

Biostatistics

Applications in Genetics, Genomics, and other 'omics data

Nuno Sepúlveda, 04.12.2023

Syllabus

1. General review

- a. What is Biostatistics?
- b. Population/Sample/Sample size
- c. Type of Data – quantitative and qualitative variables
- d. Common probability distributions
- e. Work example – Malaria in Tanzania

2. Applications in Medicine

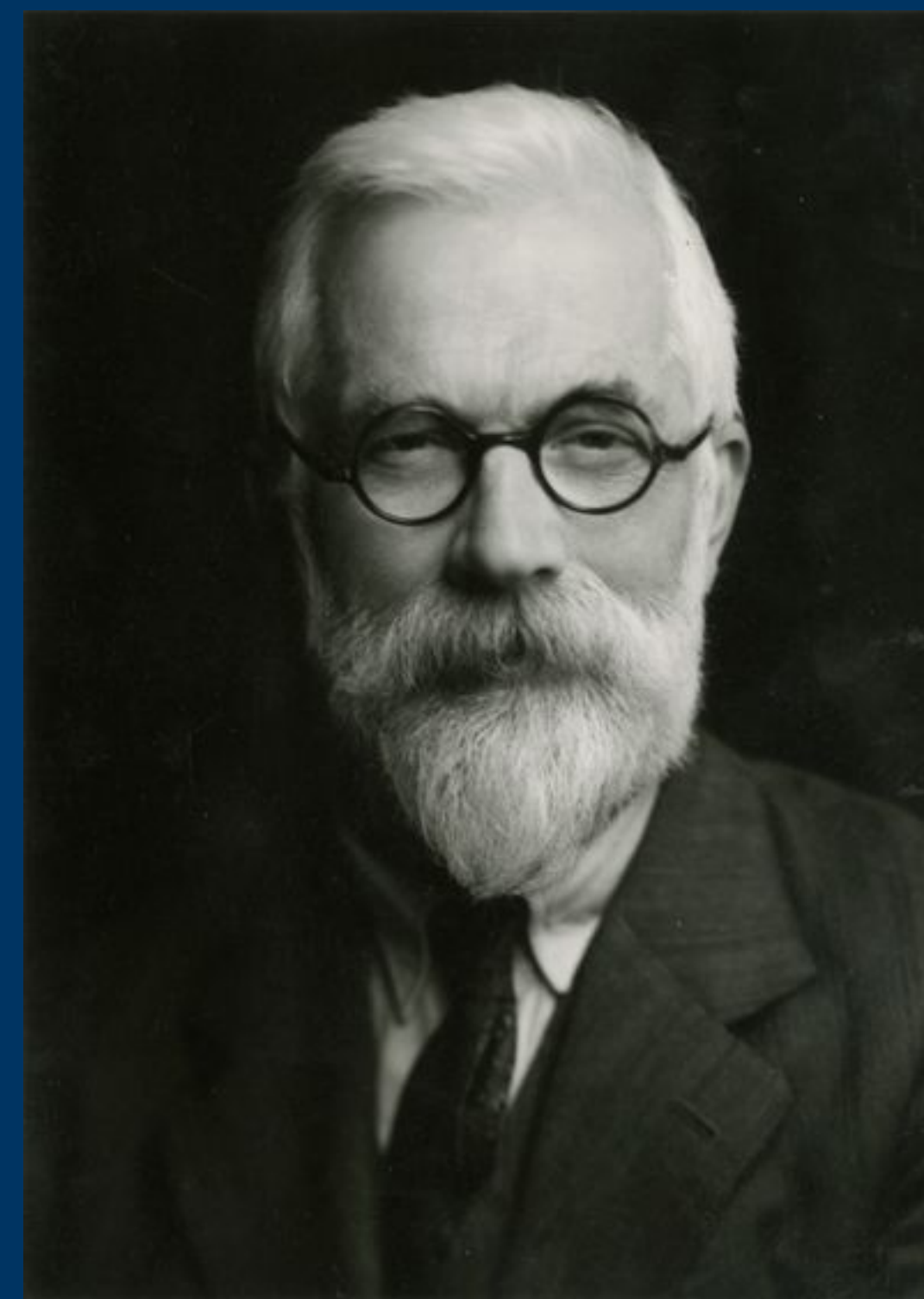
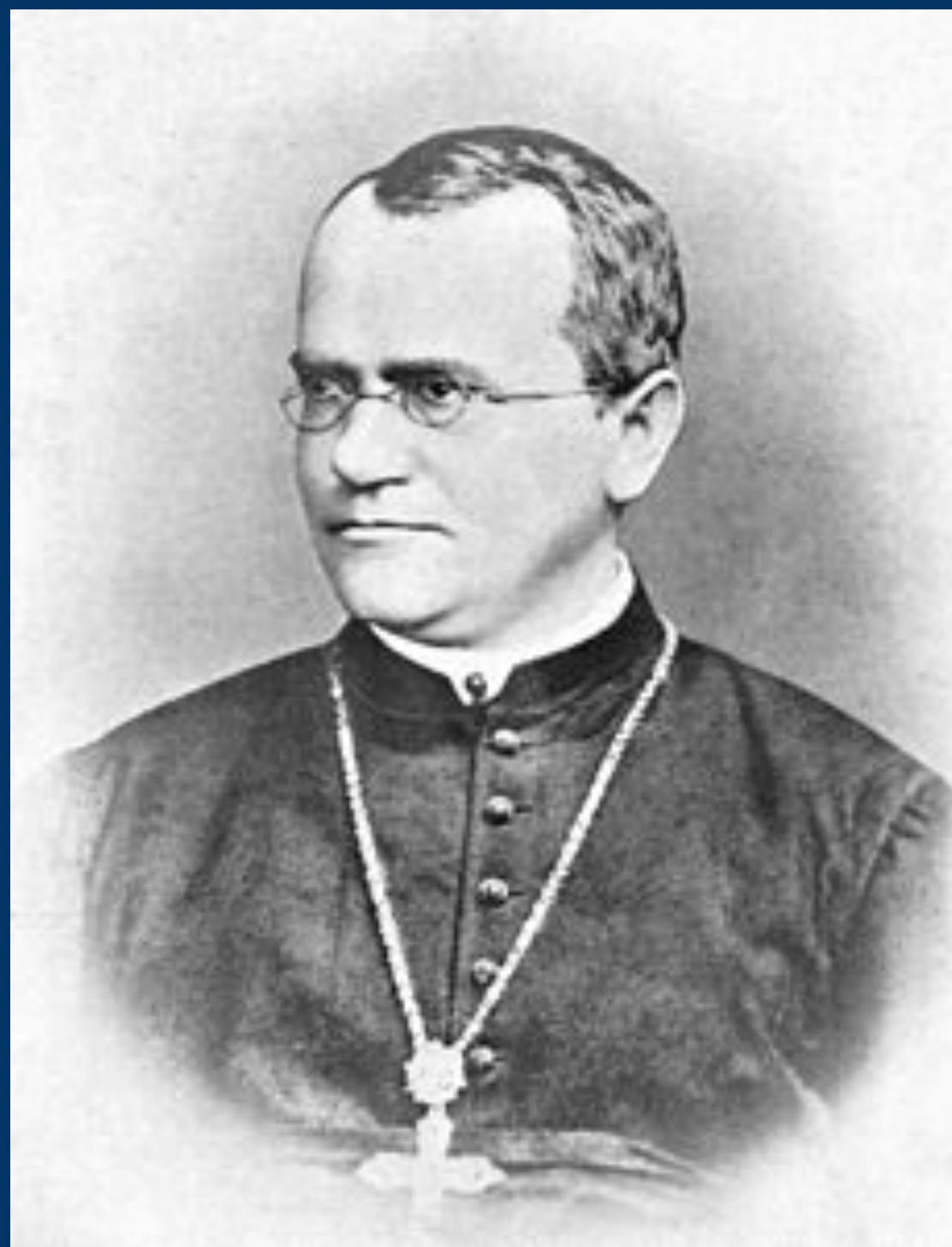
- a. Construction and analysis of diagnostic tools – Binomial distribution, sensitivity, specificity, ROC curve, Rogal-Gladen estimator
- b. Estimation of treatment effects - generalized linear models
- c. Survival analysis - Kaplan-Meier curve, log-rank test, Cox's proportional hazards model

