

# Histogram

Show the data distribution and suggests possible outliers. Its shape is similar to the population pdf/pmf, especially if the sample size is large.

**Frequency histogram:** Consists of bars, one over each bin, whose heights represent the *number* of observations in the bins.

**Relative frequency histogram:** Consists of bars, one over each bin, whose heights represent the *proportion* of observations in the bins.

**How to construct a histogram?**

- effect of number of bins (too many or too few)
  - *Don't change the defaults in R*
  - *unless there is a good reason*
- bins of unequal sizes

Find a way to  
make a rel.  
hist. freq. wrt.  
in R.

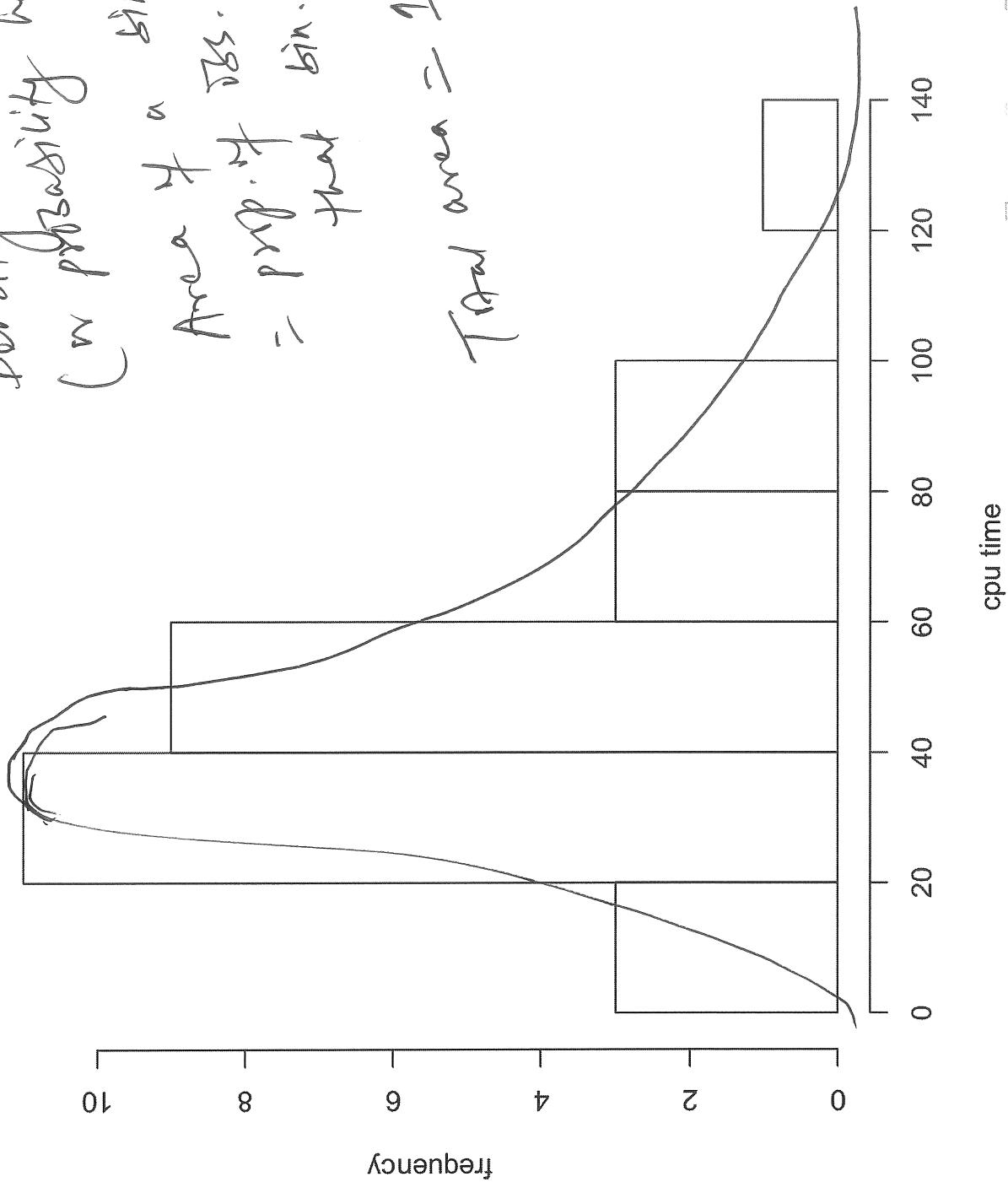
```
# frequency histogram by default
hist(cpu, xlab="cpu time", ylab="frequency", freq=TRUE,
      main="histogram of cpu data")
# probability (density) histogram
hist(cpu, freq=FALSE, xlab="cpu time",
      ylab="density", main="histogram of cpu data")
```

