Introduction to linkage and association

Dmitry Prokopenko

dprokopenko@mgh.harvard.edu

- Mon Apr 1, 2019 Introduction to linkage and association
- Wed Apr 3, 2019 Genome-Wide and Family-Based Association Studies

Contents

- Definitions
- Genetic data overview
- What is linkage
- Linkage methods
- What is association
- LD and other concepts
- Linkage versus association

DNA

- Deoxyribonucleic acid double stranded helix molecule
- Organized into long structures chromosomes
- Every chromosome carries 2 copies of information
- Length of unwrapped human DNA ~ 6ft (1.8 m) per cell



Human Genome



- 2 copies of 22 autosomes
 + 2 sex chromosomes
- ~3,200 Mb (Megabasepairs) = 3,200,000,000 basepairs
- ~20,000 protein-coding genes

https://commons.wikimedia.org/w/index.php?c urid=18594472

Central dogma of molecular biology



- DNA->RNA->Aminoacid->Protein
- Coding exons ~ 1 % of human genome
- 99.9% of human genomes same
- ~ 12M mutations (SNPs) with a frequency > 1% in population

https://commons.wikimedia.org/w/index.php?c urid=32026515

Genetic data representation

.fasta

	sam	∖/ba	۹m											
@HD	VN:1.0	50:coor	dinate	2										
@SQ	SN:chr2	0	LN:64	444167										
@PG	ID:TopH	at	VN:2.	0.14	CL:	/srv/dna_	_tools/t	ophat/1	tophat	-N 3	3read	-edit-d	list 5	read-rea
ligr	-edit-dist	2 -i 50	-I 500	90ma	x-covera	ge-intro	n 5000 -	M - O OI	ut /da	ta/us	ser446/m	apping_	tophat/i	ndex/chr
20 /	′data/user44	6/mappin	g_toph	nat/L6_	18_GTGAA	A_L007_R	L_001.fa	stq						
HWI	ST1145:74:0	101DACXX	:7:110	02:4284	:73714	16	chr	20 19	90930	3	100	1 *	0	0
	CCGTGTTT	AAAGGTGG	ATGCGG	STCACCT	TCCCAGCT	AGGCTTAG	GATTCTT	AGTTGG	CCTAGG	AAAT	CAGCTAG	ГССТСТО	TCTCAGTC	сссстст
С	BBDCCDDCCD	DDDCDDDD	DDCDCC	CDBC?D	DDDDDDDD	DDDDDDCCI	OCDDDDDD	DDDDCC	CCEDDD	C?DDI	DDDDDDDD	DDDDDDD	DDDDDBDHF	FFFDC@@
	AS:i:-15	XM:	i:3 X	(0:i:0	XG:i:0	MD:Z:550	C20C13A9	NM:i:3	3 NH:	i:2	CC:Z:=	CP:i:5	5352714	HI:i:0
HWI ·	ST1145:74:0	101DACXX	:7:111	4:2759	:41961	16	chr	20 19	93953	50	100	1 *	0	0
	TGCTGGAT	CATCTGGT	TAGTGG	SCTTCTG	ACTCAGAG	GACCTTCG	ICCCCTGG	GGCAGT	GACCT	TCCA	GTGATTCC	CCTGACA	TAAGGGGC	ATGGACGA
G	DCDDDDEDDD	DDDDDCDDD	DDDDCC	CDDDCD	DDDDEEC>	DFFFEJJJ	JJIGJJJJ	IHGBHH	201000	JJJG	JJJIJJJJ.	JIHJJJJ	JJHHHHHF	FFFFCCC
	AS:i:-16	XM:	i:3 X	(0:i:0	XG:i:0	MD:Z:600	G16T18T3	NM:i:3	3 NH:	i:1				
HWI ·	ST1145:74:0	101DACXX	:7:120	04:1476	0:4030	16	chr	20 27	70877	50	100	1 *	0	0
	GGCTTTAT	TGGTAAAA	AAGGAA	TAGCAG	ATTTAATC	AGAAATTC	CACCTGG	CCCAGC	AGCACC	AACCA	AGAAAGAA	GGGAAGA	AGACAGGA	АААААССА
С	DDDDDDDDD	CDDDDDDD	DDDEEE	EEEEFF	FEFFEGHH	HHFGDJJI	JJJIJIJJ	JIIIIG	GFJJIH	IIII	JJJJJJIG	HFAHGF	HJHFGGHF	FFDD@BB
	AS:i:-11	XM:	i:2 X	(0:i:0	XG:i:0	MD:Z:0A8	35G13	NM:i:2	2 NH:	i:1				
HWI -	ST1145:74:0	101DACXX	:7:121	0:1116	7:8699	0	chr	20 2	71218	50	50M	4700N50)M *	Θ
	Θ	GTGGCTCT	ТССАСА	GGAATG	TTGAGGAT	GACATCCA	FGTCTGGG	GTGCAC	TTGGGT	стсс	GAAGCAGA	АСАТССТ	CAAATATG	ACCTCTCG
acce	pted_hits.s	am												

