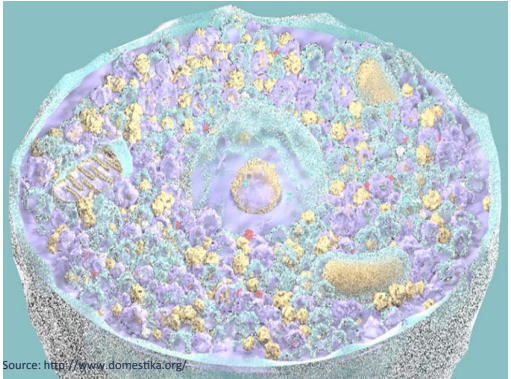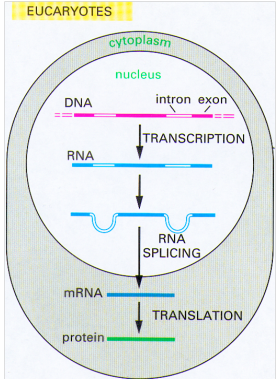# Statistical Methods for Quantitative MS-Based Proteomics:
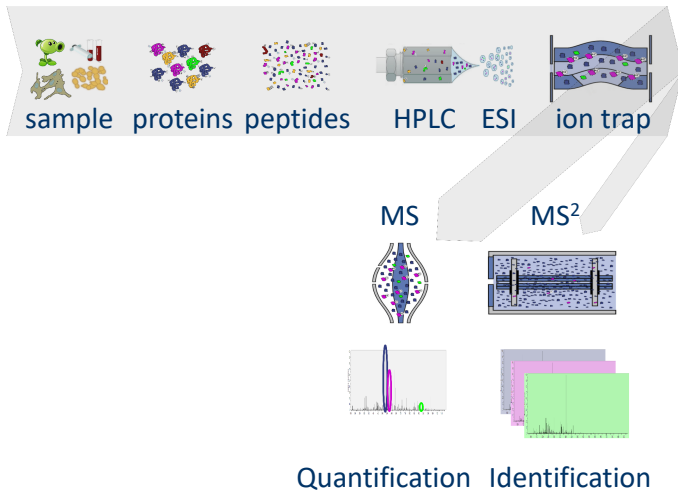## 1. Identification & False discovery rate

Lieven Clement

Proteomics Data Analysis Shortcourse
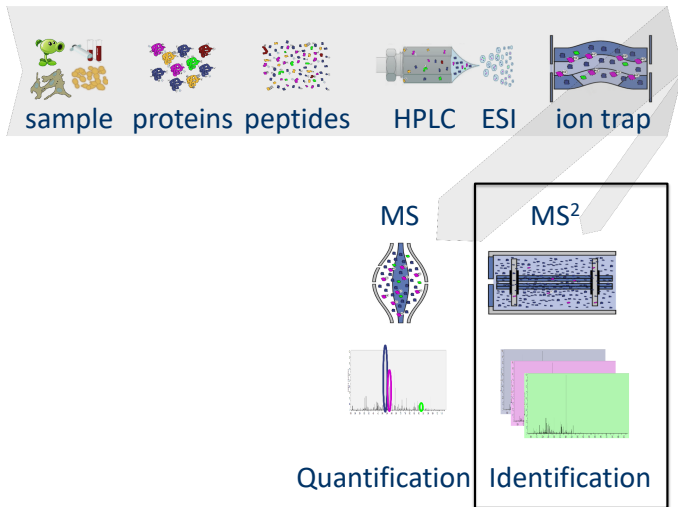
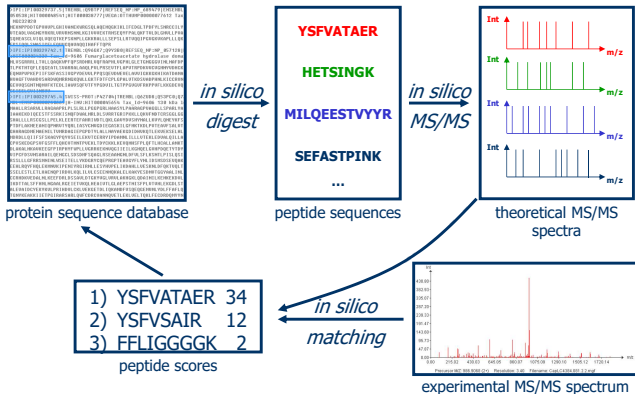# Challenges in Label Free MS-based Quantitative Proteomics



