

# COMS 4771 Machine Learning (2018 Fall)

## Homework 0

Jing Qian - jq2282@columbia.edu

September 9, 2018

### Problem 1.1

(i)

$$\begin{aligned}\Pr[X = 1] &= \sum_y \Pr[X = 1, Y = y] = 0.2 + 0.2 + 0.3 = 0.7 \\ \Pr[X = 2] &= \sum_y \Pr[X = 2, Y = y] = 0.1 + 0.1 + 0.1 = 0.3\end{aligned}\tag{1}$$

(ii)

$$\Pr[Y = 1|X = 2] = \frac{\Pr[Y = 1, X = 2]}{\Pr[X = 2]} = \frac{1}{3}\tag{2}$$

(iii)

$$\begin{aligned}\Pr[X = 1|Y = 3] &= \frac{0.3}{0.3 + 0.1} = 0.75 \\ \Pr[X = 2|Y = 3] &= \frac{0.1}{0.3 + 0.1} = 0.25\end{aligned}\tag{3}$$

$$\mathbb{E}[f(X)|Y = 3] = 1^2 * \Pr[X = 1|Y = 3] + 2^2 * \Pr[X = 2|Y = 3] = 1.75\tag{4}$$

### Problem 1.2

(i) To be a probability distribution,  $g_\theta$  has to satisfy following two conditions: 1)  $g_\theta(x) \geq 0$  for any  $x \in [0, \infty)$ ; 2)  $\int_0^\infty g_\theta(x)dx = 1$ .

Since  $\theta > 0$  and  $e^{-\frac{x}{\theta}} > 0$ ,  $g_\theta > 0$ . First condition is satisfied.

$$\int_0^\infty g_\theta(x)dx = \int_0^\infty \frac{1}{\theta} e^{-\frac{x}{\theta}} dx = \int_0^\infty e^{-y} dy = 1\tag{5}$$

where  $y = \frac{x}{\theta}$ . Second condition is satisfied too. So  $g_\theta$  is a probability distribution.

(ii)

$$\mathbb{E}[g_\theta] = \int_0^\infty x g_\theta(x) dx = \int_0^\infty \frac{x}{\theta} e^{-\frac{x}{\theta}} dx = \theta \int_0^\infty y e^{-y} dy = \theta \quad (6)$$

where  $y = \frac{x}{\theta}$ .

(iii)

$$\begin{aligned} \text{var}(g_\theta) &= \int_0^\infty (x - \mathbb{E}[g_\theta])^2 g_\theta(x) dx \\ &= \int_0^\infty \frac{(x - \theta)^2}{\theta} e^{-\frac{x}{\theta}} dx \\ &= \theta^2 \int_0^\infty y^2 e^{-y} dy - 2\theta^2 \int_0^\infty y e^{-y} dy + \theta^2 \int_0^\infty e^{-y} dy \\ &= \theta^2 \end{aligned} \quad (7)$$

### Problem 1.3

Considering that  $X$  and  $Y$  are jointly distributed Gaussian random variables,  $X + Y$  is still Gaussian distributed.

$$\begin{aligned} \mathbb{E}[X + Y] &= \sum_x \sum_y (x + y) \Pr(X = x, Y = y) \\ &= \sum_x \sum_y x \Pr(X = x, Y = y) + \sum_x \sum_y y \Pr(X = x, Y = y) \\ &= \sum_x x \Pr(X = x) + \sum_y y \Pr(Y = y) \\ &= \mathbb{E}[X] + \mathbb{E}[Y] \\ &= 0 \end{aligned} \quad (8)$$

$$\begin{aligned} \text{var}(X + Y) &= \mathbb{E}[(X + Y)^2] - \mathbb{E}[X + Y]^2 \\ &= \mathbb{E}[X^2] + \mathbb{E}[Y^2] + 2\mathbb{E}[XY] - (\mathbb{E}[X] + \mathbb{E}[Y])^2 \\ &= (\mathbb{E}[X^2] - \mathbb{E}[x]^2) + (\mathbb{E}[Y^2] - \mathbb{E}[Y]^2) + 2(\mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]) \\ &= \text{var}(X) + \text{var}(Y) + 2\text{cov}(X, Y) \\ &= 6 \end{aligned} \quad (9)$$

So  $X + Y$  has Gaussian distribution with the mean 0 and variance 6, in other words,  $N(0, 6)$ .

