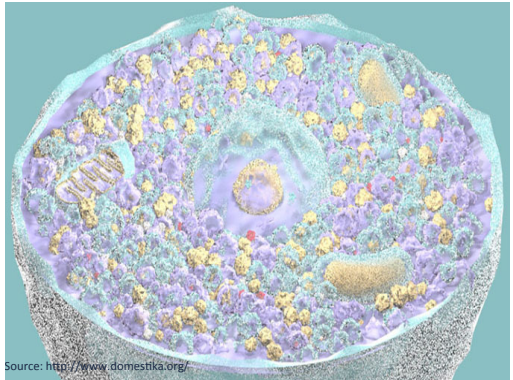
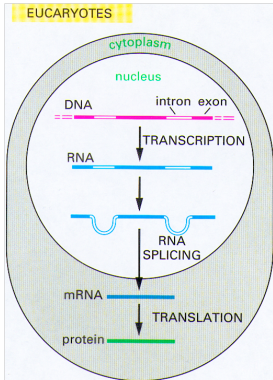


# Statistical Methods for Quantitative MS-Based Proteomics:

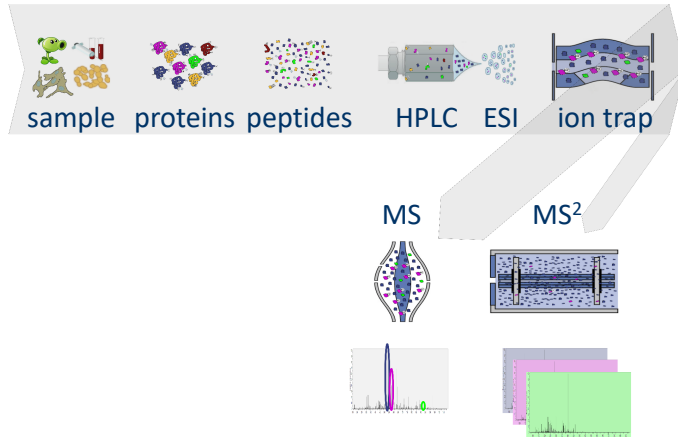
## 1. Identification & False discovery rate

Lieven Clement

Proteomics Data Analysis Shortcourse

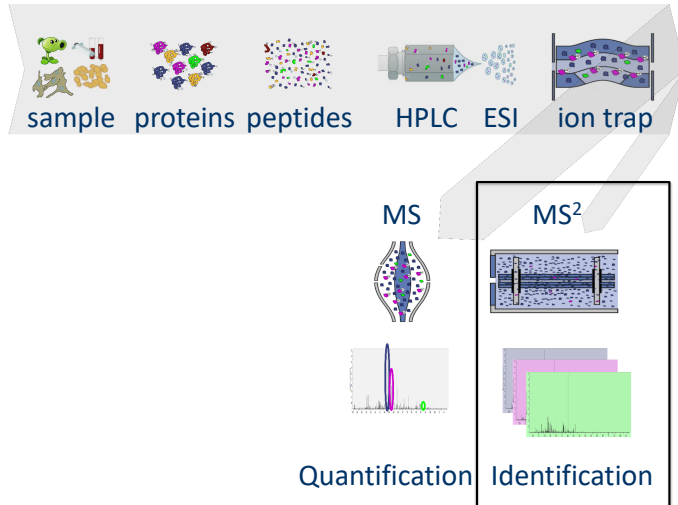


# Challenges in Label Free MS-based Quantitative Proteomics

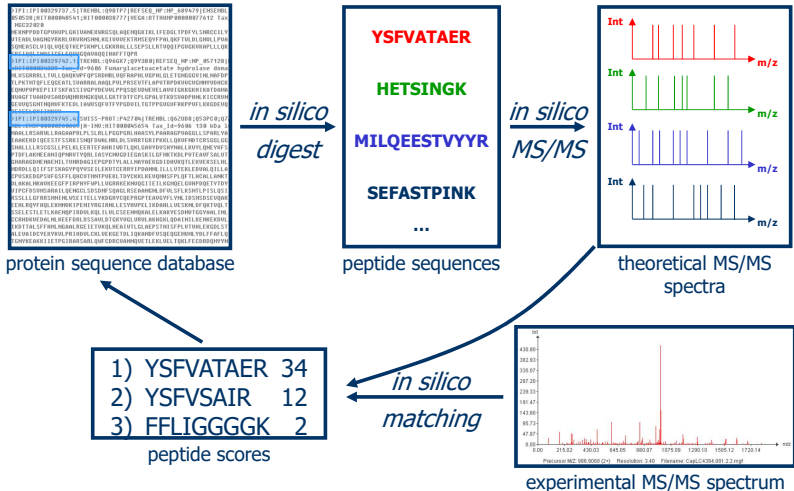


Quantification Identification

# Challenges in Label Free MS-based Quantitative Proteomics



## Identification



(slide courtesy to Lennart Martens)

# 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	$m_0$
Correct hit	FN	TP	$m_1$
Total	NR	R	$m$

Random Variables

# Table of Outcomes

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



# Table of Outcomes

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

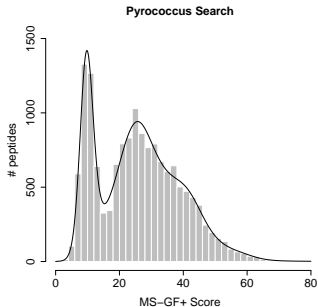
$FDP = \frac{FP}{FP+TP}$ . But is unknown! (FDP: false discovery proportion)

# Table of Outcomes

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

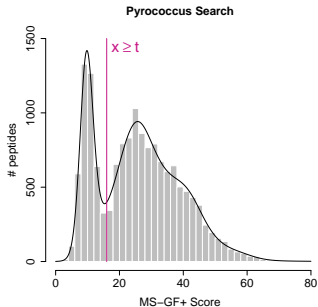
$$FDR = E \left[ \frac{FP}{FP+TP} \right]. \text{ (FDR: false discovery rate)}$$

# Search engines return score that discriminates good from bad matches

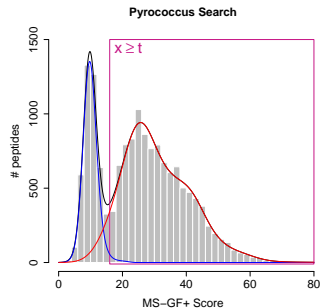


# Search engines return score that discriminates good from bad matches

Score threshold  $t$ ?



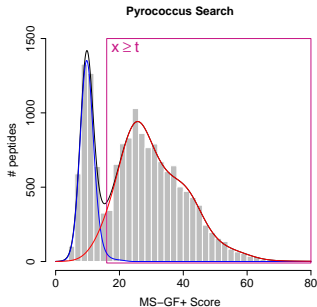
# Search engines return score that discriminates good from bad matches



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

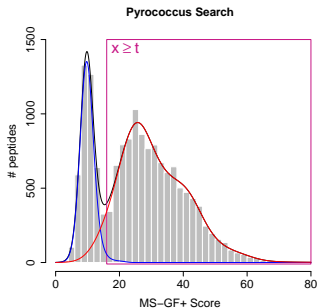


Score threshold  $t$ ?

