

高级计量经济学

理论经济学博士课程

Lecture 1-2: Review of Statistics

黄嘉平

工学博士 经济学博士

深圳大学中国经济特区研究中心 讲师

Office

粤海校区汇文楼1510

Email

huangjp@szu.edu.cn

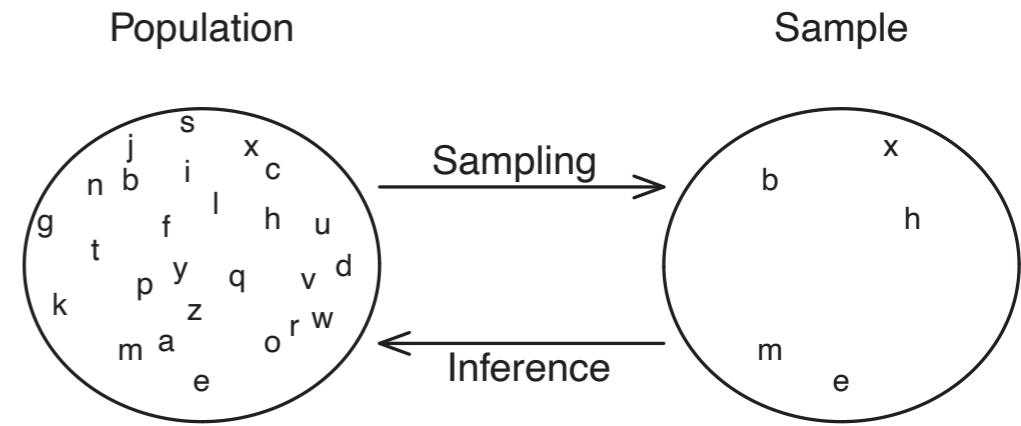
Website

<https://huangjp.com>

统计学中的基础概念

随机抽样

Random sampling



- 随机变量的向量称为**随机向量 (random vector)**，即 $Y = (Y_1, Y_2, \dots, Y_k)$ ，其中 Y_j 为随机变量

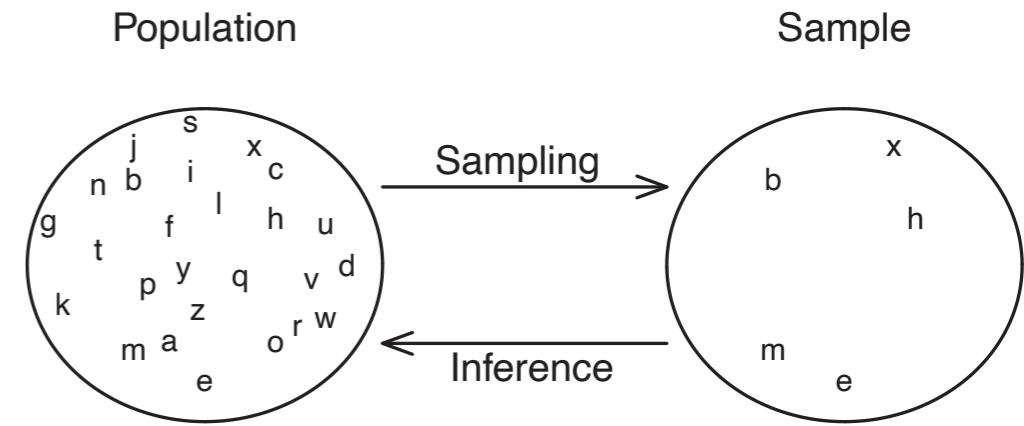
随机向量的集合 $\{X_1, X_2, \dots, X_n\}$ 为独立同分布 (independent and identically distributed/i.i.d.) 是指这些随机向量相互独立且服从同一个分布 F

独立同分布的随机向量集合 $\{X_1, X_2, \dots, X_n\}$ 称为**随机样本 (random sample)**

- F 称为**总体分布 (population distribution)** 或**总体 (population)**
- 从总体中获得随机样本的过程称为**随机抽样 (random sampling)**。随机抽样有两种解释：
 - 全体观测对象切实存在。当总体中包含 N 个个体 (N 可以是有限或无限) 时，随机抽样是从其中随机选择 n 个个体的操作。例如问卷调查或普查
“随机”代表每个个体被选中的概率相等
 - 全体观测对象不存在。此时通过特定的数据生成过程 (data generating process/DGP) 产生 n 个结果的操作。例如实验或观测性研究

样本与推断

Sample and inference



样本 (sample) 是随机抽样的结果

- 我们通常用到的数据集 (dataset) 可以看做是某个总体的样本，数据集中的每一条数据都是一个观测值 (observation)
- 当样本中有 n 个观测值时， n 被称为样本量 (sample size)

通过样本推测总体或数据生成过程的特征的行为称为推断 (inference)

- 采用哪种推断方法取决于数据是如何获得的，即抽样的方法。我们主要关注随机抽样
- 随机抽样以外还有分层抽样 (stratified sampling) 、整群抽样 (clustered sampling) 、面板数据 (panel data) 、时间序列数据 (time series data) 、空间数据 (spatial data) 等

真实数据的例子

Current Population Survey

- Current Population Survey (CPS) 是美国政府主导的关于就业的普查
- CPS 的样本量为60000个家庭
- CPS 采用两阶段分层抽样法
 1. 将全国分成 852 个地区，并按照所在地和人口特征等对这些地区进行分层，在每层中随机抽选一个地区
 2. 以地区为单位，在地区内部随机抽选对象家庭
- 右侧的表中包含2009年3月份普查中，已婚黑人女性且工作经验超过12年的子样本。Wage 为时薪（年薪/工作时长），Education 为受教育年限