.ped/.bed

	•
1	50103 - 5010004 - 5090005 - 5090004 - 2 - 2 - G - G - C - C - G - G - C - C - A - A - G - G - C - C - T - A - A - G - G - C - T - T - A - A - A - A - A - A - C - C - G - G - C - C - T - T - G - G - G - G - G - G - G - G
2	50103 · 501005 · 509005 · 509005 ·
3	50105 · 5010049 · 5090021 · 5090022 · 2 · 2 · G · G · C · C · G · G · C · C · A · A · T · T · A · A · G · G · G · G · C · T · T · A · A · G · G · C · C · G · G · C · C · T · T · G · G · C · C · T · T · G · G · G · G · G · G · G · G
4	50105 · 5010070 · 5090020 · 5090019 · 1 · 2 · G · G · C · C · G · G · C · C · A · A · G · G · A · A · G · G · C · C · T · T · A · A · G · G · T · T · A · A · A · A · A · C · C · G · G · C · C · T · T · G · G · C · C · T · T · G · G · G · G · G · G · G · G
5	50105 · 5010137 · 5090020 · 5090019 · 1 · 1 · G · G · C · C · G · G · C · C · A · A · G · G · A · A · G · G · C · C · T · T · A · A · G · G · T · T · A · A · A · A · C · C · G · G · C · C · T · T · G · G · G · G · G · G · G · G
6	50115 · 5010018 · 5090039 · 5090040 · 2 · 2 · G · G · C · C · G · G · C · C · A · A · G · G · A · A · G · G · C · C · T · T · A · A · G · G · C · C · C · A · A · A · A · C · C · G · G · C · C · T · T · T · T · G · G · C · C · T · T · T · T · G · G · G · G · G · G
7	50115 · 5010025 · 5090042 · 5090041 · 2 · 2 · G · G · C · C · G · G · C · C · A · A · G · G · A · A · G · G · C · C · T · T · A · A · G · G · C · C · C · A · A · A · A · C · C · G · G · C · C · T · T · T · T · G · G · C · C · T · T · T · T · G · G · G · G · G · G
8	50115 · 5010031 · 5090039 · 5090040 · 2 · 2 · G · G · C · C · G · G · C · C · A · A · G · G · A · A · G · G · C · C · T · T · A · A · G · G · C · C · C · A · A · A · A · C · C · G · G · C · C · T · T · T · T · G · G · C · C · T · T · T · T · G · G · G · G · G · G
9	50115 · 5010039 · 5090042 · 5090043 · 2 · 2 · G · G · C · C · G · G · C · C · A · A · G · G · A · A · G · G · C · T · T · A · A · G · G · C · C · A · A · A · A · A · C · C · G · G · C · C · T · T · G · G · C · C · T · T · T · G · G · G · G · G · G · G
10	50115 · 5010196 · 5090039 · 5090040 · 2 · 1 · G · G · C · C · G · G · C · C · A · A · G · G · A · A · G · G · C · C · T · T · A · A · G · G · C · C · C · A · A · A · A · C · C · G · G · C · C · T · T · T · T · G · G · C · C · T · T · T · T · G · G · G · G · G · G
11	50115 · 5010237 · 5090039 · 5090040 · 1 · 1 · G · G · C · C · G · G · C · C · A · A · G · G · A · A · G · G · C · C · T · T · A · A · G · G · C · C · T · A · A · A · A · C · C · G · G · C · C · T · T · T · G · G · C · C · T · T · T · G · G · G · G · G · G · G
12	50115 · 5010397 · 5090039 · 5090040 · 2 · 2 · G · G · C · C · G · G · C · C · A · A · G · G · A · A · G · G · C · C · T · T · A · A · G · G · C · C · T · A · A · A · A · C · C · G · G · C · C · T · T · T · G · G · C · C · T · T · T · G · G · G · G · G · G · G
13	50116 - 5010009 - 5410002 - 5410003 - 2 - 2 - 6 - G - C - C - G - G - C - C - A - A - G - G - A - A - G - G - C - C - T - T - A - A - G - G - T - T - A - A - A - A - C - C - G - G - C - C - T - T - G - C - C - T - T - G - G - G - G - G - G - G - G
