

COMS4721 Spring 2015: Homework 2

Daniel M. Sheehan - dms2203@columbia.edu Discussants: tudor, jade

Problem 1 - Skipped (we can pick 2 of the first 3 problems)

Problem 2

Problem 2 (Features; 10 points). It is common to pre-process the feature vectors in \mathbb{R}^d before passing them to a learning algorithm. Two simple ways to pre-process are as follows.

- **Centering:** Subtract the mean $\hat{\mu} := \frac{1}{|S|} \sum_{(x,y) \in S} x$ (of the training data) from every feature vector:

$$x \mapsto x - \hat{\mu}.$$

- **Standardization:** Perform centering, and then divide every feature by the per-feature standard deviation $\hat{\sigma}_i = \sqrt{\frac{1}{|S|} \sum_{(x,y) \in S} (x_i - \hat{\mu}_i)^2}$:

$$(x_1, x_2, \dots, x_d) \mapsto \left(\frac{x_1 - \hat{\mu}_1}{\hat{\sigma}_1}, \frac{x_2 - \hat{\mu}_2}{\hat{\sigma}_2}, \dots, \frac{x_d - \hat{\mu}_d}{\hat{\sigma}_d} \right).$$

(The same transformations should be applied to *all* feature vectors you encounter, including any future test points.)

For each of the following learning algorithms, and each of the above pre-processing transformations, explain whether or not each of the transformation can affect the resulting learned classifier.

- The classifier based on the generative model where class conditional distributions are multi-variate Gaussian distributions with a fixed covariance equal to the identity matrix I . Assume MLE is used for parameter estimation.
- The 1-NN classifier using Euclidean distance.
- The greedy decision tree learning algorithm with axis-aligned splits. (For concreteness, assume Gini index is used as the uncertainty measure, and the algorithm stops after 20 leaf nodes.)
- Empirical Risk Minimization: the (intractable) algorithm that finds the linear classifier (both the weight vector and threshold) that has the smallest training error rate.

You should assume the following: (i) the per-feature standard deviations are never zero; and (ii) there are never any “ties” whenever you do an arg max or an arg min. Also ignore computational and numerical precision issues.

- a. The classifier based on the generative model where class conditional distributions are **multi-variate Gaussian distributions with a mixed covariance** equal to the identity matrix I . Assume MLE is used for parameter estimation.
 - **Centering:** Does not have an effect on the classification rate because we're just changing the origin. Its shifted in space but relative positions are the same. We're affecting the model, but not the classifier. But the position of the sep. hyperplane has moved so the model is different.
 - **Standardization:** There is no effect because since they have covariance then the classifier is a straight line. We affect the model, distance but not the classification rate.

- b. The **1-NN classifier using Euclidean distance**.
 - **Centering**: Centering will not effect 1-NN because the distance between the points will not be altered.
 - **Standardization**: Standarization will effect the 1-NN because the distance between the features will have changed.
- c. The **greedy decision tree learning algorithm** with axis-aligned splits.
 - **Centering**: Centering will affect the classifier because it cannot return the global optimal decision tree.
 - **Standardization**: If centering affects it in this instance, standardization also does.
- d. **Empirical Risk Minimization**
 - **Centering**: Not certain this can be achieved b/c of complexity.
 - **Standardization**: ibid

Problem 3

Problem 3 (Covariance matrices; 10 points). Let \mathbf{X} be a mean-zero random vector in \mathbb{R}^d (so $\mathbb{E}(\mathbf{X}) = \mathbf{0}$). Let $\Sigma := \mathbb{E}(\mathbf{X}\mathbf{X}^\top)$ be the covariance matrix of \mathbf{X} , and suppose its eigenvalues are $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$. Let $\sigma^2 > 0$ be a positive number.

- What are the eigenvalues of $\Sigma + \sigma^2 \mathbf{I}$?
- What are the eigenvalues of $(\Sigma + \sigma^2 \mathbf{I})^{-1}$?
- Suppose $\mathbf{R} \in \mathbb{R}^{d \times d}$ is an invertible symmetric matrix that satisfies $\mathbf{R}\mathbf{R} = \Sigma + \sigma^2 \mathbf{I}$. Let $\mathbf{v}_1 \in \mathbb{R}^d$ be an eigenvector of Σ corresponding to its largest eigenvalue λ_1 . What is the variance of the random variable $\mathbf{v}_1^\top \mathbf{R}^{-1} \mathbf{X}$? Explain your answer succinctly (but precisely).

Hint: The answers can be given in terms of σ^2 and the eigenvalues of Σ .