Table 6.1: Observations From CPS Data Set

Observation	Wage	Education
1	37.93	18
2	40.87	18
3	14.18	13
4	16.83	16
5	33.17	16
6	29.81	18
7	54.62	16
8	43.08	18
9	14.42	12
10	14.90	16
11	21.63	18
12	11.09	16
13	10.00	13
14	31.73	14
15	11.06	12
16	18.75	16
17	27.35	14
18	24.04	16
19	36.06	18
20	23.08	16

统计量、参数、估计量

Statistics, parameters, and estimators

总体 F 的任意函数 θ 称为参数 (parameter)

- 例如，总体均值 $\mu = E[X]$ 是总体分布 F 的一阶矩

样本 $\{X_1, \dots, X_n\}$ 的函数称为统计量 (statistic)

- 例如，样本均值 $\bar{X} = \frac{1}{n} \sum_i^n X_i$ 是一个统计量

如果我们想用某一统计量来猜测参数 θ ，则称这个统计量是 θ 的估计量 (estimator)，记作 $\hat{\theta}$

- 例如，样本均值 \bar{X} 是总体均值 μ 的一个估计量，可以写成 $\hat{\mu} = \bar{X}$
- 当样本已经确定时，估计量 $\hat{\theta}$ 的取值被称为估计值 (estimate)

抽样分布

Sampling distribution

- 统计量是随机变量的函数，因此也是随机变量，并服从其自身的概率分布

统计量的概率分布称为为抽样分布 (**sampling distribution**)

- 估计量 $\hat{\theta}$ 的抽样分布可以帮助我们了解更多关于参数 θ 的信息
- 以样本均值 \bar{X} 为例，我们可以从以下角度了解 \bar{X} 的抽样分布
 - 偏差和方差
 - 当总体分布是正态分布时， \bar{X} 的分布
 - 大样本下 ($n \rightarrow \infty$) 的分布
 - 分布的渐进展开
 - Bootstrap 近似

$$\text{总体均值: } E[X_i] = \mu$$

$$\text{总体方差: } \text{Var}[X_i] = E[(X_i - \mu)^2] = \sigma^2$$

$$\sigma = \sqrt{\sigma^2} \text{ 称为总体标准差 (standard deviation)}$$

估计偏差

Estimation bias

估计量 $\hat{\theta}$ 的偏差 (bias) 定义为 $\text{bias}[\hat{\theta}] = E[\hat{\theta}] - \theta$

- 如果偏差为 0, 则称该估计量为非偏 (unbiased)
- 在不同的总体分布 F 下, 同一估计量可能为非偏也可能有偏
- $E[\bar{X}] = \frac{1}{n} \sum_{i=1}^n E[X_i] = \frac{1}{n} \sum_{i=1}^n \mu = \mu$, 因此当 $\mu < \infty$ 时, \bar{X} 为 μ 的非偏估计量
- μ 的其他非偏估计量包括 $X_1, \frac{\sum_{i=1}^n w_i X_i}{\sum_{i=1}^n w_i}$ 等
- $\hat{\theta}$ 是 θ 的非偏估计量 $\Rightarrow \hat{\beta} = a\hat{\theta} + b$ 是 $\beta = a\theta + b$ 的非偏估计量
- 如果 h 是非线性函数, 则 $\hat{\beta} = h(\hat{\theta})$ 不一定是 $\beta = h(\theta)$ 的非偏估计量

估计方差

Estimation variance

估计量 $\hat{\theta}$ 的方差 $\text{Var}[\hat{\theta}] = E[(\hat{\theta} - E[\hat{\theta}])^2]$ 称为抽样方差 (sampling variance)

- \bar{X} 的抽样方差：
 - X_1, \dots, X_n 为独立同分布 $\Rightarrow X_1, \dots, X_n$ 两两不相关
 - X_1, \dots, X_n 两两不相关 $\Rightarrow \text{Var}\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n \text{Var}[X_i]$ (尝试证明这个命题)
 - $\text{Var}[\bar{X}] = \text{Var}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n^2} \sum_{i=1}^n \text{Var}[X_i] = \frac{\sigma^2}{n}$ ($\text{Var}[X_i] = \sigma^2 < \infty$)
- 如果 $\hat{\beta} = a\hat{\theta} + b$, 则 $\text{Var}[\hat{\beta}] = a^2 \text{Var}[\hat{\theta}]$
- 一般情况下, 我们无法获得非线性变换 $h(\bar{X})$ 的方差的准确表达

均方误差

Mean squared error (MSE)

- 均方误差或均方误是衡量估计准确度的常用标准

估计量 $\hat{\theta}$ 的均方误差 (mean squared error) 定义为

$$\text{MSE}[\hat{\theta}] = E[(\hat{\theta} - \theta)^2]$$

- 将定义展开可得

$$\begin{aligned}\text{MSE}[\hat{\theta}] &= E[(\hat{\theta} - \theta)^2] \\ &= E[(\hat{\theta} - E[\hat{\theta}] + E[\hat{\theta}] - \theta)^2] \\ &= E[(\hat{\theta} - E[\hat{\theta}])^2] + 2E[\hat{\theta} - E[\hat{\theta}]](E[\hat{\theta}] - \theta) + (E[\hat{\theta}] - \theta)^2\end{aligned}$$

右侧第一项为 $\text{Var}[\hat{\theta}]$, 第二项为零, 第三项为 $(\text{bias}[\hat{\theta}])^2$, 因此

当 $\text{Var}[\hat{\theta}] < \infty$ 时, $\text{MSE}[\hat{\theta}] = \text{Var}[\hat{\theta}] + (\text{bias}[\hat{\theta}])^2$

- 我们倾向于选择非偏且方差小的估计量, 但两者往往无法兼得。此时需要根据研究需要进行取舍

对总体方差的估计

Estimation of variance

- 如果已知总体均值 μ , 则 $\tilde{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$ 是 σ^2 的非偏估计量

$$E[\tilde{\sigma}^2] = \frac{1}{n} \sum_{i=1}^n E[(X_i - \mu)^2] = \frac{1}{n} \sum_{i=1}^n \sigma^2 = \sigma^2$$