14	50116 - 5010063 - 5410002 - 5410003 - 1 - 2 - 6 - G - C - C - G - G - C - C - A - A - G - G - A - A - G - G - C - C - T - T - A - A - G - G - T - T - A - A - A - A - C - C - G - G - C - C - T - T - G - G - G - G - G - G - G - G
15	50116 · 5010071 · 5410002 · 5410003 · 2 · 1 · 6 · 6 · C · C · 6 · 6 · C · C · A · A · 6 · 6 · A · A · 6 · 6 · C · C · T · T · A · A · 6 · 6 · C · T · A · A · A · A · A · C · C · 6 · 6 · C · C · T · T · T · 6 · C · C · T · T · 6 · 6 · 6 · 6 · 6 · 6 · 6 · 6
16	50116 - 5010200 - 5410002 - 5410002 - 5410003 - 1 - 1 - G - G - C - C - G - G - C - C - A - A - G - G - C - C - T - T - A - A - G - G - T - T - A - A - A - A - A - C - C - G - G - C - C - T - T - G - G - G - G - G - G - G - G
17	50119 · 5010013 · 5090075 · 5090075 · 5090074 · 2 · 2 · G · G · C · C · G · G · C · C · A · A · G · G · C · C · T · T · A · A · G · G · T · T · A · A · A · A · A · C · C · G · G · C · C · T · T · T · G · C · C · T · T · G · G · G · G · G · G · G · G
18	50122 · 5010014 · 5090109 · 5090108 · 2 · 2 · G · G · C · C · G · G · C · C · A · A · G · G · A · A · G · G · C · C · T · T · A · A · G · G · C · C · T · A · A · A · A · A · C · C · G · G · C · C · T · T · T · G · G · G · G · G · G · G
19	50122 · 5010104 · 5090109 · 5090108 · 2 · 2 · G · G · C · C · G · G · C · C · A · A · G · G · A · A · G · G · C · C · T · T · A · A · G · G · C · C · T · A · A · A · A · A · C · C · G · G · C · C · T · T · T · G · G · G · G · G · G · G
20	50123 · 5010028 · 5090122 · 5090123 · 2 · 2 · G · G · C · C · G · G · C · C · A · A · G · G · A · A · G · G · C · C · T · T · A · A · G · G · C · C · T · A · A · A · A · A · C · C · G · G · C · C · T · T · T · G · G · G · G · G · G · G
21	50123 · 5010029 · 5090122 · 5090123 · 1 · 2 · G · G · C · C · G · G · C · C · A · A · G · G · A · A · G · G · C · C · T · T · A · A · G · G · C · C · C · A · A · A · A · C · C · G · G · C · C · T · T · T · T · G · G · G · G · G · G
22	50123 · 5010132 · 5090122 · 5090123 · 2 · 1 · G · G · C · C · G · G · C · C · A · A · G · T · A · A · G · G · C · T · T · A · A · G · G · C · T · A · A · A · A · C · C · G · G · C · C · T · T · T · G · G · C · C · T · T · T · G · G · G · G · G · G · G
23	50127 · 5010007 · 5090143 · 5090143 · 5090142 · 1 · 2 · G · G · C · C · G · G · C · C · A · A · G · G · G · G · G · G · G · G
24	50127 · 5010041 · 5090153 · 5090153 · 5090152 · 2 · 2 · G · G · C · C · G · G · C · C · A · A · T · T · A · A · G · G · C · C · T · T · A · A · G · G · T · T · A · A · A · G · G · C · C · T · T · G · G · C · C · T · T · G · G · G · G · G · G · G · G
25	50127 · 5010328 · 5090143 · 5090143 · 5090142 · 1 · 2 · 6 · 6 · C · C · 6 · 6 · C · C · 6 · 6
26	50130 · 5010056 · 5090159 · 5090159 · 5090158 · 2 · 2 · G · G · C · C · G · G · C · C · A · A · G · G · A · A · G · G · C · C · T · T · A · A · G · G · T · T · A · A · A · A · A · C · C · G · G · C · C · T · T · T · G · G · G · G · G · G · G
27	50130 - 5010279 - 5090159 - 5090158 - 1 - 2 - G - G - C - C - G - G - C - C - A - A - G - G - C - C - A - A - G - G - C - C - T - T - A - A - A - A - A - A - A - C - C - G - G - C - C - T - T - T - G - G - G - G - G - G - G
28	Sense - S
<	>

