

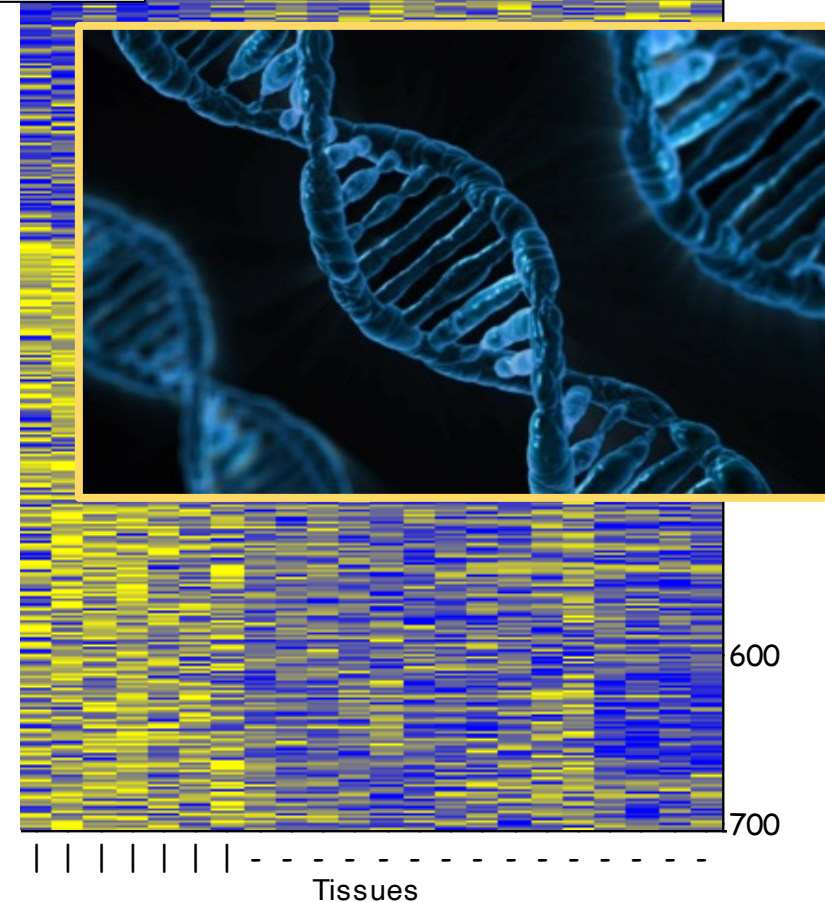
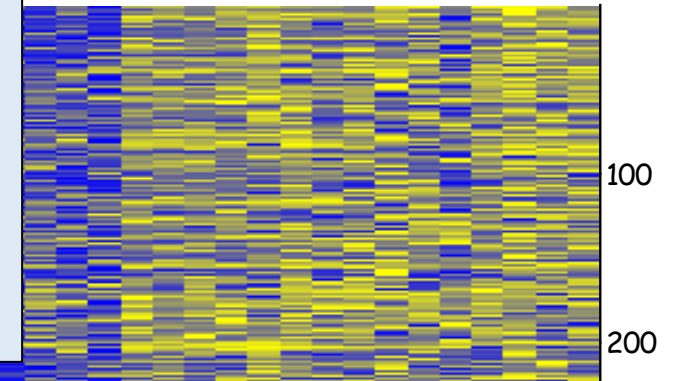
Introduction to differential gene expression, multiple testing and FDR

Zohar Yakhini, Leon Anavy,
Ben Galili – Reichman.

Partially based on slides from Zohar Yakhini,
Doron Lipson, Itai Sharon and others,
at the Technion



Breast Cancer BRCA1/BRCA2 data



The concept of p-value

Under a given NULL MODEL, what is the probability of observing a value for a MEASURED QUANTITY which is as or more extreme than the one actually measured in data?

We tossed a coin 100 times and observed 27 Hs.

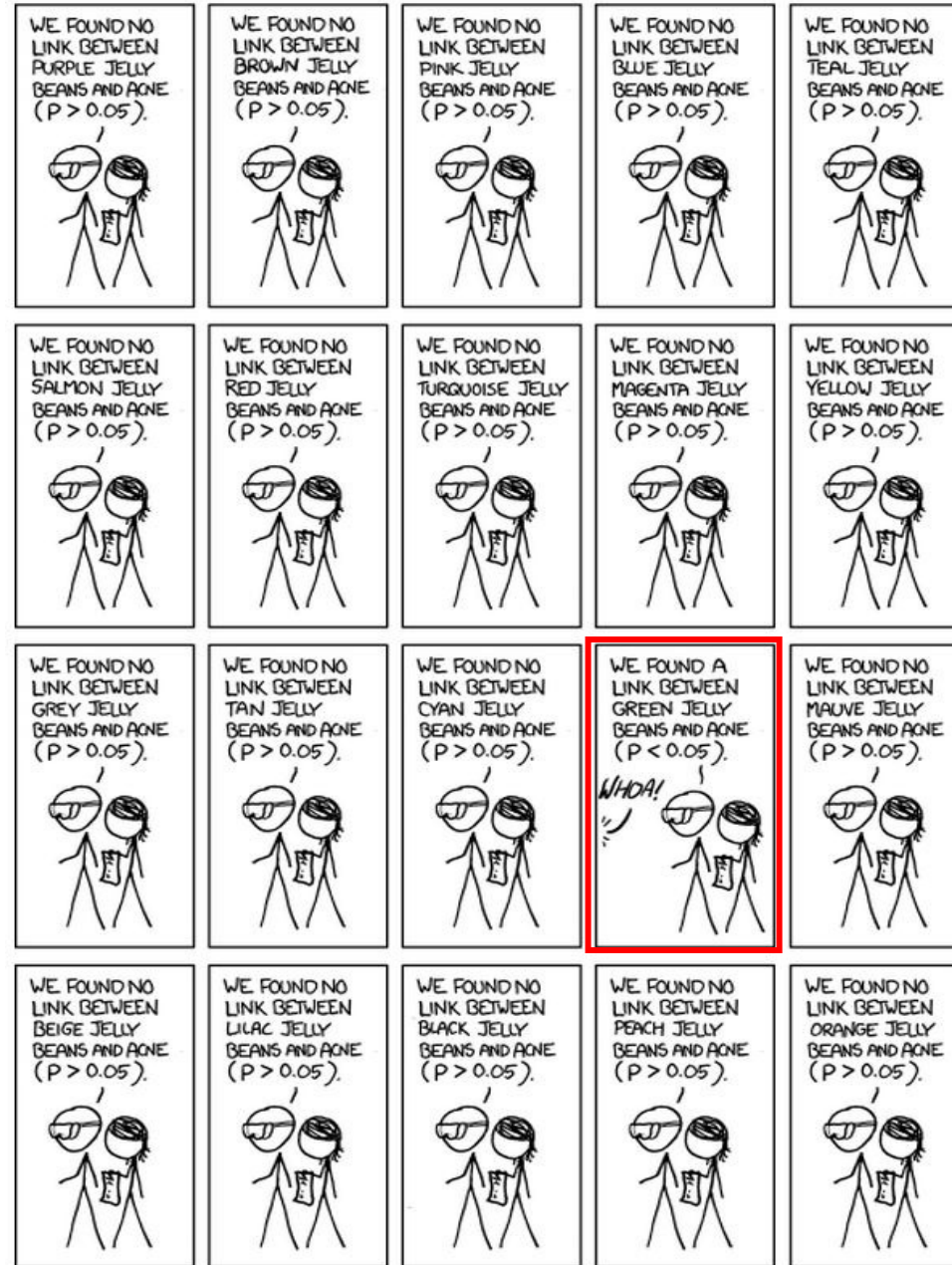
Under a FAIR INDEPENDENT COIN TOSSING MODEL, what is the probability of observing a NUMBER OF Hs which is less than 27 out of 100.

Tall people

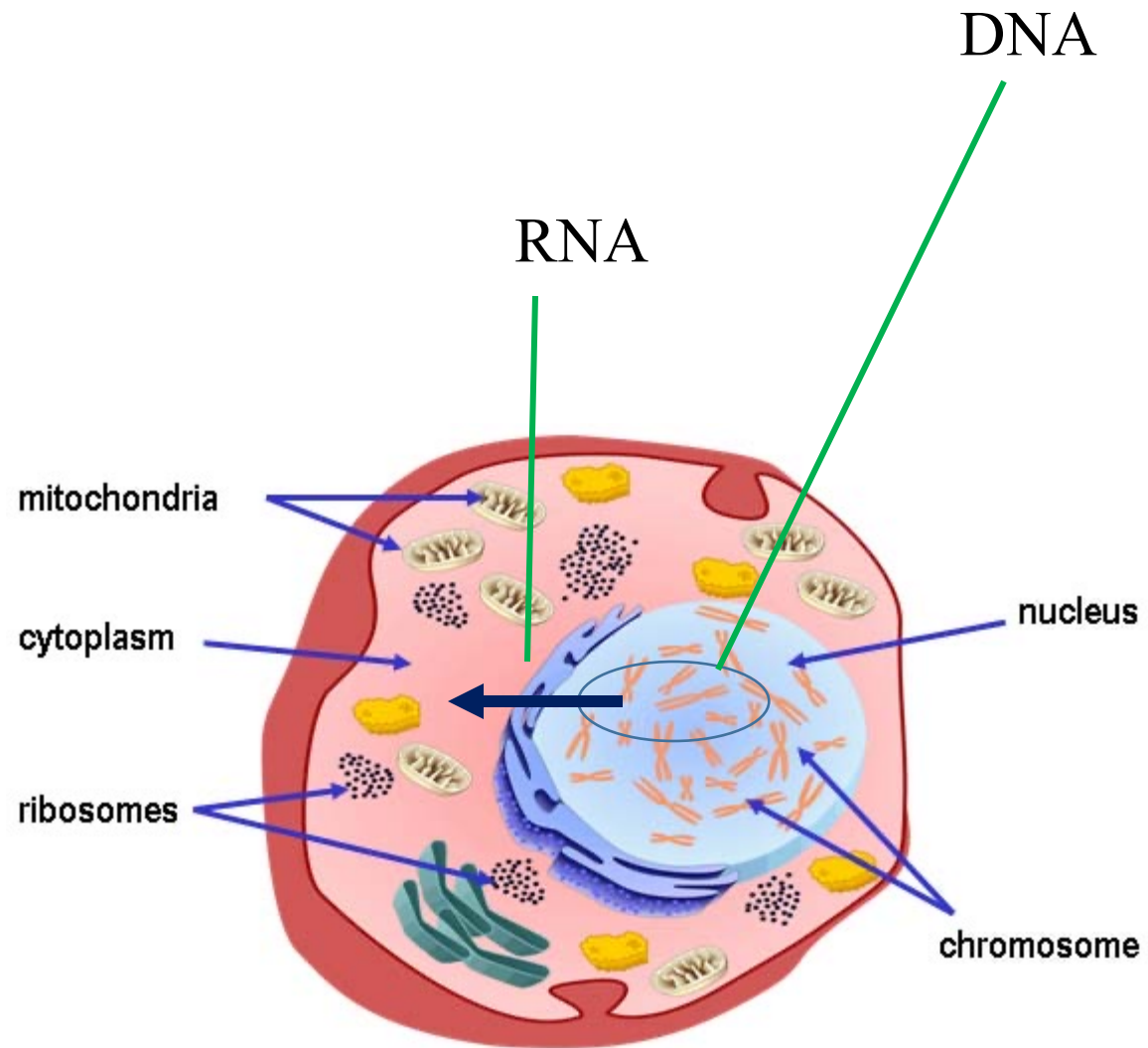
- What is the probability that the person sitting next to you on the bus going back home is $>1.90\text{m}$ tall?
- What is the probability of SOME person in the bus being $>1.90\text{m}$ tall?
- What is the probability that someone on the IDC campus is $>1.90\text{m}$ tall?

$$P(X > 1.9) = 0.01$$

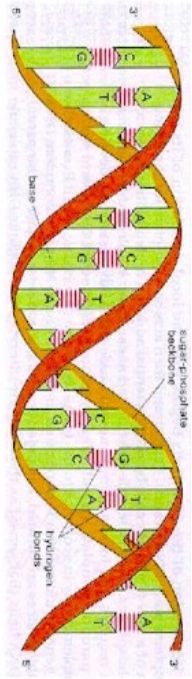
Multiple testing



A living cell

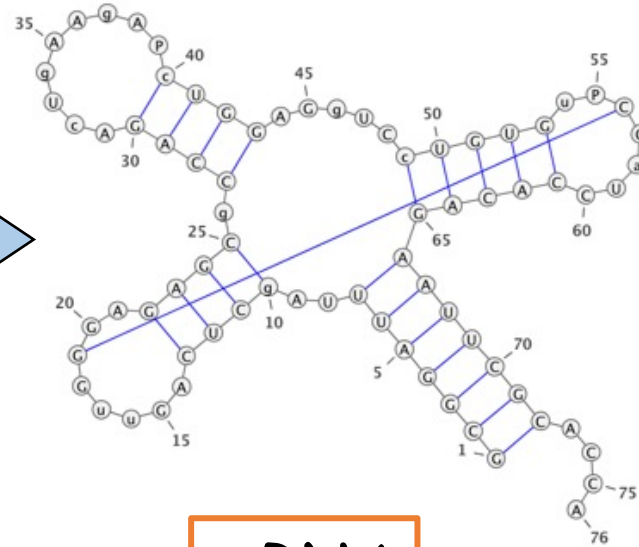


Central Dogma, information flow in living cells



Transcription

Gene
(DNA, nuclear)