## frequency histogram of cpu data

Density histogram:  
(or probability histogram):

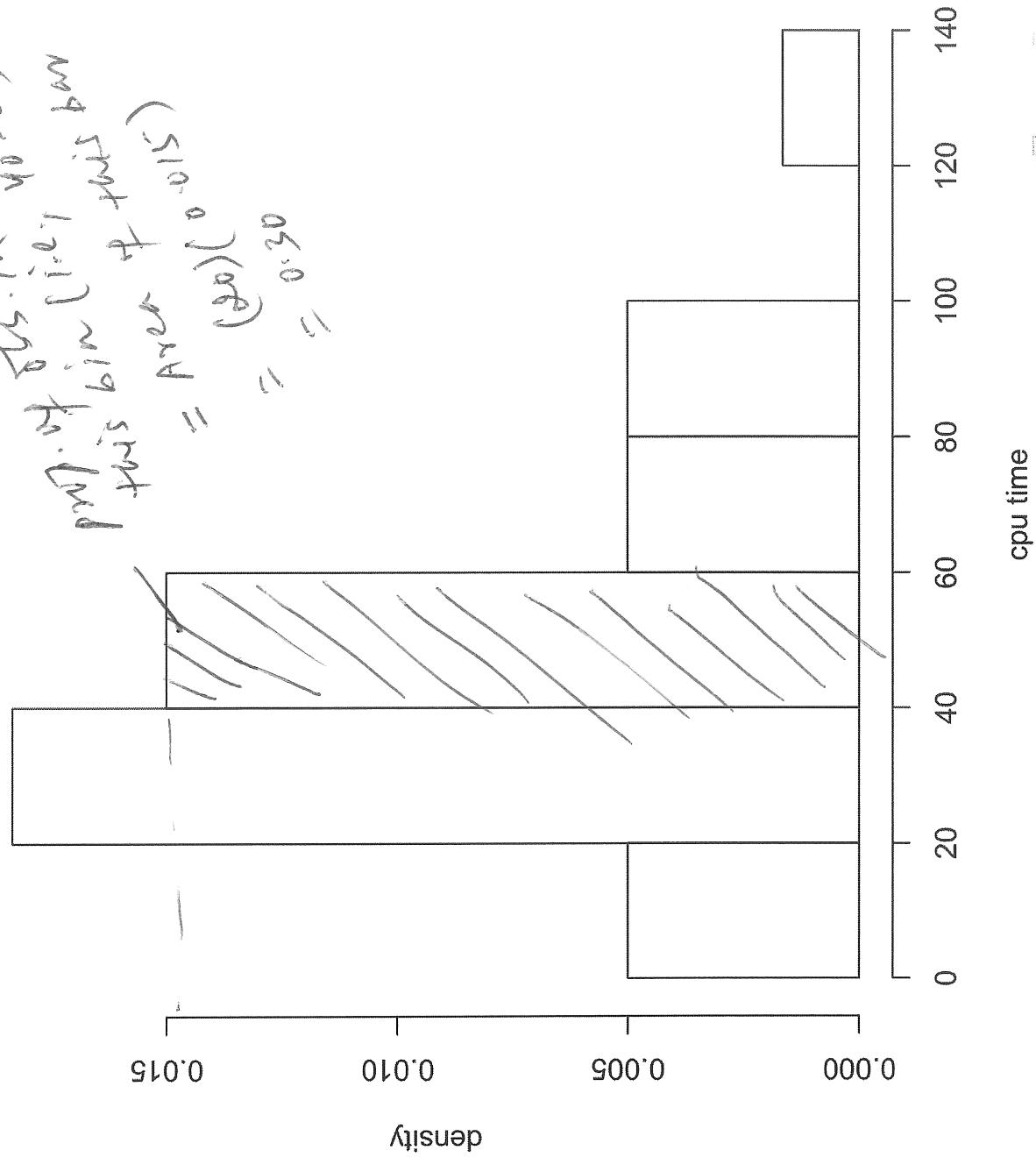
Area of a bin  
= prop. of bins.  
= prob. that bin.

