# Problem 1

**d** is the correct answer.

We have this because we can vary **w** and $b$. And since **w** has $d$ variables, we have $d+1$ variables in total.

# Problem 2

**a** is the correct answer.

I wrote some code to implement SVM with soft margin on the zip-code data by solving the equations laid out in the problem set. I used binary classification error. To help me write this code, I used the libsvm library. For this problem, I made sure to use the polynomial kernel as given in the problem. After running this code for 0 versus all, 2 versus all, 4 versus all, 6 versus all, and 8 versus all, I found that 0 versus all gave me the lowest accurary and thus the highest $E_{in}$ (ran prediction code on the training data).

# Problem 3

**a** is the correct answer.

I wrote some code to implement SVM with soft margin on the zip-code data by solving the equations laid out in the problem set. I used binary classification error. To help me write this code, I used the libsvm library. For this problem, I made sure to use the polynomial kernel as given in the problem. After running this code for 1 versus all, 3 versus all, 5 versus all, 6 versus all, and 9 versus all, I found that 1 versus all gave me the highest accurary and thus the lowest $E_{in}$ (ran prediction code on the training data).

# Problem 4

**c** is the correct answer.

Using the code I wrote for the above two problems, I got 386 support vectors for 1 versus all and 2180 support vectors for 0 versus all (with all the other parameters as specified in the problem). The difference between these two numbers is around 1800.

# Problem 5

**d** is the correct answer.

I wrote some code that runs the 1 versus 5 classifier with $Q = 2$ and $C \in \{0.001, 0.01, 0.1, 1\}$. I then observed the number of support vectors, $E_{in}$, and $E_{out}$ for each $C$ value. I found that the only statement given in the problem that is true is that the maximum $C$ (1) achieves the lowest $E_{in}$ and the highest accuracy (around 99.68% accuracy). I found that all the other statements given in the problem are false.

# Problem 6

**b** is the correct answer.

I wrote some code that runs the 1 versus 5 classifier with $Q = 2$ and $C \in \{0.0001, 0.001, 0.01, 0.1, 1\}$ and then runs the 1 versus 5 classifier with $Q = 5$ and $C \in \{0.0001, 0.001, 0.01, 0.1, 1\}$. I did this so I could compare values between all the runs. I found that the only statement given in the problem that is true is that when $C = 0.001$, the number of support vectors is lower at $Q = 5$ (76 versus 25).

# Problem 7

**b** is the correct answer.

I wrote some code that runs 10-fold cross validation for the polynomial kernel. With this code, I considered the 1 versus 5 classifier with $Q = 2$. I used $E_{cv}$ to select $C \in \{0.0001, 0.001, 0.01, 0.1, 1\}$. If there was a tie in $E_{cv}$, I selected the smaller $C$. I then ran this code for 100 random runs (so I tried 100 different partitions). With this, I discovered that $C = 0.001$ is selected most often.

# Problem 8

**c** is the correct answer.

I basically used the same code I wrote for the above problem, and for each run, I added the value of $E_{cv}$ for $C = 0.001$. I then took the average value of this list, and it was around 0.005.

# Problem 9

**e** is the correct answer.

I wrote code that runs SVM with the RBF kernel for all the values of $C$ given in the problem. I ran this code on the 1 versus 5 classifier. To help me write this code, I used the libsvm library. The value $C = 10^6$ gave me the lowest $E_{in}$ value (ran prediction code on the training data).

# Problem 10

**c** is the correct answer.

I wrote code that runs SVM with the RBF kernel for all the values of $C$ given in the problem. I ran this code on the 1 versus 5 classifier. To help me write this code, I used the libsvm library. The value $C = 100$ gave me the lowest $E_{out}$ value (ran prediction code on the test data).

