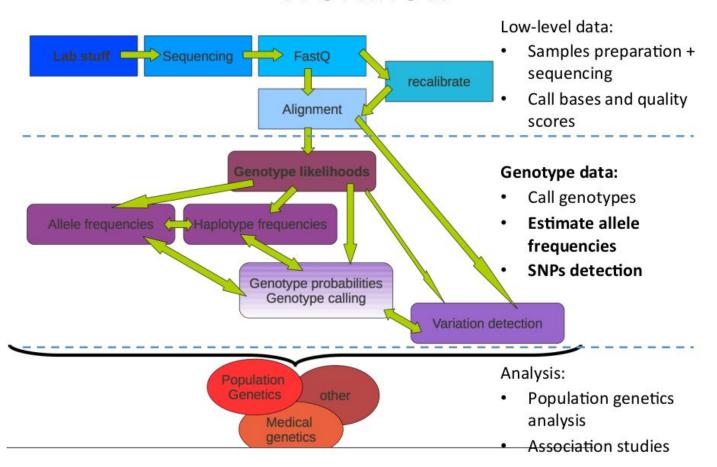
GENOME ANALYSIS & BIOINFORMATICS

<u>Week 5:</u>

Genetic variation and variant calling from NGS data

Workflow



Genome (FASTA)

>ARPM2ref|NC_000001.10|:2938046-2939467 Homo sapiens chromosome 1, GRCh37 primary reference assembly

TGGAAGAGGCCTCAGCCAGGCCAGCCACCTGGAGGGAGAGCAGACCTGCGGCTGAGGATGCAGGGCTCCCGGGCACGGTGCTAGCCCTGCGGCCCCGAGAGCTGTGGGAAGAGCTGTGGGATCCCCTATTGCACACAAAGCGGCCCTGGAGGGCTGGTCTTTATTTTGATGAGGCTGAGAAGGGAAGGCTGCGGGCATGTTTAATCCGCACGCTCTTAGACTCCCCGGCTGTGATTTTTGACAATGGCTCGGGGTTCTGCAAAGCGGGCCTGTCTGGGGGAGTTTTGGACCCCGGCACATGGTCAGCTCCATCGTGGGGCACCTGAAATTCCAGGCTCCCTCAG



Reads (FASTQ)

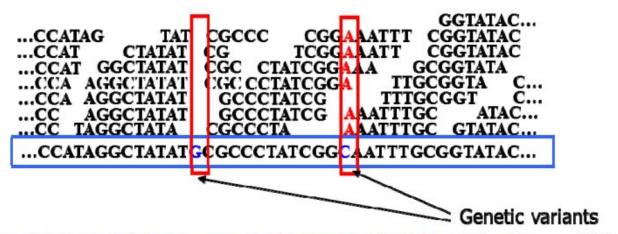
CCAATGATTTTTTTCCGTGTTTCAGAATACGGTTAA
+SRR038845.41 HWI-EAS038:6:1:0:1474 length=36
BCCBA@BB@BBBBAB@B9B@=BABA@A:@693:@B=
@SRR038845.53 HWI-EAS038:6:1:1:360 length=36
GTTCAAAAAGAACTAAATTGTGTCAATAGAAAACTC
+SRR038845.53 HWI-EAS038:6:1:1:360 length=36

Mapped Reads (mpileup, BAM)

Variants (VCF)

