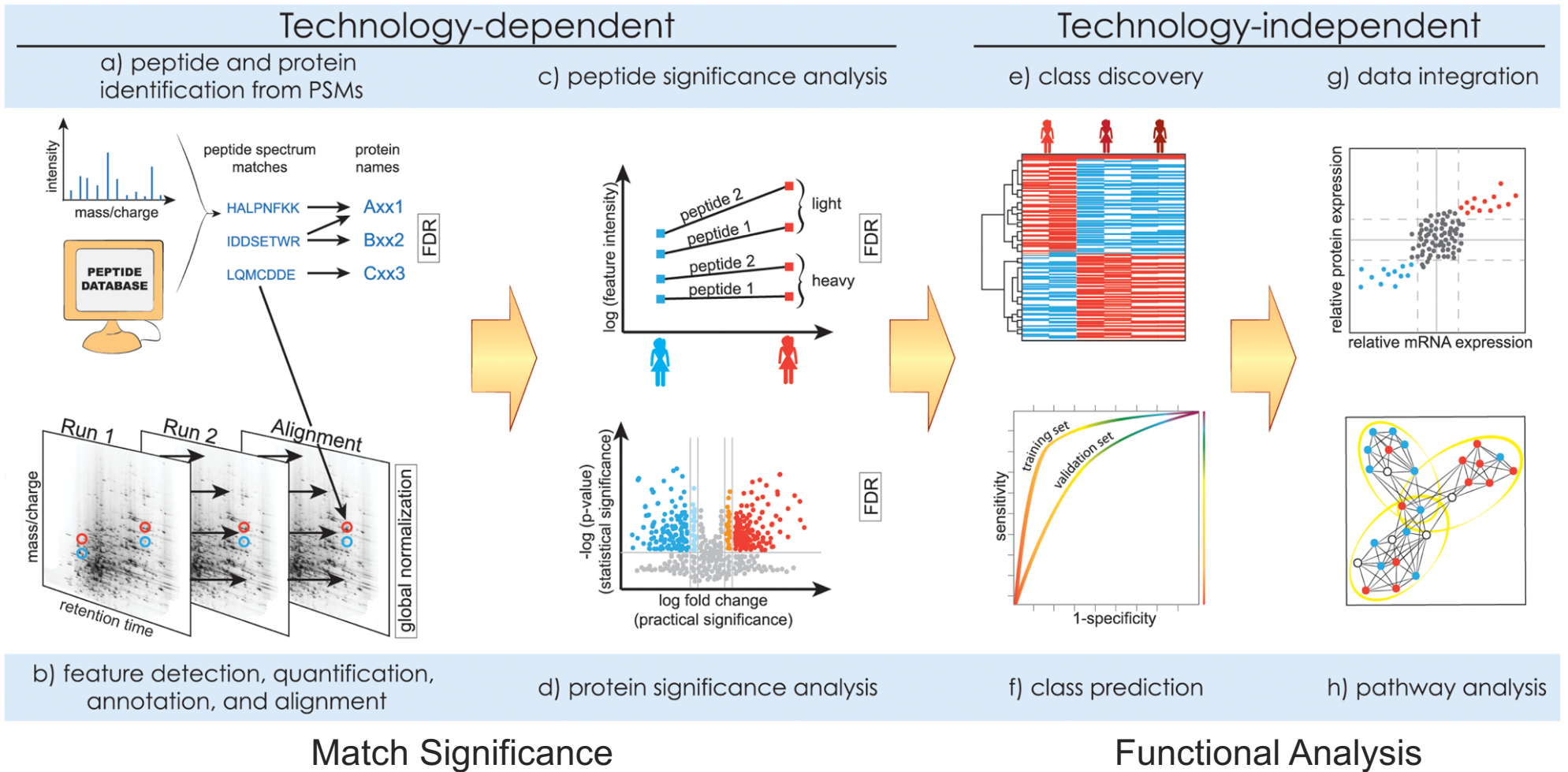# Statistics for Proteomics

**Wilson Wen Bin Goh**
*School of Biological Sciences, Nanyang Technological University*

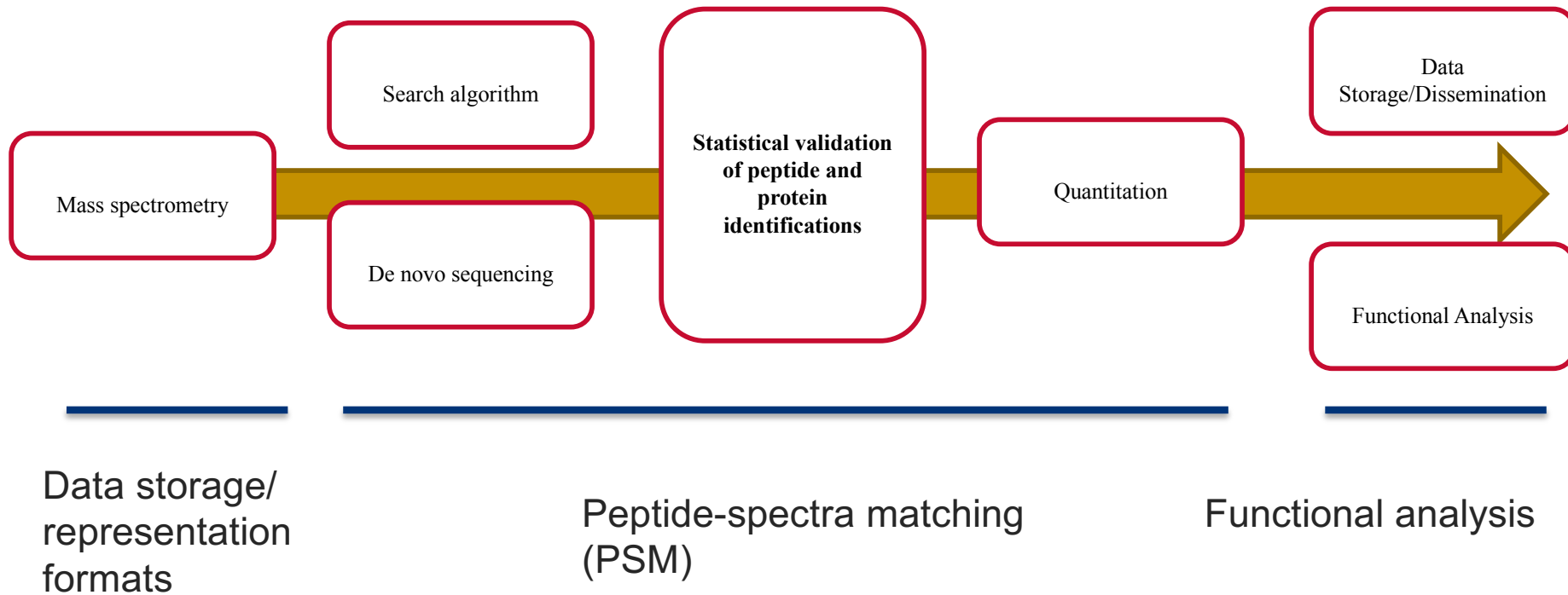*23 November 2017*

# Learning objectives

- Describe the various factors affecting PSM quality

- Describe p-values, FDR and PEP

- Describe and evaluate the various decoy library generation strategies (sequence reversal, sequence randomization) for FDR estimation

- Describe coverage and consistency issues in proteomics

# Using Proteomics in practical applications



Technology-dependent

a) peptide and protein identification from PSMs

c) peptide significance analysis

Technology-independent

e) class discovery

g) data integration

b) feature detection, quantification, annotation, and alignment

d) protein significance analysis

f) class prediction

h) pathway analysis

Match Significance

Functional Analysis

Statistics plays key roles in both areas

NANYANG TECHNOLOGICAL UNIVERSITY

# Overview of proteo-informatics



Mass spectrometry

Search algorithm

De novo sequencing

**Statistical validation of peptide and protein identifications**

Quantitation

Data Storage/Dissemination

Functional Analysis

Data storage/ representation formats

Peptide-spectra matching (PSM)

Functional analysis

NANYANG TECHNOLOGICAL UNIVERSITY

# What is PSM?

- PSM stands for Peptide-Spectra Match
- It is a pairing of sequence with spectra
- You've seen this earlier when we considered library search algorithms
- Earlier we showed a simple scenario where there is only one possible sequence match per spectra...
- But in practice...

