



Medical statistics

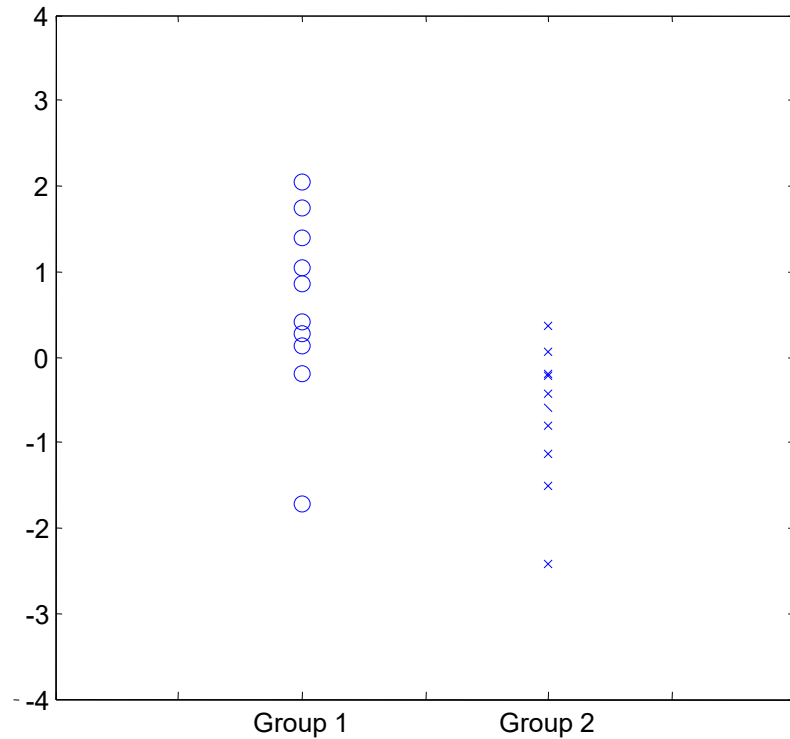
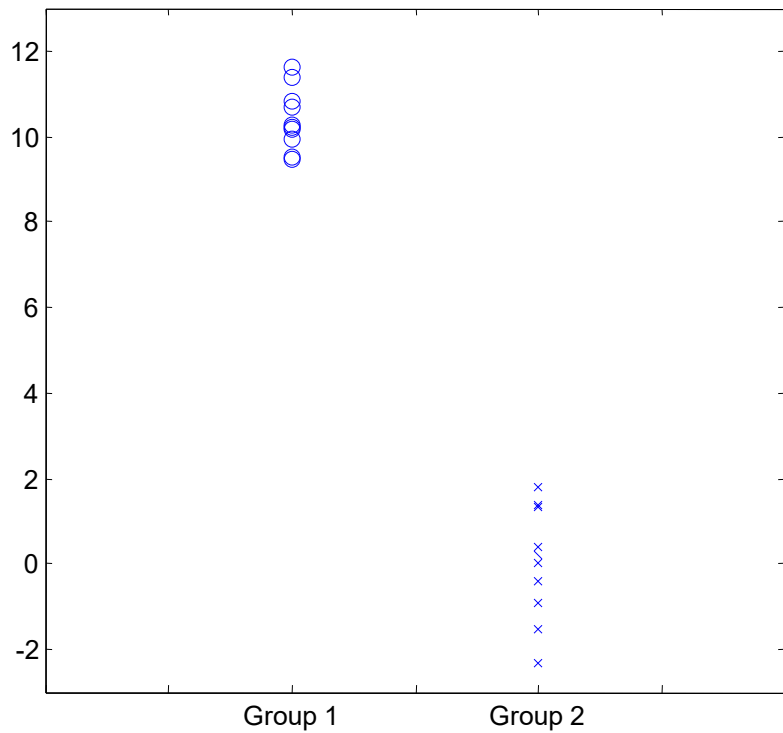
Ref: Montgomery DC, Runger GC, and Hubele NF, "Engineering statistics", 5th ed., 2010

莊子肇 副教授
中山電機系

Why statistics?

- Example: Any regional volume change of gray matter or white matter on patients with attention-deficit hyperactivity disorder (ADHD)?
 - WM/GM separation
 - Need inter-subject comparison
 - Recall the methods of image registration
 - Then?

Why Statistics?



Why statistics?

- There is always **uncertainty** in real world.
- Another example: Suppose that one of your friends, seeing that you are struggling with your course works, suggested that you quit the school, citing that “Mr. XXX made billions of dollars without even finishing elementary school!”

Would you follow the advice? Variation and risk?

