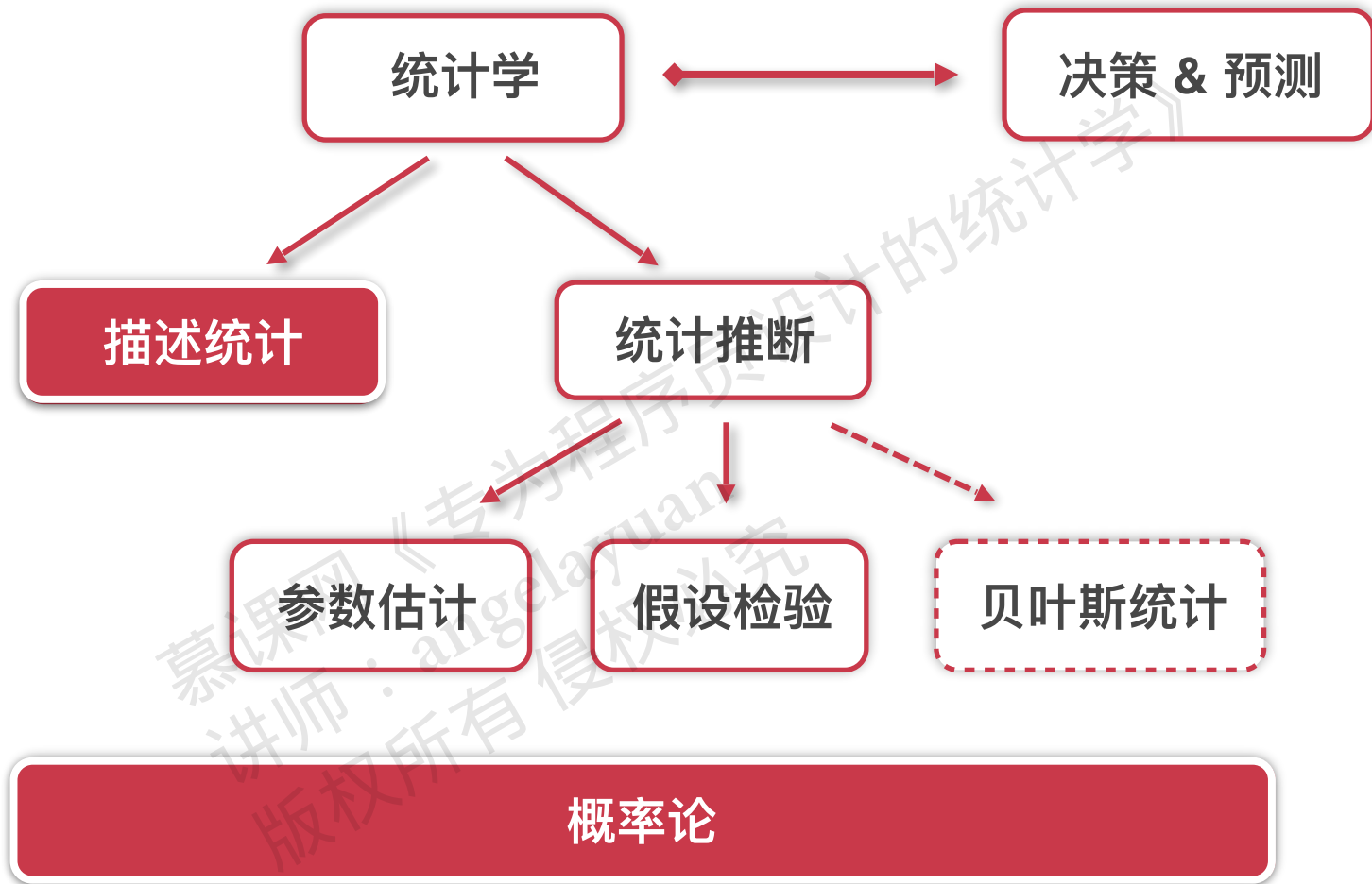


必须了解的概率论知识

慕课网《专为程序员设计的统计学》
讲师：angelayuan
版权所有 侵权必究



什么是概率论

慕课网《专为程序员设计的统计学》
讲师：angelayuan
版权所有 侵权必究

概率论

◆ 确定性现象

抛起的硬币必然下落

在标准大气压下, 水加热到100度必然沸腾

定义: 在一定条件下**必然**会发生的现象

慕课网 专为程序员设计的统计学》
讲师: angela zhan
版权所有 侵权必究

概率论

◆ 随机现象

抛一次硬币的结果{正面, 反面}

不确定,但肯定是正面或反面

抛一次骰子的结果{1, 2, 3, 4, 5, 6}

不确定,但肯定是六个点数之一

抛10000次硬币: {正, 反, 反, 正,}

正面:反面 $\approx 1:1$

抛10000次骰子: {1, 5, 3, 3, 1, 2, 2, 5, 5, 6, 4,}

1:2:3:4:5:6 $\approx 1:1:1:1:1:1$

概率论

◆ 随机现象

不确定,但肯定是正面或反面
不确定,但肯定是六个点数之一



在个别试验中
结果具有不确定性

正面:反面 $\approx 1:1$
 $1:2:3:4:5:6 \approx 1:1:1:1:1:1$



在大量重复试验中
结果呈现出固有规律性

概率论

在个别试验中
结果具有不确定性



随机试验(Experiment)

- 可以在相同条件下重复进行
- 可能的结果不止一个并且能够事先明确所有可能结果
- 进行试验前不能确定哪个结果会出现

在大量重复试验中
结果呈现出固有规律性



统计规律性

概率论

◆ 样本空间(Space)

定义: 随机试验E的所有可能结果组成的集合

试验: 扔一次硬币进行观察 → $S: \{\text{Head}, \text{Tail}\}$

试验: 扔二次硬币进行观察 → $S: \{HH, HT, TH, TT\}$

试验: 扔三次硬币进行观察 → $S: \{HHH, HHT, HTH, THH, TTH, THT, HTT, TTT\}$

概率论

◆ 随机事件

定义: 样本空间 S 的子集称为随机试验 E 的随机事件, 简称事件

$S: \{\text{Head}, \text{Tail}\} \rightarrow$ 子集: $\{H\}, \{T\}$ 基本事件

$S: \{HH, HT, TH, TT\} \rightarrow$ 子集: $\{HH\}, \{HT\}, \{TH\}, \{TT\},$
 $\{HH, HT\}, \{HT, TH\}, \dots$
 $\{HH, HT, TH\}, \{HT, TH, TT\}, \dots$