# What determines whether or not we get a good PSM?

- Search parameters
- Library quality and size
- Spectra quality
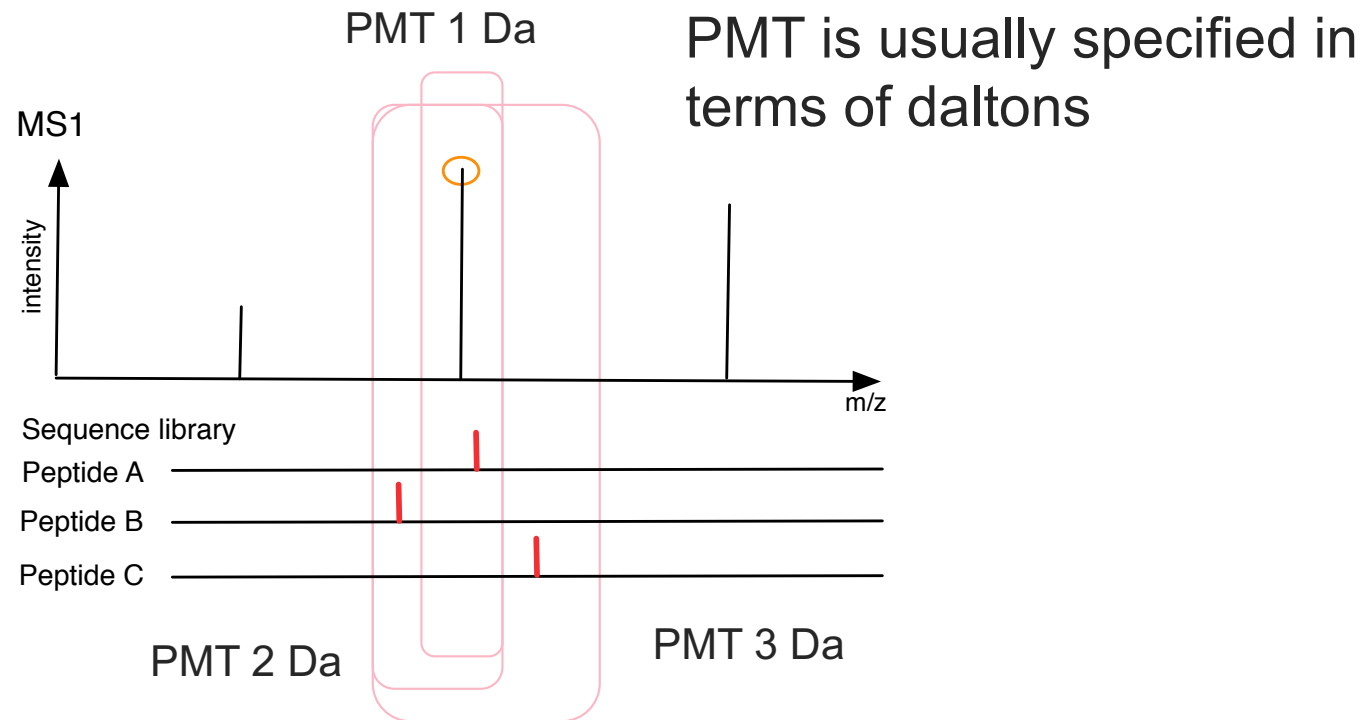- Algorithm scoring method
- Statistical evaluation

# Search parameters

- Precursor mass tolerance (PMT)
- Fragment mass tolerance (FMT)
- Post-translational modification (PTM)

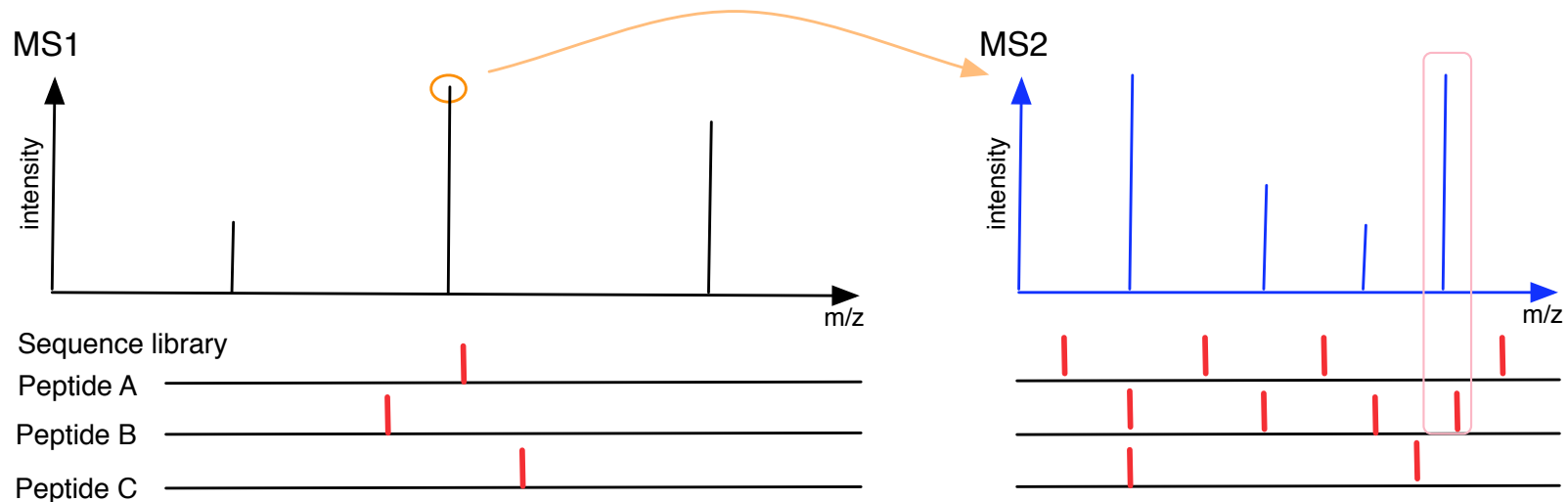# Precursor mass tolerance (PMT)

- PMT deals with MS1 (peptide level)

PMT 1 Da

PMT is usually specified in terms of daltons

MS1

intensity

m/z

Sequence library

Peptide A

Peptide B

Peptide C

PMT 2 Da

PMT 3 Da

What happens when the PMT window size is increased?

# Fragment mass tolerance (FMT)

- FMT deals with MS2 (identification level)



Which looks like the correct answer?
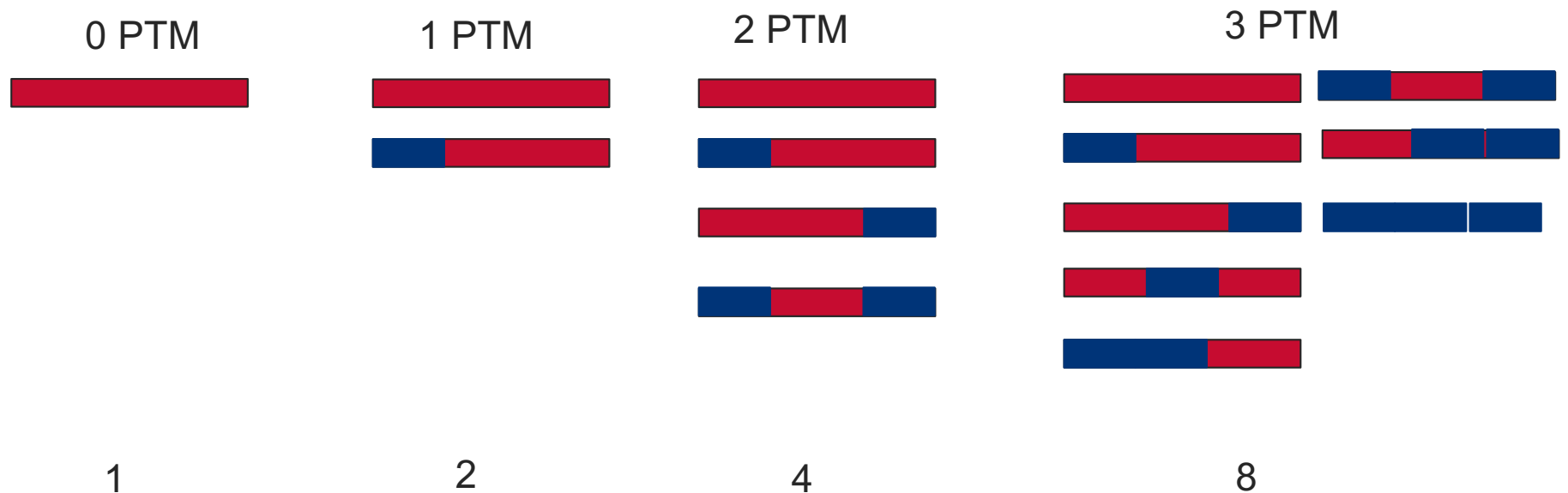
# Post-translational modification (PTM)

- Some peptide sequences undergo chemical modifications resulting in mass shifts

- If the PTMs are not specified in the search space, then the corresponding PSM may not be detectable

MS1

+PTM

MS2   +PTM

intensity

m/z

intensity

m/z

Sequence library (no PTM)

Peptide X _____ | (Missed)

Sequence library (PTM specified)

Peptide X _____ | (Found)

## Some fragments will contain the PTM in MS2

NANYANG TECHNOLOGICAL UNIVERSITY

# Search space

- The set of candidate peptides to be considered for potential match to spectra
- Without PTMs, the search space is simply the set of peptides
- With PTMs, the search space effectively doubles for every PTM to be considered.

# Examples of typical PTMs

- Phosphorylation
- Ubiquitination
- O-GlcNAcylation
- Methylation
- Acetylation

- Succinylation
- SUMOylation
- Citrullination

Some 260 000 PTM sites that have been identified in the human proteome thus far, but only a few have been assigned to key regulatory and/or other biological roles!

It is difficult to pin-point exact locations of PTMs as well. And incorporating all possibilities (where there is only 1 or few right matches)… can lead to high false positive rates (we will see how later).