- 如果 μ 为未知, 那么 $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ 是 σ^2 的非偏估计量吗?

$$\begin{aligned}\hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n (X_i - \mu + \mu - \bar{X})^2 \\&= \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 + \frac{1}{n} \sum_{i=1}^n 2(X_i - \mu)(\mu - \bar{X}) + \frac{1}{n} \sum_{i=1}^n (\mu - \bar{X})^2 \\&= \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 + 2(\bar{X} - \mu)(\mu - \bar{X}) + (\mu - \bar{X})^2 \\&= \tilde{\sigma}^2 - (\bar{X} - \mu)^2\end{aligned}$$

因此, $E[\hat{\sigma}^2] = \sigma^2 - \frac{\sigma^2}{n} = (1 - \frac{1}{n})\sigma^2$

$\Rightarrow s^2 = \frac{n}{n-1} \hat{\sigma}^2$ 是 σ^2 的非偏估计量

s^2 通常称为样本方差 (sample variance)。

也可将 $\hat{\sigma}^2$ 称为样本方差, 此时 s^2 称为偏差修正样本方差 (bias-corrected sample variance)。

注意: s 不是 σ 的非偏估计量

渐进分析

极限

Limit

- 样本量越大，样本中包含的信息越接近总体所包含的信息，因此我们希望获得 $n \rightarrow \infty$ 时估计量的极限特征

数列的极限：如果对于任意 $\varepsilon > 0$, 存在 $n_\varepsilon < \infty$, 使 $|a_n - a| \leq \varepsilon$ 针对所有的 $n > n_\varepsilon$ 成立, 则称 a 为数列 a_n 的极限 (**limit**) , 写作 $\lim_{n \rightarrow \infty} a_n = a$ 。这时我们说 a_n 在 $n \rightarrow \infty$ 时收敛 (**converges**) 于 a

随机变量的极限之概率极限：如果对于任意 $\varepsilon > 0$,

$$\lim_{n \rightarrow \infty} \Pr(|Z_n - c| \leq \varepsilon) = 1 \quad \text{或等价的 } \lim_{n \rightarrow \infty} \Pr(|Z_n - c| > \varepsilon) = 0$$

则称 c 为随机变量的序列 Z_n 的概率极限 (**probability limit**) , 写作 $\text{plim}_{n \rightarrow \infty} Z_n = c$ 或 $Z_n \xrightarrow{P} c$ 。这时我们说 Z_n 在 $n \rightarrow \infty$ 时依概率收敛 (**converges in probability**) 于 c

大数定律

Law of large numbers (LLN)

弱大数定律 (**weak law of large numbers, WLLN**) : 如果 X_i 为 i.i.d. 且 $E[X_i] < \infty$, 则当 $n \rightarrow \infty$ 时