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

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

# Search engines return score that discriminates good from bad matches



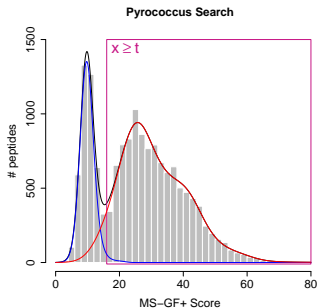
Score threshold  $t$ ?

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

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

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

# Search engines return score that discriminates good from bad matches



Score threshold  $t$ ?

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

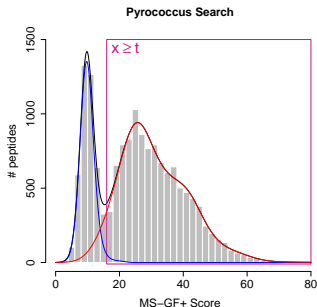
$$\text{FDR}(t) = E \left[ \frac{FP}{FP + TP} \right]$$

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

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



# Search engines return score that discriminates good from bad matches



Score threshold  $t$ ?

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

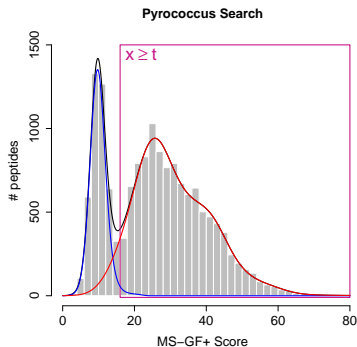
$$\text{FDR}(t) = E \left[ \frac{FP}{FP + TP} \right]$$

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

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

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

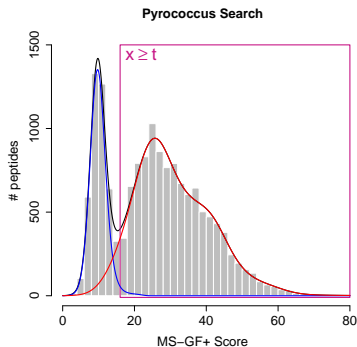
# How to estimate FDR?



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

$$\text{FDR}(t) = E \left[ \frac{FP}{FP + TP} \right] = \frac{m\pi_0 P[x \geq t | FP]}{mP[x \geq t]}$$

# How to estimate FDR?



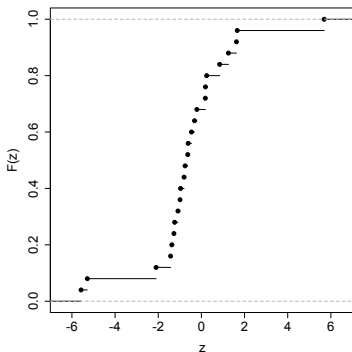
$$P.[x \geq t] = \int_t^{\infty} f.(x)$$

$$P.[x \geq t] \approx \frac{\#x \geq t}{m}$$

$$\text{FDR}(t) = E \left[ \frac{FP}{FP + TP} \right] = \frac{m\pi_0 P[x \geq t|FP]}{mP[x \geq t]}$$

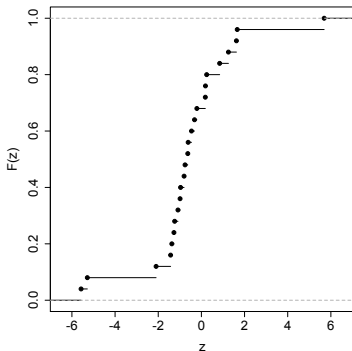
$$\widehat{\text{FDR}}(t) = \frac{m\pi_0 P[x \geq t|FP]}{\#x \geq t}$$

- $F(t) = \int_{-\infty}^t f(x)$  using the Empirical cumulative distribution function (ECDF):  $\bar{F}(t)$



$$\rightarrow \widehat{FDR}(t) = \frac{m\pi_0 P[x \geq t|FP]}{\#x \geq t} = \frac{\pi_0 P[x \geq t|FP]}{\frac{\#x \geq t}{m}} = \frac{\pi_0 [1 - F(t)]}{1 - \bar{F}(t)}$$

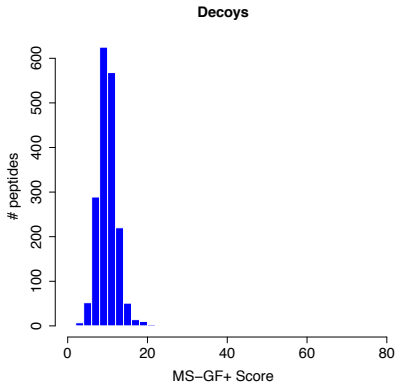
- $F(t) = \int_{-\infty}^t f(x)$  using the Empirical cumulative distribution function (ECDF):  $\bar{F}(t)$



$$\rightarrow \widehat{FDR}(t) = \frac{m\pi_0 P[x \geq t | FP]}{\#x \geq t} = \frac{\pi_0 P[x \geq t | FP]}{\frac{\#x \geq t}{m}} = \frac{\pi_0 [1 - F(t)]}{1 - \bar{F}(t)}$$

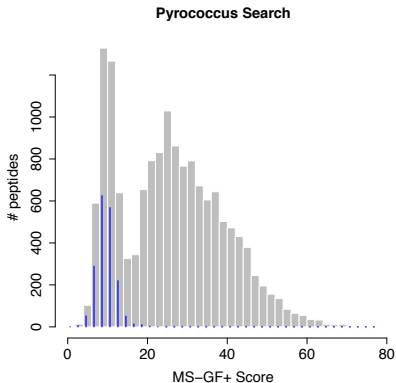
- How to characterize  $F_0(t)$  and  $\pi_0$  in proteomics?