Outline

- Data summary and presentation
- Random variables and probability
- Decision making
 - Parameter estimation & hypothesis testing
 - Z-test and T-test
 - One sample and two samples
 - Analysis of variance (ANOVA)

Population and sample

- Question: What is the average body temperature of monkeys in Taiwan?
- What we (engineers) do in most engineering applications:
 - Collect (all?) data
 - Predict what other untested samples will perform
 - **Statistical inference**

Population and sample

- Question: What is the average body temperature of monkeys in Taiwan?
 - **Population**: all monkeys in Taiwan
 - **Sample**: 10 age- and gender-matched monkeys from each county
 - How close is the sample mean to the population mean?
 - How close is the sample variation to the population variation?

Statistics

- Statistics is the science of **collecting**, **summarizing**, **presenting**, and **interpreting** data.
 - Reliable data
 - Consistent analysis
 - Be as simple as possible and make sense!

Data summary

- Sample mean (\bar{x})
 - Population mean (μ)
- The **sample mean** is a reasonable estimate of the **population mean**.

Data summary

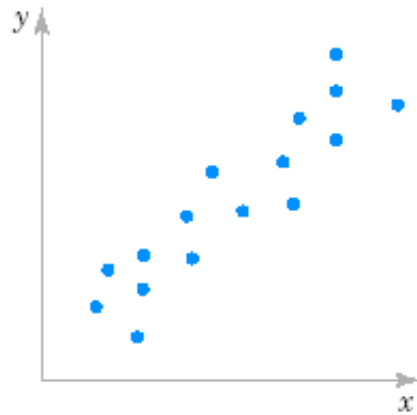
- Sample variance
 - Sample standard deviation (s)

 - Population variance
 - Population standard deviation (σ)
- The **sample variance** is a reasonable estimate of the **population variance**.

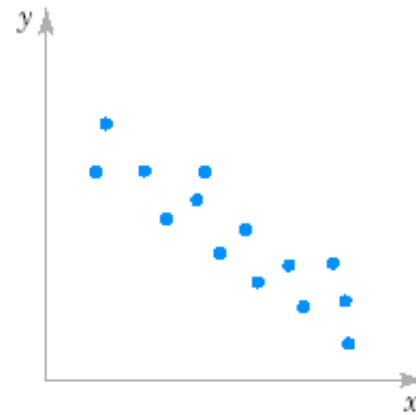
Data summary

- Coefficient of variation
- Sample correlation coefficient (r)