3. Applications in Genetics, Genomics, and other 'omics data

- a. Genetic association studies – Hardy-Weinberg test, homozygosity, minor allele frequencies, additive model, multiple testing correction
- b. Methylation association studies – M versus beta values, estimation of biological age
- c. Gene expression studies based on RNA-seq experiments – Tests based on Poisson and Negative-Binomial

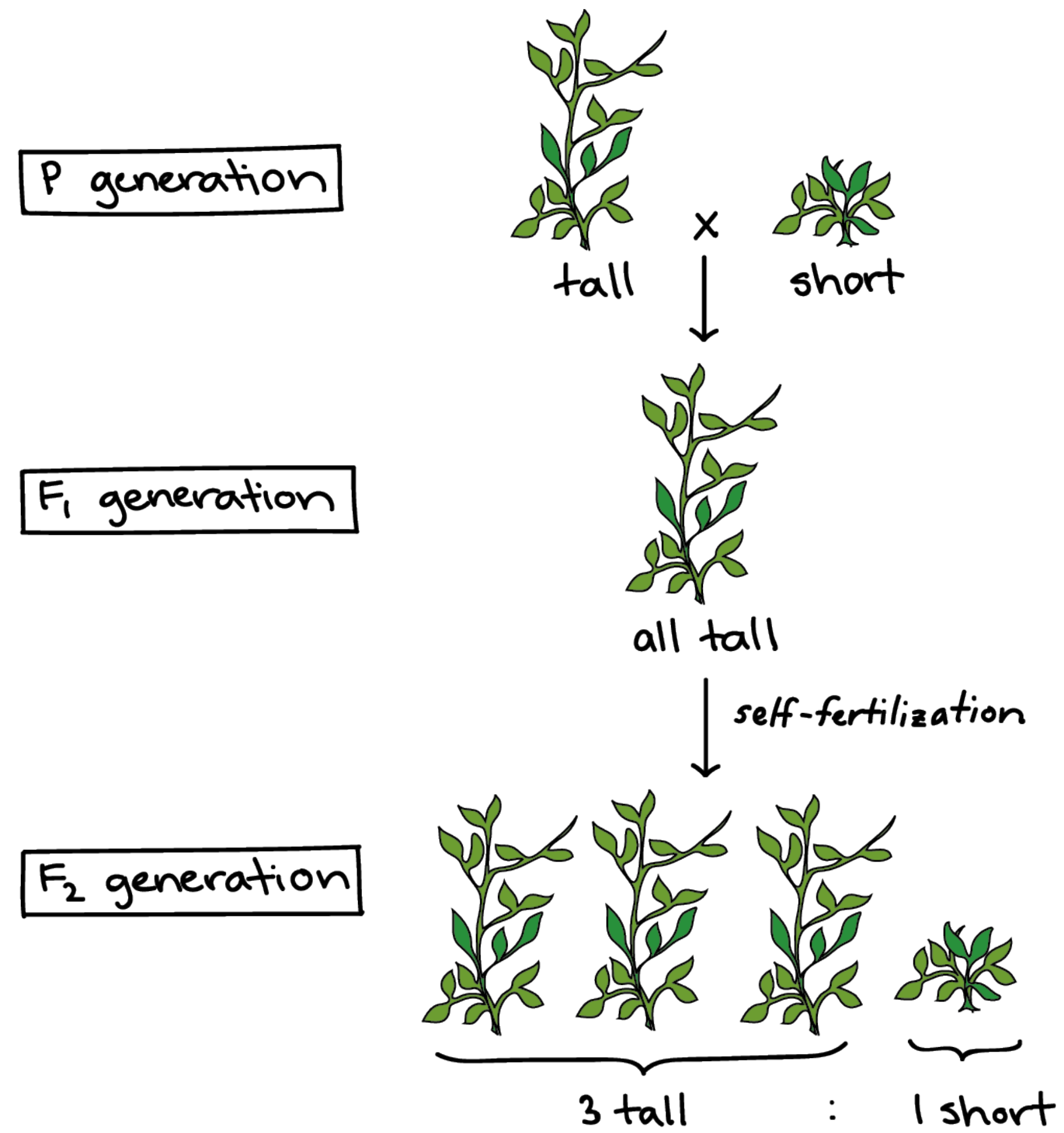
4. Other Topics

- a. Estimation of Species diversity – Diversity indexes, Poisson mixture models
- b. Serological analysis – Gaussian (skew-normal) mixture models
- c. Advanced sample size and power calculations



Do you know these people?

