# DATA2002

## Multiple testing

Garth Tarr

What is real?

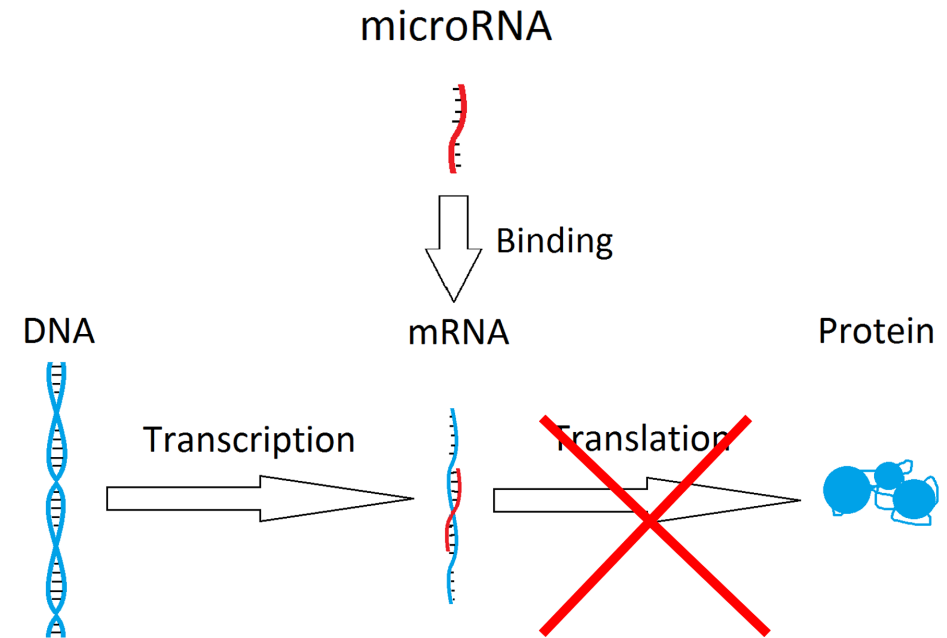Types of errors

Controlling error rates

# microRNA and Alzheimer's disease

# microRNA and Alzheimer's disease

MicroRNA are small non-coding RNA molecules that regulate gene expression.

> 💬 Is there any evidence that microRNA behaviour in the brain might be associated with Alzheimer's disease? (Patrick, Rajagopal, Wong, et al., 2017)

- Experiment by measuring the amount of 309 microRNAs in 701 subjects.

- Test for significant differences between the means of subjects with and without Alzheimer's disease for each microRNA.

# What does this microRNA data look like?

```
# raw data contains objects: AD; microRNA_Data
load("data/microRNA_full.RData")
# disease status in named vector AD
head(AD, n = 4)
```

```
## 20264936 50105725 20730959 11229148
##        1        0        1        1
```

```
disease_status = data.frame(AD) %>%
  tibble::rownames_to_column("subject")
# measurements in data frame microRNA_Data
# columns are subjects, rows are miRNAs
microRNA_Data[1:5, 1:3]
```

```
##              20264936 50105725 20730959
## hsa-let-7a 11.97670 12.61142 13.34787
## hsa-let-7b 12.24761 11.77856 11.44760
## hsa-let-7c 11.07564 11.68883 11.79726
## hsa-let-7d 11.30709 11.50804 11.37125
## hsa-let-7e 10.05073  9.80769 10.57078
```