Total area = 1.



## probability (density) histogram of cpu data

$$\begin{aligned} \text{prob. of } & \text{ this bin (i.e., } 40-60) \\ \text{this bin} & + \text{ this bin} \\ & = \text{Area} + 0.015 \\ & = 0.30 \end{aligned}$$



## Histograms of some simulated data:

---

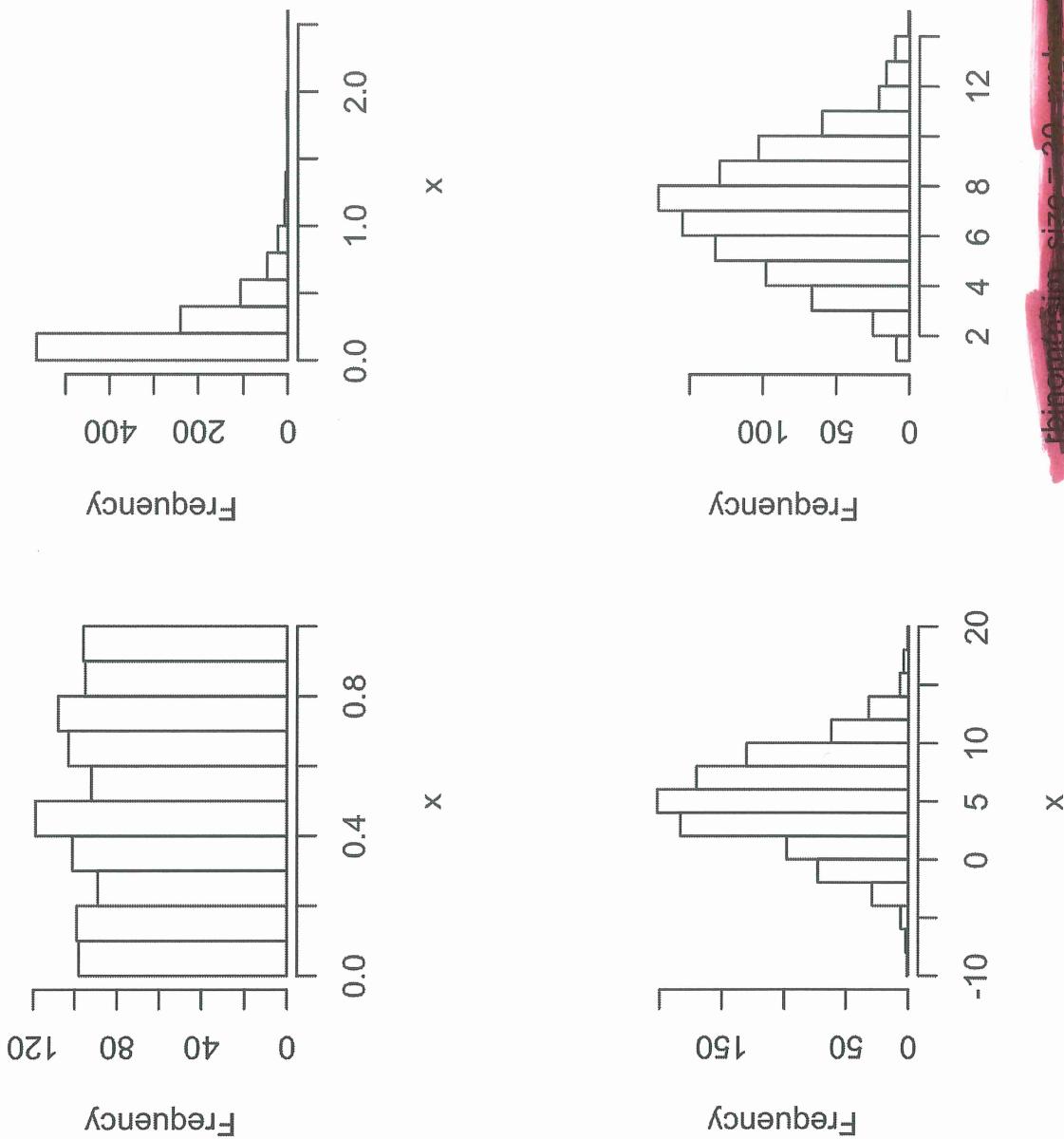
```
nsim <- 1000
# uniform (0,1) distribution
par(mfrow=c(2,2))
hist(runif(nsim), xlab="x", main="")

# exponential (lambda = 4) distribution
hist(rexp(nsim, rate=4), xlab="x", main="")

# normal (mu=5, sigma^2=16) distribution
hist(rnorm(nsim, mean=5, sd=4), xlab="x", main="")

# binomial (n=30, p=0.25)
hist(rbinom(nsim, size=30, prob=0.25), main="")

par(mfrow=c(1,1))
```