空集 \emptyset 也是 S 的子集, 它在每次试验中都不发生, 称为不可能事件

概率论

◆ 随机事件

定义: 样本空间 S 的子集称为随机试验 E 的随机事件, 简称事件

$S: \{HH, HT, TH, TT\} \rightarrow \{HH, HT, TH\} = \{\text{至少一次正面}\}$

事件发生: 在每次试验中, 当且仅当这一子集的一个样本点出现
样本空间是**必然事件**

随机现象

在个别**试验**中结果具有不确定性

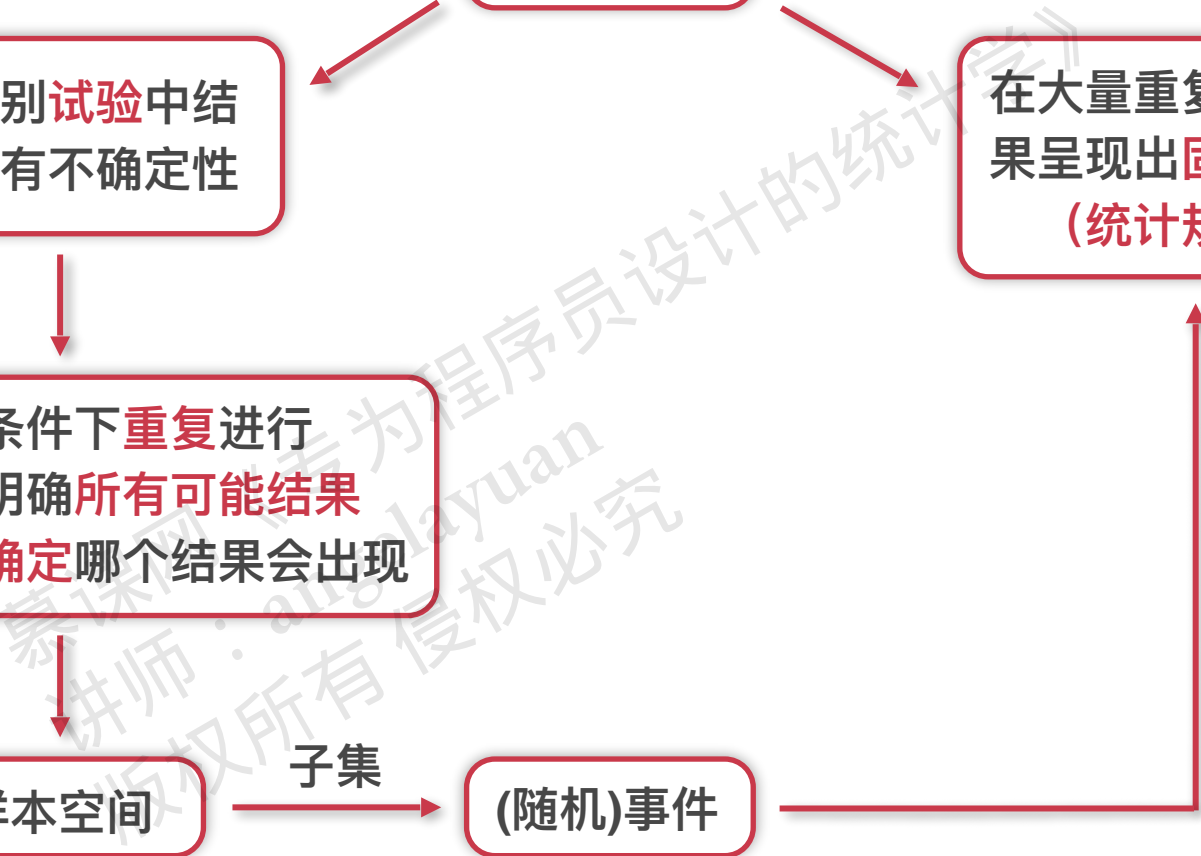
在大量重复试验中结果呈现出**固有规律性**
(统计规律性)

- 可在相同条件下**重复**进行
- 能够事先明确**所有可能结果**
- 试验前**不确定**哪个结果会出现

样本空间

子集

(随机)事件



慕课网《专为程序员设计的统计学》
讲师：angelayuan
版权所有 侵权必究

概率

个人经验/知识

明天的天气

S: {晴, 雨}

事件: {雨}

$P(\text{雨}) = 0.2$

$P(\text{雨}) = 0.7$

.....

美国总统选举

S: {希拉里赢, 希拉里输}

事件: {希拉里赢}

$P(\text{希拉里赢}) = 0.7$

$P(\text{希拉里赢}) = 0.4$

.....

- 猜测
- 个人经验/知识
- 主观
- 可信度低

频率

在相同条件下，进行了 n 次试验；在这 n 次试验中，
事件 A 发生的次数 n_A 称为事件 A 发生的**频数**
比值 n_A/n 称为事件 A 发生的**频率**，记为 $f_n(A)$

慕课网《为桂芳设计的统计学》
讲师：angelayuan
版权所有 侵权必究

频率

进行 $n = 100$ 次抛硬币

$A = \{H\}$ $n_A = 40$



$$f_n(A) = n_A/n = 40/100 = 0.4$$

性质 $0 \leq f_n(A) \leq 1$

$$f_n(S) = 1$$

若 A_1, A_2, \dots, A_k 是两两互不相容事件

$$\text{则 } f_n(A_1 \cup A_2 \cup \dots \cup A_k) = f_n(A_1) + f_n(A_2) + \dots + f_n(A_k)$$

概率

对于随机试验E的每一个事件A赋予一个实数，记为 $P(A)$ ，称为事件A的**概率**，如果集合函数 $P(\cdot)$ 满足下列条件

- 非负性: 对于每一个事件A, 有 $P(A) \geq 0$
- 规范性: 对于必然事件S, 有 $P(S) = 1$
- 可列可加性: 设 A_1, A_2, \dots 是两两互不相容事件, 有
$$P(A_1 \cup A_2 \cup \dots) = P(A_1) + P(A_2) + \dots$$

频率 vs 概率

为什么可以用频率近似概率？

$A = \{H\}$

抛 $n = 10$ 次硬币，计算频率

抛 $n = 60$ 次硬币，计算频率

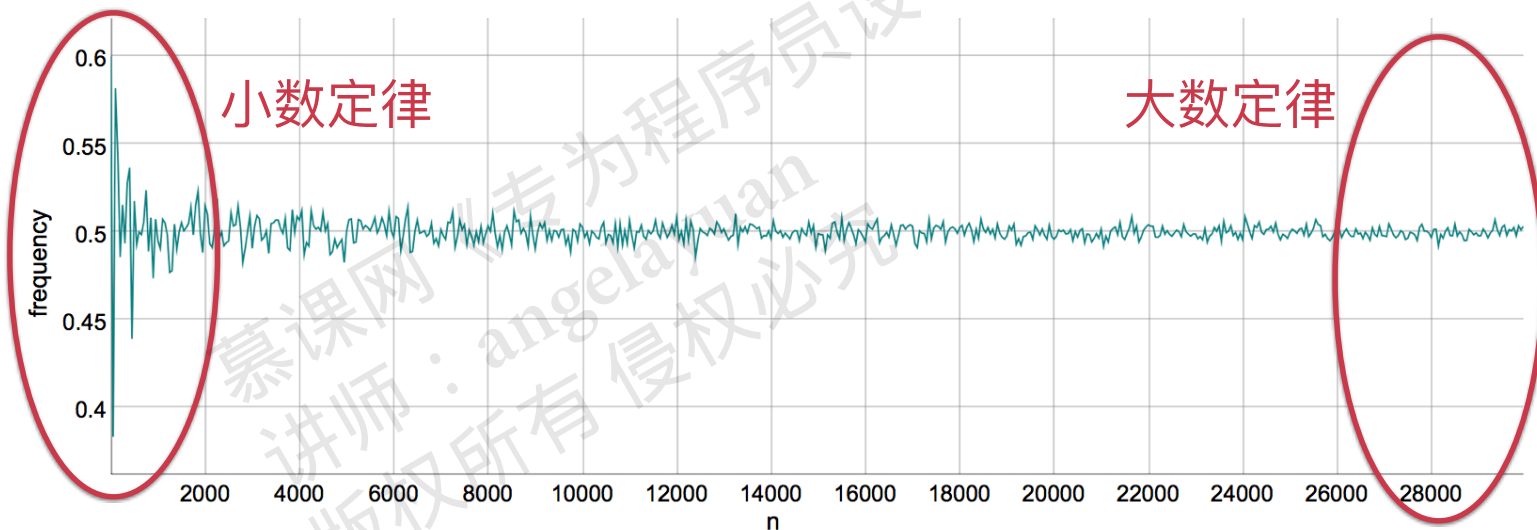
抛 $n = 110$ 次硬币，计算频率

.....