# Target-Decoy approach to establish null distribution



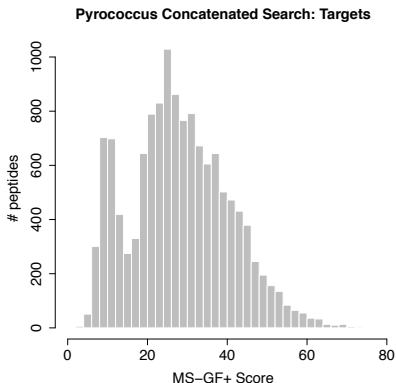
- Searching against decoy databases to generate representative bad hits
- Reversed databases are a popular choice

# Target-Decoy approach to establish null distribution



- Searching against decoy databases to generate representative bad hits
- Reversed databases are a popular choice
- Concatenated search

# Target-Decoy approach to establish null distribution

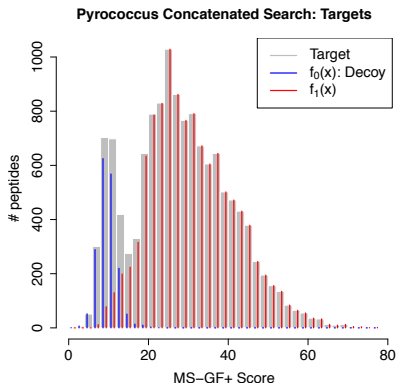


- Searching against decoy databases to generate representative bad hits
- Reversed databases are a popular choice
- Concatenated search
- Assumption that bad hits have an equal probability to map on forward (target) and reverse database (decoy)

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



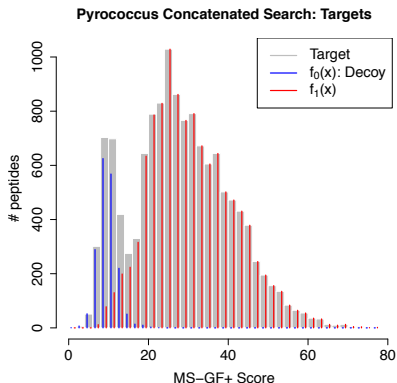
# Target-Decoy approach to establish null distribution



- Searching against decoy databases to generate representative bad hits
- Reversed databases are a popular choice
- Concatenated search
- Assumption that bad hits have an equal probability to map on forward (target) and reverse database (decoy)

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

# Target-Decoy approach to establish null distribution



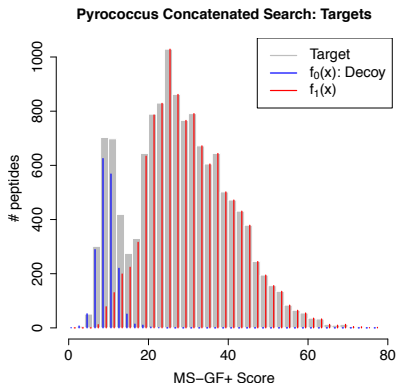
- Score cutoff?

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

- 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



- Score cutoff?

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

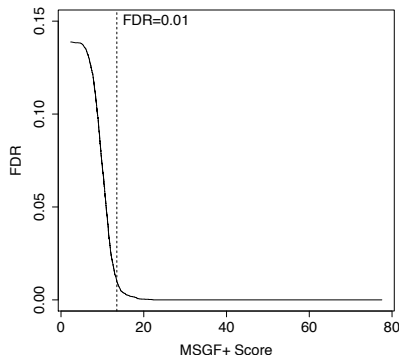
- Competitive Target - decoy:

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

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

$$\widehat{\text{FDR}}(x) = \hat{\pi}_0 \frac{1 - \bar{F}_0(x)}{1 - \bar{F}(x)}$$

# Target-Decoy approach to establish null distribution



- Score cutoff?

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

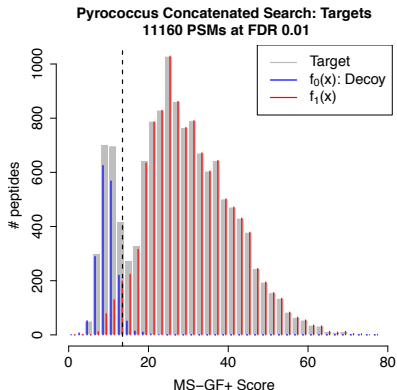
- Competitive Target - decoy:

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

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

$$\widehat{\text{FDR}}(x) = \hat{\pi}_0 \frac{1 - \bar{F}_0(x)}{1 - \bar{F}(x)}$$

# Target-Decoy approach to establish null distribution



- Score cutoff?

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

- Competitive Target - decoy:

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

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

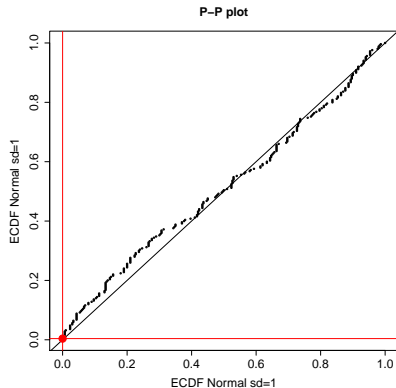
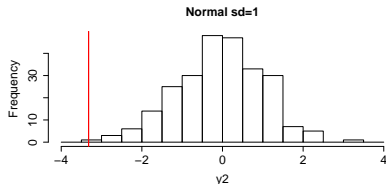
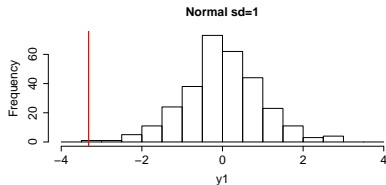
$$\widehat{\text{FDR}}(x) = \hat{\pi}_0 \frac{1 - \bar{F}_0(x)}{1 - \bar{F}(x)}$$

We have to evaluate that

- The decoys are good simulations of the targets: compare  $\bar{F}_0(x)$  with  $\bar{F}(x)$
- $\hat{\pi}_0 = \frac{\#decoys}{\#targets}$  is a good estimator for  $\pi_0$ .
- We will use Probability-Probability-plots for this purpose.
- They plot the ECDFs from two samples in function of each other.

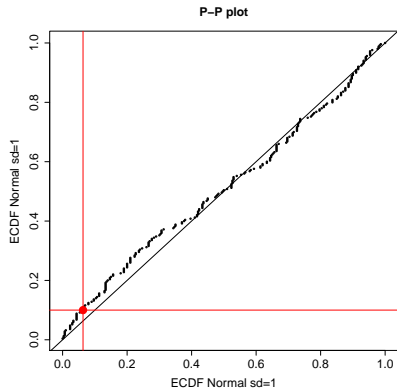
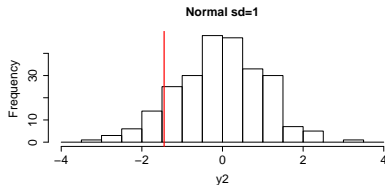
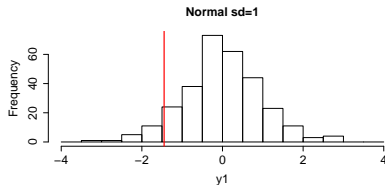
# PP-plot

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

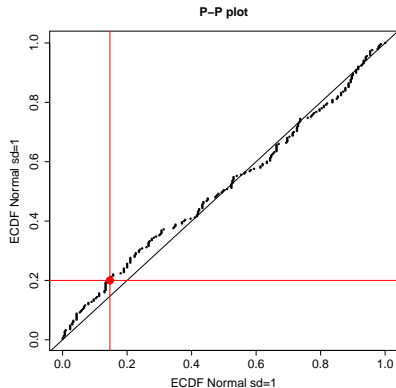
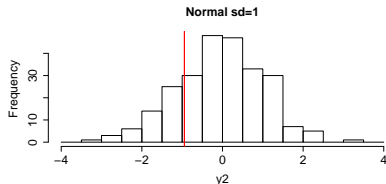
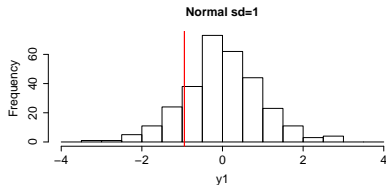
PP-plots have the property that they show a straight 45 degree line through the origin if and only if both distributions are equivalent.





