

# Variuos results for hypothesis testing

Iyar Lin

14 May, 2020

## Contents

<b>1</b>	<b>Intro</b>	<b>2</b>
<b>2</b>	<b>Notations</b>	<b>2</b>
<b>3</b>	<b>Constructing the test</b>	<b>2</b>
3.1	One sided test . . . . .	2
3.2	2 sided hypothesis test . . . . .	3
3.3	Requiring minimum lift . . . . .	4
<b>4</b>	<b>Power</b>	<b>4</b>
<b>5</b>	<b>MDE (minimum detectable effect)</b>	<b>5</b>
<b>6</b>	<b>Required minimum sample size</b>	<b>6</b>
<b>7</b>	<b>Formulas for combining min required lift, number of sides and number of hypothesis tests</b>	<b>6</b>
7.1	Test statistic . . . . .	6
7.1.1	Simulation validation . . . . .	6
7.2	Power . . . . .	8
7.2.1	Simulation validation . . . . .	8
7.3	MDE . . . . .	11
7.3.1	Simulation validation . . . . .	11
7.4	Required minimum sample size . . . . .	13
7.4.1	Simulation validation . . . . .	13
<b>8</b>	<b>Testing with several samples</b>	<b>15</b>
8.1	Notations . . . . .	15
8.2	Constructing the test . . . . .	16
8.2.1	Simulation validation . . . . .	17
8.3	Power . . . . .	17
8.3.1	Simulation validation . . . . .	17
<b>9</b>	<b>Converting all results for the continuous case</b>	<b>18</b>

*Note: All entries in the table of contents are hyperlinks*

# 1 Intro

This document contains proofs and formulas for several use cases encountered by the author in the context of experimental design. Some background in Statistics is preferable for understanding some of the results (You're welcome to contact the author in case you'd like to get more clarifications).

The first few sections contain proofs. Section Formulas for combining min required lift, number of sides and number of hypothesis tests contains unified formulae along with simulation validation.

## 2 Notations

Let a control population samples be denoted by:  $X_1, \dots, X_{n_0} \sim Ber(p_0)$  where  $Ber(p_0)$  is the Bernoulli distribution (meaning  $P(X_i = 1) = p_0$  and conversely  $P(X_i = 0) = 1 - p_0$ ).

We similarly denote the treatment population by  $Y_1, \dots, Y_{n_1} \sim Ber(p_1)$ .

We'd like to test whether applying the treatment results in any kind of lift.

Formally we'd like to test the null hypothesis

$$H_0 : p_1 - p_0 \leq 0$$

vs the alternative:

$$H_1 : p_1 - p_0 > 0$$

We estimate the population distribution parameter  $p$  using the population mean:

$$\hat{p}_0 = \frac{\sum X_j}{n_0}, \hat{p}_1 = \frac{\sum Y_j}{n_1}$$

We note that the "hat" in  $\hat{p}_i$  means this is a random variable which is aimed at estimating the unknown parameter  $p_i$ .

Given that  $\hat{p}_1$  is greater "enough" than  $\hat{p}_0$  we can conclude that we should reject the null  $H_0$  and accept the alternative  $H_1$ .

When rejecting the null we'd like to ensure that the probability we wrongly do so does not exceed some probability  $\alpha$  (often 0.05).

## 3 Constructing the test

### 3.1 One sided test

To that end we'll construct a test such that we reject the null if the difference  $\hat{p}_1 - \hat{p}_0$  exceeds some threshold  $C$  with probability  $\alpha$  when the null is in fact true:

$$P_{H_0}(\hat{p}_1 - \hat{p}_0 > C) = \alpha$$

The above probability is also called type 1 error.

We now turn to finding  $C$ .

We can subtract the mean (which is 0 under the null) from both sides of the inequality and divide by the standard deviation of  $\hat{p}_1 - \hat{p}_0$ :

$$P_{H_0} \left( \frac{\hat{p}_1 - \hat{p}_0}{SD(\hat{p}_1 - \hat{p}_0)} > \frac{C}{SD(\hat{p}_1 - \hat{p}_0)} \right) = \alpha$$

where  $SD(\hat{p}_1 - \hat{p}_0)$  denotes the standard deviation of  $\hat{p}_1 - \hat{p}_0$ .

We note that according to the central limit theorem:

$$\frac{\hat{p}_1 - \hat{p}_0}{SD(\hat{p}_1 - \hat{p}_0)} \xrightarrow[n \rightarrow \infty]{\rightsquigarrow} \mathcal{N}(0, 1)$$

This means that  $\frac{C}{SD(\hat{p}_1 - \hat{p}_0)}$  should equal the  $\alpha$  quantile of the standard normal distribution (in the case of  $\alpha = 0.05$  we have  $\Phi^{-1}(0.95) = 1.65$ ):

$$\frac{C}{SD(\hat{p}_1 - \hat{p}_0)} = \Phi^{-1}(1 - \alpha) \Rightarrow C = \Phi^{-1}(1 - \alpha) \cdot SD(\hat{p}_1 - \hat{p}_0)$$

We note that

$$SD(\hat{p}_1 - \hat{p}_0) = \sqrt{p_1(1 - p_1)/n_1 + p_0(1 - p_0)/n_0}$$

and **finally**

$$C = \Phi^{-1}(1 - \alpha) \cdot \sqrt{\hat{p}_1(1 - \hat{p}_1)/n_1 + \hat{p}_0(1 - \hat{p}_0)/n_0}$$