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} E[X_i]$$

即样本均值依概率收敛于总体均值

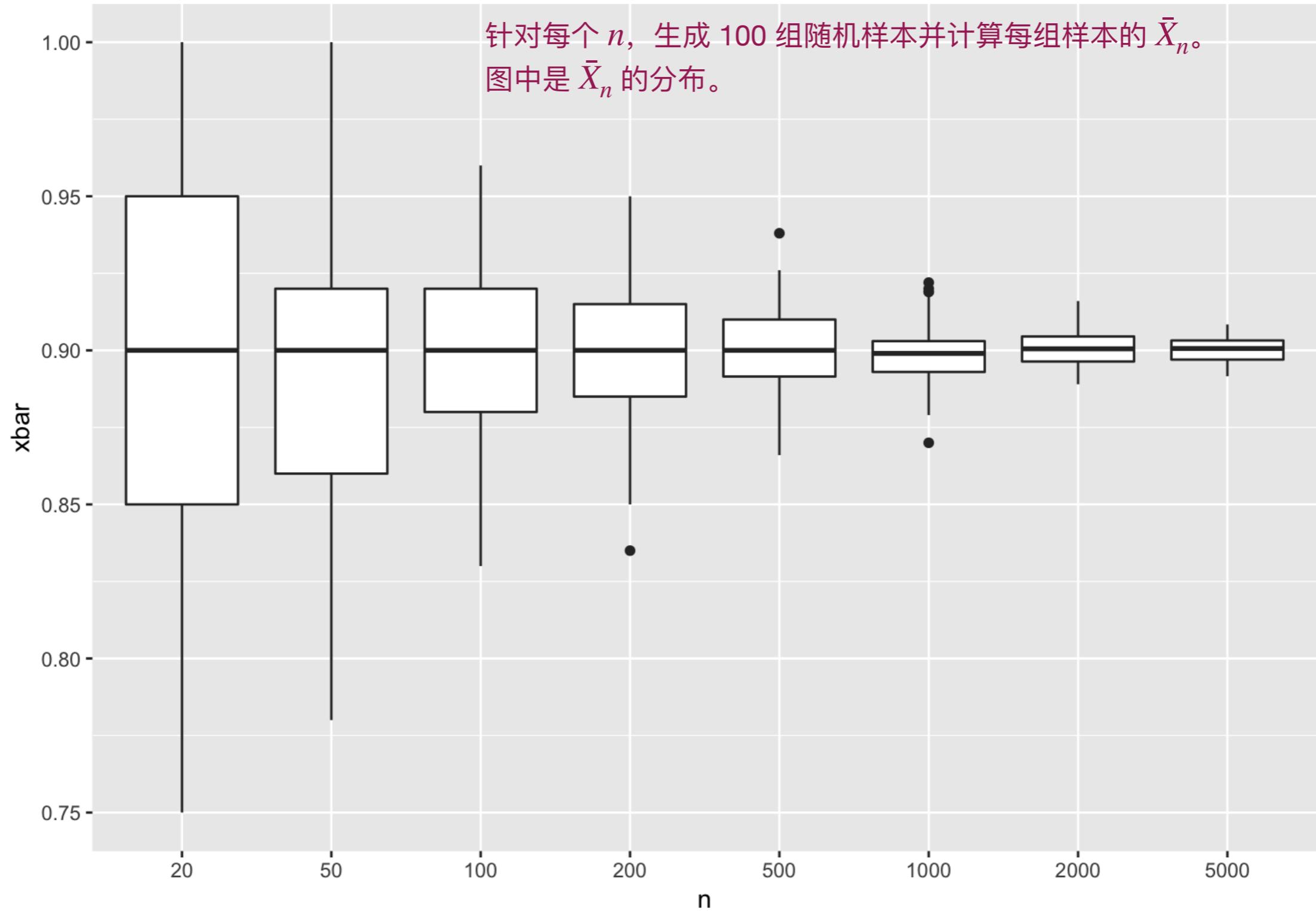
注: i.i.d. 条件可替换为独立样本
且总体方差为有限

如果估计量 $\hat{\theta}$ 在 $n \rightarrow \infty$ 时依概率收敛于 θ (即 $\hat{\theta} \xrightarrow{P} \theta$) , 则称 $\hat{\theta}$ 是 θ 的一致 (**consistent**) 估计量

- 由大数定律可知, 样本均值是总体均值的一致估计量

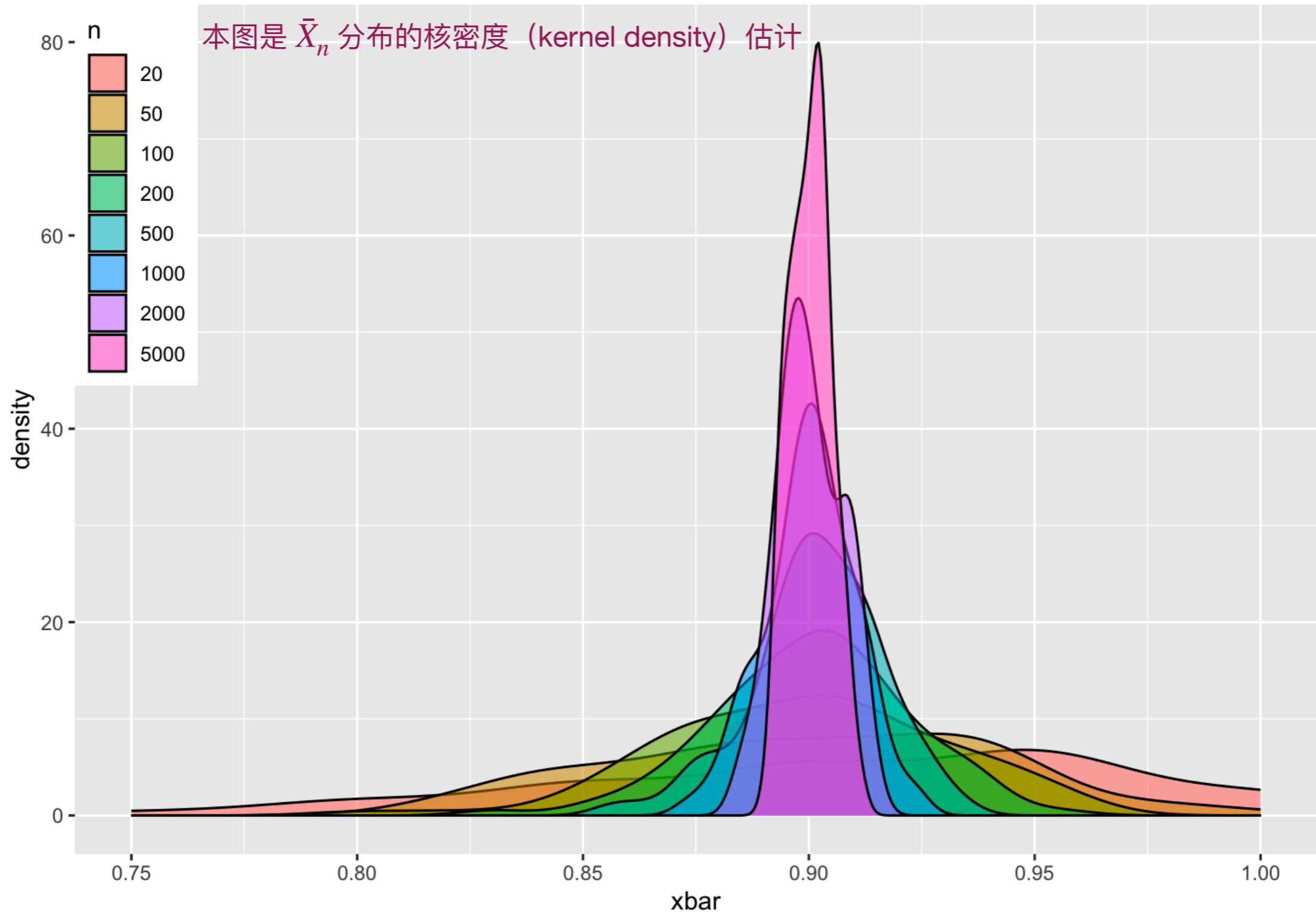
模拟大数定律

$X_i \sim \text{Bernoulli}(p = 0.9)$, $n = 20, 50, 100, 200, 500, 1000, 2000, 5000$