sample  proteins  peptides     HPLC   ESI     ion trap

MS              MS$^2$

Quantification  Identification

# Challenges in Label Free MS-based Quantitative Proteomics

# Identification



protein sequence database

peptide sequences

theoretical MS/MS spectra

*in silico* digest

*in silico* MS/MS

YSFVATAER

HETSINGK

MILQEESTVYYR

SEFASTPINK

...

1) YSFVATAER   34
2) YSFVSAIR    12
3) FFLIGGGGK    2

peptide scores

*in silico* matching

experimental MS/MS spectrum

(slide courtesy to Lennart Martens)

# E-values

Probability that a random candidate peptide produces a higher score that the observed PSM score.



protein sequence database

peptide sequences

theoretical MS/MS spectra

experimental MS/MS spectrum

1) YSFVATAER  34
2) YSFVSAIR  12
3) FFLIGGGGK  2

peptide scores

*in silico digest*

*in silico MS/MS*

*in silico matching*

YSFVATAER
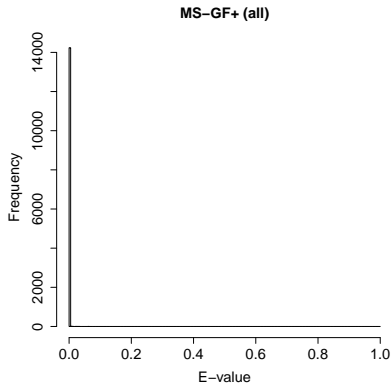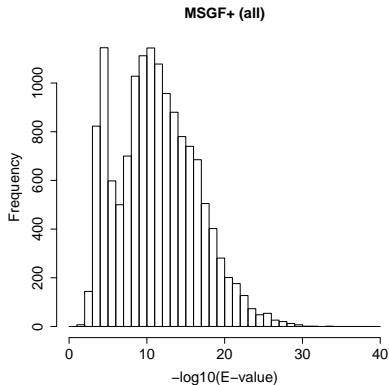
HETSINGK

MILQEESTVYYR

SEFASTPINK

...

# E-values

Probability that a random candidate peptide produces a higher score that the observed PSM score.



OMSSA all PSM

# E-values

Probability that a random candidate peptide produces a higher score that the observed PSM score.



E−values we expect for random candidate peptides

# E-values

Probability that a random candidate peptide produces a higher score that the observed PSM score.



**OMSSA decoy PSMs**

# E-values

Probability that a random candidate peptide produces a higher score that the observed PSM score.

# E-values

Probability that a random candidate peptide produces a higher score that the observed PSM score.



MSGF+ (all)

# E-values

Probability that a random candidate peptide produces a higher score that the observed PSM score.

- A bad hit is the random hit with the best score so it is also bound to have a low E-value.

# E-values

Probability that a random candidate peptide produces a higher score that the observed PSM score.

- A bad hit is the random hit with the best score so it is also bound to have a low E-value.
- If we look at E-values for all PSMs they are only useful as a score.

# E-values

Probability that a random candidate peptide produces a higher score that the observed PSM score.

- A bad hit is the random hit with the best score so it is also bound to have a low E-value.
- If we look at E-values for all PSMs they are only useful as a score.
- We should know the distribution of the maximum score of random candidate peptides when we want to do the statistics.

