

Stat 2911 Lecture Notes

Class 27, 2017

Uri Keich

© Uri Keich, The University of  
Sydney

Quantile-Quantile and  
probability plots, Extrema and  
Order Statistics (Rice 9.8,  
10.2.3, 3.7)

## Visual comparison of distributions

The Q-Q (quantile-quantile) plot of  $F_x$  and  $F_y$  is defined as the graph

$$\{(\bar{F}_x^{-1}(p), \bar{F}_y^{-1}(p)) : p \in (0,1)\}$$

2)  $y = a\bar{X} + b$        $a > 0, b \in \mathbb{R}$

Let  $p \in (0,1)$  and denote  $y_p = F_y^{-1}(p)$ . Since

$$\begin{aligned} F_y(y) &= P(a\bar{X} + b \leq y) \\ &= P(\bar{X} \leq (y-b)/a) \\ &= F_{\bar{X}}\left(\frac{y-b}{a}\right), \end{aligned}$$

$$\begin{aligned} x_p &:= F_{\bar{X}}^{-1}(p) \\ &= \frac{y_p - b}{a} \end{aligned}$$

$$\Rightarrow y_p = ax_p + b.$$

$y_p = \varphi(x_p)$  where  $Y = \varphi(\bar{X})$  this applies for any strictly ↑  $\varphi$ .

Therefore the Q-Q plot in this case is a section of the line  $y = ax + b$ .

Moreover, the y-intercept of the line is  $b = \text{location}$ , and the slope is  $a = \text{scale}$ .

In particular, recall  $Y \sim N(\mu, \sigma^2)$  can be represented as  $Y = \sigma \bar{X} + \mu$  where  $\bar{X} \sim N(0,1)$ .

⇒ The Q-Q plot of  $\bar{X}$  against  $Y$  will be the line with slope  $\sigma$  and y-intercept of  $\mu$ .

What if  $F_{\bar{X}}^{-1}$  or  $F_Y^{-1}$  do not exist?

More generally the Q-Q plot is defined as

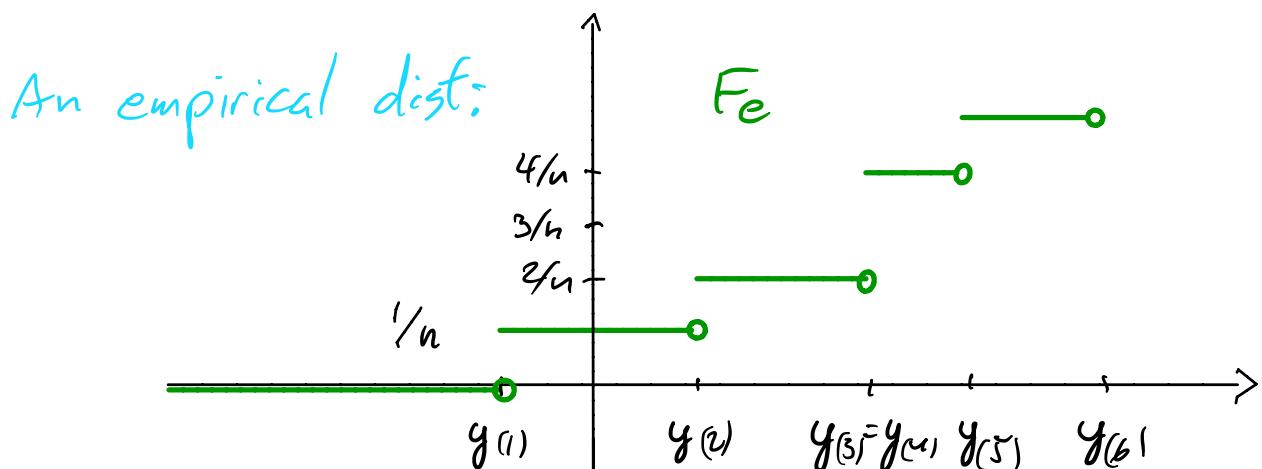
$$\{(Q_{F_X}(p), Q_{F_Y}(p)) : p \in (0, 1)\},$$

where  $Q_F(p) = \min \{x : F(x) \geq p\}$  ( $= F^{-1}(p)$  if  $F^{-1}$  exists)

Typically we wouldn't compare CDFs.