### Problem 1.4

For a fair coin, which means the possibility of tossing a head is  $\frac{1}{2}$  every time, the expected absolute difference between the number of heads  $H$  and that of tails  $T$  is:

$$\mathbb{E}[|H - T|] = \sum_{i=0}^n \binom{n}{i} \left(\frac{1}{2}\right)^i \left(\frac{1}{2}\right)^{n-i} |i - (n - i)| = \frac{1}{2^n} \sum_{i=0}^n \binom{n}{i} |2i - n| \quad (10)$$

where  $i$  is the number of heads. Since  $\binom{n}{i} |2i - n| = \binom{n}{j} |2j - n|$  while  $j = n - i$ , we have:

$$\sum_{i=0}^n \binom{n}{i} |2i - n| = \begin{cases} 2 \sum_{i=0}^{\frac{n}{2}-1} \binom{n}{i} (n - 2i) + \binom{n}{n/2} (n - n) = 2 \sum_{i=0}^{\frac{n}{2}-1} \binom{n}{i} (n - 2i), & n \text{ is even.} \\ 2 \sum_{i=0}^{\frac{n-1}{2}} \binom{n}{i} (n - 2i), & n \text{ is odd.} \end{cases} \quad (11)$$

Since the odd case and even case have slight differences, we treat them separately.

$$\begin{aligned} \mathbb{E}_{\text{even}}[|H - T|] &= \frac{1}{2^{(n-1)}} \sum_{i=0}^{\frac{n}{2}-1} \binom{n}{i} (n - 2i) = \frac{1}{2^{(n-1)}} \sum_{i=0}^{\frac{n}{2}-1} \frac{n!}{i! (n-i)!} [(n-i) - i] \\ &= \frac{1}{2^{(n-1)}} \left[ \sum_{i=0}^{\frac{n}{2}-1} \frac{n!}{i! (n-i)!} (n-i) - \sum_{i=0}^{\frac{n}{2}-1} \frac{n!}{i! (n-i)!} i \right] \\ &= \frac{1}{2^{(n-1)}} \left[ \sum_{i=0}^{\frac{n}{2}-1} \frac{n!}{i! (n-i-1)!} - \sum_{j=0}^{\frac{n}{2}-2} \frac{n!}{j! (n-j-1)!} \right] \end{aligned} \quad (12)$$

where  $j = i - 1$ . After subtraction of the two summations in Eq. (12), only the term  $i = \frac{n}{2} - 1$  remains.

$$\mathbb{E}_{\text{even}}[|H - T|] = \frac{1}{2^{(n-1)}} \frac{n!}{(\frac{n}{2}-1)! (\frac{n}{2})!} = \frac{n}{2^n} \frac{n!}{((\frac{n}{2})!)^2} \quad (13)$$

Similarly, we could get the result when  $n$  is an odd number:

$$\mathbb{E}[|H - T|] = \begin{cases} \frac{n}{2^n} \frac{n!}{((\frac{n}{2})!)^2}, & n \text{ is even.} \\ \frac{1}{2^{(n-1)}} \frac{n!}{((\frac{n-1}{2})!)^2}, & n \text{ is odd.} \end{cases} \quad (14)$$

When  $n \rightarrow \infty$ , we could get an approximate value for Eq. (14) using Stirling's approximation for factorials.