NANYANG TECHNOLOGICAL UNIVERSITY

# Library quality and size

- UniProt sequence library has 2 databases
  - SwissProt (manually curated and reviewed) - > 500K sequences
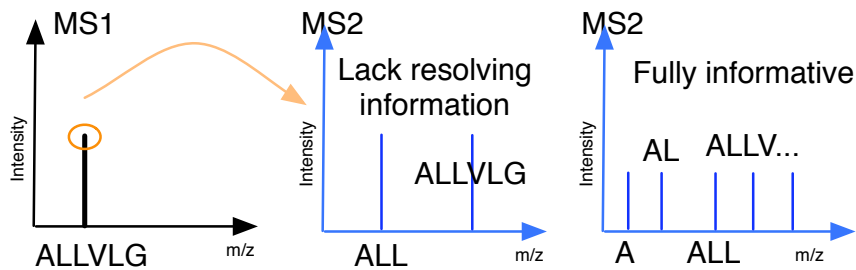  - TrEMBL (Automatic annotation, no review) - > 90M sequences

# Consider this scenario…
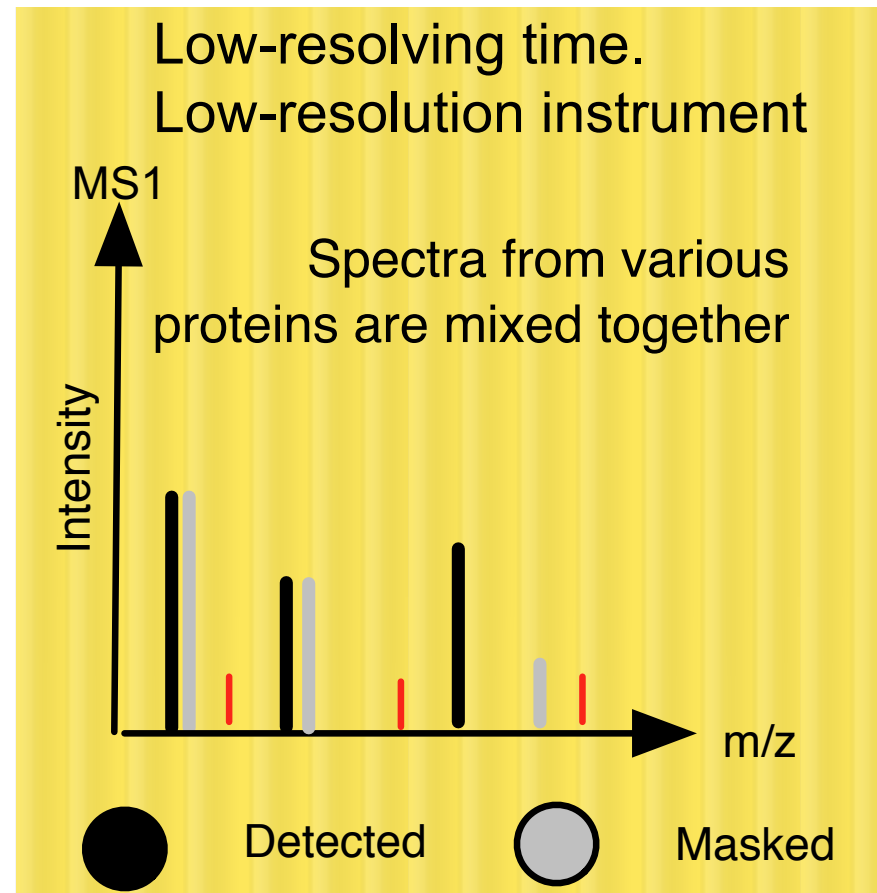


What are potential explanations for 1, 2 and 3?
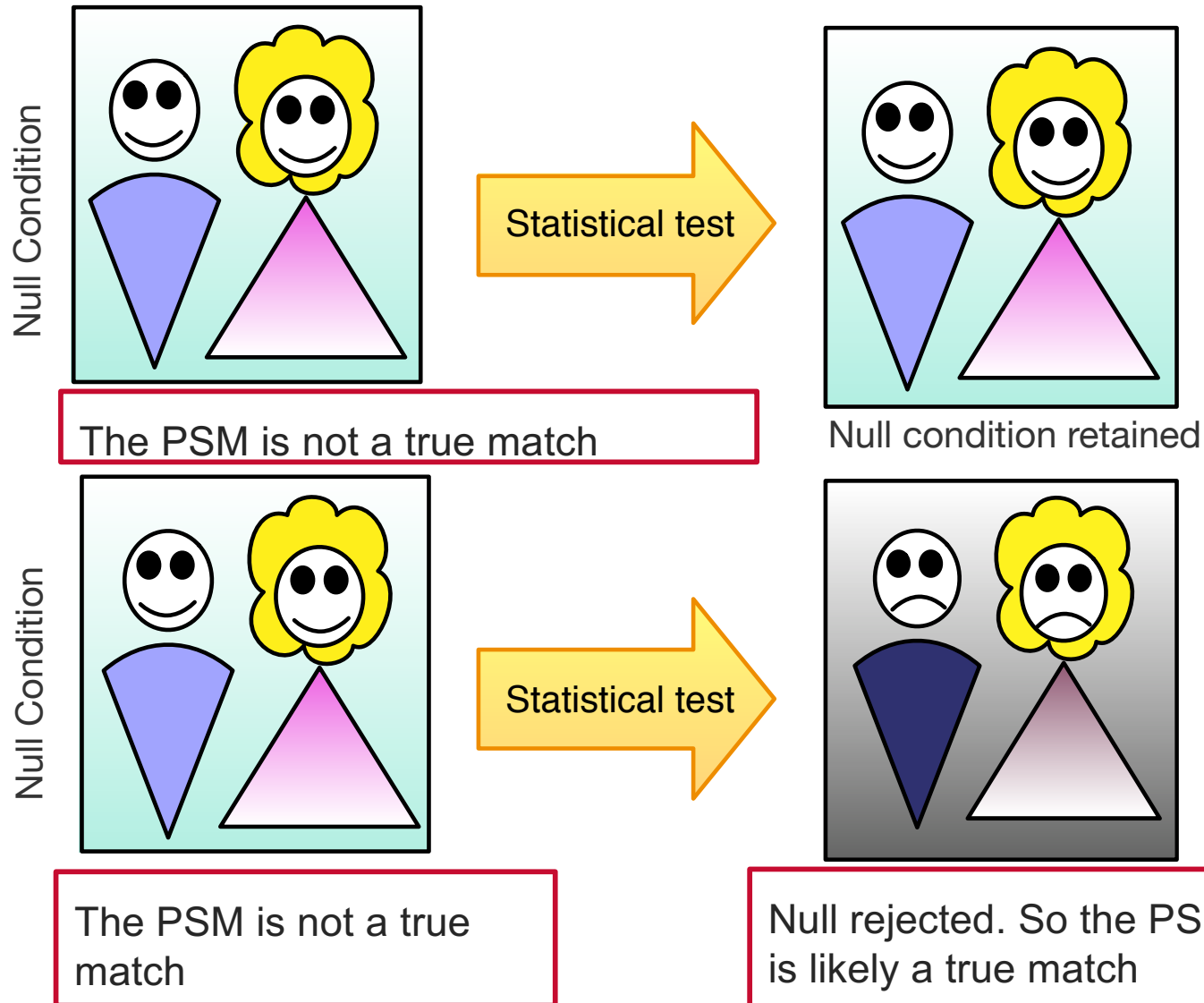
# Spectra Quality

**Incomplete fragmentation**

**Mixed signals**
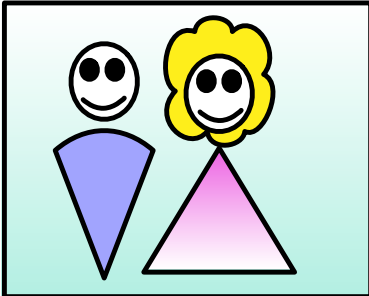


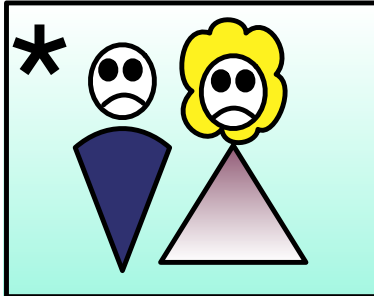Complete MS2 profile allows confident identification of spectra

# Statistical testing
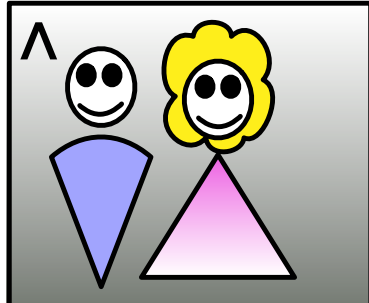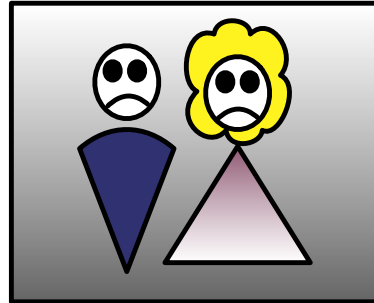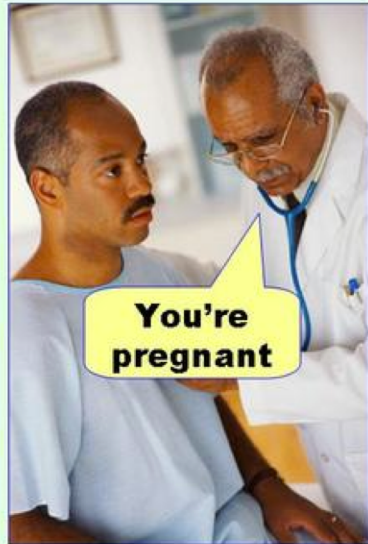
**The elements of null hypothesis statistical testing**



Null Condition

Statistical test

The PSM is not a true match

Null condition retained

Null Condition

Statistical test

The PSM is not a true match

Null rejected. So the PSM is likely a true match

**NANYANG TECHNOLOGICAL UNIVERSITY**

# Possible outcomes from a statistical test



**4 Possible outcomes**

**Predicted as**

|  | Happy (Irrelevant) | Unhappy (Relevant) |
|---|---|---|
| **Happy (Irrelevant)** | True Negative | * False Positive |
| **Unhappy (Relevant)** | ^ False Negative | True Positive |

**Reality**

Goh and Wong. Dealing with confounders in -omics analysis. Trends in Biotechnology, 2018

**NANYANG TECHNOLOGICAL UNIVERSITY**

# How to remember?

Do you recall type I and II statistical errors?

