

Problem 1: PCA and Feature Selection

SVMs and PCA

```
In [1]: import numpy as np
import matplotlib.pyplot as plt
import math
import pandas as pd
import cvxopt
from tqdm.notebook import trange
from sklearn.svm import SVC
```

```
In [2]: sonar_train = pd.read_csv('sonar_train.data', header=None)
sonar_test = pd.read_csv('sonar_test.data', header=None)
sonar_valid = pd.read_csv('sonar_valid.data', header=None)

sonar_train.loc[sonar_train[60] == 2, 60] = -1
sonar_test.loc[sonar_test[60] == 2, 60] = -1
sonar_valid.loc[sonar_valid[60] == 2, 60] = -1

def normalize(data, mean, std):
    return (data - mean) / std

def split_data(data):
    return data.iloc[:, :60].to_numpy(), data.iloc[:, 60:].to_numpy()

X, y_train = split_data(sonar_train)
train_mean = X.mean(axis=0)
train_std = X.std(axis=0)

X_train = normalize(X, train_mean, train_std)
X, y_validation = split_data(sonar_valid)
X_validation = normalize(X, train_mean, train_std)
X, y_test = split_data(sonar_test)
X_test = normalize(X, train_mean, train_std)
```

Perform PCA on the training data to reduce the dimensionality of the data set (ignoring the class labels for the moment). What are the top six eigenvalues of the data covariance matrix?

```
In [3]: def compute_covariance(norm_data):
    # Covariance matrix has dimensions (p x p)
    # Usually computed with variables as rows and observations as columns (np.cov)
    return norm_data.T.dot(norm_data)

covariance_mat = compute_covariance(X_train) #np.cov(X_train.T)
```

```
In [4]: e_val, e_vec = np.linalg.eig(covariance_mat)
i_rev = e_val.argsort()[::-1]
eig_vals = e_val[i_rev]
eig_vecs = e_vec[:, i_rev]
print("Top 6 Eigen values", eig_vals[:6])
```

```
Top 6 Eigen values [1344.02076581 1196.11970978  541.52327862  351.33131561
313.13028781
267.04094634]
```

For each $k \in \{1, 2, 3, 4, 5, 6\}$, project the training data into the best k dimensional subspace (with respect to the Frobenius norm) and use the SVM with slack formulation to learn a classifier for each $c \in \{1, 10, 100, 1000\}$. Report the error of the learned classifier on the validation set for each k and c pair.

```
In [5]: C = [1, 10, 100, 1000]
```

```
In [6]: data = {
    'k': [],
    'c': [],
    'Training Data Error': [],
    'Validation Data Error': []
}
for k in range(1,7):
    U = eig_vecs[:, :k]
    X_proj = X_train.dot(U)
    X_valid_proj = X_validation.dot(U)
    for c in C:
        data['k'].append(k)
        data['c'].append(c)
        clf = SVC(C=c, kernel='linear')
        clf.fit(X_proj, y_train.ravel())
        y_pred = clf.predict(X_proj)
        data['Training Data Error'].append(1 - np.mean(y_pred == y_train.ravel()))
        valid_pred = clf.predict(X_valid_proj)
        data['Validation Data Error'].append(1 - np.mean(valid_pred == y_validation.ravel()))
```

In [7]: `pd.DataFrame(data)`

Out[7]:

| | k | c | Training Data Error | Validation Data Error |
|----|---|------|---------------------|-----------------------|
| 0 | 1 | 1 | 0.509615 | 0.461538 |
| 1 | 1 | 10 | 0.509615 | 0.461538 |
| 2 | 1 | 100 | 0.509615 | 0.461538 |
| 3 | 1 | 1000 | 0.509615 | 0.461538 |
| 4 | 2 | 1 | 0.432692 | 0.307692 |
| 5 | 2 | 10 | 0.432692 | 0.307692 |
| 6 | 2 | 100 | 0.432692 | 0.307692 |
| 7 | 2 | 1000 | 0.432692 | 0.307692 |
| 8 | 3 | 1 | 0.269231 | 0.211538 |
| 9 | 3 | 10 | 0.269231 | 0.211538 |
| 10 | 3 | 100 | 0.269231 | 0.211538 |
| 11 | 3 | 1000 | 0.278846 | 0.211538 |
| 12 | 4 | 1 | 0.288462 | 0.211538 |
| 13 | 4 | 10 | 0.288462 | 0.211538 |
| 14 | 4 | 100 | 0.288462 | 0.211538 |
| 15 | 4 | 1000 | 0.288462 | 0.211538 |
| 16 | 5 | 1 | 0.269231 | 0.250000 |
| 17 | 5 | 10 | 0.269231 | 0.250000 |
| 18 | 5 | 100 | 0.269231 | 0.250000 |
| 19 | 5 | 1000 | 0.269231 | 0.250000 |
| 20 | 6 | 1 | 0.240385 | 0.269231 |
| 21 | 6 | 10 | 0.240385 | 0.269231 |
| 22 | 6 | 100 | 0.240385 | 0.269231 |
| 23 | 6 | 1000 | 0.240385 | 0.269231 |

How does it compare to the best classifier (with the same possible c choices) without feature selection?

```
In [8]: data = {
        'c': [],
        'Training Data Error': [],
        'Validation Data Error': []
      }
      for c in C:
        data['c'].append(c)
        clf = SVC(C=c, kernel='linear', random_state=0)
        clf.fit(X_train, y_train.ravel())
        y_pred = clf.predict(X_train)
        data['Training Data Error'].append(1 - np.mean(y_pred == y_train.ravel()))
        valid_pred = clf.predict(X_validation)
        data['Validation Data Error'].append(1 - np.mean(valid_pred == y_validation.ravel()))
      pd.DataFrame(data)
```

Out[8]:

| | c | Training Data Error | Validation Data Error |
|---|------|---------------------|-----------------------|
| 0 | 1 | 0.009615 | 0.211538 |
| 1 | 10 | 0.000000 | 0.230769 |
| 2 | 100 | 0.000000 | 0.230769 |
| 3 | 1000 | 0.000000 | 0.230769 |

What is the error of the best k/c pair on the test data? How does it compare to the best classifier (with the same possible c choices) without feature selection? Explain your observations.

Best k = 3 and c = 1, 10, 100, 1000

```
In [9]: U = eig_vecs[:, :3]
        X_proj = X_train.dot(U)
        X_test_proj = X_test.dot(U)
        res = {
            'c': [],
            'Test data error': []
        }
        for c in [1, 10, 100, 1000]:
            res['c'].append(c)
            clf = SVC(C=c, kernel='linear', random_state=0)
            clf.fit(X_proj, y_train.ravel())
            res['Test data error'].append(1 - np.mean(clf.predict(X_test_proj) == y_test.ravel()))
        pd.DataFrame.from_dict(res)
```

Out[9]:

| | c | Test data error |
|---|------|-----------------|
| 0 | 1 | 0.192308 |
| 1 | 10 | 0.192308 |
| 2 | 100 | 0.192308 |
| 3 | 1000 | 0.192308 |

```
In [10]: # SVM without feature selection
res = {
    'c': [],
    'Test data error': []
}
for c in [1]:
    res['c'].append(c)
    clf = SVC(C=c, kernel='linear', random_state=0)
    clf.fit(X_train, y_train.ravel())
    res['Test data error'].append(1 - np.mean(clf.predict(X_test) == y_test.ra
vel()))
pd.DataFrame.from_dict(res)
```

Out[10]:

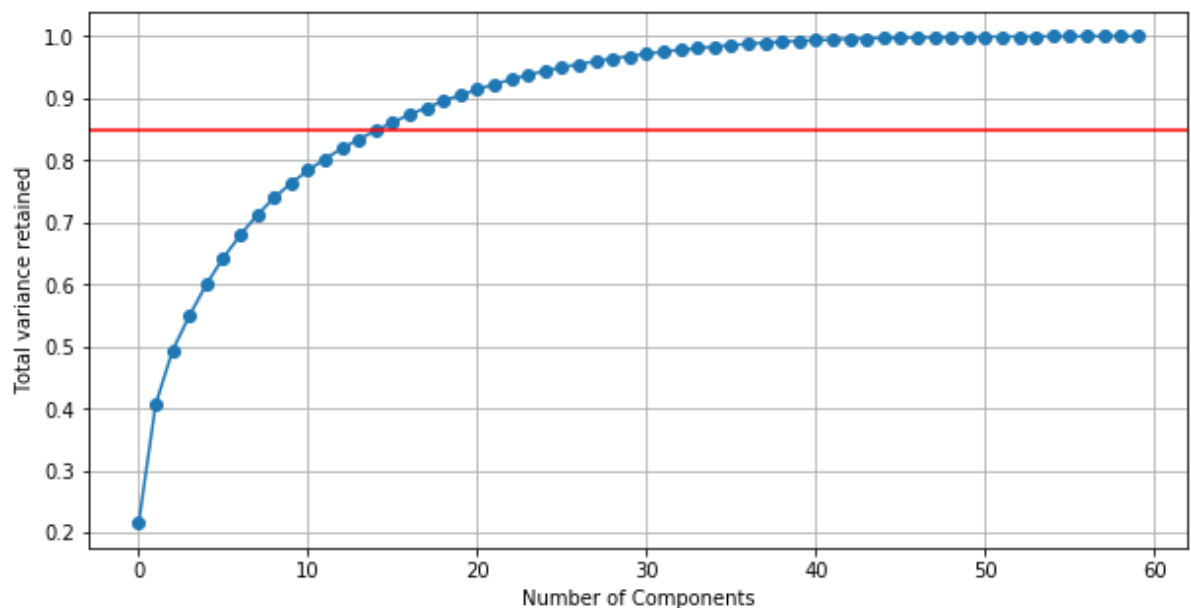
| | c | Test data error |
|---|---|-----------------|
| 0 | 1 | 0.211538 |

SVM on reduced dimension data performs better on test data compared to SVM without any feature selection

If you had to pick a value of k before evaluating the performance on the validation set (e.g., if this was not a supervised learning problem), how might you pick it?

Value of k or the number of components could be picked by heuristics. We can pick the number of components needed to explain at least 85% of the data variance. For this dataset, we can pick k = 14.

```
In [11]: fig = plt.figure(figsize=(10, 5))
plt.plot(np.cumsum(eig_vals) / np.sum(eig_vals), marker='o')
plt.xlabel('Number of Components')
plt.ylabel('Total variance retained')
plt.axhline(y=0.85, color='r', linestyle='-')
plt.grid()
```



PCA for Feature Selection

1. Compute the top k eigenvalues and eigenvectors of the covariance matrix corresponding to the data matrix omitting the labels (recall that the rows of the data matrix are the input data points).

2. Define π

3. Sample s columns independently from the probability distribution defined by π .

```
In [12]: def select_k_features(features, s, pi):  
         sel_features = np.random.choice(features, s, p=pi, replace=False)  
         return sel_features
```

Why does π define a probability distribution?

π has values ranging between 0 and 1. Also sum of all values in $\pi = 1$.

Again, using the UCI Sonar data set, for each $k \in \{1, \dots, 10\}$ and each $s \in \{1, \dots, 20\}$, report the average test error of the SVM with slack classifier over 100 experiments. For each experiment use only the s selected features (note that there may be some duplicates, so only include each feature once).

```

In [13]: data = {
    'k': [],
    's': [],
    'Average test error': []
}
datapoints, features = X_train.shape
for k in trange(1,11):
    V = eig_vecs[:k]
    pi = np.sum(V**2, axis=0) / k
    for s in trange(1,21):
        error = []
        for _ in range(100):
            sel_features = select_k_features(features, s, pi)
            clf = SVC(C=1, kernel='linear', random_state=0)
            clf.fit(X_train[:,sel_features], y_train.ravel())
            y_pred = clf.predict(X_test[:,sel_features])
            error.append(1 - np.mean(y_pred==y_test.ravel()))
        mean_error = np.mean(error)
        data['k'].append(k)
        data['s'].append(s)
        data['Average test error'].append(mean_error)
res = pd.DataFrame.from_dict(data)

