

Lecture 21 Normal Approximation To Binomial Distribution & Sampling Distribution of The Sample Proportion

BIO210 Biostatistics

Xi Chen

Spring, 2022

School of Life Sciences

Southern University of Science and Technology



南方科技大学生命科学学院
SUSTech · SCHOOL OF
LIFE SCIENCES

Clinical Infectious Diseases

BRIEF REPORT

Relationship Between the ABO Blood Group and the Coronavirus Disease 2019 (COVID-19) Susceptibility

Jiao Zhao,^{1,a} Yan Yang,^{2,a} Hanping Huang,^{3,a} Dong Li,^{4,a} Dongfeng Gu,¹ Xiangfeng Lu,⁵ Zheng Zhang,² Lei Liu,² Ting Liu,³ Yukun Liu,⁶ Yunjiao He,¹ Bin Sun,¹ Meilan Wei,¹ Guangyu Yang,^{7,b} Xinghuan Wang,^{8,b} Li Zhang,^{3,b} Xiaoyang Zhou,^{4,b} Mingzhao Xing,^{1,b} and Peng George Wang^{1,b}

¹School of Medicine, The Southern University of Science and Technology, Shenzhen,

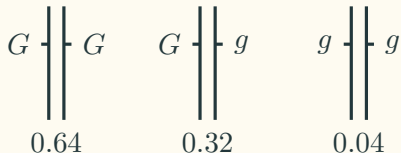
ABO Blood Types And The COVID-19

- The results showed that blood group A was associated with a higher risk for acquiring COVID-19 compared with non-A blood groups, whereas blood group O was associated with a lower risk for the infection compared with non-O blood groups.
- The ABO blood group in 3,694 normal people in Wuhan:

Total	A	B	AB	O
3,694	1,188	920	336	1,250

Data from: Xu P, Xiong Y, Cao K. Distribution of ABO and RhD blood group among Healthy Han population in Wuhan. J Clin Hematol (China). 2015(28):837.

ABO Blood Types



Allele frequency:

$$P(G) = \frac{0.64 \times 2N + 0.32N}{2N} = 0.8$$

$$P(g) = \frac{0.32N + 0.04 \times 2N}{2N} = 0.2$$

Allele	Frequency
I^A	p
I^B	q
i	r

$$p + q + r = 1$$

Genotype	Phenotype	Probability
$I^A I^A$	A	p^2
$I^A i$	A	$2pr$
$I^B I^B$	B	q^2
$I^B i$	B	$2qr$
$I^A I^B$	AB	$2pq$
ii	O	r^2

Estimation of ABO Blood Type Proportions In Wuhan

The ABO blood types in 3,694 normal people in Wuhan:

Total		A	B	AB	O
Number	3694	1,188	920	336	1,250
Proportion	1	0.32	0.25	0.09	0.34

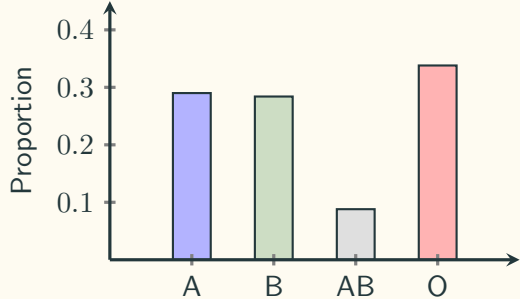
Calculation of allele frequencies
(p, q, r) of (I^A, I^B, i) using the
Hardy-Weinberg model:

$$\begin{cases} p^2 + 2pr = 0.32 \\ q^2 + 2qr = 0.25 \\ 2pq = 0.09 \\ r^2 = 0.34 \end{cases} \Rightarrow \begin{cases} p = 0.23 \\ q = 0.19 \\ r = 0.58 \end{cases}$$

ABO Blood Types Proportions In Han Chinese

Population distribution of ABO blood types in Han Chinese.

