# Project #2 - Samples and Statistics

LI YICHENG[*]

USCID:7827077047
email: l.y.c.liyicheng@gmail.com
USC Viterbi of Engineering

## I. MICROCHIPS TEST

### I. Problem Description

(a) There is a batch of 125 microchips, if the distributor find any defective microchip in the chosen microchip, the batch of microchips would be rejected. Suppose there are 6 defective units in the lot, what is the probabilty that the distributor will reject the lot if five microchips are tested.
(b) What is the fewest number of microchips that the distributor should test to reject this lot 95% of the time.

### II. Experiment(a)

A function called reject_prop is generated to show the probability that the distributor would reject the lot with test_amount microchips were tested.

*Algorithm:*

*1.Let the defective microchip be No.1 to No.6.
*2.Generate test_amount samples uniformly in 1 to 125
*3.If there exist index form 1 to 6, reject times plus one
*4.Repeat the above routine 20000 times, calculate the amount of rejections in the 20000 times, to derive the rejection probability

### III. core code: filename: reject_prop.m

```
1  function reject_p = reject_prop(test_amount)
2  %usage: reject_p = reject_prop(test_amount)
3  %given the test amount, output the reject probability
4  % input : test_amount : the test amount
5  % output: reject_p: reject probability
6  reject = 0;
7  for i = 1 : 20000 %to be accurate enough, do 20000 trial
8      samples = randperm(125,test_amount);%select test amount microchips
9      if ( sum(samples < 7) > 0) %let defective chips be No.1 to No.6 , if any of these
           number is tested just reject
10          reject = reject +1; %reject,detected such defective chip
11      end
12  end
13  reject_p = reject / 20000;
```

[*]github link: https://github.com/IAMLYCHEE/EE511-PROJECT1

## IV.   Result(a) & Analysis:

If 5 chips were tested, we query: reject_p = reject_prop(5)
*we get:*
reject_p = 0.2181
This means, in this simulation, if 5 microchips were tested, the distributor has a probability near 22% to reject the lot.
Let us compare this to the theoretical value. Because the probabilty that the distributor does not reject the lot is easier to compute, we compute that number. Using combination knowledge, we choose 5 from 199 good chips and we choose 5 from all microchips.
$\mathbb{P}(not\ reject) = \frac{\binom{119}{5}}{\binom{125}{5}} = 0.7787$
$\mathbb{P}(reject) = 1 - \mathbb{P}(not\ reject) = 0.2213$
Compare to our result 0.2181, they are very close.

## V.   Experiment(b)

We make use of the function generated in experiment(a), we just need to change the input parameter, let it increase, until the probabilty reach 95%, then we record the test_amount.

*Algorithm:*

*1.Calculate the reject probabilty under test_amount
*2.If the probability is less than 95%, increase the test_amount
*3.else record the test_amount

## VI.   core code: filename: solution1.m

```
1  test_amount = 5;
2  while(1)
3      p = reject_prop(test_amount);
4      if p < 0.95
5          test_amount = test_amount + 1;
6      else
7          break
8      end
9  end
```

## VII.   Result(b) & Analysis

It is not a function file, just run it.
*we get:* test_amount = 49
that means, if the distributor test 49 microchips, there is a probability of 95% that he would reject the lot.
Also, let us to compare this to our theoretical value, this time the equation should be:
$\mathbb{P}(reject) = 1 - \frac{\binom{119}{x}}{\binom{125}{x}} = 0.95$

Actually, in matlab, if you compute like $\binom{119}{59}$, it may not give accurate answer,(the number is too

large!), so we simplfy the formular.

$$\frac{\binom{119}{x}}{\binom{125}{x}} = \frac{\frac{119!}{(119-x)!x!}}{\frac{125!}{(125-x)!x!}} = \frac{119!(125-x)!}{125!(119-x)!} = \frac{(125-x)(124-x)(123-x)(122-x)(121-x)(120-x)}{125\cdot124\cdot123\cdot122\cdot121\cdot120} = \frac{\binom{125-x}{6}}{\binom{125}{6}}$$

$x = 50, \mathbb{P}(reject) = 0.9571$
$x = 49, \mathbb{P}(reject) = 0.9534$
$x = 48, \mathbb{P}(reject) = 0.9495$
$x = 47, \mathbb{P}(reject) = 0.9452$
Therefore, the theoretical value fit our experiment result.
If the distributor would like to reject the lot 95% of the time, he should take 48 or 49 chips to test.