.vcf/.bcf

##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">

##FORMAT=<ID=FA0,Number=A,Type=Integer,Description="Flow Evaluator Alternate allele observation count">

##FORMAT=<ID=FDP,Number=1,Type=Integer,Description="Flow Evaluator Read Depth">

##FORMAT=<ID=FRO,Number=1,Type=Integer,Description="Flow Evaluator Reference allele observation count">
##FORMAT=<ID=FSAF,Number=A,Type=Integer,Description="Flow Evaluator Alterna
##FORMAT=<ID=FSAR,Number=A,Type=Integer,Description="Flow Evaluator Alterna
##FORMAT=<ID=FSRF,Number=1,Type=Integer,Description="Flow Evaluator referend
##FORMAT=<ID=FSRR,Number=1,Type=Integer,Description="Flow Evaluator referend
##FORMAT=<ID=FSRR,Number=1,Type=Integer,Description="Flow Evaluator referend
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality, the Phi
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">

##FORMAT=<ID=RO,Number=1,Type=Integer,Description="Reference allele observations" ##FORMAT=<ID=SAF,Number=A,Type=Integer,Description="Alternate allele observa ##FORMAT=<ID=SAR,Number=A,Type=Integer,Description="Alternate allele observ. ##FORMAT=<ID=SRF, Number=1, Type=Integer, Description="Number of reference obse ##FORMAT=<ID=SRR,Number=1,Type=Integer,Description="Number of reference obse #CHROM POS ID REF ALT QUAL FILTER INFO FORMAT Sample1 Sample2 Samp 4476.14 PASS chr1 2488153 . А G AC=4;AF=1.00;AN=4;DP=648;FDP=19 chr1 2491258 . 2611.42 PASS AC=2;AF=0.500;AN=4;AO=146;DP=13 С G chr1 6528100 . GGCCCCT GGCCCTC 10278.10 PASS AC=2;AF=1.00;AN=2;A chr1 6528468 . 1859.16 PASS AC=2;AF=0.500;AN=4;AO=120;DP=89 С т chr1 6529188 . С т 11263.97 PASS AC=2;AF=0.500;AN=4;AO=606;D1 AC=2;AF=0.500;AN=4;AO=331;DP=22 chr1 6529443 . 5283.78 PASS А G chr1 6529747 . А AT 1631.35 PASS AC=1;AF=0.500;AN=2;AO=10;DP=478