## 3.2 2 sided hypothesis test

If we we'd like to test if either populations has higher  $p$  (not necessarily treatment higher than control like in the previous section) we'd be doing what's called a 2 sided hypothesis test:

$$H_0 : p_1 - p_0 = 0$$

vs

$$H_1 : p_1 - p_0 \neq 0$$

In this case we'd be constructing a test of the form:

$$P_{H_0}(|\hat{p}_1 - \hat{p}_0| > C^{two}) = \alpha$$

Which translates to

$$P_{H_0}(\hat{p}_1 - \hat{p}_0 > C^{two} \cup \hat{p}_0 - \hat{p}_1 < -C^{two}) = \alpha$$

The events  $\hat{p}_1 - \hat{p}_0 > C^{two}$  and  $\hat{p}_0 - \hat{p}_1 < -C^{two}$  are disjoint thus we have

$$P_{H_0}(\hat{p}_1 - \hat{p}_0 > C^{two} \cup \hat{p}_0 - \hat{p}_1 < -C^{two}) = P_{H_0}(\hat{p}_1 - \hat{p}_0 > C^{two}) + P_{H_0}(\hat{p}_0 - \hat{p}_1 < -C^{two}) = \alpha$$

From symmetry of the Gaussian distribution we get that

$$P_{H_0}(\hat{p}_1 - \hat{p}_0 > C^{two}) = P_{H_0}(\hat{p}_0 - \hat{p}_1 < -C^{two}) = \frac{\alpha}{2}$$

Using the same calculations from the last section we get:

$$C^{two} = \Phi^{-1} \left( 1 - \frac{\alpha}{2} \right) \cdot \sqrt{\hat{p}_1(1 - \hat{p}_1)/n_1 + \hat{p}_0(1 - \hat{p}_0)/n_0}$$

### 3.3 Requiring minimum lift

Sometimes instead of testing  $p_1 > p_0$ , we'd like to test a more demanding criteria  $p_1 > p_0 + \gamma$  where  $\gamma$  is some minimum required lift (e.g. if applying the treatment costs us money such as in the case of direct mail pieces we'd like the treatment lift to be "at least as high as  $\gamma$ ").

We'll thus test:

$$H_0 : p_1 - p_0 \leq \gamma$$

vs the alternative:

$$H_1 : p_1 - p_0 > \gamma$$

Now in order to standardize our statistic we'd subtract *gamma* instead of 0:

$$P_{H_0} \left( \frac{\hat{p}_1 - \hat{p}_0}{SD(\hat{p}_1 - \hat{p}_0)} > \frac{C - \gamma}{SD(\hat{p}_1 - \hat{p}_0)} \right) = \alpha$$

We next have

$$\frac{C - \gamma}{SD(\hat{p}_1 - \hat{p}_0)} = \Phi^{-1}(1 - \alpha) \Rightarrow C - \gamma = \Phi^{-1}(1 - \alpha) \cdot SD(\hat{p}_1 - \hat{p}_0)$$

And finally:

$$C = \Phi^{-1}(1 - \alpha) \cdot \sqrt{\hat{p}_1(1 - \hat{p}_1)/n_1 + \hat{p}_0(1 - \hat{p}_0)/n_0} + \gamma$$

## 4 Power

Let's assume that in reality the alternative is true (so  $p_1 - p_0 > 0$ ).

We denote by  $\beta$  the probability of **not** rejecting the null (also known as type 2 error):

$$\beta = P_{H_1}(\hat{p}_1 - \hat{p}_0 < C)$$

We can then again subtract the mean and divide by the standard deviation (this time under  $H_1$ ) and write:

$$\beta = P_{H_1} \left( \frac{\hat{p}_1 - \hat{p}_0 - (p_1 - p_0)}{SD(\hat{p}_1 - \hat{p}_0)} < \frac{C - (p_1 - p_0)}{SD(\hat{p}_1 - \hat{p}_0)} \right)$$

We note that according to the central limit theorem:

$$\frac{\hat{p}_1 - \hat{p}_0 - (p_1 - p_0)}{SD(\hat{p}_1 - \hat{p}_0)} \underset{n \rightarrow \infty}{\rightsquigarrow} \mathcal{N}(0, 1)$$

We can thus write the below equation:

$$\Phi^{-1}(\beta) = \frac{C - (p_1 - p_0)}{SD(\hat{p}_1 - \hat{p}_0)} \Rightarrow \beta = \Phi \left( \frac{C - (p_1 - p_0)}{SD(\hat{p}_1 - \hat{p}_0)} \right)$$

We can think about the test **power** as the probability of detecting the treatment effect (or conversely, not making the type 2 error  $\beta$ ). We thus have that power is equal to  $1 - \beta$  and:

$$1 - \beta = 1 - \Phi \left( \frac{C - (p_1 - p_0)}{SD(\hat{p}_1 - \hat{p}_0)} \right)$$

And finally the test power is:

$$1 - \beta = 1 - \Phi \left( \frac{\Phi^{-1}(1 - \alpha) \cdot \sqrt{p_1(1 - p_1)/n_1 + p_0(1 - p_0)/n_0} - (p_1 - p_0)}{\sqrt{p_1(1 - p_1)/n_1 + p_0(1 - p_0)/n_0}} \right)$$

We note that in the above equation all hats were removed, as the power calculation takes place at the design phase of an experiment before any samples are available. We can usually have a good “guess” regarding  $p_0$  based on past data.  $p_1$  is usually considered as the Minimum Detectable Effect (MDE) and is chosen by the experiment designers.

## 5 MDE (minimum detectable effect)

Often times it’s useful to choose a sample size and power and see what MDE they yield. Let’s find the formula.