## II.  NUMBER OF EVENTS IN SPECIFIED INTERVALS

### I.  Probelm Description

(a)Suppose 120 cars arrive at a freeway onramp per hour on average, simulate one hour of arrivals to the freeway ramp.Perform a Bernoulli trial to indicate a car arrival within small time interval.
(b)Sampling directly form an equivalent Possion distribution by using inverse transform method.Generate a histogram for the number of arrivals using this method.

### II.  Experiment(a)

The requirement is to derive the amount of cars at some position per hour,suppose we divide one hour into 5000 intervals, the probabilty of one car is within this interval is $\mathbb{P}(within\ the\ interval) = \frac{120}{5000}$ and we count whether there is a car in one interval by generating a random number form the u(0,1) distribution.

*Algorithm*

*1.Generate N random variables from uniform distribution [0,1], N is the number of time interval.
*2.Count the number of these variable whose value is smaller than 120/N. Record the sum as one trial result.
*3.Repeat 1,2 to get more trial results and generate the histogram.

### III.  Core Code: filename: solution2.m (part of code)

```matlab
1  function solution2(trial_budget,lambda)
2  %usage: compute the B−trial result, inverse sampling method, and theoretical value,fora
       possion distribution event
3  %input: trial times, the expected average of occurence of event (lambda)
4
5  %
6  %theoretical p.m.f of Possion distribution
7  P_the = zeros(1,lambda*1.5); %record the P theoretical value
8  x = 1 : lambda*1.5;
9  for k = 1 : lambda*1.5
10     P_the(k) = trial_budget*lambda^k * exp(−lambda)/factorial(k);
11 end
```

```
12 %————————————————————————————
13 %Bernoulli trial method
14 success_amount = zeros(trial_budget,1);
15 N = 5000;
16 P = lambda / N;
17 for i = 1: trial_budget
18     success_amount(i) = sum(rand(N,1) < P);
19 end
20 subplot(1,2,1)
21 histogram(success_amount,lambda*1.5,'BinLimits',[1,lambda*1.5]);
22 title('Bernoulli trial')
23 xlabel('number of cars')
24 ylabel('appear times')
25 %set the same x−axis [1,lambda*1.5] for comparing
26 hold on
27 plot(x,P_the,'r—o');
28 hold off
29 legend('B−trial','theoretical');
```

## IV.  Experiment(b)

The inverse tranform method.Inverse transformation sampling takes uniform samples of a number $u$ u between 0 and 1, interpreted as a probability, and then returns the largest number $x$ x from the domain of the distribution $P(X)$ such that $P(-\infty < X < x) \leq u$. For example, possion has a distribution of $\mathbb{P}(X = i) = \frac{e^{-\lambda}\lambda^i}{i!}$, we know that $\sum_{i=1}^{\infty} \mathbb{P}(X = i) = 1$,first generate a random varible u from uniform distribution [0,1], and then we just accumulate $\mathbb{P}$ from $\mathbb{P}(X = 1)$ until $\sum_{1}^{N} \mathbb{P}(X = i)$reach u, record the N.

*Algorithm:*

*1.Generate a standard uniform distribution variable u, initiate P and accumulate sum F to be P(X=1).
*2.If u > F, F + $P_{next}$ , else break, record the accumulated index i.
*3.Generate i within trial budget.
*the above generate one sample using inverse sampling method*

## V.  Core Code :

**file: gen_one_possion.m**

```
1 function i = gen_one_possion(lambda)
2 %usage function x = gen_one_possion(lambda)
3 u = rand(1);
4 i = 0;
5 p = exp(−lambda);
6 F = p;
7 while u > F
8     p = lambda * p / (i+1);
9     F = F + p; %accumulate
10     i = i+1;
11 end
```

4

**file: solution2.m (inverse sample part)**

```
1  %────────────────────────────────────
2  %inverse tranform method
3  amount_gened = zeros(trial_budget,1); %generated amount
4  for i = 1: trial_budget
5      amount_gened(i) = gen_one_possion(lambda);
6  end
7  subplot(1,2,2)
8  histogram(amount_gened,lambda*1.5,'BinLimits',[1,lambda*1.5]);
9  title('inverse-transform method');
10 xlabel('number of cars')
11 ylabel('appear times')
12 hold on
13 plot(x,P_the,'r—o');
14 hold off
15 legend('inverse-sample','theretical');
```