Mendelian genetics



Mendelian genetics

Phenotype /Trait = Biological Characteristic Under study (categorical)

Gene = Unit of Inheritance

Genotype = Composition of gene in terms of alleles

Allele = Variant of a gene

AA

Mendel's idea/interpretation

Generation F0

Phenotype A x Phenotype a



Generation F1

100% Phenotype A

F1 x F1



75%

Phenotype A



25%

Phenotype A

Generation F2

AA x aa



Aa

Aa x Aa



AA



Aa or aA



aa

First two Mendel's laws

The law of Dominance and Uniformity

Some alleles are dominant over the other alleles for a given gene

The law of Segregation

Two alleles for each gene separate from each other during gametogenesis so that the parent may only pass off one allele; thus, the offspring can only inherit one allele from each parent

Exercise 1: data_mendel_single_trait.csv

TABLE 1

Data given in Mendel (1866) for the single trait experiments. “A” (“a”) denotes the dominant (recessive) phenotype; A (a) denotes the dominant (recessive) allele; n is the total number of observations per experiment (that is, seeds for the seed trait experiments and plants otherwise); $n_{“A”}$, $n_{“a”}$, n_{Aa} and n_{AA} denote observed frequencies

	Trait	“A”	“a”	n	Obs. freq.		Theor. ratio
					$n_{“A”}$	$n_{“a”}$	“A” : “a”
F_2	Seed shape	round	wrinkled	7324	5474	1850	3 : 1
	Seed color	yellow	green	8023	6022	2001	3 : 1
	Flower color	purple	white	929	705	224	3 : 1
	Pod shape	inflated	constricted	1181	882	299	3 : 1
	Pod color	yellow	green	580	428	152	3 : 1
	Flower position	axial	terminal	858	651	207	3 : 1
	Stem length	long	short	1064	787	277	3 : 1

Test Mendel’s predictions for each trait using an appropriate statistical test.
Draw your conclusions.

Third Mendel's law

The law of Independent Assortment (law of reassortment)

Alleles of different genes segregate independently of one another during gametogenesis

Bifactorial experiments

Generation F0

Phenotypes A/B x Phenotype a/b



Generation F1

100% Phenotypes A/B

F1 x F1



9:16 Phenotype
A/B

3:16 Phenotype
A/b

3:16 Phenotype
a/B

1:16 Phenotype
a/b

Bifactorial experiments

Combined genotypes

x	BB	Bb	bb
AA	AA/BB	AA/Bb	AA/bb
Aa	Aa/BB	Aa/Bb	Aa/bb
aa	aa/BB	aa/Bb	aa/bb

Bifactorial experiments

Possibilities (n=16)

Cross	BB	Bb	bb
AA	1	2	1
Aa	2	4	2
aa	1	2	1

Bifactorial experiments

Possibilities (n=16)

Cross	BB	Bb	bb
AA	1 Phenotype A/B	2	1 Phenotype A/b
Aa	2	4	2
aa	1 Phenotype a/B	2	1 Phenotype a/b

Bifactorial experiments

Possibilities (n=16)

Cross	BB	Bb	bb
AA	1 Phenotypes A/B	2	1 Phenotypes A/b
Aa	2	4	2
aa	1 Phenotypes a/B	2	1 Phenotypes a/b

Exercise 2:

Test the third Mendel's law predictions for each trait using an appropriate statistical test.

Draw your conclusions.

Mendel-Fisher Controversy

144

HAS MENDEL'S WORK BEEN REDISCOVERED ? *

By R. A. FISHER, M.A., Sc.D., F.R.S.,

Galton Professor of Eugenics, University College, London.