Instead, given a sample  $y_1, \dots, y_n$  from an unknown dist.  $F_Y$ , we want to visually inspect the fit between the sample and a suspected dist.  $F_X$ .

We can do this by comparing  $F_X$  with  $F_E$ , the empirical (discrete) dist. of the sample.



$y(i)$  are the ordered sample points (order statistics)

so  $\{y(i)\}_i^n = \{y_i\}_i^n$  (including multiplicities!) and  $y_{(1)} \leq y_{(2)} \leq y_{(3)} \dots \leq y_{(n)}$ . Formally,

$$F_E(y) = F_E^{(n)}(y) = \frac{|\{y_i : y_i \leq y\}|}{n}.$$

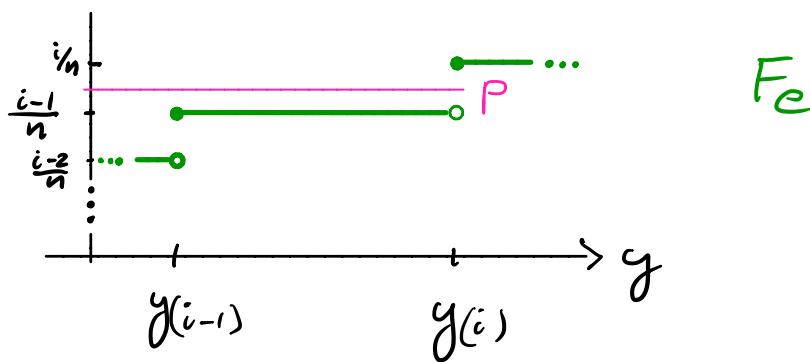
What happens as  $n \rightarrow \infty$ ? What is  $\lim_{n \rightarrow \infty} F_e^n(y)$ ?

For a fixed  $y \in \mathbb{R}$  and  $n \in \mathbb{N}$ , if  $y_1, \dots, y_n$  are iid  $F_Y$ -distributed Rvs,  $|\{y_i : Y_i \leq y\}| \sim$   
where  $p = F_Y(y)$ . Therefore,

$$F_e^n(y) = \frac{|\{y_i : Y_i \leq y\}|}{n} \xrightarrow{n} \text{by LLN.}$$

In fact  $F_e^n \xrightarrow{n} F_Y$  uniformly on  $\mathbb{R}$  with prob. 1 (a.s.)  
(Glivenko-Cantelli)

How does  $Q_{F_e}$  look like?



Clearly, for any  $P \in (\frac{i-1}{n}, \frac{i}{n}]$ ,  $Q_{F_e}(P) = y(i)$  so we can choose which  $p$  in  $(\frac{i-1}{n}, \frac{i}{n})$  we select and accordingly pair  $Q_{F_e}(p)$  with  $Q_{F_e}(p) = y(i)$ .

The Q-Q plot is commonly defined using  $p_i = \frac{i}{n+1}$ :

$$\frac{i-1}{n} < p_i = \frac{i}{n+1} < \frac{i}{n} \quad i=1, 2, \dots, n$$

so the plot consists of the  $n$ -points graph:

$$\left\{ \left( Q_{F_X}\left(\frac{i}{n+1}\right), y(i) \right) : i=1, \dots, n \right\}.$$

If  $F_{\Sigma}^{-1}$  exists then the Q-Q plot is the graph:

$$\left\{ \left( F_{\Sigma}^{-1}\left(\frac{i}{n+1}\right), y_{(i)} \right) : i=1, \dots, n \right\}.$$

The R function qqnorm which draws the Q-Q plot of a sample  $\{y_i\}$  against  $F_{\Sigma} = \Phi(N(0,1))$  uses  $p_i = \frac{i-1/2}{n}$  for  $n \geq 10$ .

We will return to the selection of  $p_i$  later.