```

```
In [16]: pd.set_option("display.max_rows", None, "display.max_columns", None)  
res
```


Out[16]:

| | k | s | Average test error |
|----|---|----|--------------------|
| 0 | 1 | 1 | 0.460769 |
| 1 | 1 | 2 | 0.401154 |
| 2 | 1 | 3 | 0.378269 |
| 3 | 1 | 4 | 0.367692 |
| 4 | 1 | 5 | 0.350192 |
| 5 | 1 | 6 | 0.320769 |
| 6 | 1 | 7 | 0.318846 |
| 7 | 1 | 8 | 0.318077 |
| 8 | 1 | 9 | 0.305192 |
| 9 | 1 | 10 | 0.302115 |
| 10 | 1 | 11 | 0.291731 |
| 11 | 1 | 12 | 0.285962 |
| 12 | 1 | 13 | 0.292308 |
| 13 | 1 | 14 | 0.286154 |
| 14 | 1 | 15 | 0.285385 |
| 15 | 1 | 16 | 0.283654 |
| 16 | 1 | 17 | 0.281731 |
| 17 | 1 | 18 | 0.267115 |
| 18 | 1 | 19 | 0.268269 |
| 19 | 1 | 20 | 0.268269 |
| 20 | 2 | 1 | 0.435962 |
| 21 | 2 | 2 | 0.389038 |
| 22 | 2 | 3 | 0.342885 |
| 23 | 2 | 4 | 0.343269 |
| 24 | 2 | 5 | 0.311154 |
| 25 | 2 | 6 | 0.299808 |
| 26 | 2 | 7 | 0.296731 |
| 27 | 2 | 8 | 0.297308 |
| 28 | 2 | 9 | 0.280962 |
| 29 | 2 | 10 | 0.265385 |
| 30 | 2 | 11 | 0.267885 |
| 31 | 2 | 12 | 0.266731 |
| 32 | 2 | 13 | 0.258269 |
| 33 | 2 | 14 | 0.252500 |
| 34 | 2 | 15 | 0.259423 |

| | k | s | Average test error |
|----|---|----|--------------------|
| 35 | 2 | 16 | 0.252885 |
| 36 | 2 | 17 | 0.240192 |
| 37 | 2 | 18 | 0.247500 |
| 38 | 2 | 19 | 0.247692 |
| 39 | 2 | 20 | 0.239423 |
| 40 | 3 | 1 | 0.428077 |
| 41 | 3 | 2 | 0.364231 |
| 42 | 3 | 3 | 0.347885 |
| 43 | 3 | 4 | 0.320769 |
| 44 | 3 | 5 | 0.293654 |
| 45 | 3 | 6 | 0.294423 |
| 46 | 3 | 7 | 0.282885 |
| 47 | 3 | 8 | 0.280769 |
| 48 | 3 | 9 | 0.265577 |
| 49 | 3 | 10 | 0.271154 |
| 50 | 3 | 11 | 0.264231 |
| 51 | 3 | 12 | 0.252692 |
| 52 | 3 | 13 | 0.258269 |
| 53 | 3 | 14 | 0.256731 |
| 54 | 3 | 15 | 0.248077 |
| 55 | 3 | 16 | 0.247692 |
| 56 | 3 | 17 | 0.247692 |
| 57 | 3 | 18 | 0.259423 |
| 58 | 3 | 19 | 0.242500 |
| 59 | 3 | 20 | 0.253077 |
| 60 | 4 | 1 | 0.430385 |
| 61 | 4 | 2 | 0.376346 |
| 62 | 4 | 3 | 0.331923 |
| 63 | 4 | 4 | 0.323077 |
| 64 | 4 | 5 | 0.305192 |
| 65 | 4 | 6 | 0.281346 |
| 66 | 4 | 7 | 0.281731 |
| 67 | 4 | 8 | 0.280385 |
| 68 | 4 | 9 | 0.273654 |
| 69 | 4 | 10 | 0.257308 |
| 70 | 4 | 11 | 0.256346 |