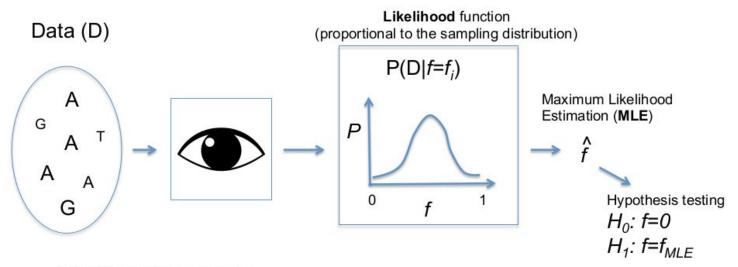
##filef	ormat=VC	Fv4.1								П
##fileD	ate=2014	0930								- 1
##sourc	e=23andm	e2vcf.pl	https	://githul	b.com/arr	ogantrobo	t/23and	lme2vcf		- 1
				hg19_re						- 1
						on="Genot	vpe">			- 1
#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	GENOTY	PE
chr1	82154	rs44772	12	а					GT	0
/0										
chr1	752566	rs30943	15	q	Α				GT	1
/1										
chr1	752721	rs31319	72	Α	G				GT	1
/1										
chr1	798959	rs11240	777	q					GT	0
/0										
chr1	800007	rs66810	49	Т	C				GT	1
/1										





- Counting high-confident, non-reference allele (i.e. Quality >= 20)
 - Freq <20% or > 80%: homozygous genotype
 - Otherwise: heterozygous
- Works well for "deeply sequenced regions" (DSR), i.e. depth > 25x
 - But suffer from under-calling of heterozygous genotypes for low-coverage regions
 - And can't give an objective measurement for reliability

Statistical inference (1)



Likelihood approach:

- All the information on the parameter is in the likelihood function (we use all the data!).
- More data leads to less bias and less variance.
- Suitable for hypothesis testing.

Calculating genotype likelihoods

$$P(X|G=bh) = \prod_{i=1}^{r} \left(\frac{L_b^{(i)}}{2} + \frac{L_h^{(i)}}{2}\right) \quad b, h \in \{A, C, G, T\}$$

Example:

Chrom6 342

A T T	Individual 1
Т	Individual 2

Individual

Calculating genotype likelihoods

$$P(X|G=bh) = \prod_{i=1}^{r} \left(\frac{L_b^{(i)}}{2} + \frac{L_h^{(i)}}{2}\right) \quad b, h \in \{A, C, G, T\}$$

Example:

Chrom6 342 ATTT

$$\begin{split} P(X|G=AC) &= (\frac{L_A^{(1)}}{2} + \frac{L_C^{(1)}}{2}) * (\frac{L_A^{(2)}}{2} + \frac{L_C^{(2)}}{2}) * (\frac{L_A^{(3)}}{2} + \frac{L_C^{(3)}}{2}) * (\frac{L_A^{(4)}}{2} + \frac{L_C^{(4)}}{2}) \\ &= (\frac{1 - \mathcal{E}}{2} + \frac{\mathcal{E}}{6}) * \frac{\mathcal{E}}{3} * \frac{\mathcal{E}}{3} * \frac{\mathcal{E}}{3} \end{split}$$

Genotype likelihoods

Genotype	Likelihood (log10)
AA	-7.44
AC	-7.74
AG	-7.74
AT	-1.22
СС	-9.91
CG	-9.91
СТ	-3.38
GG	-9.91
GT	-3.38
π	-2.49

ATTT

$$\varepsilon = 0.01$$

Likelihood Ratio:

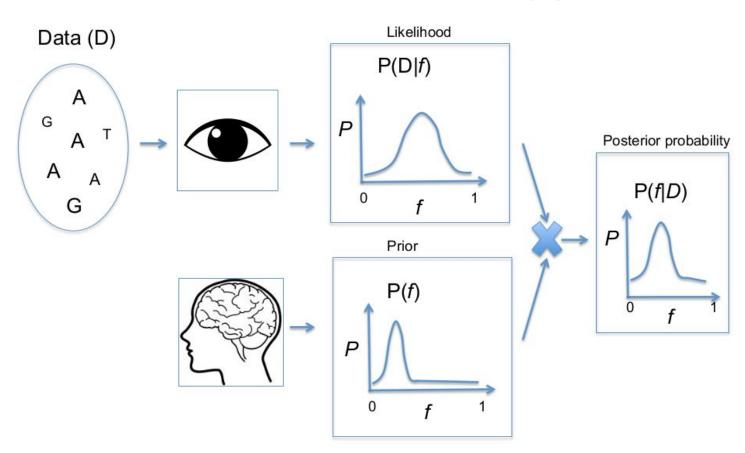
$$\log_{10} \frac{L_{G(1)}}{L_{G(2)}} > t$$

$$t = 1$$

The most likely genotype is at least **10 times** more likely than the second most likely one

