

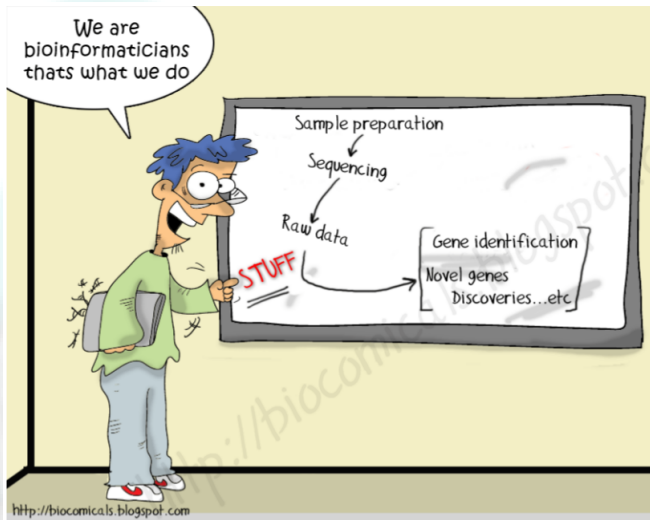
NGS data: Beauty and the Beast

How to infer population genetics parameters from messy sequencing data

Matteo Fumagalli

12th September 2017

Goal of the day

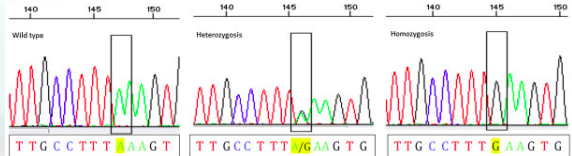
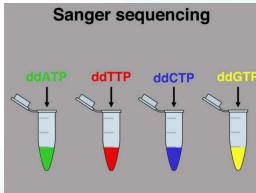


Presentation outline

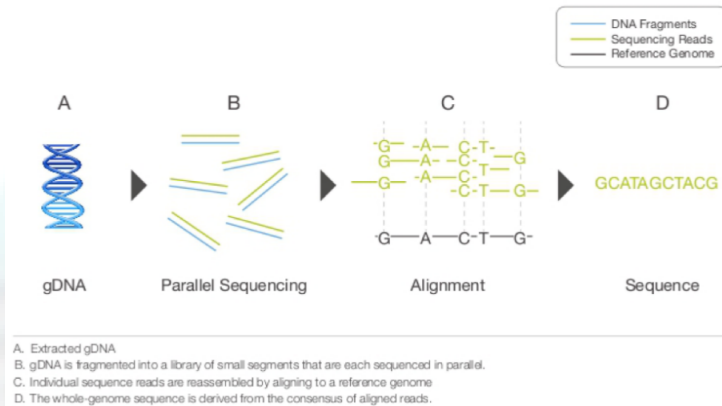
- 1 Introduction
- 2 Genotype likelihoods
- 3 Genotype calling, really?
Low-depth (unknown genotypes)
- 4 Allele frequencies
Low-depth (allele frequencies)
- 5 Experimental design

Sanger sequencing

aka first/former generation sequencing



Next Generation Sequencing



Low-level data

FASTQ

```
a'X_\Va\J'KaYJHG^]b\a^BBBBBBBBBBBBBB <-- quality score
@FC42BF1AAXX:6:1:5:732#0/1 <-- read ID
TGATTCTCTCGATATCCAGTCCTTAGTGNCATAGN <-- read (bases)
+
a^_aaaa'aa'_aaa_aaa'__'_'VBBBBBBBBB
@FC42BF1AAXX:6:1:5:492#0/1
AACAGTGGGAGGCTGCAGCAGGAGGATTNCTGAAN
+
ababb_abbbaab^'aaTaabbbaBBBBBBBBB
@FC42BF1AAXX:6:1:5:480#0/1
ACCTCCTCAGAGTTCTCGAGCTCGAGAANTCTGGN
```

Quality scores

Qscore

- The ASCII values can be interpreted as a probability
- A Q20 (ASCII 'T') score is probability of 1%
- The score is the probability, P , that the base is incorrect
-

$$Q_{\text{score}} = -10 \log_{10}(P)$$

-

$$P = 10^{\frac{-Q}{10}}$$

Quality scores

The qscores are encoded as ASCII characters, and are shifted by +33 (now the standard) or +64.

Phred Quality Score	Probability of error	Base call accuracy
0 ... 9	1 ... 0.13	!"#\$%&'()*
10 ... 19	0.1 ... 0.013	+,-./01234
20 ... 29	0.01 ... 0.0013	56789;:<=>
30 ... 39	0.001 ... 0.00013	?@ABCDEFGH
40 ... 49	0.0001 ... 0.000014	IJKLMNOPQR

Quality scores

Example

From a fastq file we observe an 'A' with a qscore encoded as '7'. What is the probability of being 'A'? What is the probability of being 'G'?

1. We find that the corresponding ASCII value of '7' is ?
(hint: <http://www.asciitable.com>)

Quality scores

Example

From a fastq file we observe an 'A' with a qscore encoded as '7'. What is the probability of being 'A'? What is the probability of being 'G'?

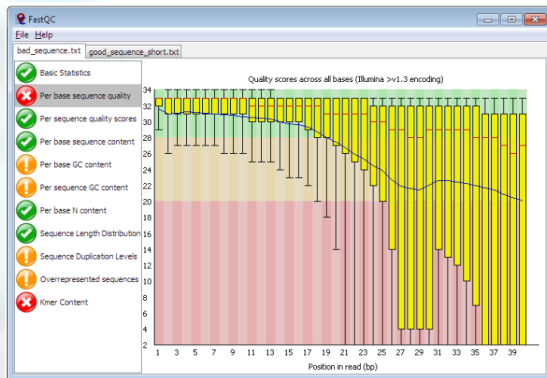
1. We find that the corresponding ASCII value of '7' is ?
(hint: <http://www.asciitable.com>)
2. We subtract 33 to get a value of ?. This is our qscore.
3. The probability of 'A' being incorrect is ? (hint: $p = 10^{\frac{-Q}{10}}$)
4. The probability of 'A' being correct is ?
5. The probability of being 'G' (or 'C' or 'T') is ?