Type I: Reject the null when the null is true
Type II: Fail to reject the null when the null is not true
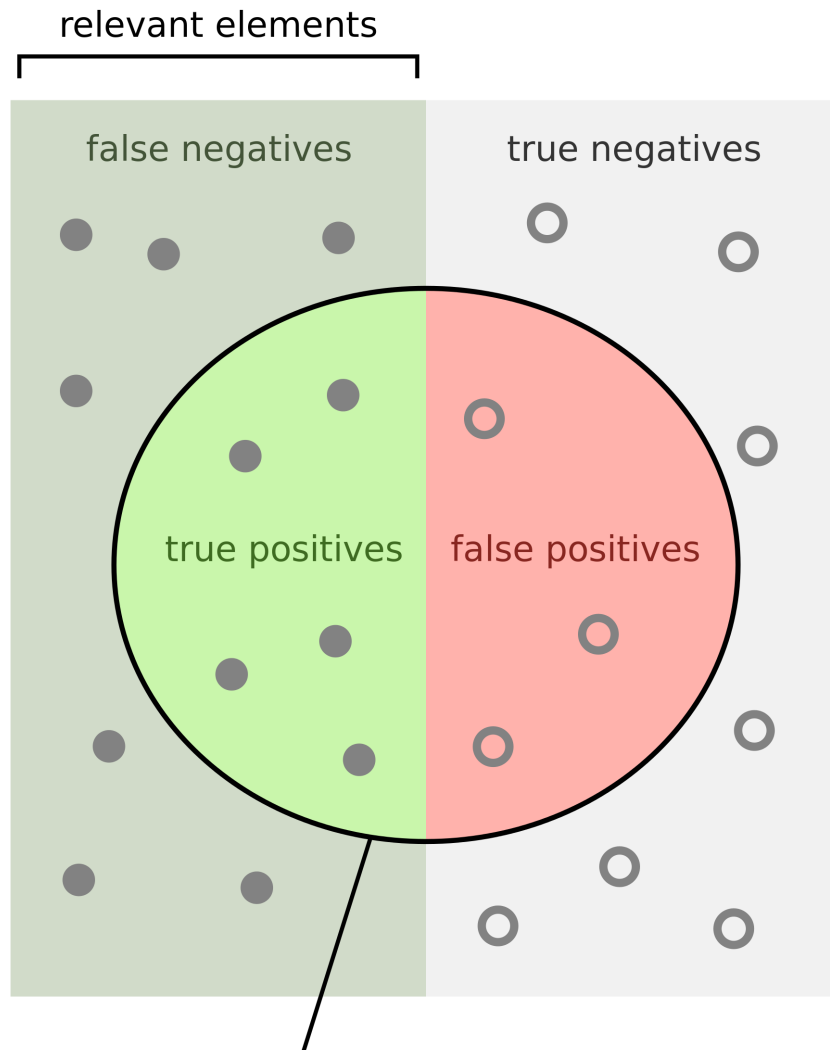
# Possible outcomes given the PSMs

The PSM is…

Predicted as

| | CORRECT | WRONG |
|---|---|---|
| CORRECT | | |
| WRONG | | |

BUT IN Reality..

Imagine we do this for every spectra…

# Recall, Precision and the F-score

relevant elements

false negatives      true negatives

true positives      false positives

selected elements

e.g. let's say we set a p-value cutoff of 0.05

How many selected items are relevant?

How many relevant items are selected?

Precision = 

Recall = 

Precision: Of the selected feature,
How many are correct?

Recall: Of the selected feature,
What is the proportion of all the correct ones we got?
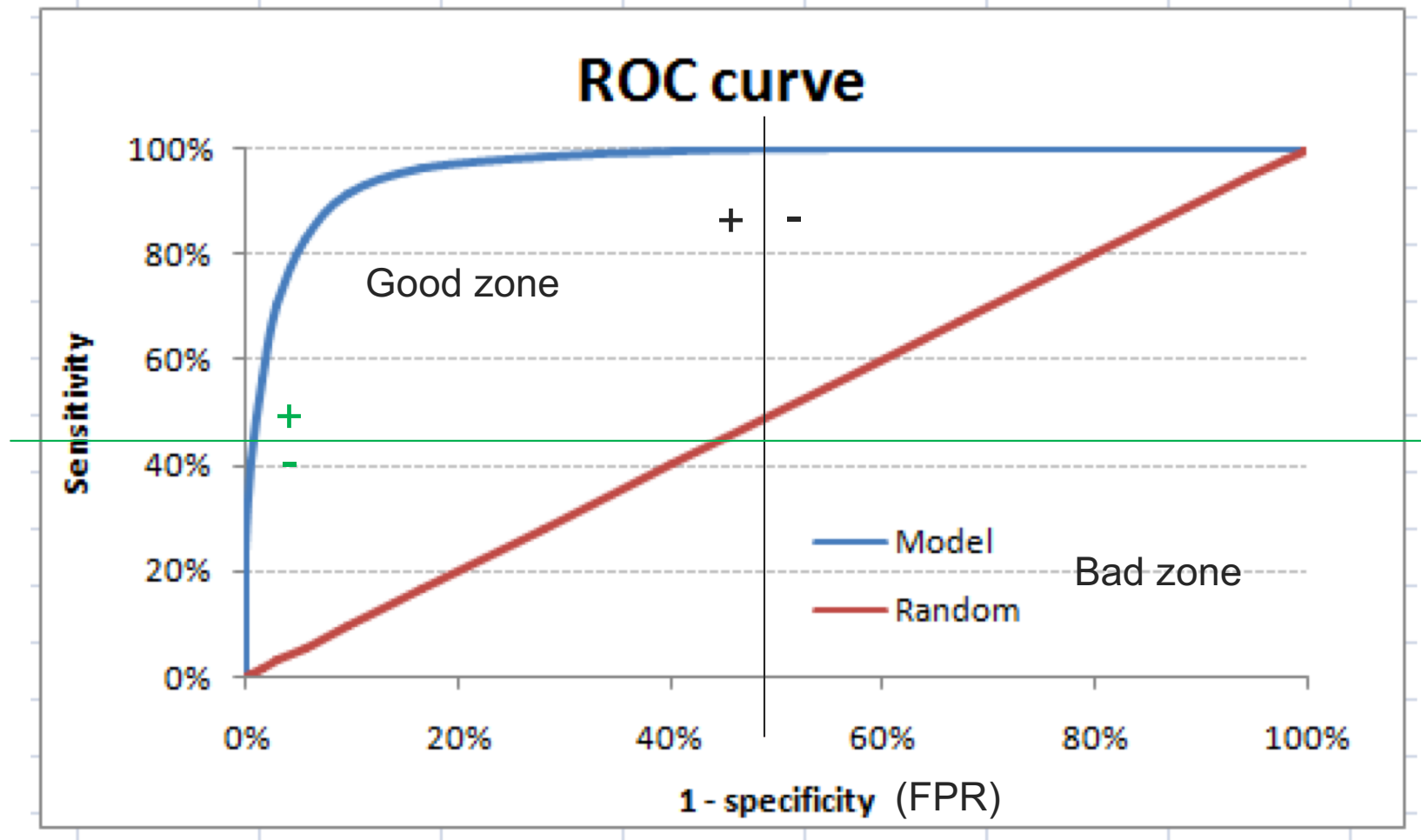
Precision and recall can be combined as:

$$F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

Recall is also often called sensitivity or True positive rate

NANYANG TECHNOLOGICAL UNIVERSITY

# Precision and recall works against each other

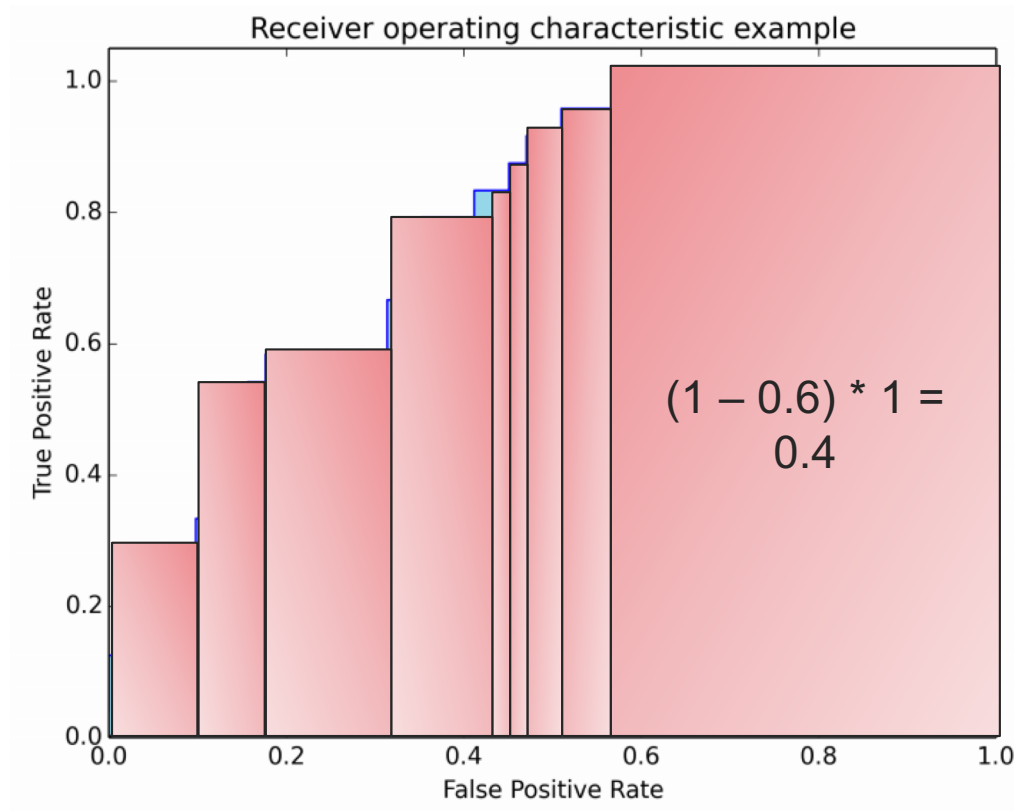# The Receiver Operator Characteristic (ROC) curve



$$\text{specificity} = \frac{\text{number of true negatives}}{\text{number of true negatives} + \text{number of false positives}}$$

Can you write down the formula for FPR yourself?

# The Area Under Curve (AUC)

More rightfully called AUROC (Area under the ROC curve)



**Receiver operating characteristic example**

Plot with y-axis "True Positive Rate" (0.0 to 1.0) and x-axis "False Positive Rate" (0.0 to 1.0).

$$(1 - 0.6) * 1 = 0.4$$

The blue area corresponds to the AUROC. The dashed line in the diagonal is expected performance due to random chance (so we have to be better than chance)
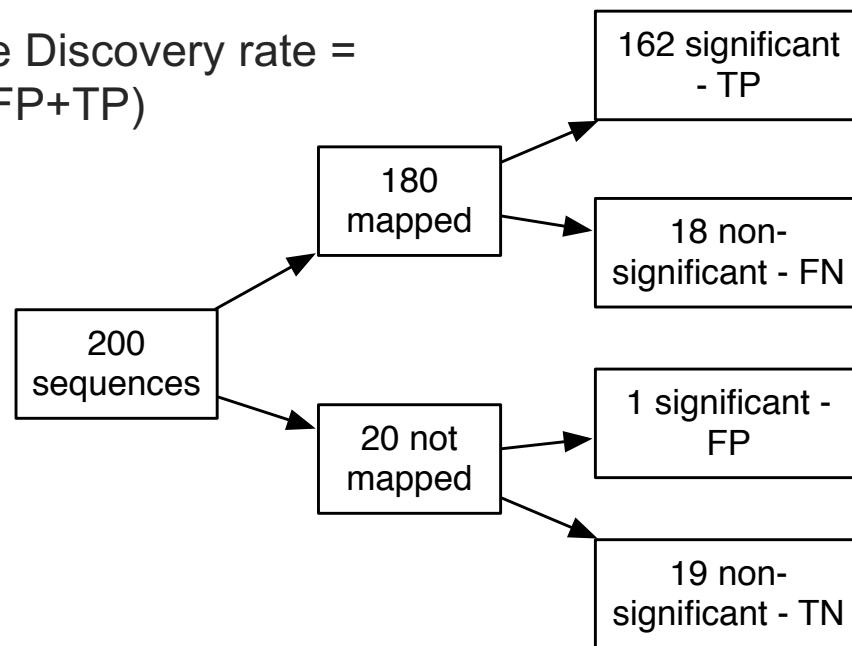
Total area = 1 x 1 = 1

Half area under the diagonal = ½ = 0.5

One simple way to get the AUROC is to simply calculate the area using simple length x breadth. But of course one may use calculus.
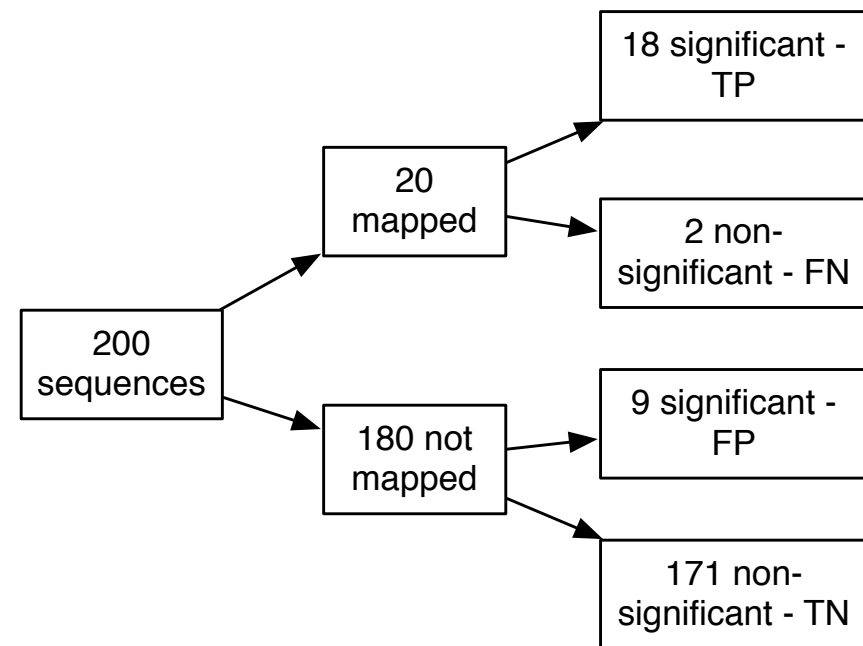
NANYANG TECHNOLOGICAL UNIVERSITY

# The False Discovery Rate

The FDR relates to the proportion of errors amongst predictions. It is equals to 1 - precision
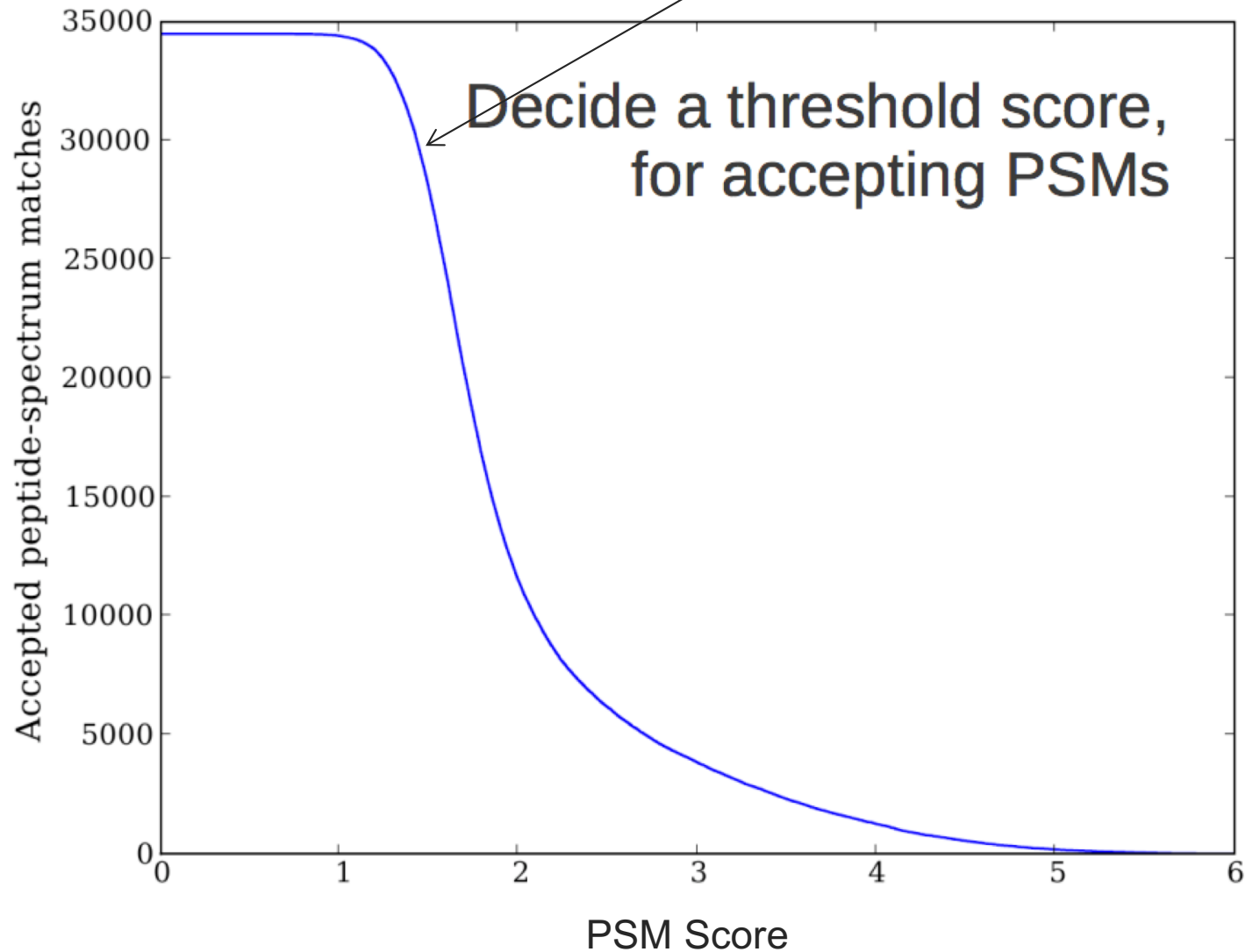
False Discovery rate = FP/(FP+TP)



1/163

9/27=1/3

The FDR is sensitive to the proportion of true features in the data.