抛 $n = 10000$ 次硬币，计算频率

频率 vs 概率

为什么可以用频率近似概率？



$$n \rightarrow \infty, f_n(A) \rightarrow P(A)$$

编程理解小数和大数定律

慕课网《专为程序员设计的统计学》
讲师：angelayuan
版权所有 侵权必究

再谈变量

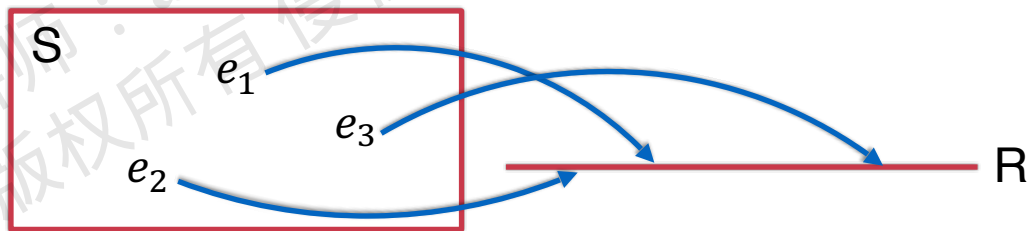
慕课网《专为程序员设计的统计学》
讲师：angelayuan
版权所有 侵权必究

变量

数据是由变量组成的; 一个变量至少包含2个不同取值

如何把数据/变量与概率论中的概念联系起来?

➡ 引入一个法则, 将样本空间 S 的每一个元素(即随机试验 E 的每一个结果)与实数对应起来



随机变量

设随机试验 E 的样本空间为 $S = \{e\}$, $X = X(e)$ 是定义在样本空间 S 上的实值单值函数, 称 $X = X(e)$ 为随机变量

试验 E : 抛硬币

样本空间 $S: \{H, T\}$

变量 $X: X(H) = 1$

$X(T) = 0$

名目尺度, 定性变量

试验 E : 观测教育程度

样本空间 $S: \{\text{小学}, \text{初中}, \text{高中}, \text{大学}\}$

变量 $X: X(\text{小学}) = 1$

$X(\text{初中}) = 2$

$X(\text{高中}) = 3$

$X(\text{大学}) = 4$

次序尺度

定性变量

随机变量

如果随机试验的结果本身就是一个数, 即 e 本身是一个数,
令 $X = X(e) = e$, X 就是一个随机变量

试验 E : 某大学学生的出生年份

样本空间 S : $\{2000, 2001, 2002, 2003\}$

变量 X : $X(2000) = 2000$

$X(2001) = 2001$

$X(2002) = 2002$

$X(2003) = 2003$

等距尺度
定量变量

随机变量

如果随机试验的结果本身就是一个数, 即 e 本身是一个数,
令 $X = X(e) = e$, X 就是一个随机变量

试验 E : 某大学学生的身高

样本空间 S : $\{1.55, 1.56, \dots, 1.90\}$

变量 X : $X(1.55) = 1.55$

$X(1.56) = 1.56$

.....

$X(1.90) = 1.90$

等比尺度
定量变量

以大写字母如 $X, Y, Z \dots$ 表示变量

以小写字母如 $x, y, z \dots$ 表示实数

随机变量

◆ 性质

随机变量的取值随试验的结果而定

试验的各个结果的出现有一定的概率

因而随机变量的取值有一定的概率

在试验之前不能预知它取什么值

随机变量的分类

离散型

取值有限个或可列无限多个

连续型

在一定区间内可以任意取值

慕课网《专为程序员设计的统计学》
讲师：angelayuan
版权所有 侵权必究

概率分布

离散型随机变量及其分布

慕课网《专为程序员设计的统计学》
讲师：angelayuan
版权所有 侵权必究

分布律

离散型随机变量: 可能取到的值是有限个或可列无限多个

随机变量X的分布律 (Probability Mass Function)

设X所有可能取的值为 x_k ($k = 1, 2, \dots$)

