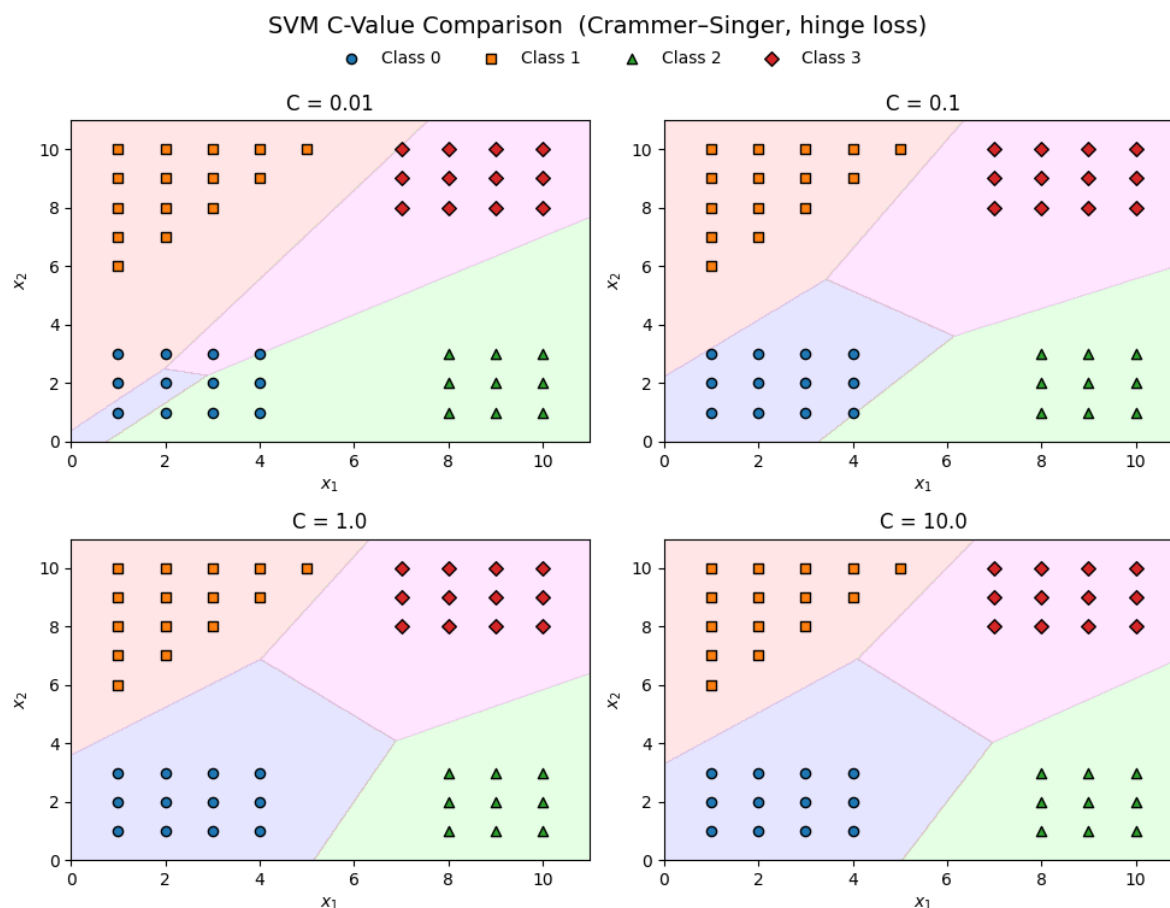


Figure 3: SVM C value comparison.

### C Value Comparison Comments

Observations across C values:

- As C increases the model partitions the data more effectively
- Since there are a limited amount of features in this dataset (*data0.txt*), we don't have to worry about minimizing the margin, since the computational cost won't skyrocket with this dataset.