mRNA

Translation



Protein

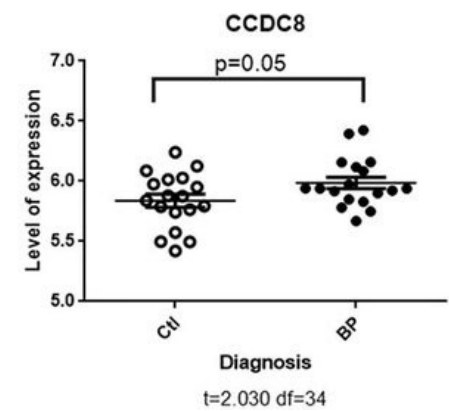
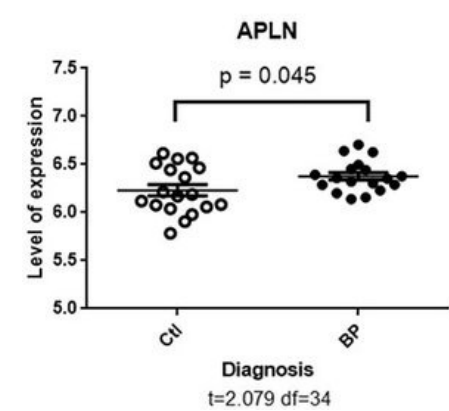
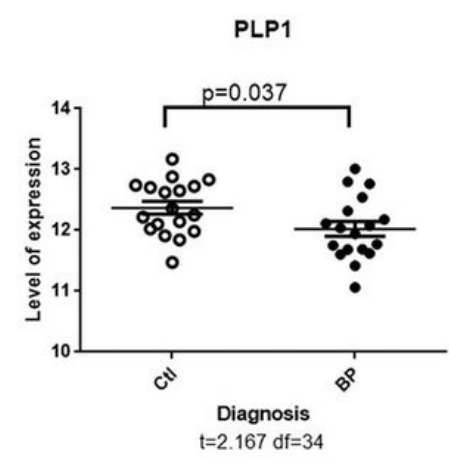
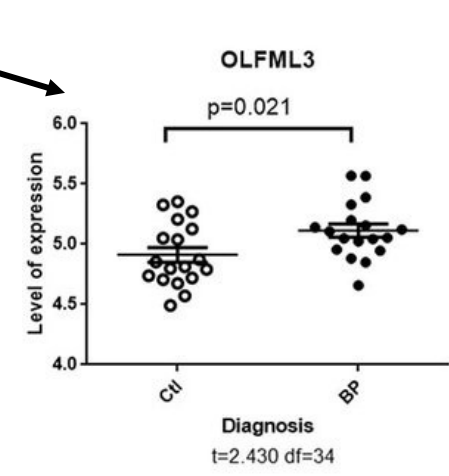
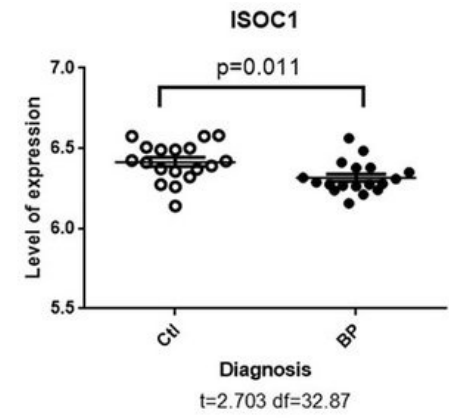
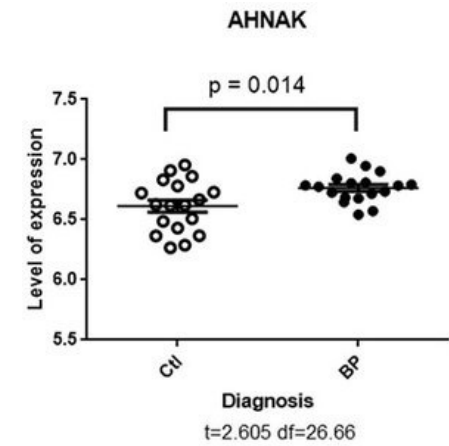
Cells express different subsets of the (organismal) genes in different tissues and under different conditions

Differentially expressed genes

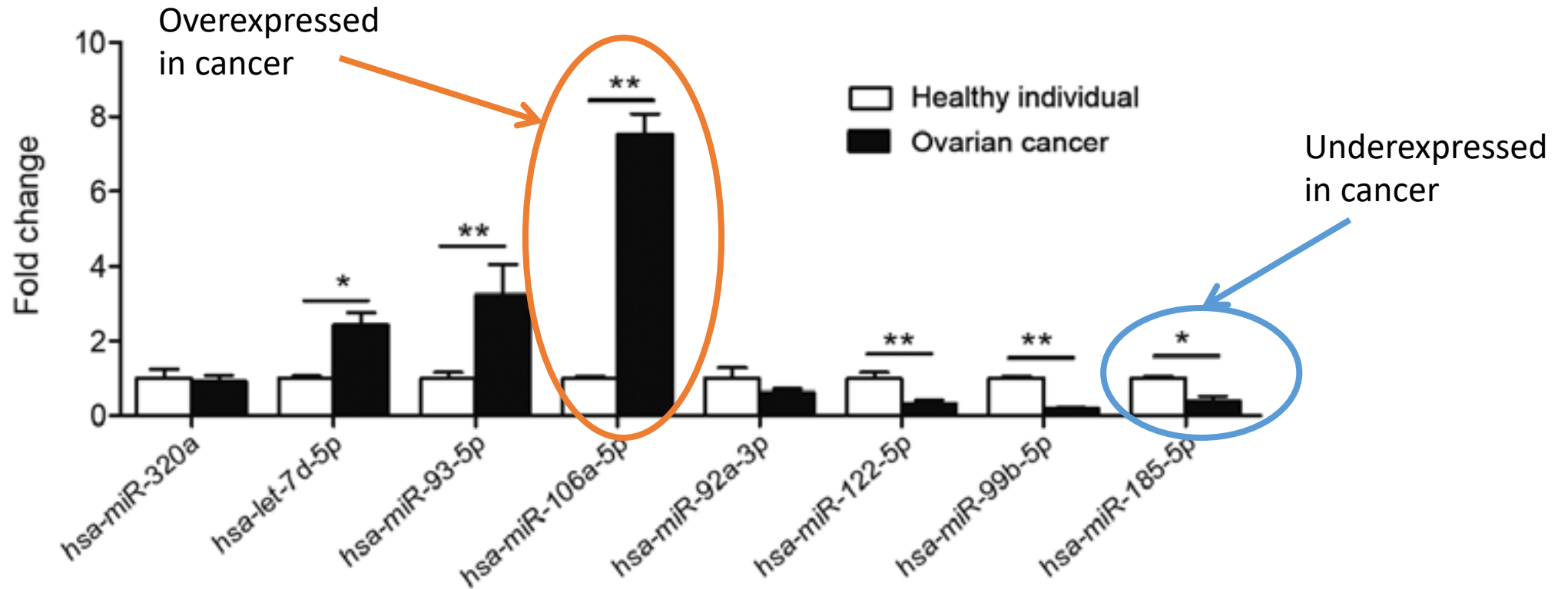
Postmortem measurement of gene expression in brain cells, Bi-Polar (BP) vs Control (Ctl)

PCR measurement

From:
Kidnapillai et al,
The use of a gene expression signature and connectivity map to repurpose drugs for bipolar disorder,
Ken Walder's Lab,
Biological Psychiatry 2018



Differentially expressed genes



Expression levels of plasma exosomal miRNA measured in patients with ovarian cancer (n=30) and healthy women (n=30) using RT-qPCR.

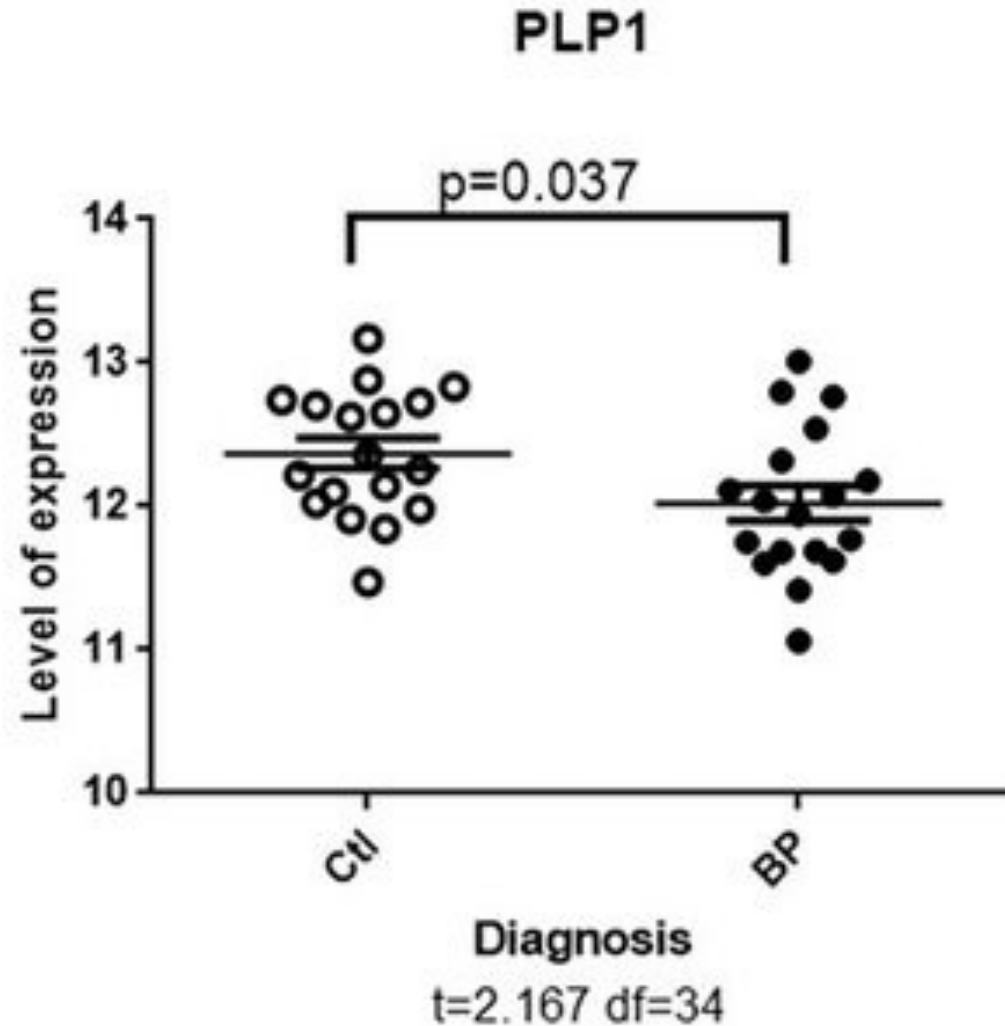
*P<0.05 and **P<0.01.

Zhang et al,
Oncology Letters 2019

Zohar Yakhini

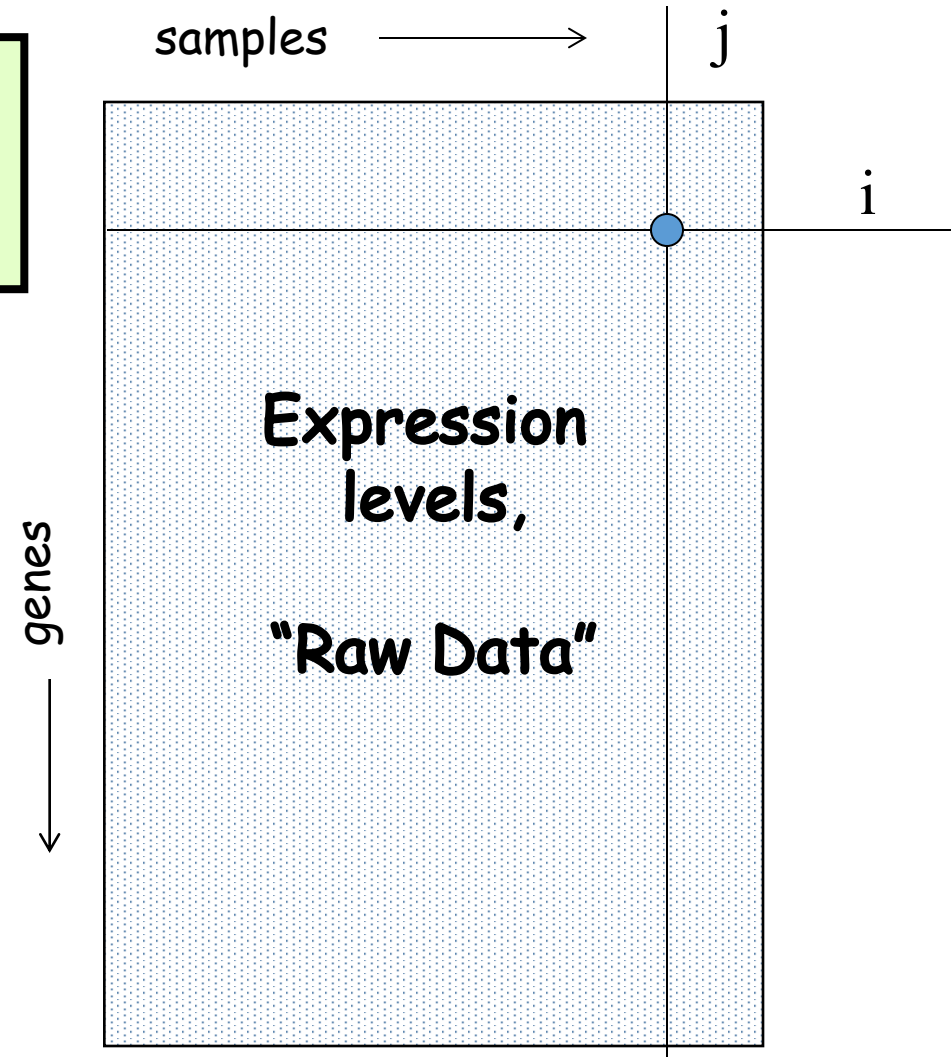
Differentially expressed genes

- t-Test
- WRS
- Other methods...



