

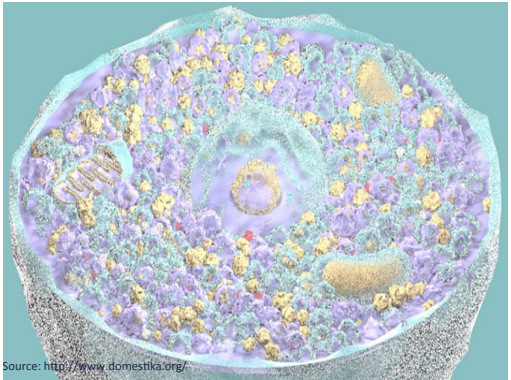
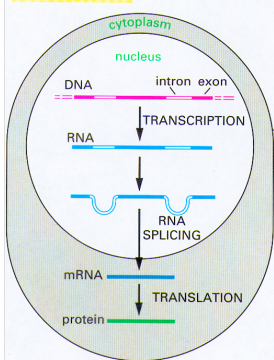
# Statistical Methods for Quantitative MS-Based Proteomics:

## 1. Identification & False discovery rate

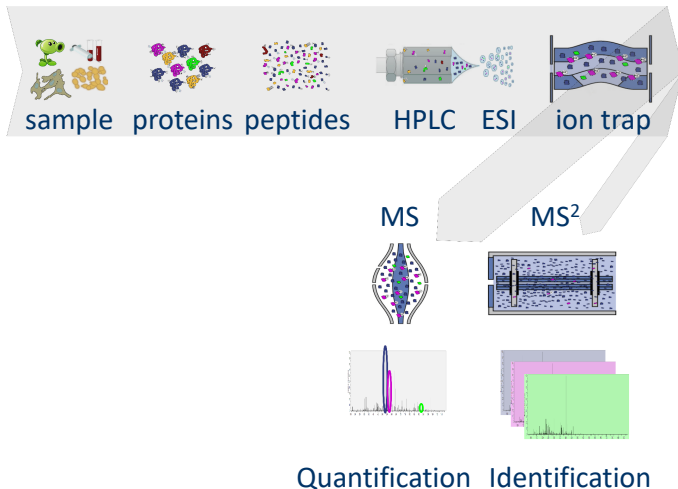
Lieven Clement

Proteomics Data Analysis Shortcourse

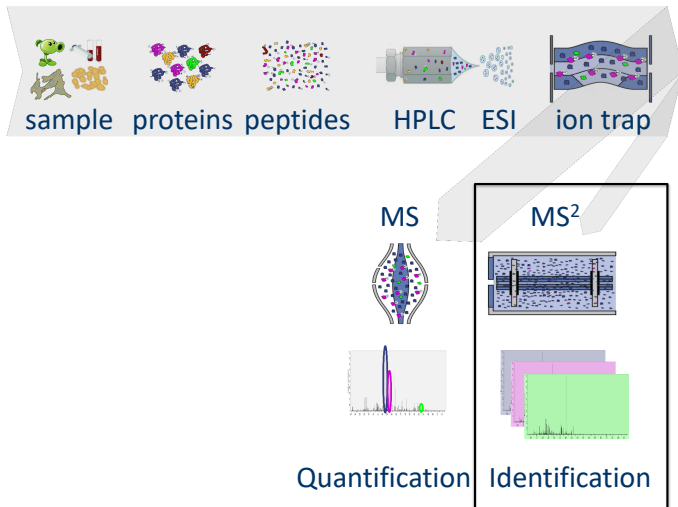
## EUCARYOTES



# Challenges in Label Free MS-based Quantitative Proteomics



# Challenges in Label Free MS-based Quantitative Proteomics



# Identification

```

>FP1:FP10029727.5|TREMBL|Q9B1P7|RET SEQ_MP_TIP_089479|ENHNL
060528|S1|0000045-1|S1|0000077|VEGA:01700P|000007|762 Tax
M022828
<...>
>FP1:FP10029740.1|TREMBL|Q9B1P7|RET SEQ_MP_TIP_089479|ENHNL
060528|S1|0000045-1|S1|0000077|VEGA:01700P|000007|762 Tax
M022828
<...>
>FP1:FP10029745.4|SWISS-PROT|P40784|TREMBL|Q42098|D56P01|Q
000000000000000000|b-102|S1|0000045 Tax_10-0606 T30 kDa 1
<...>

```

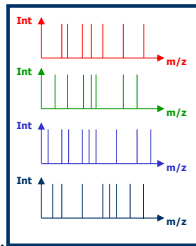
protein sequence database

*in silico*  
digest

**YSFVATAER**  
**HETSINGK**  
**MILQEESTVYR**  
**SEFASTPINK**  
...

peptide sequences

*in silico*  
MS/MS

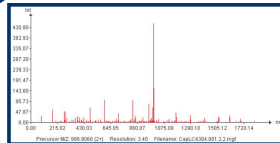


theoretical MS/MS spectra

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

peptide scores

*in silico*  
matching

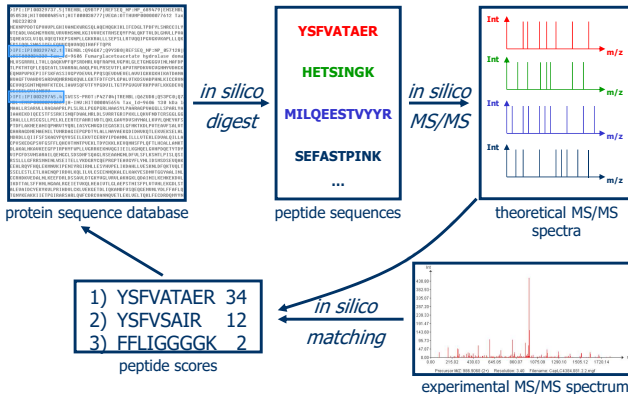


experimental MS/MS spectrum

(slide courtesy to Lennart Martens)

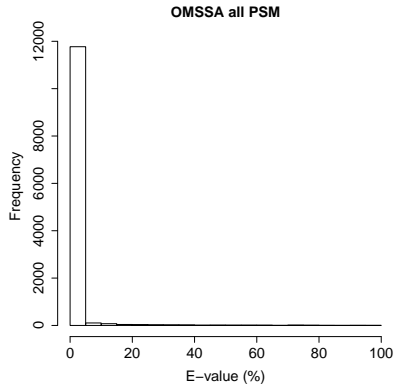
# E-values

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



## E-values

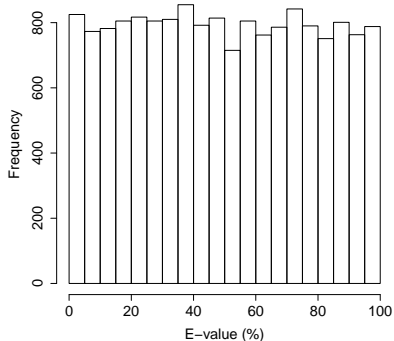
Probability that a random candidate peptide produces a higher score than the observed PSM score.



