

W1211 Introduction to Statistics

Lecture 11

Wei Wang

Oct 10th, 2012

Uniform RV

- We call a uniform rv U a **standard uniform**, if and only if $U \sim \text{uniform on } [0,1]$
- For a standard uniform rv U , we can easily calculate,

$$E(U) = \int_0^1 x \cdot 1 dx = \frac{1}{2}$$

$$E(U^2) = \int_0^1 x^2 \cdot 1 dx = \frac{1}{3}$$

$$\text{Var}(U) = E(U^2) - [E(U)]^2 = \frac{1}{12}$$

General Uniform

- Note that a general case of uniform distribution X on $[A, B]$ can be treated as a linear transform of a standard uniform, i.e., $X = (B - A)U + A$.
- Proposition:

If X is a continuous uniform rv on $[A, B]$, then
 $E(X) = (B + A)/2$, $\text{Var}(X) = (B - A)^2/12$

- R command: `dunif(x, min=0, max=1),`
`punif(q, min=0, max=1),`
`qunif(p, min=0, max=1).`

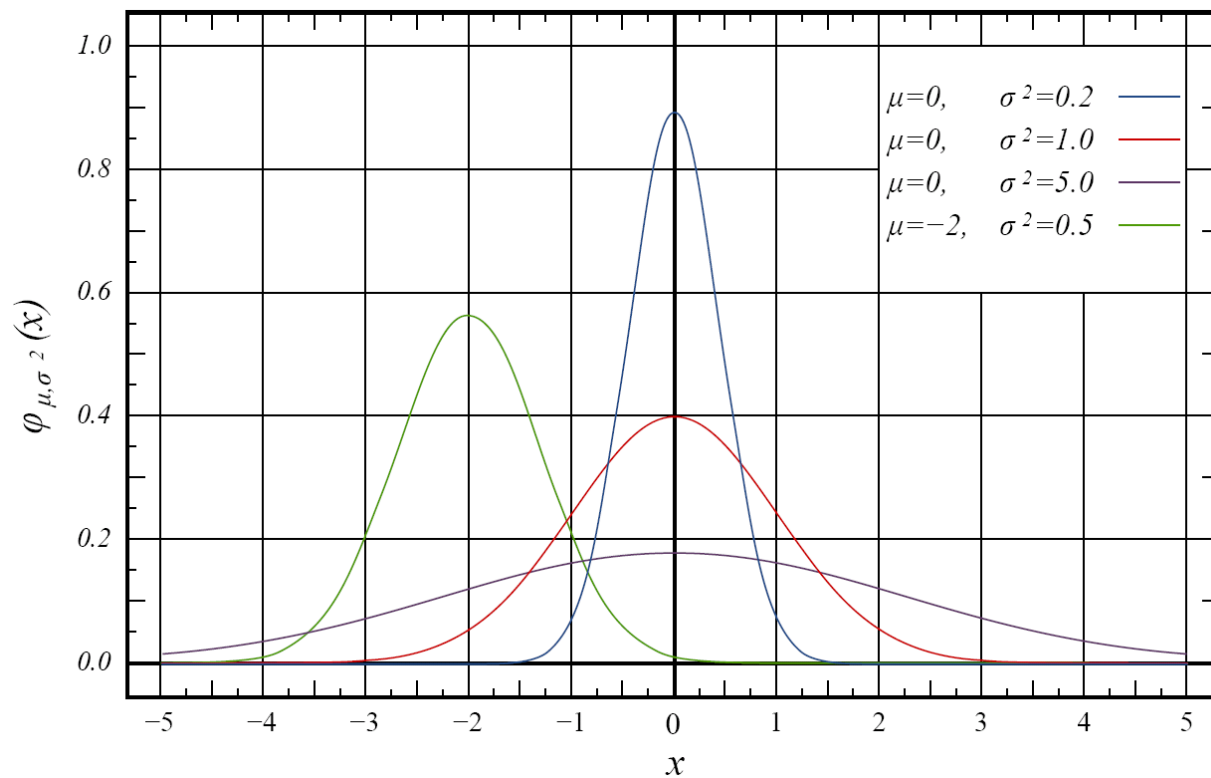
The Normal Distribution

- It's probably the most important distribution in the world!
- Many numerical populations have distributions that can be fit very closely by an appropriate normal curve. (people's height/weight; testing scores; etc.) Even when the underlying distribution is discrete, (yearly number of customers to Wal-Mart; etc.) the normal curve often gives an excellent approximation.