Starting from the equation arrived at in the power section:

$$\Phi^{-1}(\beta) = \frac{C - (p_1 - p_0)}{SD(\hat{p}_1 - \hat{p}_0)}$$

We can further develop:

$$\Phi^{-1}(\beta) \cdot \sqrt{p_1(1 - p_1)/n_1 + p_0(1 - p_0)/n_0} = \Phi^{-1}(1 - \alpha) \cdot \sqrt{p_1(1 - p_1)/n_1 + p_0(1 - p_0)/n_0} - (p_1 - p_0) \Rightarrow$$

$$p_1 - p_0 = (\Phi^{-1}(1 - \alpha) - \Phi^{-1}(\beta)) \sqrt{p_1(1 - p_1)/n_1 + p_0(1 - p_0)/n_0}$$

Finally, using the fact that under  $H_0$  we have  $p_1 = p_0$  we get:

$$MDE = p_1 - p_0 = (\Phi^{-1}(1 - \alpha) - \Phi^{-1}(\beta)) \sqrt{p_0(1 - p_0)/n_1 + p_0(1 - p_0)/n_0}$$

## 6 Required minimum sample size

Most often, an important part of an experiment design is calculating required sample sizes.

In this section we'll find the formula for calculating required minimum sample size assuming the treatment and control group samples are equal (equal samples sizes yield the highest power as demonstrated in section 2). Let's assume we'd like to obtain a test with power of  $1 - \beta$ . Using the result from section 2 we can rearrange (also denote  $n_0 = n_1 = n$ ):

$$\beta = \Phi \left( \frac{\Phi^{-1}(1 - \alpha) \cdot \sqrt{(p_1(1 - p_1) + p_0(1 - p_0)) / n} - (p_1 - p_0)}{\sqrt{(p_1(1 - p_1) + p_0(1 - p_0)) / n}} \right) \Rightarrow$$

$$\Phi^{-1}(\beta) = \frac{\Phi^{-1}(1 - \alpha) \cdot \sqrt{(p_1(1 - p_1) + p_0(1 - p_0)) / n} - (p_1 - p_0)}{\sqrt{(p_1(1 - p_1) + p_0(1 - p_0)) / n}} \Rightarrow$$

$$\Phi^{-1}(\beta) \cdot \sqrt{(p_1(1 - p_1) + p_0(1 - p_0)) / n} + (p_1 - p_0) = \Phi^{-1}(1 - \alpha) \cdot \sqrt{(p_1(1 - p_1) + p_0(1 - p_0)) / n} \Rightarrow$$

$$(p_1 - p_0) = \Phi^{-1}(1 - \alpha) \cdot \sqrt{(p_1(1 - p_1) + p_0(1 - p_0)) / n} - \Phi^{-1}(\beta) \cdot \sqrt{(p_1(1 - p_1) + p_0(1 - p_0)) / n} \Rightarrow$$

$$(p_1 - p_0) = (\Phi^{-1}(1 - \alpha) - \Phi^{-1}(\beta)) \sqrt{(p_1(1 - p_1) + p_0(1 - p_0)) / n}$$

$$\sqrt{n} = \frac{(\Phi^{-1}(1 - \alpha) - \Phi^{-1}(\beta)) \sqrt{p_1(1 - p_1) + p_0(1 - p_0)}}{(p_1 - p_0)}$$

And **finally**:

$$n = \left( \frac{(\Phi^{-1}(1 - \alpha) - \Phi^{-1}(\beta)) \sqrt{p_1(1 - p_1) + p_0(1 - p_0)}}{(p_1 - p_0)} \right)^2$$

## 7 Formulas for combining min required lift, number of sides and number of hypothesis tests

Below we can see the formulas presented above combining min required lift  $\gamma$ , number of sides  $s \in \{1, 2\}$  and number of hypothesis tests  $h \in \mathbb{N}$  (Using Bonferroni correction for multiple hypothesis testing).

### 7.1 Test statistic

The below is true for all cases **except**  $s = 2 \cap \gamma \neq 0$ .

$$C = \Phi^{-1} \left( 1 - \frac{\alpha}{(s \cdot h)} \right) \cdot \sqrt{\hat{p}_1(1 - \hat{p}_1) / n_1 + \hat{p}_0(1 - \hat{p}_0) / n_0} + \gamma$$

#### 7.1.1 Simulation validation

```

rm(list = ls())
C <- function(p_1_hat, n_1, p_0_hat, n_0, alpha, s, h, gamma) {
  if (gamma != 0 & s == 2) stop("using gamma != 0 with s = 2 is not supported yet")
  qnorm(1 - alpha / (s * h)) *
    sqrt(p_1_hat * (1 - p_1_hat) / n_1 +
          p_0_hat * (1 - p_0_hat) / n_0) + gamma
}

p_0 <- 0.2
p_1 <- 0.2 # Null is true
n_0 <- n_1 <- 10000
alpha <- 0.05

s <- 1:2
h <- 1:5
gamma <- seq(-0.03, 0.02, 0.01)
M <- 10000
rejected_null <- array(
  dim = c(length(h), length(gamma), length(s)),
  dimnames = list(paste0("#tests = ", h), paste0("gamma = ", gamma), paste0("sides = ", s))
)

for (i in 1:dim(rejected_null)[1]) {
  for (j in 1:dim(rejected_null)[2]) {
    for (k in 1:dim(rejected_null)[3]) {
      if (gamma[j] != 0 & s[k] == 2) next
      p_1_hat <- as.matrix(replicate(
        n = h[i],
        rbinom(M, size = n_1, prob = p_1 + gamma[j]) / n_1
      ), ncol = h[i])

      p_0_hat <- as.matrix(replicate(
        n = h[i],
        rbinom(M, size = n_0, prob = p_0) / n_0
      ), ncol = h[i])
      const <- mapply(
        function(p_1_hat, p_0_hat) {
          C(
            p_1_hat = p_1_hat, n_1 = rep(n_1, M),
            p_0_hat = p_0_hat, n_0 = rep(n_0, M),
            alpha = alpha, s = s[k], h = h[i],
            gamma = gamma[j]
          )
        },
        split(p_1_hat, rep(1:ncol(p_1_hat),
          each = nrow(p_1_hat)
        )),
        split(p_0_hat, rep(1:ncol(p_0_hat),
          each = nrow(p_0_hat)
        ))
      )

      if (s[k] == 2) {