## E-values

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

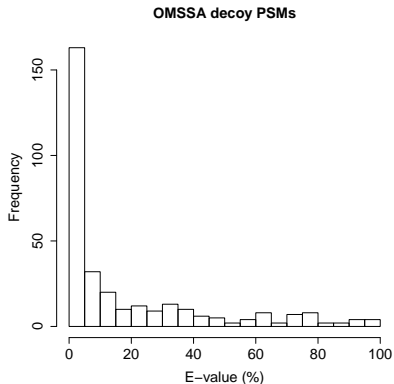
E-values we expect for random candidate peptides





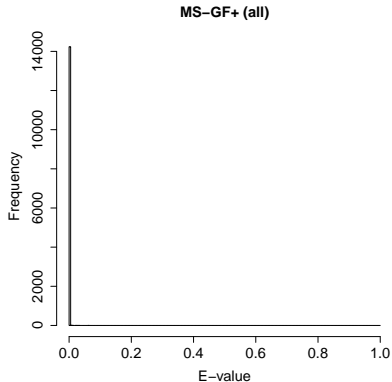
## E-values

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



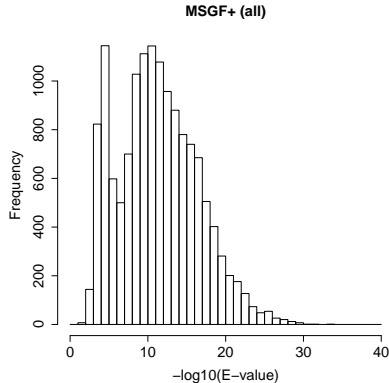
## E-values

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



## E-values

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



## E-values

Probability that a random candidate peptide produces a higher score than 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 than 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 than 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	
Unobservable	Bad hit	TN	FP	$m_0$
	Correct hit	FN	TP	$m_1$
Observable	Total	NR	R	$m$

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

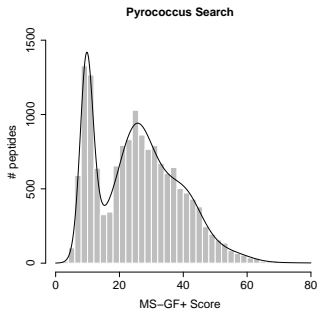


## Table of Outcomes

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

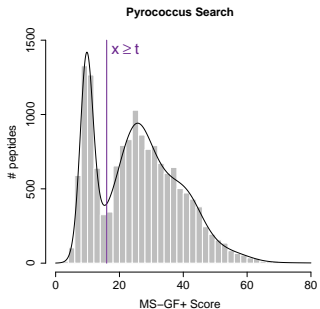
$$FDR = E \left[ \frac{FP}{FP+TP} \right]. \text{ (FDR: false discovery rate)}$$

# Search engines return score that discriminates good from bad matches



# Search engines return score that discriminates good from bad matches

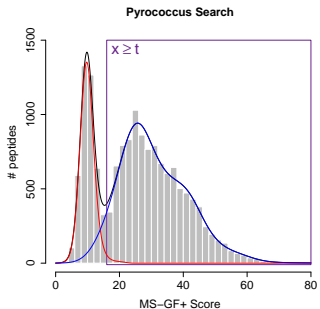
Score threshold  $t$ ?



# Search engines return score that discriminates good from bad matches

Score threshold  $t$ ?

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

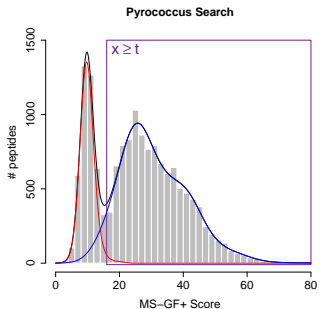


# Search engines return score that discriminates good from bad matches

Score threshold  $t$ ?

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

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



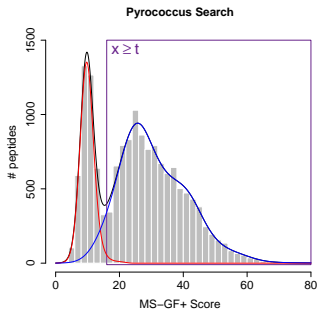
# Search engines return score that discriminates good from bad matches

Score threshold  $t$ ?

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

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

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



# Search engines return score that discriminates good from bad matches

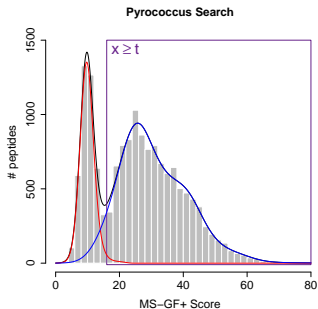
Score threshold  $t$ ?

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

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

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

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



# Search engines return score that discriminates good from bad matches

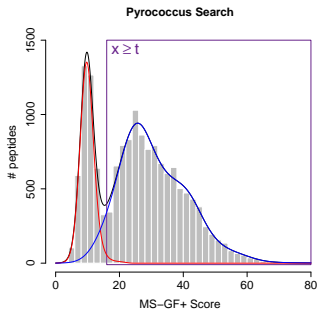
Score threshold  $t$ ?

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

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

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

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



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



# Search engines return score that discriminates good from bad matches

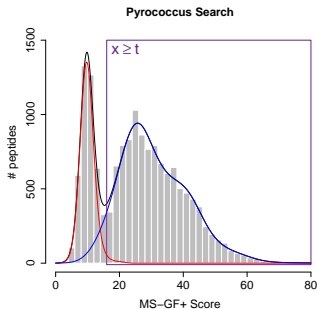
Score threshold  $t$ ?

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

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

$$\begin{aligned} FDR(t) &= \frac{m_0 P[x \geq t | FP]}{m P[x \geq t]} \\ &= \frac{m P[FP] P[x \geq t | FP]}{m P[x \geq t]} \end{aligned}$$

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



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

# Search engines return score that discriminates good from bad matches

Score threshold  $t$ ?

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

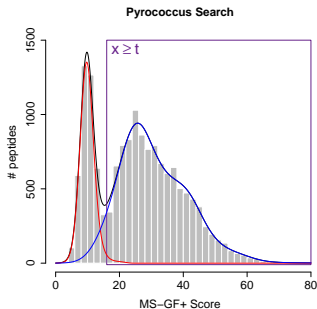
$$\text{FDR}(t) = E \left[ \frac{FP}{FP+TP} \right]$$

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

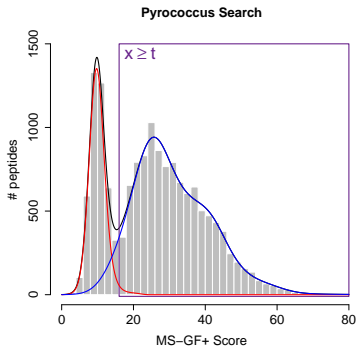
$$\text{FDR}(t) = \frac{\pi_0 P_0[x \geq t]}{P[x \geq t]}$$

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?