Reshape the microRNA data and merge in the disease status information.

```
mirna = microRNA_Data %>%
  rownames_to_column("microRNA") %>%
  pivot_longer(cols = -1,
               names_to = "subject",
               values_to = "value") %>%
  left_join(disease_status)
head(mirna)
```

```
## # A tibble: 6 × 4
##   microRNA   subject   value    AD
##   <chr>      <chr>     <dbl> <int>
## 1 hsa-let-7a 20264936  12.0      1
## 2 hsa-let-7a 50105725  12.6      0
## 3 hsa-let-7a 20730959  13.3      1
## 4 hsa-let-7a 11229148  12.5      1
## 5 hsa-let-7a 20151388  12.3      1
## 6 hsa-let-7a 11259716  12.4      1
```
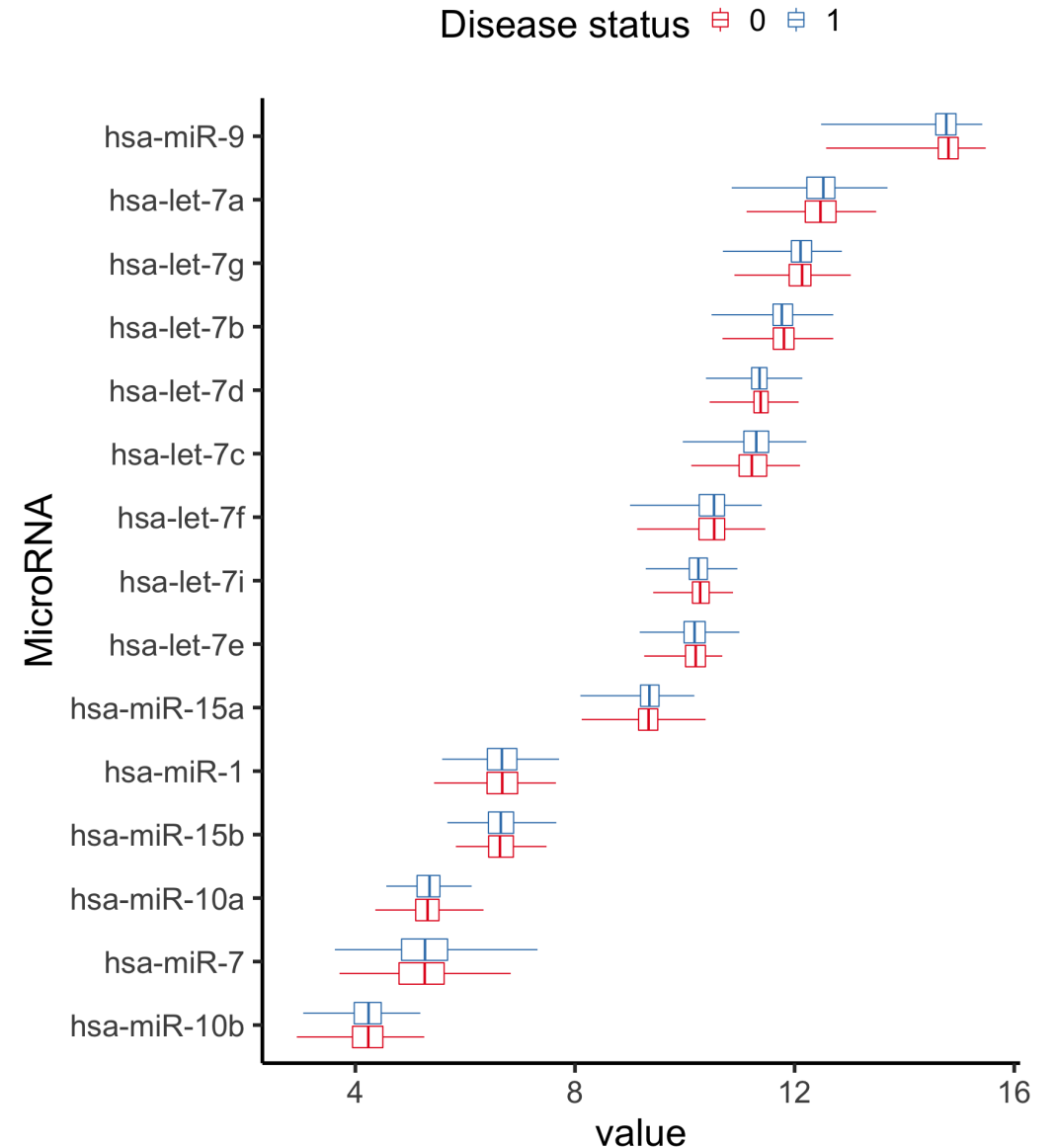
# How many subjects have with Alzheimer's?

```r
mirna %>% select(subject, AD) %>%
  distinct() %>%
  janitor::tabyl(AD) %>%
  janitor::adorn_pct_formatting()
```

```
##  AD   n percent
##   0 273   38.9%
##   1 428   61.1%
```
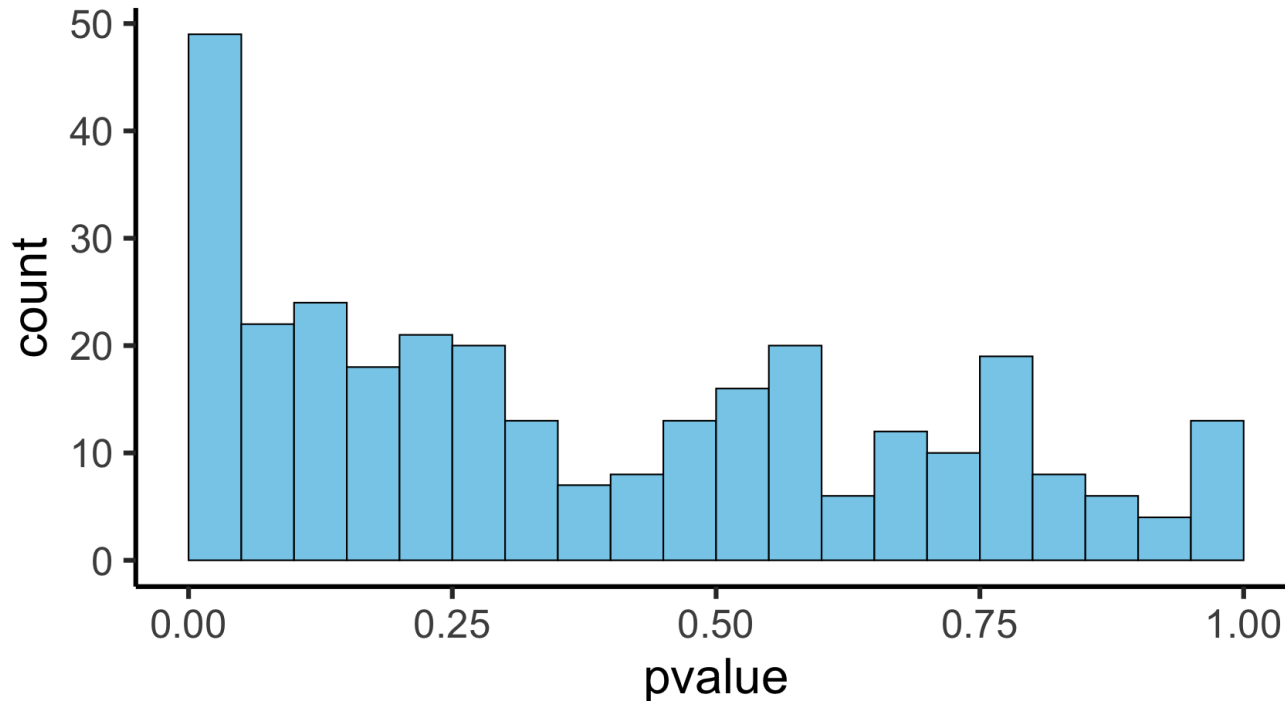
```r
mirna %>%
  group_by(microRNA) %>%
  nest() %>%
  ungroup() %>%
  slice(1:15) %>% # extract first 15 groups
  unnest(cols = everything()) %>%
  ggplot() +
  aes(y = reorder(microRNA, value),
      x = value, colour = factor(AD)) +
  geom_boxplot(coef = 10) +
  scale_color_brewer(palette = "Set1") +
  theme(legend.position = "top") +
  labs(colour = "Disease status",
       y = "MicroRNA")
```