- A continuous rv is said to have a normal (Gaussian) distribution with parameters μ and σ , where $-\infty < \mu < \infty$, and $0 < \sigma$, if the pdf of X is

$$f(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/(2\sigma^2)} \quad -\infty < x < \infty$$

The Normal pdf

- Normal distribution is a **bell-shaped**, **single peaked** and **symmetric** distribution.



Parameters

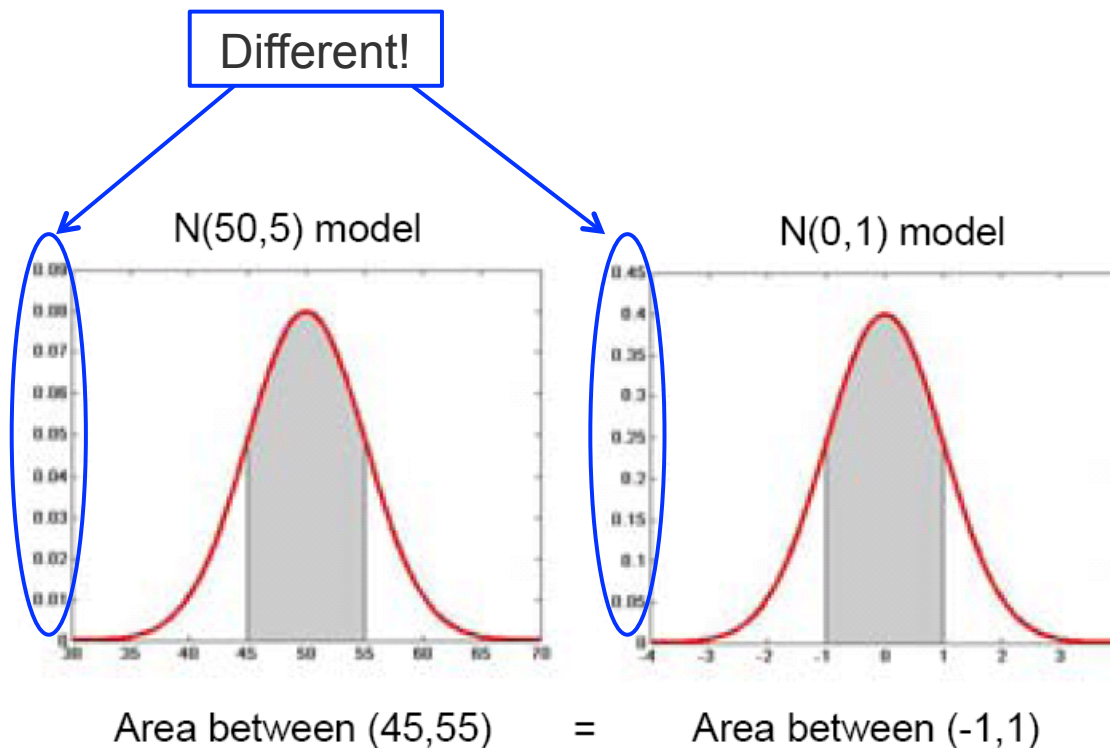
- Clearly $f(x; \mu, \sigma) \geq 0$, but a somewhat complicated calculus argument must be used to verify that

$$\int_{-\infty}^{\infty} f(x; \mu, \sigma) dx = 1.$$

- Parameter μ , stands for the **expected value** of the normal distribution.
Exercise: show that if $X \sim N(\mu, \sigma^2)$, then $E(X) = \mu$.
- Parameter σ , stands for the **standard deviation** of the normal distribution.
Exercise: show that if $X \sim N(\mu, \sigma^2)$, then $\text{Var}(X) = \sigma^2$.