Total	A	B	AB	O
592,243	171,473	168,040	52,088	200,642
1	0.290	0.284	0.088	0.338



中国汉族人ABO血型的分布

空军成都医院* 彭德仁

研究ABO血型的分布在医学、法医学及人类学等方面都有重要意义。关于国人的ABO血型分布早在1918年就有人报道过^[1]，现已积累了大量的资料。1963年，尚书颂等曾统计分析了15万多中国人ABO血型的分布资料，提出将我国各省区的ABO血型分布划分为4种类型^[2]。1982年，陈雅勇等收集了1920~1979年国内外发表的中国人的ABO血型分布资料共28万多例，通过计算各群体间的遗传距离，将全国分为4个组^[3]。但这两篇文献都包含了少数民族的资料。Mourant等所著的《人类血型分布》中也仅收集到18万多中国人的ABO血型分布资料^[4]。由于中国是一个多民族的国家，就ABO血型而言，不同的民族可能有不同的分布特点，即使是同一民族，其分布特点因地域等原因也可能不尽相同。为了给医学、法医学及人类学等研究提供一些基本数据，本文收集了1920~1988年国内外发表的有关汉族的ABO血型分布资料共59万多人，并对其进行统计分析。

材料与方法

〈一〉资料来源

国内发表的资料主要取自1963~1966年的《天津医药杂志输血及血液学附刊》、1978~

1979年的《输血及血液学》杂志、1980~1988年的《中华血液学杂志》、1981~1988年的《中华医学检验杂志》等的162篇文献；国外发表的资料主要取自《人类血型分布》^[4]。所收集的资料仅限于汉族，每份资料的人数均多于30人且注明了居住地区，全部资料共计1 022 237人。

〈二〉基因频率的计算与Hardy-Weinberg吻合度测验

对所有收集到的有关资料用赵冠茂推荐的方法来计算ABO基因频率^[5]， p, q, r 分别代表A、B、O基因频率。为了估计调查资料的可靠性，对每份原始资料均作了显著性测验^[6]，如果 $|D/B| \leq 2$ ，表示观察值与期望值无显著性差异，此时 $P \geq 0.05$ （即Hardy-Weinberg吻合度测验观察值与期望值吻合度很好）；如果 $|D/B| > 2$ ，则 $P < 0.05$ ，表示观察值与期望值有显著性差异。AB型的观察值小于期望值， D/B 为正值，反之则为负值。所收集到的资料中除去 $|D/B| > 2$ 的81份外，最后所选择的327份的数据按地区合并，并根据Hirschfeld等提出的 $(A + AB)/(B + AB)$ 公式计算民族指数。但该指数只反映A和B基因的比例，并不能反映差异程度^[7]。

〈三〉遗传距离

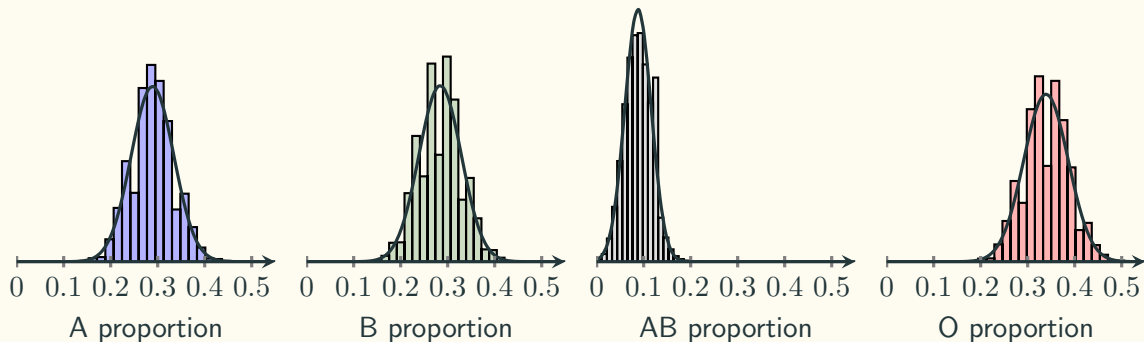
为比较ABO血型分布在各地区间的差异，使用遗传距离 d ，其公式为 $d = 4(1 - \cos \theta)/\pi$

* 邮政编码 610061

Sampling Distribution of ABO Blood Type Proportions

Population distribution of ABO blood types in Han Chinese.