FastQ files

```
@HWI-ST397_0000:2:1:2248:2126#CTTGTA/1
TTGGATCTGAAAGATGAATGTGAGAGACACAATCCAAGTCATCTCTC
ATG
+HWI-ST397_0000:2:1:2248:2126#CTTGTA/1
eeee\dZddaddddddeeeeeeeedaed_ec_ab_\NSRNRcdddc[_c^d
```

- sequencer
- flowcell
- lane/cell/tile
- position within the tile
- barcode id for pooling/multiplexing
- pair
- ...

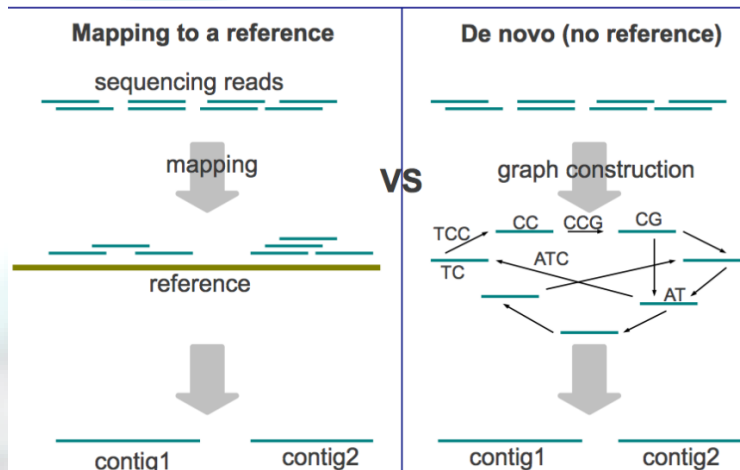
Quality check of fastq files: fastQC



- distribution of qscores over read
- overrepresented kmers
- ...

<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

Alignment of reads



Mapping to a reference

Issues:

- (i) Millions of short reads.
- (ii) Blat/blast is too slow.

But:

Mate-pair gives additional information.

Many aligners exists:

- Soap/soap2 (BGI)
- Maq/bwa (Heng Li)
- Bowtie/bowtie2 (Langmead B)
- Eland, SSAHA2, RMAP, shrimp, zoom, GEM, snap, novoAlign.

Most are based on burrows wheeler transform BTW
(BWA, Bowtie, Soap2, ...)

Mapping to a reference

Many different aligners, what's the difference?

- Memory usage
- Speed
- Gapped (indels)
- Using qscores (bwa pssm)
- Estimating a **mappability score**
- Multiple best hits
- Paired end data
- Output **formats** (SAM,...)

Alignment file

an alignment file includes

reads TTTGTTCTTTCTTTCTCTCTAGTCTTCTT ...

Qscore NVFVN]^['^_]^^U]] '['[_VS[_^Z]_ ...

start position chr4 53351385

multiple best hits 1

Number of mismatch 2

sequence strand -

read quality* V

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

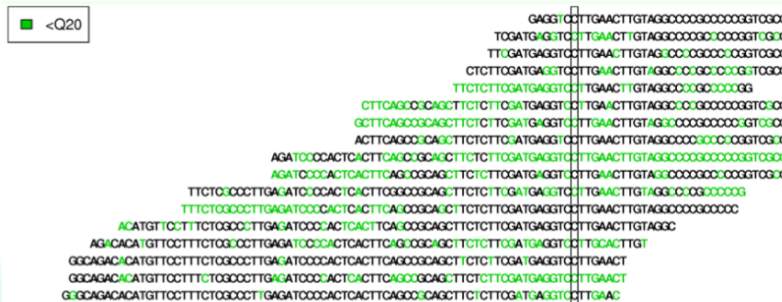
TGGAAGAGGCTCAGCAGGCCAGGCCACCTGGAGGGAGAGCAGACCTGCGGCTGAGGATGCAGGGCTCC
 CGGGCACGGTGCTAGCCCTGCCCTTGAGACACCCGAGAGCTGTGGGAAGAGCTGTGGGATCCCCATTATGC
 ATCACAAAGCGGCCCTTGAGGGCTGGTCTTTATTTTGATGAGGCTGAGAAGGGAAGGCTGCGGGCATGTT
 TAATCCGACAGCTTTAGACTCCCGGGCTGTGATTTTGCAGTACGCTGGGGTCTGCAAAGCGGGCTG
 TCTGGGGAGTTTGGAGCCCGGACATGGTCACTCCATCTGGGGGACCTGAAGAAATCAGGCTCCCTCAG

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

[illegible]

```
##fileformat=VCFv4.0
##fileDate=20140930
##source=23andme2vcf.pl https://github.com/arrogantrobot/23andme2vcf
##reference=file:///23andme_v3 hg19_ref.txt.gz
##FORMAT=
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT GENOTYPE
chr1 82154 rs4477212 a . . . . . GT
/0
chr1 752566 rs3094315 g A . . . . . GT
/1
chr1 752721 rs3131972 A G . . . . . GT
/1
chr1 798959 rs11240777 g . . . . . GT
/0
chr1 800087 rs6681049 T C . . . . . GT
```

Our data: mapped reads with quality scores



- Coverage: fraction of the genome with data
- Depth: number of reads mapped to a position
- Counts: number of different alleles mapped to a position
- Effective Base Depth: similar to the counts, but weighing for qscores and mapping quality