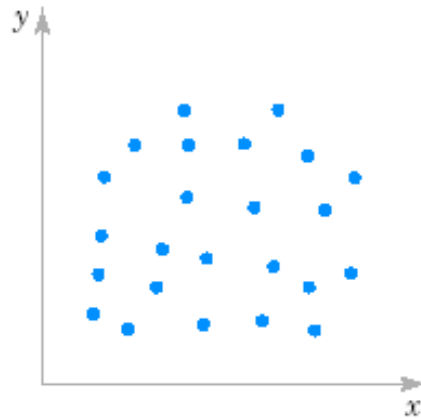
Scatter diagram of multivariate



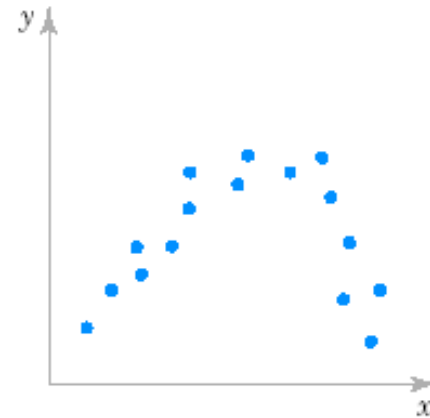
(a) r is near +1



(b) r is near -1



(c) r is near 0, y and x are unrelated



(d) r is near 0, y and x are nonlinearly related

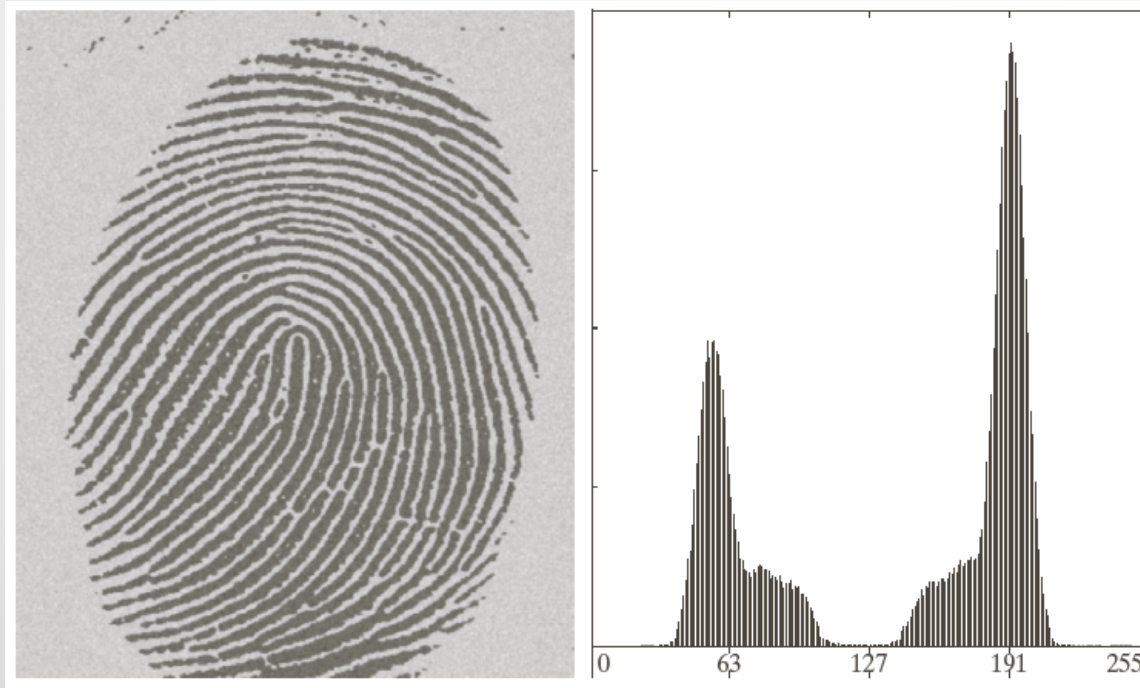
Data display: Stem-and-Leaf diagram

Barry Bonds				
Year	Age	OBP	SLG	OPS
1986	22	.330	.416	.746
1987	23	.329	.492	.821
1988	24	.368	.491	.859
1989	25	.351	.426	.777
1990	26	.406	.565	.971
1991	27	.410	.514	.924
1992	28	.456	.624	1.080
1993	29	.458	.677	1.135
1994	30	.426	.647	1.073
1995	31	.431	.577	1.008
1996	32	.461	.615	1.076
1997	33	.446	.585	1.031
1998	34	.438	.609	1.047
1999	35	.389	.617	1.006
2000	36	.450	.707	1.157
2001	37	.515	.863	1.378
2002	38	.582	.799	1.381
2003	39	.529	.749	1.278
2004	40	.609	.812	1.422
2005	41	.404	.667	1.071
2006	42	.454	.545	.999
2007	43	.480	.565	1.045
22 (avg)		.444	.607	1.051

Stemplot of Bonds' OPS ($N = 22$)

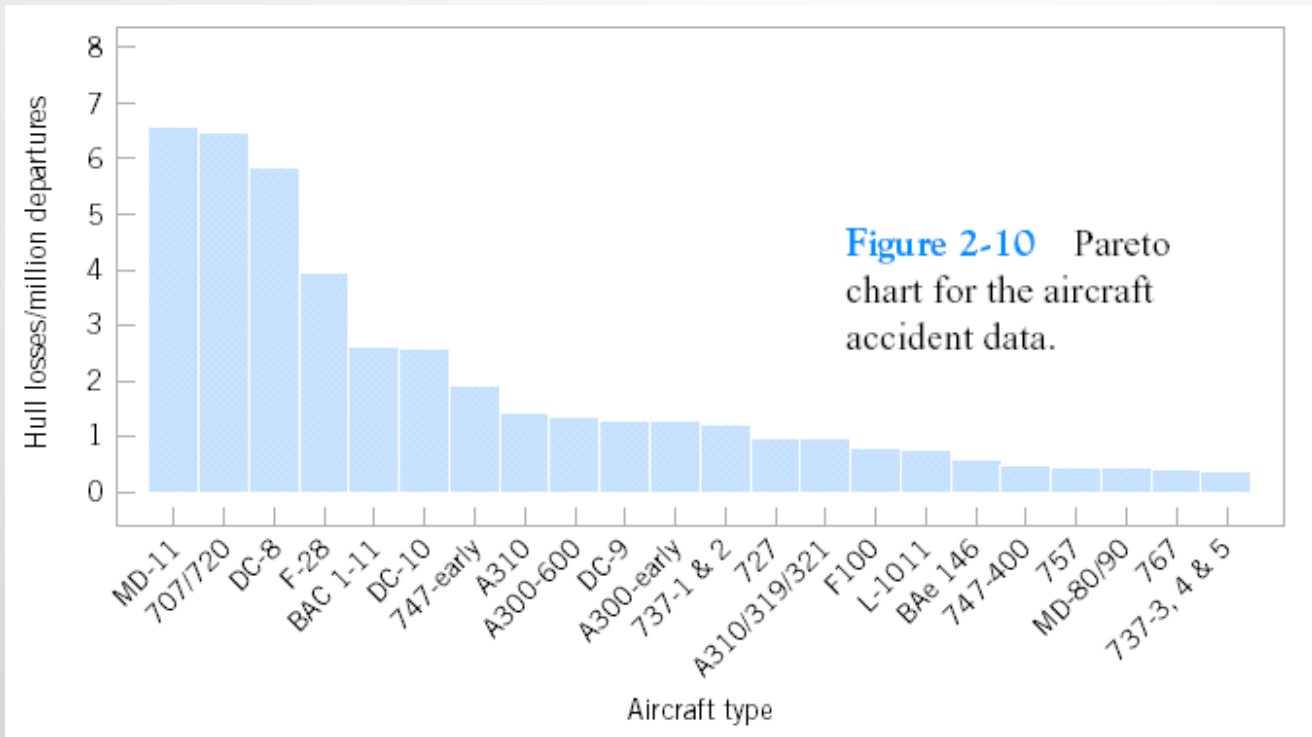
<i>Stem</i>	<i>Leaf</i>
.7	5 8
.8	2 6
.9	2 7
1.0	0 1 1 3 5 5 7 7 8 8
1.1	4 6
1.2	8
1.3	8 8
1.4	2