Why does the last histogram have a "normal shape?"

This is because of CLT.

$$\text{Bin}(n, p) \approx \text{Normal}(\mu = np, \sigma^2 = np(1-p)),$$

provided  $n$  is large.

$$np \geq 5 \text{ and } n p(1-p) \geq 5$$

Rule of thumb:

# QQ Plot

Plot quantiles of one dataset against quantiles of another dataset (from a known distribution with cdf  $F$ ). If the points fall on a straight line, the distribution  $F$  may be a good fit to the data — allows a graphical check of how well  $F$  fits the data.

**Data:**  $x_1, \dots, x_n$  (a random sample)

**Sorted data:**  $x_{(1)}, \dots, x_{(n)}$

- These are sample quantiles or “order statistics.”
- They estimate population quantiles of the distribution  $F$ .  
*(Is this correct? Is this model for data?)*

**Q:** What are the associated probabilities?

- Each sample observation has  $1/n$  probability weight under the empirical distribution.
- The sample quantiles  $x_{(1)}, x_{(2)}, \dots, x_{(n-1)}, x_{(n)}$  are associated with probabilities  $1/n, 2/n, \dots, (n-1)/n, n/n$ .
- $x_{(i)}$  estimates  $(i/n)$ th population quantile, i.e.,  $F^{-1}(i/n)$ ,  
 $i = 1, \dots, n$ .

$$\frac{p}{n} \quad x_{(np)} \quad F^{-1}(p)$$

$\uparrow$   
Sample quantile with  $p = \frac{1}{n}$

$$\frac{2}{n} \quad x_{(2)} \quad \dots \quad p = \frac{2}{n} \quad F^{-1}(p = \frac{2}{n})$$

$$\frac{n-1}{n} \quad x_{(n-1)} \quad \dots \quad p = \frac{n-1}{n} \quad F^{-1}(p = \frac{n-1}{n})$$

↓  
 $F^{-1}(p=1)$  is very large.

$$x_{(n)} \approx p = \frac{n}{n} = 1$$

QQ plot: Plot the following pairs of points:  $(x(i), F^{-1}(i/n)),$   
 $i = 1, \dots, n.$

**Problem:**  $F^{-1}(1)$  may be  $\infty$ .

**Solution:** Consider an offset  $a$ .

Old probabilities:  $\frac{1}{n}, \frac{2}{n}, \dots, \frac{n-1}{n}, \frac{n}{n}$

New probabilities:  $\frac{1-a}{n+1-2a}, \frac{2-a}{n+1-2a}, \dots, \frac{n-1-a}{n+1-2a}, \frac{n-a}{n+1-2a}$

- Default in R:  $a = 3/8$  if  $n \leq 10$  and  $a = 1/2$  if  $n > 10$ .