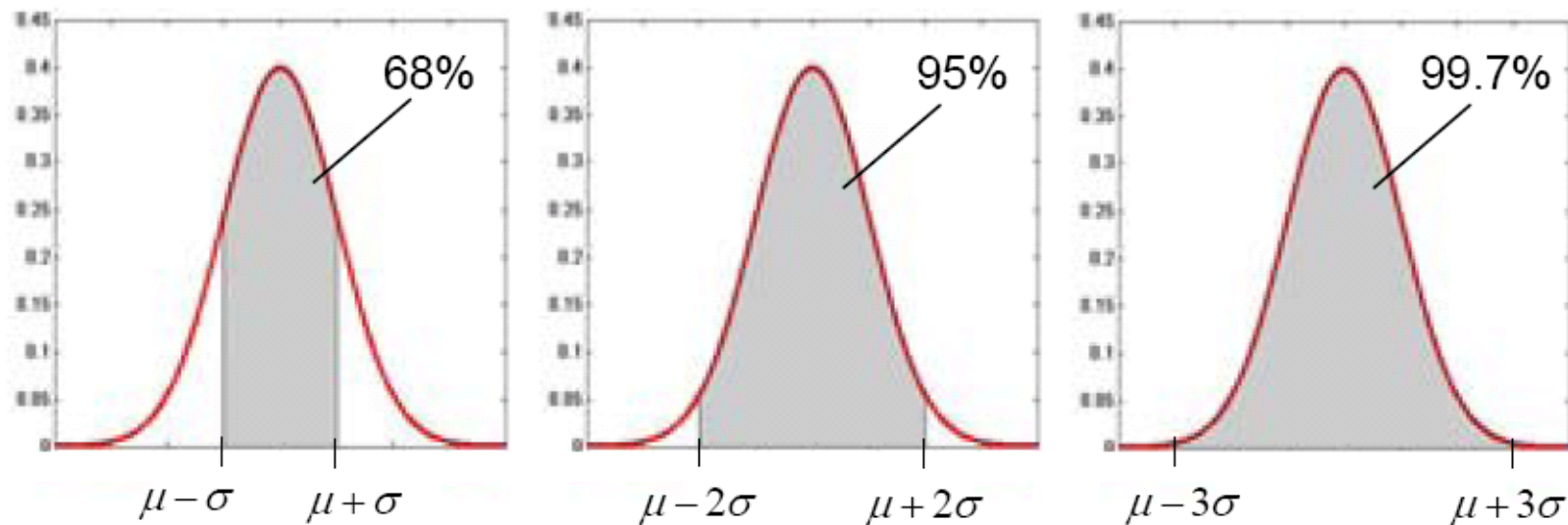
Basic Properties

- All normal models have the same shape and the same area within x standard deviations of its mean.



The 68-95-99.7 Rule

- For any normal distribution, we have the following result:



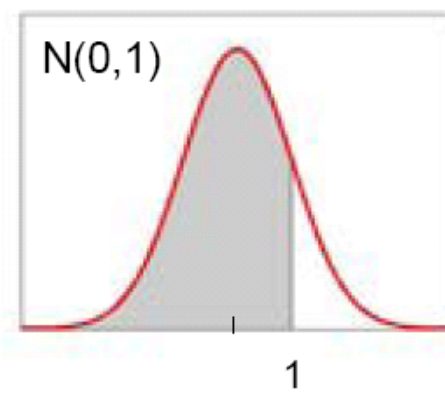
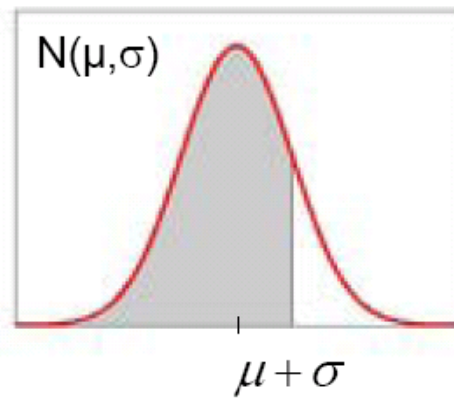
Example

Ex. On an exam the scores followed an approximate normal model with $\mu = 72$ and $\sigma = 8$.