gen/.bgen

21 · rs240444 · 11002011 · T · C · 1 · 0 · 0 · 1 · 0 · 0 · 1 · 0 · 0 · 1 · 0 · 0	^
21.rs2260895.14642464 A.G.O.O.I.I.O.O.O.O.I.I.O.O.O.O.I.O.I.O.I	
21 · rs2821796 · 14649798 · C · A · 0 · 0 · 1 · 1 · 0 · 0 · 0 · 1 · 0 · 0	
21.rs2742182.14665973.C.T.0.1.0.1.0.1.0.1.0.1.0.0.0.0.1.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0	
21.rrs3119488.14693233.4.G.0.1.0.0.0.1.1.0.0.0.1.1.0.0.0.1.1.0.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.0.0.0.0.0.0.0.0.0.0.0.0.0.0.0.0.0.0	
21.rs2259403.14693381.G.A.0.1.0.0.0.1.0.0.0.1.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.0.1.0.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0	
21.rs3126387.14736458.C.A.1.0.0.0.1.0.0.0.000000	
21.rs2821847.14767569.T.C.0.1.0.0.0.1.0.1.0.0.1.0.1.0.1.0.0.1.0.1.0.0.1.0.1.0.0.1.0.1.0.0.1.0.1.0.0.1.0.1.0.0.1.0.1.0.0.0.1.0.0.0.1.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0	
21 · rs7410107 · 14775729 · G · A · 0 · 1 · 0 · 0 · 0 · 1 · 0 · 1 · 0 · 1 · 0 · 1 · 0 · 0	
21.rs2747351.14780336 G A 0 1 0 1 0 1 0 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 0 1 0 0 0 0 1 0 0 0 0 1 0 0 0 0 1 0 0 0 0 1 0 0 0 0 1 0 0 0 0 1 0 0 0 0 0 1 0	
21.rs3132414.14790733 A.T.1.0.0.0.1.0.0.0.0.000000	
21.rs59682.14800853.C.T.0.1.0.1.0.1.0.0.0.1.0.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.1.0.0.1.0.1.0.0.1.0.0.1.0.0.1.0.0.0.1.0.0.1.0.0.0.1.0	
21.rs3119486.14818933.T.C.1.0.0.0.1.0.0.0.0.0.0.000000	
21. rs468843.14820408.A.G.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0	
21 · rs3126409 · 14832310 · G · A · 1 · 0 · 0 · 1 · 0 · 0 · 1 · 0 · 0 · 1 · 0 · 0	
21 · rs469388 · 1483885 · A·G · 1 · 0 · 0 · 1 · 0 · 0 · 1 · 0 · 0 · 1 · 0 · 0	
21.rs10482849.14916055.T.A.1.0.0.0.1.0	
21 · rs9984128 · 15203801 · G · A · 1 · 0 · 0 · 1 · 0 · 0 · 1 · 0 · 0 · 1 · 0 · 0	
21.rs9984294.15206017.T.A.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0	
21.rs3115511.15214708.G.A.1.0.0.0.1.0.1.0.1.0.1.0.0.0.0.1.0.1.0	
21 · rs400304 · 15280652 · T · C · 1 · 0 · 0 · 1 · 0 · 0 · 1 · 0 · 0 · 1 · 0 · 0	
21 · rs7280714 · 15318381 · C · G · 1 · 0 · 0 · 0 · 1 · 0 · 1 · 0 · 1 · 0 · 0	
21.rs8134986.15331478.C.A.O.O.1.O.1.O.0.O.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.1.O.0.1.O.0.1.O.0.1.O.0.1.O.0.1.O.0.0.1.0.0.0.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.0.1.O.0.0.0.1.O.0.0.0.1.O.0.0.0.0	
21 · rs447479 · 15336419 · T · G · 1 · 0 · 0 · 1 · 0 · 0 · 1 · 0 · 0 · 1 · 0 · 0	
21 · rs447807 · 15337085 · T · C · 0 · 0 · 1 · 0 · 1 · 0 · 0 · 0 · 1 · 0 · 0	
21 · rs1297083 · 15372509 · C · T · 0 · 0 · 1 · 0 · 1 · 0 · 0 · 0 · 1 · 0 · 0	
21 · rs28222391 · 15412399 · G · A · 0 · 1 · 0 · 0 · 1 · 0 · 0 · 1 · 0 · 0	
21.7*48129930.15412893.0.15.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0.0.0.1.0	×

Some definitions

- A *locus L* is a well-defined position along a chromosome.
- single nucleotide variant (SNV) variation in a single nucleotide
- In the population, different variants can exist at a locus. These variants are called *alleles*. A locus *L* is described by the set of alleles {A,G}, or {A,G,T}, {G,ACTAA}, {AGT,A}. One allele is fixed as the "reference allele".
- A *genotype* at a locus consists of a pair of alleles, one inherited from the father and one from the mother. A person is said to be *homozygous* at locus *L* if both alleles are the same (i.e., *A/A or G/G*) and is said to be *heterozygous* if the alleles are different (i.e., *A/G*).

Genetic data representation



What is linkage?

Gregor Mendel

- Moravian Augustinian monk who founded the modern science of genetics
- Derived the principles of Mendelian inheritance (Mendel's laws)
- Used only dichotomous traits and assumed full penetrance, one disease one locus



Mendel's first law

 Mendel's First Law (Segregation): One allele of each parent is randomly and independently selected, with probability 1/2, for transmission to the offspring; the alleles unite randomly to form the offspring's genotype.