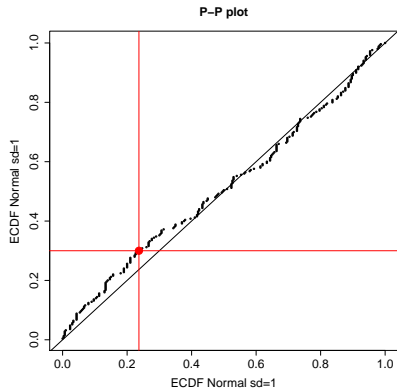
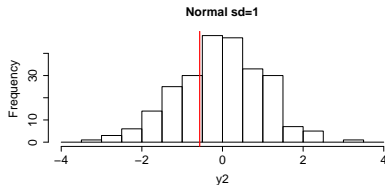
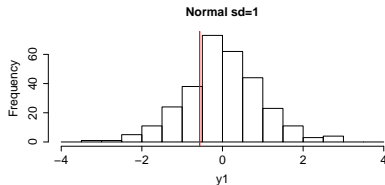
# PP-plot

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



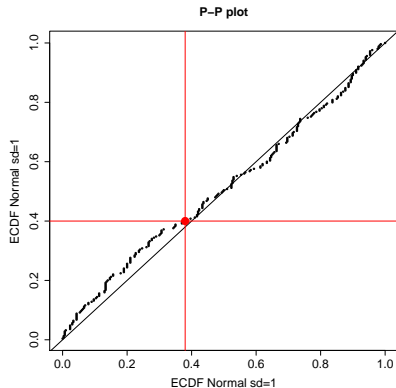
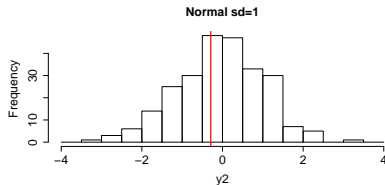
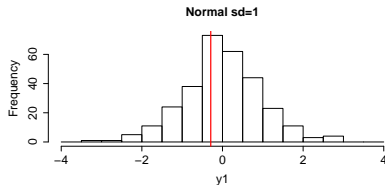
# PP-plot

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



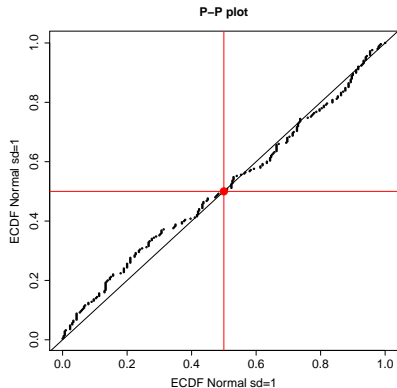
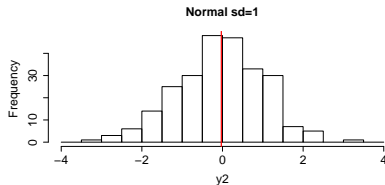
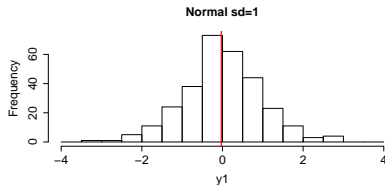
# PP-plot

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



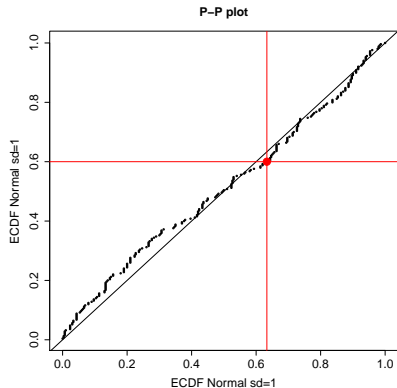
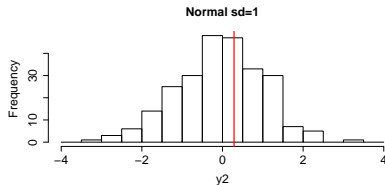
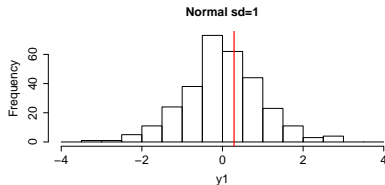
# PP-plot

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



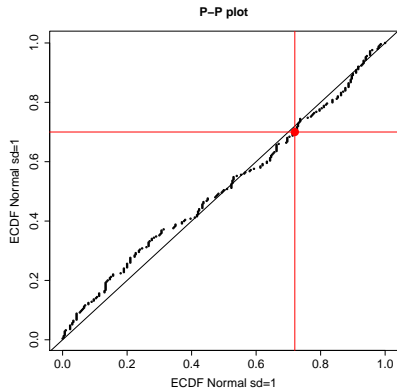
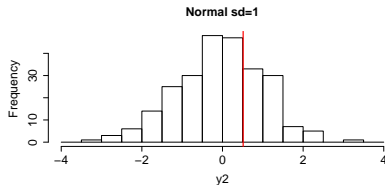
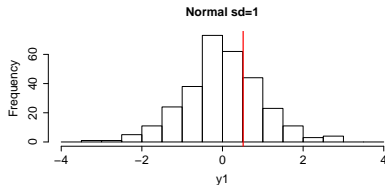
# PP-plot

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



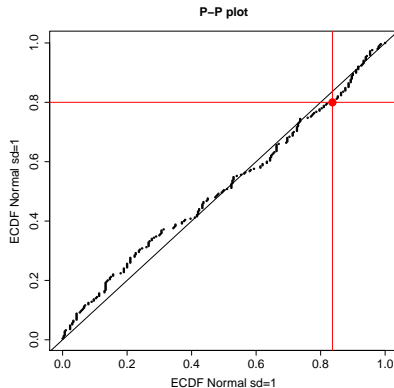
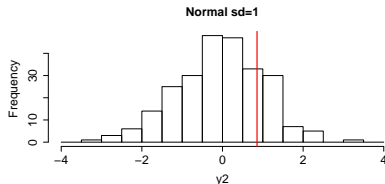
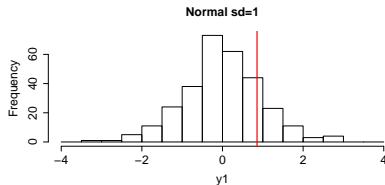
# PP-plot