Problem 3

$$\sum \lambda v = \lambda v$$

a) eigenvalues of $\Sigma + \sigma^2 \mathbf{I}$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{bmatrix}$$

$$\sigma^2 I + \Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{bmatrix} + \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}$$

$\Sigma \leftarrow \text{original covariance}$

$$\sigma^2 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}$$

$$\Sigma + \sigma^2 I = \begin{bmatrix} \sigma_1^2 + \sigma^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 + \sigma^2 \end{bmatrix}$$

$$(\Sigma + \sigma^2 I)v' = \lambda' v'$$

$$\Sigma v' + \sigma^2 I v' = \lambda' v'$$

$$\Sigma v' = \lambda' v' - \sigma^2 I v'$$

$$\Sigma v' = (\lambda' - \sigma^2)v'$$

v' eigenvector of Σ

$\lambda' - \sigma^2$ eigenvalue of Σ

$$\lambda' - \sigma^2 = \lambda$$

$$\lambda' = \lambda + \sigma^2$$

b) eigenvalues of $(\Sigma + \sigma^2 I)^{-1}$

$$(\Sigma + \sigma^2 I)^{-1} v' = \lambda' v'$$

$$v' = \lambda' (\Sigma + \sigma^2 I) v'$$

$$\frac{1}{\lambda'} v' = (\Sigma + \sigma^2 I) v'$$

$$= \Sigma v' + \sigma^2 v'$$

$$\sum v' = \frac{1}{\lambda'} v' - \sigma^2 v'$$

$$\sum v' = (\frac{1}{\lambda'} - \sigma^2) v', \leftarrow \text{this here is } \lambda$$

$$\lambda = \frac{1}{\lambda'} - \sigma^2$$

$$\lambda + \sigma^2 = \frac{1}{\lambda'}$$

////////////////////////////////////

$$\lambda' = \frac{1}{\lambda + \sigma^2}$$

////////////////////////////////////

c)

$$RR = \sum + \sigma^2 I$$

$$R^{-1}R = I, \quad RR^{-1} = I$$

$$\sum v_1 = \lambda_1 v_1$$

$$E(X) = 0$$

$$Y = v_1^T R^{-1} X$$

$$\mu_Y = E(Y) = E(v_1^T, R^{-1} X)$$

$$= v_1^T R^{-1} E(X), \leftarrow \text{which is } 0$$

$$Var(Y) = E[Y^2]$$

$$= E[(v_1^T R^{-1} X)(v_1^T R^{-1} X)^T]$$

$$= E[v_1^T R^{-1} X X^T R^{-T} v_1]$$

$$= v_1^T R^{-1} E(X X^T) R^{-T} v_1$$

$$= v_1^T R^{-1} \sum R^{-T} v_1$$

$$RR = \sum + \sigma^2 I, \quad R = R^T \rightarrow R^{-1} = (R^{-1})^T$$

$$R^{-1}RR = R^{-1} \sum + \sigma^2 R^{-1} I$$

^ above is IR , which is R

$$R = R^{-1} \sum + \sigma^2 R^{-1}, \quad x(R^{-1})^T = R^{-1}$$

$$RR^{-1} = R^{-1} \sum R^{-T} + \sigma^2 R^{-1} R^{-1}$$

$$I = R^{-1} \sum R^{-T} + \sigma^2 R^{-1} R^{-1}$$

$$R^{-1} \sum R^{-T} = I - \sigma^2 R^{-1} R^{-1}$$

$$\text{Var}(Y) = v_1^T (I - \sigma^2 R^{-1} R^{-1}) v_1$$

$$= v_1^T v_1 - \sigma^2 v_1^T R^{-1} R^{-1} v_1$$

$$= 1 - \sigma^2 v_1^T R^{-1} R^{-1} v_1$$

$$\sum v_1 = \lambda_1 v_1, \quad \leftarrow \sum = RR - \sigma^2 I$$

$$RR = \sum + \sigma^2 I$$

$$RRv_1 = (\lambda_1 + \sigma^2)v_1 \quad \text{multiply both sides by } R^{-1}$$

$$R^{-1}RRv_1 = R^{-1}(\lambda_1 + \sigma^2)v_1$$

^ is I

$$Rv_1 = (\lambda_1 + \sigma^2)R^{-1}v_1$$

$$R^{-1}v_1 = \frac{Rv_1}{\lambda_1 + \sigma^2}$$

transpose

$$(R^{-1}v_1)^T = V_1^T R^{-T} = v_1^T R^{-1} = \frac{v_1^T R^T}{\lambda_1 + \sigma^2}$$

$$\rightarrow v_1^T R^{-1} = \frac{v_1^T R}{\lambda_1 + \sigma^2}$$

$$\text{var}(Y) = 1 - \sigma^2 v_1^T R^{-1} R^{-1} v_1$$

$$= 1 - \sigma^2 \left(\frac{v_1^T R}{\lambda_1 + \sigma^2} \right) \left(\frac{R v_1}{\lambda_1 + \sigma^2} \right)$$