X取各个可能值的概率, 即事件 $\{X = x_k\}$ 的概率为

$$P\{X = x_k\} = p_k \quad (k = 1, 2, \dots)$$

分布律

$$P\{X = x_k\} = p_k \quad (k = 1, 2, \dots)$$

抛硬币 X所有可能的取值为 $x_1 = 0, x_2 = 1$

X取各个可能值的概率为 $P\{X = 0\} = 0.5; P\{X = 1\} = 0.5$



掷 X所有可能的取值为

$$x_1 = 1, x_2 = 2, x_3 = 3, x_4 = 4, x_5 = 5, x_6 = 6$$

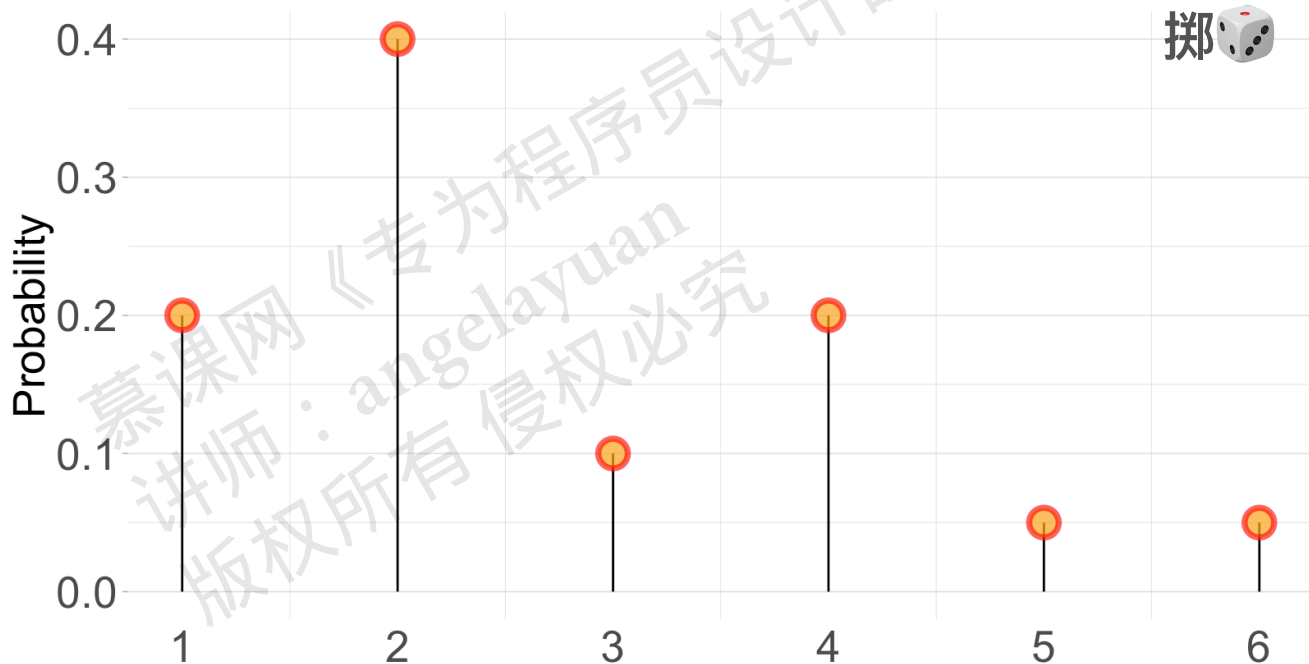
X取各个可能值的概率为

$$P\{X = 1\} = 0.2; P\{X = 2\} = 0.4; P\{X = 3\} = 0.1$$

$$P\{X = 4\} = 0.2; P\{X = 5\} = 0.05; P\{X = 6\} = 0.05$$

分布律

$$P\{X = x_k\} = p_k \quad (k = 1, 2, \dots)$$



两种重要的分布

- ◆ (0-1)分布/两点分布 (Bernoulli distribution)
- ◆ 伯努利试验，二项分布 (Binomial distribution)

慕课网《专为程序员设计的统计学》
讲师：angelayuan
版权所有 侵权必究

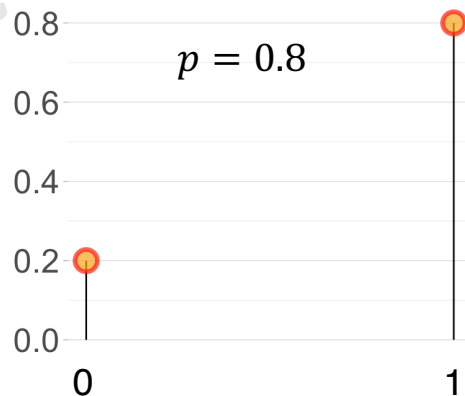
(0-1)分布/两点分布 (Bernoulli distribution)

设随机变量 X 只可能取0与1两个值, 它的分布律是

$$P\{X = k\} = p^k(1 - p)^{1-k}, (k = 0, 1; 0 < p < 1)$$

则称 X 服从以 p 为参数的(0-1)分布或两点分布

X	0	1
$P\{X = k\}$	$1 - p$	p



(0-1)分布/两点分布 (Bernoulli distribution)

X	0	1
$P\{X = k\}$	$1 - p$	p

对于随机试验E, 如果其样本空间S只包含两个元素, 总能够在S上定义一个服从(0-1)分布的随机变量来描述这个随机试验的结果

性别: $X=0$ 当性别为女; $P\{\text{女}\} = 1 - p$

$X=1$ 当性别为男; $P\{\text{男}\} = p$

伯努利试验，二项分布 (Binomial distribution)

设试验E只有两个可能的结果: A, \bar{A} , 则称E为伯努利试验

抛硬币 性别

设 $P(A) = p$ ($0 < p < 1$), 此时 $P(\bar{A}) = 1 - p$

将E独立重复地进行n次, 则称这一串重复的独立试验为n重伯努利试验

抛n次硬币 记录n个人的性别

伯努利试验，二项分布 (Binomial distribution)

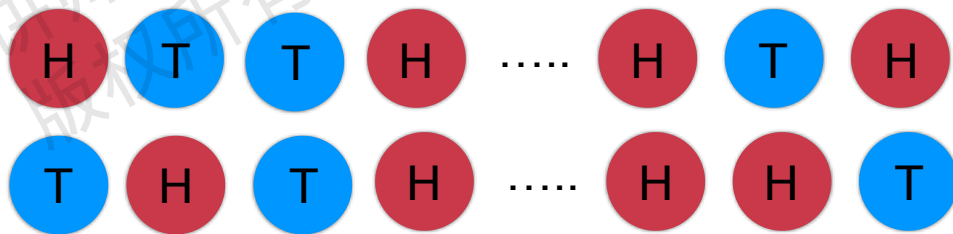
以 X 表示 n 重伯努利试验中事件 A 发生的次数, X 是一个随机变量,
 X 的所有可能的取值为 $0, 1, 2, \dots, n$, 其分布律为

$$P\{X = k\} = \binom{n}{k} p^k (1-p)^{n-k}, (k = 0, 1, 2, \dots, n; 0 < p < 1)$$



前 k 次正面朝上

后 $(n-k)$ 次反面朝上



伯努利试验，二项分布 (Binomial distribution)

以 X 表示 n 重伯努利试验中事件 A 发生的次数, X 是一个随机变量, X 的所有可能的取值为 $0, 1, 2, \dots, n$, 其分布律为

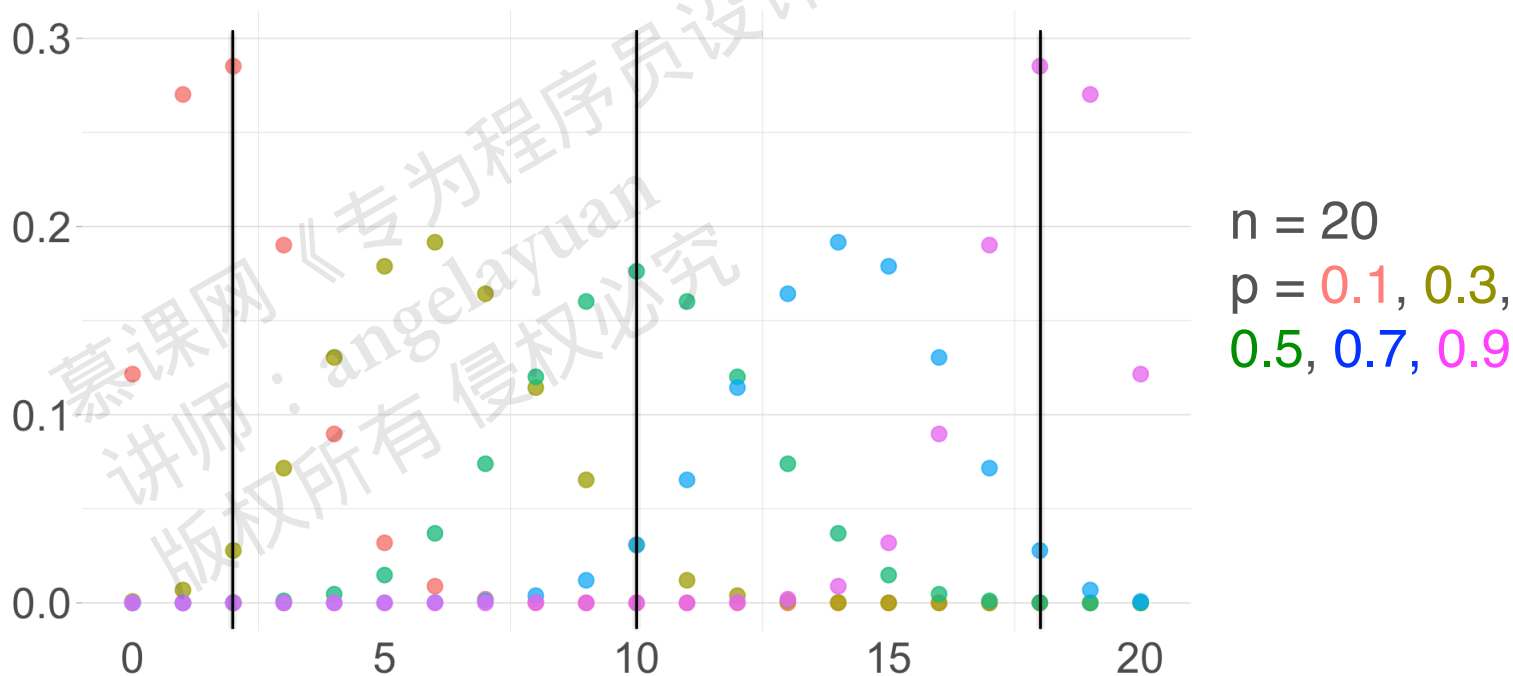
$$P\{X = k\} = \binom{n}{k} p^k (1-p)^{n-k}, (k = 0, 1, 2, \dots, n; 0 < p < 1)$$