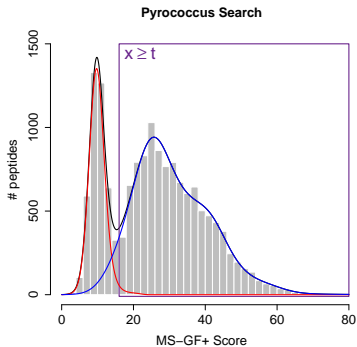


$$\text{FDR}(t) = E \left[ \frac{FP}{FP+TP} \right]$$

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

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

# How to estimate FDR?



$$\hat{P}[x \geq t] = \frac{\#x \geq t}{m} \Rightarrow$$

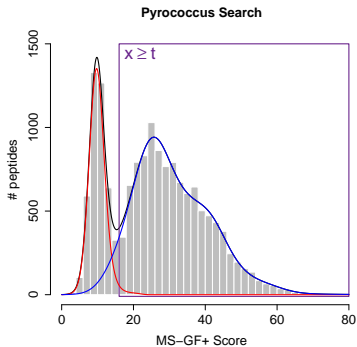
$$\text{FDR}(t) = E \left[ \frac{FP}{FP+TP} \right]$$

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

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

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

# How to estimate FDR?



$$\hat{P}[x \geq t] = \frac{\#x \geq t}{m} \Rightarrow$$

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

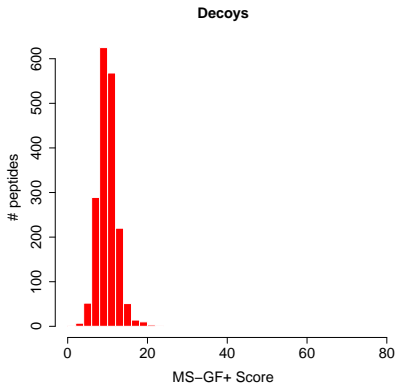
$$= \frac{\pi_0 P_0[x \geq t]}{P[x \geq t]}$$

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

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

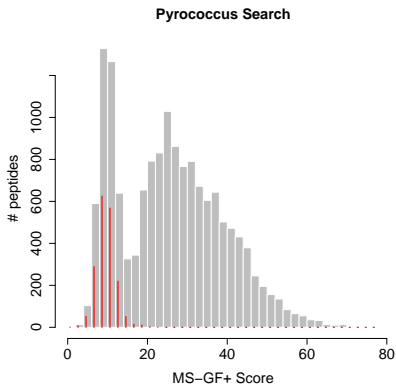
How to characterize  $f_0(t)$  and  $\pi_0$  in proteomics?

# Target-Decoy approach to establish null distribution



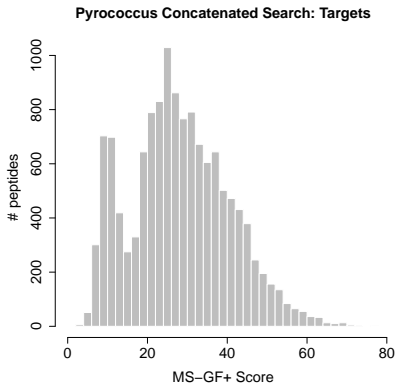
- Search against decoy database to generate representative bad hits
- Reversed databases are popular

# Target-Decoy approach to establish null distribution



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

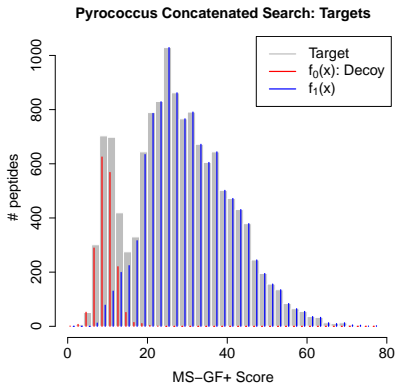
# Target-Decoy approach to establish null distribution



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



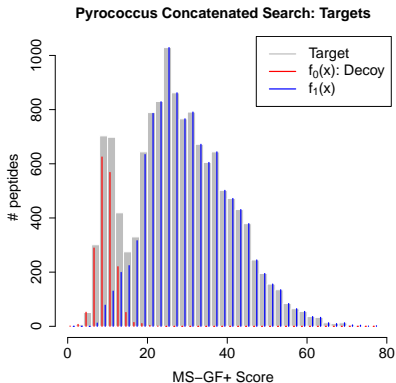
# Target-Decoy approach to establish null distribution



- 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



- 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

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

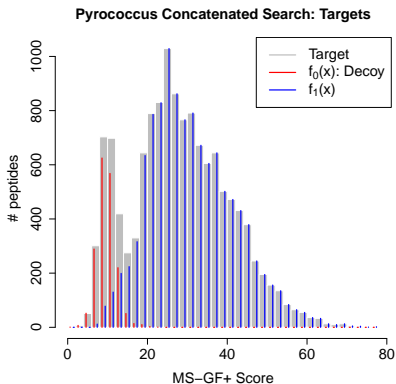
- Score cutoff:  

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

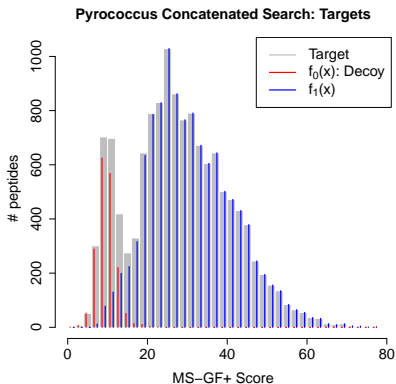
# Target-Decoy approach to establish null distribution

- Competitive Target - decoy:

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



# Target-Decoy approach to establish null distribution



- Competitive Target - decoy:

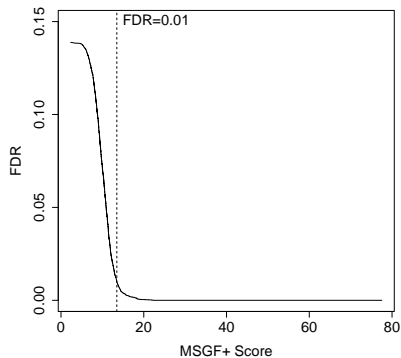
$$\widehat{\text{FDR}}(x) = \frac{\# \text{decoys} | X \geq x}{\# \text{targets} | X \geq x}$$

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

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

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

# Target-Decoy approach to establish null distribution



- Competitive Target - decoy:

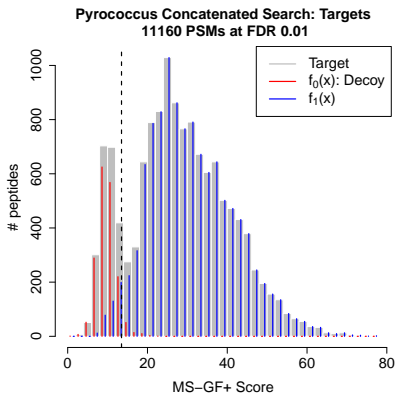
$$\widehat{\text{FDR}}(x) = \frac{\# \text{decoys} | X \geq x}{\# \text{targets} | X \geq x}$$

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

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

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

# Target-Decoy approach to establish null distribution



- Competitive Target - decoy:

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

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

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

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

## Assess TDA assumptions

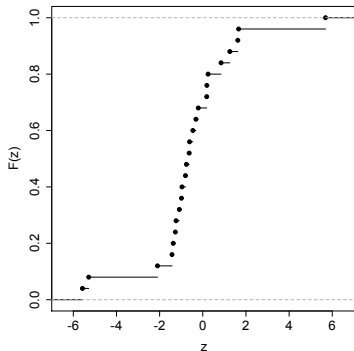
We have to evaluate that

- 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

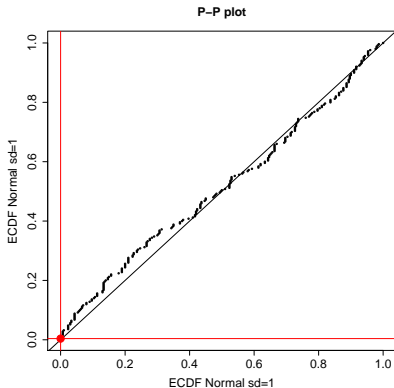
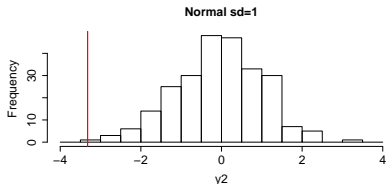
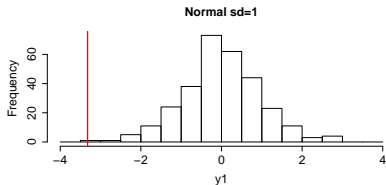


$$\hat{F}_D(x) = \frac{\#\text{decoys} | X \leq x}{\#\text{decoys}}, \quad \hat{F}(x) = \frac{\#\text{targets} | X \leq x}{\#\text{targets}}$$