# How to get all the good quality matches?
## What if we just use our eye power?



**PSM quality score**

Can do it by eye?

Decide a threshold score, for accepting PSMs

*Accepted peptide-spectrum matches* (y-axis: 0 to 35000)

*PSM Score* (x-axis: 0 to 6)

Statistics provides a more objective manner of evaluation

# Statistics help us determine the best match

- p-values
- False Discovery Rate (FDR)
- Posterior Error Probability (PEP)

The p-value

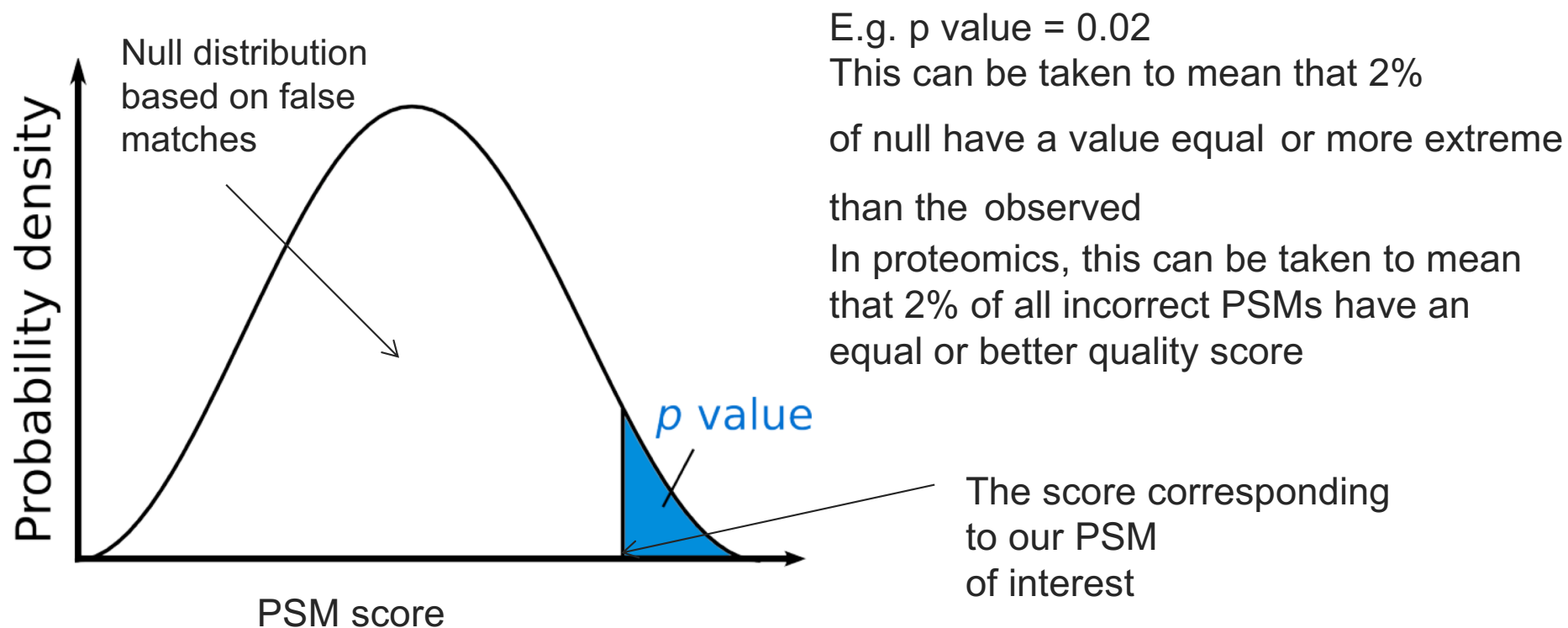We can set up our PSMs as a statistical test (based on the following hypothesis statements)

H0: A PSM is incorrect          H1: A PSM is correct

We may reject the null hypothesis with a certain degree of error:
Type I: Falsely reject the null hypothesis (False positive)
Type II: Falsely accept the null hypothesis (False negative)

E.g. p value = 0.02
This can be taken to mean that 2%

of null have a value equal or more extreme

than the observed

In proteomics, this can be taken to mean that 2% of all incorrect PSMs have an equal or better quality score

Null distribution based on false matches

Probability density

*p* value

The score corresponding to our PSM of interest

PSM score

The probability of obtaining an equal or more extreme result, assuming the null hypothesis is true (Type I error)

# The p-value

- We are comparing the observed score against a distribution of "null scores"
- The null distribution are comprised of the natural distribution of values when there is no signal i.e., when a PSM is incorrect (or the null statement is true)
- Why does this make sense?
- Because under this setup, a small p-value would imply that the observed PSM score is very significant (unlikely to arise due to chance).

# The p-value

- For each hypothesis tested. Suppose we use a statistical cutoff at 0.01, we should therefore expect 1 in 100 times the result is a false positive

- Suppose 100 tests are performed, then we should expect 100 * 0.01 = 1 false positive

- To control for this, a multiple test correction can be used. For example, to maintain 0.01 FPs given 100 tests, the cutoff can be reduced from 0.01 to 0.01/100 = 0.0001

# The p-value (in proteomics)

- Is a **local** measure, meaning that it is confined specifically to the particular PSM under consideration (it is therefore self-contained)

- **Global** measures on the other hand, considers all PSMs scores concurrently and relative to each other (they are therefore not self-contained).

- Lets say we observe a PSM with a score of 1, we can build an empirical reference distribution of similar false/random sequences and find out what are their respective PSM scores. If the observed PSM does better than at least some alpha threshold, then we can say that this PSM is statistically significant, and so we reject the null hypothesis for the alternative.

- This is computationally very intensive. ALSO… what is a reasonable null?

# False Discovery Rate (FDR)

False Discovery Rate (FDR): The expected fraction of false positives among the significant test statistics. (FP/FP+TP)
Compare this against the false positive rate which is FP/(FP+TN)

| score | type |
|-------|------|
| 7.5 | correct |
| 7.2 | correct |
| 6.9 | correct |
| 6.8 | correct |
| 6.7 | incorrect |
| 6.5 | correct |
| 6.4 | correct |
| 6.4 | correct |
| 6.3 | incorrect |
| 6.1 | correct |
| 6.0 | incorrect |
| 5.9 | correct |
| 5.7 | incorrect |
| ... | ... |

Threshold

So how do we look at this?
Let's say we have a set of PSM scores and decide to draw the line at 6, i.e., we accept all PSMs with scores > 6.
Let's also assume we have perfect knowledge of correct and wrong matches.
We note that 10 PSMs are retained.
Of these, 2 are wrong. So the FDR is therefore
$$FDR = 2/10 = 20\% = 0.2$$

This seems great. But in reality, we don't know which ones are wrong. This is similar to the null problem in p-value generation. So how do we create something which we know to be wrong or sure?

Benjamini & Hochberg, JSTOR, 1995; Storey & Tibshirani, PNAS, 2003

NANYANG TECHNOLOGICAL UNIVERSITY

False Discovery Rate (FDR)