As mentioned,  $\hat{F}_c(y) \xrightarrow{n} F_y(y)$  and it follows that with  $i_n \rightarrow \infty$  s.t.  $p_{i_n} = \frac{i_n}{n+1} \xrightarrow{n} p$ , with prob. 1

$$\lim_{n \rightarrow \infty} Y_{(i_n)}^n = \lim_{n \rightarrow \infty} Q_{\hat{F}_c^n}(p_{i_n}) = F_y^{-1}(p),$$

provided  $F_y^{-1}$  exists and is cont. at  $p$  (e.g.  $F_y$  is cont. and strictly ↑).

It follows that as  $n \rightarrow \infty$  the (empirical) Q-Q plot will converge to the distributional Q-Q plot:

$$\left\{ \left( F_{\Sigma}^{-1}(p), F_y^{-1}(p) \right) : p \in (0,1) \right\}$$

In particular, if  $y = a\Sigma + b$ , we expect the empirical Q-Q plot to roughly follow a straight line. Useful for gauging if your sample came from a normal dist.

Q-Q plots compare quantiles or  $F^{-1}$ .

Probability (P-P) plots compare the CDFs ( $F$ ):

If  $F$  and  $G$  are CDFs then the **probability plot** comparing  $F$  and  $G$  is defined by the graph

$$\{(G(x), F(x)) : x \in \mathbb{R}\}$$

Advantage over Q-Q plots: the scale is constant.

If  $G^{-1}$  exists then this graph is essentially the same as

$$\{(p, F(G^{-1}(p))) : p \in (0, 1)\},$$

which suggests the following alternative definition of a **probability plot**:

$$\{(p, F(Q_G(p))) : p \in (0, 1)\}.$$

And for comparing a sample  $y_1, \dots, y_n$  with  $F$

$$\left\{ \left( \frac{i}{n+1}, F(y_{(i)}) \right) : i=1, \dots, n \right\}$$

In summary:

Given a sample  $y_1, \dots, y_n$  we can visually gauge how well it fits the CDF  $F$  by plotting

$$\text{Q-Q: } F^{-1}\left(\frac{i}{n+1}\right) \quad \text{vs.} \quad y_{(i)}$$

$$\text{P-P: } \frac{i}{n+1} \quad \text{vs.} \quad F(y_{(i)})$$

# Theoretical PP & QQ against normal

User: Anon

May 15, 2017

The following three distributions are visually compared with the standard normal distribution using QQ plots and PP.

- $N(10, \sigma^2 = 25)$
- A gamma distribution with shape  $a = 4$  and rate  $\lambda = 0.4$  which has the same mean (location) and standard deviation (scale) as the normal  $N(10, \sigma^2 = 25)$  distribution
- Cauchy distribution with location  $\mu = 10$  and scale  $\sigma = 5$  again:

$$f(x) = \frac{1}{\pi} \cdot \frac{1}{1 + \left(\frac{x-\mu}{\sigma}\right)^2} \cdot \frac{1}{\sigma}$$