- qqplot gives a general QQ plot
- qqnorm gives normal QQ plot — it uses  $N(0, 1)$  distribution as  $F$

(with offset)

**Q:** What are the probabilities for  $n = 30?$

$$\frac{1 - \frac{1}{2}}{30} = \frac{0.5}{30}, \quad \frac{1.5}{30}, \quad \frac{2.5}{30}$$

with offset  
 $(30 + 1 - 1)$

Recall:

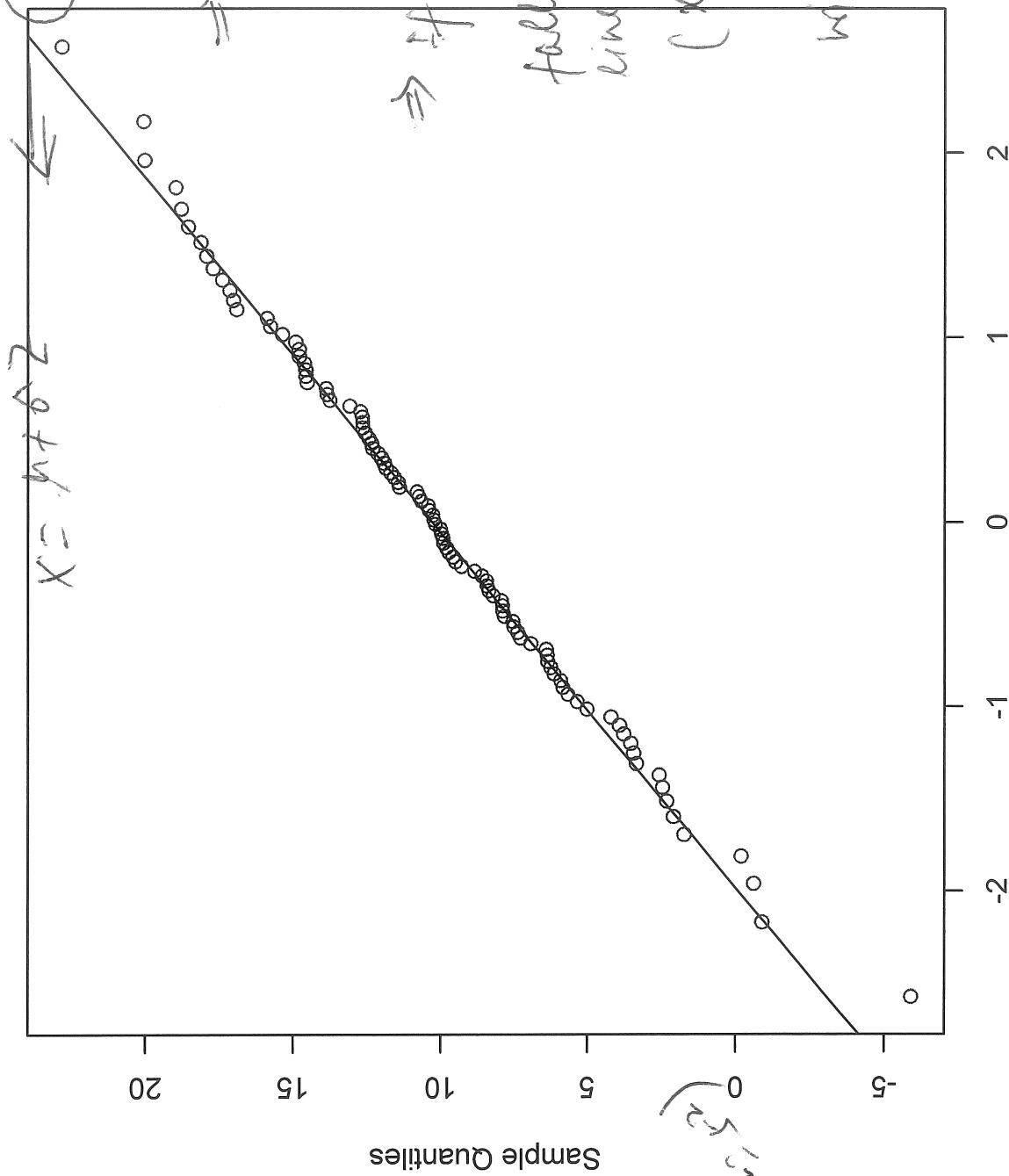
$$\text{If } X \sim N[\mu, \sigma^2],$$
$$Z = \frac{X - \mu}{\sigma} \sim N[0, 1].$$

$$x_p = \mu + z_p \sigma$$

pop. Quantiles  
( $x_{(i)}$ , of  $z$ )

fall on a straight line, then pop. Quantiles ( $x_{(i)}$ , of  $z$ ) will also fall in a straight line.