# Code

```
1  from __future__ import division
2  from svm import *
3  from svmutil import *
4  from collections import Counter
5  import sys
6  import random
7  import numpy as np
8
9  class MySvm:
10     def __init__(self, file_train, file_test):
11         self.training_labels = self.load_labels(file_train)
12         self.training_data = self.load_data(file_train)
13         self.training_data_curr = self.training_data
14         self.scores_tr = []
15
16         self.test_labels = self.load_labels(file_test)
17         self.test_data = self.load_data(file_test)
18         self.test_data_curr = self.test_data
19         self.scores_test = []
20
21         self.cross_data_curr_list = []
22         self.scores_cross_list = []
23         self.scores_tr_list = []
24         self.training_data_curr_list = []
25
26         self.model = None
27         self.model_list = []
28         self.num_support_vectors = 0
29
30     def load_labels(self, file_name):
31         file_obj = open(file_name)
32         labels = []
33         for line in file_obj:
34             labels.append(float(line.split()[0]))
35         return labels
36
37     def load_data(self, file_name):
38         file_obj = open(file_name)
39         data = []
40         for line in file_obj:
41             item = []
42             line_split = line.split()
43             for i in range(1, len(line_split)):
44                 item.append(float(line_split[i]))
45             data.append(item)
46         return data
47
48     def one_versus_all(self, num_one, kernel_type, error_const, poly_degree,
49             is_cross, is_repeat = False):
50         if not is_repeat:
51             self.scores_tr, self.training_data_curr = self.get_one_versus_all_lists(
52                     num_one, self.training_labels, self.training_data)
53             self.scores_test, self.test_data_curr = self.get_one_versus_all_lists(
54                     num_one, self.test_labels, self.test_data)
55             if is_cross:
56                 self.scores_tr, self.training_data_curr, self.scores_cross, self.cross_data
57
58         prob = svm_problem(self.scores_tr, self.training_data_curr)
59         param_str = '-t %d -r 1 -g 1 -c %f -d %d' %(kernel_type, error_const, poly_degree)
60         param = svm_parameter(param_str)
61         model = svm_train(prob, param)
```

```
62             self.model = model
63             self.num_support_vectors = len(model.get_SV())
64
65       def one_versus_one(self, num_one, num_other, kernel_type, error_const,
66                 poly_degree, is_cross, is_repeat = False):
67           if not is_repeat:
68               self.scores_tr, self.training_data_curr = self.get_one_versus_one_lists(
69                       num_one, num_other, self.training_labels, self.training_data)
70               self.scores_test, self.test_data_curr = self.get_one_versus_one_lists(
71                       num_one, num_other, self.test_labels, self.test_data)
72               if is_cross:
73                   self.scores_tr_list, self.training_data_curr_list, self.scores_cross_list,
74
75           if not is_cross:
76               prob = svm_problem(self.scores_tr, self.training_data_curr)
77               param_str = '-t %d -r 1 -g 1 -c %f -d %d' %(kernel_type, error_const, poly_degr
78               param = svm_parameter(param_str)
79               model = svm_train(prob, param)
80               self.model = model
81               self.num_support_vectors = len(model.get_SV())
82           else:
83               self.model_list = []
84               sum_num_support_vectors = 0
85               for i in range(0, len(self.scores_tr_list)):
86                   prob = svm_problem(self.scores_tr_list[i], self.training_data_curr_list[i])
87                   param_str = '-t %d -r 1 -g 1 -c %f -d %d' %(kernel_type, error_const, poly_
88                   param = svm_parameter(param_str)
89                   model = svm_train(prob, param)
90                   self.model_list.append(model)
91                   sum_num_support_vectors += len(model.get_SV())
92               self.num_support_vectors = sum_num_support_vectors / len(self.scores_tr_list)
93
94       def get_one_versus_all_lists(self, num_one, labels, data_items):
95           scores = []
96           for label in labels:
97               if label == num_one:
98                   scores.append(1)
99               else:
100                  scores.append(-1)
101          return (scores, data_items)
102
103      def get_one_versus_one_lists(self, num_one, num_other, labels, data_items):
104          scores = []
105          data_curr = []
106          for i in range(0, len(labels)):
107              label = labels[i]
108              item = data_items[i]
109              if label == num_one:
110                  scores.append(1)
111                  data_curr.append(item)
112              elif label == num_other:
113                  scores.append(-1)
114                  data_curr.append(item)
115          return (scores, data_curr)
116
117      def get_cross_val_lists(self):
118          z_list = list(zip(self.scores_tr, self.training_data_curr))
119          random.shuffle(z_list)
120          shuffled_scores_tr, shuffled_training_data_curr = zip(*z_list)