Data display: Histogram



- More compact

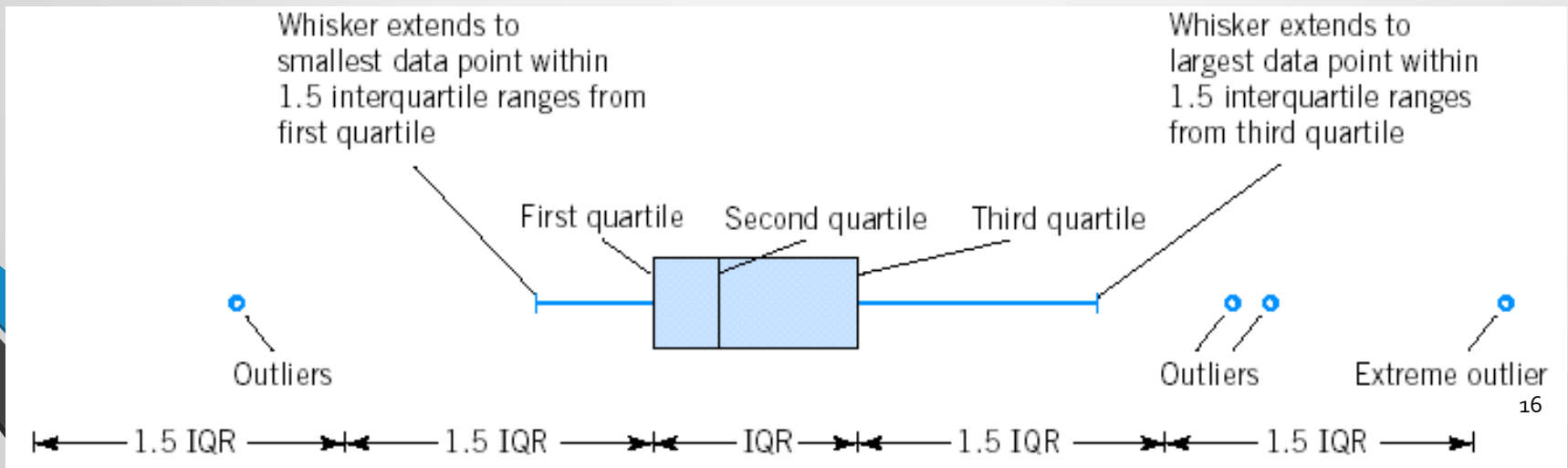
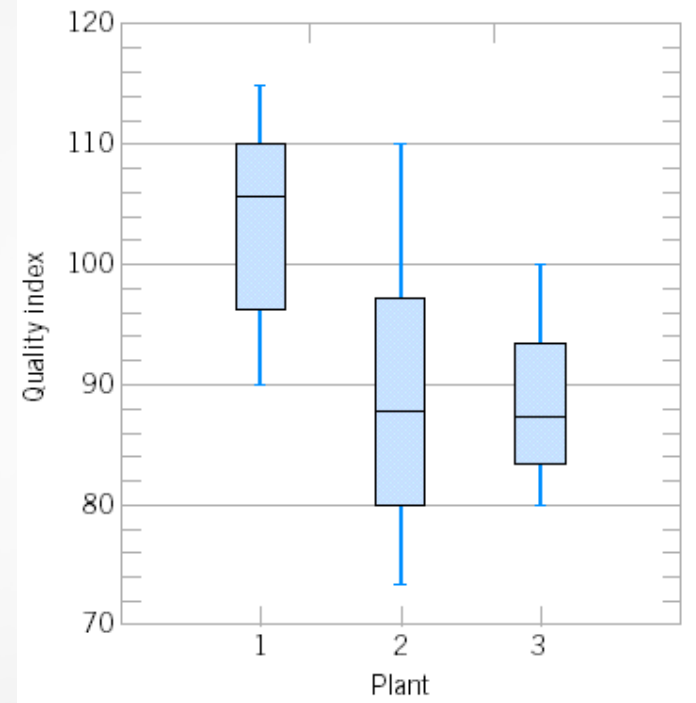
Data display: Pareto chart