其中 $\binom{n}{k} = \frac{n!}{k!(n-k)!}$ $n! = n \times (n-1) \times \dots \times 2 \times 1$

则称 X 服从参数为 n, p 的二项分布, 记为 $X \sim b(n, p)$

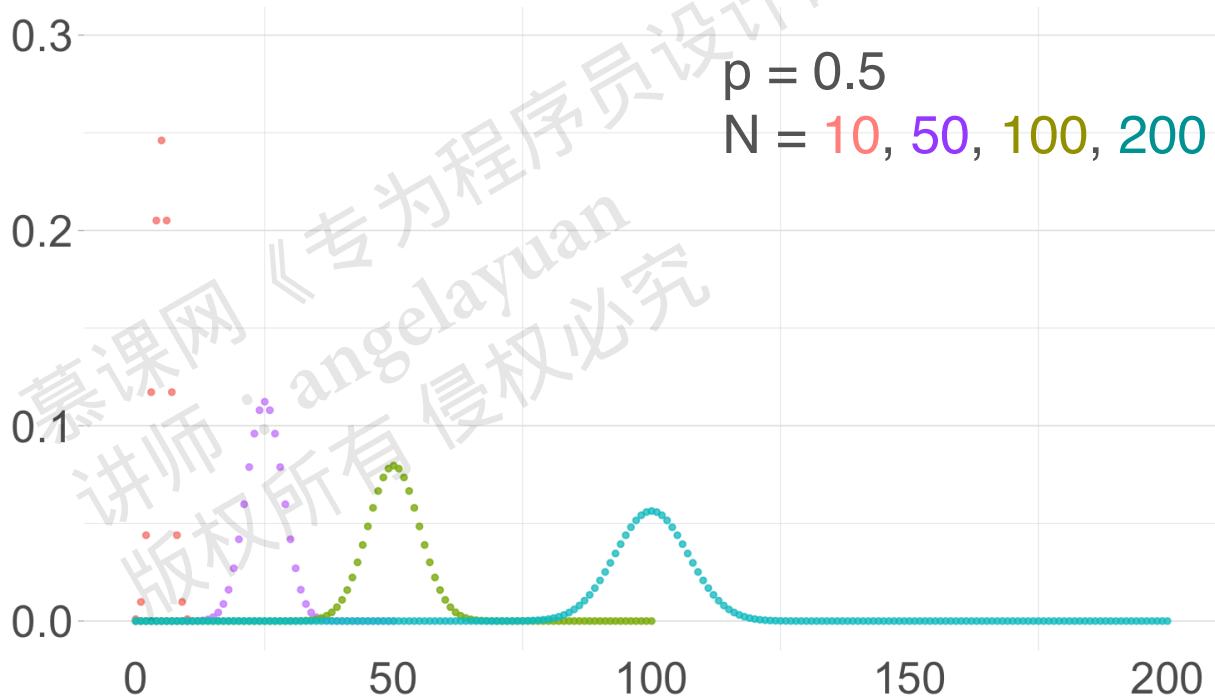
伯努利试验，二项分布 (Binomial distribution)

$$P\{X = k\} = \binom{n}{k} p^k (1 - p)^{n-k}, (k = 0, 1, 2, \dots, n; 0 < p < 1)$$



伯努利试验，二项分布 (Binomial distribution)

$$P\{X = k\} = \binom{n}{k} p^k (1 - p)^{n-k}, (k = 0, 1, 2, \dots, n; 0 < p < 1)$$



伯努利试验，二项分布 (Binomial distribution)

$$P\{X = k\} = \binom{n}{k} p^k (1 - p)^{n-k}, (k = 0, 1, 2, \dots, n; 0 < p < 1)$$

每次抛硬币，正面朝上的概率是 p

抛 n 次硬币，其中有 k 次正面朝上的概率是多少？

某地区，一个新生儿是女孩的概率是 p

现有 n 个新生儿，其中有 k 个是女孩的概率是多少？

连续型随机变量及其分布

慕课网《专为程序员设计的统计学》
讲师：angelayuan
版权所有 侵权必究

分布函数

设 X 是一个随机变量， x 是任意实数，函数

$$F(x) = P\{X \leq x\}, -\infty < x < \infty$$

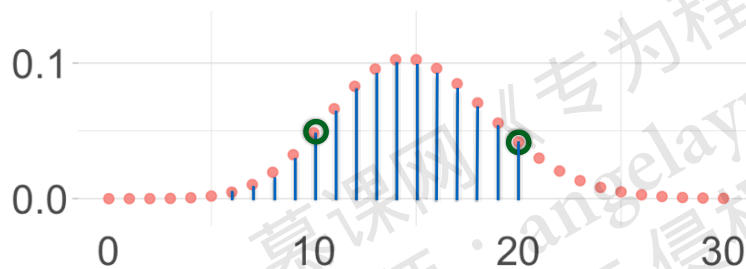
称为 X 的分布函数(Cumulative Distribution Function)