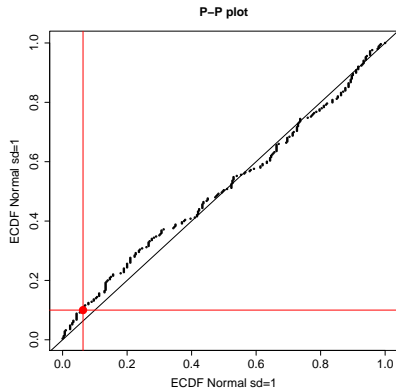
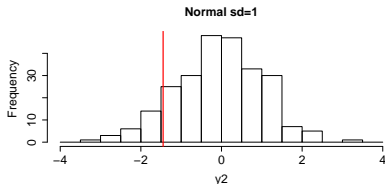
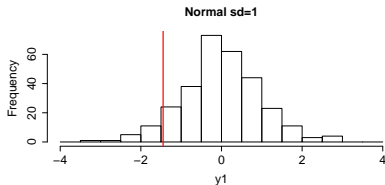
# PP-plot

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



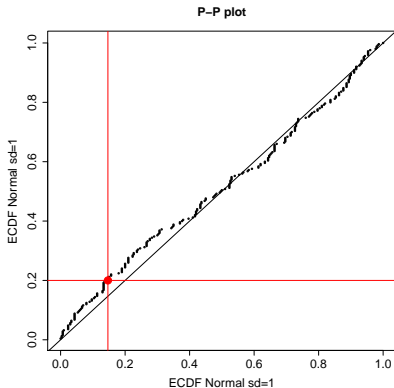
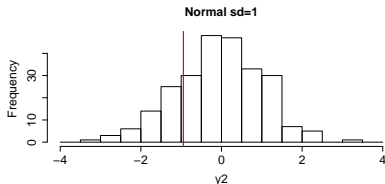
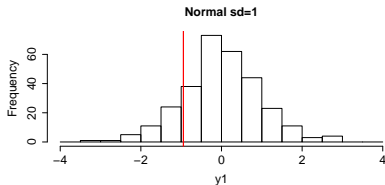
# PP-plot

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



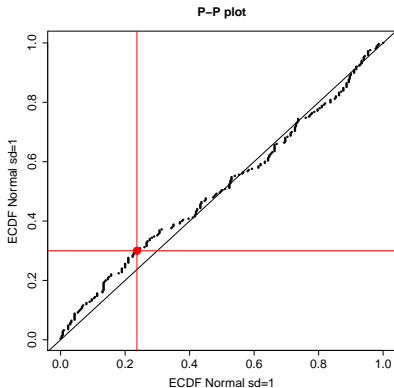
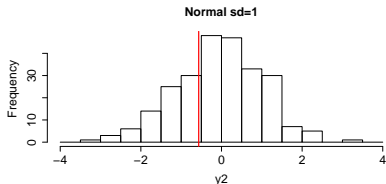
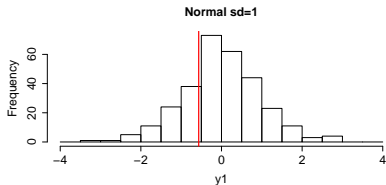
# PP-plot

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



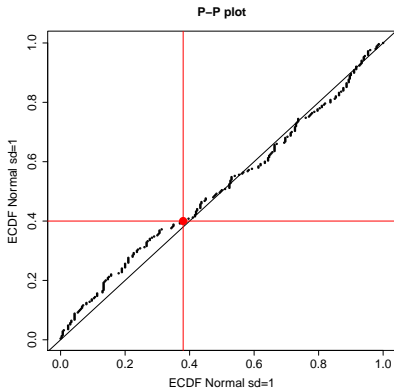
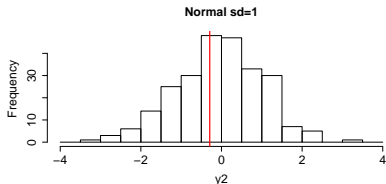
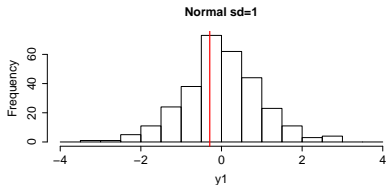
# PP-plot

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



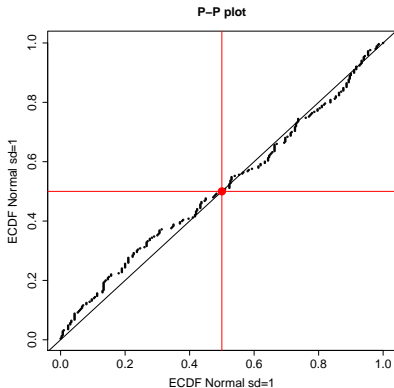
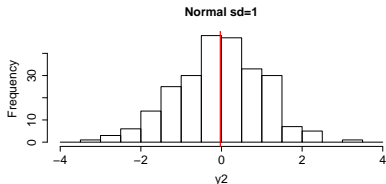
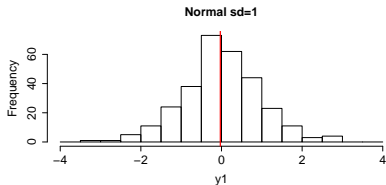
# PP-plot

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



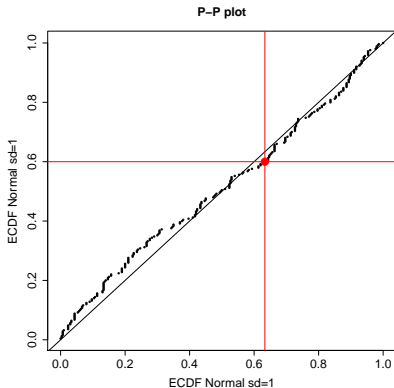
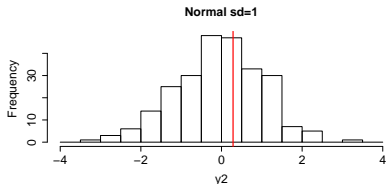
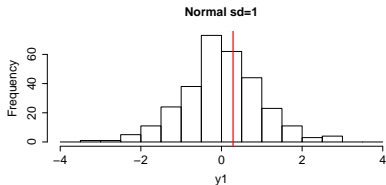
# PP-plot