Let's visualise the distribution of p-values for all 309 microRNA.

```
mirna_pval = mirna %>%
  group_by(microRNA) %>%
  summarise(pvalue = t.test(value ~ AD)$p.value)
mirna_pval %>% ggplot() + aes(x = pvalue) +
  geom_histogram(boundary = 0, binwidth = 0.05, fill = "skyb
```



```
mirna_pval %>%
  summarise(
    n_sig = sum(pvalue < 0.05)
  )
```
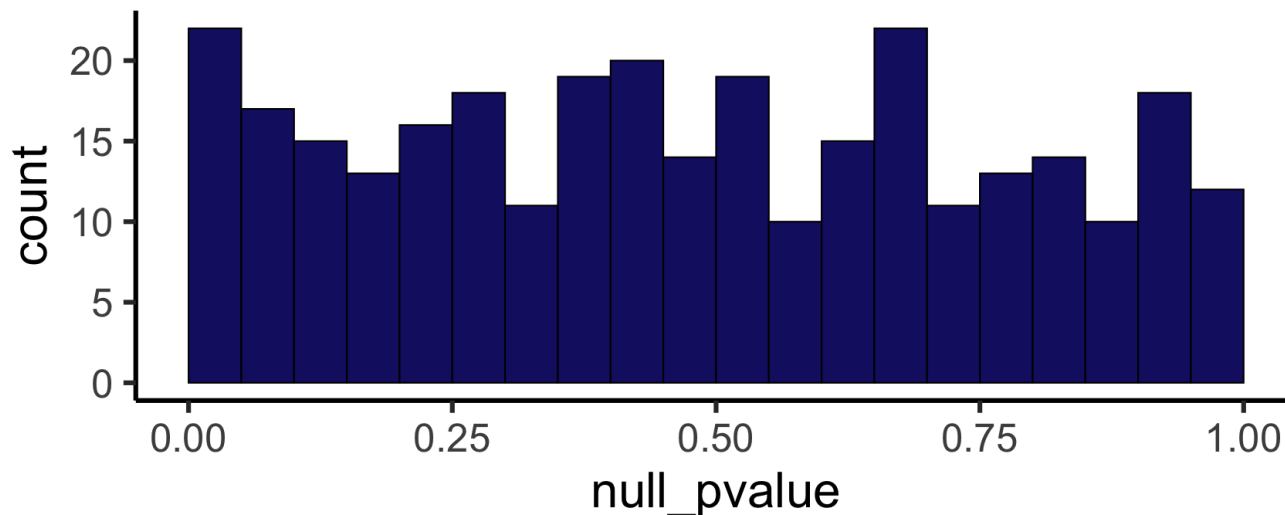
```
## # A tibble: 1 × 1
##   n_sig
##   <int>
## 1    49
```

- Of the 309 microRNA tested, 49 with p-values less than 0.05.

- Are all of these important?