Mendel first law





Mendel's first law

				Offspring's				
Father's	Mother's	Genotype						
Genotype	Genotype	dd	dD	DD				
dd	dd	1	0	0				
dd	dD	$\frac{1}{2}$	$\frac{1}{2}$	0				
dd	DD	Õ	$\tilde{1}$	0				
dD	dd	$\frac{1}{2}$	$\frac{1}{2}$	0				
dD	dD	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$				
dD	DD	$\overset{4}{0}$	$\frac{1}{2}$	$\frac{1}{2}$				
DD	dd	0	$\frac{1}{1}$	$\tilde{0}$				
DD	dD	0	$\frac{1}{2}$	$\frac{1}{2}$				
DD	DD	0	$\tilde{0}$	$\frac{1}{1}$				

Table 2.1: Distribution of offspring's genotype conditional upon parental genotypes

Mendel's second law

- Question: are the genes underlying two traits transmitted independently?
- Mendel's Second Law (Independent Assortment): The alleles underlying two or more different traits are transmitted to offspring independently of each other; the transmission of each trait separately follows the first law of segregation.
- Linkage lack of independent transmission, when the second law does not hold.

Mendelian diseases versus complex diseases

- Monogenic aka Mendelian diseases
 - Caused by one or few genes (highly penetrant)
 - Rare disease
 - Example: sickle-cell anemia, cystic fibrosis
- Complex disorder
 - Caused by several factors, including multiple genes and environment
 - Variants can increase/decrease risk to a particular disease (variable penetrance)
 - Example: COPD, Alzheimer's disease, Schizophrenia

Penetrance functions

- P(Y=1|G) penetrance function
- f_{DD} :=P(1|DD); f_{Dd} :=P(1|Dd); f_{dd} :=P(1|dd)
- Mendel's penetrance functions (complete penetrance):
 - Dominant mode of inheritance: $f_{DD} = f_{Dd} = 1$; $f_{dd} = 0$
 - Recessive mode of inheritance: $f_{DD}=1$; $f_{Dd}=f_{dd}=0$

Realistic penetrance for complex diseases

- P(Y|G)<1
- For dominant model: $f_{DD} = f_{Dd} < 1$; $f_{dd} > 0$
- Other models
 - Recessive
 - Additive

Recombination

Genetic recombination happens as a result of the separation of genes that occurs during gamete formation in meiosis



Each gamete gets one copy of the chromosome, each with a unique combination of alleles.

Recombination rearranges chromosomes, generating new allele combinations. While just one homologous chromosome pair is shown above, the same process happens for all of them.

• θ – recombination fraction between two loci

• $\theta = \frac{1}{2}$ under independent assortment

Courtesy of https://learn.genetics.utah.edu/content/pigeons/geneticlinkage/

What is Linkage

- Genetic linkage is a physical concept.
- Two loci are linked, if they are physically near to each other on the chromosome and hence unlikely to be separated during recombination. Θ ≠1/2.
- Which markers are perfectly unlinked?

Linkage example

Not Linked

Gene 1 and Gene 2 are far apart on the same chromosome.



Linked



Courtesy of https://learn.genetics.utah.edu/content/pigeons/geneticlinkage/

No new allele combinations

Linkage analysis (mapping)



• Approaches that use joint transmission of affection status and alleles at the observed marker locus to localize the DSL (disease sucseptability locus)

Linkage of early onset AD to APP gene



Advantages and disadvantages

- Few genome-wide markers required
- Requires being able to infer relationship between variant at DSL and disease trait. Easy with Mendelian diseases and known mode of inheritance
- - may implicate large regions with many genes
- - Requires family data, which might be difficult to collect

Linkage methods

- Parametric linkage analysis
 - H0: $\theta = \frac{1}{2}$, i.e. marker and DSL are unlinked
 - H1: θ <1/2

• LOD score:
$$Z(\theta) = \log_{10} \frac{\prod_{i=1}^{n} L_i(\theta + y_i)}{\prod_{i=1}^{n} L_i(\theta = 1/2 \mid y_i)} = \sum_{i=1}^{n} \log_{10} \frac{L_i(\theta \mid y_i)}{L_i(\theta = 1/2 \mid y_i)}$$

 $\prod L_{i}(A \mid u_{i})$

- Requires specification
 - Disease model
 - Marker parameters (allele frequencies)
- Nonparametric linkage analysis
 - Based on IBS and IBD (allele sharing)

What is association

What is association

- Allelic association (linkage disequilibrium)
 - Allelic values at DSL and marker are associated in the population
 - Population-based concept

- Association (broad term)
 - Testing statistical independence between a variant and a trait of interest.
 - Can use population-based (unrelated) or family-based designs
 - GWAS genome-wide association study
 - WGAS whole genome association study