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



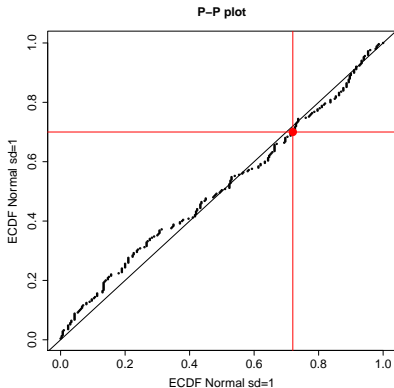
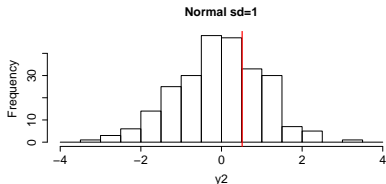
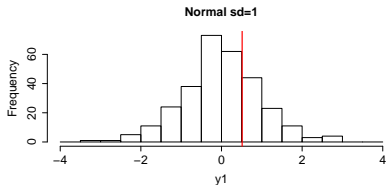
# PP-plot

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

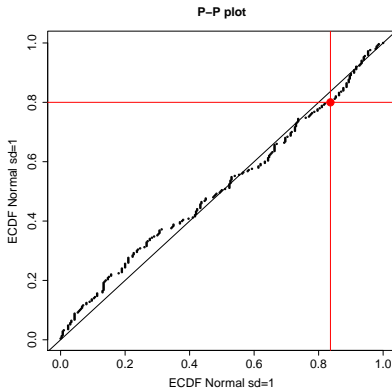
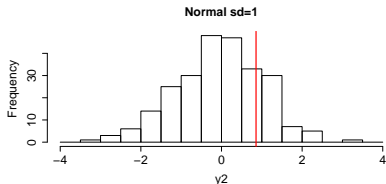
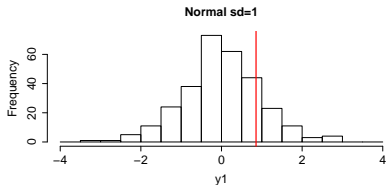
PP-plots have the property that they show a straight 45 degree line through the origin if and only if both distributions are equivalent.





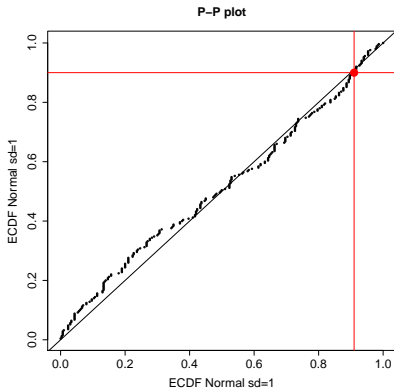
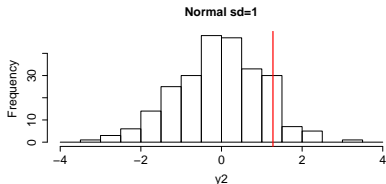
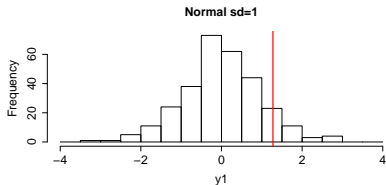
# PP-plot

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



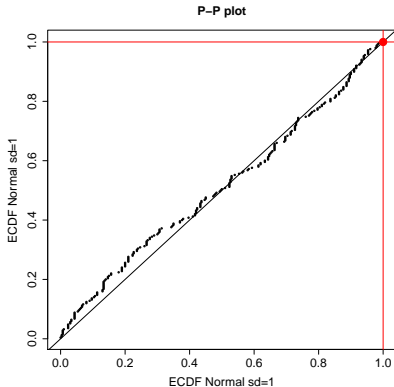
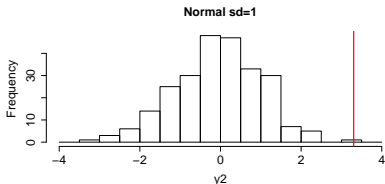
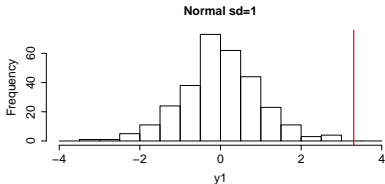
# PP-plot

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

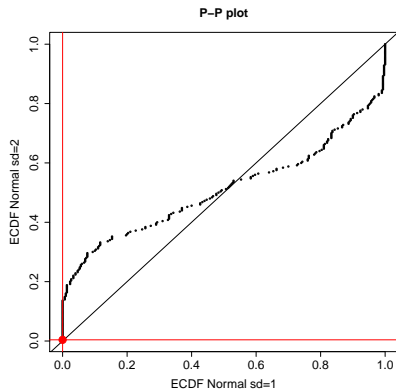
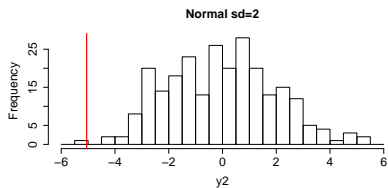
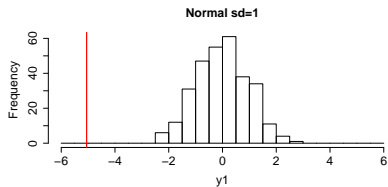


# PP-plot

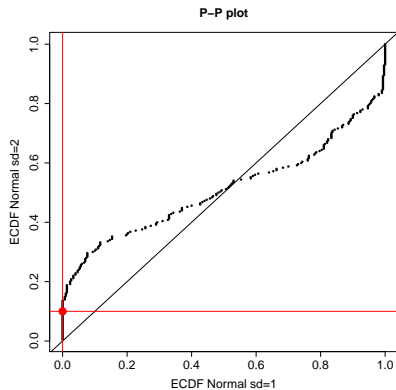
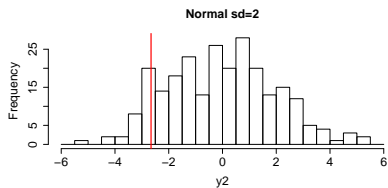
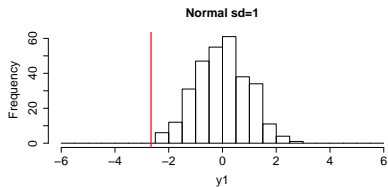
PP-plots have the property that they show a straight 45 degree line through the origin if and only if both distributions are equivalent.



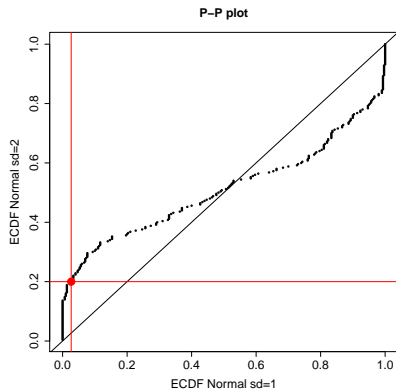
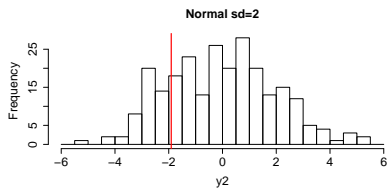
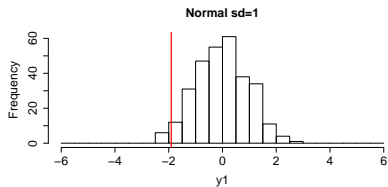
## PP-plot