(in our example t=1.27)

Statistical inference (2)



Bayesian inference

$$P(\theta|X) = \frac{P(X|\theta)P(\theta)}{P(X)} = \frac{P(X|\theta)P(\theta)}{\sum_{\theta} P(X|\theta)P(\theta)}$$

$$P(X|\theta) \longleftarrow \text{ Likelihood of } \theta$$

$$P(\theta) \longleftarrow \text{ Prior probability distribution of } \theta$$

$$P(\theta|X) \longleftarrow \text{ Posterior probability distribution of } \theta$$

- Parameter is not fixed (point estimate) but rather has a probability distribution
- We update our "belief" on the parameter after performing the experiment
- As P(f|D) is a proper probability distribution, we can easily derive credible intervals

Genotype	Likelihood (log10)	Prior	Posterior probability	
AA	-7.44	1/10	~0	
AC	-7.74	1/10	~ 0	Ge
AG	-7.74	1/10	~ 0	
AT	-1.22	1/10	0.94	
CC	-9.91	1/10	~ 0	
CG	-9.91	1/10	~ 0	
СТ	-3.38	1/10	0.006	
GG	-9.91	1/10	~ 0	
GT	-3.38	1/10	0.006	
π	-2.49	1/10	0.05	

Genotype posterior probabilities

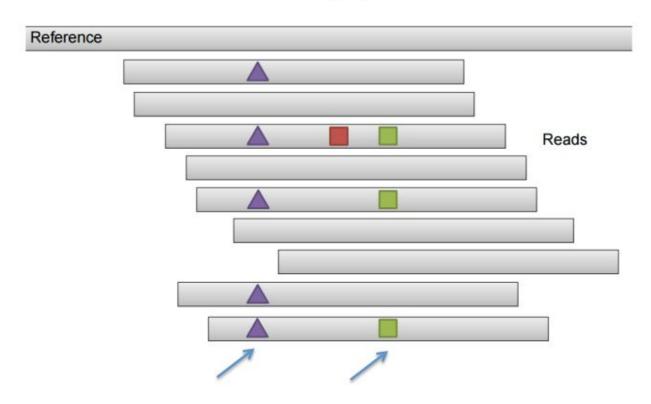
But **only** call the genotype if the largest probability is above a threshold (e.g. > 0.95)

Gei	notype	Likelihood (log10)	Prior	Posterior probability	
AA		-7.44	0.16	~0	
AC		-7.74	0	0	Ge
AG		-7.74	0	0	
AT		-1.22	0.48	0.96	
CC		-9.91	0	0	
CG		-9.91	0	0	
СТ		-3.38	0	0	
GG		-9.91	0	0	
GT		-3.38	0	0	
П		-2.49	0.36	0.38	

Genotype posterior probabilities

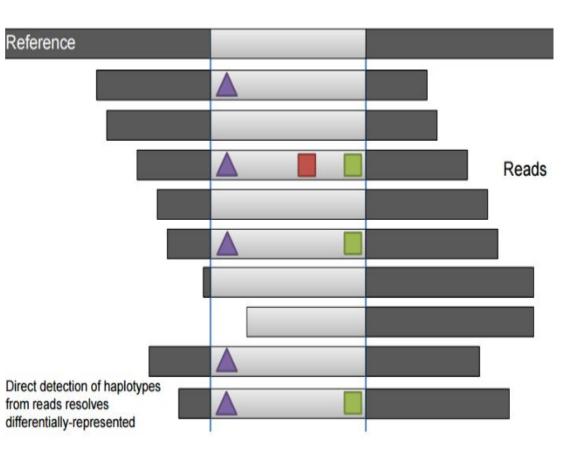
P(A) = f=0.6

SNP calling procedures



We completely rely on how reads have been mapped

Haplotype-based caller



- 1. Define active regions
- 2. Determine haplotypes by assembly of active regions
- 3. Determine LH of haplotypes given read data
- 4. Assign sample genotypes
- 5. Determine allele frequency at all sites
- 6. Determine presence of SNPs

