

温兆和 10205501432, 数据科学算法作业1.

1. 解: 抽样间距为: $160/20 = 8$

设起始编号为 x , 则 $x + 15 \times 8 = 126$.

解得: $x = 6$.

4. 解: $14/4 = 3.5$ 故抽样间距为 4.

第1个: 1. 第2个: $1+4=5$. 第3个: $5+4=9$ 第4个: $9+4=13$

5. 解: 应使用分层抽样法. 由样本/总体 $= 100/500 = \frac{1}{5}$.

则 A 车间以系统抽样抽取 $280 \times \frac{1}{5} = 56$ 样本.

B 车间以系统抽样抽取 $95 \times \frac{1}{5} = 19$ 样本.

C 车间以系统抽样抽取 $125 \times \frac{1}{5} = 25$ 样本.

8. 解: 由于要研究血型与色弱的关系, 最好抽取的样本中每种血型人数相等. 且若采用等比分配法, 会造成 AB 型血样本过少. 故最好采用分层抽样法 + 等额分配法. 从四种血型中各随机抽取 $20/4 = 5$ 样本.

13. 解: 应该使用水库抽样法. 维护一个大小为 k 的数组, 先把链表中遍历到的前 k 个元素全部放入数组中. 此后, 对于被遍历到的第 m 个元素 ($m > k$), 在 $1 \sim m$ 中随机地生成一个随机数 i . 若 $i \leq k$, 就把数组中的第 i 个元素替换成当前被遍历到的元素; 否则就把当前遍历到的元素丢弃. 如此循环往复, 直到链表中的所有元素都被遍历过一次.

14. 证: 这种抽样算法等价于: 对每个元素赋予一个随机数后, 对这批总体使用水库排序算法。前 k 个元素直接入数组; 当 $m > k$ 时, 对遍历到的第 m 个元素, 若其随机数大于前数组中元素的随机数的最小值, 就把这个带有 " k 个元素中的最小值" 的元素从数组中替换掉。否则, 它就被丢弃。现在研究当 $m > k$ 时, 遍历到的第 m 个元素的值大于数组中元素的概率是否仍为 $\frac{k}{m}$ 。

对第 $k+1$ 个元素: 前 k 个元素按随机数从小到大排, 间隔为 $(k+1)$ 个缝:

① □ ② □ ... □ ④ $(k+1)$

元素 I_{k+1} 要留在样本中, 只能落入缝 ② ~ 缝 ④, $p = \frac{k}{k+1}$

其余元素要留在样本中: 则不能成为最小值, $p = \frac{k}{k+1}$

假设, 前 m 个元素想留在样本, $p = \frac{k}{m}$, 对元素 I_{m+1} :

I_{m+1} 欲留在样本中: 则要落在前 m 个元素之间最右边的 k 个缝里, $p = \frac{k}{m+1}$

其余元素欲留在样本: ① 它是前 k 大元素中最小者, 且未被替换 ② 其已留在样本中且非最小者。

$$p = \frac{k}{m} \left(\frac{1}{k} \times \frac{m+1-k}{m+1} + \frac{k-1}{k} \times 1 \right) = \frac{k}{m+1}$$

由数学归纳法可知, 引进第 m 个元素时, 所有 m 个元素留在样本中概率均为 $\frac{k}{m}$ 。

故: 该算法为等概率抽样。