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



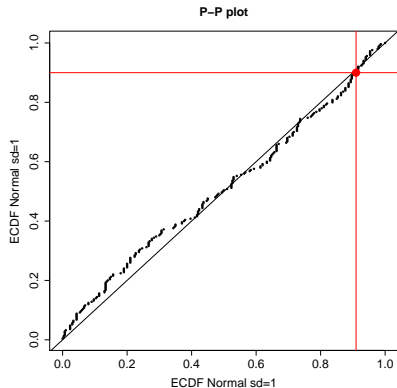
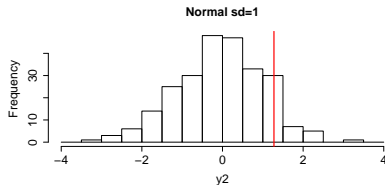
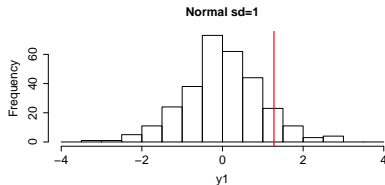
# PP-plot

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

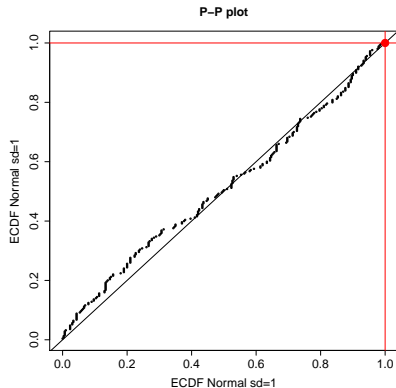
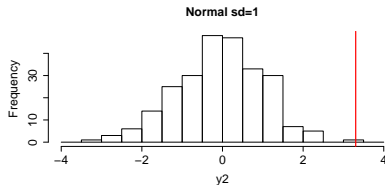
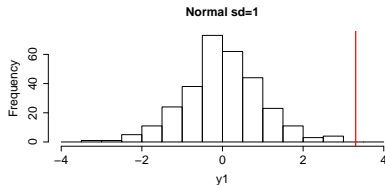
PP-plots have the property that they show a straight 45 degree line through the origin if and only if both distributions are equivalent.



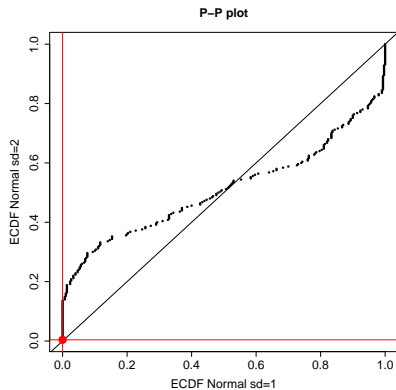
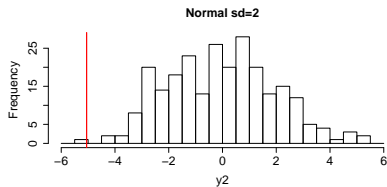
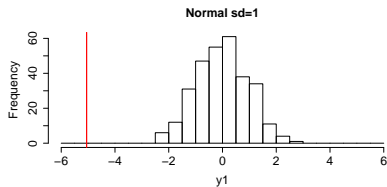


# PP-plot

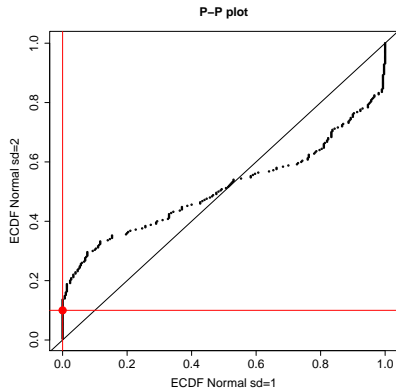
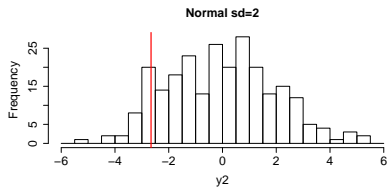
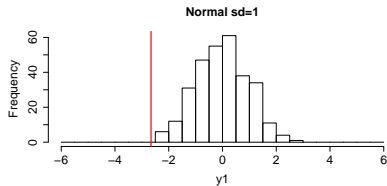
PP-plots have the property that they show a straight 45 degree line through the origin if and only if both distributions are equivalent.



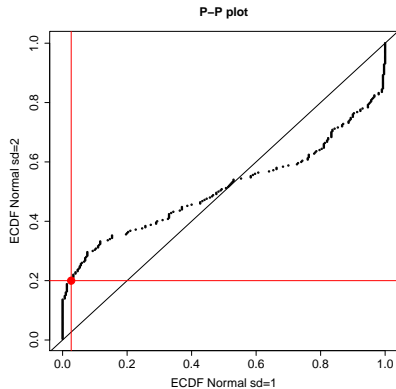
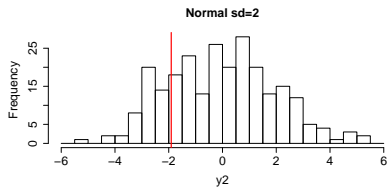
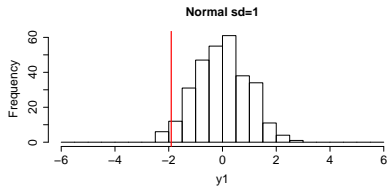
# PP-plot



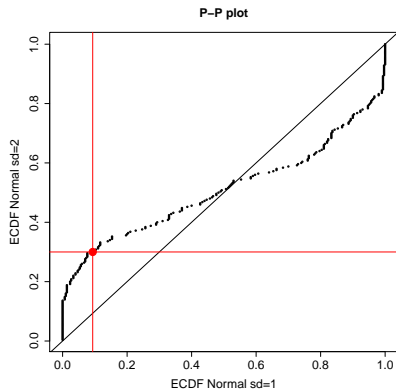
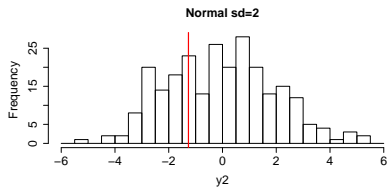
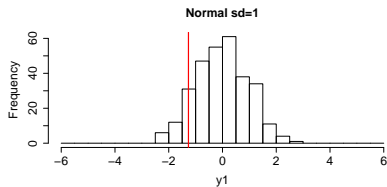
# PP-plot



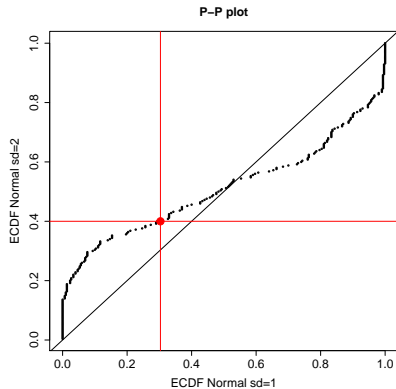
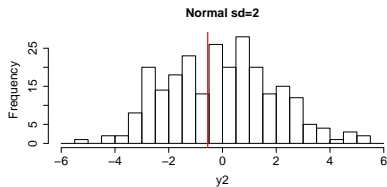
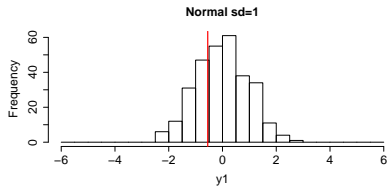
# PP-plot