```

```

    diff <- abs(p_1_hat - p_0_hat)
  } else {
    diff <- p_1_hat - p_0_hat
  }

  rejected_null[i, j, k] <- mean(apply(diff - const, 1, function(row) any(row > 0)))
}
}
}

```

Below we can see the results:

Table 1: sides = 1

	gamma = -0.03	gamma = -0.02	gamma = -0.01	gamma = 0	gamma = 0.01	gamma = 0.02
#tests = 1	0.0523	0.0507	0.0489	0.0502	0.0528	0.0508
#tests = 2	0.049	0.0471	0.0497	0.0504	0.0461	0.0502
#tests = 3	0.0504	0.0503	0.0476	0.0481	0.0498	0.0533
#tests = 4	0.0509	0.0535	0.05	0.0495	0.0493	0.0506
#tests = 5	0.0476	0.0504	0.0482	0.0474	0.0487	0.048

Table 2: sides = 2

#tests = 1	0.05
#tests = 2	0.0467
#tests = 3	0.0503
#tests = 4	0.0489
#tests = 5	0.0468

## 7.2 Power

$$1 - \beta = 1 - \Phi \left( \frac{\Phi^{-1}(1 - \frac{\alpha}{s \cdot h}) \cdot \sqrt{p_1(1 - p_1)/n_1 + p_0(1 - p_0)/n_0} - (p_1 - (p_0 + \gamma))}{\sqrt{p_1(1 - p_1)/n_1 + p_0(1 - p_0)/n_0}} \right)$$

### 7.2.1 Simulation validation

```

power <- function(p_1, n_1, p_0, n_0, alpha, s, h, gamma) {
  if (gamma != 0 & s == 2) stop("using gamma != 0 with s = 2 is not supported yet")
  1 - pnorm((qnorm(1 - alpha / (s * h)) * sqrt(p_1 * (1 - p_1) / n_1 + p_0 * (1 - p_0) / n_0) -
    (p_1 - (p_0 + gamma))) / sqrt(p_1 * (1 - p_1) / n_1 + p_0 * (1 - p_0) / n_0))
}

p_0 <- 0.2
n_0 <- n_1 <- 10000
alpha <- 0.05
coef <- 1.05
s <- 1:2

```



```

h <- 1:5
gamma <- seq(-0.03, 0.02, 0.01)
M <- 10000

rejected_null_calc_power <- array(
  dim = c(length(h), length(gamma), length(s)),
  dimnames = list(paste0("#tests = ", h), paste0("gamma = ", gamma), paste0("sides = ", s))
)

for (i in 1:dim(rejected_null)[1]) {
  for (j in 1:dim(rejected_null)[2]) {
    for (k in 1:dim(rejected_null)[3]) {
      if (gamma[j] != 0 & s[k] == 2) next
      p_1 <- (p_0 + gamma[j]) * coef # alternative is true
      p_1_hat <- as.matrix(replicate(
        n = h[i],
        rbinom(M, size = n_1, prob = p_1) / n_1
      ), ncol = h[i])

      p_0_hat <- as.matrix(replicate(
        n = h[i],
        rbinom(M, size = n_0, prob = p_0) / n_0
      ), ncol = h[i])
      const <- mapply(
        function(p_1_hat, p_0_hat) {
          C(
            p_1_hat = p_1_hat, n_1 = rep(n_1, M),
            p_0_hat = p_0_hat, n_0 = rep(n_0, M),
            alpha = alpha, s = s[k], h = h[i],
            gamma = gamma[j]
          )
        },
        split(p_1_hat, rep(1:ncol(p_1_hat),
          each = nrow(p_1_hat)
        )),
        split(p_0_hat, rep(1:ncol(p_0_hat),
          each = nrow(p_0_hat)
        ))
      )

      if (s[k] == 2) {
        diff <- abs(p_1_hat - p_0_hat)
      } else {
        diff <- p_1_hat - p_0_hat
      }

      rejected_null_calc_power[i, j, k] <- paste0(
        "calc = ",
        round(power(
          p_0 = p_0,
          p_1 = p_1,
          n_1 = n_1,
          n_0 = n_0,

```

```

    alpha = alpha,
    s = s[k],
    h = h[i],
    gamma = gamma[j]
  ), 3),
  ", actual = ",
  round(mean(apply(
    diff - const, 1,
    function(row) {
      row[1] > 0
    }
  )), 3)
)
}
}
}

```

Below we can see the results:

Table 3: sides = 1

	gamma = -0.03	gamma = -0.02	gamma = -0.01	gamma = 0	gamma = 0.01	gamma = 0.02
<b>#tests</b> <b>= 1</b>	calc = 0.456, actual = 0.458	calc = 0.485, actual = 0.491	calc = 0.514, actual = 0.509	calc = 0.543, actual = 0.546	calc = 0.571, actual = 0.57	calc = 0.598, actual = 0.601
<b>#tests</b> <b>= 2</b>	calc = 0.335, actual = 0.34	calc = 0.362, actual = 0.367	calc = 0.39, actual = 0.395	calc = 0.418, actual = 0.418	calc = 0.445, actual = 0.446	calc = 0.473, actual = 0.472
<b>#tests</b> <b>= 3</b>	calc = 0.277, actual = 0.278	calc = 0.302, actual = 0.302	calc = 0.327, actual = 0.326	calc = 0.353, actual = 0.349	calc = 0.38, actual = 0.383	calc = 0.407, actual = 0.398
<b>#tests</b> <b>= 4</b>	calc = 0.24, actual = 0.238	calc = 0.263, actual = 0.264	calc = 0.287, actual = 0.292	calc = 0.312, actual = 0.314	calc = 0.338, actual = 0.341	calc = 0.364, actual = 0.362
<b>#tests</b> <b>= 5</b>	calc = 0.214, actual = 0.213	calc = 0.236, actual = 0.239	calc = 0.259, actual = 0.257	calc = 0.283, actual = 0.282	calc = 0.307, actual = 0.309	calc = 0.332, actual = 0.329

Table 4: sides = 2

<b>#tests = 1</b>	calc = 0.418, actual = 0.415
<b>#tests = 2</b>	calc = 0.312, actual = 0.316
<b>#tests = 3</b>	calc = 0.26, actual = 0.255
<b>#tests = 4</b>	calc = 0.228, actual = 0.223
<b>#tests = 5</b>	calc = 0.205, actual = 0.206

## 7.3 MDE

$$MDE = p_1 - p_0 = \left( \Phi^{-1}\left(1 - \frac{\alpha}{s \cdot h}\right) - \Phi^{-1}(\beta) \right) \sqrt{p_0(1 - p_0)/n_1 + p_0(1 - p_0)/n_0} + \gamma$$

### 7.3.1 Simulation validation

```
MDE <- function(power, n_1, p_0, n_0, alpha, s, h, gamma) {
  if (gamma != 0 & s == 2) stop("using gamma != 0 with s = 2 is not supported yet")
  (qnorm(1 - alpha / (s * h)) - qnorm(1 - power)) *
    sqrt(p_0 * (1 - p_0) / n_1 + p_0 * (1 - p_0) / n_0) + gamma
}

p_0 <- 0.2
n_0 <- n_1 <- 10000
alpha <- 0.05
s <- 1:2
h <- 1:5
gamma <- seq(-0.03, 0.02, 0.01)
M <- 10000
pow <- 0.8
rejected_null_calc_MDE <- array(
  dim = c(length(h), length(gamma), length(s)),
  dimnames = list(paste0("#tests = ", h), paste0("gamma = ", gamma), paste0("sides = ", s))
)

for (i in 1:dim(rejected_null)[1]) {
  for (j in 1:dim(rejected_null)[2]) {
    for (k in 1:dim(rejected_null)[3]) {
      if (gamma[j] != 0 & s[k] == 2) next
      mde <- MDE(
        power = pow, n_1 = n_1, n_0 = n_0, p_0 = p_0, alpha = alpha,
        s = s[k], h = h[i], gamma = gamma[j]
      )
      p_1 <- p_0 + mde # MDE is true

      p_1_hat <- as.matrix(replicate(
        n = h[i],
        rbinom(M, size = n_1, prob = p_1) / n_1
      ), ncol = h[i])

      p_0_hat <- as.matrix(replicate(
        n = h[i],
        rbinom(M, size = n_0, prob = p_0) / n_0
      ), ncol = h[i])
      const <- mapply(
        function(p_1_hat, p_0_hat) {
          C(
            p_1_hat = p_1_hat, n_1 = rep(n_1, M),
            p_0_hat = p_0_hat, n_0 = rep(n_0, M),
            alpha = alpha, s = s[k], h = h[i],
            gamma = gamma[j]
          )
        }
      )
    }
  }
}
```

```

    },
    split(p_1_hat, rep(1:ncol(p_1_hat),
      each = nrow(p_1_hat)
    )),
    split(p_0_hat, rep(1:ncol(p_0_hat),
      each = nrow(p_0_hat)
    ))
  )

  if (s[k] == 2) {
    diff <- abs(p_1_hat - p_0_hat)
  } else {
    diff <- p_1_hat - p_0_hat
  }

  rejected_null_calc_MDE[i, j, k] <- paste0(
    "MDE = ",
    round(mde, 4),
    ", rejected = ",
    round(mean(apply(
      diff - const, 1,
      function(row) {
        row[1] > 0
      }
    )), 3)
  )
}
}
}

```

Below we can see the results:

Table 5: sides = 1

	gamma = -0.03	gamma = -0.02	gamma = -0.01	gamma = 0	gamma = 0.01	gamma = 0.02
<b>#tests</b> <b>= 1</b>	MDE = -0.0159, rejected = 0.803	MDE = -0.0059, rejected = 0.804	MDE = 0.0041, rejected = 0.802	MDE = 0.0141, rejected = 0.79	MDE = 0.0241, rejected = 0.786	MDE = 0.0341, rejected = 0.783
<b>#tests</b> <b>= 2</b>	MDE = -0.0142, rejected = 0.806	MDE = -0.0042, rejected = 0.793	MDE = 0.0058, rejected = 0.789	MDE = 0.0158, rejected = 0.789	MDE = 0.0258, rejected = 0.79	MDE = 0.0358, rejected = 0.777
<b>#tests</b> <b>= 3</b>	MDE = -0.0132, rejected = 0.81	MDE = -0.0032, rejected = 0.8	MDE = 0.0068, rejected = 0.789	MDE = 0.0168, rejected = 0.784	MDE = 0.0268, rejected = 0.787	MDE = 0.0368, rejected = 0.779
<b>#tests</b> <b>= 4</b>	MDE = -0.0126, rejected = 0.814	MDE = -0.0026, rejected = 0.808	MDE = 0.0074, rejected = 0.799	MDE = 0.0174, rejected = 0.79	MDE = 0.0274, rejected = 0.778	MDE = 0.0374, rejected = 0.773

	gamma = -0.03	gamma = -0.02	gamma = -0.01	gamma = 0	gamma = 0.01	gamma = 0.02
#tests = 5	MDE = -0.0121, rejected = 0.805	MDE = -0.0021, rejected = 0.803	MDE = 0.0079, rejected = 0.794	MDE = 0.0179, rejected = 0.793	MDE = 0.0279, rejected = 0.774	MDE = 0.0379, rejected = 0.771

Table 6: sides = 2

#tests = 1	MDE = 0.0158, rejected = 0.793
#tests = 2	MDE = 0.0174, rejected = 0.782
#tests = 3	MDE = 0.0183, rejected = 0.786
#tests = 4	MDE = 0.0189, rejected = 0.779
#tests = 5	MDE = 0.0193, rejected = 0.781

## 7.4 Required minimum sample size

$$n = \left( \frac{(\Phi^{-1}(1 - \frac{\alpha}{s \cdot h}) - \Phi^{-1}(\beta)) \sqrt{p_1(1 - p_1) + p_0(1 - p_0)}}{(p_1 - (p_0 + \gamma))} \right)^2$$

### 7.4.1 Simulation validation

```

min_sample_size <- function(p_1, p_0, alpha, s, h, gamma, power) {
  if (gamma != 0 & s == 2) stop("using gamma != 0 with s = 2 is not supported yet")
  round((((qnorm(1 - alpha / (s * h)) - qnorm(1 - power)) *
    sqrt(p_1 * (1 - p_1) + p_0 *
      (1 - p_0))) / (p_1 - (p_0 + gamma)))^2)
}

p_0 <- 0.2
alpha <- 0.05
coef <- 1.05
s <- 1:2
h <- 1:5
gamma <- seq(-0.03, 0.02, 0.01)
M <- 10000
pow <- 0.8

rejected_null_calc_min_n <- array(
  dim = c(length(h), length(gamma), length(s)),
  dimnames = list(paste0("#tests = ", h), paste0("gamma = ", gamma), paste0("sides = ", s))
)

for (i in 1:dim(rejected_null)[1]) {
  for (j in 1:dim(rejected_null)[2]) {
    for (k in 1:dim(rejected_null)[3]) {
      if (gamma[j] != 0 & s[k] == 2) next
    }
  }
}

```

```

p_1 <- (p_0 + gamma[j]) * coef # alternative is true
n_0 <- n_1 <- min_sample_size(
  p_1 = p_1, p_0 = p_0, alpha = alpha,
  s = s[k], h = h[i], gamma = gamma[j], power = pow
)
p_1_hat <- as.matrix(replicate(
  n = h[i],
  rbinom(M, size = n_1, prob = p_1) / n_1
), ncol = h[i])

p_0_hat <- as.matrix(replicate(
  n = h[i],
  rbinom(M, size = n_0, prob = p_0) / n_0
), ncol = h[i])
const <- mapply(
  function(p_1_hat, p_0_hat) {
    C(
      p_1_hat = p_1_hat, n_1 = rep(n_1, M),
      p_0_hat = p_0_hat, n_0 = rep(n_0, M),
      alpha = alpha, s = s[k], h = h[i],
      gamma = gamma[j]
    )
  },
  split(p_1_hat, rep(1:ncol(p_1_hat),
    each = nrow(p_1_hat)
  )),
  split(p_0_hat, rep(1:ncol(p_0_hat),
    each = nrow(p_0_hat)
  ))
)

if (s[k] == 2) {
  diff <- abs(p_1_hat - p_0_hat)
} else {
  diff <- p_1_hat - p_0_hat
}

rejected_null_calc_min_n[i, j, k] <- paste0(
  "n = ",
  n_0,
  ", rejected = ",
  round(mean(apply(
    diff - const, 1,
    function(row) {
      row[1] > 0
    }
  )), 3)
)
}
}
}

```

Below we can see the results:

Table 7: sides = 1

	gamma = -0.03	gamma = -0.02	gamma = -0.01	gamma = 0	gamma = 0.01	gamma = 0.02
<b>#tests</b> <b>= 1</b>	n = 26240, rejected = 0.794	n = 23912, rejected = 0.801	n = 21901, rejected = 0.797	n = 20149, rejected = 0.8	n = 18611, rejected = 0.802	n = 17252, rejected = 0.793
<b>#tests</b> <b>= 2</b>	n = 33312, rejected = 0.801	n = 30357, rejected = 0.797	n = 27804, rejected = 0.809	n = 25579, rejected = 0.798	n = 23627, rejected = 0.803	n = 21902, rejected = 0.806
<b>#tests</b> <b>= 3</b>	n = 37429, rejected = 0.794	n = 34108, rejected = 0.806	n = 31240, rejected = 0.8	n = 28741, rejected = 0.801	n = 26547, rejected = 0.794	n = 24608, rejected = 0.798
<b>#tests</b> <b>= 4</b>	n = 40341, rejected = 0.796	n = 36762, rejected = 0.805	n = 33670, rejected = 0.796	n = 30977, rejected = 0.794	n = 28613, rejected = 0.807	n = 26523, rejected = 0.801
<b>#tests</b> <b>= 5</b>	n = 42594, rejected = 0.802	n = 38816, rejected = 0.805	n = 35551, rejected = 0.806	n = 32707, rejected = 0.801	n = 30211, rejected = 0.795	n = 28005, rejected = 0.798