Challenges

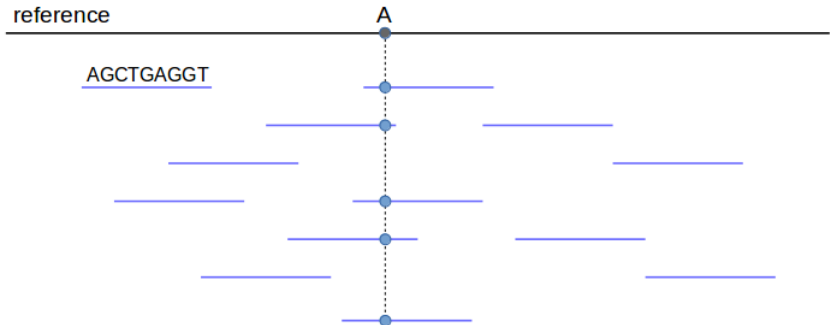


- Variable and low depth
- High sequencing and mapping errors



Quality control filters

The data



- is a **nucleotide**/base/allele with a certain **quality** score

Genotype likelihoods

Likelihood

$$P(D|G = \{A_1, A_2, \dots, A_n\})$$

with

$A_i \in \{A, C, G, T\}$ and n being the ploidy

How many genotypes likelihoods do we need to calculate for each individual at each site?

Base quality in Phred scale



Calculating genotype likelihoods

Likelihood function

$$P(D|G = \{A_1, A_2, \dots, A_N\}) = \prod_{i=1}^R \sum_{j=1}^N \frac{L_{A_j,i}}{N}$$

- $L_{A_j,i} = P(D|A_G = A_j)$
- $A_i \in \{A, C, G, T\}$
- R is the depth (nr. of reads)
- N is the ploidy (nr. of chromosomes)

Example:

A
A
A
G

with all quality scores equal to 20 (in phred score)

$P(D|G = AC) = ?$

Calculating genotype likelihoods

Likelihood function

$$P(D|G = \{A_1, A_2, \dots, A_N\}) = \prod_{i=1}^R \sum_{j=1}^N \frac{L_{A_j,i}}{N}$$

A
A
A
G
& Q=20

$$P(D|G = \{A, C\}) = \dots$$

Calculating genotype likelihoods

Likelihood function

$$P(D|G = \{A_1, A_2, \dots, A_N\}) = \prod_{i=1}^R \sum_{j=1}^N \frac{L_{A_j,i}}{N}$$

A

A

A

G

& Q=20

$N = 2; i = 1; A_1 = A; A_2 = C$

$$P(D|G = \{A, C\}) = \left(\frac{L_{A,1}}{2} + \frac{L_{C,1}}{2}\right) \times \dots$$

What are $L_{A,1}$ and $L_{C,1}$?

Calculating genotype likelihoods

Likelihood function

$$P(D|G = \{A_1, A_2, \dots, A_N\}) = \prod_{i=1}^R \sum_{j=1}^N \frac{L_{A_j,i}}{N}$$

A
A
A
G
& Q=20

$$L_{C,1} = \frac{\epsilon_1}{3}$$

$$L_{A,1} = 1 - \epsilon_1$$

$$P(D|G = \{A, C\}) = \left(\frac{1 - \epsilon_1}{2} + \frac{\epsilon_1}{6}\right) \times \dots$$

Calculating genotype likelihoods

Likelihood function

$$P(D|G = \{A_1, A_2, \dots, A_N\}) = \prod_{i=1}^R \sum_{j=1}^N \frac{L_{A_j,i}}{N}$$

A

A

A

G

& Q=20

$$L_{C,1} = \frac{\epsilon_1}{3}$$

$$L_{A,1} = 1 - \epsilon_1$$

$$P(D|G = \{A, C\}) = \left(\frac{1 - \epsilon_1}{2} + \frac{\epsilon_1}{6}\right) \times \left(\frac{1 - \epsilon_2}{2} + \frac{\epsilon_2}{6}\right) \times \left(\frac{1 - \epsilon_3}{2} + \frac{\epsilon_3}{6}\right) \times \frac{\epsilon_4}{3}$$

What are $\epsilon_1, \epsilon_2, \dots$?

Genotype likelihoods

Calculating genotype likelihoods

Genotype	Likelihood (log10)	
AA	-2.49	
AC	-3.38	
AG	-1.22	A
AT	-3.38	A
CC	-9.91	A
CG	-7.74	G
CT	-9.91	& $\epsilon = 0.01$
GG	-7.44	
GT	-7.74	
TT	-9.91	

Genotype calling

Genotype	Likelihood (log10)
AA	-2.49
AC	-3.38
AG	-1.22
AT	-3.38
CC	-9.91
CG	-7.74
CT	-9.91
GG	-7.44
GT	-7.74
TT	-9.91

AAAG & $\epsilon = 0.01$

What is the genotype here?

Genotype calling

Genotype	Likelihood (log10)
AA	-2.49
AC	-3.38
AG	-1.22
AT	-3.38
CC	-9.91
CG	-7.74
CT	-9.91
GG	-7.44
GT	-7.74
TT	-9.91

AAAG & $\epsilon = 0.01$

What is the genotype?

AG.

Maximum Likelihood

The simplest genotype caller:
choose the genotype with the
highest likelihood.

Major and minor alleles

Likelihood function

$$\log P(D|G = A) = \sum_{i=1}^R \log L_{A_j,i}$$

AAAG & $\epsilon = 0.01$

Allele	log-Likelihood
A	-2.49
C	-3.38
G	-1.22
T	-3.38

We can reduce the genotype space to 3 entries (from 10).

Genotype likelihoods

AAAG & '5555' & A,G alleles

Genotype	log-Likelihood
AA	-5.73
AG	-2.80
GG	-17.12

Examples varying qualities and reads...