## VI.    Result (a),(b) & Analysis
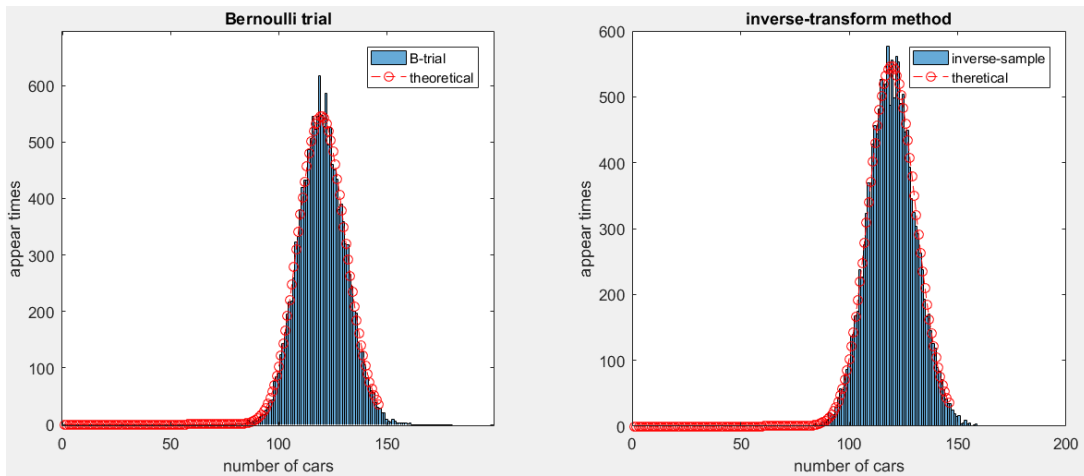
**query: solution2(10000,120)**



**Figure 1:** *outcome histogram*

From the figure, it is easy to observe that both the Bernoulli trial and the inverse sampling method reached the theoretical outcome. Both all even had more probability to get 120 cars.

For Bernoulli trial, theoretically,

$\binom{n}{k}p^k(1-p)^{n-k} \to e^{-r}\frac{r^k}{k!}$ when $n \to \infty$

In these case, the n is 5000, and the p is 12/5000, nearly zero. So the above convergence is suitable in this case.

For the inverse sampling method.

Let $F$ be a continuous cumulative distribution function, and let $F^{-1}$ be its inverse function.

$F^{-1}(u) = \inf \{x \mid F(x) \geq u\}$      $(0 < u < 1)$.

Proof:

$\Pr(F^{-1}(U) \le x)$

$= \Pr(U \le F(x))$    (applying $F$, to both sides)

$= F(x)$          (because $\Pr(U \le y) = y$)

*reference:https://en.wikipedia.org/wiki/Inverse_transform_sampling*

Therefore the inverse sampling method is good for generating sample numbers at random from any probability distribution given its cumulative distribution function.

## III.    AMOUNT OF S.T.D UNIFORM SAMPLE NEEDED TO REACH SUM OF 4

### I.    Problem Description

**(a):** Generate a random variable as the smallest number of uniform random samples whose sum is greater than four.

**(b):** Generate a histogram using 100,1000, and 10000 samples for N.

### II.    Experiment

To generate random variable N is straigtforward:

Algorithm:

*1:Initiate a zero sum.

*2:Generate a random number from standard uniform distribution.

*3:Add the number to the sum. loop time increase 1

*4: if the sum is larger than four, record the loop times and break,

else goto 2.

**filename: generateN.m**

```
1  function N = generateN(thre)
2  %usage N = generateN(thre)
3  %input : thre : the threshold for sum, in this experiment , it is 4
4  %output: the smallest number of sum (u. random samples) > 4
5  sum = 0;
6  N = 0;
7  while sum < thre
8      sum = sum + rand(1);
9      N = N+1;
10 end
```

To generate histogram with defined number of N samples: and derive the expectation, just repeat the above function specific times.Use the following function to generate one histogram with specific amount of N.