分布函数适用于离散型随机变量和连续型随机变量

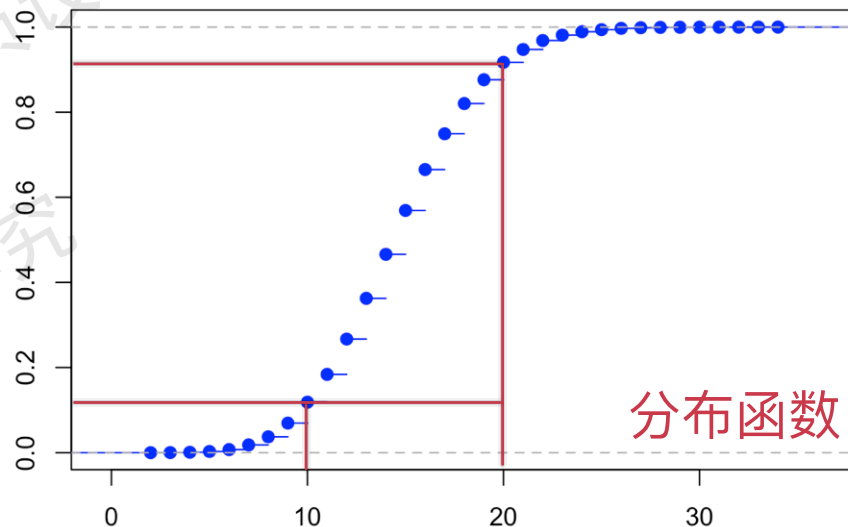
分布函数

$$P\{X = x_k\} = p_k \quad (k = 1, 2, \dots)$$

$$F(x) = P\{X \leq x\}, -\infty < x < \infty$$



分布律



分布函数

若已知 X 的分布函数，就可知 X 落在任一区间 $(x_1, x_2]$ 的概率

分布函数完整地描述了随机变量的统计规律性

慕课网《为程序员设计的统计学》
讲师：angelayuda
版权所有 侵权必究

概率密度

◆ 连续型随机变量

如果对于随机变量 X 的分布函数 $F(x)$,存在非负函数 $f(x)$

使对于任意实数 x 有
$$F(x) = \int_{-\infty}^x f(t) d(t)$$

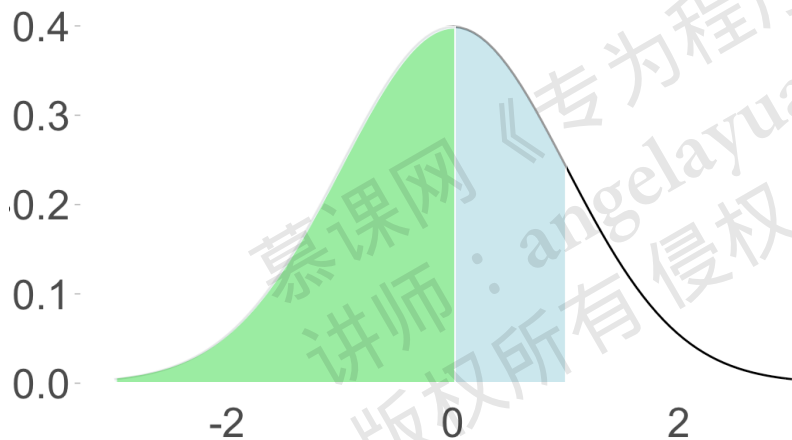
则称 X 为连续型随机变量

函数 $f(x)$ 称为 X 的**概率密度函数**(Probability Density Function)

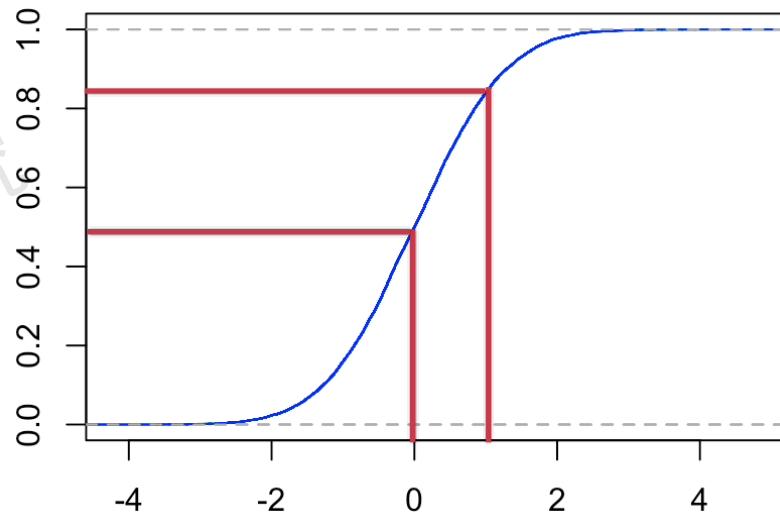
概率密度

$$F(x) = \int_{-\infty}^x f(t) d(t)$$

概率密度 $f(x)$



分布函数 $F(x)$



概率密度

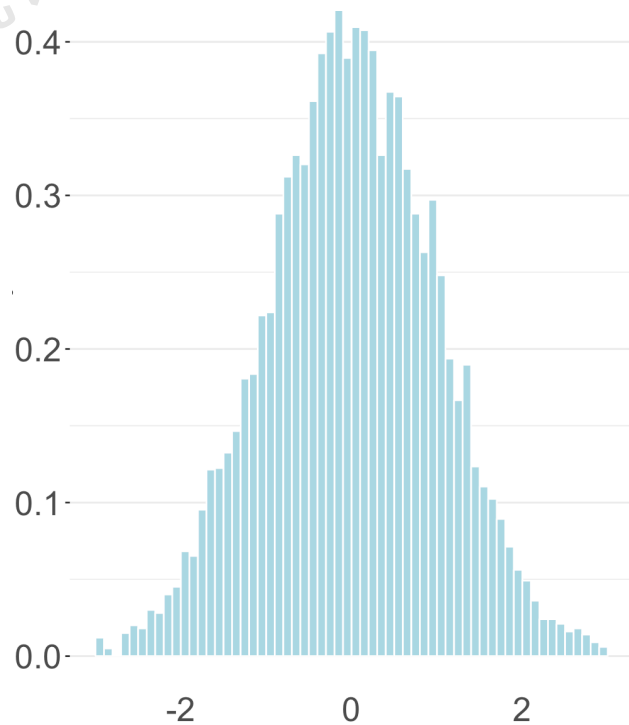
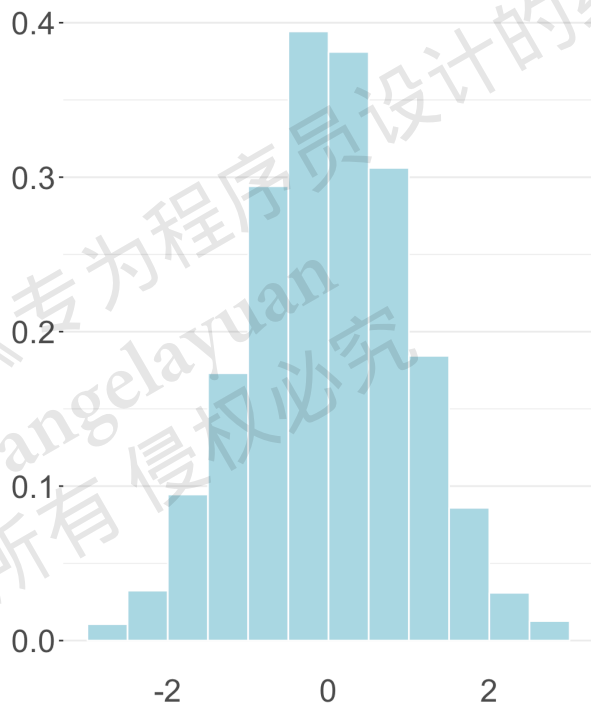
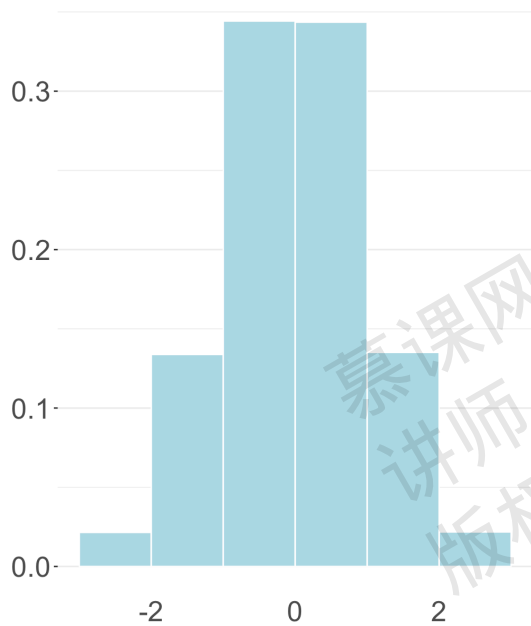
对于连续型随机变量 X 而言, 它取任一指定实数值 a 的概率均为0

$$P\{X = a\} = 0$$

若 A 是不可能事件, 则有 $P(A) = 0$

若 $P(A) = 0$, 并不意味着 A 是不可能事件

概率密度



正态分布

慕课网《专为程序员设计的统计学》
讲师：angelayuan
版权所有 侵权必究

正态分布 (Normal distribution)