| | k | s | Average test error |
|-----|---|----|--------------------|
| 71 | 4 | 12 | 0.258654 |
| 72 | 4 | 13 | 0.263077 |
| 73 | 4 | 14 | 0.255385 |
| 74 | 4 | 15 | 0.258654 |
| 75 | 4 | 16 | 0.247692 |
| 76 | 4 | 17 | 0.250385 |
| 77 | 4 | 18 | 0.245769 |
| 78 | 4 | 19 | 0.250962 |
| 79 | 4 | 20 | 0.247115 |
| 80 | 5 | 1 | 0.434615 |
| 81 | 5 | 2 | 0.365577 |
| 82 | 5 | 3 | 0.345577 |
| 83 | 5 | 4 | 0.302115 |
| 84 | 5 | 5 | 0.295192 |
| 85 | 5 | 6 | 0.286731 |
| 86 | 5 | 7 | 0.283077 |
| 87 | 5 | 8 | 0.271538 |
| 88 | 5 | 9 | 0.266923 |
| 89 | 5 | 10 | 0.273269 |
| 90 | 5 | 11 | 0.265000 |
| 91 | 5 | 12 | 0.258269 |
| 92 | 5 | 13 | 0.255385 |
| 93 | 5 | 14 | 0.261538 |
| 94 | 5 | 15 | 0.250769 |
| 95 | 5 | 16 | 0.254808 |
| 96 | 5 | 17 | 0.240962 |
| 97 | 5 | 18 | 0.246731 |
| 98 | 5 | 19 | 0.242308 |
| 99 | 5 | 20 | 0.242500 |
| 100 | 6 | 1 | 0.427692 |
| 101 | 6 | 2 | 0.372308 |
| 102 | 6 | 3 | 0.341923 |
| 103 | 6 | 4 | 0.317500 |
| 104 | 6 | 5 | 0.294231 |
| 105 | 6 | 6 | 0.282308 |
| 106 | 6 | 7 | 0.282115 |

| | k | s | Average test error |
|-----|---|----|--------------------|
| 107 | 6 | 8 | 0.267500 |
| 108 | 6 | 9 | 0.270385 |
| 109 | 6 | 10 | 0.265385 |
| 110 | 6 | 11 | 0.260000 |
| 111 | 6 | 12 | 0.255577 |
| 112 | 6 | 13 | 0.248462 |
| 113 | 6 | 14 | 0.248654 |
| 114 | 6 | 15 | 0.244231 |
| 115 | 6 | 16 | 0.252308 |
| 116 | 6 | 17 | 0.252115 |
| 117 | 6 | 18 | 0.239231 |
| 118 | 6 | 19 | 0.239038 |
| 119 | 6 | 20 | 0.242885 |
| 120 | 7 | 1 | 0.430577 |
| 121 | 7 | 2 | 0.352885 |
| 122 | 7 | 3 | 0.344808 |
| 123 | 7 | 4 | 0.313654 |
| 124 | 7 | 5 | 0.289423 |
| 125 | 7 | 6 | 0.282885 |
| 126 | 7 | 7 | 0.274038 |
| 127 | 7 | 8 | 0.265577 |
| 128 | 7 | 9 | 0.260385 |
| 129 | 7 | 10 | 0.254808 |
| 130 | 7 | 11 | 0.258654 |
| 131 | 7 | 12 | 0.259038 |
| 132 | 7 | 13 | 0.259038 |
| 133 | 7 | 14 | 0.242885 |
| 134 | 7 | 15 | 0.250192 |
| 135 | 7 | 16 | 0.242500 |
| 136 | 7 | 17 | 0.243846 |
| 137 | 7 | 18 | 0.241538 |
| 138 | 7 | 19 | 0.244615 |
| 139 | 7 | 20 | 0.245192 |
| 140 | 8 | 1 | 0.428269 |
| 141 | 8 | 2 | 0.384038 |
| 142 | 8 | 3 | 0.336731 |