Genotype likelihoods - example

AAAG & '5550' & A,G alleles

Genotype	log-Likelihood
----------	----------------

Genotype likelihoods - example

AAAG & '5550' & A,G alleles

Genotype	log-Likelihood
AA	-4.58
AG	-2.81
GG	-17.14

Genotype likelihoods - example

AAAG & '555K' & A,G alleles

Genotype	log-Likelihood
----------	----------------

Genotype likelihoods - example

AAAG & '555K' & A,G alleles

Genotype	log-Likelihood
AA	-10.80
AG	-2.80
GG	-17.11

Genotype likelihoods - example

AAAAAAAAAG & '555555550' & A,G alleles

Genotype	log-Likelihood
----------	----------------

Genotype likelihoods - example

AAAAAAAAAG & '555555550' & A,G alleles

Genotype	log-Likelihood
AA	-4.64
AG	-7.01
GG	-51.37

NGS data uncertainty

Issue

There is a notable amount of statistical uncertainty in assigning individual genotypes depending on the number of reads and their quality. How can we deal with that when doing population genetics analysis (e.g. estimating genetic diversity)?

Solutions

1. let's pretend we don't have such uncertainty

NGS data uncertainty

Issue

There is a notable amount of statistical uncertainty in assigning individual genotypes depending on the number of reads and their quality. How can we deal with that when doing population genetics analysis (e.g. estimating genetic diversity)?

Solutions

1. let's pretend we don't have such uncertainty
2. genotype filtering
3. (a third way)

Genotype likelihood ratio

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

i.e. $t = 1$ meaning that the most likely genotype is 10 times more likely than the second most likely one

Genotype posterior probability

AAAG & $\epsilon = 0.01$ & A,G alleles

Genotype	Likelihood (log)	Prior	Posterior
AA	-5.73	1/3	0.05
AG	-2.80	1/3	0.95
GG	-17.12	1/3	0

Only call genotypes if the largest probability is above a certain threshold (e.g. 0.95).

Pros and cons?

Genotype posterior probability

AAAG & $\epsilon = 0.01$ & A,G alleles

Genotype	Likelihood (log)	Prior	Posterior
AA	-5.73	1/3	0.05
AG	-2.80	1/3	0.95
GG	-17.12	1/3	0

Only call genotypes if the largest probability is above a certain threshold (e.g. 0.95).

Pros and cons?

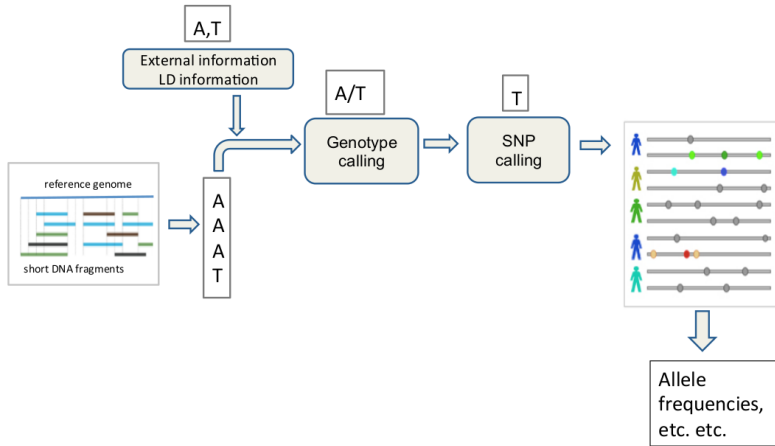
- Yes: genotype are called with higher **confidence**
- No: more **missing** data

Exercise 1

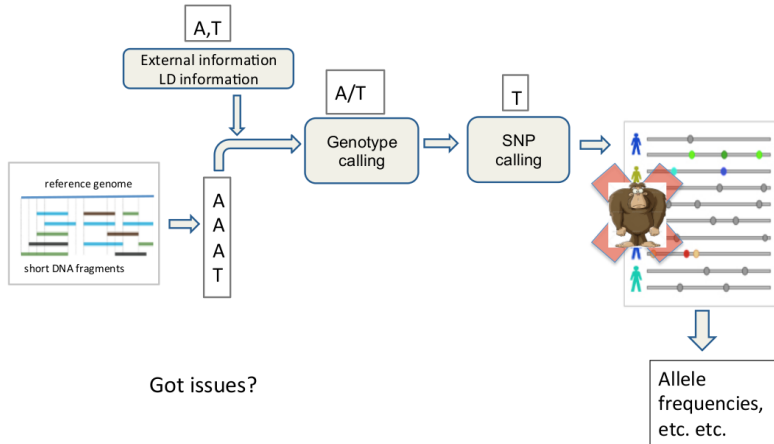
Simulate NGS data and calculate genotype likelihoods and probabilities.

Assess the amount of uncertainty and data missingness.

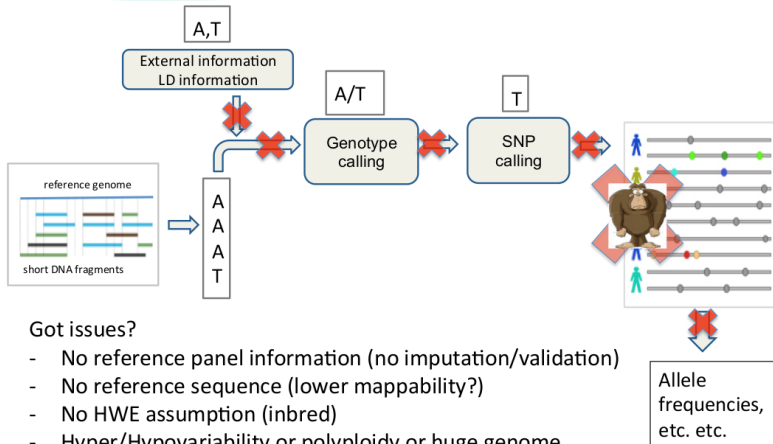
NGS data processing in the **model** world