### Normal Q-Q Plot

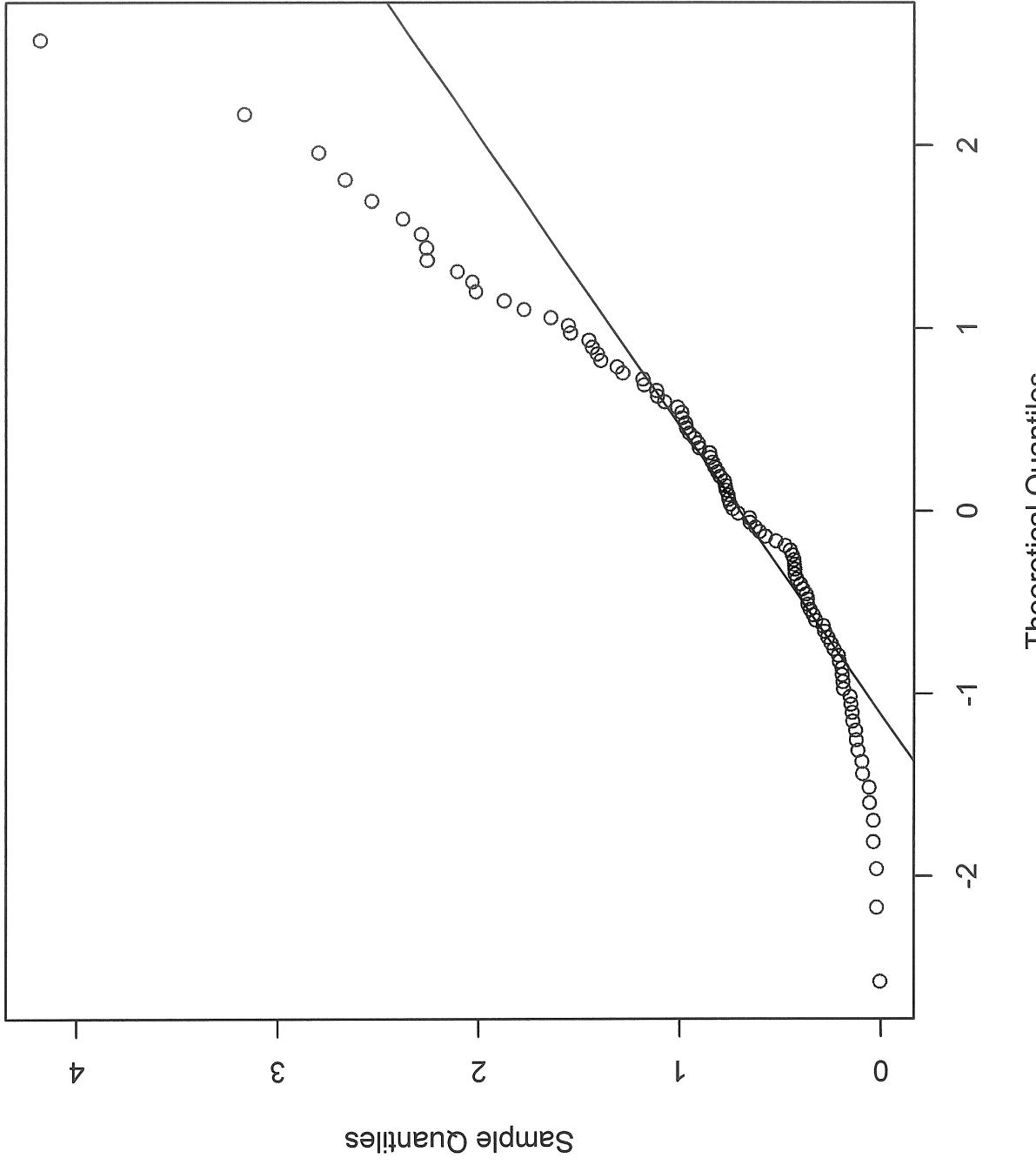


Quantiles  
Quantile Function  
Cumulative Distribution Function  
CDF

Theoretical Quantiles

\ quantiles of  $F_N(z)$

## Normal Q-Q Plot



# R Code for QQ plots

```
# QQ plot 1  
x <- rnorm(100, 10, 5)  
qqnorm(x)  
qqline(x)  
  
# QQ plot 2  
x <- rexp(100, 1)  
qqnorm(x)  
qqline(x)
```