| | k | s | Average test error |
|-----|---|----|--------------------|
| 143 | 8 | 4 | 0.330000 |
| 144 | 8 | 5 | 0.301346 |
| 145 | 8 | 6 | 0.289038 |
| 146 | 8 | 7 | 0.269615 |
| 147 | 8 | 8 | 0.280385 |
| 148 | 8 | 9 | 0.274231 |
| 149 | 8 | 10 | 0.256731 |
| 150 | 8 | 11 | 0.258269 |
| 151 | 8 | 12 | 0.253846 |
| 152 | 8 | 13 | 0.262500 |
| 153 | 8 | 14 | 0.251346 |
| 154 | 8 | 15 | 0.245769 |
| 155 | 8 | 16 | 0.246154 |
| 156 | 8 | 17 | 0.248462 |
| 157 | 8 | 18 | 0.245192 |
| 158 | 8 | 19 | 0.236538 |
| 159 | 8 | 20 | 0.241154 |
| 160 | 9 | 1 | 0.415577 |
| 161 | 9 | 2 | 0.373846 |
| 162 | 9 | 3 | 0.336154 |
| 163 | 9 | 4 | 0.318462 |
| 164 | 9 | 5 | 0.297500 |
| 165 | 9 | 6 | 0.289615 |
| 166 | 9 | 7 | 0.287692 |
| 167 | 9 | 8 | 0.268846 |
| 168 | 9 | 9 | 0.275192 |
| 169 | 9 | 10 | 0.276731 |
| 170 | 9 | 11 | 0.251731 |
| 171 | 9 | 12 | 0.247115 |
| 172 | 9 | 13 | 0.261731 |
| 173 | 9 | 14 | 0.251346 |
| 174 | 9 | 15 | 0.248846 |
| 175 | 9 | 16 | 0.250192 |
| 176 | 9 | 17 | 0.234038 |
| 177 | 9 | 18 | 0.245769 |
| 178 | 9 | 19 | 0.247308 |

| | k | s | Average test error |
|-----|----|----|--------------------|
| 179 | 9 | 20 | 0.252308 |
| 180 | 10 | 1 | 0.431731 |
| 181 | 10 | 2 | 0.366154 |
| 182 | 10 | 3 | 0.333269 |
| 183 | 10 | 4 | 0.310000 |
| 184 | 10 | 5 | 0.300192 |
| 185 | 10 | 6 | 0.286154 |
| 186 | 10 | 7 | 0.285769 |
| 187 | 10 | 8 | 0.266923 |
| 188 | 10 | 9 | 0.262115 |
| 189 | 10 | 10 | 0.256731 |
| 190 | 10 | 11 | 0.256731 |
| 191 | 10 | 12 | 0.264423 |
| 192 | 10 | 13 | 0.257885 |
| 193 | 10 | 14 | 0.242308 |
| 194 | 10 | 15 | 0.247692 |
| 195 | 10 | 16 | 0.242115 |
| 196 | 10 | 17 | 0.244038 |
| 197 | 10 | 18 | 0.235385 |
| 198 | 10 | 19 | 0.231154 |
| 199 | 10 | 20 | 0.231731 |

```
In [15]: res.loc[res['Average test error'] == res['Average test error'].min()]
```

Out[15]:

| | k | s | Average test error |
|-----|----|----|--------------------|
| 198 | 10 | 19 | 0.231154 |

Does this provide a reasonable alternative to SVM with slack formulation without feature selection on this data set? What are the pros and cons of this approach?

SVM with feature selection ($k = 7$, $s = 20$, $c = 1$) gives lowest average test error of 23% whereas SVM without feature selection ($c = 1$) gives lowest test error of 21%. There is not much improvement with feature selection.

Pros:

- Training individual models is computationally faster as number of features is less

Cons:

- Choosing k and s is difficult. Grid search like above takes long time