1. THE POLEMIC USE OF THE REDISCOVERY.

THE tale of Mendel's discovery of the laws of inheritance, and of the sensational rediscovery of his work thirty-four years after its publication and sixteen after Mendel's death, has become traditional in the teaching of biology. A careful scrutiny can but strengthen the truth in such a tradition, and may serve to free it from such accretions as prejudice or hasty judgment may have woven into the story. Few statements are so free from these errors as that which I quote from H. F. Roberts' valuable book *Plant Hybridisation before Mendel* (p. 286) :

"The year 1900 marks the beginning of the modern period in the study of heredity. Despite the fact that there had been some development of the idea that a living organism is an aggregation of characters in the form of units of some description, there had been no attempts to ascertain by experiment, how such supposed units might behave in the offspring of a cross. In the year above mentioned the papers of Gregor Mendel came to light, being quoted almost simultaneously in the scientific contributions of three European botanists, De Vries in Holland, Correns in Germany, and Von Tschermak in Austria. Of Mendel's two papers, the important one in this connection, entitled 'Experiments in Plant Hybridization', was read at the meetings of the Natural History Society of Brünn in Bohemia (Czecho-Slovakia) at the sessions of February 8 and March 8, 1865. This paper had passed entirely unnoticed by the scientific circles of Europe, although it appeared in 1866 in the Transactions of the Society. From its publication until 1900, Mendel's paper appears to have been completely overlooked, except for the citations in Focke's 'Pflanzenmischlinge', and the single citation of Hoffmann, elsewhere referred to."

* For further commentary on Mendel's work written by Fisher in 1955, see *Experiments in Plant Hybridisation: Gregor Mendel*. (Ed. J.H. Bennett) Edinburgh: Oliver & Boyd, 1965. As indicated there, all of the years given in Fisher's (1936) reconstruction of the timing of Mendel's experimental programme must be reduced by one.

detail by his paper as a whole. Although no explanation can be expected to be satisfactory, it remains a possibility among others that Mendel was deceived by some assistant who knew too well what was expected. This possibility is supported by independent evidence that the data of most, if not all, of the experiments have been falsified so as to agree closely with Mendel's expectations.

Mendel-Fisher Controversy

Statistical Science
2010, Vol. 25, No. 4, 545–565
DOI: 10.1214/10-STS342
© Institute of Mathematical Statistics, 2010

A Statistical Model to Explain the Mendel–Fisher Controversy

Ana M. Pires and João A. Branco

Abstract. In 1866 Gregor Mendel published a seminal paper containing the foundations of modern genetics. In 1936 Ronald Fisher published a statistical analysis of Mendel's data concluding that "*the data of most, if not all, of the experiments have been falsified so as to agree closely with Mendel's expectations.*" The accusation gave rise to a controversy which has reached the present time. There are reasonable grounds to assume that a certain unconscious bias was systematically introduced in Mendel's experimentation. Based on this assumption, a probability model that fits Mendel's data and does not offend Fisher's analysis is given. This reconciliation model may well be the end of the Mendel–Fisher controversy.

Key words and phrases: Genetics, ethics, chi-square tests, distribution of p -values, minimum distance estimates.



American
Genetic
Association

Journal of Heredity, 2016, 635–646

doi:10.1093/jhered/esw058

Perspective

Advance Access publication August 30, 2016

OXFORD

Perspective

Are Mendel's Data Reliable? The Perspective of a Pea Geneticist

Norman F. Weeden

From the Department of Plant Sciences and Plant Pathology, Montana State University, Bozeman, MT 59717

Address correspondence to N. F. Weeden at the address above or e-mail: nweeden@montana.edu

Received April 11, 2016; First decision May 26, 2016; Accepted August 24, 2016.

Corresponding Editor: John Stommel

Exercise 3:

$$X \rightsquigarrow F \Rightarrow \begin{cases} Y = F(X) \rightsquigarrow \text{Uniform}(0,1) \\ Y = 1 - F(X) \rightsquigarrow \text{Uniform}(0,1) \end{cases}$$

Create a pooled sample of the p-values from exercises 1 and 2 and test whether the p-values are coming from an Uniform distribution.

Draw your conclusions.

First creation of Genotype-Mapping



If you know the genotype, then you know the phenotype

One gene that controls a single binary trait

Mendel triggered the scientific curiosity

What is actually a gene and an allele?