Gene Expression Data, Matrix Representation

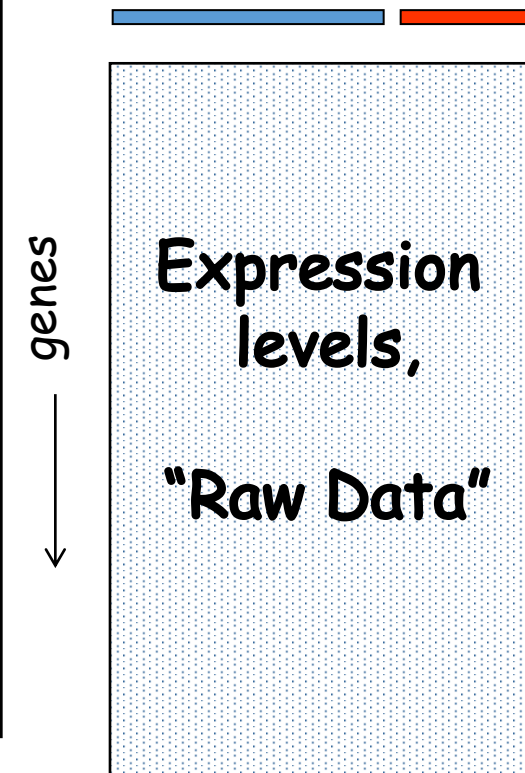
$E(i, j)$ represents the expression level of the gene indexed by i in the experiment/sample indexed by j .



Classified Gene Expression Data, Informative Genes

- Examples:
 - Tumor vs Normal
 - Subtypes of a pathology
 - Prognosis (responders vs non-responders)
 - etc
- Informative genes: genes that sharply separate the two classes. They are ***DIFFERENTIALLY EXPRESSED***
- How do we assign DE scores to genes?

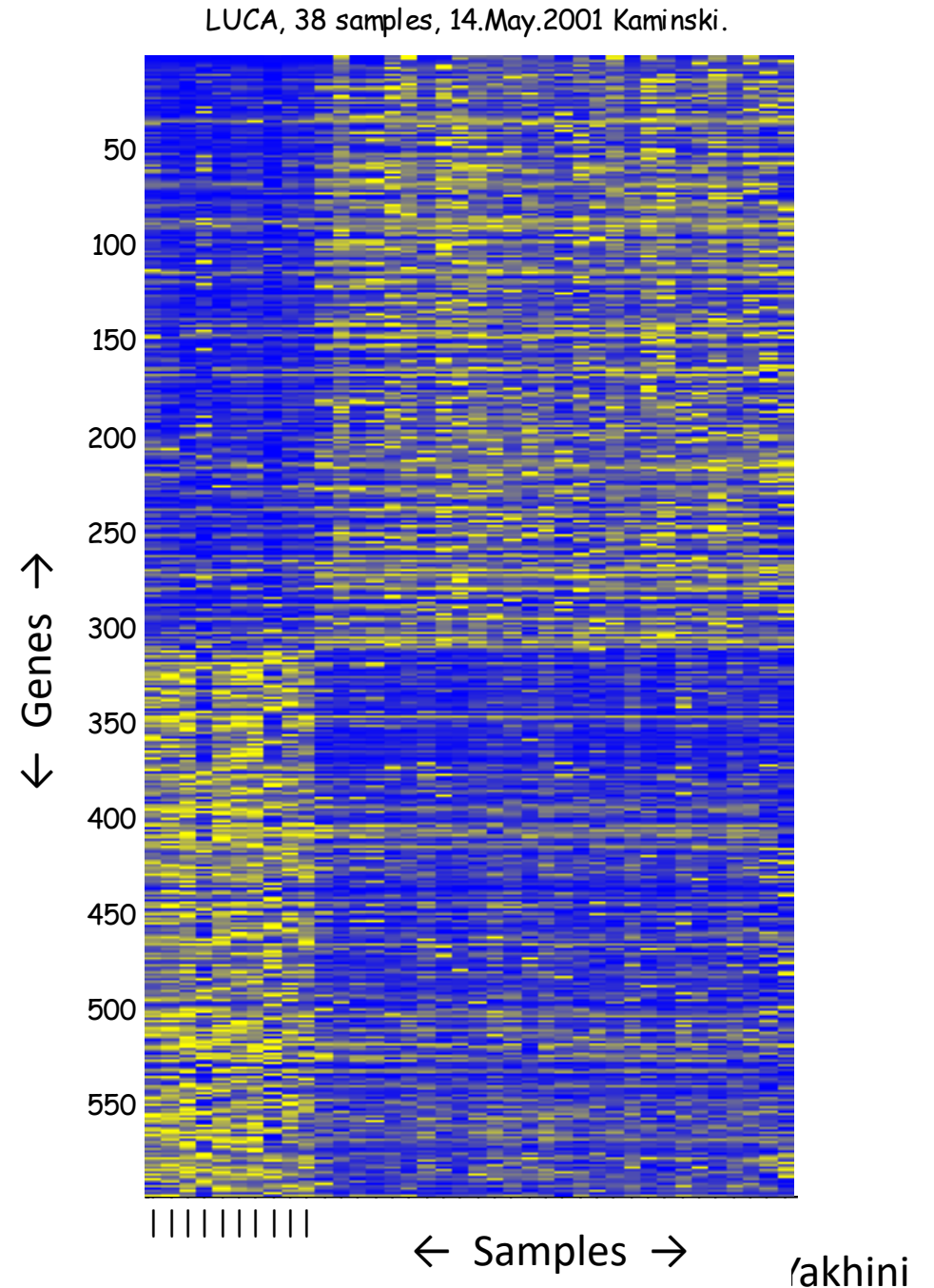
Samples, coming from several classes
→



Lung Cancer Informative Genes

Dehan et al, (Kaminski Lab at Sheba),
Lung Cancer 2007

- 24 tumors (various types and origins)
- 10 normals (normal edges and normal lung pools)



Wilcoxon Rank Sum test

- Compute the sum of the ranks of the +s:

+	+	-	-	+	+	+	-	-	-	-	-	+	-	-
a1	a2	a3	a4	a5	a6	a7	a8	a9	a10	a11	a12	a13	a14	a15

- In this case:

$$RS(+) = 1+2+5+6+7+13 = 34$$

- The null model: all +/- configurations are equiprobable
- The mean value under the null model is $E(RS(+)) = ?$
- The p-value of the deviation can be computed using normal approximations

Wilcoxon example

+ + + + + - - - - - - - - - -

$RS(+)$ = ?, p-val = ?

+ + + - - - - - + - - - + - -

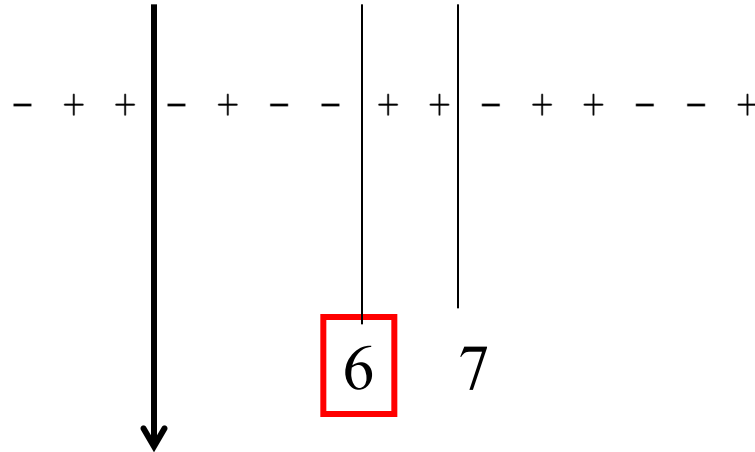
$RS(+)$ = ? , p-val = ?

$E(RS(+))$ = ?

Threshold Error Rate (TNoM) Score

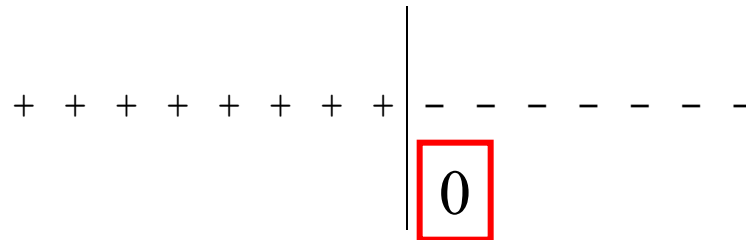
Find the threshold that best separates tumors from normals, count the number of errors committed there.

Ex 1:



of errors = $\min(7, 8) = 7$.

Ex 2: A perfect single gene classifier gets a score of 0.



Threshold Error Rate (TNoM) Score

Consider the ordered vectors of labels below:

$v1 = 0\ 0\ 0\ 0\ 1\ 0\ 0\ 0\ 1\ 1\ 1\ 0$

$v2 = 0\ 0\ 1\ 0\ 0\ 1\ 0\ 0\ 0\ 1\ 1\ 0$

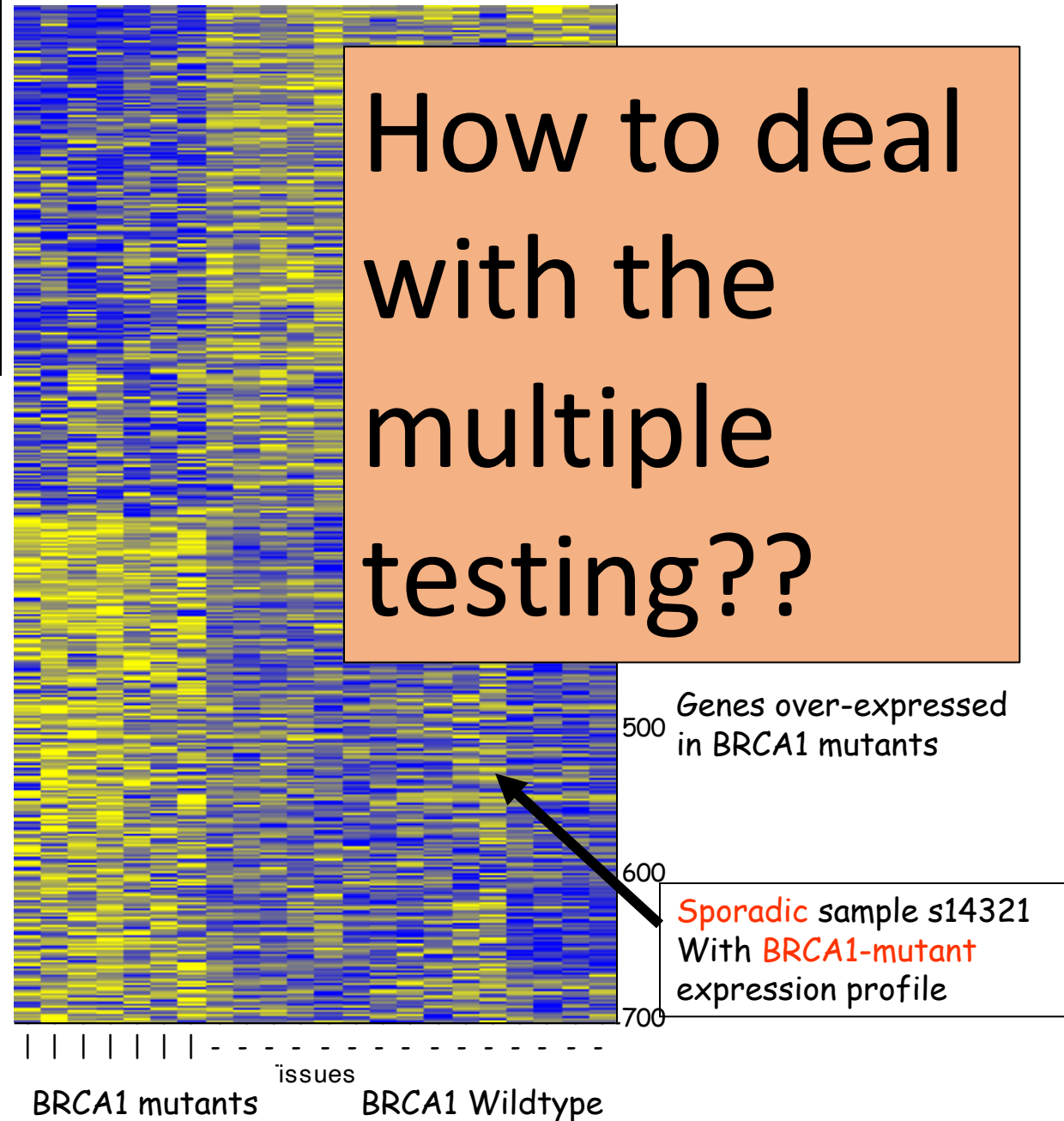
Compute the TNoM score for these.