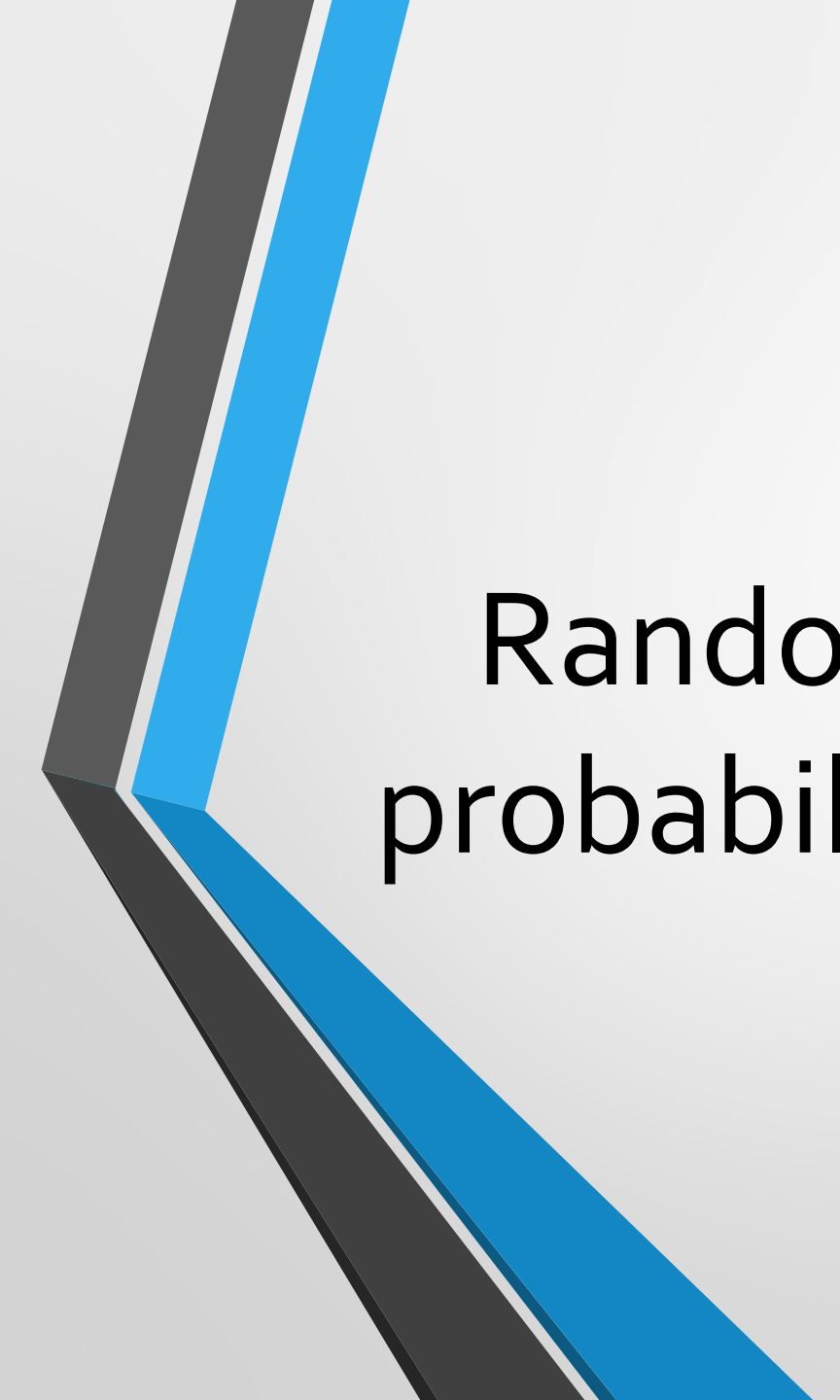


- Widely used in quality and process improvement studies

Data display: Box plot

- Main features
 - Q_1 , Median (Q_2), Q_3
 - IQR: interquartile range
 - Whisker
 - Outlier