If there was no association between any microRNAs and Alzheimer's disease we would expect our p-values to follow a uniform distribution. We can generate a set of p-values *knowing* that there is no association and visualise this:

```r
set.seed(2021)
mirna_pval = mirna_pval %>%
  mutate(null_pvalue = runif(n = n(), min = 0, max = 1))
mirna_pval %>% ggplot() + aes(x = null_pvalue) +
  geom_histogram(boundary = 0, binwidth = 0.05,
                 fill = "midnightblue", colour = "black")
```



```r
mirna_pval %>%
  summarise(
    n_sig = sum(null_pvalue < 0.05)
  )
```

```
## # A tibble: 1 × 1
##   n_sig
##   <int>
## 1    22
```
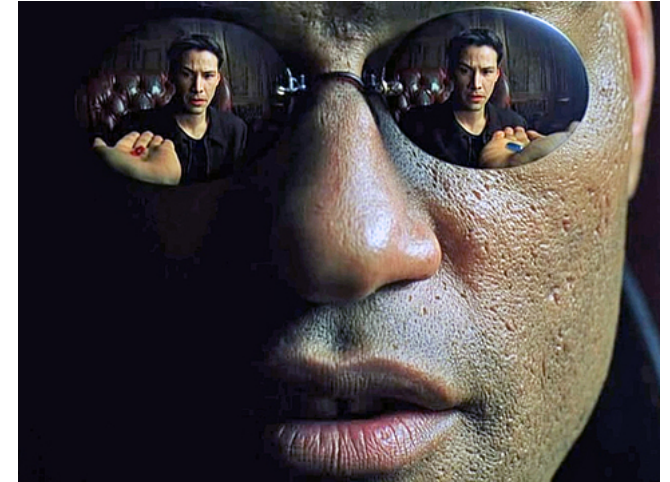
- When we **know** that there are no truly important microRNAs, we still see 22 "significant" p-values in this simulated example.

# Deciding what's real

> ### 💬 Statistical thinking
>
> How do you know if a significant association you find is **real** and not just random chance?
>
> 1. Because someone else published it
>
> 2. Because p-value is less than 0.05
>
> 3. Because of all the tests I ran, that one had the lowest p-value
>
> 4. Because it makes biological sense

# The reality of the situation

- We never really know what is a *real* association.

- A small p-value provides some evidence against the null bit it could still be a false positive.

- Type 1 error ( $\alpha = 0.05$ )

For every model we evaluate at $\alpha = 0.05$, we accept that there is a 5% chance that we reject the null hypothesis when the null hypothesis is actually true.

## 💬 Statistical Thinking

Are jelly beans associated with acne?

## 💬 Statistical Thinking

If we performed 20 tests with $\alpha = 0.05$, how many tests are likely to be significant by chance alone?

## 💬 Statistical Thinking

Does this make you question the conclusions from any of the tests you've done in the labs?

# Types of errors

Suppose you are testing a hypothesis that a parameter $\theta$ equals zero versus the alternative that it does not equal zero.

- **Type I error or false positive** ($V$) Conclude that $\theta$ does not equal zero when it does
- **Type II error or false negative** ($T$) Conclude that $\theta$ equals zero when it doesn't

Possible outcomes from a series of $m$ hypothesis tests.

| Outcomes | Truth: $\theta = 0$ | Truth: $\theta \neq 0$ | Number of tests |
|---|---|---|---|
| **Conclusion:** $\theta = 0$ | $U$ | $T$ | $m - R$ |
| **Conclusion:** $\theta \neq 0$ | $V$ | $S$ | $R$ |
| **Number of tests** | $m_0$ | $m - m_0$ | $m$ |

## Error rates

**False positive rate**: the rate at which null results ($\theta = 0$) are called significant: $\mathrm{E}\left[\dfrac{V}{m_0}\right]$

**Family wise error rate (FWER)**: the probability of at least one false positive $P(V \geq 1)$

# Accounting for multiple testing

- If p-values are correctly calculated calling all p-values less than $\alpha$ significant will control the false positive rate at level $\alpha$, on average.

- Suppose that you perform 10,000 tests and the reality is that $\theta = 0$ for all of them.

> 💬 In what sort of situation might you be doing 10,000 hypothesis tests?

- Suppose that you call all p-values less than 0.05 significant.

- The expected number of false positives is: $10,000 \times 0.05 = 500$ false positives.

- **How do we avoid so many false positives?**

Consider two approaches:

- Controlling the **Family-Wise Error Rate (FWER)**

- Controlling the **False Discovery Rate (FDR)**

# Controlling the family-wise error rate

# Family-wise error rate

**Family wise error rate (FWER)**: the probability of at least one false positive.

Let $T_1, \ldots, T_m$ be $m$ test statistics for null hypothesis $H_{01}, \ldots, H_{0m}$

$$
\begin{aligned}
\text{FWER} &= P(\text{falsely rejecting one or more } H_{0i}) \\
&= P(V \geq 1)
\end{aligned}
$$

If the **null hypothesis is always true** but we conduct $m$ tests each at significance level $\alpha$ then...

- The probability of at least one false positive is $1 - (1 - \alpha)^m$

- e.g. if $m = 20$ then the FWER is 64%.

Possible outcomes from a series of $m$ hypothesis tests.

| Outcomes | Truth: $\theta = 0$ | Truth: $\theta \neq 0$ | Number of tests |
|---|---|---|---|
| **Conclusion:** $\theta = 0$ | $U$ | $T$ | $m - R$ |
| **Conclusion:** $\theta \neq 0$ | $V$ | $S$ | $R$ |
| **Number of tests** | $m_0$ | $m - m_0$ | $m$ |

```
m = 20
alpha = 0.05
1-(1-alpha)^m
```

```
## [1] 0.6415141
```

# Bonferroni correction

The Bonferroni correction is the oldest multiple testing correction.