**filename: Nhistogram.m**

```
1  function E = Nhistogram(eval_budget,thre)
2  %use: E = Nhistogram(eval_budget,thre)
3  % plot the histogram using eval_budget for N, threshold thre
4  result = zeros(eval_budget,1);
```

```
5  for i = 1 : eval_budget
6      result(i)= generateN(thre);
7  end
8  histogram(result);
9  xlabel('number of uniform samples')
10 ylabel('appear times')
11 E = mean(result);
```

To get the result just use the following querys:

```
1  subplot(1,3,1)
2  E_100=Nhistogram(100,4);
3  title('100')
4  subplot(1,3,2)
5  E_1000=Nhistogram(1000,4);
6  title('1000')
7  subplot(1,3,3)
8  E_10000=Nhistogram(10000,4);
9  title('10000')
```

## III.  Result(a)(b) & Analysis
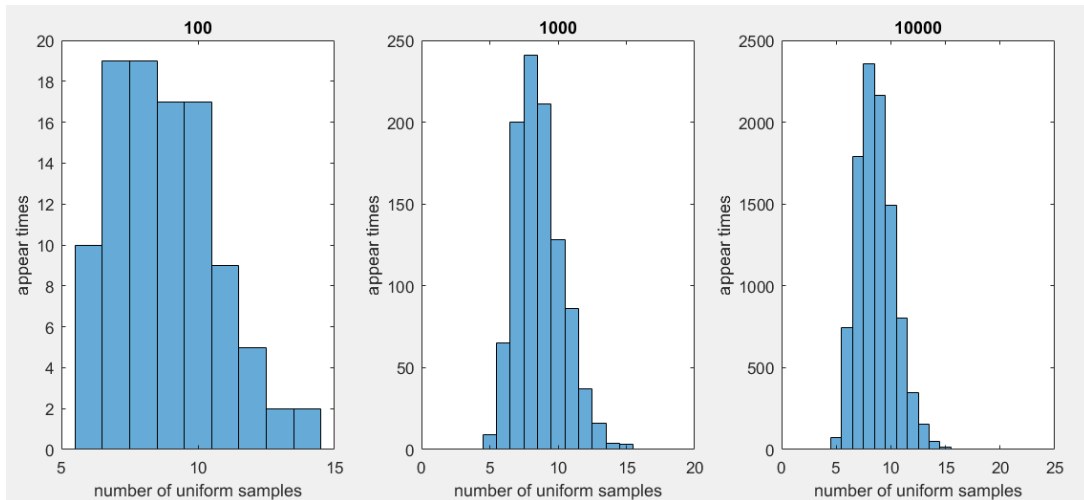
**Graph:**



**Figure 2:** *Histograms for 100,1000,10000 N samples*

**Expectations**
E_100 = 8.750, E_1000 = 8.6570, E_10000 = 8.6652
**Analysis**
When it comes to a larger number of N, there tends to exist a convergence in the figure because it is obviously that the 1000 graph and the 10000 graph looks similar. So let us we analyse theoretical why the shape is look like this.
Actually, it is really difficult to generate a curve which can be described in a continous function to fit the histogram, because it only has ten more bars and it is really a discrete distribution. My thought to analyse the problem is as follows:

*1. Assume N numbers from standard uniform distribution have the sum larger than 4

*2. So k numbers is less then 4/N, while others larger than 4/N, but their average is larger than 4/N(else we can not get sum larger than 4).

*3. Therefore, the probabilty for N amount uniform numbers sum larger than 4 is $\mathbb{P}(N) = \binom{N}{k}(\frac{4}{N})^k(1 - \frac{4}{N})^{N-k}$

*4. Assume k = N/2, that is half of these number is larger than 4/N, another half is smaller than 4/N

*Take a look at this distribution: $\mathbb{P}(N) = \binom{N}{N/2}(\frac{4}{N})^{N/2}(1 - \frac{4}{N})^{N-N/2}$

The following code are generated to analyse the above distribution function:

**filename: P_n.m**

```
1  function p = P_n(N)
2  p = nchoosek(N,floor(N/2))*(4/N)^(N/2)
       *(1-4/N)^(N-(N/2)));
3  end
```