NGS data processing in the **non-model** world



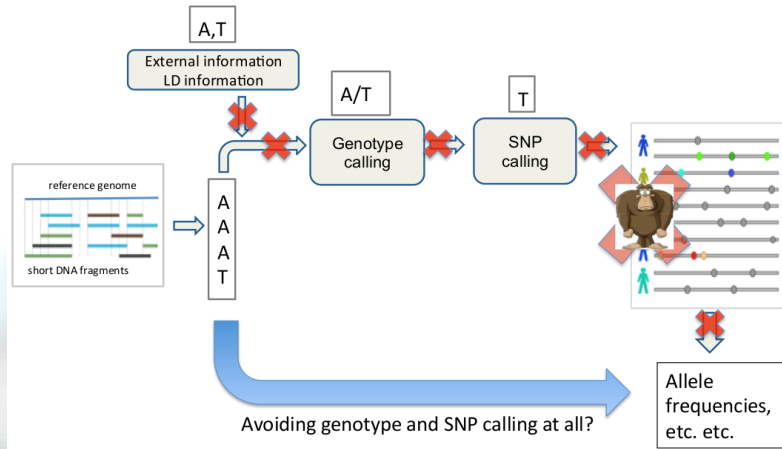
NGS data processing in the **non-model** world



Got issues?

- No reference panel information (no imputation/validation)
- No reference sequence (lower mappability?)
- No HWE assumption (inbred)
- Hyper/Hypovariability or polyploidy or huge genome
- No money (?)
- **Your inferences will be wrong!**

NGS data processing in the **non-model** world



Summary statistics

Aim: estimate the number of heterozygotes (H) from **unknown** genotypes.

Data:

Sample	Data	$P(G = AA D)$	$P(G = AG D)$	$P(G = GG D)$
1	A	0.66	0.33	0.01
2	AAAG	0.14	0.86	0.00
3	AGG	0.00	0.92	0.08
4	GG	0.00	0.20	0.80

with $Q = 15$

Solutions:

1. call genotypes: $H = 2$

Summary statistics

Aim: estimate the number of heterozygotes (H) from **unknown** genotypes.

Data:

Sample	Data	$P(G = AA D)$	$P(G = AG D)$	$P(G = GG D)$
1	A	0.66	0.33	0.01
2	AAAG	0.14	0.86	0.00
3	AGG	0.00	0.92	0.08
4	GG	0.00	0.20	0.80

with $Q = 15$

Solutions:

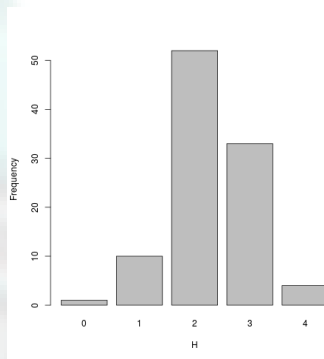
1. call genotypes: $H = 2$
2. sample genotypes: H is defined by a distribution

Summary statistics

Sample	Data	$P(G = AA D)$	$P(G = AG D)$	$P(G = GG D)$
1	A	0.66	0.33	0.01
2	AAAG	0.14	0.86	0.00
3	AGG	0.00	0.92	0.08
4	GG	0.00	0.20	0.80

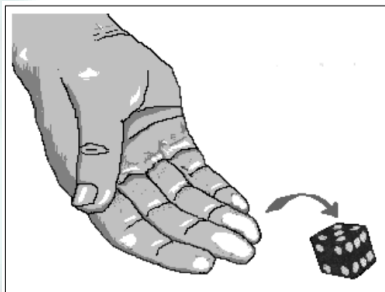
with $Q = 15$

Sampling genotypes:



Genotype calling, really?

Expected value



- What are the possible outcomes of this experiment?
- With what probability?

Expected value

The expected value of a discrete random variable is the probability-weighted average of all possible values

$$E[X|D] = \sum_{i=1}^N x_i p(X = x_i|D)$$

Expected value

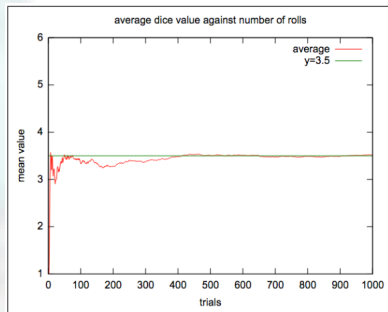
The expected value of a discrete random variable is the probability-weighted average of all possible values

$$E[X|D] = \sum_{i=1}^N x_i p(X = x_i|D)$$

$$E[X|D] = 1 \cdot \frac{1}{6} + 2 \cdot \frac{1}{6} + 3 \cdot \frac{1}{6} + 4 \cdot \frac{1}{6} + 5 \cdot \frac{1}{6} + 6 \cdot \frac{1}{6} = \frac{(1+2+3+4+5+6)}{6} = \frac{21}{6} = 3.5$$

Expected value

It is the average value if you perform the same experiment many times.



Summary statistics

Sample	Data	$P(G = AA D)$	$P(G = AG D)$	$P(G = GG D)$
1	A	0.66	0.33	0.01
2	AAAG	0.14	0.86	0.00
3	AGG	0.00	0.92	0.08
4	GG	0.00	0.20	0.80

with $Q = 15$

1. call genotypes: $H = 2$
2. sample genotypes: H is defined by a distribution
3. expected value: $\hat{H} =$

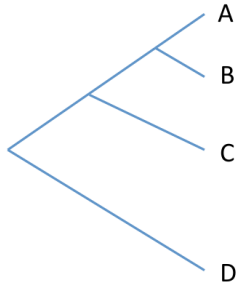
Summary statistics

Sample	Data	$P(G = AA D)$	$P(G = AG D)$	$P(G = GG D)$
1	A	0.66	0.33	0.01
2	AAAG	0.14	0.86	0.00
3	AGG	0.00	0.92	0.08
4	GG	0.00	0.20	0.80

with $Q = 15$

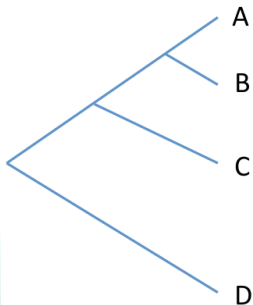
1. call genotypes: $H = 2$
2. sample genotypes: H is defined by a distribution
3. expected value: $\hat{H} =$
 $p(G_1 = AG) + p(G_2 = AG) + p(G_3 = AG) + p(G_4 = AG) = 0.33 + 0.86 + 0.92 + 0.20 = 2.31$

Genetic distances



Genotype 1	Genotype 2	Distance
aa	aa	0
aa	aA	1
aa	AA	2
aA	aa	1
aA	aA	0
aA	AA	2
...