- 68% of the students scored between 72 ± 8 or (64, 80).
- 95% of the scores were between $72 \pm 2 \cdot 8$ or (56, 88).
- 99.7% of the scores were between $72 \pm 3 \cdot 8$ or (48, 96).

- What proportion scored below 84?

Key Result



$$area\{y < \mu + \sigma\} = area\{z < 1\}$$

Standard Normal

- If $Z \sim N(0, 1)$, i.e., if Z is a normal random variable with $\mu=0$, $\sigma=1$. Then Z is said to have a **standard normal distribution**.
- Any normally distributed rv's could be obtained by using standard normal rv's. To put it more mathematically, if $X \sim N(\mu, \sigma^2)$, then X could be written as

$$X = \mu + \sigma \cdot Z$$

where Z is a standard normal rv.

- Conversely, if $X \sim N(\mu, \sigma^2)$, then

$$Z = (X - \mu) / \sigma$$

has a **standard normal distribution**. And Z is often called the “**z-score**” of X .

Example cont.

Ex. The exam scores followed a $N(72,8)$ model.

What proportion of the students scored below 84?

$$z = \frac{y - \mu}{\sigma} = \frac{84 - 72}{8} = 1.5$$

Answer: 93.32%

[illegible]

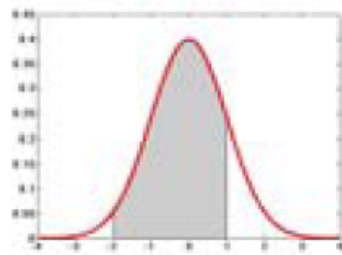
Simplification

- Thus, any problem about any normal rv $X \sim N(\mu, \sigma^2)$, can be **translated** to a problem about a standard normal rv Z .

Ex. $P(a \leq X \leq b) = P[(a-\mu)/\sigma \leq (X-\mu)/\sigma \leq (b-\mu)/\sigma] = P[(a-\mu)/\sigma \leq Z \leq (b-\mu)/\sigma]$.

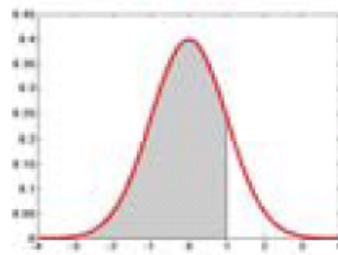
- The cumulative distribution function of standard normal distribution, that is $\Phi(z) = P(Z \leq z)$, is already known! (Appendix Table.)
- Check Table A.3 to determine $P(Z \leq 0.76)$; $P(Z > 0.76)$; $P(-1.32 \leq Z \leq 0.76)$.
- **Question:** How to get the p -th percentile of the standard normal from A.3?

Using the Normal Table



0.8185

=

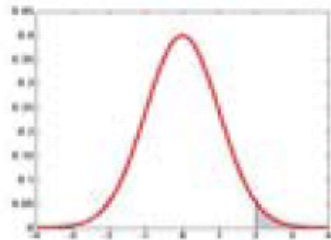


0.8413

-

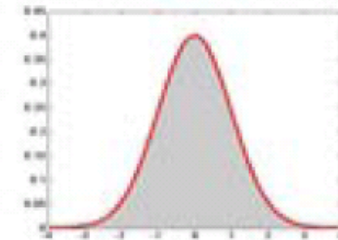


0.0228



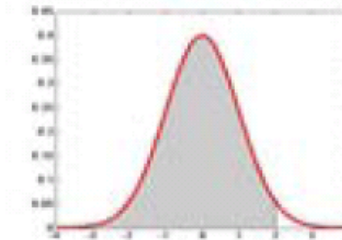
0.0228

=



1.00

-



0.9772

R instead of tables

- R command: `dnorm(x, mean = 0, sd = 1),`
`pnorm(q, mean = 0, sd = 1),`
`qnorm(p, mean = 0, sd = 1) .`