Given that the number of false positives for $m$ tests is $m\alpha$ then consider defining a new threshold for significance:

$$\alpha^* = \frac{\alpha}{m}$$

- This is conservative but keeps FWER $< \alpha$.

- e.g. for $m = 20$

$$1 - (1 - \alpha^*)^m = 1 - (1 - 0.05/20)^{20} = 0.0488$$

# Bonferroni correction

**Basic idea**

- Suppose you do $m$ tests

- You want to control FWER at level $\alpha$ so $P(V \geq 1) < \alpha$

- Calculate p-values in the usual way

- Set $\alpha^* = \alpha/m$ (or alternatively calculate adjusted p-values: $p^* = \text{p-value} \times m$)

- Call all p-values less than $\alpha^*$ significant (or all adjusted p-values less than $\alpha$ significant)

**Pros**: Easy to calculate, conservative

**Cons**: May be very conservative

💬 What does conservative mean in this context?

# Controlling family-wise error rate (FWER)

**Mathematically**

Let $p_1, p_2, \ldots, p_m$ be the p-values from $m$ hypothesis tests.

$$
\begin{aligned}
\text{FWER} &= P(\text{rejecting at least one true null hypothesis}) \\
&= P(V \geq 1) \\
&= P\left\{ \bigcup_{i=1}^{m_0} \left( p_i \leq \frac{\alpha}{m} \right) \right\} \\
&\leq \sum_{i=1}^{m_0} \left\{ P\left( p_i \leq \frac{\alpha}{m} \right) \right\} \\
&= m_0 \frac{\alpha}{m} \\
&\leq m \frac{\alpha}{m} \\
&= \alpha.
\end{aligned}
$$

This is not examinable.

# 10 microRNA p-values: Bonferroni method

For the sake of illustration, we're going to control the error rates at $\alpha = 0.2$. Let's also say that we are only interested in $m = 10$ of the microRNA.

```
alpha = 0.2
m = 10
M = nrow(mirna_pval)
set.seed(123, sample.kind = "Rounding")
sample_rows = sample(1:M, size = m)
mirna10 = mirna_pval %>%
  select(microRNA, pvalue) %>%
  slice(sample_rows) %>%
  mutate(p_bonferroni = pmin(pvalue*m, 1)) %>%
  arrange(pvalue)
```

We compare the original p-values to $0.2/m = 0.02$ or the adjusted p-values p_bonferroni to $\alpha = 0.2$. We get the same conclusion either way.

```
mirna10 %>% knitr::kable(digits = 4) %>%
  kableExtra::kable_styling(font_size = 16)
```

| microRNA | pvalue | p_bonferroni |
|----------|--------|--------------|
| hsa-miR-874 | 0.0022 | 0.0220 |
| hsa-miR-185 | 0.0054 | 0.0544 |
| hsa-miR-548a-3p | 0.0532 | 0.5323 |
| hsa-miR-221 | 0.1730 | 1.0000 |
| hsa-let-7e | 0.2594 | 1.0000 |
| hsa-miR-345 | 0.3666 | 1.0000 |
| hsa-miR-27b | 0.3766 | 1.0000 |
| hsa-miR-337-3p | 0.5807 | 1.0000 |
| hsa-miR-639 | 0.6607 | 1.0000 |
| hsa-miR-640 | 0.6912 | 1.0000 |

We would typically choose the error rate to be 0.05 or perhaps even smaller, but using a larger error rate of 0.2 makes it a bit easier to outline the concepts.

# Controlling the false discovery rate

# False Discovery Rate (FDR)

**Aim:** to keep the expected proportion of false positives in your rejected tests (FDR) close to $\alpha$

Let...

- $R$ = total number of $H_{0i}$ rejected
- $V$ = number of $H_{0i}$ falsely rejected
- FDR = $\mathrm{E}\left(\dfrac{V}{R}\right)$

Possible outcomes from a series of $m$ hypothesis tests.

| Outcomes | Truth: $\theta = 0$ | Truth: $\theta \neq 0$ | Number of tests |
|---|---|---|---|
| **Conclusion:** $\theta = 0$ | $U$ | $T$ | $m - R$ |
| **Conclusion:** $\theta \neq 0$ | $V$ | $S$ | $R$ |
| **Number of tests** | $m_0$ | $m - m_0$ | $m$ |

# Controlling false discovery rate (FDR)

The Benjamini–Hochberg procedure is the most popular correction when performing *lots* of tests say in genomics, imaging, astronomy, or other signal-processing disciplines.