BRCA1 Differential Expression

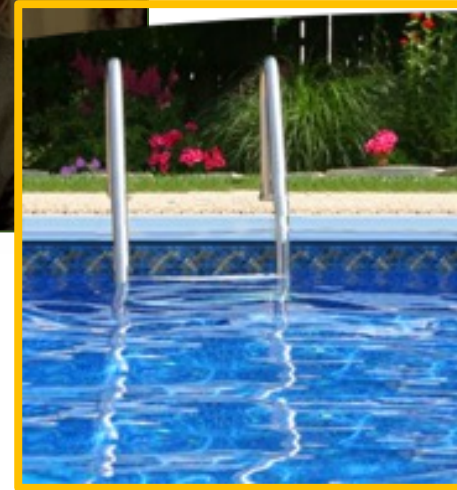
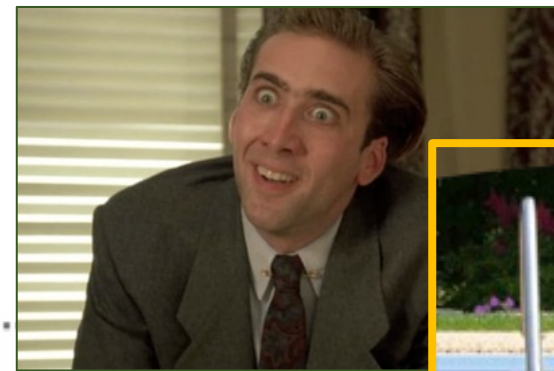
Collab with NIH;
Hedenfalk et al,
NEJM 2001



Breast Cancer BRCA1/BRCA2 data



Spurious Correlations



Number of people who drowned by falling into a pool
correlates with
Films Nicolas Cage appeared in

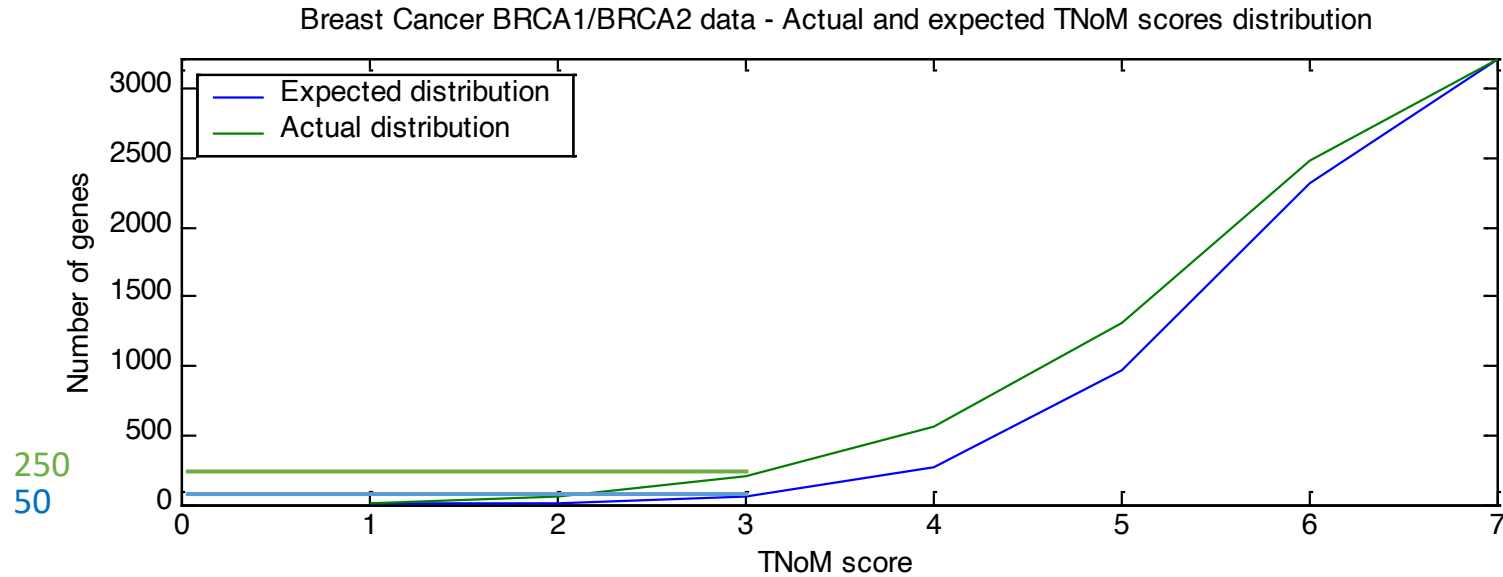
Correlation: 66.6% ($r=0.666004$)



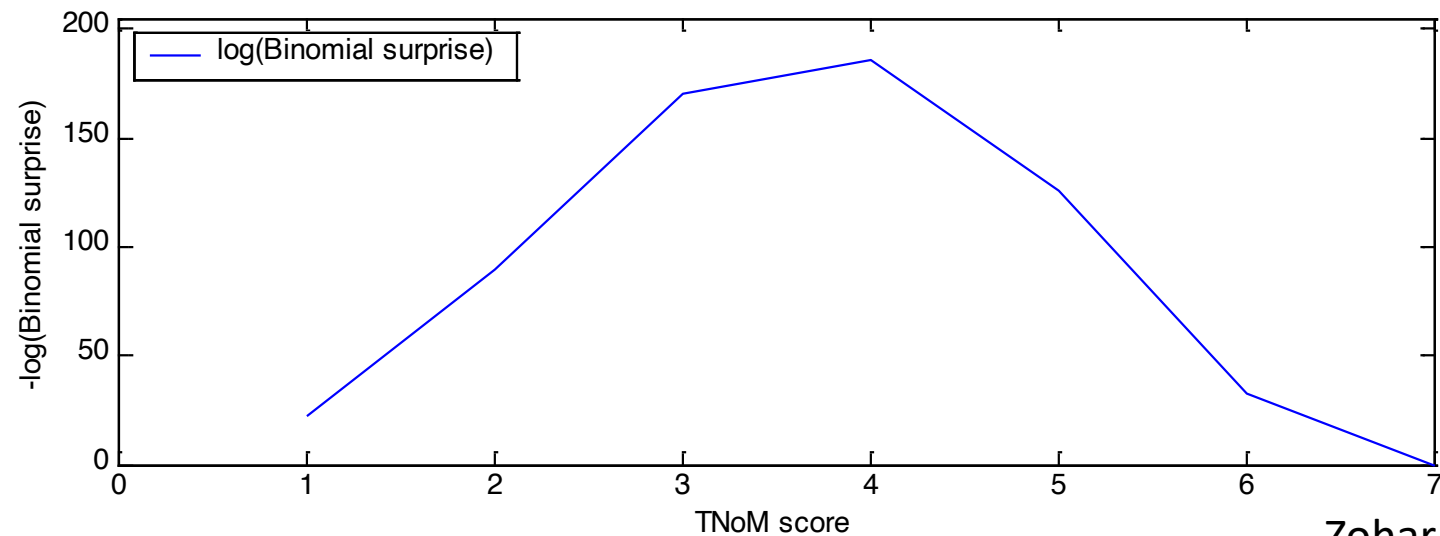
Data sources: Centers for Disease Control & Prevention and Internet Movie Database

tylervigen.com

BRCA1 Differential Expression: Overabundance Analysis



$$TNoM(3) \sim \text{Binomial}(pval(3), 3000)$$
$$P(TNoM(3) \geq 250) = ?$$



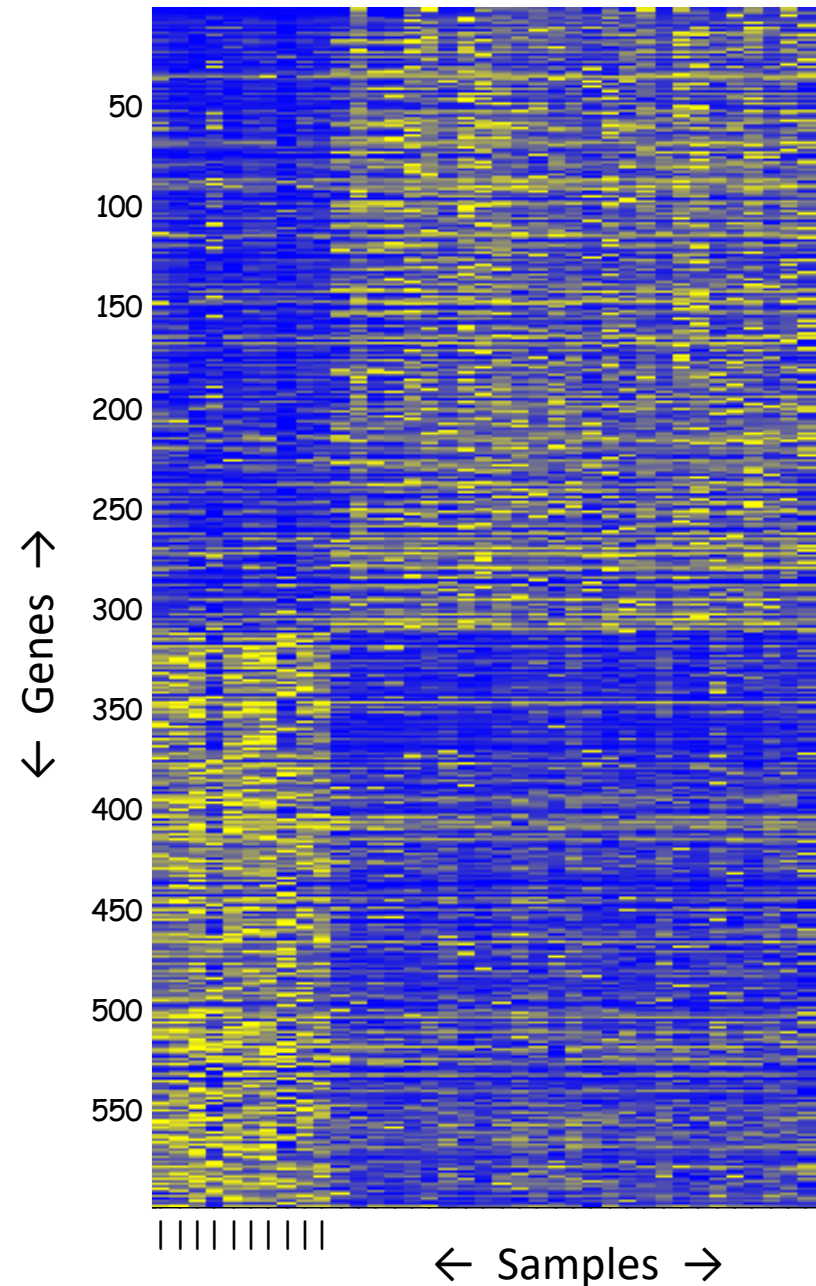
Lung Cancer Informative Genes

Dehan et al,
(Kaminski Lab at Sheba),
Lung Cancer 2007

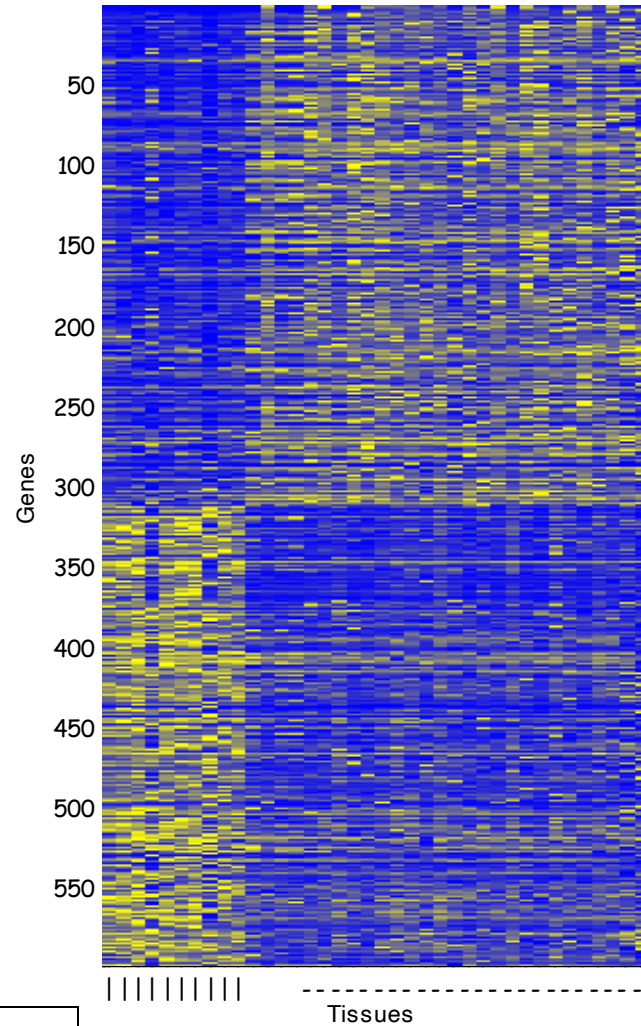
- 24 tumors (various types and origins)
- 10 normals (normal edges and normal lung pools)



LUCA, 38 samples, 14.May.2001 Kaminski.

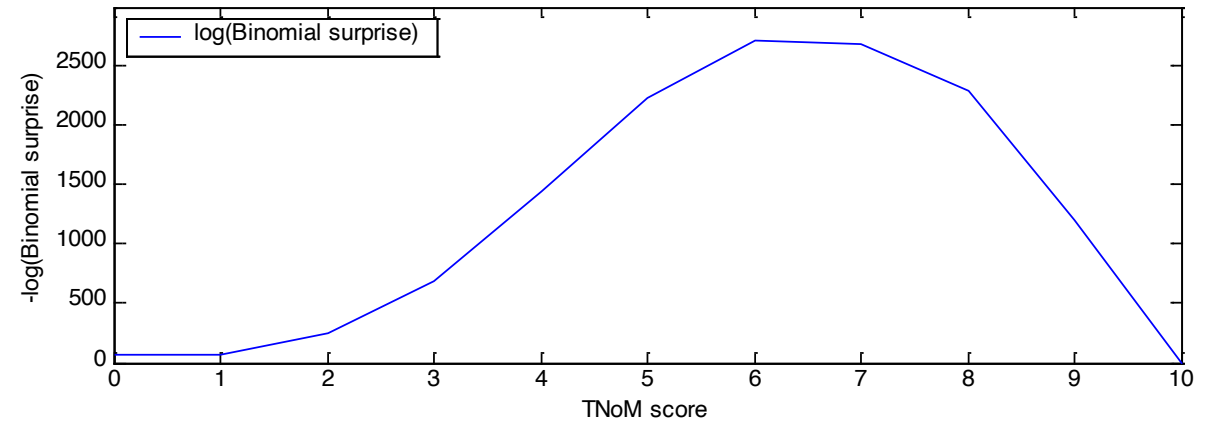
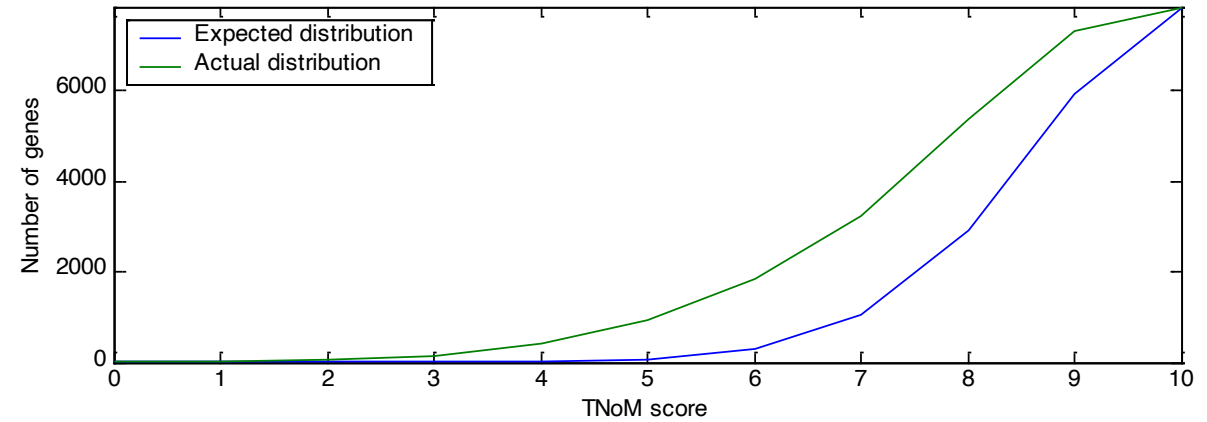


LUCA, 38 samples, 14.May.2001 Kaminski.