Table 8: sides = 2

<b>#tests = 1</b>	n = 25579, rejected = 0.797
<b>#tests = 2</b>	n = 30977, rejected = 0.809
<b>#tests = 3</b>	n = 34119, rejected = 0.804
<b>#tests = 4</b>	n = 36341, rejected = 0.803
<b>#tests = 5</b>	n = 38062, rejected = 0.802

## 8 Testing with several samples

### 8.1 Notations

Sometimes we'd like to collect several samples and do hypothesis testing over them.

Formally speaking let's assume we collect samples in 2 periods (this will be further generalized to T periods later). We have the populations:

First period control:

$$X_1^1, \dots, X_{n_0^1}^1 \sim Ber(p_0^1)$$

First period treatment:

$$Y_1^1, \dots, Y_{n_1^1}^1 \sim Ber(p_1^1)$$

Second period control:

$$X_1^2, \dots, X_{n_0^2}^2 \sim Ber(p_0^2)$$

First period treatment:

$$Y_1^2, \dots, Y_{n_1^2}^2 \sim Ber(p_1^2)$$

While we don't assume  $p_0^1 = p_0^2$  or  $p_1^1 = p_1^2$  we do assume that the lift between treatment and control is the same such that  $p_1^1 - p_0^1 = p_1^2 - p_0^2 = \delta$ .

So far we've really tested

$$H_0 : \delta \leq 0$$

vs

$$H_1 : \delta > 0$$

We can estimate the lift using the average of both periods estimates:

$$\hat{\delta} = \frac{\hat{\delta}_1 + \hat{\delta}_2}{2}$$

where  $\hat{\delta}_1 = \hat{p}_1^1 - \hat{p}_0^1$  and  $\hat{\delta}_2 = \hat{p}_1^2 - \hat{p}_0^2$ .

Let's assume that  $var(\delta_1) < var(\delta_2)$ . The question arises: when should we use  $\hat{\delta}$  instead of  $\hat{\delta}_1$ ?

We'd do so when  $var(\hat{\delta}_1) > var(\hat{\delta})$ . This happens when:

$$var(\hat{\delta}_1) > var(\hat{\delta}) = var\left(\frac{\hat{\delta}_1 + \hat{\delta}_2}{2}\right) = \frac{1}{4} (var(\hat{\delta}_1) + var(\hat{\delta}_2))$$

re-arranging we get

$$3 \cdot var(\hat{\delta}_1) > var(\hat{\delta}_2)$$

Given  $T$  periods we can use this logic to choose periods in an inductive way.

## 8.2 Constructing the test

Again, we'd like to construct a test of the form

$$P_{H_0}(\hat{\delta} > C) = \alpha$$

We note that the standard deviation of our estimator for this case is:

$$SD(\hat{\delta}) = \frac{1}{2} \sqrt{var(\hat{\delta}_1) + var(\hat{\delta}_2)} = \frac{1}{2} \sqrt{p_1^1(1-p_1^1)/n_1^1 + p_0^1(1-p_0^1)/n_0^1 + p_1^2(1-p_1^2)/n_1^2 + p_0^2(1-p_0^2)/n_0^2}$$

We thus have that our constant is:

$$C = \Phi^{-1}(1 - \alpha) \cdot \frac{1}{2} \sqrt{p_1^1(1-p_1^1)/n_1^1 + p_0^1(1-p_0^1)/n_0^1 + p_1^2(1-p_1^2)/n_1^2 + p_0^2(1-p_0^2)/n_0^2}$$

Given  $T$  periods used in the test we have:

$$C = \Phi^{-1}(1 - \alpha) \cdot \frac{1}{T} (var(\delta^1) + var(\delta^2) + \dots + var(\delta^T))$$

The rest of the results in this document can be found in a similar manner by swapping the  $SD$  term.



### 8.2.1 Simulation validation

Let's validate using 3 periods, different sample sized and base conversion rates.

```
rm(list = ls())
p_01 <- p_11 <- 0.2 # Null is true
p_02 <- p_12 <- 0.24
p_03 <- p_13 <- 0.23

n_01 <- 5000
n_11 <- 7000
n_02 <- 4000
n_12 <- 10000
n_03 <- 3000
n_13 <- 5000

alpha <- 0.05

M <- 100000 # number of simulations
p_01_hat <- rbinom(M, size = n_01, prob = p_01) / n_01
p_11_hat <- rbinom(M, size = n_11, prob = p_11) / n_11
p_02_hat <- rbinom(M, size = n_02, prob = p_02) / n_02
p_12_hat <- rbinom(M, size = n_12, prob = p_12) / n_12
p_03_hat <- rbinom(M, size = n_03, prob = p_03) / n_03
p_13_hat <- rbinom(M, size = n_13, prob = p_13) / n_13

var_1_hat <- p_11_hat * (1 - p_11_hat) / n_11 + p_01_hat * (1 - p_01_hat) / n_01
var_2_hat <- p_12_hat * (1 - p_12_hat) / n_12 + p_02_hat * (1 - p_02_hat) / n_02
var_3_hat <- p_13_hat * (1 - p_13_hat) / n_13 + p_03_hat * (1 - p_03_hat) / n_03

C <- qnorm(1 - alpha) * (1 / 3) * sqrt(var_1_hat + var_2_hat + var_3_hat)

diff <- ((p_11_hat - p_01_hat) + (p_12_hat - p_02_hat) + (p_13_hat - p_03_hat)) / 3
rejected_null <- mean(diff > C)
```

We rejected 0.05095 of simulations. Cool.