What is the gene involved?

Is it possible to derive genotype-phenotype rules for other type of traits such as the occurrence of a given disease or height?

Some useful concepts

Gene = a stretch of DNA located in a chromosome. The stretch encodes a protein

Allele = variant in the DNA sequence of a gene

Chromosome = a long DNA molecule that contains genetic information of an organism

Genome = the set of all the chromosomes that enables the creation of life

Human Genome = 1-23 autosomal chromosomes, X and Y sexual chromosomes

Beyond Mendelian

Genotype	→	Phenotype	Probabilities
AA		A or a	π_{AA}
Aa		A or a	π_{Aa}
aa		A or a	π_{aa}

Complete penetrance versus incomplete penetrance

What is the gene responsible for this trait?

Genetic mapping: general principle

What is the gene responsible for this trait?

Use experimental cross a la Mendel.

Use genetic markers (with alleles a and b) at known location in the genome.

Test for association between the genotypes and the trait.

Example: Genetic mapping of Type 1 diabetes in mice

data_todd_1991.csv

ARTICLES

Genetic analysis of autoimmune type 1 diabetes mellitus in mice

John A. Todd, Timothy J. Aitman, Richard J. Cornall, Soumitra Ghosh, Jennifer R. S. Hall, Catherine M. Hearne, Andrew M. Knight^{*}, Jennifer M. Love, Marcia A. McAleer, Jan-Bas Prins, Nanda Rodrigues, Mark Lathrop[†], Alison Pressey[‡], Nicole H. DeLarato[‡], Laurence B. Peterson[§] & Linda S. Wicker[‡]

Nuffield Department of Surgery, John Radcliffe Hospital, Headington, Oxford OX3 9DU, UK

^{*} Transplantation Biology, Clinical Research Centre, Watford Road, Harrow, Middlesex HA1 3UJ, UK

[†] CEPH, 27 rue Juliette Dodu, Paris 75010, France

[‡] Autoimmune Diseases Research and [§] Department of Cellular and Molecular Pharmacology, Merck Sharp & Dohme Research Laboratories, Rahway, New Jersey 07065, USA

Genetic mapping: type 1 diabetes in mice

High incidence

Low incidence

NOD

x

(B10.H-2g x NOD) F1

NN

NB

Progeny

NN

NB

Each genetic marker

50%

50%

53 genetic markers across the genome

Exercise 4:

Use an appropriate statistical test and test whether second Mendel’s law apply to the data.

Which the genetic marker has the highest association with the trait?