# The target-decoy analysis

Estimating FDR:
How to purposely create your incorrect PSMs



Target
database

Decoy
database

## Target database
Protein sequences of the studied organism.
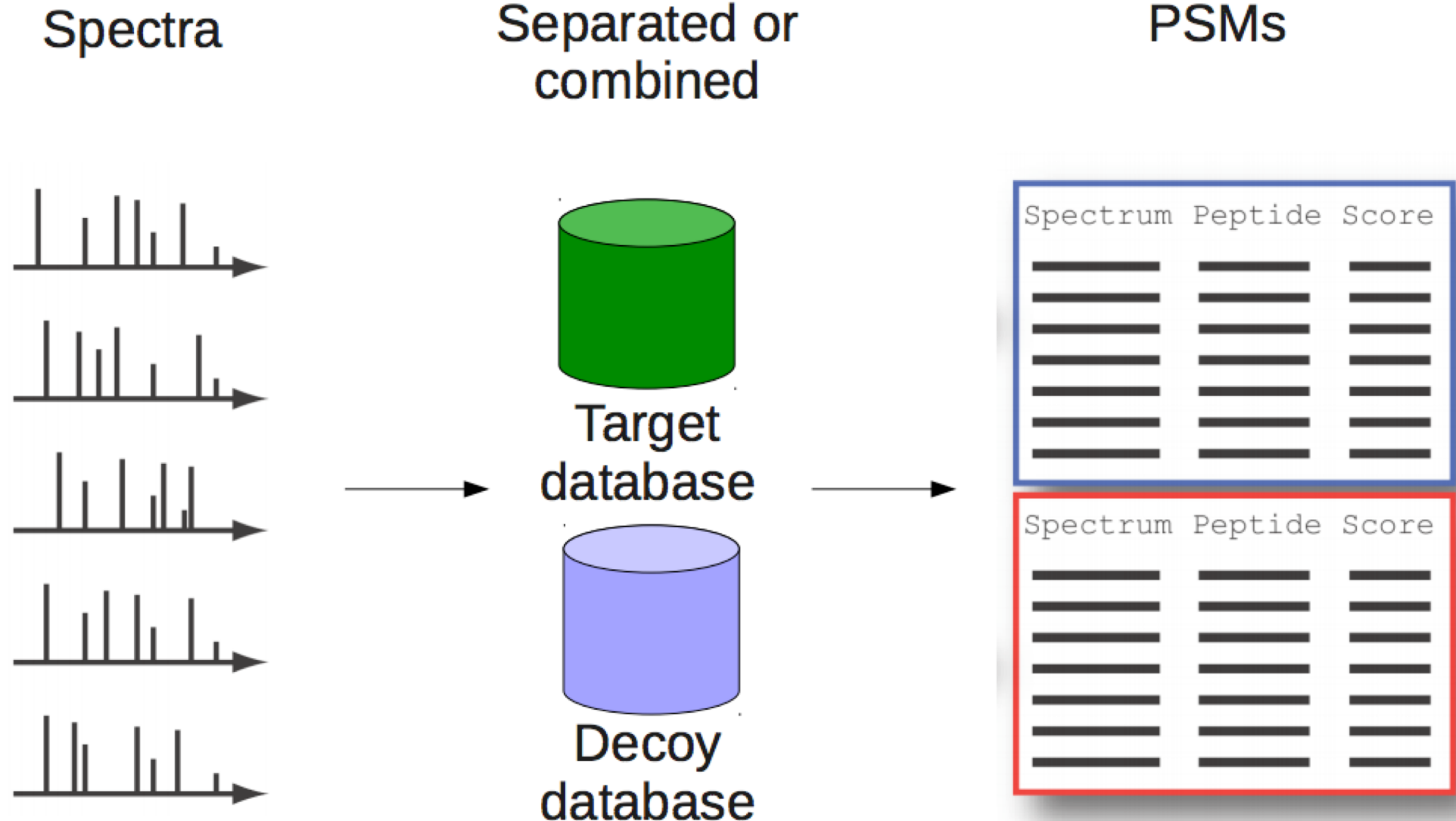
## Decoy database
Reversed or shuffled sequences.

## Assumption
Spectra matched to the decoy database are good models of **incorrect** matches to the target database.

In other words, all matches to decoy are false positives

False Discovery Rate (FDR)

# The target-decoy analysis

| Spectra | Separated or combined | PSMs |
|---|---|---|



Target database

Decoy database

| Spectrum | Peptide | Score |
|---|---|---|

| Spectrum | Peptide | Score |
|---|---|---|

Elias and Gygi. Target-Decoy Search Strategy for Mass Spectrometry-Based Proteomics. Methods Mol Biol. 2010.

NANYANG
TECHNOLOGICAL
UNIVERSITY

## Target-decoy searching steps

- Construct a concatenated target-decoy sequence list, marking decoy sequences with a text flag in their annotation.

- Use a MS/MS search engine to interpret input MS/MS spectra using target-decoy sequence list.

- Evaluate the relative proportion of target and decoy sequences in the search space to derive the multiplicative factor required to estimate false positives, if necessary.

- Estimate false positive-related statistics.

- Use decoy hits to guide the establishment of filtering criteria.

- Report statistics for filtered data set.

## Decoy construction rules

- **Similar** amino acid distributions as target protein sequences.

- **Similar** protein length distribution as target protein sequence list.

- **Similar** numbers of proteins as target protein list.

- **Similar** numbers of predicted peptides as target protein list.

- **No** predicted peptides **in common** between target and decoy sequence lists.

# False Discovery Rate (FDR)

## Reversal

**Advantages**
- Simple
- Preserve general features of the target sequence list e.g. same inter-protein redundancy
- Defined transformation therefore repeatable

**Disadvantages**
- Non-random transformation is less statistically rigorous
- Cannot be used for peptides with low sequence complexity

## Shuffling

**Advantages**
- Simple
- Has desired stochastic properties

**Disadvantages**
- Redundancies and homologies between protein entries will not be preserved, so many more decoy peptides than originally present in the target sequence list

## Random Proteins

**Advantages**
- Has desired stochastic properties
- Can preserve amino acid bias and protein length distribution

**Disadvantages**
- Redundancies and homologies between protein entries will not be preserved, so many more decoy peptides than originally present in the target sequence list

# FDR estimation based on decoy

## No decoy



$$FDR = \frac{B}{A}$$

**Combined searches**
Target and decoy database are searched together

$$\widehat{FDR} = \frac{\{\#decoys\ over\ threshold\}}{\{\#targets\ over\ threshold\}}$$

I.e., $\pi_0$ is 1
Simpler. Since estimating $\pi_0$ can be tricky.

## With decoy



decoy PSMs    target PSMs

$$FDR = \frac{B}{A} = \frac{\pi_0\ B'}{A}$$

$\pi_0$ is the fraction of incorrect target PSMs among target PSMs

Target-decoy analysis



$$FDR = \frac{B}{A} = \frac{\pi_0\ B'}{A}$$

$$\pi_0 = \frac{N_i}{N}$$

$\pi_0$ is the fraction of incorrect target PSMs among target PSMs

# Posterior Error Probability (PEP)

Posterior Error Probability (PEP): The probability that the null hypothesis is true for a particular test statistic
In proteomics, it can be taken to mean the probability that a given PSM is wrong.
PEP is sometimes called "local FDR"



$$PEP = \frac{b}{a}$$

A PEP is the probability that a PSM scoring x is incorrect

# Statistical evaluation: when to use what?

| **P-value** | **FDR/q value** | **PEP** |
|---|---|---|
| In experiments with one, or very few, spectra | For the confidence in sets, of PSMs, peptides or proteins | For the confidence in a particular PSM, peptide or proteins |



Optimal

Conservative ← Bonferroni adjusted p–value — Posterior error probability — False discovery rate or q–value — Unadjusted p–value → Anti–conservative

**Figure 3.** Methods for assigning statistical significance

Käll et al. Assigning significance to peptides identified by tandem mass spectrometry using decoy databases. JPR, 2008

NANYANG TECHNOLOGICAL UNIVERSITY

# Overview of proteo-informatics



Mass spectrometry

Search algorithm

De novo sequencing

Statistical validation of peptide and protein identifications

Quantitation

Data Storage/Dissemination

**Functional Analysis**

Data storage/ representation formats

Peptide-spectra matching (PSM)

Functional analysis

# Functional analysis 1 (Comparative analysis)

- The process of creating knowledge and insight from biological data

- Comparative analysis (group A vs B)
  - Assumption: the differences between two groups are phenotypically relevant and can be used to construct mechanistic explanation
  - This is a fallacious assumption.

# The Anna Karenina Principle

- Happy families are all alike, every unhappy family is unhappy in its own way --- Leo Tolstoy

- Interpreted as: There are many ways to violate the null hypothesis, but only one way that is pertinent to the outcome of interest

# Happiness does not come easy



Lack of $

No leisure time

No communication

Awful in-laws

......

Scandals

Incompatibility

Happiness requires positive fulfilment of all possible categories. Failure in any leads to unhappiness

Goh and Wong. Dealing with confounders in -omics analysis. Trends in Biotechnology, 2018

NANYANG
TECHNOLOGICAL
UNIVERSITY

# Anna Karenina in comparative proteomics



**False Dichotomy**
**Null: Gene does not cause disease**
**Alternative: Gene causes disease**

Wrong test construction

Batch effect

Wrong null distribution

**The gene is relevant**

Chance association

Non-causal association

Goh and Wong. Dealing with confounders in -omics analysis. Trends in Biotechnology, 2018

# Dealing with the Anna Karenina

**Causes**

- A careless null/alternative hypothesis due to forgotten assumptions:
  - Distributions of the feature of interest in the two samples are identical to the two corresponding populations
  - Features not of interest are equalized/controlled for in the two samples
  - No other explanation for the significance of the test
  - Null distribution models the real world
- These make it easy to reject the carelessly stated null hypothesis and accept an incorrect alternative hypothesis.

NANYANG
TECHNOLOGICAL
UNIVERSITY

# Dealing with the Anna Karenina

**Good Practices to Avoid Wrong Conclusions and Get Deeper Insight**

- Check for sampling bias:
  - Are the distributions of the feature of interest in the two samples same as that in the two populations?
- Check for exceptions:
  - Are there large subpopulations for which the test outcome is opposite?
  - Are there large subpopulations for which the test outcome becomes much more significant?
- Check for validity of the null distribution:
  - Is there evidence that suggests the null distribution is inappropriate?
- Check the hypothesis statement construction
  - Are the hypothesis statements being framed correctly (as opposed to a statement that is prone to being rejected for the wrong reasons)?