# Table of Outcomes

|  | Called Bad | Called Correct |  |
|---|---|---|---|
| Bad hit | TN | FP | $m_0$ |
| Correct hit | FN | TP | $m_1$ |
| Total | NR | R | $m$ |

- TN: number of true negatives
- FP: number of false positives
- FN: number of false negatives
- TP: number of true positives
- NR: number of non-rejections, R: number of rejections

# Table of Outcomes

|  |  | Called Bad | Called Correct |  |
|---|---|:---:|:---:|:---:|
|  | Bad hit | TN | FP | $m_0$ |
| Unobservable |  |  |  |  |
|  | Correct hit | FN | TP | $m_1$ |
| Observable | Total | NR | R | $m$ |

$FDP = \frac{FP}{FP+TP}$. But is unkown! (FDP: false discovery proportion)

# Table of Outcomes

|  |  | Called Bad | Called Correct |  |
|---|---|---|---|---|
|  | Bad hit | TN | FP | $m_0$ |
| Unobservable |  |  |  |  |
|  | Correct hit | FN | TP | $m_1$ |
| Observable | Total | NR | R | $m$ |

$FDR = E\left[\frac{FP}{FP+TP}\right]$. (FDR: false discovery rate)

# Search engines return score that discriminates good from bad matches



Pyrococcus Search

# Search engines return score that discriminates good from bad matches

Score threshold $t$?

# Search engines return score that discriminates good from bad matches



Score threshold $t$?

$$f(x) = \pi_0 f_0(x) + (1 - \pi_0) f_1(x)$$

# Search engines return score that discriminates good from bad matches



Score threshold $t$?

$$f(x) = \pi_0 f_0(x) + (1 - \pi_0) f_1(x)$$

$$\text{FDR}(t) = E\left[\frac{FP}{FP + TP}\right]$$

# Search engines return score that discriminates good from bad matches



**Pyrococcus Search**

Score threshold $t$?

$$f(x) = \pi_0 f_0(x) + (1 - \pi_0) f_1(x)$$

$$\mathrm{FDR}(t) = E\left[\frac{FP}{FP + TP}\right]$$

$$\mathrm{FDR}(t) = \frac{m_0 P[x \geq t | FP]}{m P[x \geq t]}$$

$$= \frac{m P[FP] P[x \geq t | FP]}{m P[x \geq t]}$$

# Search engines return score that discriminates good from bad matches



**Pyrococcus Search**

Score threshold $t$?

$$f(x) = \pi_0 f_0(x) + (1 - \pi_0) f_1(x)$$

$$\text{FDR}(t) = E\left[\frac{FP}{FP + TP}\right]$$

$$\text{FDR}(t) = \frac{m_0 P[x \geq t | FP]}{m P[x \geq t]}$$

$$= \frac{m P[FP] P[x \geq t | FP]}{m P[x \geq t]}$$

$$\text{FDR}(t) = \frac{\pi_0 P_0[x \geq t]}{P[x \geq t]}$$

# Search engines return score that discriminates good from bad matches



**Pyrococcus Search**

Score threshold $t$?

$$f(x) = \pi_0 f_0(x) + (1 - \pi_0)f_1(x)$$

$$\text{FDR}(t) = E\left[\frac{FP}{FP+TP}\right]$$

$$\text{FDR}(t) = \frac{m_0 P[x \geq t | FP]}{m P[x \geq t]}$$

$$= \frac{m P[FP] P[x \geq t | FP]}{m P[x \geq t]}$$

$$\text{FDR}(t) = \frac{\pi_0 P_0[x \geq t]}{P[x \geq t]}$$

$$P[x \geq t] = \int\limits_{x=t}^{+\infty} f(x)dx$$

# Search engines return score that discriminates good from bad matches



**Pyrococcus Search**

Score threshold $t$?

$$f(x) = \pi_0 f_0(x) + (1 - \pi_0) f_1(x)$$

$$\text{FDR}(t) = E\left[\frac{FP}{FP + TP}\right]$$

$$\text{FDR}(t) = \frac{m_0 P[x \geq t | FP]}{m P[x \geq t]}$$

$$= \frac{m P[FP] P[x \geq t | FP]}{m P[x \geq t]}$$

$$\text{FDR}(t) = \frac{\pi_0 P_0[x \geq t]}{P[x \geq t]}$$

FDR is a set property: $FDR(t) = \dfrac{\pi_0 \int\limits_{x=t}^{+\infty} f_0(x)dx}{\int\limits_{x=t}^{+\infty} f(x)dx}$