**filename: analysis3.m**

```
1  p = zeros(1,20);
2  for n = 5: 17
3      p(n) = P_n(n);
4  end
5  stem(p)
6  title('analysis')
7  xlabel('N')
8  ylabel('probability')
```
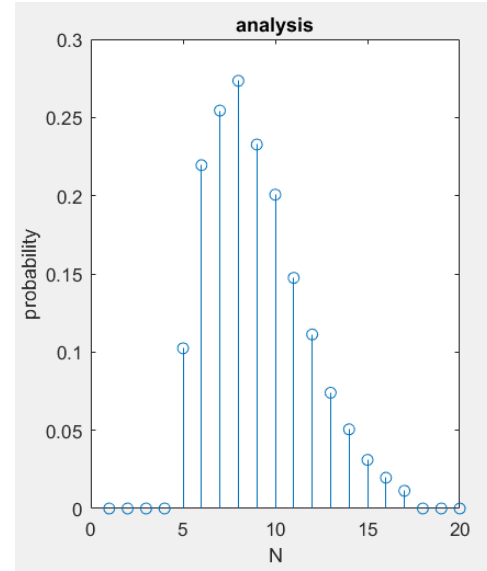
**Theoretical graph**



**Figure 3:** *Theoretical analysis*

    See, the simulation in figure 2 looks like this kind of distribution, the right side of this distribution is actually binomial distribution. Because from the formular, it is easy to derive it has the same decreasing velocity as the right side of binomial distribution.

## IV.  Some kind of sequence analysis

### I.  Problem Description

**(a):** we have a sequence $X_k$ where $p_j = \frac{p}{j}$ for $j = 1, 2, 3, \cdots 60$.Generate a historam showing the behaviour of this sequence.

**(b):**Define the random variable $N_j = min\{k : X_k = j\}$.Simulate sampling from $N_60$.

### II.  Experiment(a):

This is like the inverse sampling method. First, [0,1] is devide into 60 parts with length $p/60, p/59, \cdots, p$ seperately. Then we generate a vairable from standard uniform distribution. And then find out what region does this varible locate.It is easier to be understand by the following

8

algorithm.

*Algorithm:*

*1.Normalize p, let $\sum\limits_{1}^{60} p_j = 1$

*2.Generate a number u from standard uniform distribution
*3.Set accumulate number acc = p/X, X= 60
*4.if u > acc , stop, record X,
else X = X-1,acc=acc+p/X;
goto 4 again.
**filename: generateX.m**

```
1  function X = generateX(p)
2  %usage: result = generateN_4(p)
3  %input : the constant p
4  X = 60;
5  randNumber = rand(1);
6  acc = p/X;
7  while (randNumber > acc && X >=1)
8      X = X−1;
9      acc=acc+p/X;
10 end
```
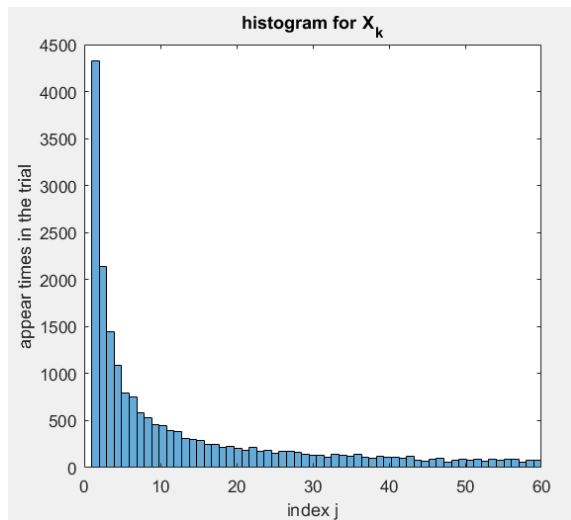
**filename: solution4.m (X histogram part)**

```
1  %———————————————————————————
2  %generate the constant p to make it normalized
3  sum =0;
4  for i = 1 : 60
5      sum = sum +1/i;
6  end
7  p = 1/sum;
8  %———————————————————————————
9  %generate the histogram, do 10000 trials
10 %basic idea is kind of similar to inverse sampling
11 result = zeros(10000,1);
12 for k = 1 : 10000
13     result(k) = generateX(p);
14 end
15 histogram(result)
16 title('histogram for {X_k}')
17 xlabel('index j')
18 ylabel('appear times in the trial')
```

III.   Result(a)&Analysis

**figure:**

**Figure 4:** *X distribution*

**Analysis** Using inverse sampling method works well in this situation, because we can see the second bar has half height as the first bar and the third bar has 1/3 height to the first bar which is corresponding to p,p/2,p/3....p/60. Although there seems to be some 'noise' when it come to p/45 to p/60(not smmoth enough), given their relatively low probability, such noise can be ignored.