$$\mathbb{E}_{\text{even}}[|H - T|] = \frac{n}{2^n} \frac{n!}{((\frac{n}{2})!)^2} \approx \frac{n}{2^n} \frac{\sqrt{2\pi n} (\frac{n}{e})^n}{\pi n (\frac{n}{2e})^n} = \sqrt{\frac{2n}{\pi}} \quad (15)$$

And

$$\begin{aligned} \mathbb{E}_{\text{odd}}[|H - T|] &= \frac{1}{2^{(n-1)}} \frac{n!}{((\frac{n-1}{2})!)^2} \approx \frac{1}{2^{(n-1)}} \frac{\sqrt{2\pi n} (\frac{n}{e})^n}{\pi (n-1) (\frac{n-1}{2e})^{(n-1)}} \\ &= \sqrt{\frac{2n}{\pi}} \frac{1}{e} \frac{1}{(1 - \frac{1}{n})^n} \\ &\approx \sqrt{\frac{2n}{\pi}} \end{aligned} \quad (16)$$

So when independently toss a fair coin  $n$  times, the expected absolute difference between the number of heads  $H$  and the number of tails  $T$  is about  $\sqrt{\frac{2n}{\pi}}$ .

## Problem 2.1

(i) Let  $\mathbf{A}$  be a matrix contains vectors in the question and be transformed into row echelon form:

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 8 & 10 \\ 3 & 3 & 6 \\ 4 & 2 & 6 \end{bmatrix} \iff \mathbf{V} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad (17)$$

Since  $\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{V}) = 2$ , the dimension of the subspace  $S$  is 2.

(ii) To get the orthogonal linear projection of the point  $\begin{pmatrix} 6 \\ 5 \\ 9 \\ 2 \end{pmatrix}$  onto the subspace  $S$ , we need to calculate the orthonormal basis of subspace  $S$  first. From (i), we could get the orthonormal basis  $(\mathbf{v}_1, \mathbf{v}_2)$  of subspace  $S$ :

$$(\mathbf{v}_1, \mathbf{v}_2) = \left( \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \right) \quad (18)$$

Gram-Schmidt orthogonalization also works but the given orthonormal basis are much more complex.

So the projection  $P$  of given point  $\mathbf{u} = \begin{pmatrix} 6 \\ 5 \\ 9 \\ 2 \end{pmatrix}$  onto  $S$  is:

$$\mathbf{P} = \sum_i \frac{\langle \mathbf{u}, \mathbf{v}_i \rangle}{\langle \mathbf{v}_i, \mathbf{v}_i \rangle} \mathbf{v}_i = \langle \mathbf{u}, \mathbf{v}_1 \rangle \mathbf{v}_1 + \langle \mathbf{u}, \mathbf{v}_2 \rangle \mathbf{v}_2 = \begin{pmatrix} 6 \\ 5 \\ 0 \\ 0 \end{pmatrix}. \quad (19)$$

## Problem 2.2

A square matrix is invertible if and only if none of its eigenvalues is zero. So to prove  $\mathbf{A}^T \mathbf{A} + \rho \mathbf{I}$  is invertible, we could try to prove that all its eigenvalues are non-zero. Let  $\lambda$  be one eigenvalue of matrix  $\mathbf{A}^T \mathbf{A} + \rho \mathbf{I}$ . Then we have:

$$\det(\mathbf{A}^T \mathbf{A} + \rho \mathbf{I} - \lambda \mathbf{I}) = 0 \quad (20)$$

So,

$$\lambda = \det(\mathbf{A}^T \mathbf{A}) + \rho \quad (21)$$

Let  $\tilde{\lambda}$  and  $\mathbf{x}$  be one eigenvalue and corresponding eigenvector of matrix  $\mathbf{A}^T \mathbf{A}$ . In other words,  $\mathbf{A}^T \mathbf{A} \mathbf{x} = \tilde{\lambda} \mathbf{x}$ . Then,

$$\|\mathbf{A} \mathbf{x}\|^2 = (\mathbf{A} \mathbf{x})^T \mathbf{A} \mathbf{x} = \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} = \tilde{\lambda} \mathbf{x}^T \mathbf{x} = \tilde{\lambda} \quad (22)$$

Since  $\det(\mathbf{A}^T \mathbf{A}) = \tilde{\lambda} = \|\mathbf{A} \mathbf{x}\|^2 \geq 0$  and  $\rho > 0$ ,  $\lambda = \det(\mathbf{A}^T \mathbf{A}) + \rho > 0$ .