若连续型随机变量 X 的概率密度为

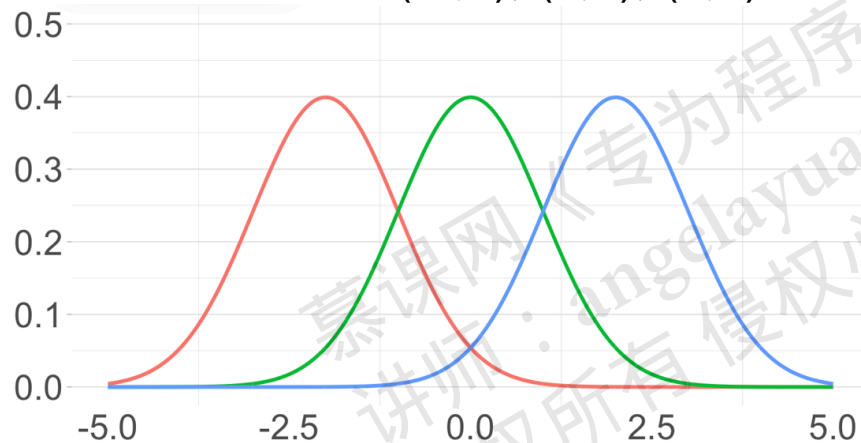
$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, (-\infty < x < \infty)$$

其中 $\mu, \sigma (\sigma > 0)$ 为常数, 则称 X 服从参数为 μ, σ

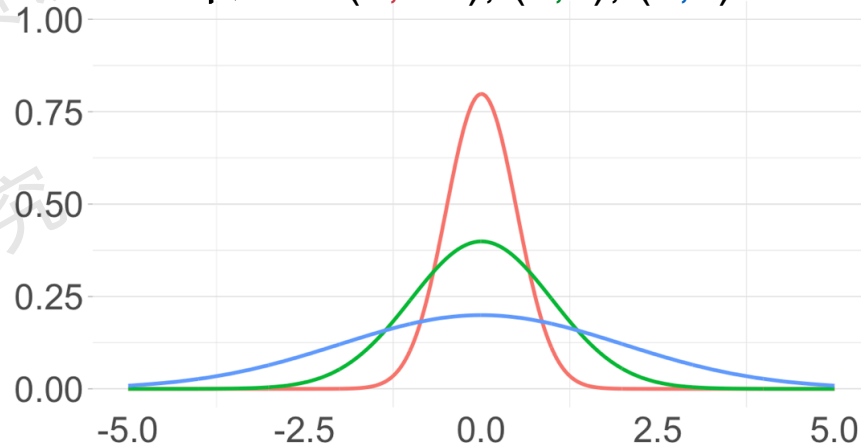
的正态分布或高斯(Gauss)分布, 记为 $X \sim N(\mu, \sigma^2)$

正态分布 (Normal distribution)

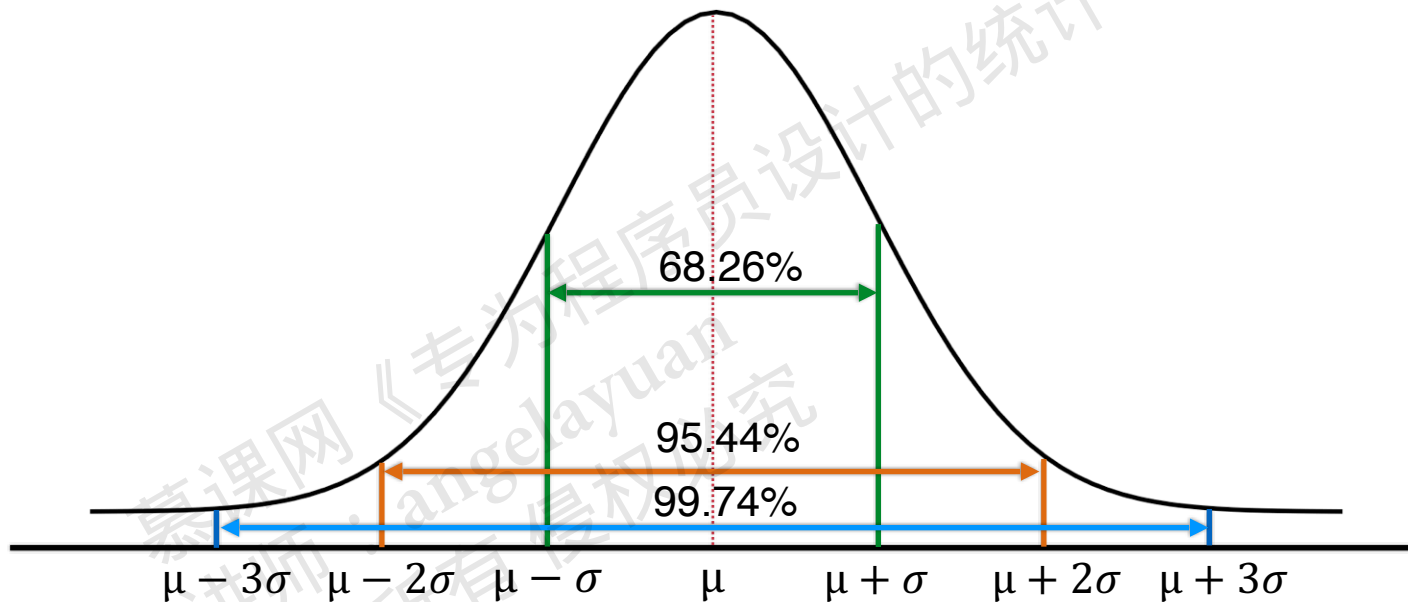
$\mu, \sigma = (-2, 1), (0, 1), (2, 1)$



$\mu, \sigma = (0, 0.5), (0, 1), (0, 2)$



正态分布 (Normal distribution)

