

# Stat-581: Probability and Statistical Inference for Data Science

## Quiz - 3

### Sequential Probability Ratio Test - Bernoulli

Prakruti Joshi (phj15)

Keya Desai (kd706)

Twisha Naik (tn268)

#### 1. SPRT function for Bernoulli

```
strp_bernoulli<-function(alpha0 = 0.01, alpha1 = 0.01, p1 = 0.45, p2 = 0.55, bern_p =0.3){
  # alpha0 is Type1 error  alpha1 is Type2 error
  # p<=p1: NULL Hypothesis  p>=p1: Alternate Hypothesis
  S = 0
  log_likelihood = 0
  ## to keep track of number of steps required for convergence ##
  n_converge = 0
  # calculating threshold for stopping
  A = log(alpha1/(1 - alpha0))
  B = log((1 - alpha1)/alpha0)
  hypo_accepted = -1

  while(TRUE){
    n_converge = n_converge + 1
    # generating bernoulli RV with p = bern_p
    data_point = rbinom(1, 1, bern_p)
    # Log-likelihood ratio
    log_likelihood = (data_point*p2 + (1-data_point)*(1-p2)) -
      (data_point*p1 + (1-data_point)*(1-p1))
    # cumulative sum of the log-likelihood ratio
    S = S + log_likelihood
    # Stopping Rule #
    if(S>=B){
      #Accept H1
      hypo_accepted = 1
      break
    }
    if(S<=A){
      #Accept H0
      hypo_accepted = 0
      break
    }
  }
  return(list(n_converge = n_converge, hypo_accepted = hypo_accepted))
}
```

Define a pair of hypotheses:

1. Null hypothesis  $H_0$ : The  $p$  value of underlying Bernoulli random variable is  $p_1=0.45$
2. Alternate hypothesis  $H_1$ : The  $p$  value of underlying Bernoulli random variable is  $p_2=0.55$

For a given sample from Bernoulli distribution, compute the log likelihood  $\log\lambda_i$ :

- Likelihood function for a Bernoulli sample  
 $f_p(x) = p^x(1-p)^{1-x}$
- Log-Likelihood function for a Bernoulli sample  
 $\log(f_p(x)) = xp + (1-x)(1-p)$

$$\begin{aligned}\log\lambda(x) &= \log(f_{p_2}(x)/f_{p_1}(x)) \\ &= \log(f_{p_2}(x)) - \log(f_{p_1}(x)) \\ &= [xp_2 + (1-x)(1-p_2)] - [xp_1 + (1-x)(1-p_1)]\end{aligned}$$

Define thresholds:

$$\begin{aligned}\alpha_0 &= \text{desired type I error} \\ \alpha_1 &= \text{desired type II error}\end{aligned}$$

$$\begin{aligned}A &= (\alpha_1/(1-\alpha_0)) \\ B &= ((1-\alpha_1)/\alpha_0)\end{aligned}$$

Given  $\alpha_0 = 0.1$  and  $\alpha_1 = 0.1$ ,  $A = -4.59512$  and  $B = 4.59512$

Cumulative sum of log likelihoods =  $S_i$

Start with  $S_0 = 0$

Generate a new bernoulli random variable  $x_i$

$$S_i = S_{i-1} + \log\lambda(x_i)$$

Stopping conditions:

1.  $A < S_i < B$  : continue monitoring (critical inequality)
2.  $S_i \leq B$  : Accept  $H_1$
3.  $S_i \geq A$  : Accept  $H_0$

## 2. Simulation function for Bernoulli

```
simulate_strp_bernoulli <- function(bern_p = 0.3, nsim = 100){  
  
  sum_n = 0  
  H0_count = 0  
  H1_count = 0  
  
  ## Averaging over the STPR function  
  for(i in c(1:nsim)){  
    # calling the STPR function which generates bernoulli variables with p = bern_p #  
    strp_result = strp_bernoulli(bern_p=bern_p)  
  
    ## if H0 is accepted  
    if(strp_result$hypo_accepted == 0){  
      H0_count = H0_count + 1  
    }  
    ## if H1 is true  
    if(strp_result$hypo_accepted == 1){  
      H1_count = H1_count + 1  
    }  
  
    sum_n = sum_n + strp_result$n_converge  
  }  
  avg_steps = sum_n/nsim  
  
  return(list(avg_steps=avg_steps, H0_count=H0_count, H1_count=H1_count))  
}
```