Therefore, the matrix  $\mathbf{A}^T \mathbf{A} + \rho \mathbf{I}$  is invertible because its eigenvalues are all positive.

## Problem 3.1

(i) The function could be expanded as following:

$$\begin{aligned} f(\mathbf{x}) &= \|\mathbf{Ax} - \mathbf{b}\|^2 + \|\mathbf{x}\|^2 = (\mathbf{Ax} - \mathbf{b})^T(\mathbf{Ax} - \mathbf{b}) + \mathbf{x}^T\mathbf{x} \\ &= \mathbf{x}^T\mathbf{A}^T\mathbf{Ax} + \mathbf{x}^T\mathbf{x} - \mathbf{x}^T\mathbf{A}^T\mathbf{b} - \mathbf{b}^T\mathbf{Ax} + \mathbf{b}^T\mathbf{b} \end{aligned} \quad (23)$$

Using the derivatives of matrices:

$$\begin{aligned} \frac{\partial \mathbf{x}^T\mathbf{B}}{\partial \mathbf{x}} &= \frac{\partial \mathbf{B}^T\mathbf{x}}{\partial \mathbf{x}} = \mathbf{B} \\ \frac{\partial \mathbf{x}^T\mathbf{B}\mathbf{x}}{\partial \mathbf{x}} &= (\mathbf{B} + \mathbf{B}^T)\mathbf{x} \\ \frac{\partial \mathbf{x}^T\mathbf{B}\mathbf{x}}{\partial \mathbf{x}} &= 2\mathbf{x} \end{aligned} \quad (24)$$

We could get:

$$\begin{aligned} \nabla f(\mathbf{x}) &= \frac{\partial (\mathbf{x}^T\mathbf{A}^T\mathbf{Ax} + \mathbf{x}^T\mathbf{x} - \mathbf{x}^T\mathbf{A}^T\mathbf{b} - \mathbf{b}^T\mathbf{Ax} + \mathbf{b}^T\mathbf{b})}{\partial \mathbf{x}} \\ &= 2\mathbf{A}^T\mathbf{Ax} + 2\mathbf{x} - 2\mathbf{A}^T\mathbf{b} \end{aligned} \quad (25)$$

(ii) To find  $\mathbf{x}$  minimizes  $f$ , first, we need to find the  $\mathbf{x}$  where  $\nabla f(\mathbf{x})$  equals to  $\mathbf{0}$ . In other words,  $\nabla f(\mathbf{x}) = 2\mathbf{A}^T\mathbf{Ax} + 2\mathbf{x} - 2\mathbf{A}^T\mathbf{b} = \mathbf{0}$ . So  $(\mathbf{A}^T\mathbf{A} + \mathbf{I})\mathbf{x} - \mathbf{A}^T\mathbf{b} = \mathbf{0}$ . Then  $\mathbf{x} = (\mathbf{A}^T\mathbf{A} + \mathbf{I})^{-1}\mathbf{A}^T\mathbf{b}$ .

Because when  $\|\mathbf{x}\| \rightarrow \infty$ ,  $f(\mathbf{x}) \rightarrow \infty$ ,  $\min(f(\mathbf{x}))$  appears when  $\nabla f(\mathbf{x}) = \mathbf{0}$ . Hence  $\mathbf{x} = (\mathbf{A}^T\mathbf{A} + \mathbf{I})^{-1}\mathbf{A}^T\mathbf{b}$  minimizes  $f$ .