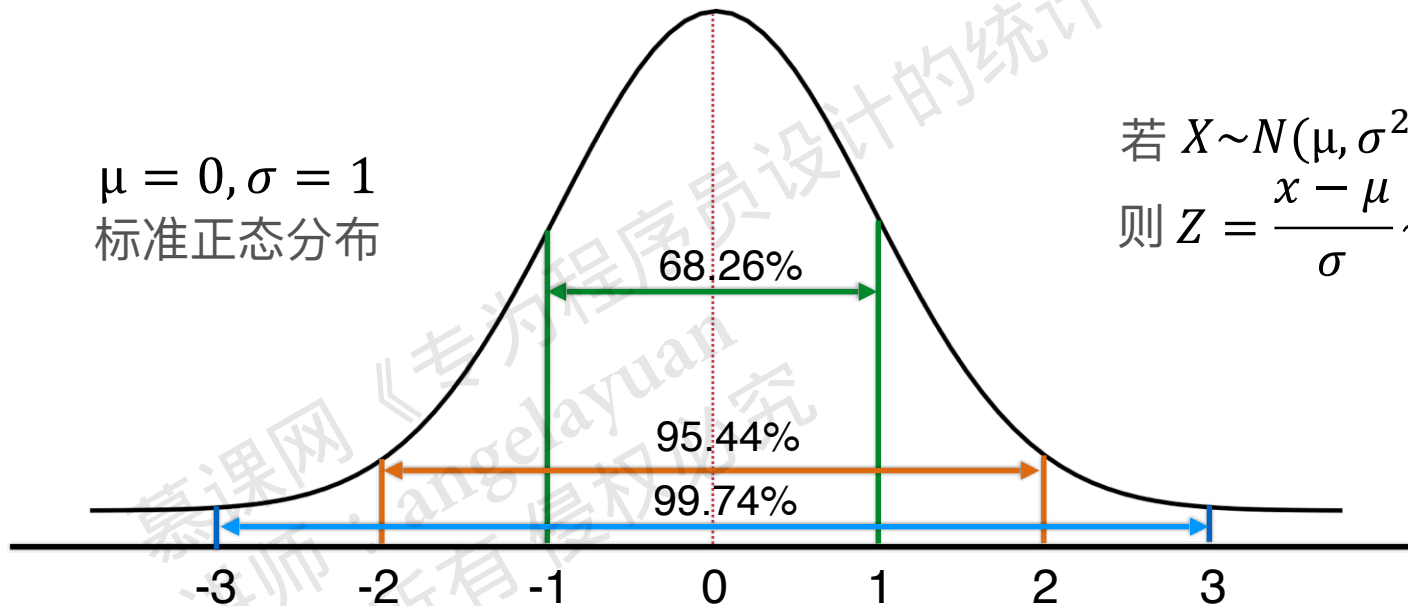


3 σ 法则: 正态变量的取值范围是正负无穷, 但它的值几乎肯定落在 $(\mu - 3\sigma, \mu + 3\sigma)$

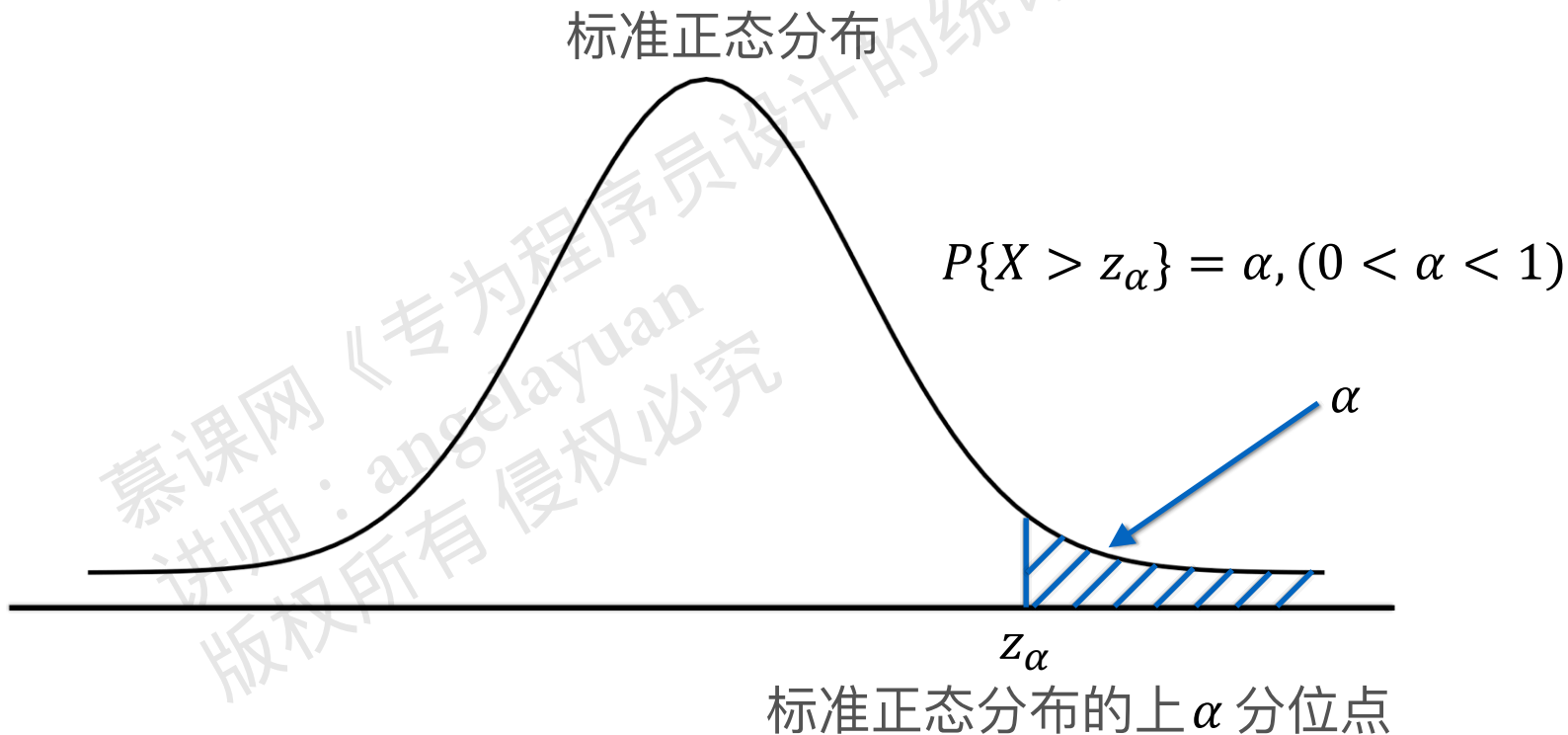
正态分布 (Normal distribution)

$\mu = 0, \sigma = 1$
标准正态分布

若 $X \sim N(\mu, \sigma^2)$
则 $Z = \frac{x - \mu}{\sigma} \sim N(0, 1)$



正态分布 (Normal distribution)



慕课网《专为程序员设计的统计学》
讲师：angelayuan
版权所有 侵权必究

本章小结

必须了解的概率论知识

```
graph TD; A[必须了解的概率论知识] --> B[随机现象<br/>样本空间<br/>随机事件]; A --> C[概率分布<br/>• 两点分布、二项分布<br/>• 正态分布]; A --> D[把数据/变量和<br/>概率论中的概念<br/>建立联系]; A --> E[概率<br/>• 频率vs概率<br/>• 小数/大数定律];
```

- 随机现象
- 样本空间
- 随机事件

- 概率分布
 - 两点分布、二项分布
 - 正态分布

- 把数据/变量和概率论中的概念建立联系

- 概率
 - 频率vs概率
 - 小数/大数定律