## 8.3 Power

General formula for  $T$  periods is

$$1 - \beta = 1 - \Phi \left( \frac{\Phi^{-1}(1 - \alpha) \cdot \frac{1}{T} \sqrt{\text{var}(\hat{\delta}^1) + \text{var}(\hat{\delta}^2) + \dots + \text{var}(\hat{\delta}^T)} - \delta}{\frac{1}{T} \sqrt{\text{var}(\hat{\delta}^1) + \text{var}(\hat{\delta}^2) + \dots + \text{var}(\hat{\delta}^T)}} \right)$$

Let's validate that!

### 8.3.1 Simulation validation

```
rm(list = ls())
delta <- 0.003
```

```

p_01 <- 0.2
p_11 <- p_01 + delta
p_02 <- 0.24
p_12 <- p_02 + delta
p_03 <- 0.23
p_13 <- p_03 + delta

n_01 <- 5000
n_11 <- 7000
n_02 <- 4000
n_12 <- 10000
n_03 <- 3000
n_13 <- 5000

var_1 <- p_11 * (1 - p_11) / n_11 + p_01 * (1 - p_01) / n_01
var_2 <- p_12 * (1 - p_12) / n_12 + p_02 * (1 - p_02) / n_02
var_3 <- p_13 * (1 - p_13) / n_13 + p_03 * (1 - p_03) / n_03

alpha <- 0.05

power_calc <- 1 -
  pnorm((qnorm(1 - alpha) * (1 / 3) * sqrt(var_1 + var_2 + var_3) - delta) /
        ((1 / 3) * sqrt(var_1 + var_2 + var_3)))

M <- 100000 # number of simulations
p_01_hat <- rbinom(M, size = n_01, prob = p_01) / n_01
p_11_hat <- rbinom(M, size = n_11, prob = p_11) / n_11
p_02_hat <- rbinom(M, size = n_02, prob = p_02) / n_02
p_12_hat <- rbinom(M, size = n_12, prob = p_12) / n_12
p_03_hat <- rbinom(M, size = n_03, prob = p_03) / n_03
p_13_hat <- rbinom(M, size = n_13, prob = p_13) / n_13

var_1_hat <- p_01_hat * (1 - p_01_hat) / n_11 + p_01_hat * (1 - p_01_hat) / n_01
var_2_hat <- p_02_hat * (1 - p_02_hat) / n_12 + p_02_hat * (1 - p_02_hat) / n_02
var_3_hat <- p_03_hat * (1 - p_03_hat) / n_13 + p_03_hat * (1 - p_03_hat) / n_03

C <- qnorm(1 - alpha) * (1 / 3) * sqrt(var_1_hat + var_2_hat + var_3_hat)

diff <- ((p_11_hat - p_01_hat) + (p_12_hat - p_02_hat) + (p_13_hat - p_03_hat)) / 3
rejected_null <- mean(diff > C)

```

The calculated power is 0.151632. The fraction of rejected simulations is 0.15435. Pretty close.

## 9 Converting all results for the continuous case

All results in this document are derived for the binary case.

In case our distributions of interest are continuous with means  $\mu_0, \mu_1$  and variances  $\sigma_0^2, \sigma_1^2$  (note that they don't have to be necessarily Gaussian) then all results in this paper can be used, converting in all formulas:

$$p \rightarrow \mu$$

$$p \rightarrow \mu$$

$$SD(\hat{\mu}_1 - \hat{\mu}_0) \rightarrow \sqrt{\sigma_0^2/n_1 + \sigma_0^2/n_0}$$

$$SD(\hat{\mu}_1 - \hat{\mu}_0) \rightarrow \sqrt{\sigma_1^2/n_1 + \sigma_0^2/n_0}$$

So for example in the continuous case the test constant  $C$  would be

$$C = \Phi^{-1}(1 - \alpha) \cdot \sqrt{\hat{\sigma}_0^2/n_1 + \hat{\sigma}_0^2/n_0}$$

```
rm(list = ls())
mu_0 <- mu_1 <- 10 # Null is true
sigma_0 <- sigma_1 <- 4
n_0 <- n_1 <- 1000
alpha <- 0.05

M <- 100000 # number of simulations
mu_0_samples <- replicate(n = M, rnorm(n = n_0, mean = mu_0, sd = sqrt(sigma_0)))
mu_1_samples <- replicate(n = M, rnorm(n = n_1, mean = mu_1, sd = sqrt(sigma_1)))

mu_0_hat <- apply(mu_0_samples, 2, mean)
sigma_0_hat <- apply(mu_0_samples, 2, var)
mu_1_hat <- apply(mu_1_samples, 2, mean)
C <- qnorm(1 - alpha) * sqrt(2 * sigma_0_hat / n_0)

diff <- mu_1_hat - mu_0_hat
rejected_null <- mean(diff > C)
```

**9.0.0.1 Simulation validation** The fraction of simulations we rejected the null was 0.05064, pretty close to our chosen  $\alpha = 0.05$ .