Lung Cancer Overabundance Analysis

Lung Cancer Data - Actual and expected TNoM scores distribution



Tall people, Bonferroni

- What is the probability that the person sitting next to you on the bus this evening is $>1.90\text{m}$ tall?
- What is the probability of SOME person in the bus being $>1.90\text{m}$ tall?
- What is the probability that someone on the IDC campus is $>1.90\text{m}$ tall?
- What is a possible naïve (and not very tight) correction? p/N

Bonferroni - cont

What is the probability of rejecting at least one TRUE hypothesis (each test is one hypothesis testing)?
(FWER – Family-Wise Error Rate)

The probability of rejecting a specific TRUE hypothesis is at most p/N .

$$\text{FWER} = P \left\{ \bigcup_{i=1}^{N_0} \left(p_i \leq \frac{\alpha}{N} \right) \right\} \leq \sum_{i=1}^{N_0} \left\{ P \left(p_i \leq \frac{\alpha}{N} \right) \right\} = N_0 \frac{\alpha}{N} \leq \alpha$$

Bonferroni - cont

Assuming that we profiled 20K genes.

Is a DE p-value of $5 \cdot 10^{-6}$ significant (at 0.05)?

How many genes with a p-value at most $5 \cdot 10^{-6}$ are we expecting to see?

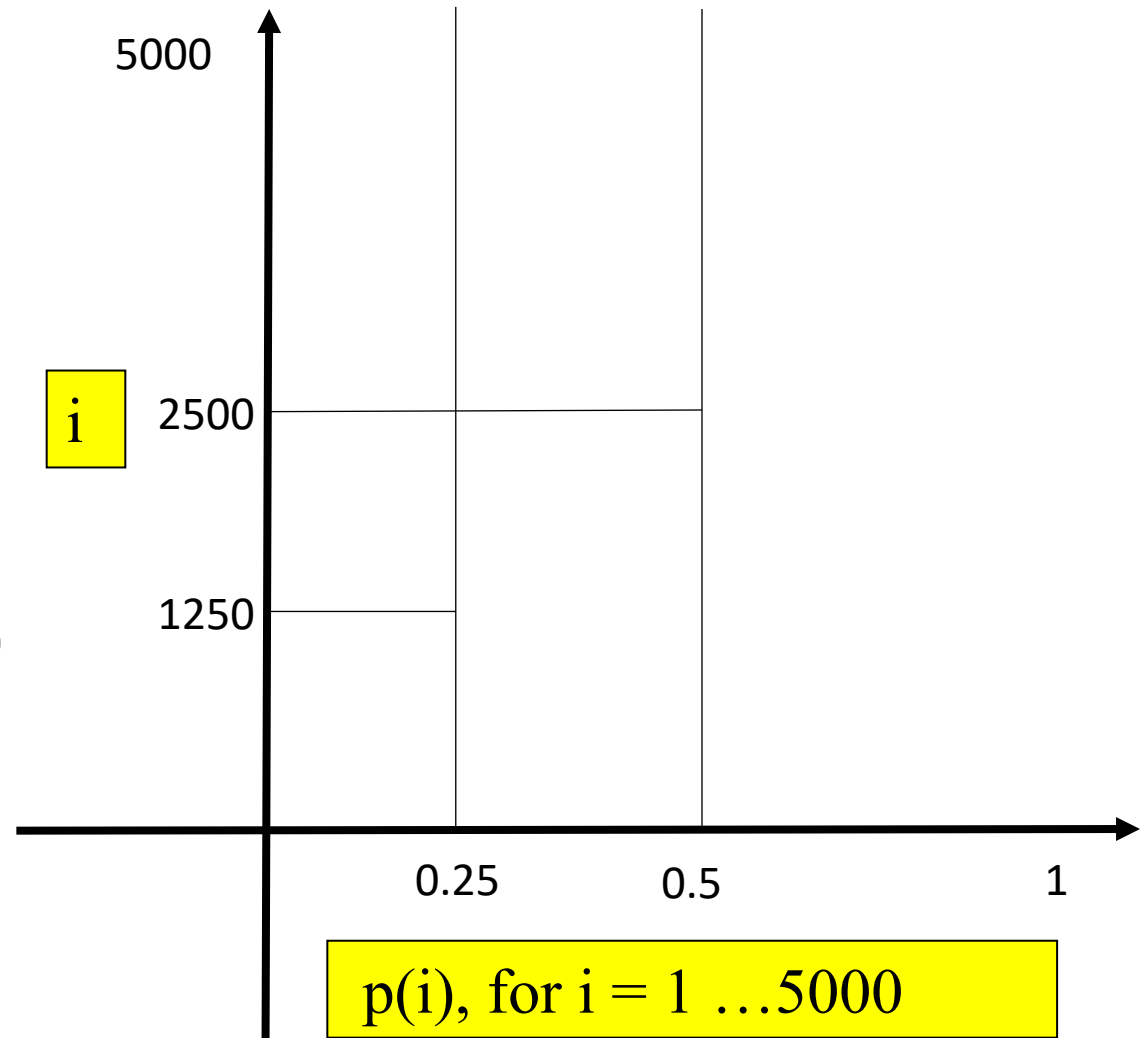
How are p-values distributed under the null model???

- Single Experiment:

1. Toss a fair coin 100 times
2. Compute $\hat{\theta} = \frac{\#1}{100}$
3. Compute the (one sided **left**) p-value of $\hat{\theta}$ under the fair coin null model $P(X \leq \#1)$
 $X \sim \text{Binom}(0.5, 100)$

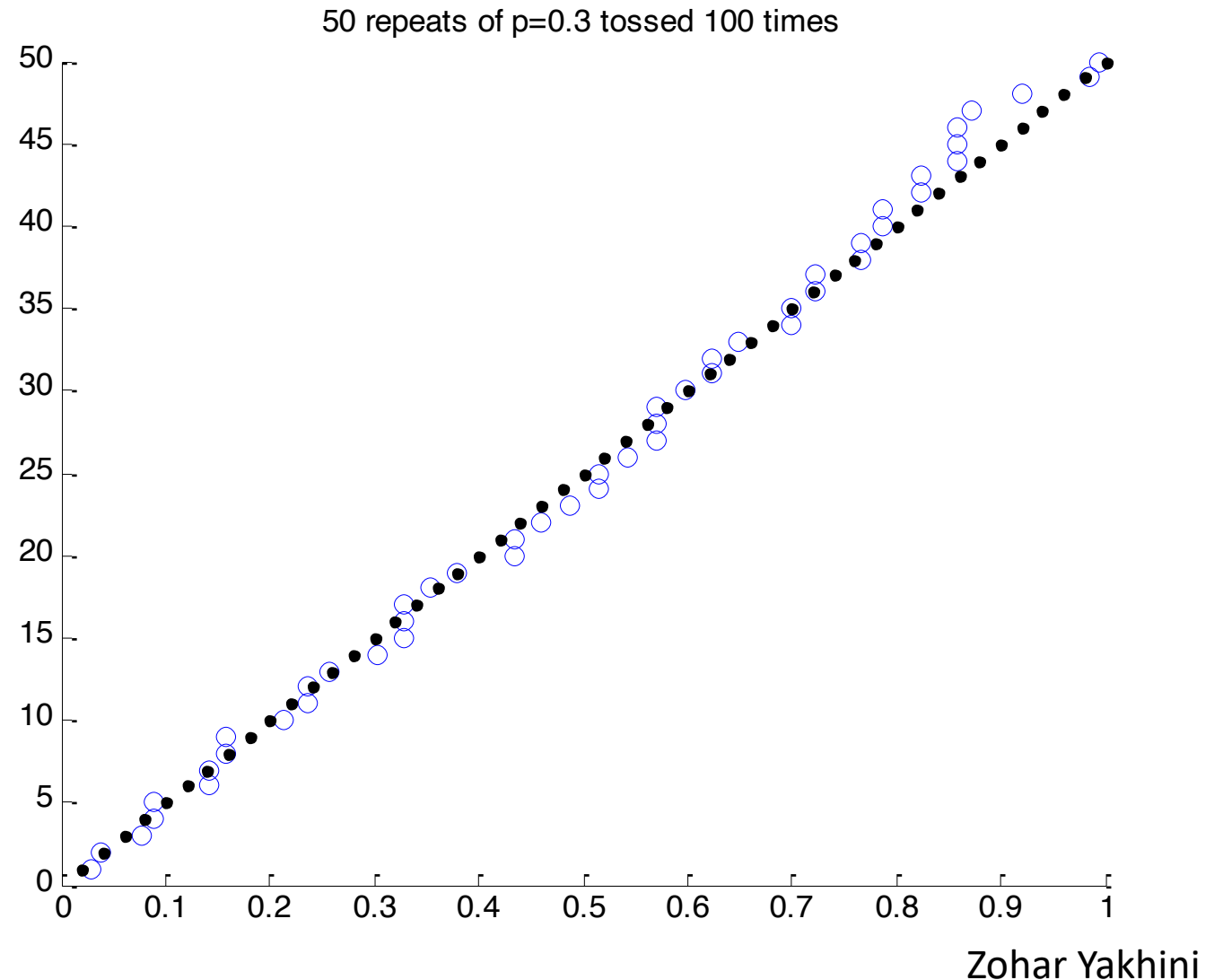
- Repeat 5000 times:

$$(\hat{\theta}(1), p(1)), (\hat{\theta}(2), p(2)), \dots, (\hat{\theta}(5000), p(5000))$$
$$p(1) \leq p(2) \leq \dots \leq p(5000)$$



Distribution of p-values under the null model

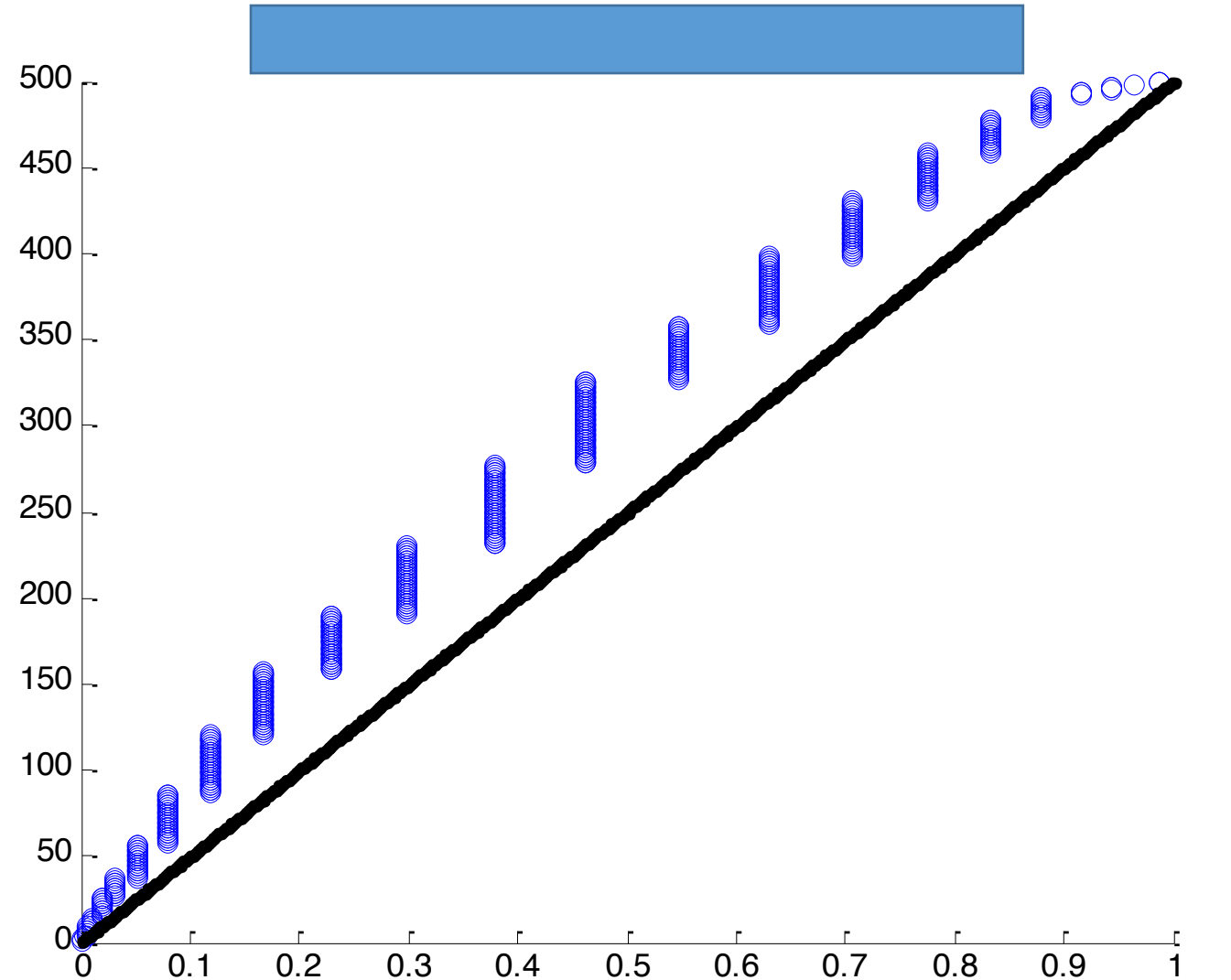
- Generating data using a coin with $p = 0.3$
- under a null model of $p = 0.3$