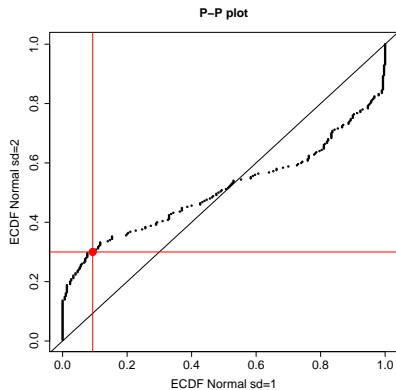
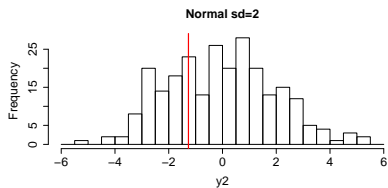
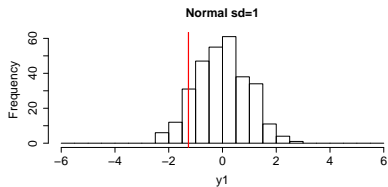
## PP-plot



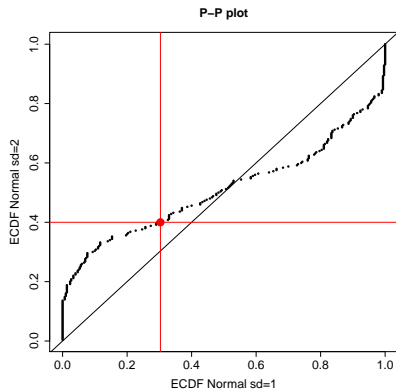
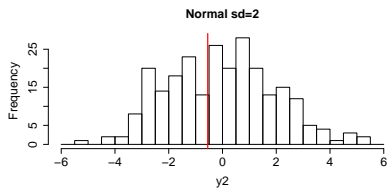
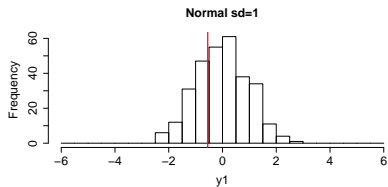
## PP-plot



## PP-plot

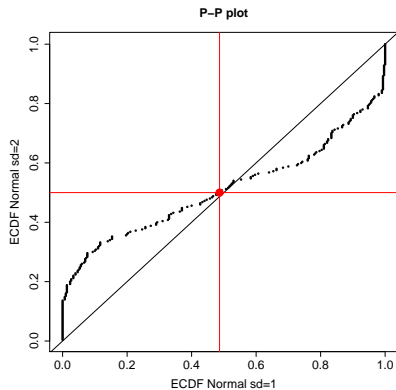
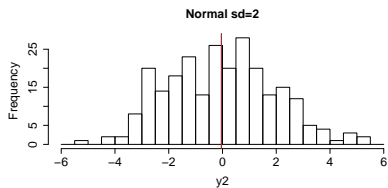
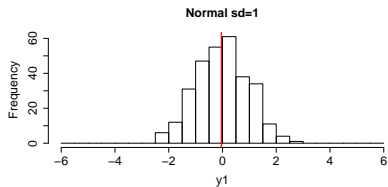


## PP-plot

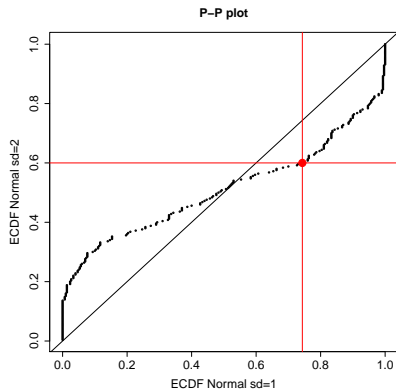
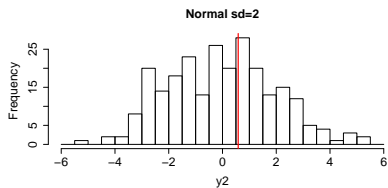
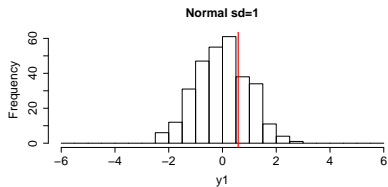




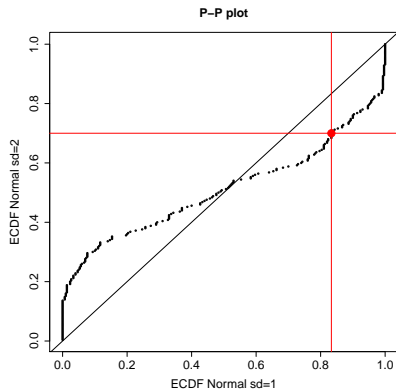
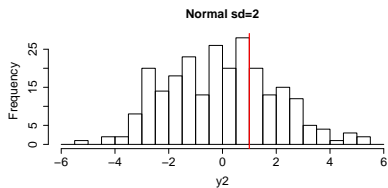
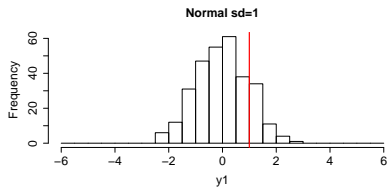
## PP-plot



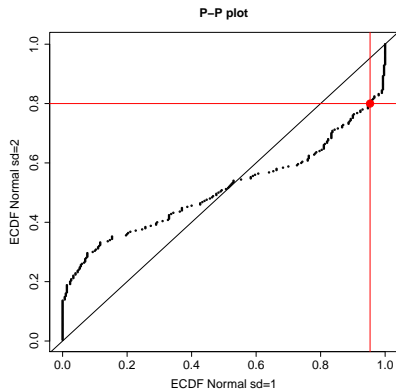
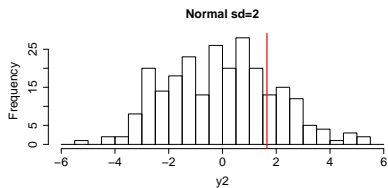
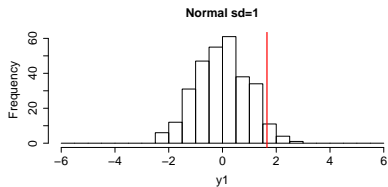
## PP-plot



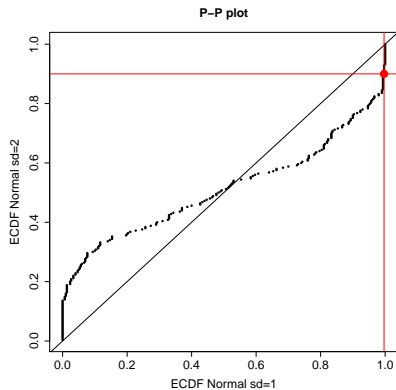
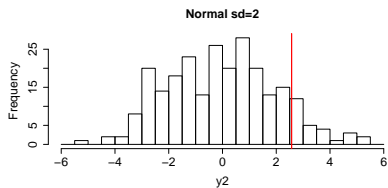
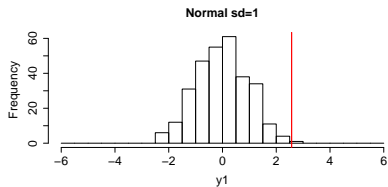
## PP-plot