Genetic distances



Genotypes are {aa, aA, AA} as {0, 1, 2}

For individuals i and j and N sites:

$$d(i, j) = -\log \left(1 - \frac{1}{N} \sum_{s=1}^N \frac{|g(i, s) - g(j, s)|}{2} \right)$$

genotype of i at site s

e.g. $G(i=A, s=1)=0$ and $G(j=B, s=1)=1$ then $d(i, j)=1$

Genetic distances from known genotypes

Genotypes are {aa, aA, AA} as {0, 1, 2}

For individuals i and j and N sites:

$$d(i, j) = -\log \left(1 - \frac{1}{N} \sum_{s=1}^N \frac{|g(i, s) - g(j, s)|}{2} \right)$$

$$d(i, j) = 1 * 1.00 = 1.00/2$$

		B		
		0	1	2
A	0	0	1	0
	1	0	0	0
	2	0	0	0

Genetic distances from **unknown** genotypes

Genotypes are {aa, aA, AA} as {0, 1, 2}

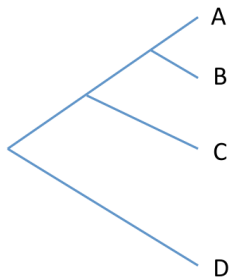
For individuals i and j and N sites:

$$d(i,j) = -\log \left(1 - \frac{1}{N} \sum_{s=1}^N \frac{|g(i,s) - g(j,s)|}{2} \right)$$

$$E[d(i,j)] = 0*0.30 + 1*0.50 + 2*0.10 + 1*0.10 + \dots = 0.80/2$$

		B		
		0	1	2
A	0	0.30	0.50	0.10
	1	0.10	0	0
	2	0	0	0

Genetic distances from **unknown** genotypes



Genotypes are {aa, aA, AA} as {0, 1, 2}

For individuals i and j and N sites:

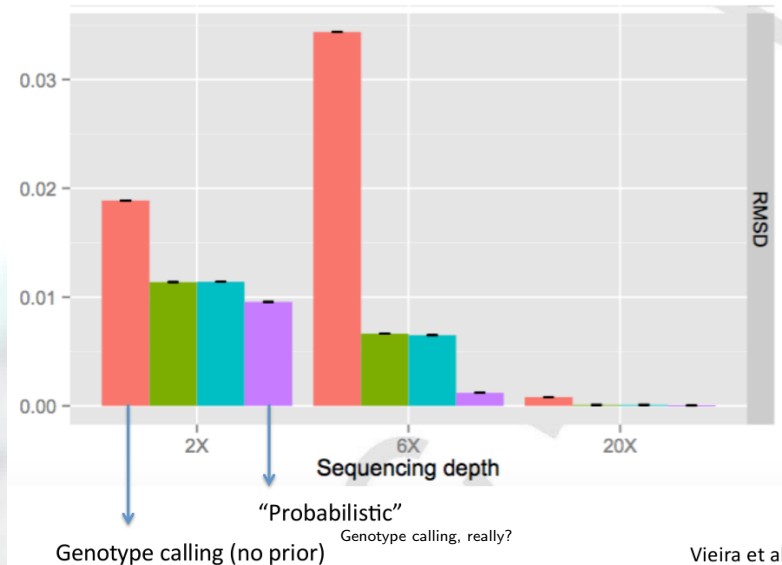
$$d(i, j) = -\log \left(1 - \frac{1}{N} \sum_{s=1}^N \frac{|g(i, s) - g(j, s)|}{2} \right)$$

Iterate across all possible genotypes

Genotypes probability

$$d(i, j) = -\log \left(1 - \frac{1}{N} \sum_{s=1}^N \sum_{g(i, s)=0}^2 \sum_{g(j, s)=0}^2 \frac{|g(i, s) - g(j, s)|}{2} * P(g(i, s), g(j, s)) \right)$$

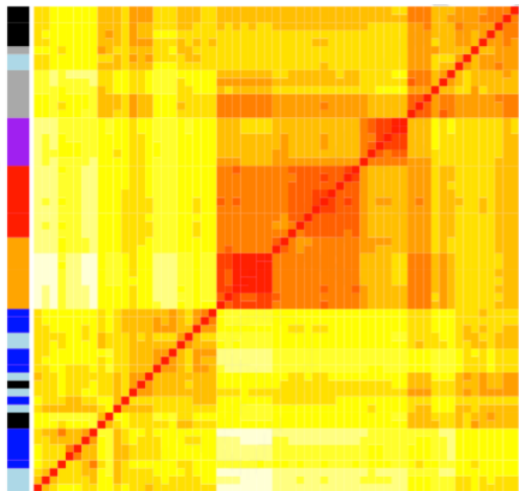
Genetic distances from **unknown** genotypes



Clustering from **unknown** genotypes

- nivara
- rufipogon
- Chinese rufipogon
- Indica
- aromatic
- tropical japonica
- temperate japonica

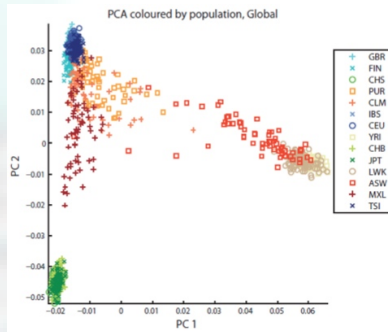
Original data: ~2M SNPs
Here: 5.4M SNPs at 2X



Population structure

Principal Component Analysis (PCA) is a data reduction method for:

- visualisation
- correction for population stratification
- information on population history and differentiation?