Distribution of p-values under the null model

$$P(X \leq \#1)$$

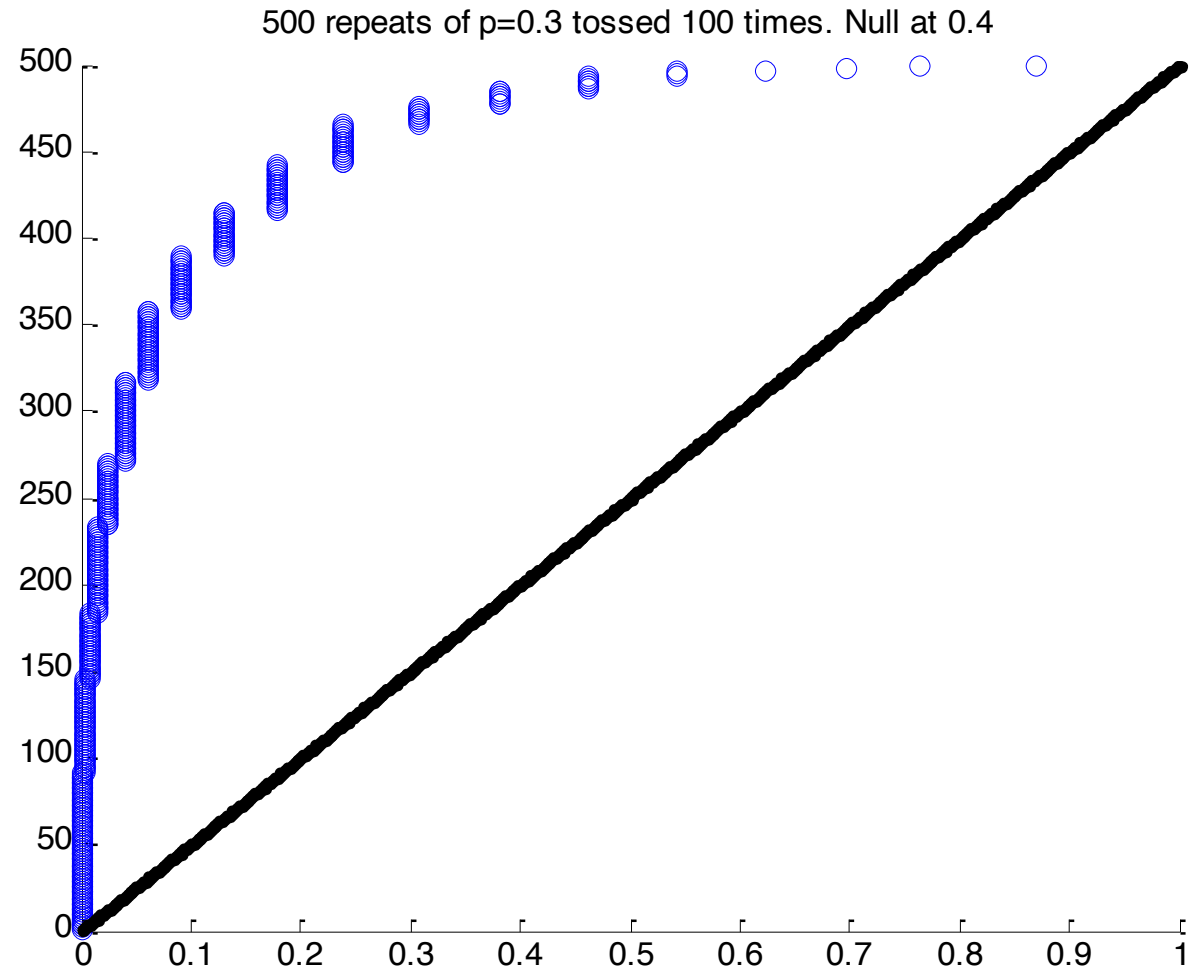
- Generating data using a coin with $p = 0.3$
- under a null model of $p = 0.32$



Distribution of p-values under the null model

$$P(X \leq \#1)$$

- Generating data using a coin with $p = 0.3$
- under a null model of $p = 0.4$



False Discovery Rate (FDR)

What fraction of the observed DE is expected at random (under a null-model)?

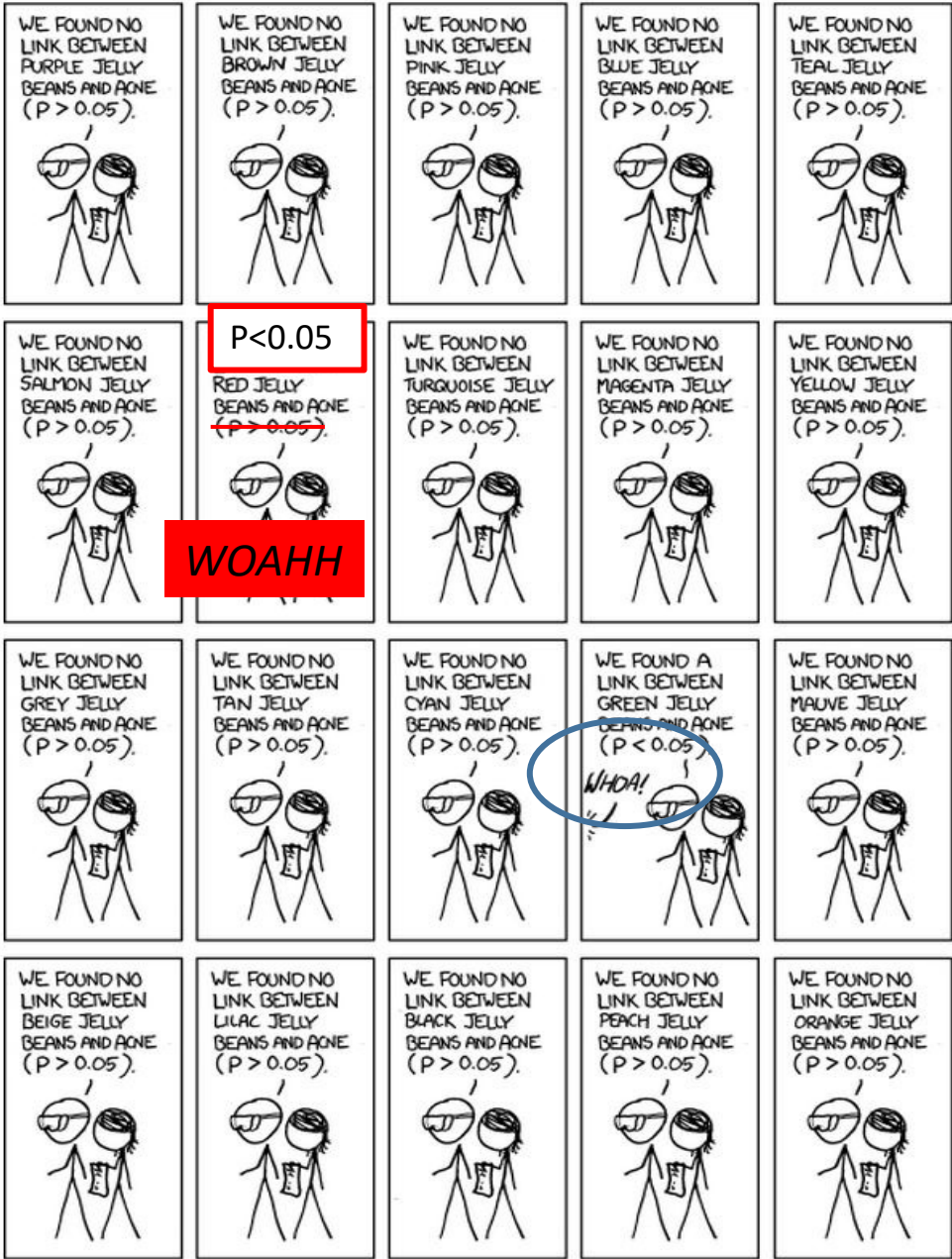
Expected

$$FDR(p) = \frac{pN}{O(p)}$$

Observed: # of tests for which $p\text{-val} \leq p$

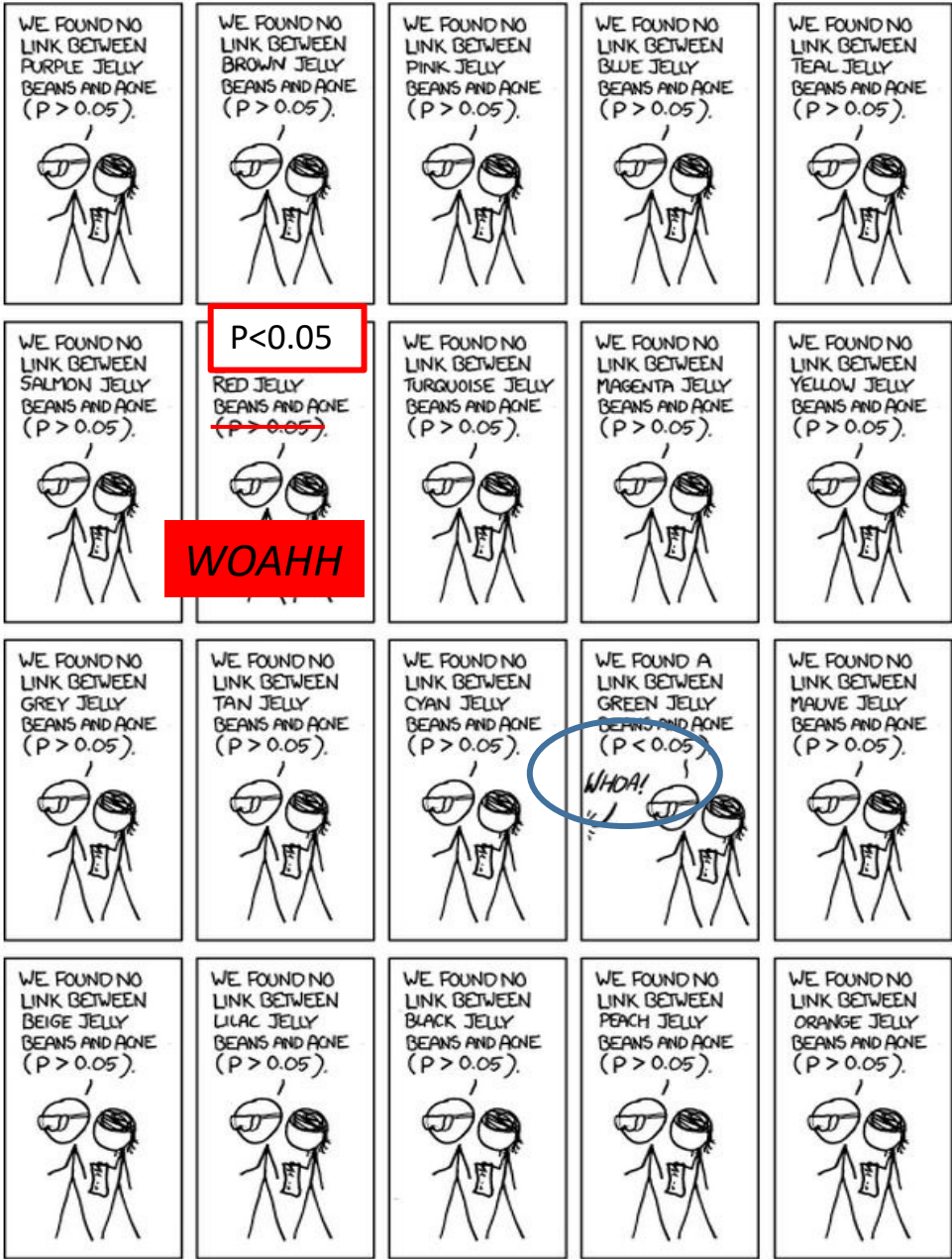
More events

What if two colors
were linked at 0.05?
Three?