**Basic idea**:

- Suppose you do $m$ tests

- You want to control FDR at level $\alpha$

- Calculate p-values normally

- Order the p-values from smallest to largest $p_{(1)} \leq p_{(2)} \leq \cdots \leq p_{(m)}$

- Find $j^* = \max j$ such that $p_{(j)} \leq \dfrac{j}{m}\alpha$

- Reject all $H_{0i}$ where $p_{(i)} \leq \dfrac{j^*}{m}\alpha$

**Pros**: Still pretty easy to calculate, less conservative (maybe much less)

**Cons**: Allows for more false positives, may behave strangely under dependence

# 10 microRNA p-values: BH method

Controlling the FDR rates at $\alpha = 0.2$. Take a the same sample of 10 p-values from the microRNA experiment.

```r
alpha = 0.2
m = 10
p_vals = sort(mirna10$pvalue)
# BH procedure
# j=1: smallest p-value < 1*alpha/m?
p_vals[1] < 1*alpha/m
```

```
## [1] TRUE
```

```r
# j=2: second smallest p-value < 2*alpha/m?
p_vals[2] < 2*alpha/m
```

```
## [1] TRUE
```

```r
# j=3: third smallest p-value < 3*alpha/m?
p_vals[3] < 3*alpha/m
```

```
## [1] TRUE
```

```r
# j=4: fourth smallest p-value < 4*alpha/m?
p_vals[4] < 4*alpha/m
```

```
## [1] FALSE
```

```r
# j=5: fifth smallest p-value < 5*alpha/m?
p_vals[5] < 5*alpha/m
```

```
## [1] FALSE
```

```r
# and so on ...
```

# 10 microRNA p-values: BH method

```r
# in general...
result = vector(length = length(p_vals))
p_vals = sort(p_vals) # we already did this but just emphasising it
for(j in seq(p_vals)) { # seq(p_vals) is the same as 1:length(pvals)
  result[j] = p_vals[j] < j*alpha/m
}
result
```

```
##  [1]  TRUE  TRUE  TRUE FALSE FALSE FALSE FALSE FALSE FALSE
## [10] FALSE
```

```r
largest_true = max(which(result == TRUE))
largest_true
```

```
## [1] 3
```

```r
significant_pvals = p_vals[1:largest_true]
significant_pvals
```
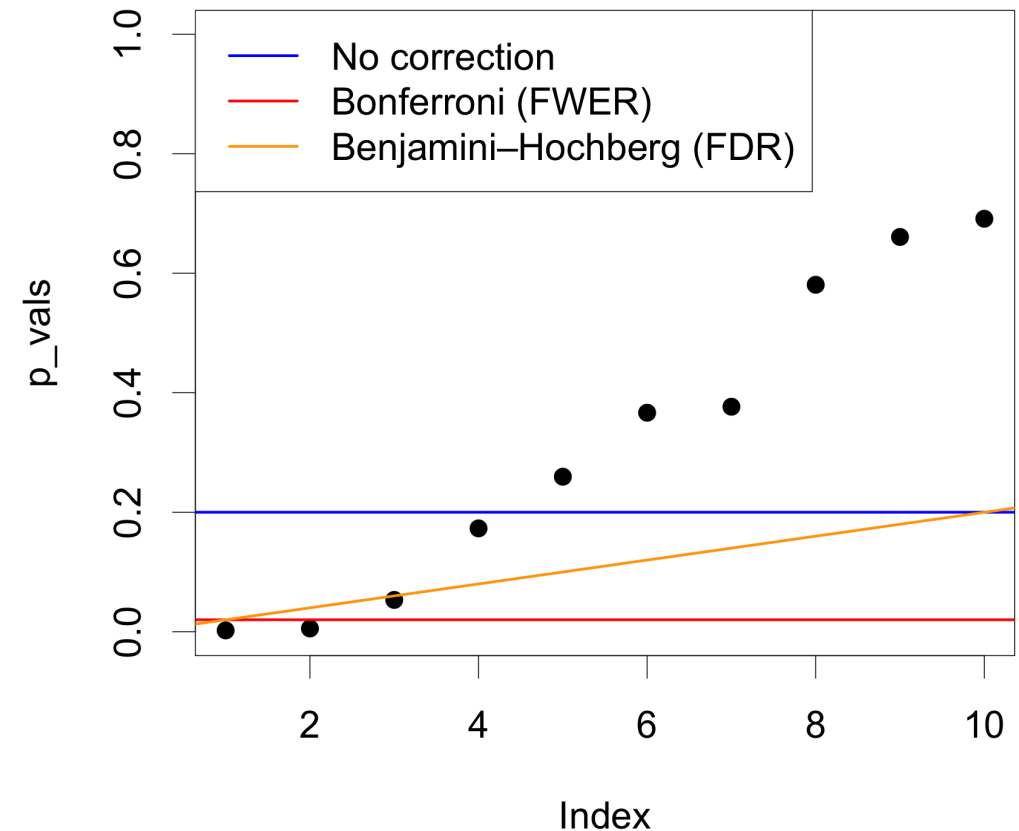
```
## [1] 0.002196115 0.005442173 0.053227665
```

# 10 microRNA p-values: BH vs Bonferonni

Controlling all error rates at $\alpha = 0.20$ and using our sample of 10 microRNAs.

```
alpha = 0.2
m = 10
par(cex = 2.6, mar = c(4,4,1,1))
plot(p_vals,
     ylim = c(0,1), pch=19)
abline(h = alpha, col = "blue", lwd = 3)
abline(h = alpha/m, col = "red", lwd = 3)
abline(a = 0, b=alpha*1/m, col = "orange",
       lwd = 3)
legend("topleft",
       legend = c("No correction",
                  "Bonferroni (FWER)",
                  "Benjamini-Hochberg (FDR)"),
       lty = 1, lwd = 3,
       col = c("blue","red","orange"))
```