## Problem 3.2

Suppose  $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ , then:

$$\begin{aligned} g(x) &= x^T Ax - b^T x + c \\ &= [x_1 \ x_2] \begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - [1 \ 2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + 3 \\ &= x_1^2 + 6x_1x_2 + x_2^2 - x_1 - 2x_2 + 3 \end{aligned} \quad (26)$$

The first partial derivatives are

$$\begin{aligned} P_{x_1}(x) &= P_{x_1}(x_1, x_2) = 2x_1 + 6x_2 - 1 \\ P_{x_2}(x) &= P_{x_2}(x_1, x_2) = 6x_1 + 2x_2 - 2 \end{aligned} \quad (27)$$

So  $\nabla g((1, 1)) = (2 + 6 - 1, 6 + 2 - 2) = (7, 6)$ .

## Problem 4.1

The output is:

(ii) 2

(iii) [85 87 19 21 3 60 62 69] [15 16 22 3 9 90 91 97 78 84]

(iv) 50.5

(vi) [1022.0677061701972, 263.62801935563306, 102.83034619975379]

\* If "the top three eigenvalues" refers to the eigenvalues with largest absolute value, then the output would be: [1022.0677061701972, 263.62801935563306, -235.74325551868949].

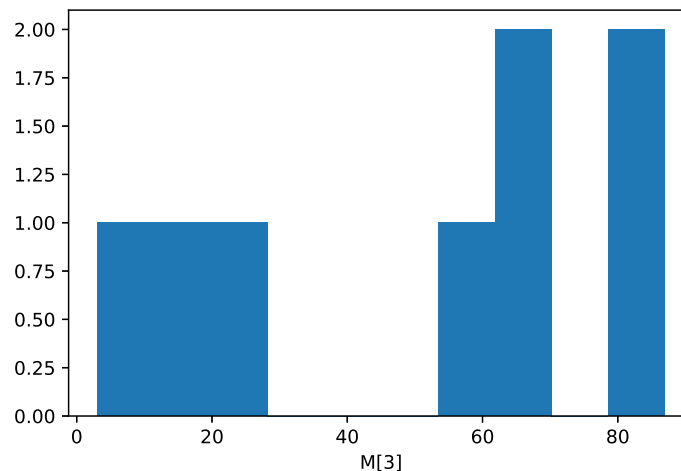


Figure 1: Histogram of the 4th row of  $\mathbf{M}$ .

HW0\_4.1.py

```

1 import numpy as np
2 import scipy.io as sio
3 import matplotlib.pyplot as plt
4 #load data
5 m = sio.loadmat('hw0data.mat')['M']
6 #print the dimensions of M.
7 print(np.ndim(m))
8 #print the 4th row and 5th column entry of M
9 print(m[3], m[:, 4])
10 #print the mean value of the 5th column of M
11 print(np.mean(m[:, 4]))
12 #compute the histogram of the 4th row of M and show the figure
13 plt.hist(m[3])
14 plt.show()

```

```

15 #compute and print the top three eigenvalues of the matrix MTM
16 evals, evecs = np.linalg.eig(np.dot(m.transpose(), m))
17 print(sorted(evals, reverse = True)[:3])