# Search engines return score that discriminates good from bad matches



**Pyrococcus Search**

Score threshold $t$?

$$f(x) = \pi_0 f_0(x) + (1 - \pi_0) f_1(x)$$

$$\text{FDR}(t) = E\left[\frac{FP}{FP+TP}\right]$$

$$\text{FDR}(t) = \frac{m_0 P[x \geq t | FP]}{m P[x \geq t]}$$

$$= \frac{m P[FP] P[x \geq t | FP]}{m P[x \geq t]}$$

$$\text{FDR}(t) = \frac{\pi_0 P_0[x \geq t]}{P[x \geq t]}$$

local fdr (posterior error probability, PEP): $fdr(x) = \frac{\pi_0 f_0(x)}{f(x)}$

Probability that an individual PSM is a bad hit.

# How to estimate FDR?



$$\mathrm{FDR}(t) = E\left[\frac{FP}{FP+TP}\right]$$

$$= \frac{\pi_0 P_0[x \geq t]}{P[x \geq t]}$$

$$P_.[x \geq t] = \int\limits_{t}^{\infty} f_.(x)\,dx$$

# How to estimate FDR?



$$\mathrm{FDR}(t) \;=\; E\left[\frac{FP}{FP+TP}\right]$$

$$\;=\; \frac{\pi_0 P_0[x \ge t]}{P[x \ge t]}$$

$$P.[x \ge t] \;=\; \int\limits_{t}^{\infty} f.(x)\,dx$$

$$\hat{P}[x \ge t] = \frac{\#x \ge t}{m} \qquad \Rightarrow \qquad \widehat{\mathrm{FDR}}(t) = \frac{\pi_0 P_0[x \ge t]}{\frac{\#x \ge t}{m}}$$

# How to estimate FDR?



$$\text{FDR}(t) = E\left[\frac{FP}{FP+TP}\right]$$

$$= \frac{\pi_0 P_0[x \geq t]}{P[x \geq t]}$$

$$P.[x \geq t] = \int_t^\infty f.(x)dx$$

$$\hat{P}[x \geq t] = \frac{\#x \geq t}{m} \qquad \Rightarrow \qquad \widehat{\text{FDR}}(t) = \frac{\pi_0 P_0[x \geq t]}{\frac{\#x \geq t}{m}}$$

How to characterize $f_0(t)$ and $\pi_0$ in proteomics?

# Target-Decoy approach to establish null distribution



- Search against decoy database to generate representative bad hits
- Reversed databases are popular

# Target-Decoy approach to establish null distribution



- Search against decoy database to generate representative bad hits
- Reversed databases are popular
- Concatenated search

# Target-Decoy approach to establish null distribution



Pyrococcus Concatenated Search: Targets

- Search against decoy database to generate representative bad hits
- Reversed databases are popular
- Concatenated search

# Target-Decoy approach to establish null distribution



**Pyrococcus Concatenated Search: Targets**

- Search against decoy database to generate representative bad hits
- Reversed databases are popular
- Concatenated search
- Assumption: bad hits has equal probability to map on target and decoy

$$\hat{\pi}_0 = \frac{\#decoys}{\#targets}$$

# Target-Decoy approach to establish null distribution



**Pyrococcus Concatenated Search: Targets**

- Search against decoy database to generate representative bad hits
- Reversed databases are popular
- Concatenated search
- Assumption: bad hits has equal probability to map on target and decoy

$$\hat{\pi}_0 = \frac{\#decoys}{\#targets}$$

- Score cuttoff:
$$\text{FDR}(x) = E\left[\frac{FP}{FP+TP}\right]$$

# Target-Decoy approach to establish null distribution



Pyrococcus Concatenated Search: Targets

- Competitive Target - decoy:

$$\widehat{\mathrm{FDR}}(x) = \frac{\#decoys | X \geq x}{\#targets | X \geq x}$$

# Target-Decoy approach to establish null distribution



Pyrococcus Concatenated Search: Targets

- Competitive Target - decoy:

$$\widehat{\text{FDR}}(x) = \frac{\#decoys | X \geq x}{\#targets | X \geq x}$$

$$\widehat{\text{FDR}}(x) = \frac{\#decoys}{\#targets} \frac{\frac{\#decoys | X \geq x}{\#decoys}}{\frac{\#targets | X \geq x}{\#targets}}$$

$$\widehat{\text{FDR}}(x) = \hat{\pi}_0 \frac{\overset{+\infty}{\underset{t}{\int}} \widehat{f_0(x)dx}}{\overset{+\infty}{\underset{t}{\int}} \widehat{f(x)dx}}$$

$$\widehat{\text{FDR}}(x) = \hat{\pi}_0 \frac{\hat{P}_0[X \geq x]}{\hat{P}[X \geq x]}$$

# Target-Decoy approach to establish null distribution



- Competitive Target - decoy:

$$\widehat{FDR}(x) = \frac{\#decoys | X \geq x}{\#targets | X \geq x}$$

$$\widehat{FDR}(x) = \frac{\#decoys}{\#targets} \frac{\frac{\#decoys | X \geq x}{\#decoys}}{\frac{\#targets | X \geq x}{\#targets}}$$

$$\widehat{FDR}(x) = \hat{\pi}_0 \frac{\overbrace{\int_t^{+\infty} f_0(x)dx}}{\int_t^{+\infty} f(x)dx}$$

$$\widehat{FDR}(x) = \hat{\pi}_0 \frac{\hat{P}_0[X \geq x]}{\hat{P}[X \geq x]}$$

# Target-Decoy approach to establish null distribution



**Pyrococcus Concatenated Search: Targets**
**11160 PSMs at FDR 0.01**

Legend: Target, $f_0(x)$: Decoy, $f_1(x)$

x-axis: MS–GF+ Score
y-axis: # peptides

- Competitive Target - decoy:

$$\widehat{FDR}(x) = \frac{\#decoys|X \geq x}{\#targets|X \geq x}$$

$$\widehat{FDR}(x) = \frac{\#decoys}{\#targets} \frac{\frac{\#decoys|X \geq x}{\#decoys}}{\frac{\#targets|X \geq x}{\#targets}}$$

$$\widehat{FDR}(x) = \hat{\pi}_0 \frac{\overbrace{\int_t^{+\infty} f_0(x)dx}}{\underbrace{\int_t^{+\infty} f(x)dx}}$$

$$\widehat{FDR}(x) = \hat{\pi}_0 \frac{\hat{P}_0[X \geq x]}{\hat{P}[X \geq x]}$$

# Assess TDA assumptions

We have to evaluate that

- The decoys are good simulations of the bad target hits: compare distributions $F_D(x)$ with $F(x)$

$$F_D(x) = \int_{-\infty}^{t} f_D(x)dx \quad \leftrightarrow \quad F(x) = \int_{-\infty}^{t} f(x)dx$$

- $\hat{\pi}_0 = \frac{\#decoys}{\#targets}$ is a good estimator for $\pi_0$.

- We will use Probability-Probability-plots (PP-plot) for this purpose.

- To make PP-plots we need estimates for $F_D(x)$ and $F(x)$.
- The empirical cumulative distribution (ECDF) is used for that purpose



$$\hat{F}_D(x) = \frac{\#decoys | X \leq x}{\#decoys}, \quad \hat{F}(x) = \frac{\#targets | X \leq x}{\#targets}$$

# PP-plot

PP-plots have the property that they show a straight 45 degree line through the origin if and only if both distributions are equivalent.

# PP-plot

PP-plots have the property that they show a straight 45 degree line through the origin if and only if both distributions are equivalent.

# PP-plot

PP-plots have the property that they show a straight 45 degree line through the origin if and only if both distributions are equivalent.