Genotype calling, really?

Covariance matrix

Genotype (0,1,2) Allele frequency

$$\text{cov}(i, j) = \frac{1}{(m-1)} \frac{\sum_{s=1}^m (G_s^{(i)} - 2\hat{p}_s)(G_s^{(j)} - 2\hat{p}_s)}{\sqrt{\hat{p}_s(1-\hat{p}_s)}}$$

Covariance matrix

Genotype (0,1,2) \swarrow Allele frequency

$$\text{cov}(i, j) = \frac{1}{(m-1)} \frac{\sum_{s=1}^m (G_s^{(i)} - 2\hat{p}_s)(G_s^{(j)} - 2\hat{p}_s)}{\sqrt{\hat{p}_s(1-\hat{p}_s)}}$$

Genotype (0,1,2) \swarrow Allele frequency

$$\text{cov}(i, j) = \frac{1}{(m-1)} \frac{\sum_{s=1}^m (G_s^{(i)} - 2\hat{p}_s)(G_s^{(j)} - 2\hat{p}_s)}{\sqrt{\hat{p}_s(1-\hat{p}_s)}}$$

Iterate across all genotypes \quad Weight by their probability

$$\text{cov}\hat{v}_{(i,j)} := \frac{1}{(\sum_{s=1}^m P_{\text{var},s}) - 1} \frac{\sum_{s=1}^m \sum_{G_s^{(i)}=0}^2 \sum_{G_s^{(j)}=0}^2 (G_s^{(i)} - 2\hat{p}_s)(G_s^{(j)} - 2\hat{p}_s) P(G_s^{(i)} | X_s^{(i)}) P(G_s^{(j)} | X_s^{(j)}) P_{\text{var},s}}{\sqrt{\hat{p}_s(1-\hat{p}_s)}}$$

Probability of the site being variable
(to avoid SNP calling)

Does it really work?

Genotype calling, really?

Exercise 2

Perform a PCA with called genotypes and using genotype probabilities.

Assess their relative performance.

Estimating allele frequencies

Assuming 2 alleles (A,G) with true allele frequency of 0.50

Sample	True genotype	Reads allele A	Read 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

What is the simplest estimator of allele frequencies?

Estimating allele frequencies

Assuming 2 alleles (A,G) with true allele frequency of 0.50

Sample	True genotype	Reads allele A	Read 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
Total		41	14

$$\hat{f} = \frac{\sum_{i=1}^N n_{A,i}}{\sum_{i=1}^N (n_{A,i} + n_{G,i})}$$

$$\hat{f} = 0.75$$

What is wrong with this estimator?

Estimating allele frequencies

Assuming 2 alleles (A,G) with true allele frequency of 0.50

Sample	True genotype	Reads allele A	Read 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
Total		41	14

$$\hat{n}_A = \sum_{i=1}^N (1 - \epsilon)n_{A,i} + \epsilon n_{G,i} - \epsilon n_{A,i} - (1 - \epsilon)n_{G,i}$$

$$\hat{f} = 0.77$$

Estimating allele frequencies

Maximum Likelihood estimator

$$P(D|f) = \prod_{i=1}^N \sum_{g \in \{0,1,2\}} P(D|G = g)P(G = g|f)$$

Estimating allele frequencies

Maximum Likelihood estimator

$$P(D|f) = \prod_{i=1}^N \sum_{g \in \{0,1,2\}} P(D|G = g)P(G = g|f)$$

$P(D|G = g)$ is the genotype likelihood and $P(G = g|f)$ is given by HWE (for instance).

In our previous example, $\hat{f} = 0.46$ which is much closer to the true value than previous estimators.

SNP calling

Challenges

- If high levels of missing data, then genotypes can be lost.
- Rare variants are hard to detect.
- Trade off between false positive and false negative rates.

How to call SNPs?

- If at least one heterozygous genotype has been called.
- If the estimated allele frequency is above a certain threshold.

Call a SNP if

$$\hat{f} \geq t$$

where t can be the minimum sample allele frequency detectable (e.g. $t = 1/2N$ with N diploids).

Likelihood Ratio Test

A Likelihood Ratio Test (LRT) compares the goodness of fit between the null and the alternative model:

- Null model: $f = 0$
- Alternative model: $f \neq 0$

$$T = -2 \log \frac{L(f = 0)}{L(f = \hat{f}_{MLE})}$$

where T is χ^2 distributed with 1 degree of freedom.

Exercise 3

Estimate allele frequencies and call SNPs.
Assess the accuracy using different thresholds.
Perform a new PCA with data filtering.

Sample allele frequency likelihoods

$$P(D|f) = \prod_{i=1}^N \sum_{g \in \{0,1,2\}} P(D|G = g)P(G = g|f)$$

$P(D f = 0)$	$P(D f = 1)$	$P(D f = 2)$	\dots	$P(D f = 2k)$
--------------	--------------	--------------	---------	---------------

with k diploids.

Sample allele frequency probabilities

If unfolded, $2k+1$ entries

$p_0=0$	$p_1=0$	$p_2=1$	$p_3=0$...	$p_{2k}=0$
---------	---------	---------	---------	-----	------------

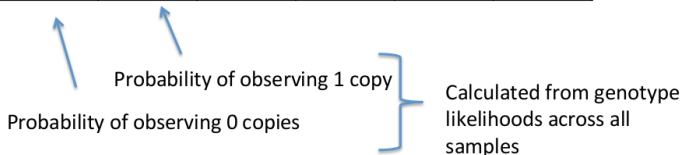
e.g. A is ancestral, G is derived (alternate)

AA AA AG AA AG AA AA AA AA

If genotypes are unknown and counting is not possible?