```

## Problem 4.2

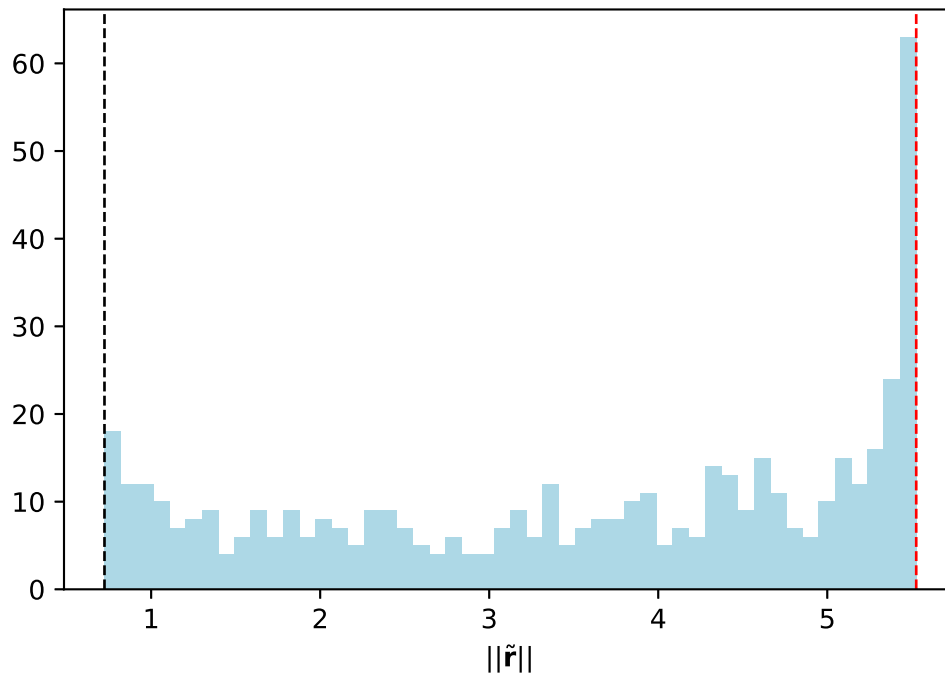


Figure 2: Histogram of values of  $\|\tilde{\mathbf{r}}\|$ . Black dashed line is  $\lambda_{\min}$  and red dashed line is  $\lambda_{\max}$ .

(vii) From Fig. (2), we could see that  $\|\tilde{\mathbf{r}}\|$  lie between  $\lambda_{\min}$  and  $\lambda_{\max}$ , which means  $\lambda_{\min}$  is the lower bound of  $\|\tilde{\mathbf{r}}\|$  while  $\lambda_{\max}$  is the upper bound.