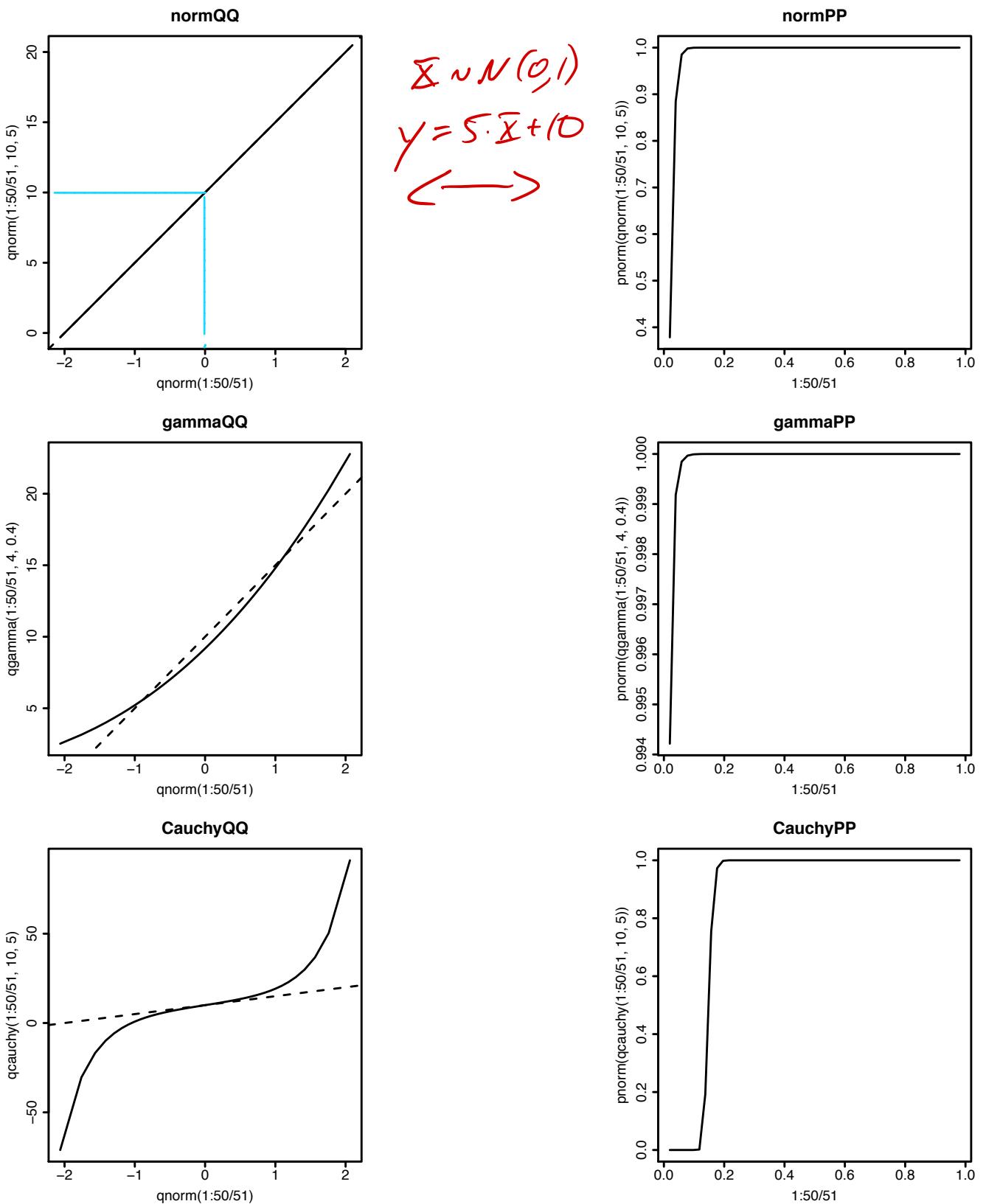
We first compare all three theoretical distributions with the  $N(0, 1)$  distribution.

```
> par(mfrow=c(3,2), pty="s", ps=20, cex=0.25)
> plot(qnorm(1:50/51), qnorm(1:50/51, 10, 5), type="l", main="normQQ")
> abline(c(10, 5), lty="dashed")
> plot(1:50/51, pnorm(qnorm(1:50/51, 10, 5))), type="l", main="normPP")
> plot(qnorm(1:50/51), qgamma(1:50/51, 4, .4), type="l", main="gammaQQ")
> abline(c(10, 5), lty="dashed")
> plot(1:50/51, pnorm(qgamma(1:50/51, 4, .4))), type="l", main="gammaPP")
> abline(c(10, 5), lty="dashed")
> plot(qnorm(1:50/51), qcauchy(1:50/51, 10, 5), type="l", main="CauchyQQ")
> abline(c(10, 5), lty="dashed")
> plot(1:50/51, pnorm(qcauchy(1:50/51, 10, 5))), type="l", main="CauchyPP")
> abline(c(10, 5), lty="dashed")
>
```

$$\{(F(\rho), G(\rho)) : \rho \in \mathbb{R}\}$$

$$\{(\rho, F(G(\rho))) : \rho \in \mathbb{R}\}$$

$$F(z) = \Phi(z) = \int_{-\infty}^z e^{-t^2/2} dt$$



# P-P and Q-Q plots: sample size $n = 50, 1000$

User: Anon

May 15, 2017

Samples of size  $n = 50$  are taken from the same 3 distributions as before:

- $N(10, \sigma^2 = 25)$
- A gamma distribution with shape  $a = 4$  and rate  $\lambda = 0.4$  which has the same mean (location) and standard deviation (scale) as the normal  $N(10, \sigma^2 = 25)$  distribution
- Cauchy distribution with location  $\mu = 10$  and scale  $\sigma = 5$  again:

$$f(x) = \frac{1}{\pi} \cdot \frac{1}{1 + \left(\frac{x-\mu}{\sigma}\right)^2} \cdot \frac{1}{\sigma}$$

All samples are then visually compared against the standard normal distribution and the normal sample is also compared against the correct  $N(10, \sigma^2 = 25)$  distribution. These plots are then followed by those generated from samples of size  $n = 1000$ .

```
> n = 50
> x = rnorm(n, 10, 5)
> d = sqrt(12*25)/2
> z = rgamma(n, 4, .4)
> y = rcauchy(n, 10, 5)
> par(mfrow=c(3,2), pty="s", ps=20, cex=0.25)
> plot(1:n/(n+1), pnorm(sort(x), 10, 5), main="PPnormal (correct mean/sd)")
> abline(c(0,1), lty="dashed")
> plot(1:n/(n+1), pnorm(sort(x), 0, 1), main="PPnormal (incorrect mean/sd)")
> abline(c(0,1), lty="dashed")
```

```
> plot(qnorm(1:n/(n+1), 10, 5), sort(x), main="QQnormal (correct mean/sd)")
> abline(c(0,1), lty="dashed")
> plot(qnorm(1:n/(n+1)), sort(x), main="QQnormal (incorrect mean/sd)")
```