## IV.    Experiment(b)

The above we have already define a function to generate one X, the requirement in this part is straight forward, if the X is not 60, keep generate X, else stop and record the generate times.

*Algorithm*

*1.Init times be 0
*2.Generate one X
*3.X is 60, stop, else times plus one, goto 2
*4.Repeat 1,2,3 several times.
**filename:solution4.m (N_60 part)**

```
1  %
2  % N_60 experiment
3  result_N60 = zeros(1000,1);
4  for i = 1 : 1000
5      N = 0;
6      while generateX(p) ~= 60
7          N = N +1;
8      end
9      result_N60(i) = N;
10 end
11 figure
12 histogram(result_N60)
13 title('histogram for {N_{60}}')
14 xlabel('the least amount of trials to reach N_{60}')
15 ylabel('the number of such amount of trials')
16 Expectation = mean(result_N60);
17 variance = var(result_N60);
18 %
```
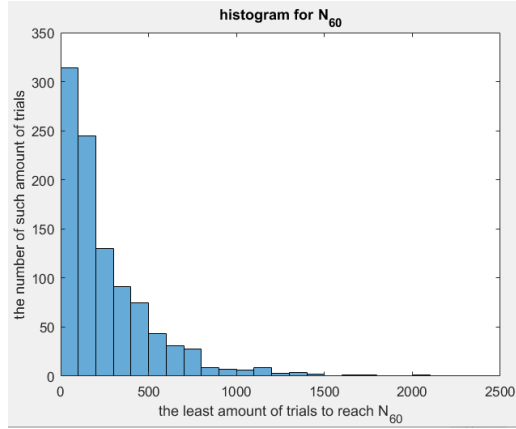
## V.   Result(b) & Analysis

**Figure & Data**



**Figure 5:** *N_60 Experiment*

Expectation = 280.2020
variance = 7.0729e+04.
**Analysis**
This event is actually conform to a geomerty distribution with the probability hitting 60 is $p/60$, not hitting, $\frac{60-p}{60}$. Therefore, theretical expectation is $\frac{1}{p}$, and the variance $\frac{(1-p)}{p^2}$. In this problem, $p = \frac{p_0}{60}$, and $p_0$ is the normalized one we generated in Experiment(a), which is 0.2137. Therefore, $\mathbb{P}(hit\ 60) = 0.2137/60$.
**Theoretical expectation:** $1/\mathbb{P}(hit\ 60) = 60/0.2137 = 280.7674$
**Experiment expectation: 280.2020.**
**Theoretical variance**: $(1 - \mathbb{P}(hit\ 60))/\mathbb{P}(hit\ 60)^2 = 7.8563e + 04$
**Experiment variance: 7.0729e+04**

To sum up, in this experiment, the simulation data conform the theoretical data. Therefore we are doing a experiment which conforms the geometry distribution.

## V.   Accept-reject method

### I.   Problem Description

Use accept-reject method to sample from $p_j$ by sampling from another uniform distribution$q_j$. $p_j$ conforms $p_{1\ to\ 5} = 0.06, p_6 = p_9 = 0.15, p_7 = p_10 = 0.13, p_8 = 0.14$.

### II.   Experiment:

**(a):**Understanding the accept-reject method. The rejection sampling method generates sampling values from a target X with arbitrary probability density function f(x) by using a proposal distribution Y with probability density g(x). The idea is that one can generate a sample value from X by instead sampling from Y and accepting the sample from Y with probability $f(x)/(Mg(x))$, repeating the draws from Y until a value is accepted.
**(b):**First, we derive the constant M, M need to be make sure that $Mg(x) > f(x)$ should be true for every x, therefore, we take the max of $(f(x)/g(x))$ to determine the M. Then we using the accept-reject method to derive $f(x)$.
**(c):**Overlay the target distribution according to the accepted amount, for example we accept 200 samples in Experiment(b), then the target distribution is [200*0.06,200*0.06,...200*0.14], and compute the efficiency.