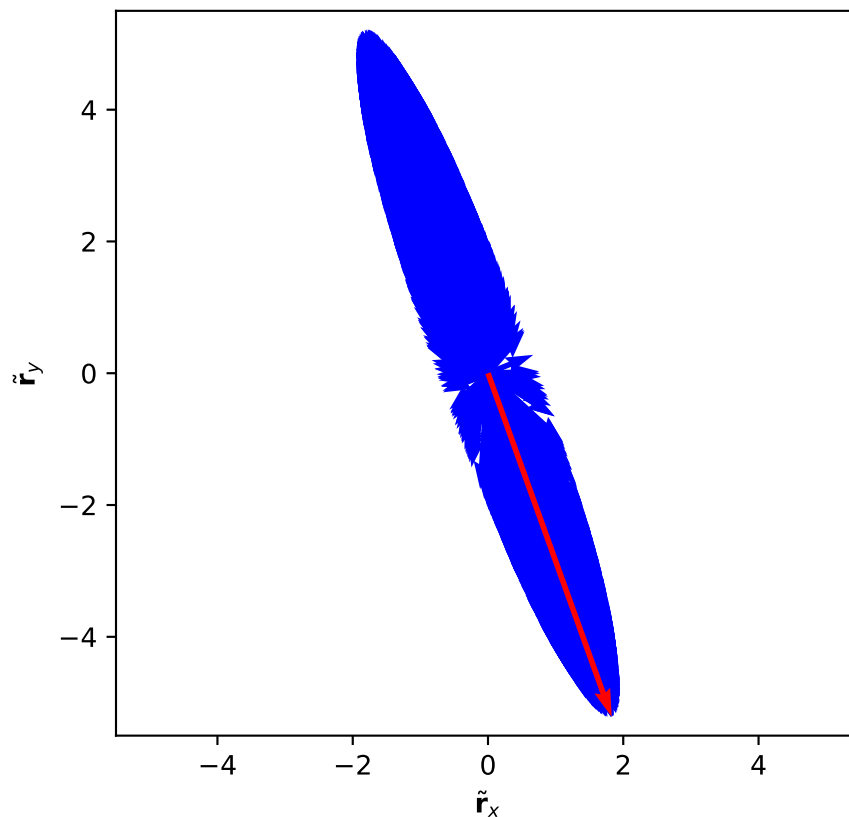


Figure 3: 2-D plot of all the distorted vectors  $\tilde{\mathbf{r}}$  (in blue color) and  $\mathbf{L}\mathbf{v}_{\max}$  (in red color).

(x) From Fig. (3), we could see that  $\mathbf{L}\mathbf{v}_{\max}$  is the semi-major axis of the distribution of all the distorted vectors  $\tilde{\mathbf{r}}$ . Since eigenvectors point in directions that are stretched by the transformation from matrix  $\mathbf{L}$ , the eigenvector  $\mathbf{v}_{\max}$ , which corresponds to the largest eigenvalue, points to the direction that is stretched most severely. As a result, after the transformation of matrix  $\mathbf{L}$ , the set of vectors  $\mathbf{r}$  which originally have randomly distributed directions became a set of vectors  $\tilde{\mathbf{r}}$  which have distorted distribution of directions. Most of them tend to point to the direction (or the reverse direction) of  $\mathbf{v}_{\max}$ .

HW0.4.2.py

```

1 import numpy as np
2 import random
3 import matplotlib.pyplot as plt
4 nV = 500
5 #create a 2*2 matrix L
6 L = np.matrix([[5/4, -3/2], [-3/2, 5]])
7 #create a 2*500 array R to represent 500 random, unit length,
   2-d vectors, in which R[:, i] is the i-th vector

```



```

8 R = np.zeros((2, nV))
9 for i in range(nV):
10     tmp = np.array([random.gauss(0,1), random.gauss(0,1)])
11     tmp = tmp/(sum(tmp * tmp) ** 0.5)
12     R[:, i] = tmp
13 #compute R2 = LR, R2[:, i] is the distorted R[:, i]
14 R2 = np.dot(L,R)
15 #compute the eigenvalues of L and denote the minimum eigenvalue
    with lmax and lmin.
16 evals, evecs = np.linalg.eig(L)
17 [lMax, lMin] = sorted(evals, reverse = True)
18 #lr[i] is the length of the i-th vector in R2
19 lr = np.zeros((nV))
20 for i in range(nV):
21     lr[i] = (np.dot(R2[:, i].transpose(), R2[:, i]))[0, 0] **
        0.5
22 #create a histogram of values of lr (use 50 bins) and compare
    it to lMax and lMin.
23 plt.hist(lr, bins = 50)
24 plt.axvline(lMax, color='r', linestyle='dashed', linewidth=1)
25 plt.axvline(lMin, color='k', linestyle='dashed', linewidth=1)
26 plt.show()
27 #compute the eigenvectors of L and Let vmax denote the
    eigenvector corresponding to the maximum eigenvalue lmax.
28 for i in range(2):
29     if evals[i] == lMax:
30         vMax = evecs[:, i]
31         break
32 #make a two-dimensional plot of R2 (in blue color) and the
    eigenvector vmax (in red color).
33 Lv = np.asarray(np.dot(L, vMax))
34 fig = plt.figure(figsize = (5, 5))
35 plt.quiver([0], [0], np.asarray(R2[0, :]), np.asarray(R2[1, :])
    , color = 'b', angles='xy', scale_units='xy', scale=1)
36 plt.quiver([0], [0], Lv[0], Lv[1], color = 'r', angles='xy',
    scale_units='xy', scale=1)
37 plt.xlim(-5.5, 5.5)
38 plt.ylim(-5.5, 5.5)
39 plt.show()

```