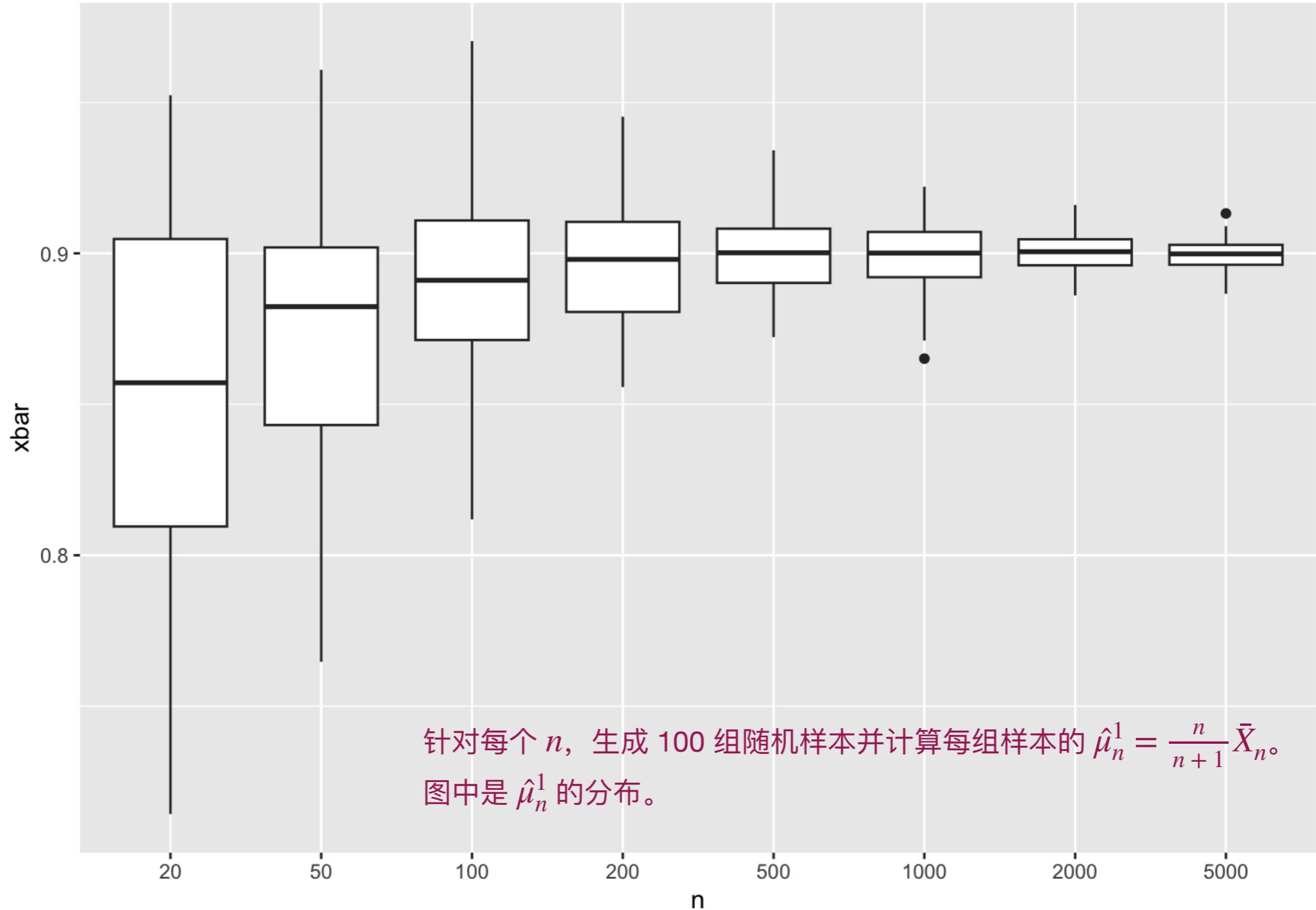
模拟大数定律

$X_i \sim \text{Bernoulli}(p = 0.9)$, $n = 20, 50, 100, 200, 500, 1000, 2000, 5000$