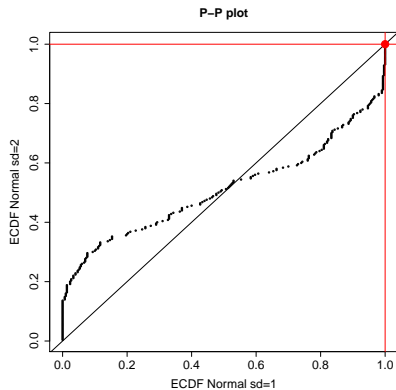
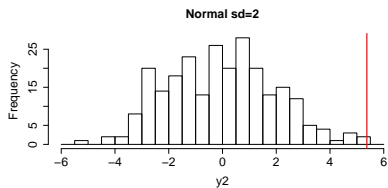
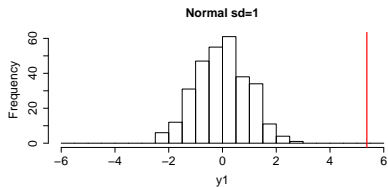
## PP-plot



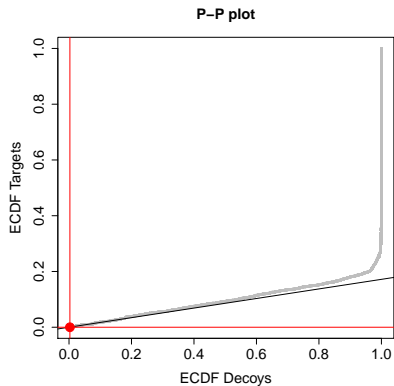
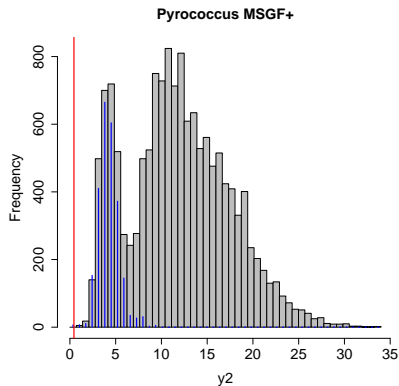
## PP-plot



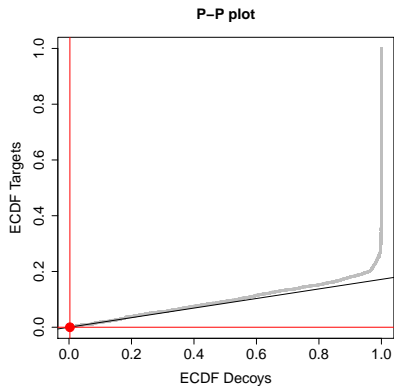
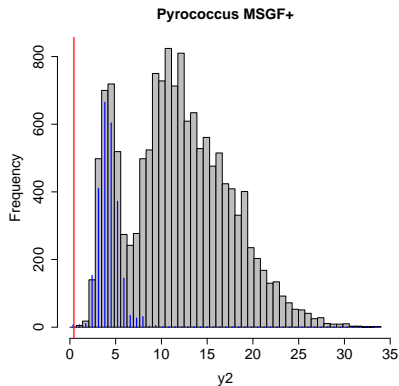
## PP-plot



# PP-plot: pyrococcus



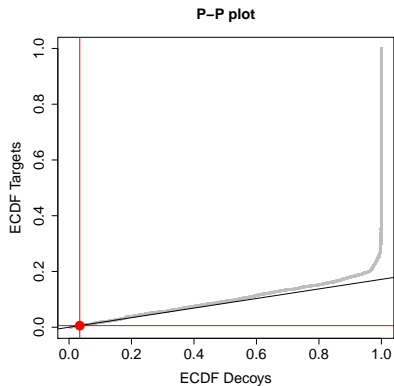
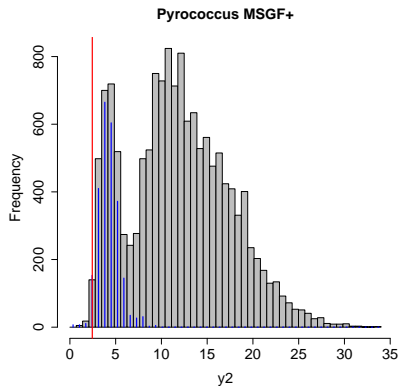
# PP-plot: pyrococcus



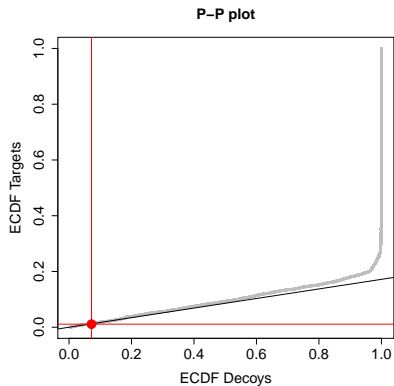
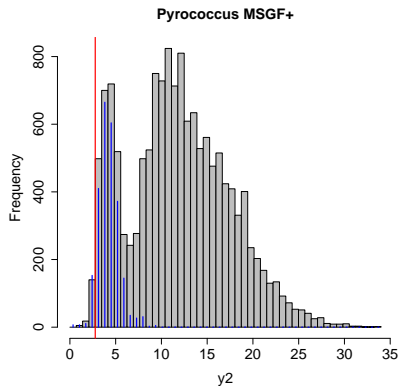
What about  $\hat{\pi}_0$ ?



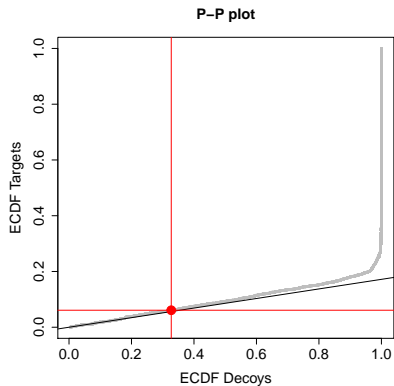
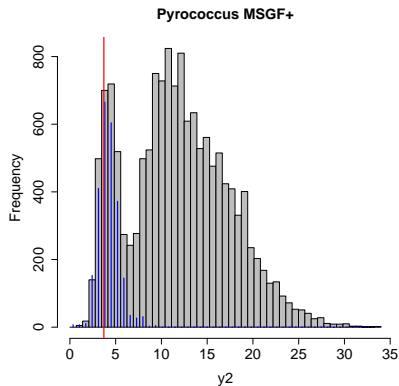
## PP-plot: pyrococcus



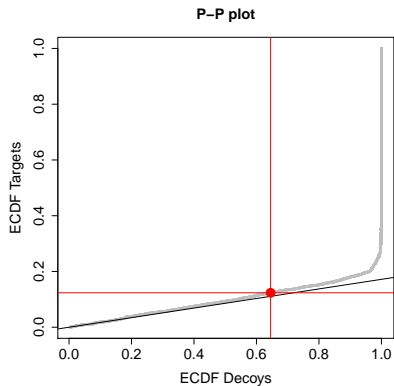
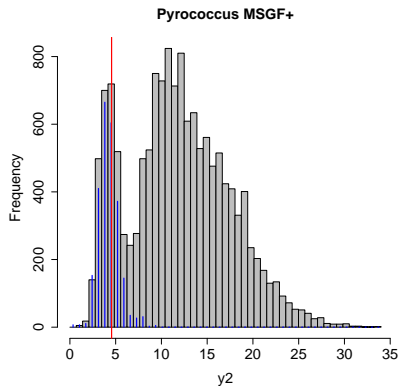
# PP-plot: pyrococcus