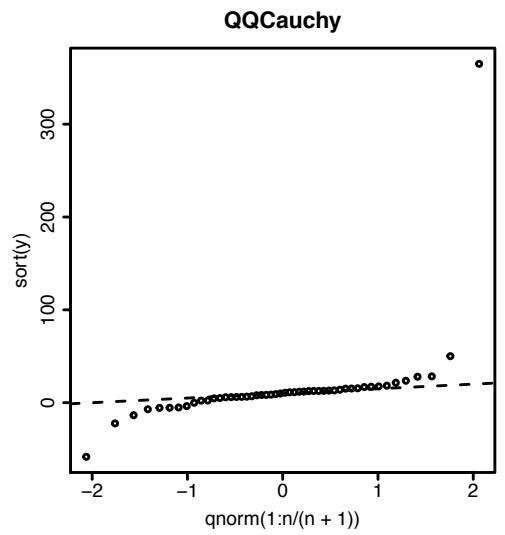
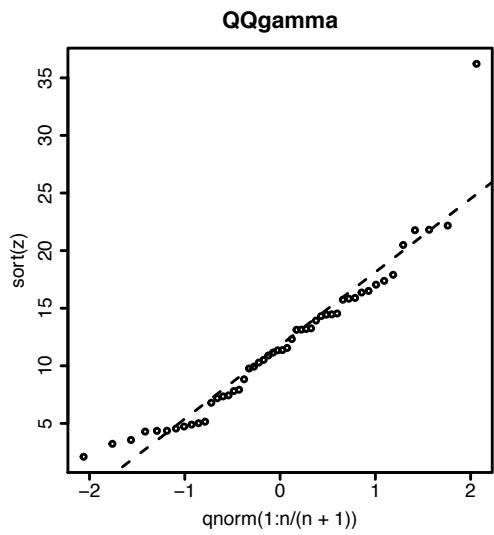
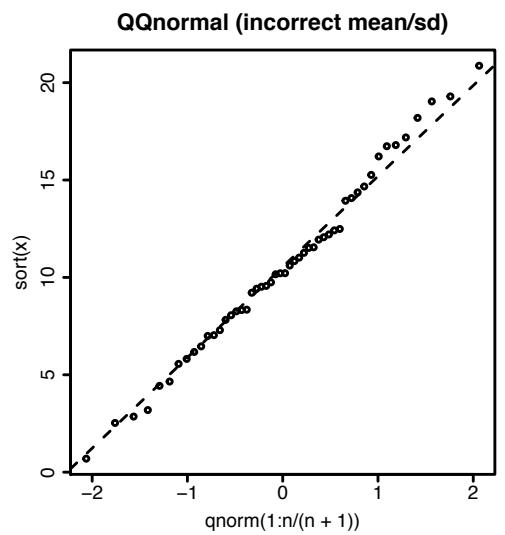
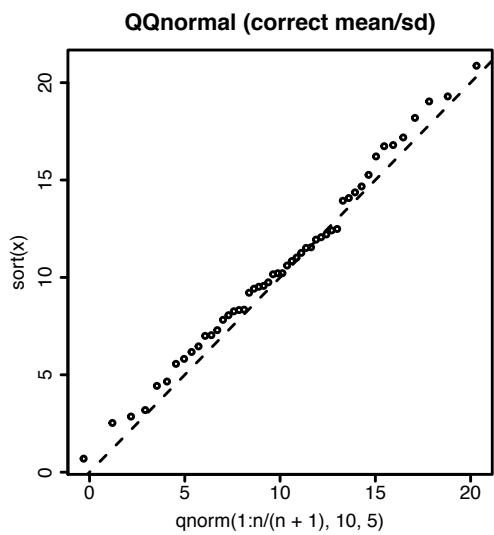
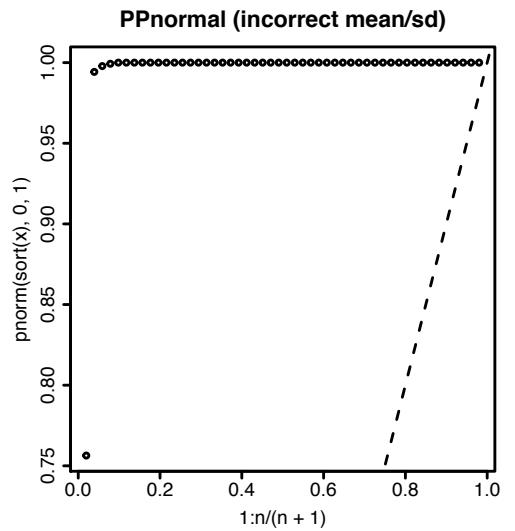
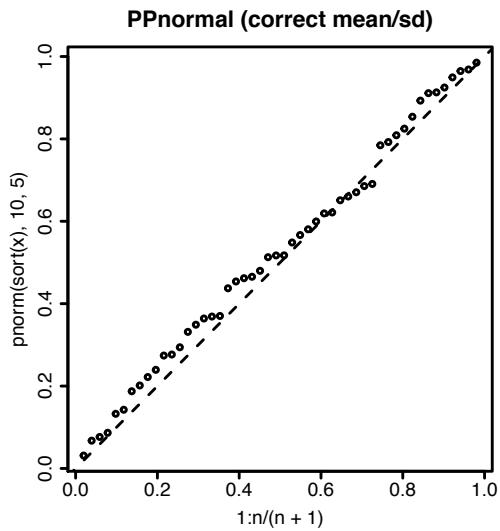
```
> abline(c(mean(x), sd(x)), lty="dashed")
> plot(qnorm(1:n/(n+1)), sort(z), main="QQgamma")
> abline(c(mean(z), sd(z)), lty="dashed")
```

```
> plot(qnorm(1:n/(n+1)), sort(y), main="QQCauchy")
> abline(c(10, 5), lty="dashed")
>
```