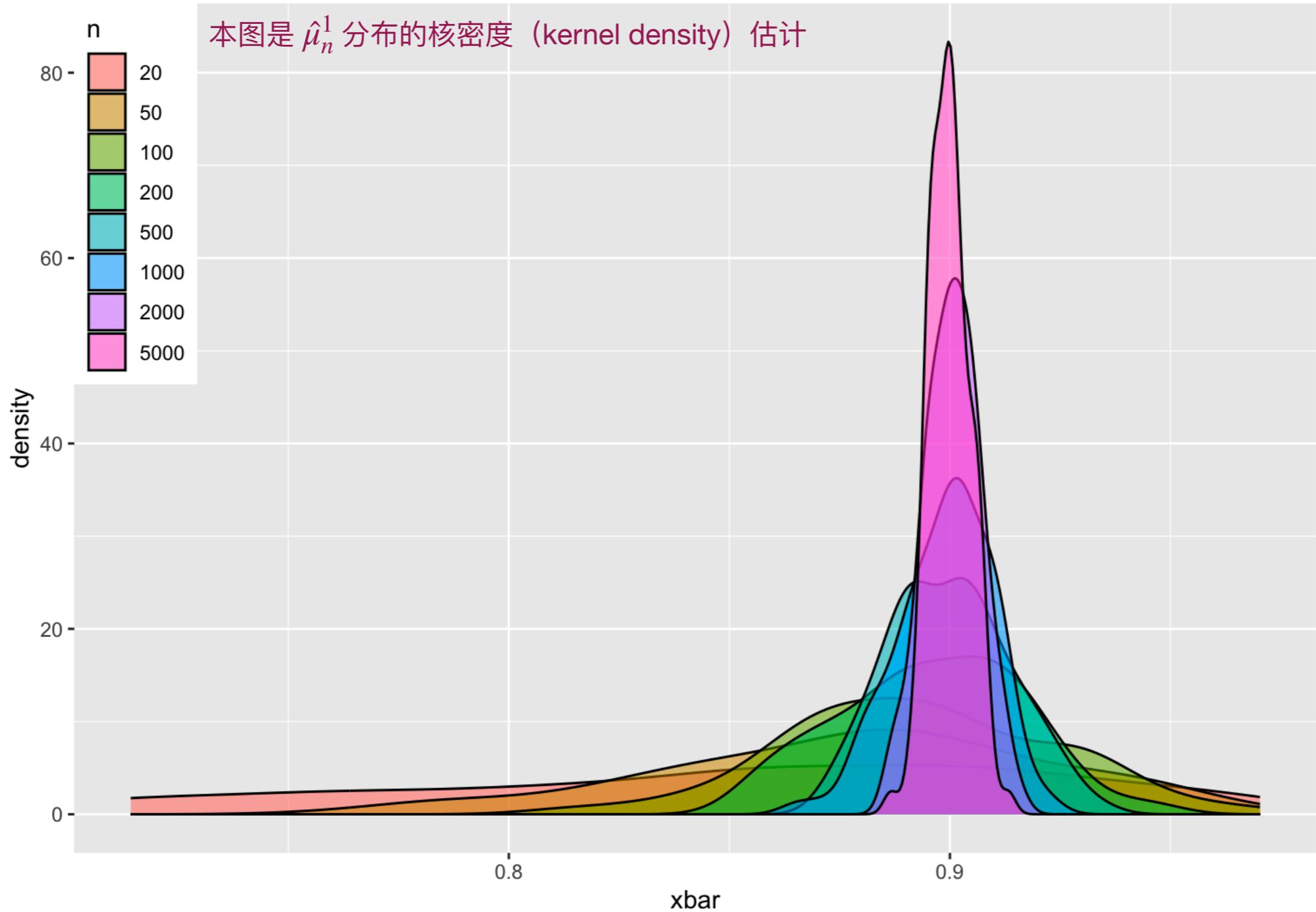
模拟大数定律

$X_i \sim \text{Bernoulli}(p = 0.9)$, $n = 20, 50, 100, 200, 500, 1000, 2000, 5000$



模拟大数定律

$X_i \sim \text{Bernoulli}(p = 0.9)$, $n = 20, 50, 100, 200, 500, 1000, 2000, 5000$



中心极限定理

Central limit theorem (CLT)

随机变量的极限之分布极限： Z_n 服从分布函数 $G_n(u) = \Pr(Z_n \leq u)$ 。如果对任意 u , $G(u) = \Pr(Z \leq u)$ 是连续函数且当 $n \rightarrow \infty$ 时 $G_n(u) \rightarrow G(u)$, 则称 Z_n 依分布收敛 (converges in distribution) 于 Z , 写作 $Z_n \xrightarrow{d} Z$

- 称 $G(u)$ 为 Z_n 的渐进分布 (asymptotic distribution) 或大样本分布 (large sample distribution) 或极限分布 (limit distribution)

Lindeberg-Lévy 中心极限定理：如果 X_i 为 i.i.d. 且 $E[X_i^2] < \infty$, 则当 $n \rightarrow \infty$ 时

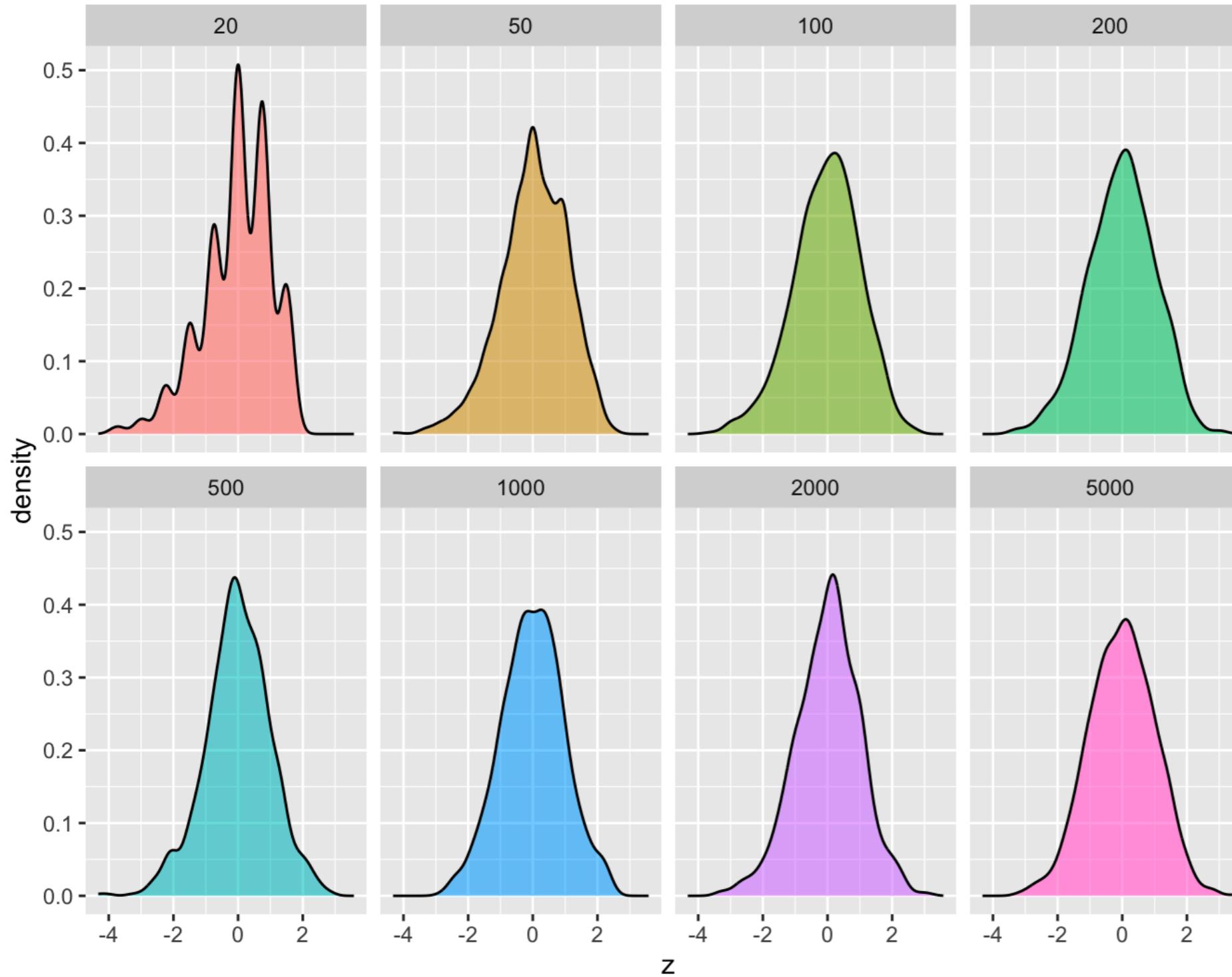
$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} Z \sim N(0, \sigma^2)$$