**Time series plot:** Plot of a data on a variable against time — shows how the variable changes over time.

```
# Data from Exercise 8.5
year <- seq(from=1790, to=2010, by=10)
# > year
# [1] 1790 1800 1810 ...
# >
uspop <- c(3.9, 5.3, 7.2, 9.6, ..., 281.4, 308.7)

plot(year, uspop, ylab="Population (in millions)",
     main="US population since 1790")
```

**Scatterplot:** Plot of one variable ( $X$ ) against another variable  $Y$  — shows the relationship between the two variables. See Figure 8.11 of the textbook.



## Point estimation (Chapter 9)

**Problem:**  $X \sim f_\theta(x)$ , where  $\theta$  is an unknown parameter. This  $\theta$  may be a vector.

**Data:**  $X_1, \dots, X_n$  — a random sample of  $X$ .

We have seen a number of descriptive statistics and what they estimate. But the choice of an  $\hat{\theta}$  of  $\theta$  may not be obvious.

**A general method of parameter estimation:** Method of maximum likelihood. It has generally good properties.

Population:  $X \sim f_\theta(x)$  ↗ pdf  
 $\theta$  = unknown parameter

↓

$X_1, X_2, \dots, X_n$  (ps)  
Estimate  $\theta$  by  $\hat{\theta}$  which is called the sample data  
Goal:

# Method of Maximum Likelihood

Likelihood function of data: Joint pdf or pmf of sample data considered as a function of  $\theta$  with data held fixed at the observed values  $X_1 = x_1, X_2 = x_2, \dots, X_n = x_n$ .

$$L(\theta) = L(\theta; x_1, x_2, \dots, x_n) = \prod_{i=1}^n f_\theta(x_i)$$

indep.

joint pdf/pmf:  $\theta$  is fixed.  
joint of this a fn.  
of  $x_1, \dots, x_n$   
likelihood fn:  $x_1, \dots, x_n$  fixed,  
think of this as  
a fn. of  $\theta$ .

- A function of  $\theta$  — the data are held fixed.

Maximum likelihood estimator (MLE) of  $\theta$ : The value  $\hat{\theta}$  of  $\theta$  that maximizes the likelihood function as a function of  $\theta$ .

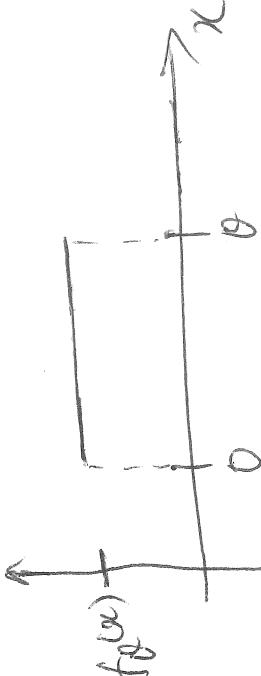
- Can think of MLE as the value of  $\theta$  that is “most likely” to have led to the observed data.
- Essentially a calculus problem.

# How to find MLE?

Direct approach: Directly maximize the likelihood function.

**Ex:** Let  $X_1, X_2, \dots, X_n$  represent a random sample from a Uniform  $(0, \theta)$  distribution where  $\theta > 0$ . Find the MLE of  $\theta$ .

Recall:

$$f_{\theta}(x) = \begin{cases} \frac{1}{\theta}, & 0 < x \leq \theta \\ 0, & \text{otherwise} \end{cases}$$


PDF of  $X_i$ , which represents the population

which is max. in the population.  
will see that:  $\hat{\theta} = \max_i X_i$  in the sample is MLE of  $\theta$ .

Note:

$$\mathbb{E}[X] = \frac{\theta}{2} \Rightarrow \theta = \frac{2\mathbb{E}[X]}{\theta} \quad \text{also } \hat{\theta} = 2\bar{X} \text{ can be used as an estimate of } \theta.$$

MLE =  $\hat{\theta}_{MLE} = \bar{X}(n)$  which one is better?  
 $\hat{\theta}_{MLE} = 2\bar{X}$