Example

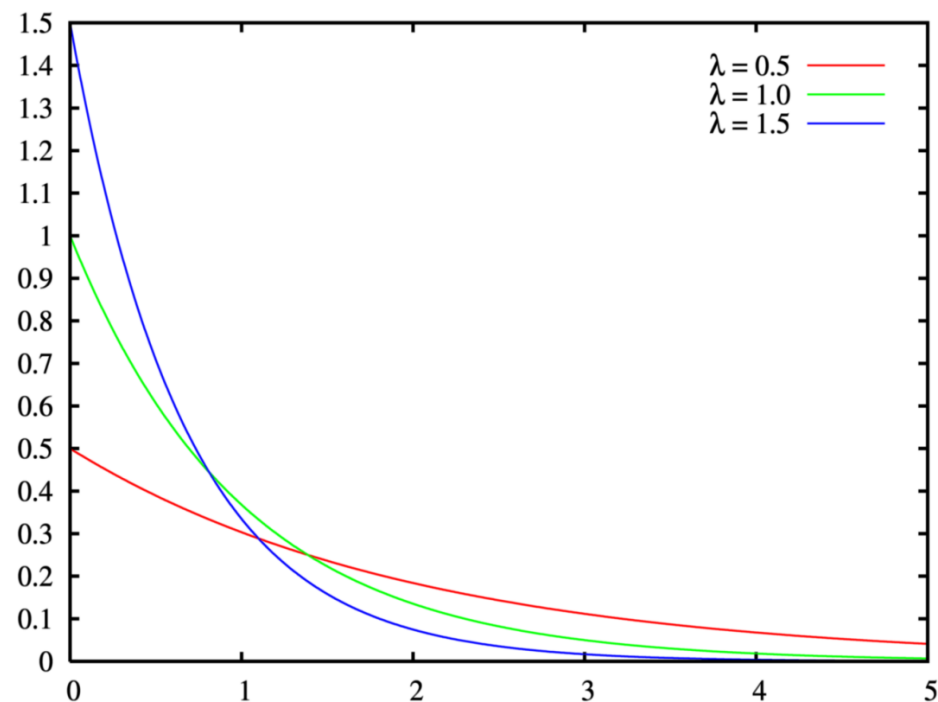
Ex. Suppose the height of all Columbia students can be described by a $N(68, 4)$ model.

1. What proportion of students is shorter than 74 inches?
2. What proportion of students is taller than 74 inches?
3. How tall does a student have to be to be among the 10% tallest students?

The Exponential Distribution

- X is said to have an exponential distribution with parameter λ ($\lambda > 0$) if the pdf of X is

$$f(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$



More on Exponential

- Note that an exponential rv X can only take positive values. And the cdf of X is

$$F(x; \lambda) = \begin{cases} \int_0^x \lambda e^{-\lambda y} dy = 1 - e^{-\lambda x} & x \geq 0 \\ 0 & x < 0 \end{cases}$$

- Thus $P(X > x) = 1 - F(x; \lambda) = e^{-\lambda x}$
- **Proposition:** (proof?)

If X is an exponential rv with parameter λ , then
 $E(X) = 1/\lambda$, $\text{Var}(X) = 1/\lambda^2$

- R command: `dexp(x, lamda=1)`,
`pexp(q, lamda=1)`,
`qexp(p, lamda=1)`.

Exponential Distribution and Poisson Distribution

- ▶ Suppose the number of events occurring in a time interval of length t has Poisson Distribution with parameter αt , and the numbers of occurrences in non-overlapping intervals are independent of one another. Then the distribution of elapsed time between the occurrence of two successive events is exponential with parameter $\lambda = \alpha$.