TABLE 1 A linkage map of the mouse genome and associations of markers with type 1 diabetes													
Chromosome (location, cM)	Locus	Diabetics He	Diabetics Ho	Non- diabetics He	Non- diabetics Ho	$\chi^2 > 4$	Chromosome (location, cM)	Locus	Diabetics He	Diabetics Ho	Non- diabetics He	Non- diabetics Ho	$\chi^2 > 4$
1 (3)	<i>D1Nds4</i>	38	58	57	37	8.4	9 (24)	<i>Thy-1</i>	41	56	49	47	
1 (41)	<i>Bcl-2</i>	45	52	49	47		9 (29)	<i>Ncam</i>	39	58	48	49	
1 (42)	<i>D1Nds2</i>	44	49	50	47		9 (33)	<i>Cyp1a2</i>	39	58	49	48	
1 (48)	<i>D1Nds1</i>	50	47	49	46		9 (44)	<i>D9Nds2</i>	44	53	45	52	
1 (73)	<i>Crp</i>	57	40	45	51		9 (46)	<i>D9Nds1</i>	44	52	46	48	
2 (35)	<i>D2Nds1</i>	40	46	18	31		10 (29)	<i>D10Nds1</i>	47	50	49	46	
2 (46)	<i>B2m</i>	39	44	19	29		11 (10)	<i>Glns</i>	46	51	51	46	
							11 (42)	<i>Acrb</i>	31	66	51	46	8.4
3 (32)	<i>Il-2</i>	33	64	49	48	5.4	11 (47)	<i>D11Nds1</i>	30	67	51	46	9.3
3 (53)	<i>D3Nds1</i>	17	80	54	43	30.4	11 (52)	<i>Mpo</i>	36	61	53	44	6.0
3 (67)	<i>Tshb</i>	24	73	50	47	14.8	11 (68)	<i>Gfap</i>	36	61	51	46	4.7
3 (86)	<i>Adh-1</i>	28	69	49	48	9.5	11 (71)	<i>Myla</i>	36	61	50	47	4.1
4 (18)	<i>D4Nds3</i>	44	40	9	10		12 (4)	<i>Odc</i>	54	43	47	50	
4 (29)	<i>Mup-1</i>	43	43	21	28		12 (45)	<i>Mtv-9</i>	36	41	13	16	
4 (30)	<i>Orm-1</i>	44	42	21	28		13 (20)	<i>Hist1</i>	42	32	14	21	
4 (62)	<i>D4Nds2</i>	40	56	55	40		13 (39)	<i>D13Nds1</i>	58	39	36	60	9.6
4 (69)	<i>Lck</i>	44	53	52	41		13 (68)	<i>P198-13</i>	45	27	?	?	
4 (95)	<i>Pnd</i>	11	20	25	15		14 (8)	<i>Plau</i>	68	29	43	54	13.2
							14 (27)	<i>Tcra</i>	61	36	42	52	6.4
5 (10)	<i>D5Nds1</i>	50	47	51	45		14 (38)	<i>Nfl</i>	52	45	45	47	
5 (30)	<i>D5Nds2</i>	51	43	53	43		14 (42)	<i>Hpg</i>	55	42	47	49	
5 (46)	<i>Afp</i>	49	37	22	27		15 (18)	<i>Myc</i>	38	59	48	46	
5 (94)	<i>Zp-3</i>	41	36	28	20		15 (24)	<i>D15Nds1</i>	37	60	50	45	4.1
6 (32)	<i>Ly-3</i>	38	48	22	27		15 (27)	<i>Ly-6C</i>	38	59	51	45	
6 (68)	<i>Prp</i>	36	60	30	24	4.6	15 (49)	<i>Gdc-1</i>	32	48	23	19	
							15 (53)	<i>Hox-3</i>	40	57	52	44	
7 (6)	<i>Ckmm</i>	64	33	41	55	10.5	16 (42)	<i>D16Nds2</i>	26	31	16	18	
7 (27)	<i>Ngfg</i>	53	44	37	54		18 (24)	<i>Fim-2</i>	37	40	21	17	
7 (48)	<i>D7Nds2</i>	42	53	38	58		18 (29)	<i>Ii</i>	38	46	20	18	
7 (64)	<i>Hbb</i>	39	55	44	50		19 (35)	<i>Cyp2c</i>	43	37	17	21	
8 (0)	<i>Polb</i>	43	43	21	28		X (23)	<i>Hprt</i>	29	28	25	14	
8 (35)	<i>Mt-2</i>	37	49	26	23		X (39)	<i>DXNds3</i>	28	29	27	16	

Discussion:

What is the statistical challenge of genetic mapping?

Multiple testing problem

In absence of association

$\alpha = 0.05$ $m =$ number of genetic markers (statistical tests)

$Y =$ Number of significant tests $| H_0, \alpha = 0.05 \rightsquigarrow \text{Binomial}(m, p = \alpha)$

Expected number of false positive associations

$$E[Y | H_0, \alpha] = m \times \alpha$$

$$E[Y | H_0, \alpha] = 53 \times 0.05 = 2.65$$

Dealing with multiple testing (classical methods)

Redefine the type I error for the overall analysis

$$P[Y \geq 1 | H_0, \alpha^*] = \alpha$$

$$E[Y | H_0, \alpha^*] = \alpha$$

$$1 - (\alpha^*)^m = \alpha \Leftrightarrow \alpha^* = 1 - (1 - \alpha)^{1/m}$$

$$m \times \alpha^* = \alpha \Leftrightarrow \alpha^* = \frac{\alpha}{p}$$

Sidak-Dunn correction

Bonferroni correction

$$\alpha^* = 1 - (1 - \alpha)^{1/53} \approx 0.00084$$

$$\alpha^* = \frac{0.05}{53} = 0.00082$$

In the previous exercise, was the strongest association statistically significant controlling for multiple testing?