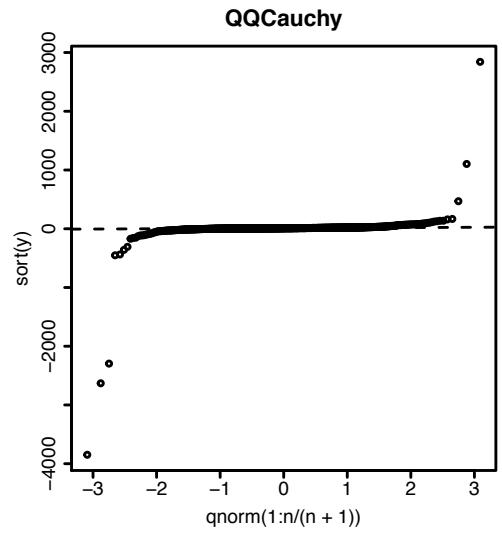
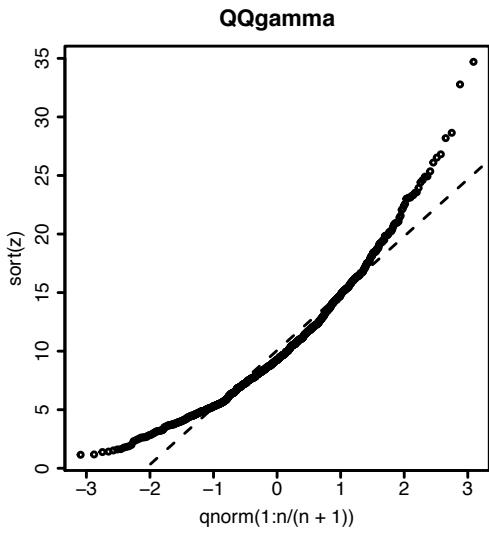
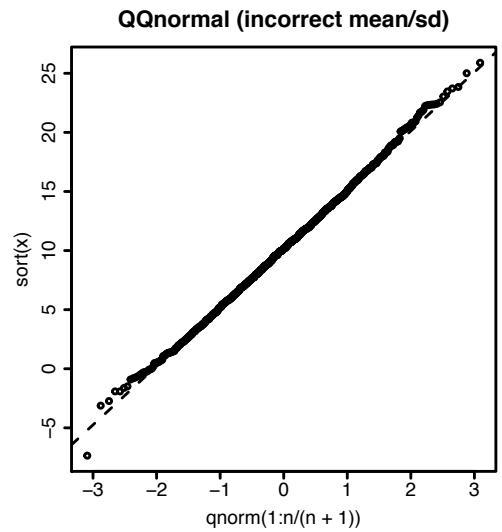
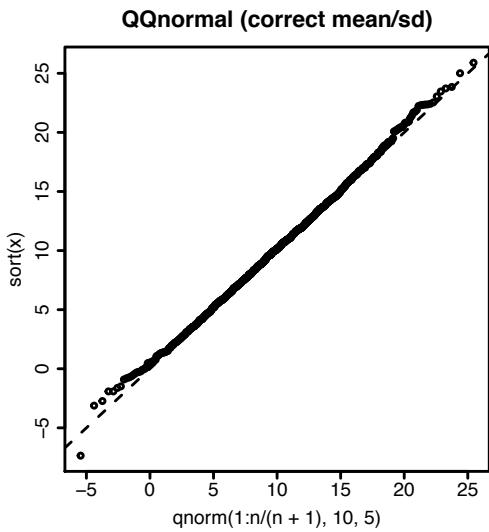
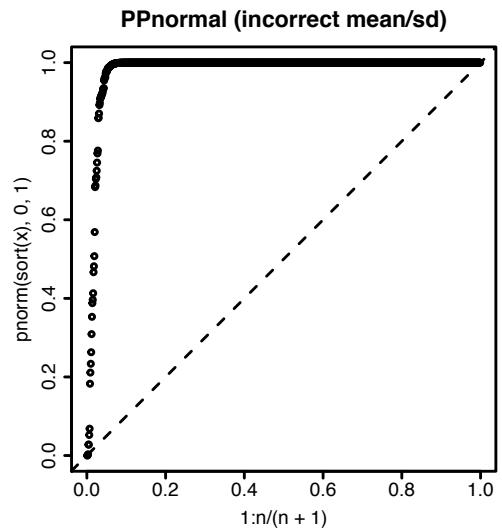
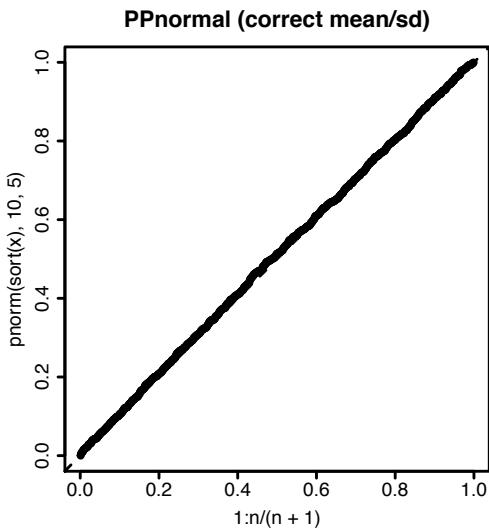
$$\left\{ \left( \frac{i}{n+1}, F(Y_{(i)}) \right) : i=1, \dots, n \right\}$$