Goh and Wong. Dealing with confounders in -omics analysis. Trends in Biotechnology, 2018

NANYANG
TECHNOLOGICAL
UNIVERSITY

# Dealing with the Anna Karenina

**Good Practices to Avoid Wrong Conclusions and Get Deeper Insight**

- Check your assumptions
  - Are the right assumptions being made (e.g. the independence of measured variables)?

- Check if appropriate summary statistics are used
  - If an event is extremely rare, then using mean/median-based statistics will miss it; ditto if many similar events are present, but only one is relevant

- **Note:** Even if all (or most) of the above points are addressed, it still does not ensure phenotypic relevance, only correlation.

Goh and Wong. Dealing with confounders in -omics analysis. Trends in Biotechnology, 2018

NANYANG
TECHNOLOGICAL
UNIVERSITY

# Functional analysis 2 (Missing proteins)



Cannot do analysis properly with such data

# Missing proteins

- Any gene sequence whose respective protein has never been observed is an MP.

- Alongside various initiatives---e.g. GPMDB, PeptideAtlas and neXtProt --- the goal is to establish a genome-proteome bridge.

# Missing proteins (Tiers 1 to 5)

| PE Tier | Inclusion criteria | Percentage of proteome against 20,055 proteins | Notes |
|---|---|---|---|
| 1 | Evidence at proteome level | ~82.0% (16,518) | *At least 2 unique non-overlapped peptides  at least 9 amino acid residues long |
| 2 | Evidence at transcript level only | ~11.5% (2,290) | *The transcript must be confidently detected, but no corresponding protein evidence |
| 3 | Homology inference only | ~3.0% (565) | *Inferred homologues without protein or transcript support |
| 4 | Predicted | ~0.5% (94) | *Predicted coding sequence, without homology, transcript or protein support |
| 5 | Dubious | ~3.0% (588) | *The sequence may not fully meet the criteria for a predicted coding sequence<br>*Uncertainty over the veracity of the coding sequence (i.e., we do not know the sequence is correct)<br>*Some studies do not consider PE5 as MPs |

NANYANG TECHNOLOGICAL UNIVERSITY

# Missing proteins and relations to coverage, consistency problems

- An MP may be one of the following
  - sequence is known but hard to detect,
  - sequence is known but never detected in MS
  - sequence is not known but evidence exists for it e.g. via gene prediction or in raw spectra.
- The "missing-protein problem" (MPP)---viz. the difficulty in detecting certain proteins despite transcript or theoretical evidence---should more rightfully be considered a narrow manifestation of the more general coverage (the inability to survey the entire proteome) and consistency (the inability to consistently detect a protein) problems

# Coverage and consistency

# Why do proteins go missing?



**Biological**

**Low Abundance**
The MP's spectra have low intensities and barely distinguishable from background

**Unknown variants**
Known sequence
Unknown mods
Unknown variant
The MP has unknown splice forms or PTMs

**Sequence ambiguity**
MP
P1
P2
The MP lacks unique sequences

**Technical**

**Low-Instrument Resolution**
Detected    Undetected
Low-resolution instruments cannot capture all peptide signals concurrently

**Instrument Bias**
MS1    MS2
Selected
Missed
Random precursor selection

**Cross-interference**
The MP's spectra are often mixed up with other proteins
Detected    Masked

**Informatics**

**Large reference libraries**
Library
statistical cutoff
True match    False match
Cutoff result in MPs

**Not found in library**
Library    ?
The MP sequence is absent

**Parameterization**
Match tolerance setting
Detected
Theoretical
Search settings unsuitable

Zhou et al. Understanding missing proteins: A functional perspective. Drug Discovery Today, 2018

**NANYANG TECHNOLOGICAL UNIVERSITY**

# Missing value imputation (MVI)

- A few strategies:
  - We fill in 0s or a fixed value based on the average of all protein expression
  - For each missing value per protein, we fill in the average value based on all observed values for that same protein
  - We estimate the missing value based on proteins that are known to be correlated

# Missing value imputation (MVI)



Limitation: It mostly only works well for inconsistency issue

# MVI methods really don't work very well



**Figure 2.**
Boxplot of the average $\log_{10}$ CV(RMSE) for the imputed dilution series datasets (Table 1) at the (A) peptide and (B) protein levels. The lower line represents the 25th percentile, the upper line of the box represents the 75th percentile, and the inner line corresponds to the median $\log_{10}$ CV(RMSE).

NANYANG TECHNOLOGICAL UNIVERSITY

**High abundance has lower % of MPs. However, low abundance is not a solely explanation. The MPs are equally distributed across the horizontal median.**



**Figure 1.**
Average $\log_{10}$ intensity as measured by peptide peak area in the control group versus fraction of missing values and peptide counts associated with bins corresponding to the fraction of missing data comparing phenotypes and exposures for datasets from (A) human plasma and (B) mouse lung. The control group for the human plasma is the normal glucose tolerant (NGT) samples, and the sham group for the mouse lung is the regular weight mice with no lipopolysaccharide (LPS) exposure. The vertical red line represents median average intensity, and the horizontal red line represents the point that 50% of the values are missing.

TECHNOLOGICAL UNIVERSITY

# How about we use the idea of "guilt-by-association?"

- **Postulate**: The chance of a protein complex being present in a sample is proportional to the fraction of its constituent proteins being correctly reported in the sample

- Suppose proteomics screen has 75% reliability; a complex comprises proteins A, B, C, D, E; and screen reports A, B, C, D only but not E.

$\Rightarrow$ Complex has 60% (= 0.75 * 4 / 5) chance to be present

$\Rightarrow$ The unreported protein E also has $\geq$ 60% chance to be present, as presence of the complex implies presence of all its constituents

$\Rightarrow$ **improving coverage (recover missing proteins)**

$\Rightarrow$ Each of the reported proteins (A, B, C, and D) individually has 90% (= 100% * 0.6 + 75% * 0.4) chance of being true positive, whereas a reported protein that is isolated has a lower 75% chance of being true positive

$\Rightarrow$ **removing noise**

Goh and Wong. Integrating networks and proteomics: moving forward. Trends in Biotechnology, 2016

Goh and Wong. Design principles for clinical network-based proteomics. Drug Discovery Today, 2016

NANYANG
TECHNOLOGICAL
UNIVERSITY

# How about we use the idea of "guilt-by-association?"

## The functional class scoring (FCS) algorithm



Goh et al. Comparative network-based recovery analysis and proteomic profiling of neurological changes in valporic acid-treated mice. *JPR*, 2013

NANYANG
TECHNOLOGICAL
UNIVERSITY

# Does context really work?

| Method | Novel Suggested Proteins | Recovered proteins | Recall | Precision |
|---|---|---|---|---|
| PEP | 1037 | 158 | 0.317 | 0.152 |
| Maxlink | 822 | 226 | 0.454 | 0.275 |
| FCS (predicted) | 638 | 224 | 0.450 | 0.351 |
| FCS (complexes) | 895 | 477 | 0.958 | 0.533 |

- Looks like running FCS on real complexes is able to recover more proteins and more accurately

But we can't rank the individual proteins simply based on p-values. Can we do better? This is a story for another time. Or simply refer to https://www.comp.nus.edu.sg/~wongls/talks/wls-incob2017.pdf

NANYANG TECHNOLOGICAL UNIVERSITY

Goh et al. Comparative network-based recovery analysis and proteomic profiling of neurological changes in valporic acid-treated mice. *JPR*, 2013

# What have we learnt?

- Getting good quality PSMs requires consideration of a large number of factors
- The p-value, FDR and PEP are used as statistical approaches for different purposes
- There are 3 strategies for creating decoy libraries in FDR estimation
- Proteomics is plagued with coverage and consistency issues, requiring various rescue analysis

## You should be able to

- Describe the various factors affecting PSM quality

- Describe p-values, FDR and PEP

- Describe and evaluate the various decoy library generation strategies (sequence reversal, sequence randomization) for FDR estimation

- Describe coverage and consistency issues in proteomics

# Readings

- Elias and Gygi. Target-Decoy Search Strategy for Mass Spectrometry-Based Proteomics. Methods Mol Biol. 604: 55–71, 2010.

- Goh et al. Comparative network-based recovery analysis and proteomic profiling of neurological changes in valporic acid-treated mice. JPR. 12 (5), 2116-2127, 2013

- Goh and Wong. Advanced bioinformatics methods for practical applications of proteomics. Briefings in Bioinformatics, 2017 (https://doi.org/10.1093/bib/bbx128)

- Käll et al. Assigning significance to peptides identified by tandem mass spectrometry using decoy databases. JPR. 7(1):29-34, 2008

# Readings

- Goh and Wong. Integrating networks and proteomics: moving forward. Trends in Biotechnology, 34(12):951-959, 2016

- Goh and Wong. Design principles for clinical network-based proteomics. Drug Discovery Today, 21(7):1130-1138, 2016

- Goh and Wong. Dealing with confounders in omics analysis. Trends in Biotechnology, S0167-7799(18)30047-7, 2018

- Goh and Wong. Advanced bioinformatics methods for practical applications in proteomics. Briefings in Bioinformatics, 2017.

- Zhou et al. Understanding missing proteins: A functional perspective. Drug Discovery Today, 23(3):644--651, March 2018.

NANYANG
TECHNOLOGICAL
UNIVERSITY