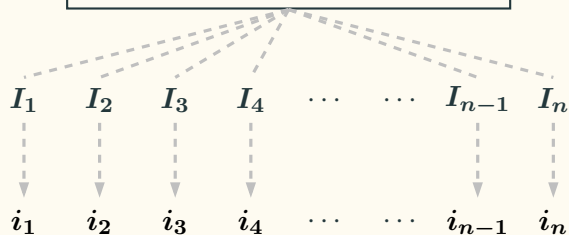
Total	A	B	AB	O
1	0.290	0.284	0.088	0.338



Sampling Distribution of The Sample Proportion

Fraction of type A blood in the population: π

A, A, B, O, O, AB, O, A, O, B, ...



$$Y = \sum_{i=1}^n I_i$$

$$\bar{I} = \frac{1}{n}Y$$

y

\bar{i}

1	0	1	1	1	0	$0.3n$	0.3	sample 1
0	1	1	0	1	1	$0.4n$	0.4	sample 2
0	1	0	0	0	1	$0.35n$	0.35	sample 3

The Sum And The Mean of Indicator Variables

I: Indicator Variable

i.i.d.

$$I_1 \sim \text{Bernoulli}(\pi)$$

$$I_2 \sim \text{Bernoulli}(\pi)$$

$$I_3 \sim \text{Bernoulli}(\pi)$$

\vdots

\vdots

$$I_{n-1} \sim \text{Bernoulli}(\pi)$$

$$I_n \sim \text{Bernoulli}(\pi)$$

Meaning of Y : number of people with blood type A per n people.

Meaning of \bar{I} : The proportion of people with blood type A.

$$Y = \sum_{i=1}^n I_i \sim ? \quad \bar{I} = \frac{1}{n} Y \sim ?$$

By definition:

$$Y \sim B(n, \pi)$$

By The Central Limit Theorem:

$$\bar{I} \dot{\sim} \mathcal{N} \left(\mu = \pi, \sigma^2 = \frac{\pi(1-\pi)}{n} \right)$$

$$Y = n\bar{I} \dot{\sim} \mathcal{N} (\mu = n\pi, \sigma^2 = n\pi(1-\pi))$$

Normal Approximation To A Binomial Distribution

Our knowledge about Han Chinese (Peng, 1991):

Total	A	B	AB	O
1	0.290	0.284	0.088	0.338

A sample from Wuhan (Xu *et al.*, 2015): 1,188 out of 3,694 people have blood type A.

Questions:

1. When draw a random sample ($n = 3,694$), what is the probability of getting 1,100 – 1,200 people with blood type A?
2. When draw a random sample ($n = 3,694$), what is the probability of getting 1,188 people with blood type A?

Normal Approximation To A Binomial Distribution

Question 1:

Use the Binomial probability :

$$\sum_{k=1100}^{1200} \binom{3694}{k} 0.29^k 0.71^{3694-k} = 0.152949$$

Use the Normal probability :

$$P(1100 \leq x \leq 1200) = P\left(\frac{1100 - 1017.26}{27.58} \leq z \leq \frac{1200 - 1017.26}{27.58}\right) = 0.148681$$

Use the Normal probability with **continuity correction** :

$$P(1100 - 0.5 \leq x \leq 1200 + 0.5) = 0.152923$$

Normal Approximation To A Binomial Distribution

Question 2:

Use the Binomial probability :

$$\binom{3694}{1188} 0.29^{1188} 0.71^{3694-1188} = 2.16 \times 10^{-6}$$

Use the Normal probability with continuity correction :

$$P(1188 - 0.5 \leq x \leq 1188 + 0.5) = 1.86 \times 10^{-6}$$

Sampling Distribution of The Sample Proportion

- $\bar{I} \sim$ **Sampling Distribution of The Sample Proportion**
- Generally, we used $p = \frac{x}{n}$ to represent the sample proportion, which is an estimate for the population parameter π .
- According to the **Central Limit Theorem**, when the sample size n is large enough, we have:

$$p \sim \mathcal{N}(\mu_p, \sigma_p^2), \text{ where } \mu_p = \pi, \sigma_p^2 = \frac{\pi(1 - \pi)}{n}$$

Sampling Distribution of The Sample Proportion