此处 $\mu = E[X_i]$, $\sigma^2 = E[(X_i - \mu)^2]$, $N(a, b^2)$ 为均值为 a 方差为 b^2 的正态分布

- 在有限样本下, $Z_n = \sqrt{n}(\bar{X}_n - \mu)$ 的均值为 0, 方差为 σ^2 。LLN 告诉我们 Z_n 的渐进分布是正态分布
- CLT 也可以写成 $\bar{X}_n \xrightarrow{a} N\left(\mu, \frac{\sigma^2}{n}\right)$, 即可以用 $N\left(\mu, \frac{\sigma^2}{n}\right)$ 作为 \bar{X}_n 分布的近似

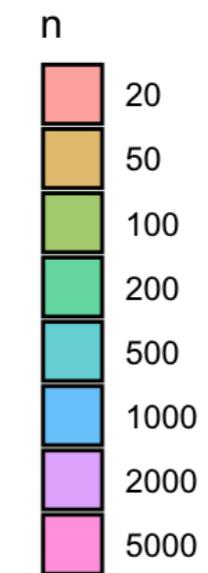
模拟中心极限定理

$X_i \sim \text{Bernoulli}(p = 0.9)$, 针对每个 n 生成 1000 组随机样本



左图是 Z_n/σ 的分布的
核密度 (kernel density) 估计

根据 LLN, 渐进分布为 $N(0,1)$

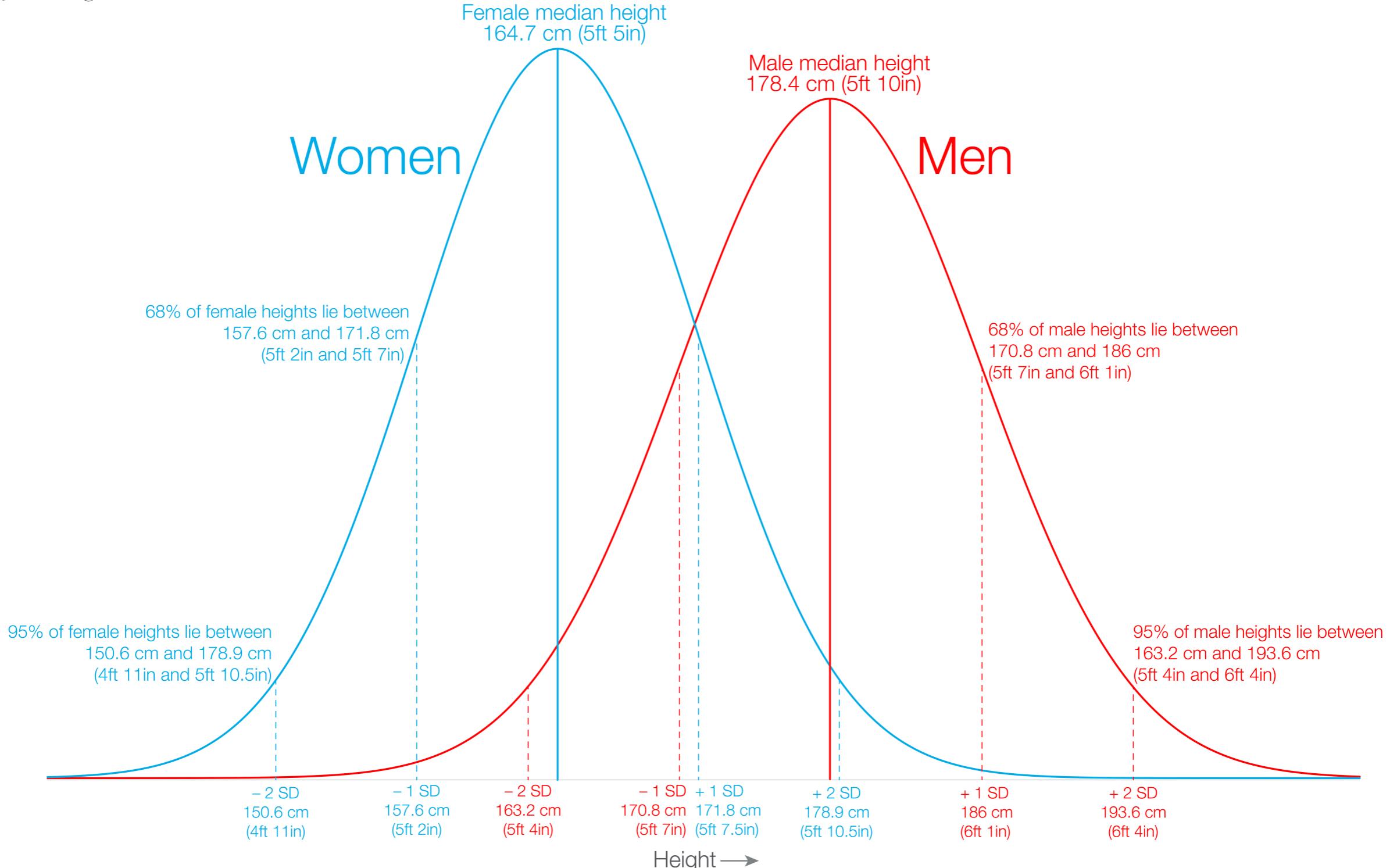


The distribution of male and female heights

The distribution of adult heights for men and women based on large cohort studies across 20 countries in North America, Europe, East Asia and Australia. Shown is the sample-weighted distribution across all cohorts born between 1980 and 1994 (so reaching the age of 18 between 2008 and 2012).