Random variables and probability distributions

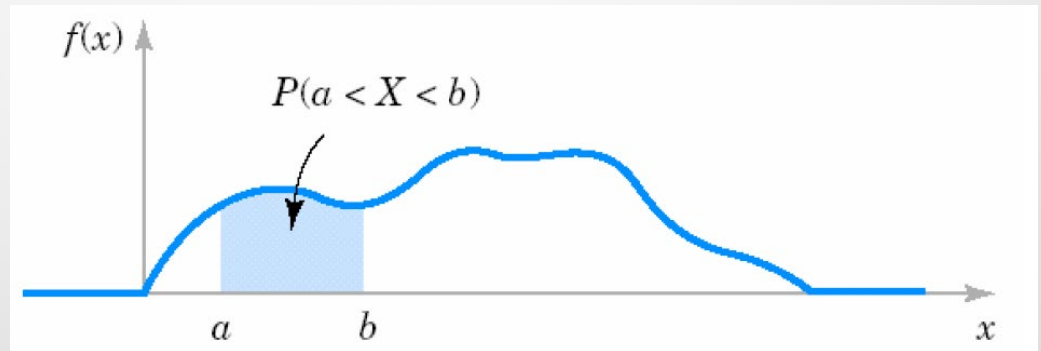
Random variable

- A **random variable** is a numerical variable whose measured value can change from one replicate of the experiment to another.
 - Example: body temperature, electric current, number of transmitted bits received in error
- **Probability**: the likelihood that particular values occur

Probability

- Probability density function (PDF, $f(x)$)
 - Properties of the PDF*

$$P(a < X < b) = \frac{\text{event}}{\text{event}}$$



- Approximated by histogram.

Probability

- Cumulative distribution function (CDF, $F(x)$)
 - $F(x) = P(X \leq x)$
 - $P(a < X < b) = ?$

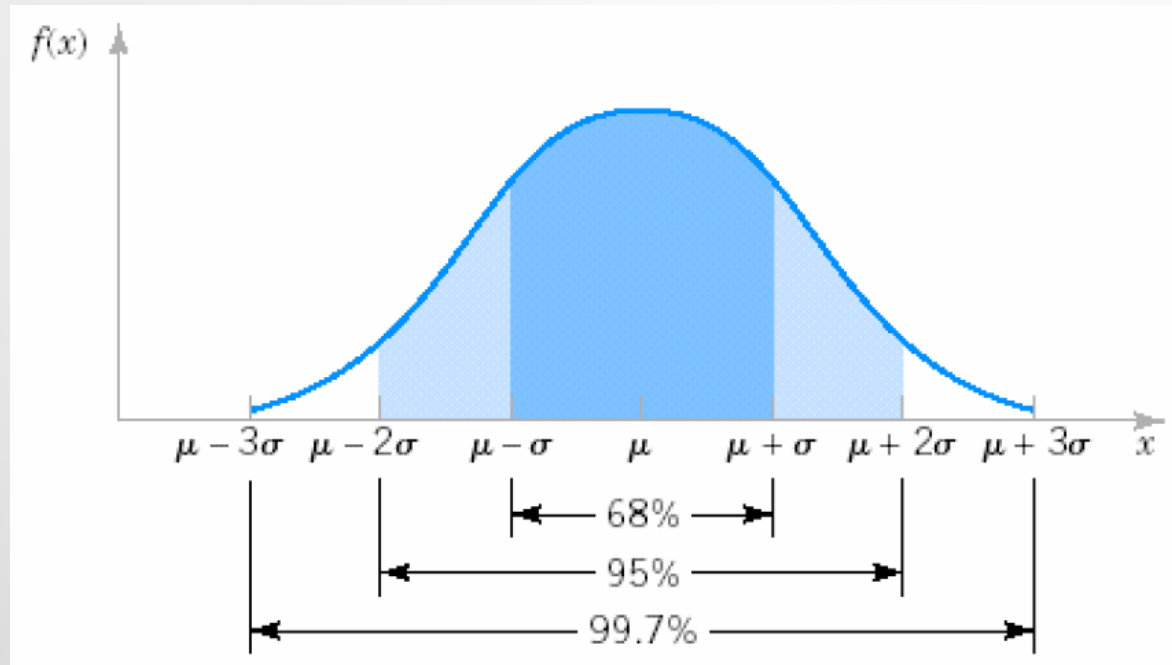
Continuous random variables

- Suppose X is a continuous random variable...
 - Mean (expected value, $E(x)$)
 - Variance ($V(x)$)