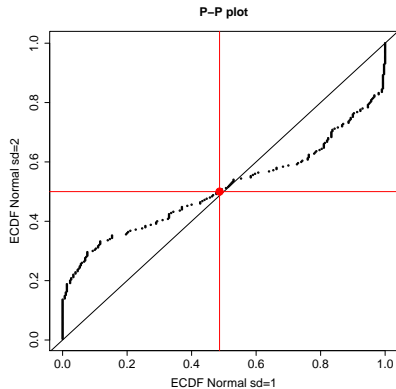
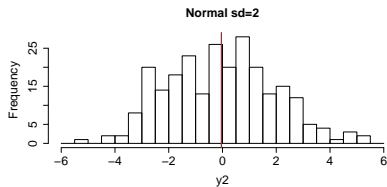
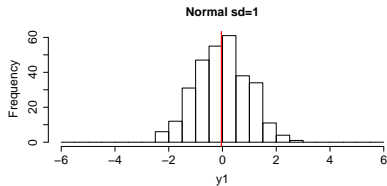
# PP-plot



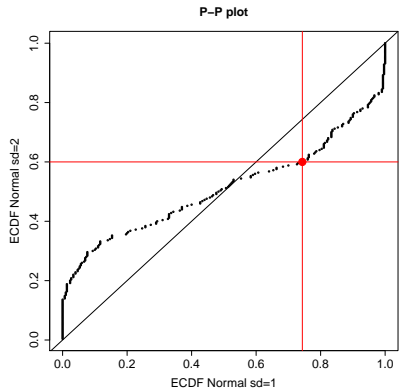
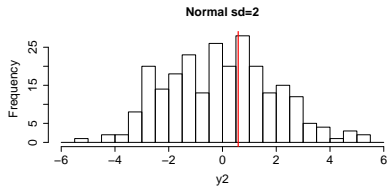
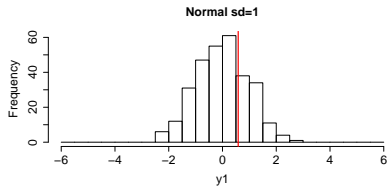
# PP-plot



# PP-plot

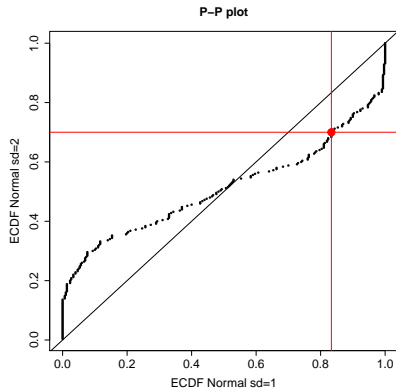
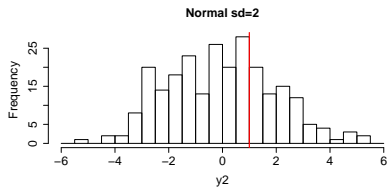
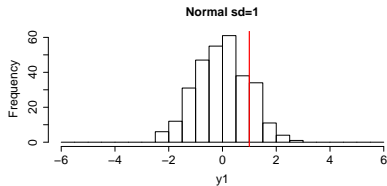


# PP-plot

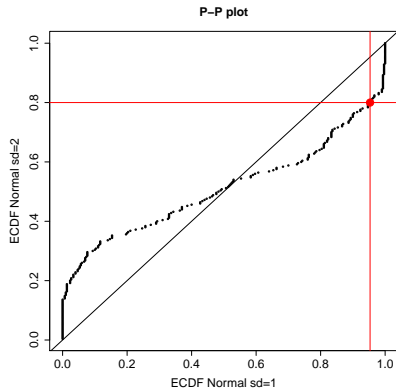
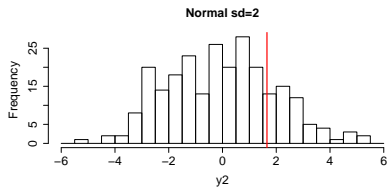
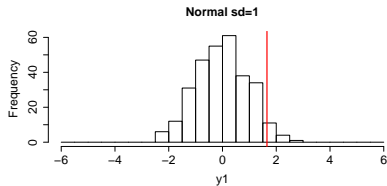




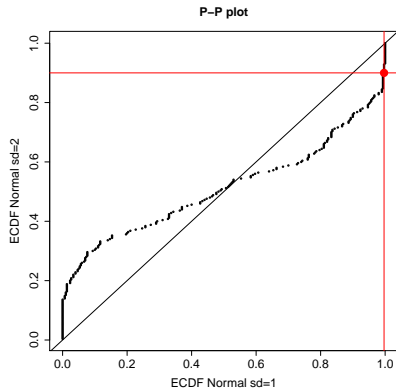
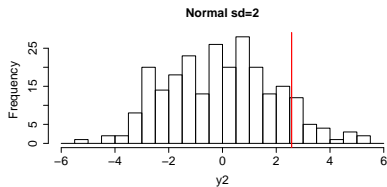
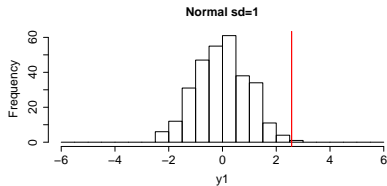
# PP-plot



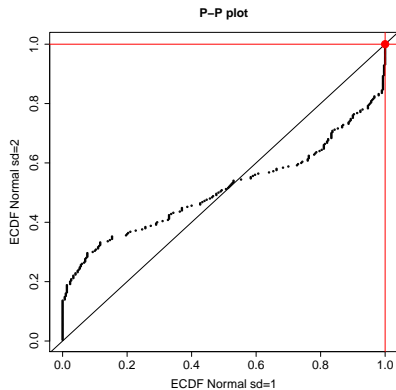
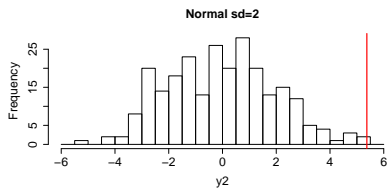
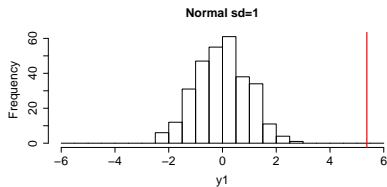
# PP-plot