Since human heights within a population typically form a normal distribution:

- 68% of heights lie within 1 standard deviation (SD) of the median height;
- 95% of heights lie within 2 SD.



Note: this distribution of heights is not globally representative since it does not include all world regions due to data availability.

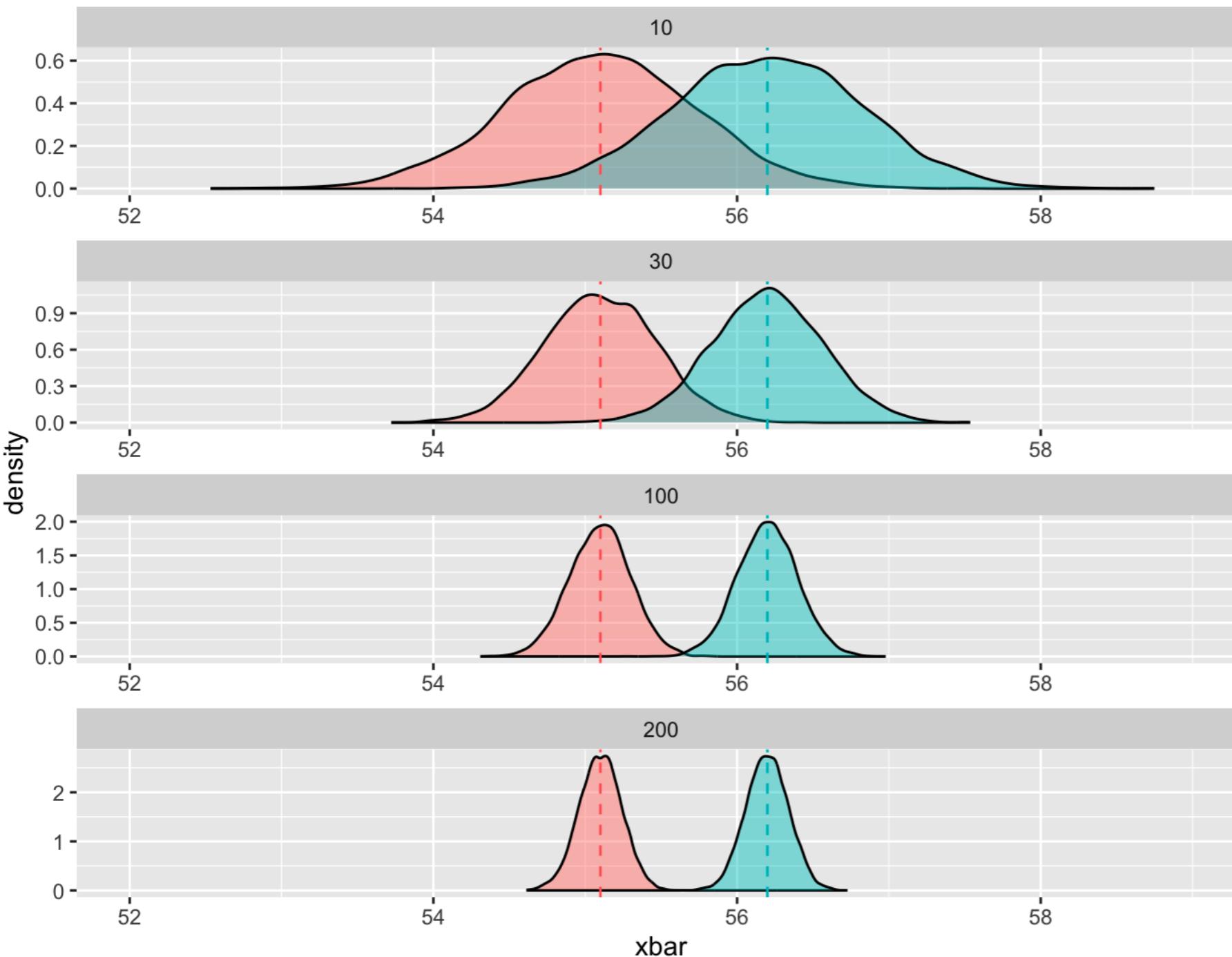
Data source: Jelenkovic et al. (2016). Genetic and environmental influences on height from infancy to early adulthood: An individual-based pooled analysis of 45 twin cohorts.

This is a visualization from OurWorldInData.org, where you find data and research on how the world is changing.

Licensed under CC-BY by the author Cameron Appel.

两个总体的抽样分布

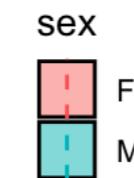
男婴与女婴的体长分布



参考 WHO 婴儿体长分布标准,
出生后 6 个月的婴儿体长为：

男婴： $(\mu, \sigma) = (56.2, 2)$
女婴： $(\mu, \sigma) = (55.1, 2)$

假设二者都服从正态分布，左图为
不同样本量 ($n = 10, 30, 100, 200$) 时
样本均值 \bar{X}_n 的抽样分布的核密度估计



男婴与女婴的实际身长均值之差为

$$1.1 \text{ cm} > 0$$

你如何理解统计结果的统计学意义与
现实意义？