Normal distribution

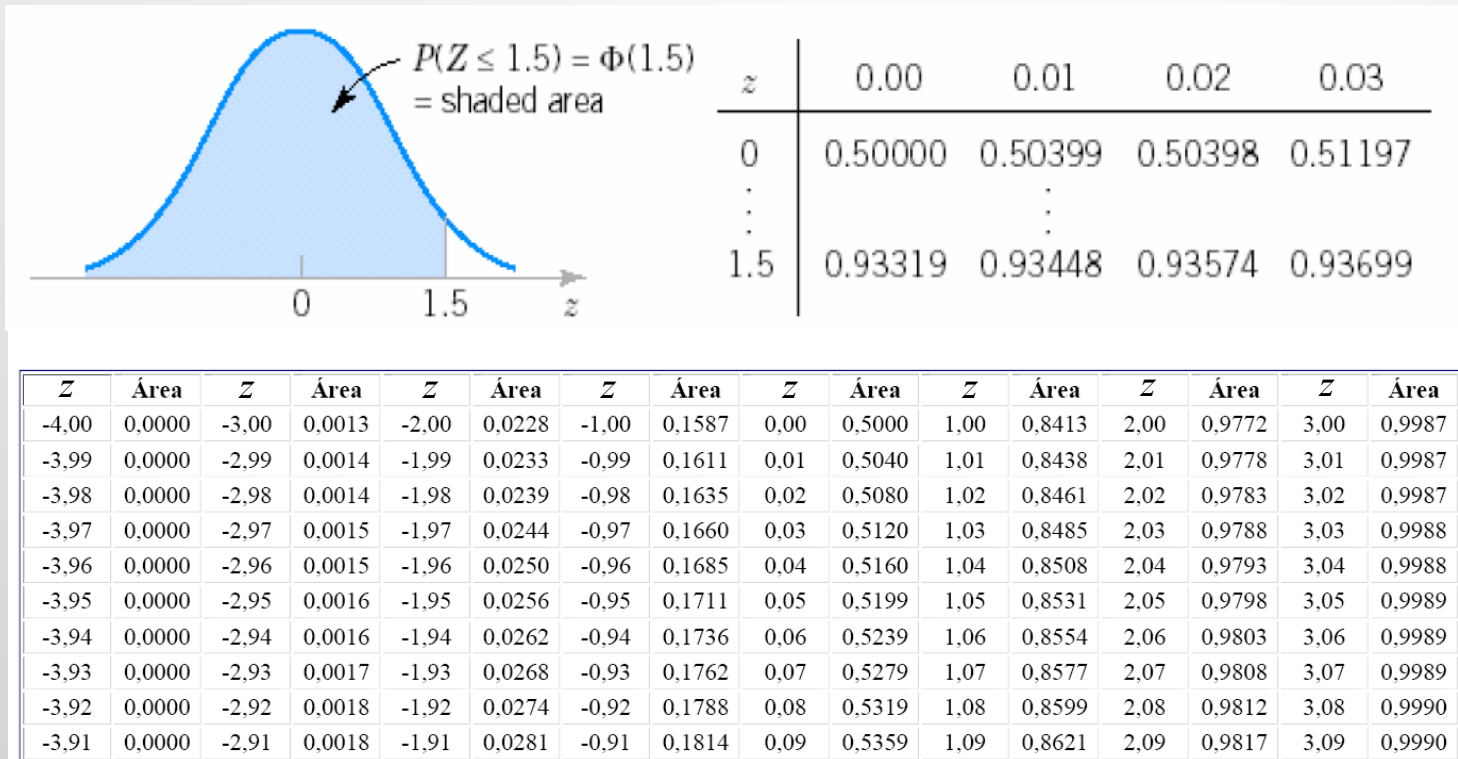
- Although distributions can have different shapes, the most widely used model for a random variable (e.g. blood pressure, height) is **normal distribution**.
 - Also referred to as a Gaussian distribution.
 - $f(x) = ?$

Normal distribution



- Standard normal distribution (Z)*
 - A normal random distribution with $\mu = 0$ and $\sigma^2 = 1$

Standard normal distribution



CDF of a standard normal distribution

Find $P(-1.92 < Z < 1.05) = ?$

Example

- The diameter of a shaft in an optical storage drive is normally distributed with mean 0.2508 inch and standard deviation 0.0005 inch. The specifications on the shaft are 0.2500 ± 0.0013 inch. What proportion of the shaft conforms to the specification?

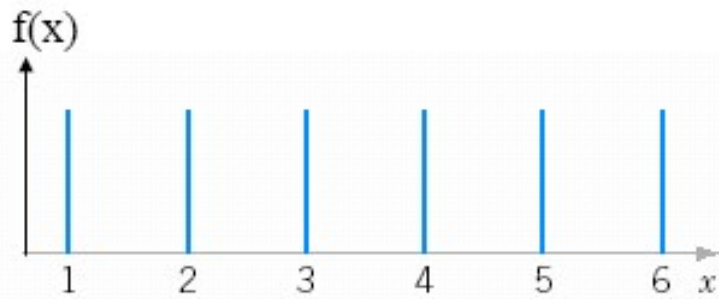
Normal distribution sampling theorem

- If a variable X is normally distributed with a mean μ and a standard deviation σ , the sampling distribution of the mean \bar{X} , based on random samples of size n , will also be normally distributed and have a mean μ , and a standard deviation $\sigma_{\bar{X}}$.