# PP-plot

PP-plots have the property that they show a straight 45 degree line through the origin if and only if both distributions are equivalent.

# PP-plot

PP-plots have the property that they show a straight 45 degree line through the origin if and only if both distributions are equivalent.

# PP-plot

PP-plots have the property that they show a straight 45 degree line through the origin if and only if both distributions are equivalent.

# PP-plot

PP-plots have the property that they show a straight 45 degree line through the origin if and only if both distributions are equivalent.

# PP-plot

PP-plots have the property that they show a straight 45 degree line through the origin if and only if both distributions are equivalent.

# PP-plot

PP-plots have the property that they show a straight 45 degree line through the origin if and only if both distributions are equivalent.

# PP-plot

PP-plots have the property that they show a straight 45 degree line through the origin if and only if both distributions are equivalent.

# PP-plot

PP-plots have the property that they show a straight 45 degree line through the origin if and only if both distributions are equivalent.

# PP-plot

# PP-plot

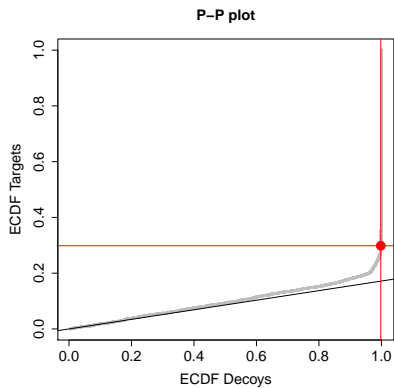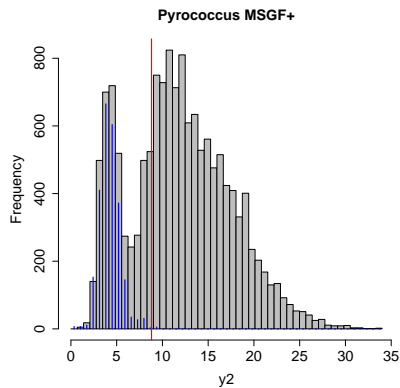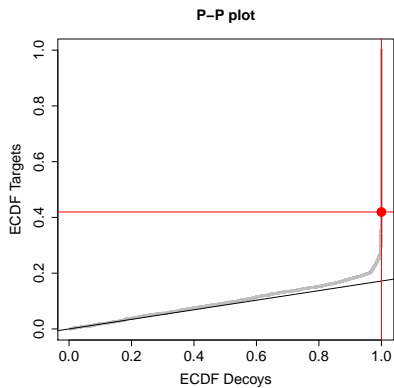# PP-plot

# PP-plot

# PP-plot
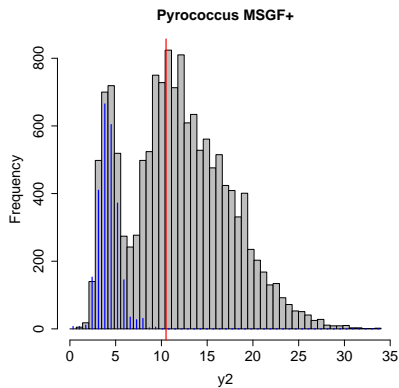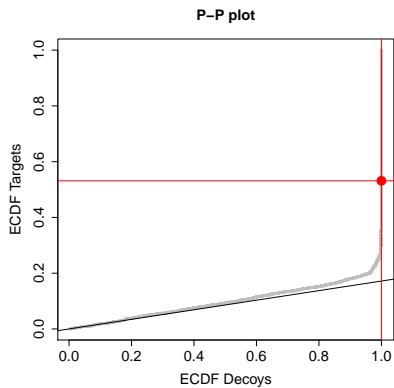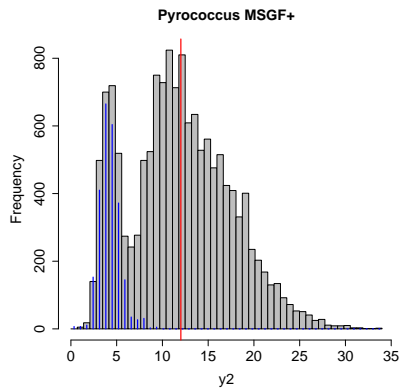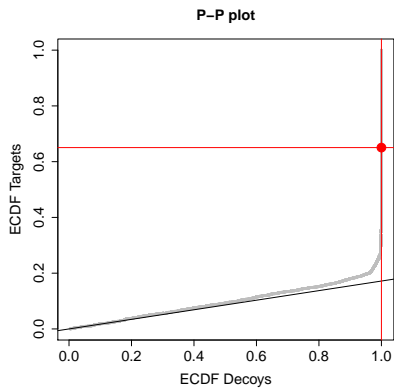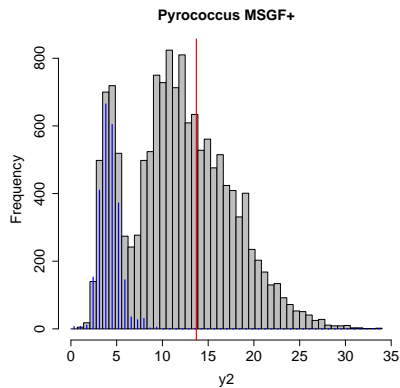
# PP-plot