### 3. Test results including all output

#### a. Run it on a sequence of x's distributed Bernoulli .3, and .56

p_values	Average Steps for convergence	H0_count	H1_count
0.3	112.92	100	0
0.46	602.4	100	0
0.5	2021.58	55	45
0.54	597	0	100
0.56	484.64	0	100

#### b. What do you think it would do for .54 ? Try it. Why does it give the result you got?

Intuitively setting  $p = 0.54$  should generate random bernoulli variables having higher likelihood to Hypothesis 1 (p value of underlying Bernoulli random variable is 0.55) than Hypothesis 0 (The p value of underlying Bernoulli random variable is 0.45). This can be determined and verified by the SPRT function for Bernoulli defined above. We tested this by running the simulation 100 times and the results are shown in Table-1. Thus, for  $p = 0.54$ , hypothesis 1 is true as seen.

For  $p = 0.5$ , it has equal likelihood to both hypothesis 0 and hypothesis 1, since it is the midpoint of (0.45,0.55). Thus, the ratio of acceptance of hypothesis 0 to acceptance of hypothesis 1 is close to 1 i.e. both are accepted nearly equal number of times. For 10 trails, hypothesis 1 might be true 5 times and hypothesis 0 maybe true 5 times. The ratio of hypothesis 0 : hypothesis 1 might be 4:6 or the inverse. One of the sample results are shown in (Table 1) on simulating the SPRT bernoulli 100 times.

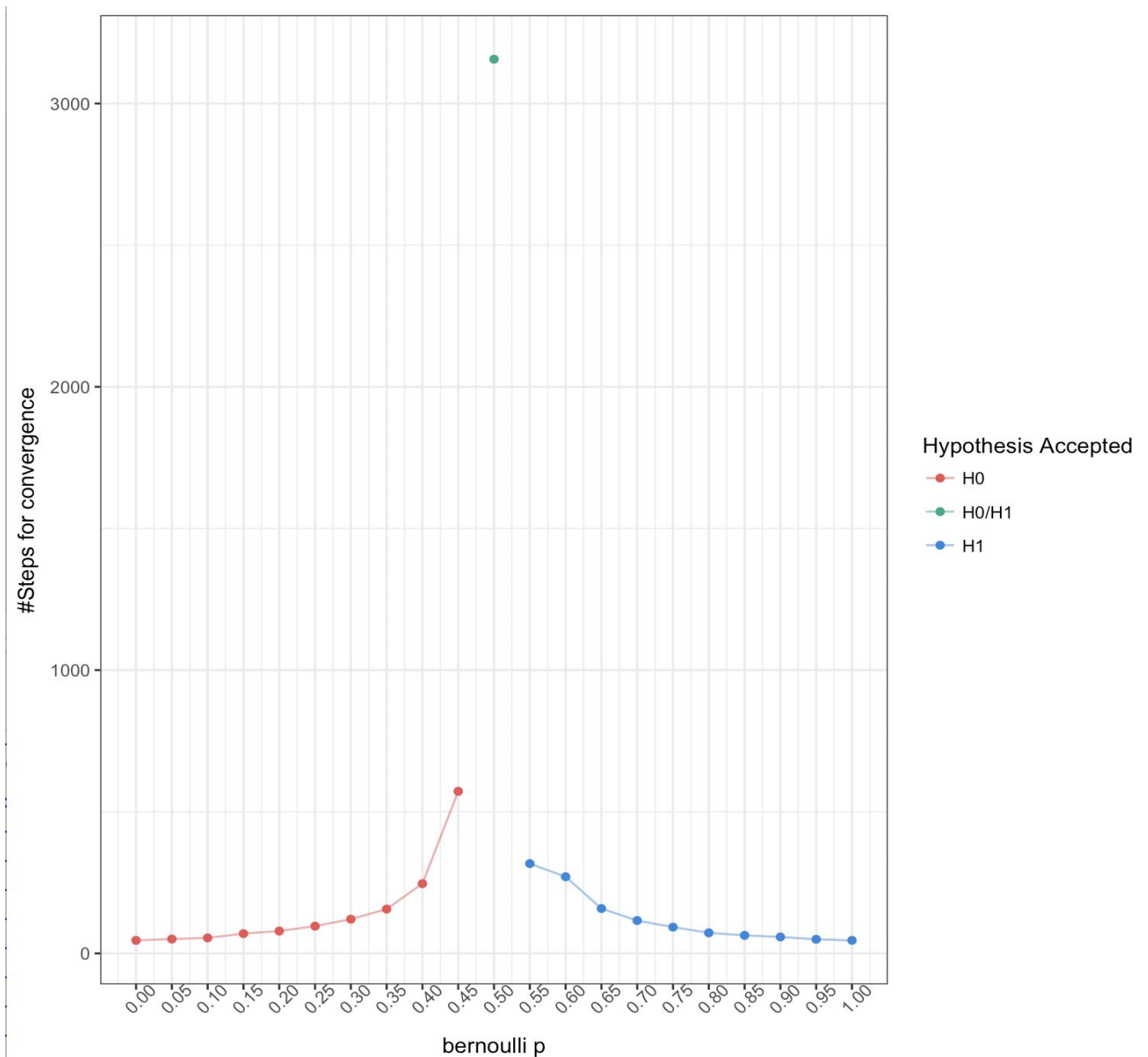
### Extended Analysis:

Running the simulation for all the values of  $p$  ranging from 0 to 1 and running the simulation 100 times for each  $p$  value.

#### Data Frame:

▲	bern_p ▲	avg_steps_to_converge ▲	count_H0 ▲	count_H1 ▲	final_hypothesis_accepted ▲
1	0.00	46.0	10	0	H0
2	0.05	50.6	10	0	H0
3	0.10	54.8	10	0	H0
4	0.15	69.8	10	0	H0
5	0.20	79.0	10	0	H0
6	0.25	96.2	10	0	H0
7	0.30	120.8	10	0	H0
8	0.35	156.4	10	0	H0
9	0.40	246.2	10	0	H0
10	0.45	572.2	10	0	H0
11	0.50	3156.4	5	5	H0/H1
12	0.55	317.0	0	10	H1
13	0.60	270.6	0	10	H1
14	0.65	158.4	0	10	H1
15	0.70	116.0	0	10	H1
16	0.75	92.8	0	10	H1
17	0.80	72.6	0	10	H1
18	0.85	63.6	0	10	H1
19	0.90	58.2	0	10	H1
20	0.95	49.8	0	10	H1
21	1.00	46.0	0	10	H1

Plot:



**Key Observations:**

1. The number of simulations is maximum for  $p = 0.5$ . This is because the Bernoulli random variable with  $p = 0.5$  is equidistant from the set of both the Bernoulli random variable with  $p = 0.45$  and  $p = 0.55$ . Thus, the value of  $S_i$  will keep on fluctuating until it converges to either H0 or H1. This drastically increases the number of simulations needed for convergence.

2. In general, the  $p$  values ranging from  $(0.45, 0.55)$  require greater number of simulations. This can be evidently seen from the graph.
3. For  $p=0$  and  $p=1$ , the bernoulli random variable becomes deterministic.
  - a.  $p=0: x=0 \Rightarrow \log \text{likelihood} = (1-0.55) - (1-0.45) = -0.1$
  - b.  $p=1: x=1 \Rightarrow \log \text{likelihood} = 0.55 - 0.45 = 0.1$

Hence the  $S_i$  will uniformly go down below  $A = -4.5$  or rises above  $B = 4.5$  by taking a constant 46 number of steps.

4. Hypothesis result:

For  $p < 0.5$  : hypothesis 0 is accepted

For  $p > 0.5$  : hypothesis 1 is accepted

For  $p = 0.5$  : hypothesis 0 and hypothesis 1 are accepted approximately in 1:1 ratio (accepted equal number of times)