Estimating allele frequencies

Individual	True genotype	Reads allele A	Reads allele G
1	AA	7	0
2	AA	25	1
3	AG	5	3
4	AG	4	4
5	GG	0	2
6	GG	0	4
Tot.		41	14

Maximum Likelihood (ML) estimator (Kim et al. 2011)

$$L = \prod_{i=1}^{N} p(D_i \mid f)$$
 Genotype likelihoods

$$p(D_i | f) = \sum_{g \in \{0,1,2\}} p(D | G = g) p(G = g | f)$$

If we assume HWE:
$$p(G = AA \mid f) = f^2$$

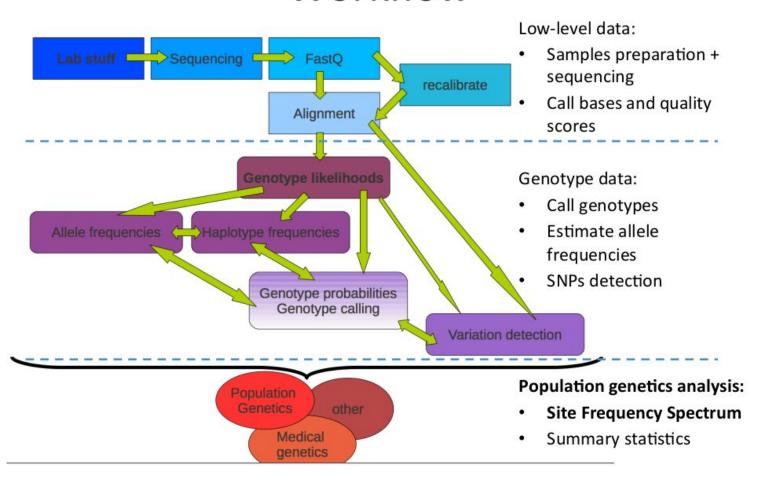
 $p(G = AG \mid f) = 2 f(1 - f)$

$$p(G = GG | f) = (1-f)^2$$

$$\hat{f} = \operatorname{arg\,max}_{p} \prod_{i=1}^{N} p(D_{i} \mid f)$$

$$\hat{f} = 0.46$$

Workflow



Possible measures of genetic variation

Total raw data

the numbers of the different kinds of polymorphic sites and their distribution along the sequence/genome

Pros

most detailed information available

Cons

computationally intensive not practical often times impossible

SFS is the most information rich summary statistics

If the mutation rate per site is constant along the sequence, the number of segregating sites should, on average, be proportional to the length of the sequences. S should be proportional mutation rate.

S should increase with the sample size

S & Pi: reduces all information contained in the data into a single number!

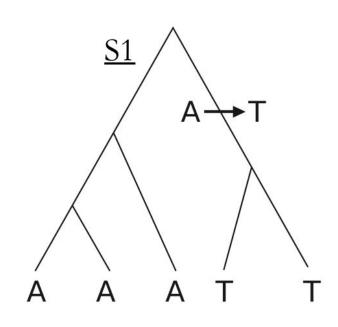
SFS: an intermediate measure, between the total data and the extreme summaries.

Both S & Pi can be derived from the SFS

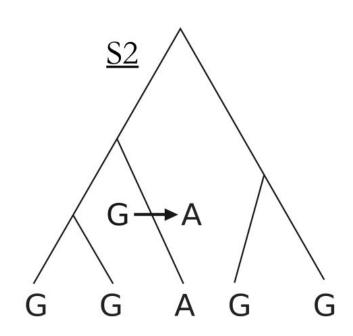
What is the SFS?