More events

$$FDR(0.05) = ?$$



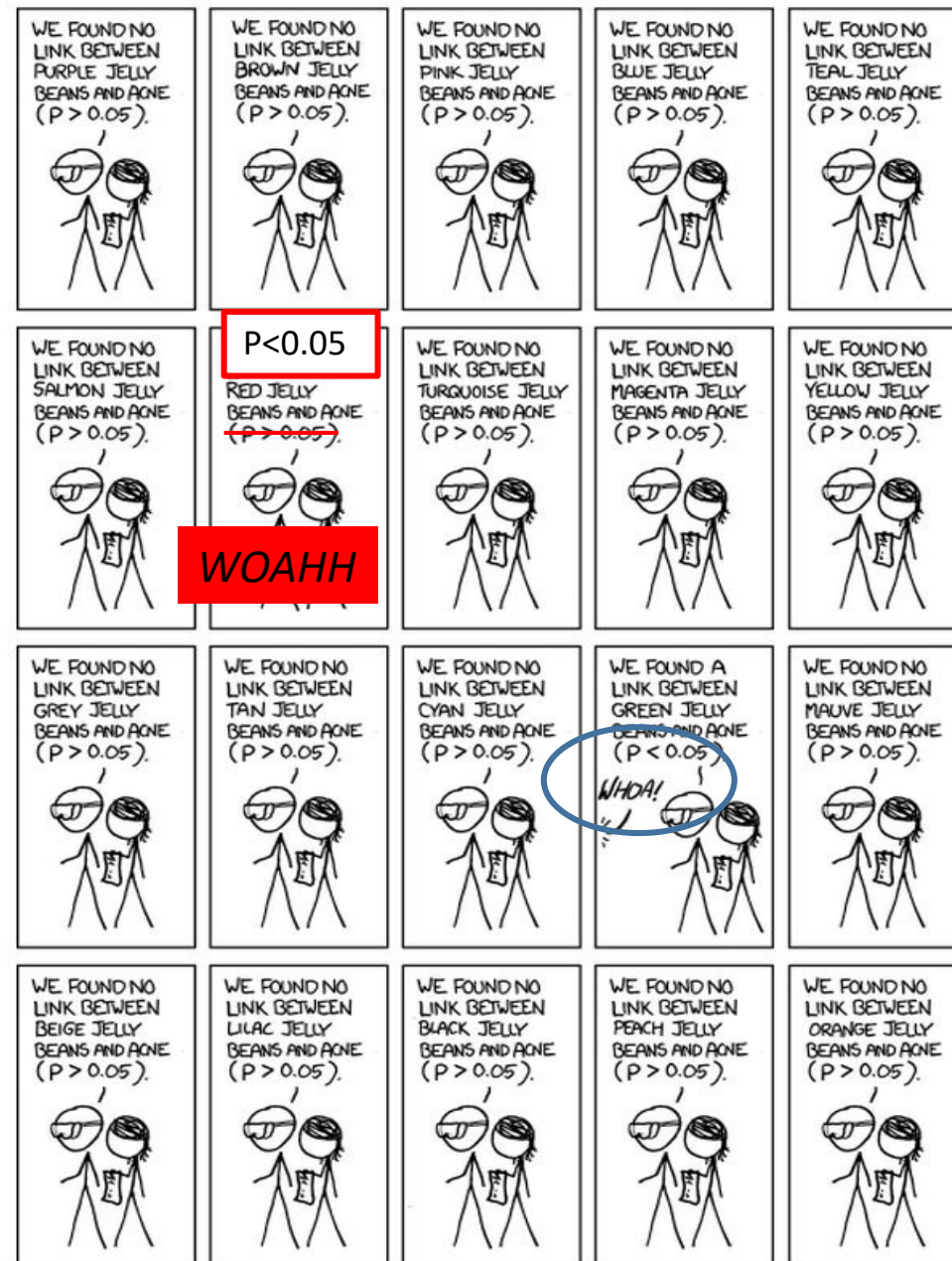
FDR at any i

$$FDR(p) = \frac{pN}{O(p)}$$

$$FDR(i) = \frac{p(i)N}{i}$$

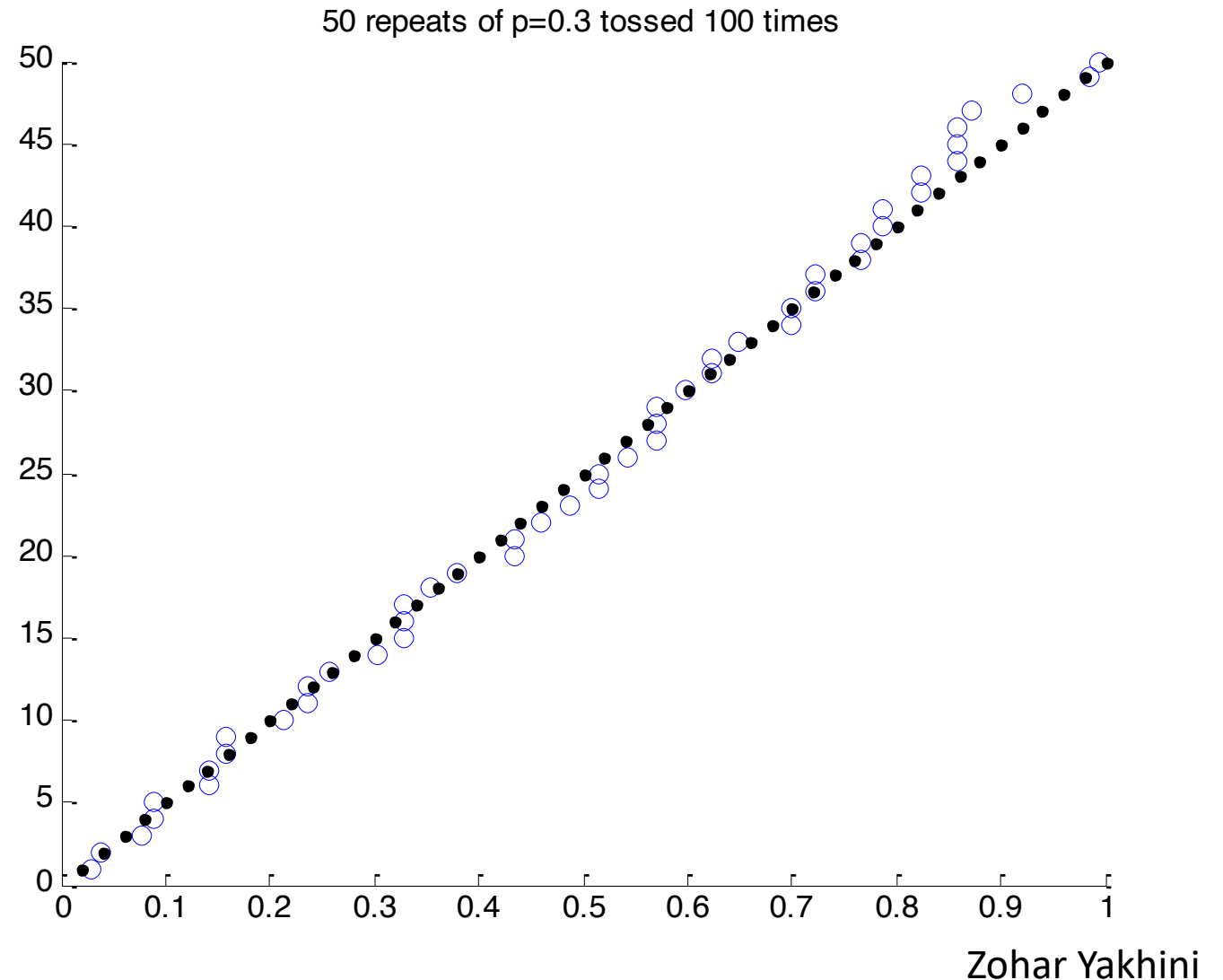
More events

$$FDR(2) = ?$$



Distribution of p-values under the null model

- Generating data using a coin with $p = 0.3$
- under a null model of $p = 0.3$

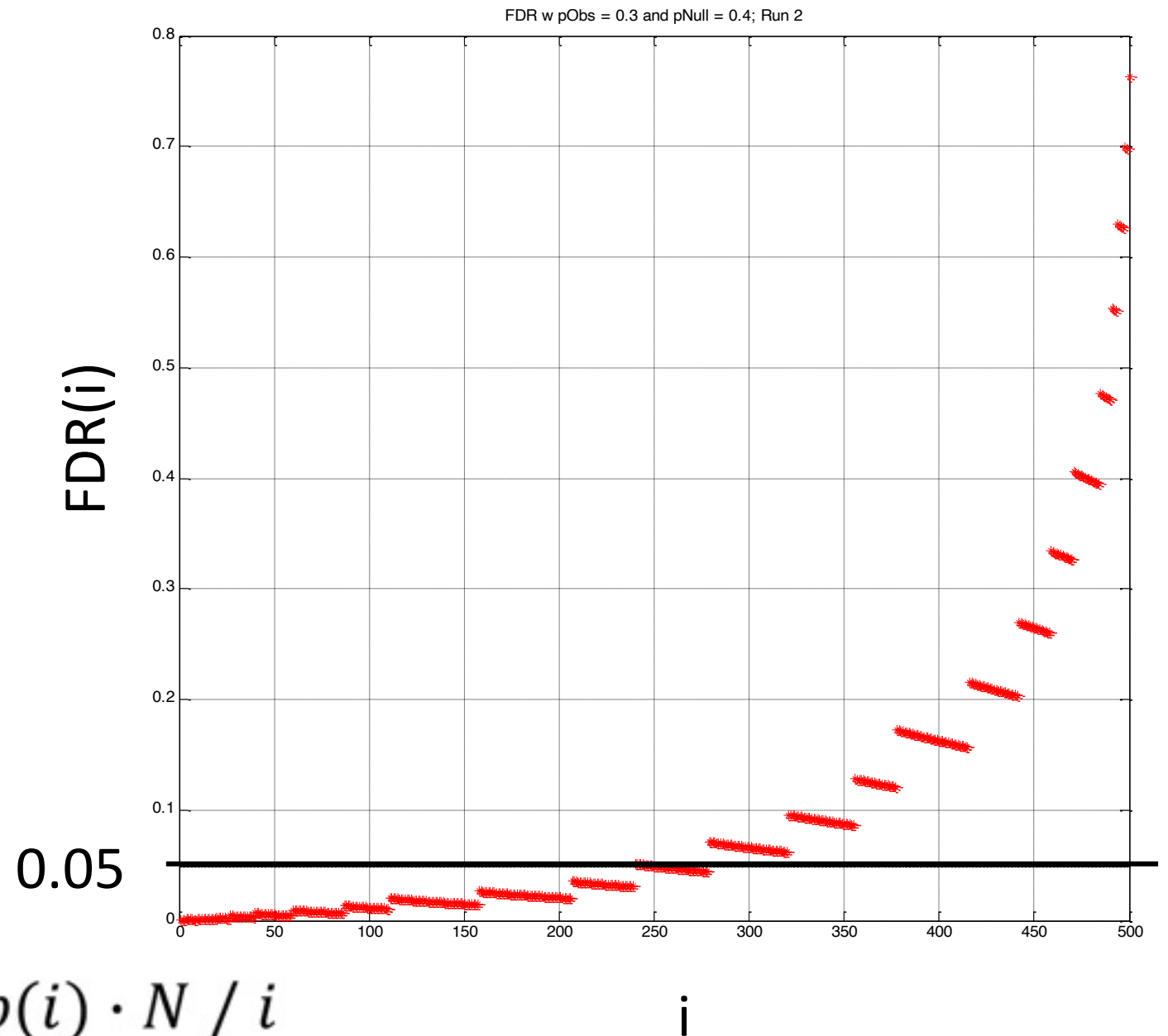


FDR

- Generating data using a coin with $p = 0.3$
- under a null model of $p = 0.4$



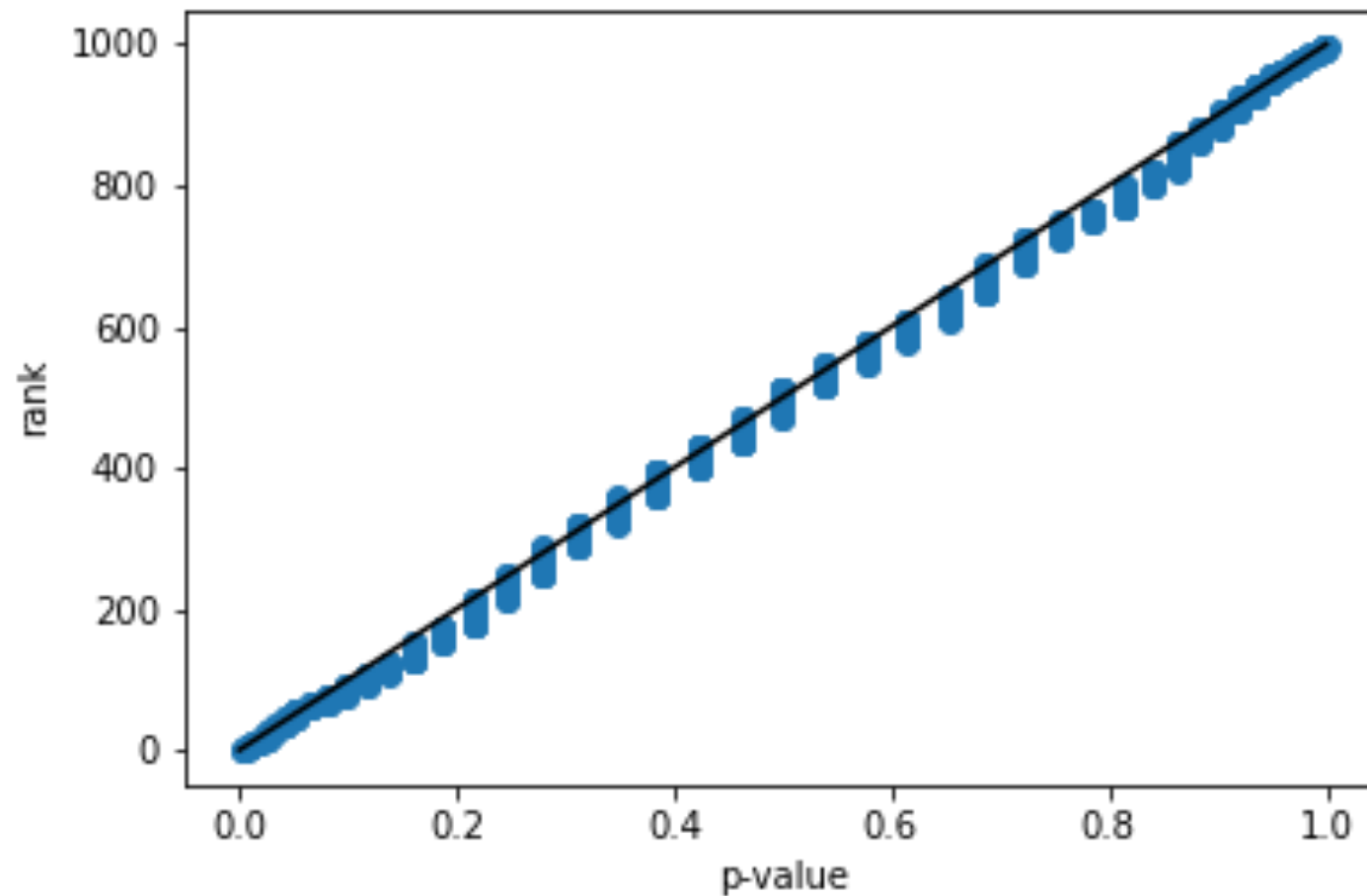
$$FDR(i) = p(i) \cdot N / i$$



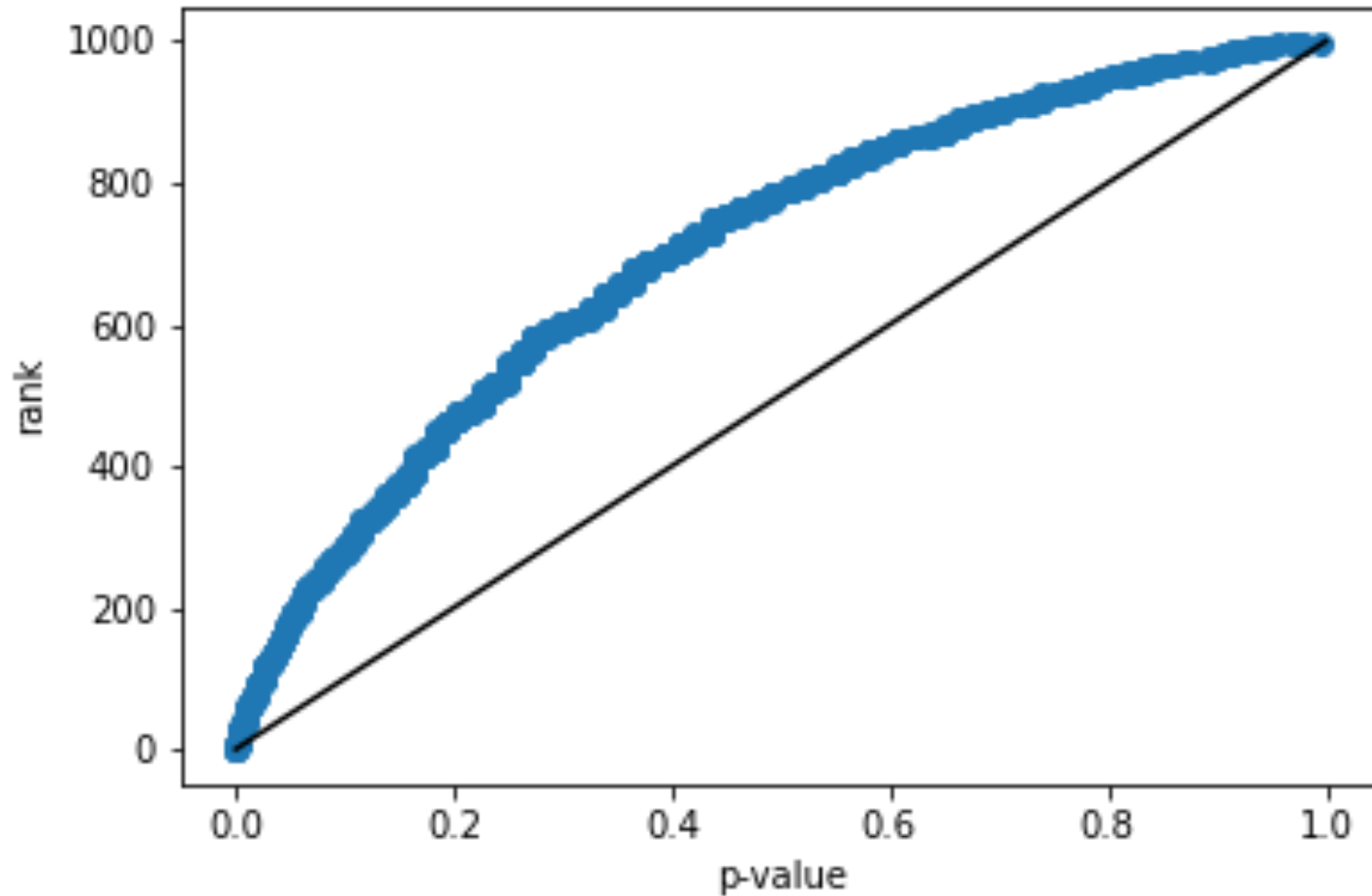
A more subtle correction

$$\text{RobFDR}(i) = \min_{j \geq i} \left(\frac{p(j) \cdot N}{j} \right)$$

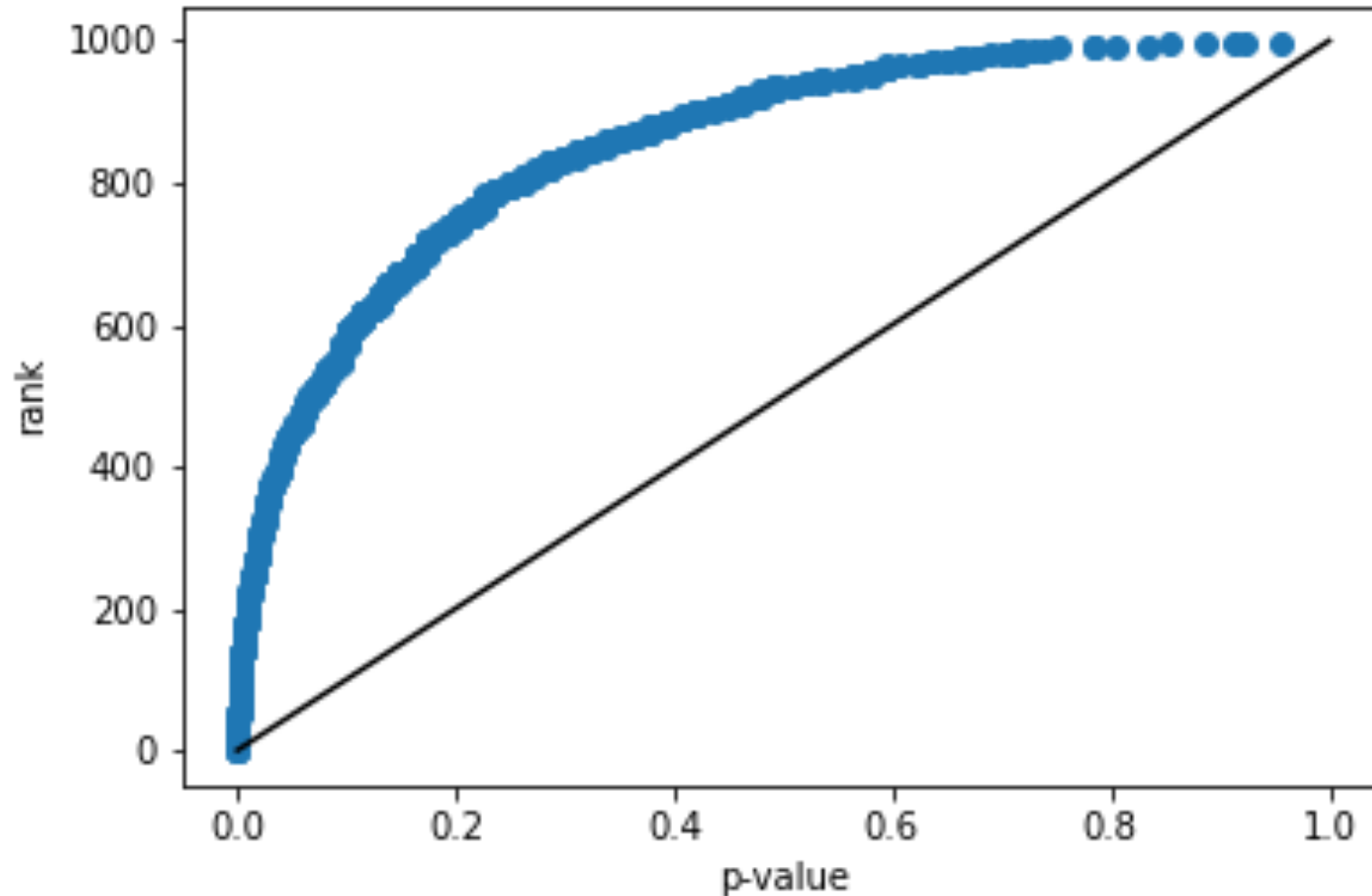
WRS FDR: $A \sim N(0,1)$ and $B \sim N(0,1)$, $n=20,20$



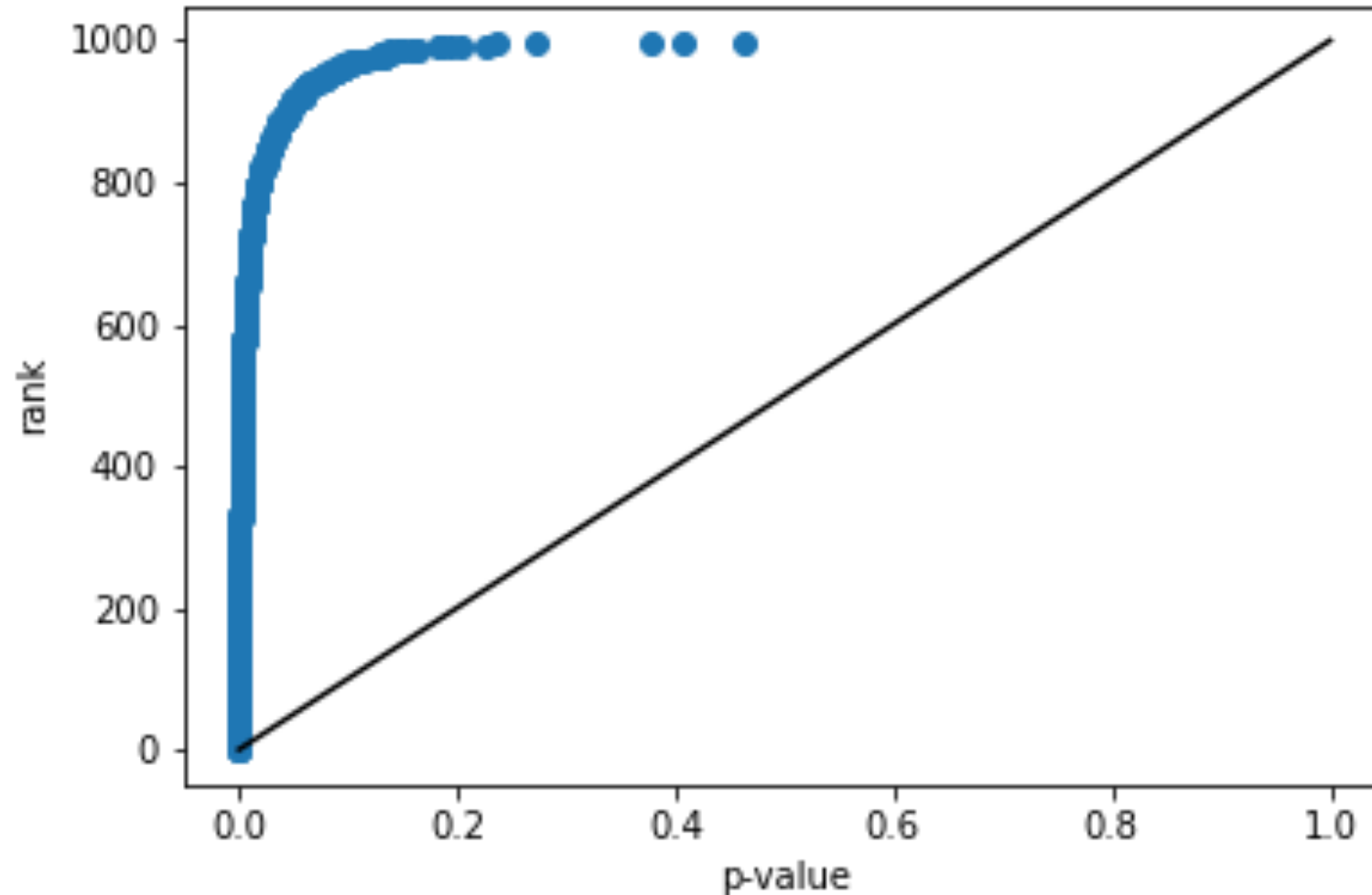
WRS FDR: $A \sim N(0,1)$ and $B \sim N(0.25,1)$, $n=20,20$



WRS FDR: $A \sim N(0,1)$ and $B \sim N(?,1)$, $n=20,20$



WRS FDR: $A \sim N(0,1)$ and $B \sim N(1,1)$, $n=20,20$



FDR – the procedure