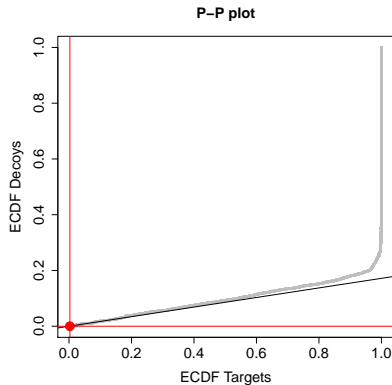
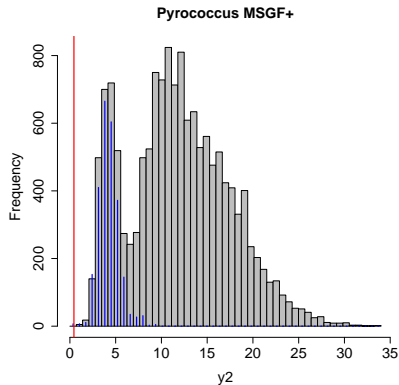
# PP-plot



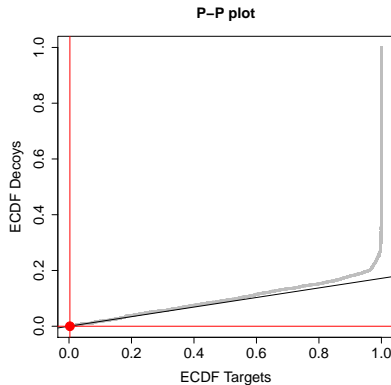
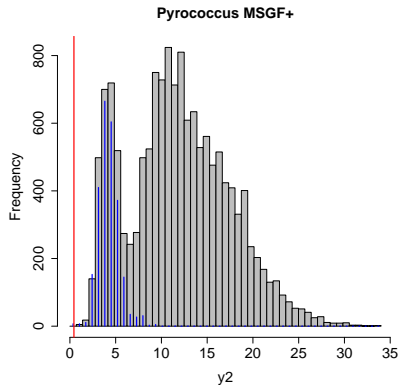
# PP-plot



# PP-plot: pyrococcus

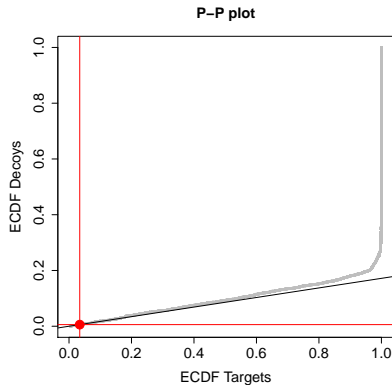
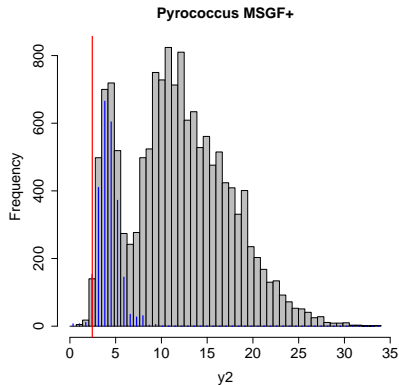


# PP-plot: pyrococcus

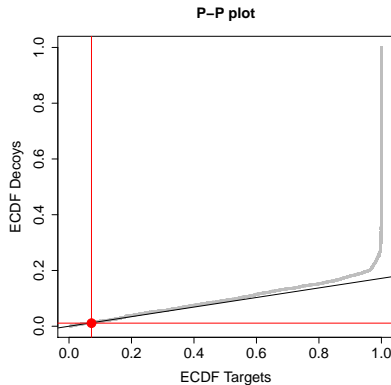
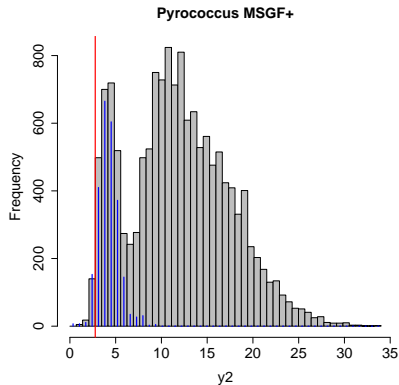


What about  $\hat{\pi}_0$ ?

# PP-plot: pyrococcus

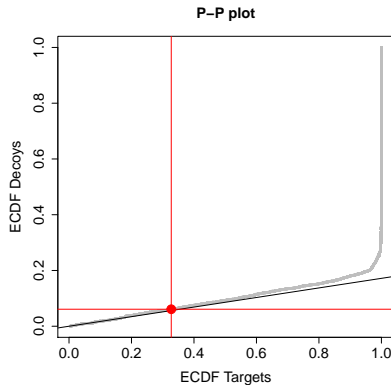
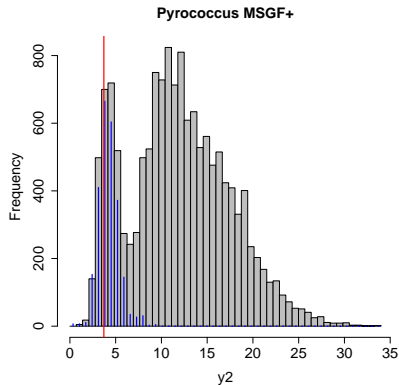


# PP-plot: pyrococcus

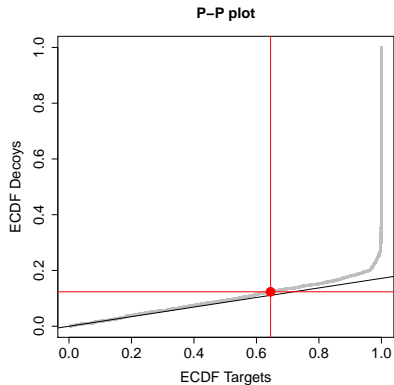
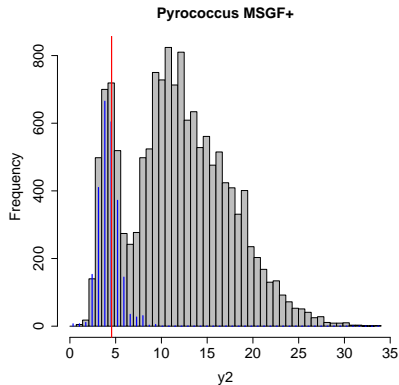




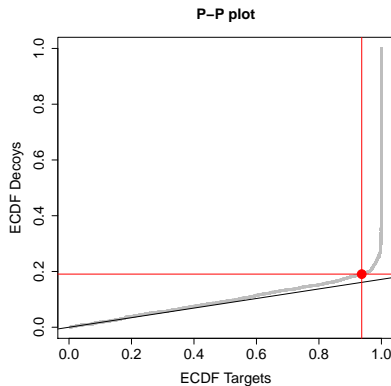
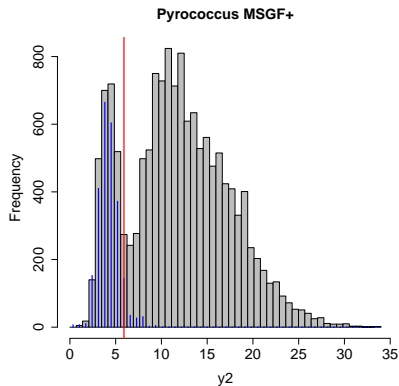
# PP-plot: pyrococcus