$$= 1 - \frac{\sigma^2}{(\lambda_1 + \sigma^2)^2} v_1^T (RR) v_1$$

$$\text{note : } RR = \sum + \sigma^2 I$$

$$= 1 - \frac{\sigma^2}{(\lambda_1 + \sigma^2)^2} [v_1^T \sum v_1 + \sigma^2 v_1^T v_1]$$

$$v_1^T v_1 \text{ is } 1$$

$$= 1 - \frac{\sigma^2}{(\lambda_1 + \sigma^2)^2} (\lambda_1 + \sigma^2)$$

$$= 1 - \frac{\sigma^2}{\lambda_1 + \sigma^2}$$

Problem 4

Problem 4 (Linear/quadratic classifiers; 40 points). Download the spam data set `spam_fixed.mat` from Courseworks. This is the data set described in the *ESL* text for a binary classification problem of predicting whether an e-mail is spam or not. The training data and test data are in the usual format. You can read about the original features in *ESL* (Chapter 9.1, pages 300–301); in this version, the features are (approximately) standardized.

Write a MATLAB script or Python script that, using only the training data, tries out six different methods for learning linear/quadratic classifiers, and ultimately selects (via ten-fold cross validation) and uses one of these methods. (Think of the learning method itself as a “hyperparameter”!) The six methods to try are the following.

1. Averaged-Perceptron with 64 passes through the data.

Averaged-Perceptron is like Voted-Perceptron (which uses Online Perceptron), except instead of forming the final classifier as a majority vote over the various linear classifiers, you simply form a single linear classifier by averaging the weight vectors and thresholds of the various linear classifiers. Important: Before each pass of Online Perceptron, you should randomly shuffle the order of the training examples.

2. Logistic regression model with MLE for parameter estimation.

3. Generative model classifier where class conditional distributions are multivariate Gaussian distributions *with shared covariance matrix for all classes*. Use MLE for parameter estimation.

4. Same as above, except arbitrary Gaussians (i.e., each class with its own covariance matrix).

- 5&6. Averaged-Perceptron and logistic regression as above, with feature map $\phi: \mathbb{R}^{57} \rightarrow \mathbb{R}^{1710}$ given by $\phi(\mathbf{x}) := (x_1, x_2, \dots, x_{57}, x_1^2, x_2^2, \dots, x_{57}^2, x_1x_2, x_1x_3, \dots, x_1x_{57}, \dots, x_{56}x_{57})$. No need to use the kernel trick here.

Note that we want to learn linear classifiers that are possibly non-homogeneous: so you should learn a weight vector \mathbf{w} in \mathbb{R}^{57} (or \mathbb{R}^{1710}) *and* a threshold value $t \in \mathbb{R}$.

You should write your own code for implementing Online Perceptron, Averaged-Perceptron, and cross-validation (put these code in separate files). The code (and the driver script that runs everything) should be easy-to-read, commented as necessary. For the other learning methods, you can use existing library implementations, but you are responsible for the correctness of these implementations (as per the specifications from the course lectures).

Report the cross-validation error rates for all methods, the training error rate of the classifier learned by the selected method (and state which method was chosen), and the test error rate for the learned classifier. Submit all codes & scripts you write yourself, and give references to any existing software you use for the other learning algorithms.

```
1 | from scipy.io import loadmat
2 | import numpy as np
3 | from sklearn.linear_model import LogisticRegression
4 | from sklearn.ensemble import RandomForestClassifier
5 |
6 | data = loadmat('spam_fixed.mat')
7 |
```



```

8 Ytrain = data['labels'].flatten()
9 Xtrain = data['data']
10 Ytest = data['testlabels'].flatten()
11 Xtest = data['testdata']
12
13 class Perceptron(object):
14     def fit(self, X, Y):
15         N, D = X.shape
16         V = [np.zeros(D)]
17         C = [0]
18         for t in xrange(64):
19             for i in xrange(N):
20                 v = V[-1]
21                 y_hat_i = np.sign(v.dot(X[i]))
22                 if y_hat_i == Y[i]:
23                     C[-1] = C[-1] + 1
24                 else:
25                     new_v = v + Y[i]*X[i]
26                     V.append(new_v)
27                     C.append(1)
28
29             self.v = np.zeros(D)
30             total_votes = 0
31             for c, v in zip(C[1:], V[1:]):
32                 self.v += c*v
33                 total_votes += c
34             self.v /= total_votes
35
36     def score(self, X, Y):
37         P = np.sign(X.dot(self.v))
38         return np.mean(P == Y)
39
40 class BigPerceptron(object):
41     def fit(self, X, Y):
42         X2 = transform(X)
43         self.model = Perceptron()
44         self.model.fit(X2, Y)
45
46     def score(self, X, Y):
47         X2 = transform(X)
48         return self.model.score(X2, Y)
49
50 class Gauss1(object):
51     # Ax = b
52     # x = inv_A * b