121          index = int(len(self.scores_tr) / 10)
122
123          scores_tr_list, training_data_curr_list, scores_cross_list, cross_data_curr_list =
124          start = 0
```

```
125          end = index
126          while end < len(self.scores_tr):
127              if end + index >= len(self.scores_tr):
128                  end = len(self.scores_tr)
129              scores_cross_list.append(list(shuffled_scores_tr[start:end]))
130              cross_data_curr_list.append(list(shuffled_training_data_curr[start:end]))
131              scores_tr_list.append(list(shuffled_scores_tr[0:start]) + list(shuffled_scores_
132              training_data_curr_list.append(list(shuffled_training_data_curr[0:start]) + lis
133              start += index
134              end += index
135          return (scores_tr_list, training_data_curr_list, scores_cross_list, cross_data_curr
136
137      def get_error_in(self):
138          p_labels, p_acc, p_vals = svm_predict(self.scores_tr,
139                  self.training_data_curr, self.model)
140          return p_acc[0]
141
142      def get_error_out(self):
143          p_labels, p_acc, p_vals = svm_predict(self.scores_test,
144                  self.test_data_curr, self.model)
145          return p_acc[0]
146
147      def get_error_cv(self):
148          sum_acc = 0
149          for i in range(0, len(self.scores_cross_list)):
150              p_labels, p_acc, p_vals = svm_predict(self.scores_cross_list[i],
151                      self.cross_data_curr_list[i], self.model_list[i])
152              sum_acc += p_acc[0]
153          print '========LENGTH = ', len(self.scores_cross_list)
154          return sum_acc / len(self.scores_cross_list)
155
156 if __name__ == '__main__':
157      prob1_3 = False
158      prob5 = False
159      prob6 = False
160      prob7_8 = True
161      prob9_10 = False
162      print 'Kshit is gay'
163      my_svm = MySvm('features.train', 'features.test')
164
165      if prob1_3:
166          for i in range(1, 11, 2):
167              print 'i = ', i
168              my_svm.one_versus_all(i, 1, .01, 2, False)
169              p_acc = my_svm.get_error_in()
170              print 'Num support vectors = ', my_svm.num_support_vectors
171              print '\n\n'
172
173      if prob5:
174          print 'Q = ', 2
175          for c in [.0001, .001, .01, .1, 1]:
176              print 'C = ', c
177              my_svm.one_versus_one(1, 5, 1, c, 2, False)
178              print 'E in'
179              my_svm.get_error_in()
180              print 'E out'
181              my_svm.get_error_out()
182              print 'Num support vectors = ', my_svm.num_support_vectors
183              print '\n'
184          print '\n\n\n'
185
186      if prob6:
187          print 'Q = ', 5
```

```
188          for c in [.0001, .001, .01, .1, 1]:
189              print 'C = ', c
190              my_svm.one_versus_one(1, 5, 1, c, 5, False)
191              print 'E in'
192              my_svm.get_error_in()
193              print 'E out'
194              my_svm.get_error_out()
195              print 'Num support vectors = ', my_svm.num_support_vectors
196              print '\n'
197
198      if prob7_8:
199          c_list = []
200          e_dict = {.0001 : [], .001 : [], .01 : [], .1 : [], 1 : []}
201          for i in range(0, 100):
202              min_error = sys.maxint
203              c_val = 1
204              repeat = False
205              for c in [.0001, .001, .01, .1, 1]:
206                  my_svm.one_versus_one(1, 5, 1, c, 2, True, repeat)
207                  repeat = True
208                  # Get error, NOT accurary
209                  error = (100 - my_svm.get_error_cv()) / 100
210                  e_dict[c].append(error)
211                  if error < min_error:
212                      min_error = error
213                      c_val = c
214              c_list.append(c_val)
215          c_data = Counter(c_list)
216          print 'Counts = ', c_data.most_common()
217          print 'Mode = ', c_data.most_common(1)
218          for key, val in e_dict.iteritems():
219              new_val = []
220              val_max = max(val)
221              val_min = min(val)
222              val_mean = sum(val) / len(val)
223              new_val = [val_mean, val_max, val_min]
224              e_dict[key] = new_val
225          print e_dict
226
227      if prob9_10:
228          min_error_in = sys.maxint
229          min_error_out = sys.maxint
230          c_val_in = 0
231          c_val_out = 0
232          for c in [.01, 1, 100, 10000, 1000000]:
233              my_svm.one_versus_one(1, 5, 2, c, 2, False, False)
234              error_in = (100 - my_svm.get_error_in()) / 100
235              if error_in < min_error_in:
236                  min_error_in = error_in
237                  c_val_in = c
238              error_out = (100 - my_svm.get_error_out()) / 100
239              if error_out < min_error_out:
240                  min_error_out = error_out
241                  c_val_out = c
242          print 'Min error in = ', min_error_in
243          print 'C value in = ', c_val_in
244          print 'Min error out = ', min_error_out
245          print 'C value out = ', c_val_out
```