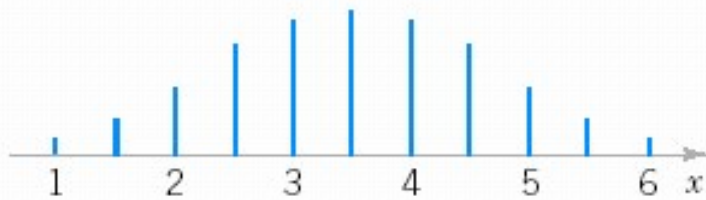
Central limit theorem

- If a variable X has a **distribution** $\sim(\mu, \sigma^2)$, the sampling distribution of the mean \bar{X} , based on random samples of size n , will have a mean equal to μ , and a standard deviation $\sigma_{\bar{X}}$, and the shape will tend to be a **normal distribution** when $n \rightarrow \infty$.

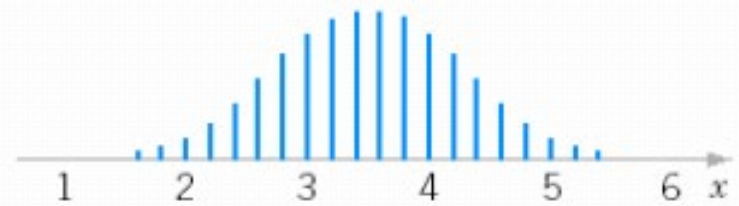
Central limit theorem



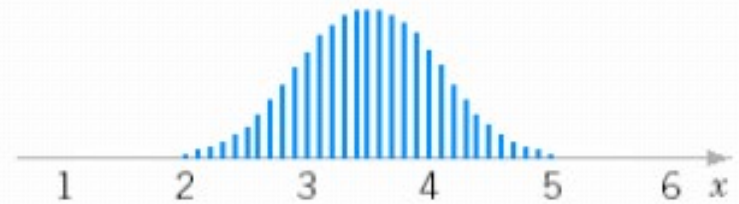
(a) One die



(b) Two dice



(d) Five dice

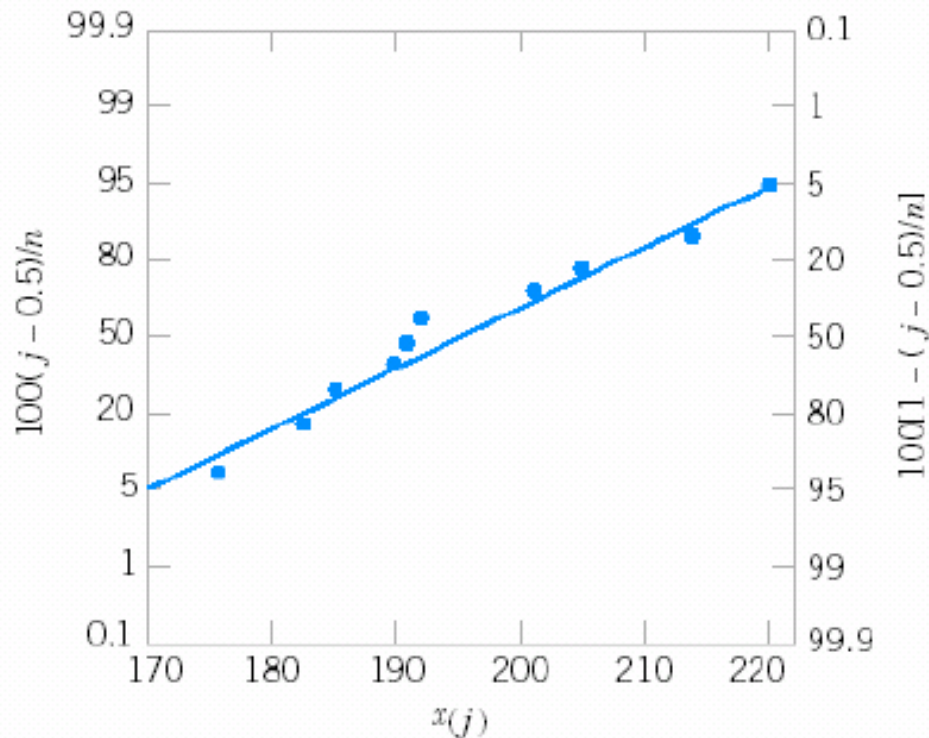


(e) Ten dice

Probability plotting

- How do we know if a normal distribution is a reasonable model for data?
- **Probability plotting** is a graphical method for determining whether sample data conform to a hypothesized distribution based on a subjective visual examination of the data.

Probability plotting



j	$x_{(j)}$	$(j - 0.5)/10$
1	176	0.05
2	183	0.15
3	185	0.25
4	190	0.35
5	191	0.45
6	192	0.55
7	201	0.65
8	205	0.75
9	214	0.85
10	220	0.95

Figure 3-23 Normal probability plot for the battery life.

• **Rule of thumb:** draw the straight line between the 25th and 75th percentile points