Yoav Benjamini and
Yosef Hochberg
1995



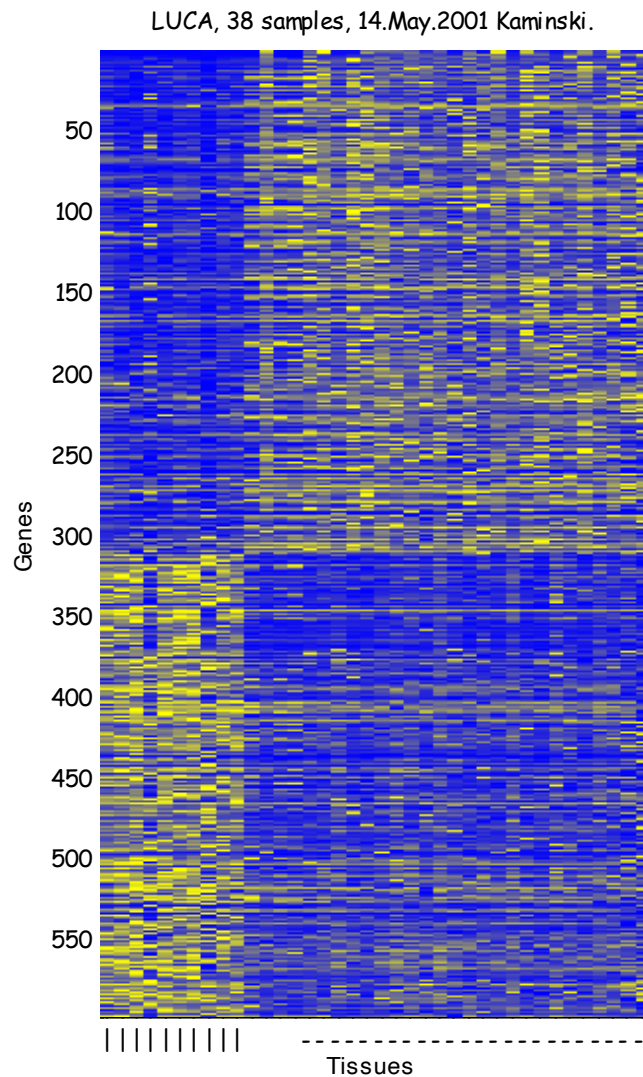
Replace by
the robust
version

- Assume that we performed N measurements (comparisons, observations etc)
- Rank the computed significance of the findings:
 $p(1) \leq p(2) \leq p(3) \leq p(4) \leq \dots \leq p(N-1) \leq p(N)$
- Under the null model, the expected number of observations with p-value better than $p(i)$ is $p(i) \cdot N$
- The false discovery rate at i is therefore:

$$\text{RobFDR}(i) = \min_{j \geq i} \left(\frac{p(j) \cdot N}{j} \right)$$

- A corrected hypothesis testing in this case would be to find the max i that satisfies $\text{FDR}(i) \leq \tau$, where τ (e.g 0.05) is the required false discovery rate.

In the context of DGE ...



We observed (say) 200 genes
at FDR=0.05, using a WRS test

Multiple testing and FDR

- Greater data volumes require more careful inferential statistics approaches.
- Approaches to addressing multiple testing:
 - Bonferroni correction
 - Report FDR results
 - Simulations under a null