Example

Ex. Suppose you are waiting for a bus at a bus station. And the distribution of the length of the time you have to wait to get on the bus after you arrive at the bus station is exponentially distributed with parameter λ . Assume you have already waited for s minutes, how much longer do you expect to wait?

First, we have to figure out the conditional probability distribution of the additional waiting time given we have waited for s minutes. For any $t > 0$

$$\begin{aligned} P(X \geq s + t \mid X \geq s) &= P[(X \geq s + t) \cap (X \geq s)] / P(X \geq s) \\ &= P(X \geq s + t) / P(X \geq s) \\ &= e^{-\lambda t} \end{aligned}$$

which is again an **exponential distribution**! Thus the expected additional waiting time is $1/\lambda$.

Memoryless Property

- From the previous example, we know that if a waiting time (or lifetime of something) follows an exponential distribution, *the distribution of additional waiting time (lifetime) is exactly the same as the original distribution of waiting time (lifetime)*. In other words, the exponentially distributed waiting time does NOT remember how much time you have waited, it starts afresh at any time!
- It is popular to model the distribution of component lifetime using the exponential distribution. However, the memoryless property may not be realistic in many applied problems. More general lifetime models can be furnished by the gamma, Weibull, and lognormal distributions. (Book: p159 – p168).

Normal Probability Plot

- ▶ Because of the important role that Normal Distribution plays in statistical inference, we often want to assess whether a given sample is roughly normal distributed. Normal Probability Plot is used for this purpose.

Normal Probability Plot

- ▶ Because of the important role that Normal Distribution plays in statistical inference, we often want to assess whether a given sample is roughly normal distributed. Normal Probability Plot is used for this purpose.
- ▶ The basic strategy is to compare sample features with population features. In probability plot, we are using sample percentile(quantile) and population percentile(quantile), so it is also known as Q-Q plot.

Normal Probability Plot

- ▶ Because of the important role that Normal Distribution plays in statistical inference, we often want to assess whether a given sample is roughly normal distributed. Normal Probability Plot is used for this purpose.
- ▶ The basic strategy is to compare sample features with population features. In probability plot, we are using sample percentile(quantile) and population percentile(quantile), so it is also known as Q-Q plot.
- ▶ The definition of a normal probability plot

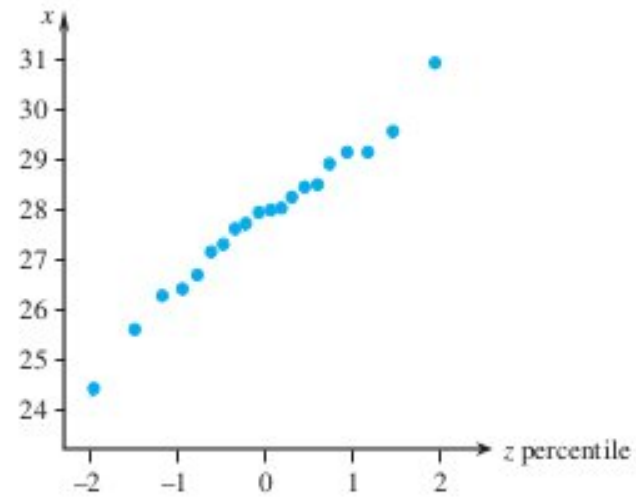
A plot of the n pairs

$([100(i - .5)/n]\text{th } z \text{ percentile}, i\text{th smallest observation})$

on a two-dimensional coordinate system is called a **normal probability plot**. If the sample observations are in fact drawn from a normal distribution with mean value μ and standard deviation σ , the points should fall close to a straight line with slope σ and intercept μ . Thus a plot for which the points fall close to some straight line suggests that the assumption of a normal population distribution is plausible.