LD

- Non-random sharing of combinations of genetic variants within a population due to the history of recombination, mutation, and selection in a genomic region
- Linkage disequilibrium population concept, exists when two loci are close (up to 500kb)
- There are regions of long range LD conserved genetic regions
- 3 measures of LD



LD

• Consider 2 diallelic loci {A,a} and {B,b}. If the occurrence of alleles A and B are independent events, then $pAB = pA \cdot pB$. And the alleles are said to be in linkage equilibrium.

under Linkage Equilibrium							
	B Locus						
A Locus	В	b	Total				
A	$p_{AB} = p_A p_B$	$p_{Ab} = p_A p_b$	p_A				
a	$p_{aB} = p_a p_B$	$p_{ab} = p_a p_b$	p_a				
Column Total	p_B	p_b					

Population Allele Frequencies at Two Loci under Linkage Equilibrium

Coefficient D

- Measures the departure from independence
- $D = p_{AB} p_A p_B$
- D=0 linkage equilibrium
- Highly sensitive to marginal values (allele frequencies)

Coefficient D'

• "Normalized" D

•
$$D' = \begin{cases} \frac{D}{D_{max}}, & \text{if } D < 0\\ \frac{D}{D_{min}}, & \text{if } D > 0 \end{cases}$$

- Dmax=min($p_A p_b, p_a p_B$)
- Dmin= $-max(p_A p_B, p_a p_b)$
- Ranges from 0 to 1

$R \text{ or } R^2$

- Correlation coefficient between 2 markers
- $r = p_{AB} pApB / \sqrt{p_A p_B p_a p_b} = D / \sqrt{p_A p_B p_a p_b}$

		B Locus		Row	
	A Locus	В	b	Total	
	A	43	27	70	
	a	2	28	30	
Column Total		45	55	100	

$$D = (43 - 70 * 45/100)/100 = 0.115$$

$$D_{\text{max}} = \min(70 * 55, 30 * 45)/10,000 = 0.135$$

$$D' = 0.115/0.135 = 0.8581$$

$$r = 0.115/\sqrt{0.7 * 0.3 * 0.55 * 0.45} = .5044$$

Laird and Lange, 2010 "The fundamentals of modern statistical genetics"



Allele frequencies

• Allele frequency – proportion of chromosomes in population/dataset having the allele of interest.

Example:

• {A,a}, n_{AA}, n_{Aa}, n_{aa}

• $\hat{p} = (n_{Aa} + 2n_{AA})/n$

• $\hat{q} = 1 - \hat{p}$

AAAaaa \hat{p} \hat{q} Dataset 1904204900.30.7(n=1000) \cdot \cdot \cdot \cdot \cdot

Strength of genetic effect versus allele frequency



HWE (Hardy-Weinberg equilibrium)

- A population is said to be in HWE if the genotypes satisfy:
- $P(AA) = p^2;$
- P(Aa) = 2pq;
- $P(aa) = q^2;$
- Implication: Allele frequencies does not change in a population from generation to generation.
- Useful to identify genotyping errors, population structure





By Johnuniq - Own work, CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=6045237

Using LD in association studies



Test for association between a phenotype Y and a marker G



Advantages and disadvantages of association studies

- Uses LD to localize DSL
- Better coverage through large number of markers
- Appropriate for complex traits
- A large number of individuals required -> different sources of bias
- Large number of markers -> multiple testing problem

Association methods (population-based)

• Linear and logistic regression models

 $Y \sim \beta G + \gamma Z + \varepsilon$

- Covariates include principal components to account for population stratification
- Mixed models (LMM and GLMM)
 - Allows to increase the sample size by including related individuals
 - Also efficiently accounts for population structure

Association methods (family-based)

- Test for both linkage and association
- Robust to population substructure: admixture, stratification, failure of HWE
- Requires genotyped families (parent-child, or siblings)
- TDT test for trio design (affected offspring)
- FBAT generalization to general phenotypes, general pedigrees, missing parental genotypes, and multiple variants (Lake and Laird 2001, Laird and Lange 2006,...)

Linkage versus association

• Linkage is a physical concept: The two loci are "close' together on the same chromosome. There is hardly any recombination between disease locus and marker locus

• Association is a statistical / population-based concept: A particular marker allele tends to be present with disease allele in a population and is associated with the phenotype.

Next lecture

- Genome-wide association studies
 - Population-based association studies
 - Family-based association studies
 - Common genetic variant methods
 - Rare genetic variant methods
 - Problems (QC, POPSTRAT)