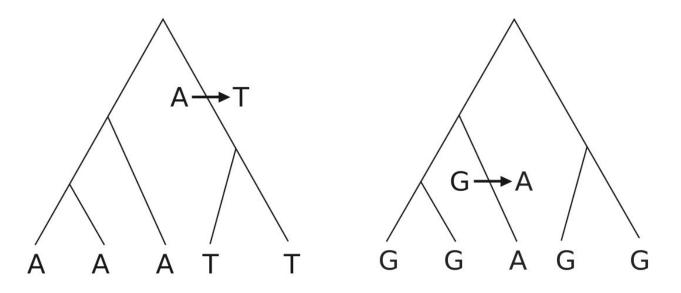
Is the number of counts of each possible kind of polymorphic site (distinguishing between sites that partition the sample in different ways), based on the correspondence between genealogical structure and observable DNA data under the infinite-sites mutation model.

	<u>S1</u>	<u>S2</u>	
Seq 1 _	Ą	Ģ	
Seq 2 _	Ţ	Ģ	
Seq 3 _	Ą	Ģ	
Seq 4 _	Ą	Ą	
Seq 5 _	Ţ	Ģ	

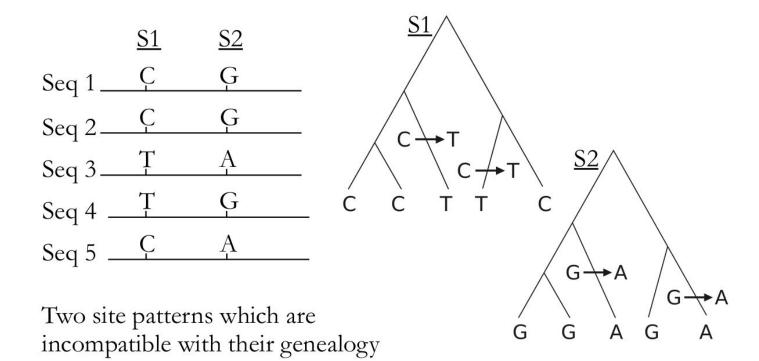


	<u>S1</u>	<u>S2</u>	
Seq 1 _	Ą	Ģ	
Seq 2 _	Ţ	Ģ	
Seq 3 _	Ą	Ģ	
Seq 4 _	Ą	Ą	
Seq 5 _	Ţ	Ģ	

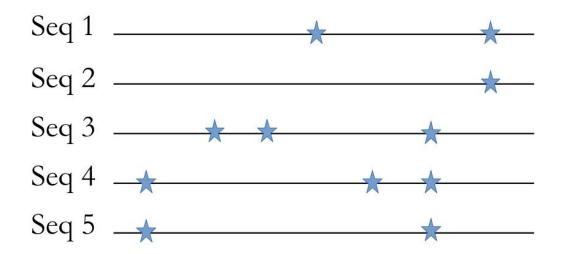




Under the infinite sites model there is a branch in the genealogy that partitions the sample identically to the pattern at that site. If this is the case we say that the site pattern is compatible with the genealogy