```

```

53     # solve(A, b)
54     def fit(self, X, Y):
55         # w = inv_cov * (mu1 - mu0)
56         # b = 0.5*(mu1 + mu0) * inv_cov * (mu0 - mu1) + log(pi1/pi0)
57         cov = np.cov(X.T)
58         idx1 = np.where(Y == 1)[0]
59         idx0 = np.where(Y == -1)[0]
60         mu0 = X[idx0, :].mean(axis=0)
61         mu1 = X[idx1, :].mean(axis=0)
62         self.w = np.linalg.solve(cov, mu1 - mu0)
63         N = len(Y)
64         pi1 = float(len(idx1)) / N
65         pi0 = float(len(idx0)) / N
66         self.b = -0.5*(mu0 + mu1).dot(self.w) + np.log(pi1/pi0)
67
68     def score(self, X, Y):
69         P = np.sign(X.dot(self.w) + self.b)
70         return np.mean(P == Y)
71
72 class Gauss2(object):
73     def fit(self, X, Y):
74         idx1 = np.where(Y == 1)[0]
75         idx0 = np.where(Y == -1)[0]
76
77         x0 = X[idx0, :]
78         x1 = X[idx1, :]
79         mu0 = x0.mean(axis=0)
80         mu1 = x1.mean(axis=0)
81         cov0 = np.cov(x0.T)
82         cov1 = np.cov(x1.T)
83
84         self.A = 0.5*(np.linalg.inv(cov0) - np.linalg.inv(cov1))
85         icov0mu0 = np.linalg.solve(cov0, mu0)
86         icov1mu1 = np.linalg.solve(cov1, mu1)
87         self.w = icov1mu1 - icov0mu0
88         N = len(Y)
89         pi1 = float(len(idx1)) / N
90         pi0 = float(len(idx0)) / N
91         self.b = 0.5*(mu0.dot(icov0mu0) - mu1.dot(icov1mu1)) + np.log(pi1/pi0)
92
93     def score(self, X, Y):
94         P = np.sign((X.dot(self.A)*X).sum(axis=1) + X.dot(self.w) + self.b)
95         return np.mean(P == Y)
96
97     def transform(X):

```

```

98     N, D = X.shape
99     X2 = np.zeros((N, 1710))
100    X2[:, :D] = X
101    j = D
102    for i in xrange(D):
103        for k in xrange(D):
104            if i <= k:
105                X2[:, j] = X[:, i]*X[:, k]
106                j += 1
107    # mu = X2.mean(axis=0)
108    # std = X2.std(axis=0)
109    # X2 = (X2 - mu) / std
110    return X2
111
112    class BigLogistic(object):
113        def fit(self, X, Y):
114            X2 = transform(X)
115            self.model = LogisticRegression()
116            self.model.fit(X2, Y)
117
118        def score(self, X, Y):
119            X2 = transform(X)
120            return self.model.score(X2, Y)
121
122    def cross_validation(model, X, Y):
123        # split the data into 10 parts
124        N = len(Y)
125        batchsize = N / 10 + 1
126        scores = []
127        for i in xrange(10):
128            # test on i-th part, train on other 9 parts
129            # (i + 1)*batchsize
130            start = i*batchsize
131            end = (i*batchsize + batchsize)
132            Xvalid = X[start:end]
133            Yvalid = Y[start:end]
134
135            Xtrain = np.concatenate([ X[:start] , X[end:] ])
136            Ytrain = np.concatenate([ Y[:start] , Y[end:] ])
137
138            model.fit(Xtrain, Ytrain)
139            scores.append(model.score(Xvalid, Yvalid))
140    return np.mean(scores)
141
142

```

```
143 models = {
144     '2. Logistic': LogisticRegression(),
145     'TEST for. RandomForest': RandomForestClassifier(),
146     '1. Perceptron': Perceptron(),
147     '6. Big Logistic': BigLogistic(),
148     '5. Big Perceptron': BigPerceptron(),
149     '3. Gauss 1': Gauss1(),
150     '4. Gauss 2': Gauss2(),
151 }
152
153 def main():
154     for name, model in models.iteritems():
155         print "Model:", name, "accuracy:", cross_validation(model, Xtrain, Ytrain)
156
157 if __name__ == '__main__':
158     main()
159
160
161 Model: 5. Big Perceptron accuracy: 0.912894492741
162 Model: 2. Logistic accuracy: 0.918437344953
163 Model: TEST for. RandomForest accuracy: 0.938654356408
164 Model: 3. Gauss 1 accuracy: 0.868188191643
165 Model: 4. Gauss 2 accuracy: 0.806716353517
166 Model: 1. Perceptron accuracy: 0.918415773238
167 Model: 6. Big Logistic accuracy: 0.923985590094
```