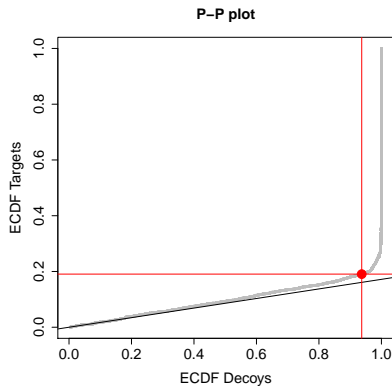
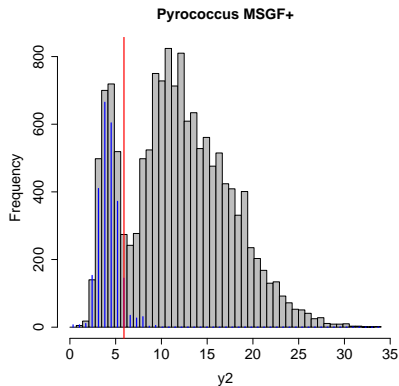
## PP-plot: pyrococcus



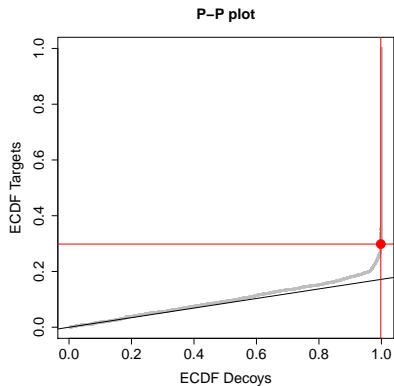
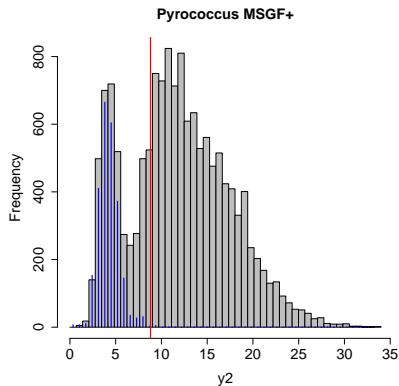
# PP-plot: pyrococcus



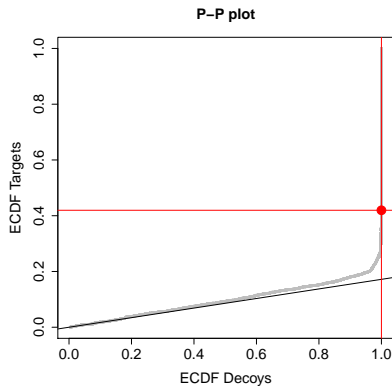
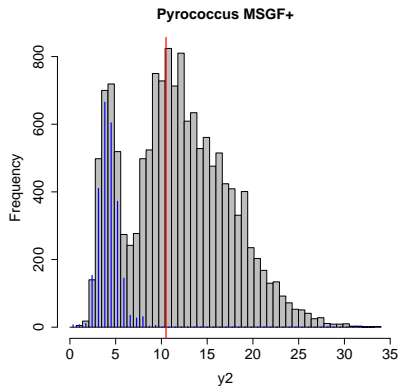
## PP-plot: pyrococcus



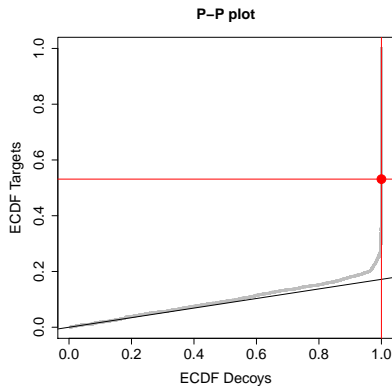
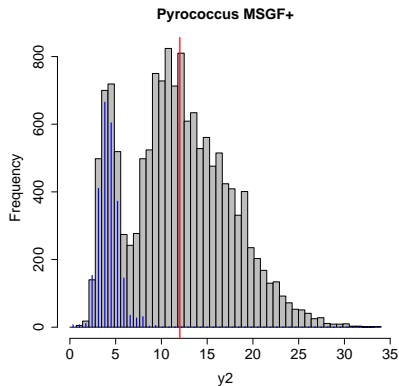
## PP-plot: pyrococcus



## PP-plot: pyrococcus

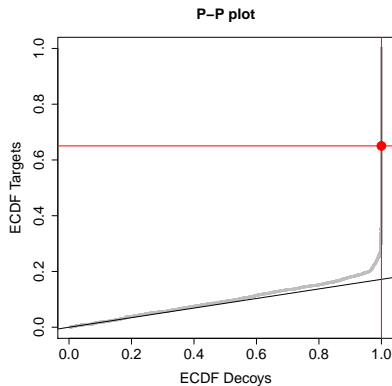
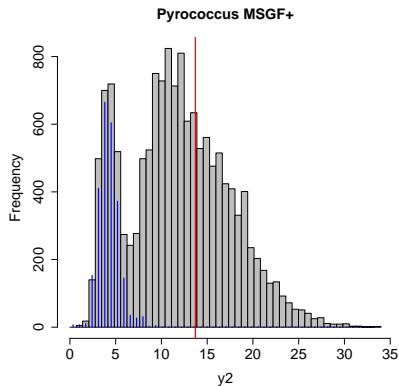


## PP-plot: pyrococcus

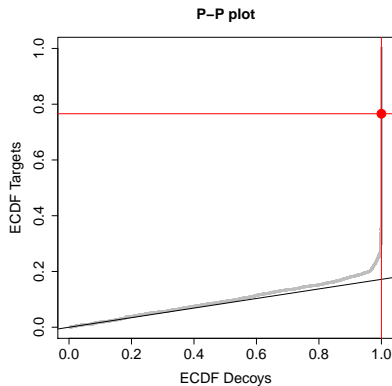
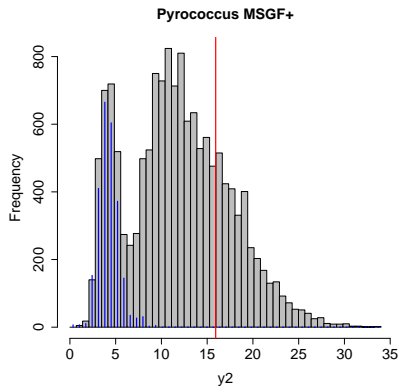




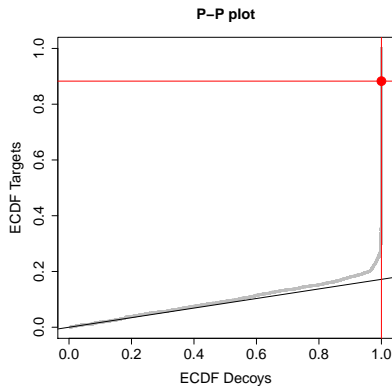
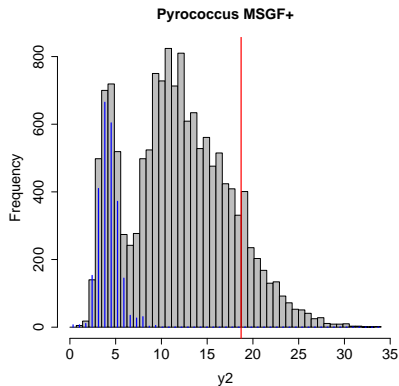
# PP-plot: pyrococcus



## PP-plot: pyrococcus



## PP-plot: pyrococcus



## PP-plot: pyrococcus