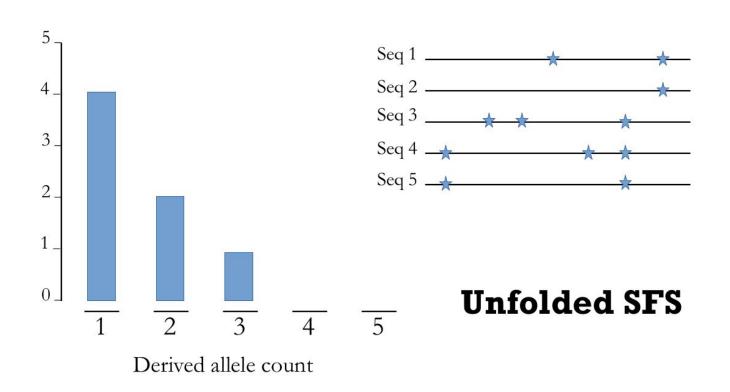


Calculating the SFS

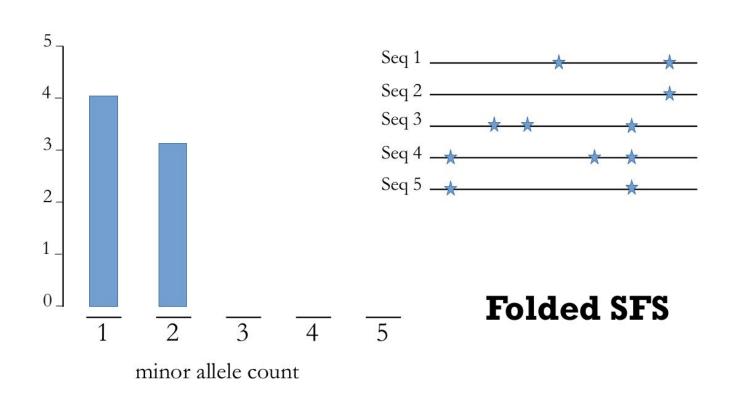


 N_i = The number of sites " ξi " at which the mutant base is present in i copies and the ancestral base is in n - i copies

Calculating the SFS



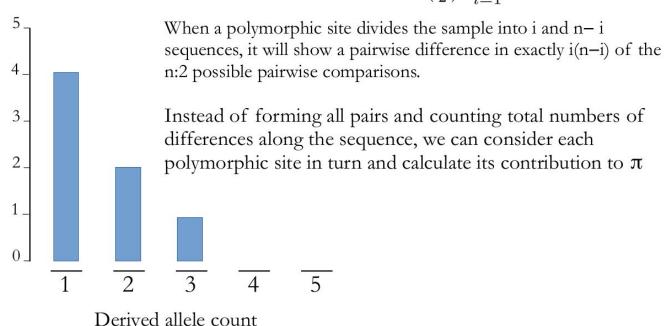
Calculating the SFS



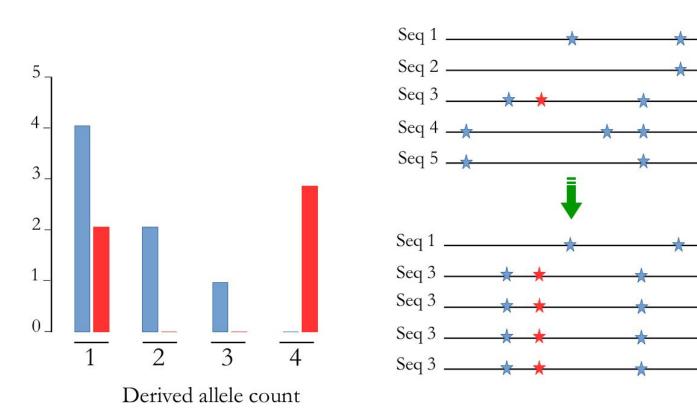
Deriving summary stats from the SFS

For S both simply count the number of poly sites

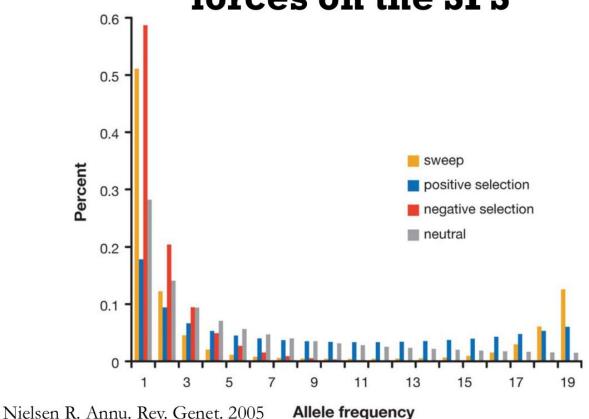
$$\pi = \frac{1}{\binom{n}{2}} \sum_{i=1}^{[n/2]} i(n-i)\eta_i.$$



Positive selection & SFS



Influence of evolutionary forces on the SFS



Effect of errors on the SFS