# PP-plot

# PP-plot

# PP-plot

# PP-plot

# PP-plot

# PP-plot: pyrococcus
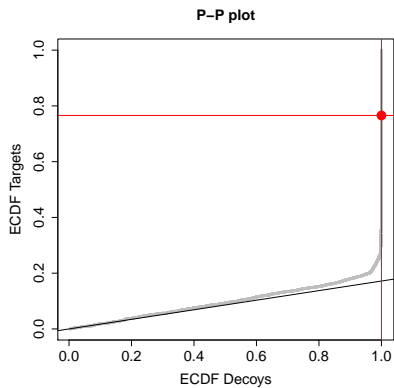
# PP-plot: pyrococcus
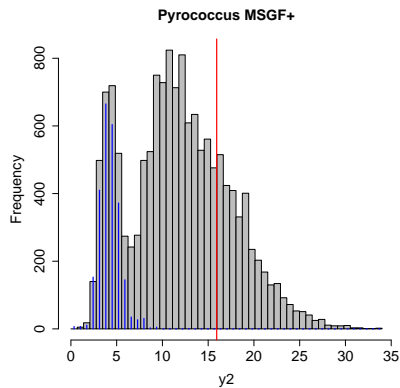


What about $\hat{\pi}_0$?

# PP-plot: pyrococcus

# PP-plot: pyrococcus

# PP-plot: pyrococcus

# PP-plot: pyrococcus

# PP-plot: pyrococcus

# PP-plot: pyrococcus

# PP-plot: pyrococcus

# PP-plot: pyrococcus

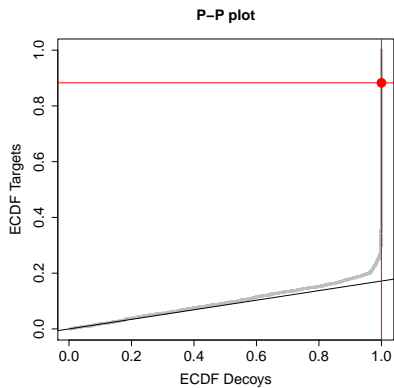# PP-plot: pyrococcus

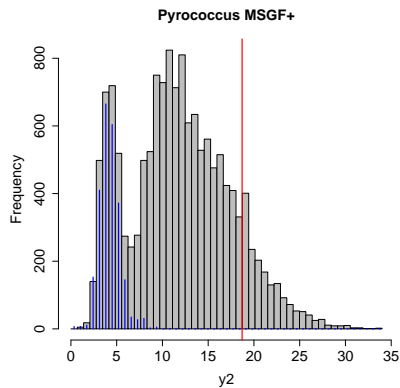# PP-plot: pyrococcus

# PP-plot: pyrococcus

# PP-plot: pyrococcus