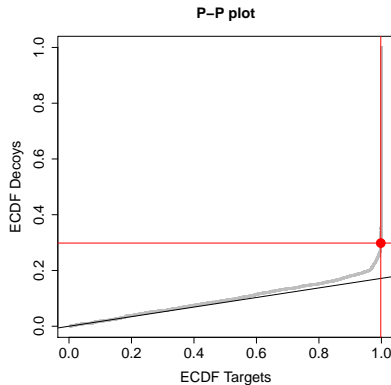
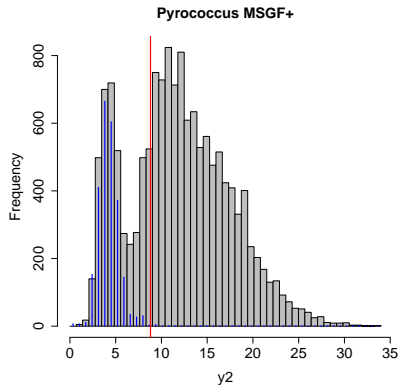
# PP-plot: pyrococcus



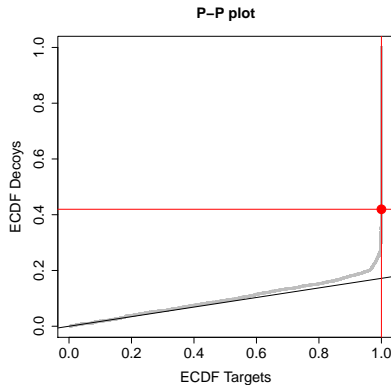
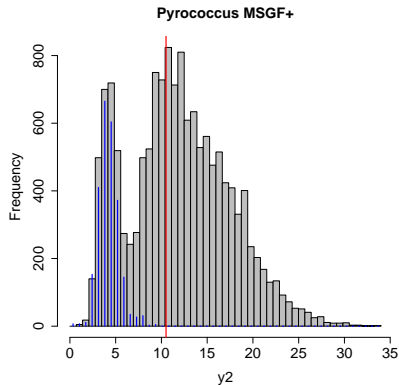
# PP-plot: pyrococcus



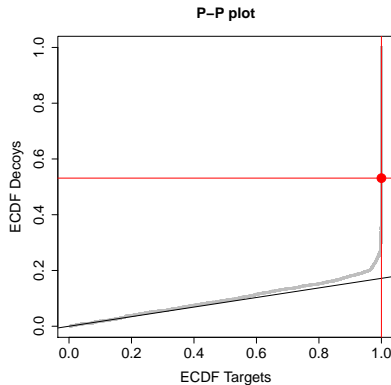
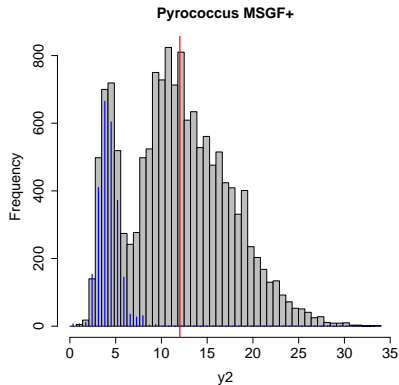
# PP-plot: pyrococcus



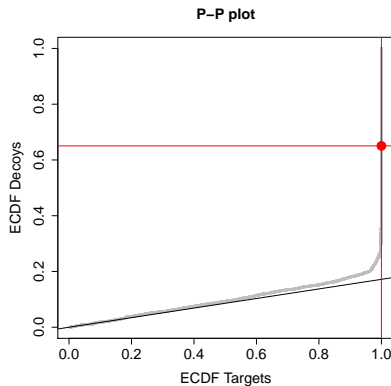
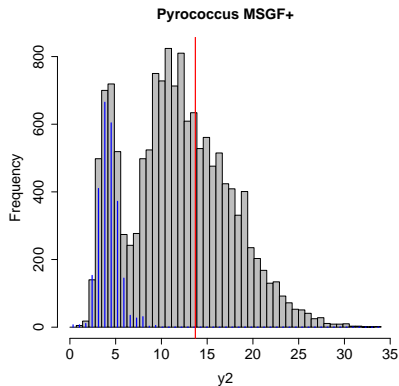
# PP-plot: pyrococcus



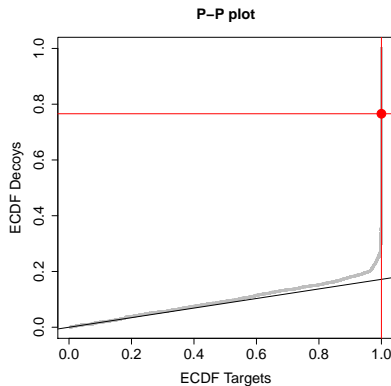
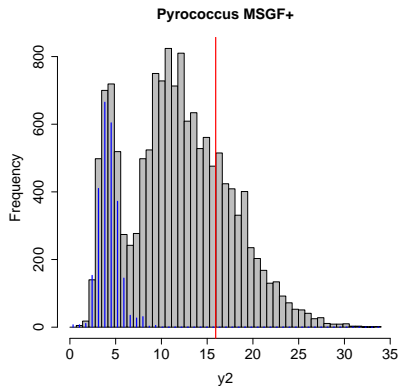
# PP-plot: pyrococcus



# PP-plot: pyrococcus

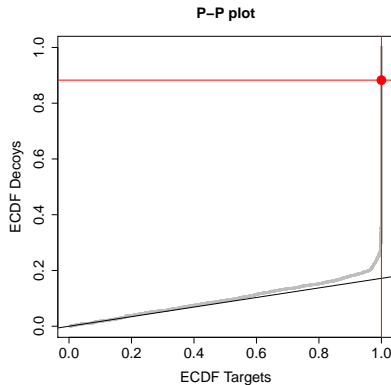
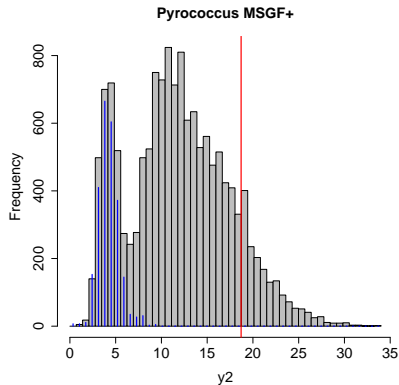


# PP-plot: pyrococcus





# PP-plot: pyrococcus



# PP-plot: pyrococcus