The lines are the significance thresholds for the three methods. If a point is below the line, the method would consider it "significant".
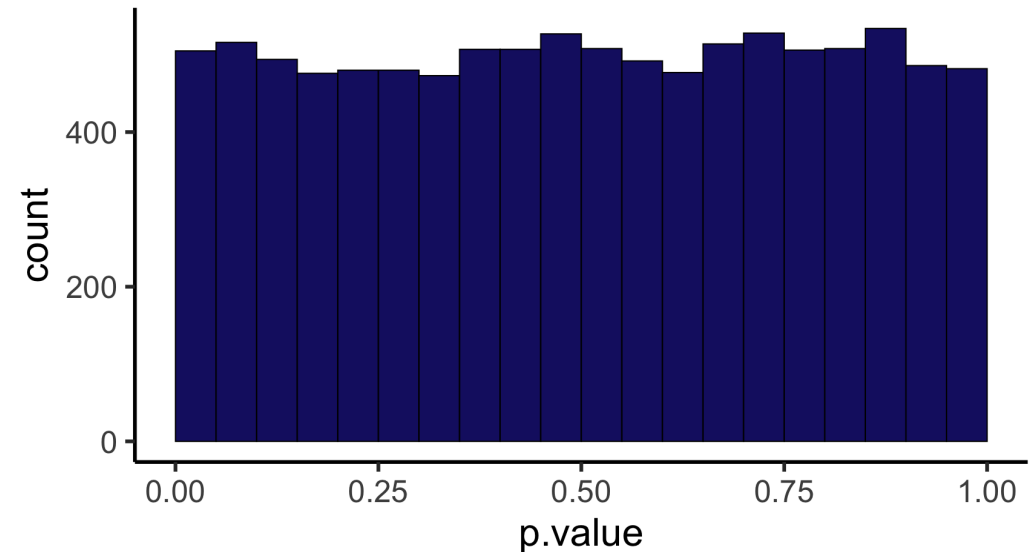
# Simulation experiments

# Case study I: no true positives

```r
set.seed(1234)
p_vals = rep(NA, 1000)
B = 10000
case1 = tibble(experiment = 1:B) %>%
  group_by(experiment) %>%
  summarise(x_sample = rnorm(20),
            y_sample = rnorm(20)) %>%
  nest() %>%
  mutate(
    test = map(data,
               ~t.test(.$x_sample,
                       .$y_sample,
                       var.equal = TRUE) %>%
                 broom::tidy())) %>%
  unnest(test) %>%
  ungroup()
mean(case1$p.value < 0.05)
```

```
## [1] 0.0505
```

```r
case1 %>%
  ggplot() + aes(x = p.value) +
  geom_histogram(boundary = 0,
                 binwidth = 0.05,
                 fill = "midnightblue",
                 colour = "black")
```

# Case study I: no true positives

Get R to do the corrections for us using the `p.adjust()` function:

```
case1 = case1 %>%
  mutate(
    p_bonf = p.adjust(p.value, "bonferroni"),
    p_bh = p.adjust(p.value, "BH")
  )
case1 %>% select(experiment, p.value,
                 p_bonf, p_bh) %>%
  head()
```

```
## # A tibble: 6 × 4
##   experiment p.value p_bonf  p_bh
##        <int>   <dbl>  <dbl> <dbl>
## 1          1   0.703      1 0.995
## 2          2  0.0319      1 0.962
## 3          3  0.0539      1 0.962
## 4          4   0.778      1 0.995
## 5          5   0.480      1 0.995
## 6          6   0.990      1 0.999
```

Proportion of "significant" results

```
case1 %>% ungroup() %>%
  summarise(
    original = mean(p.value < 0.05),
    bonferroni = mean(p_bonf < 0.05),
    bh = mean(p_bh < 0.05)
  )
```

```
## # A tibble: 1 × 3
##   original bonferroni    bh
##      <dbl>      <dbl> <dbl>
## 1   0.0505          0     0
```

# Case study II: 50% true positives

```r
set.seed(1234)
B = 10000
case2 = tibble(experiment = 1:B) %>%
  group_by(experiment) %>%
  summarise(x_sample = rnorm(20),
            y_sample = rnorm(20)) %>%
  rowwise() %>%
  mutate(truth = if_else(experiment<=B/2, "mu1 - mu2 = 0", "mu1 - mu2 = 2"),
         y_sample = if_else(truth == "mu1 - mu2 = 2", y_sample + 2, y_sample)) %>%
  ungroup() %>%
  nest(data = c(x_sample, y_sample)) %>%
  mutate(test = map(data, ~t.test(.$x_sample, .$y_sample, var.equal = TRUE) %>% broom::tidy())) %>%
  unnest(test) %>%
  ungroup() %>%
  mutate(
    prediction = if_else(p.value < 0.05, "reject H0", "don't reject H0"),
    p_bonf = p.adjust(p.value, method = "bonferroni"),
    p_bh = p.adjust(p.value, method = "BH"),
    pred_bonf = if_else(p_bonf < 0.05, "reject H0", "don't reject H0"),
    pred_bh = if_else(p_bh < 0.05, "reject H0", "don't reject H0")
  )
```