*Algorithm:*

*1. Set distribution q to be $q_1 = 0.1, q_2 = 0.1, \cdots, q_{10} = 0.1$

*2. Compute the constant M

*3. Generate a rand number i from 1 to 10, generate a number u from std U[0,1]

*4. if u< p(i)/Mq(i) , accept and record, else reject.

*5. Repeat 3,4 for trail_budget times.

**filename: solution5.m**

```
 1 function [mean_exp, mean_the, var_exp, var_the, efficiency] = solution5(trial_budget)
 2 %————————————————————————————————————————————————————————
 3 %usage:[mean_exp, mean_the, var_exp, var_the, efficiency] = solution5(trial_budget)
 4 %trial_budget: the number of times you would like to query for generating
 5 %the sample using accept−reject method
 6 %generate the sample mean and variance using accept−reject method
 7 %generate the theoretical mean and variance
 8 %generate the efficiency of the accept−reject method
 9 %Li Yicheng
10 %————————————————————————————————————————————————————————
11
12 %accept−reject method
13 p = [0.06*ones(1,5),0.15,0.13,0.14,0.15,0.13];
14 %the distribution of p_j
15 q = ones(1,10)*0.1;
16 %the distribution to sample from
17 M = max(p./q);
18 %derive the coeffience to let all Mq >= p
19 i=1;
20 for k = 1 : trial_budget
21     y = randi(10);
22     u = rand(1);
23     if u < p(y)/(M*q(y))
24         accept(i) = y;
25         i=i+1;
26     end
27 end
28 %plot part
29 histogram(accept');
30 xlabel('index 1 to 10')
31 ylabel('appear times')
32 %————————————————————————————————————————————————————————
33 %the target distribution
34 amount = length(accept);
35 %to compare the two graph we need to set them have same total sample amount
36 index = 1;
37 for i = 1 : 10
38     std_distri(index:index + round(amount*p(i))−1)= i*ones(round(amount*p(i)),1);
39     %add the 'i's into the sequence, for example in the first loop add
40     %(amount * 0.06 ) ones into the sequence from 1 to amount * 0.06;
41     index = index + round(amount*p(i));
42 end
43 %overlay with the theoretical
44 hold on
45 histogram(std_distri);
46 title('accept−reject sampling')
47 legend('accept−reject','target')
48
49 %the sample mean,variance and the theoretical values
50 mean_exp = mean(accept);
51 mean_the = (1:10)*p';
```
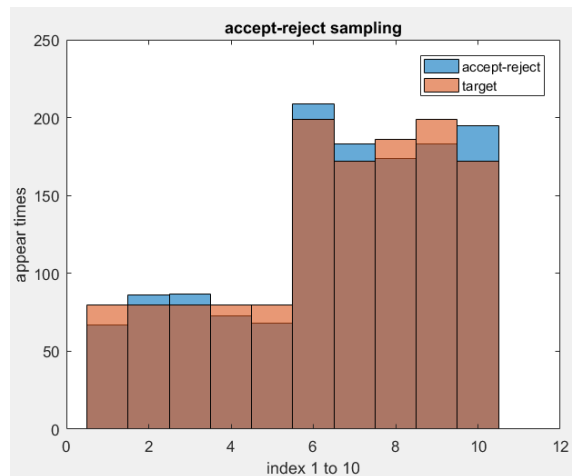
```
52  var_exp = var(accept);
53  var_the = var(std_distri);
54  %the efficiency
55  efficiency = amount/trial_budget;
```

## III.   Result & Analysis

**query: [m_exp,m_the,v_exp,v_the,efficiency]= solution5(2000)**



**m_exp=6.3273**
**m_the=6.4800**
**v_exp=7.7341**
**v_the=7.2164**
**efficiency=0.6630**
**Theoretical efficience** $1/M = 1/(max(p_j/q_j)) = 1/(0.15/0.1) = 1/1.5 = 0.6666$

**Figure 6:** *Experiment result & Target distribution*

All the data derived from experiment is close to theoretical value. Therefore, accept-reject sampling is also a good method to sample from some ditribution by sampling from another distribution. The only drawback may happen is just the efficiency. Because sometimes we may get a lot of unwanted data. So the choice of M should be considered if using this method.