Sample allele frequency probabilities

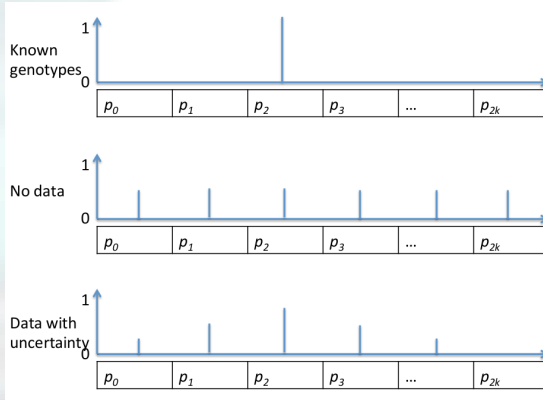
$p_0=0.05$	$p_1=0.15$	$p_2=0.70$	$p_3=0.10$...	p_{2k}
------------	------------	------------	------------	-----	----------



e.g. A is ancestral, G is derived (alternate)

If genotypes are unknown and counting is not possible.

Sample allele frequency probabilities



Sample allele frequency probabilities

Summary statistics

$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
--------------	--------------	--------------	-----	---------------

with k diploids.

$$\hat{f} =$$

Sample allele frequency probabilities

Summary statistics

$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
--------------	--------------	--------------	-----	---------------

with k diploids.

$$\hat{f} = \sum_{i=0}^{2k} \left(\frac{i}{2k}\right) \cdot P(f = i|D)$$

Sample allele frequency probabilities

Summary statistics

$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
--------------	--------------	--------------	-----	---------------

with k diploids.

$$P_{var} =$$

Sample allele frequency probabilities

Summary statistics

$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
--------------	--------------	--------------	-----	---------------

with k diploids.

$$P_{var} = 1 - P(f = 0|D) - P(f = 2k|D)$$

Sample allele frequency probabilities

Summary statistics - number of segregating sites

site 1	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
site 2	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
site 3	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
...	
site M	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$

$$E[S] =$$

Sample allele frequency probabilities

Summary statistics - number of segregating sites

site 1	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
site 2	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
site 3	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
...
site M	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$

$$E[S] = \sum_{j=1}^M$$

Sample allele frequency probabilities

Summary statistics - number of segregating sites

site 1	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
site 2	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
site 3	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
...
site M	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$

$$E[S] = \sum_{j=1}^M (1 - P(f_j = 0|D) - P(f_j = 2k|D))$$

Sample allele frequency probabilities

Summary statistics - nucleotide diversity $D = 2 \cdot f \cdot (1 - f)$

site 1	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
site 2	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
site 3	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
...
site M	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$

$$E[D] =$$

Sample allele frequency probabilities

Summary statistics - nucleotide diversity $D = 2 \cdot f \cdot (1 - f)$

site 1	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
site 2	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
site 3	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$
...
site M	$P(f = 0 D)$	$P(f = 1 D)$	$P(f = 2 D)$...	$P(f = 2k D)$

$$E[D] = \sum_{j=1}^M \sum_{i=0}^{2k} \binom{i}{2k} \cdot \binom{2k-i}{2k} \cdot P(f_j = i|D)$$

Site frequency spectrum (SFS)

The SFS is a summary of population genetic data. Each bin represents the proportion of sites with a particular derived (or minor) allele frequency.

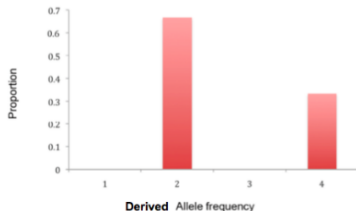
Ancestral sequence

a

c

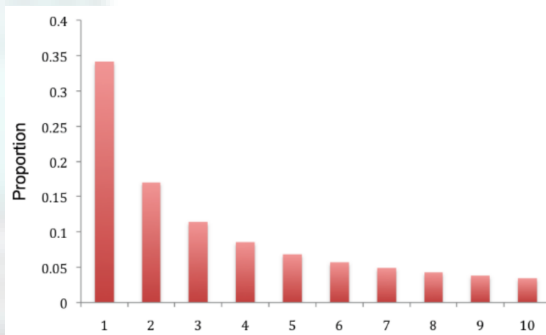
a

Sequence 1	aggtatgcta	gaa c cctaga	aagacacaga	gatagacaag
Sequence 2	aggtatgcta	gaa a cctaga	ta g acacaga	gatagacaag
Sequence 3	aggtatgcta	gaa a cctaga	ta g acacaga	gatagacaag
Sequence 4	aggtatgct g	gaa c cctaga	ta g acacaga	gatagacaag
Sequence 5	aggtatgct g	gaa c cctaga	ta g acacaga	gatagacaag



Site frequency spectrum (SFS)

SFS under the neutral coalescent model for a sample of $n = 11$ haploid individuals.



Site frequency spectrum (SFS)

From low-depth data:

- Let X be the sequencing data for our entire genome (all sites with ancestral and/or derived reads).
- X_s is the number of ancestral and derived reads at a particular site s .
- For 1 population, the SFS is a 1-dimensional vector $\vec{\gamma}$ with entries γ_i :
- $L(X|\gamma) = \prod_{s=1}^N L(X_s|\vec{\gamma}) = \prod_{s=1}^N \sum_{i=0}^{2n} \gamma_i P[X_s|D=i]$
- Then, we can use likelihood maximization algorithms to find a maximum likelihood estimate for each entry of the SFS (the values γ_i)

Exercise - 4

Calculate the SFS for each population and compare them. Estimate summary statistics (e.g. π , θ_W) and make some comments on their different values.

Advanced:

Calculate the joint (2D) SFS and F_{ST} in sliding windows.

Examples of optimal experimental design.

Thank you for your attention