# Case study II: 50% true positives

```
# no adjustment
case2 %>% janitor::tabyl(prediction, truth)
```

```
##        prediction mu1 - mu2 = 0 mu1 - mu2 = 2
##   don't reject H0          4742             0
##        reject H0           258          5000
```

```
# Bonferroni: controls FWER
case2 %>% janitor::tabyl(pred_bonf, truth)
```

```
##        pred_bonf mu1 - mu2 = 0 mu1 - mu2 = 2
##   don't reject H0          5000           979
##        reject H0             0          4021
```

```
# BH: controls FDR
case2 %>% janitor::tabyl(pred_bh, truth)
```

```
##          pred_bh mu1 - mu2 = 0 mu1 - mu2 = 2
##   don't reject H0          4883             0
##        reject H0           117          5000
```

| Outcomes | Truth: $\theta = 0$ | Truth: $\theta \neq 0$ | Number of tests |
|---|---|---|---|
| **Conclusion:** $\theta = 0$ | $U$ | $T$ | $m - R$ |
| **Conclusion:** $\theta \neq 0$ | $V$ | $S$ | $R$ |
| **Number of tests** | $m_0$ | $m - m_0$ | $m$ |

- Bonferroni is controlling FWER

$$\mathrm{FWER} = P(\text{falsely rejecting one or more } H_{0i})$$
$$= P(V \geq 1)$$

- BH is controlling FDR = $E\left(\dfrac{V}{R}\right)$

# Case study II: 50% true positives

```r
pval2 = case2 %>% select(experiment, p.value, p_bonf, p_bh) %>%
  pivot_longer(cols = c(p.value, p_bonf, p_bh), names_to = "method", values_to = "p_value") %>%
  mutate(method = recode(method, "p.value" = "Original", "p_bh" = "BH", "p_bonf" = "Bonferroni"))
pval2 %>% ggplot() + aes(x = p_value, fill = method) +
  geom_histogram(boundary = 0, binwidth = 0.05, colour = "black") +
  facet_grid(~method) + scale_fill_brewer(palette = "Set1") + scale_x_continuous(breaks = c(0,1)) +
  theme(legend.position = "none")
```

# MicroRNA revisited

# All microRNA p-values

```r
mirna_pval = mirna_pval %>%
  mutate(
    p_bonf = p.adjust(pvalue, method = "bonferroni"),
    p_bh = p.adjust(pvalue, method = "BH")
  )
mirna_pval %>%
  summarise(original_n_sig = sum(pvalue < 0.05),
            bonf_n_sig = sum(p_bonf < 0.05),
            bh_n_sig = sum(p_bh < 0.05))
```

```
## # A tibble: 1 × 3
##   original_n_sig bonf_n_sig bh_n_sig
##            <int>      <int>    <int>
## 1             49          4        7
```

# All microRNA p-values

```
mirna_pval %>%  arrange(pvalue) %>%
  select(-null_pvalue, `Original p-value` = pvalue, `Bonferroni p-value` = p_bonf,  `BH p-value` = p_
  DT::datatable(rownames = FALSE, options = list(dom = 'tp', pageLength = 7)) %>%
  DT::formatSignif(2:4, digits = 2)
```

| microRNA ⇕ | Original p-value ⇕ | Bonferroni p-value ⇕ | BH p-value ⇕ |
|---|---|---|---|
| hsa-miR-132 | 1.3e-12 | 4.0e-10 | 4.0e-10 |
| hsa-miR-129-5p | 3.3e-7 | 0.00010 | 0.000051 |
| hsa-miR-1260 | 0.000063 | 0.020 | 0.0065 |
| hsa-miR-200a | 0.00013 | 0.041 | 0.010 |
| hsa-miR-34c-5p | 0.00065 | 0.20 | 0.040 |
| hsa-miR-744 | 0.00087 | 0.27 | 0.040 |
| hsa-miR-129-3p | 0.00091 | 0.28 | 0.040 |

Previous   1   2   3   4   5   …   45   Next

# Final comments

- Multiple testing is an entire subfield of statistics

- A basic Bonferroni/BH correction is usually enough

- If there is strong dependence between tests there may be problems

# Further reading

- Statistical significance for genome-wide studies

- Introduction to multiple testing

# References

Patrick, E., S. Rajagopal, H. A. Wong, C. McCabe, J. Xu, A. Tang, S. H. Imboywa, J. A. Schneider, N. Pochet, A. M. Krichevsky, L. B. Chibnik, D. A. Bennett, and P. L. De Jager (2017). "Dissecting the role of non-coding RNAs in the accumulation of amyloid and tau neuropathologies in Alzheimer's disease". In: *Molecular Neurodegeneration* 12.51, pp. 1-13. DOI: 10.1186/s13024-017-0191-y.