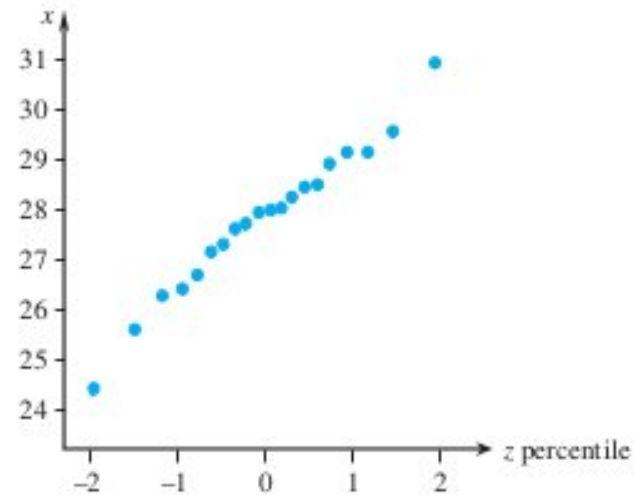
Examples of Normal Probability Plot

- ▶ A Normal Sample



Examples of Normal Probability Plot

► A Normal Sample



► Two Non-normal Samples

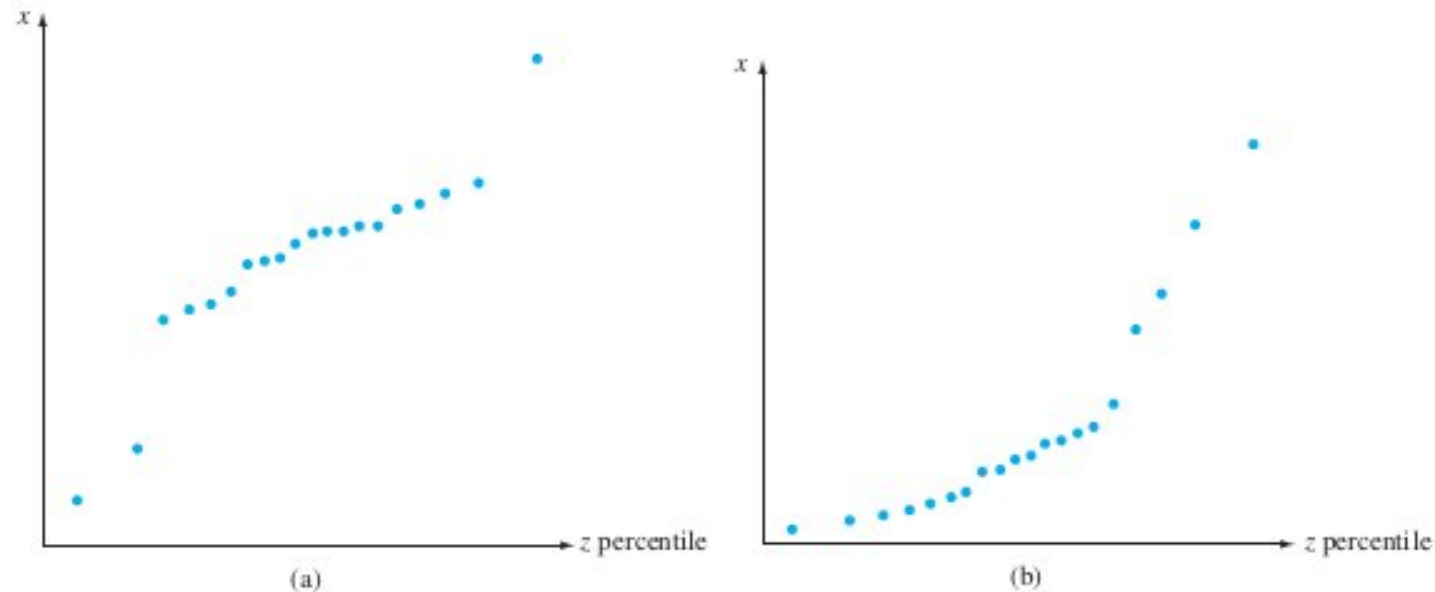


Figure 4.37 Probability plots that suggest a nonnormal distribution: (a) a plot consistent with a heavy-tailed distribution; (b) a plot consistent with a positively skewed distribution