$$\left\{ \left( F^{-1}\left( \frac{i}{n+1} \right), Y_{(i)} \right) : i=1, \dots, n \right\}$$

$$F(t) = \Phi(t) = \int_{-\infty}^t e^{-t^2/2} dt$$

*n=50*

(omitted repeated code)

 $n=1000$ 

## Extrema and Order Statistics

Let  $X_1, \dots, X_n$  be iid RVs with CDF  $F$ .

Let  $\cup = X_{(n)} = \max \{X_1, \dots, X_n\}$

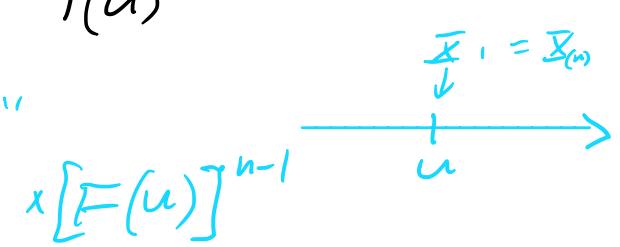
$$V = \underline{X}_{(1)} = \min \{X_1, \dots, X_n\}$$

What are the dists of  $\cup$  and  $V$ ?

$$\begin{aligned} F_\cup(u) &= P(\cup \leq u) \\ &= P(X_1 \leq u, \dots, X_n \leq u) \\ &= \prod_i^n P(X_i \leq u) \\ &= [F(u)]^n \end{aligned}$$

If  $F$  has a density  $f$ , then by chain rule so does  $F_\cup$  and

$$f_\cup(u) = n [F(u)]^{n-1} f(u)$$

Intuitively : " $f_\cup(u) du$ " = " $f(u) du$ " 

$$\times [F(u)]^{n-1}$$

$$\times n \quad (\text{which } i=(n))$$

Similarly ,  $F_V(v) = 1 - P(V > v)$

$$\begin{aligned} &= 1 - P(X_1 > v, \dots, X_n > v) \\ &= 1 - \prod_i^n P(X_i > v) \\ &= 1 - [1 - F(v)]^n \end{aligned}$$

If  $F$  has a density  $f$ , then by chain rule

$$f_V(v) = n [1 - F(v)]^{n-1} f(v) .$$

More generally, let  $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$  be the sorted values of  $X_1, \dots, X_n$ .

The  $X_{(i)}$  are called the *order statistics*, so  $X_{(1)}$  is the first,  $X_{(i)}$  is the i<sup>th</sup>, and  $X_{(n)}$  is the last order statistic.

Claim. If  $F$  is a cont. function (so  $P(X=x)=0$ ) then with prob. 1  $X_{(1)} < X_{(2)} < \dots < X_{(n)}$ .

Proof. exercise (easy if  $F$  has a density)

What is the dist. of  $X_{(k)}$ ?