

Strategies to work with HLA data in human populations: a tutorial to format HLA data and make basic analyses

Prof. Alicia Sanchez-Mazas & Dr José Manuel Nunes

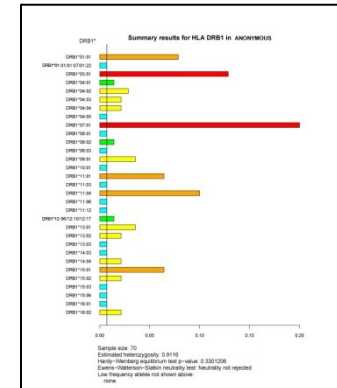
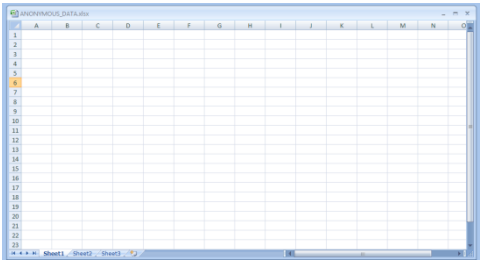
Dpt Genetics & Evolution - Anthropology Unit

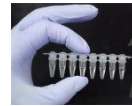
University of Geneva



Anonymous population sample

How to format my data and compute basic statistics?





Anonymous population sample

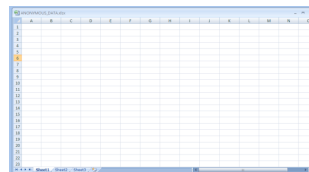


The data

Most of the times, your original table is an Excel file with genotypes, let's work with it

[http://geneva.../Anonymous .xls](http://geneva.../Anonymous.xls)

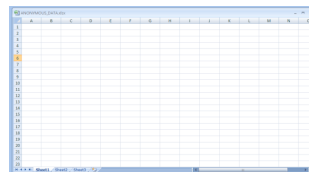




Data formatting

This is an example of original table of HLA genotypes for one locus: It contains a sample number, the population name, and pairs of possible genotypes

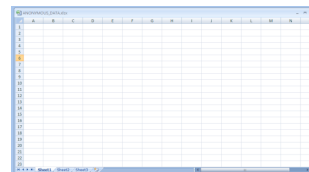
Sample Number	Population	DRB1-allele1	DRB1-allele2	DRB1-allele1	DRB1-allele2	DRB1-allele1	DRB1-allele2
ANON_1	Anonymous	*1404	*1404				
ANON_2	Anonymous	*03010101/03010102/030108	*110401				
ANON_3	Anonymous	*03010101/03010102/030108	*110401				
ANON_4	Anonymous	*15010101/15010102	*150201				
ANON_5	Anonymous	*07010101/07010102	*07010101/07010102				
ANON_6	Anonymous	*010101	*160201				
ANON_7	Anonymous	*010101	*03010101/03010102/030108				
ANON_8	Anonymous	*010101	*03010101/03010102/030108				
ANON_9	Anonymous	*010101/0107/0122	*07010101/07010102/070102				
ANON_10	Anonymous	*010101	*110101/110108				
ANON_11	Anonymous	*07010101/07010102	*07010101/07010102				
ANON_12	Anonymous	*07010101/07010102	*111201				
ANON_13	Anonymous	*07010101/07010102	*160201				
ANON_14	Anonymous	*0402	*07010101/07010102				
ANON_15	Anonymous	*110401	*15010101/15010102				
ANON_16	Anonymous	not done					
ANON_17	Anonymous	*07010101/07010102	*130201				
ANON_18	Anonymous	*03010101/03010102/030108	*130201				
ANON_19	Anonymous	*040301	*160201				
ANON_20	Anonymous	*010101	*110101/110108				
ANON_21	Anonymous	*07010101/07010102	*090102	*07010101/07010102	*0906	*07010101/07010102	*0909
ANON_22	Anonymous	*040501	*100101	*040503	*100101	*040504	*100101
ANON_23	Anonymous	*090102	*110401				
ANON_24	Anonymous	*07010101/07010102	*1506				



Data formatting

The data may include untyped samples (here for sample number ANON_16, in green): These should be removed. Select the line(s) and remove it (them).

Sample Number	Population	DRB1-allele1	DRB1-allele2	DRB1-allele1	DRB1-allele2	DRB1-allele1	DRB1-allele2
1							
2	ANON_1	Anonymous	*1404	*1404			
3	ANON_2	Anonymous	*03010101/03010102/030108	*110401			
4	ANON_3	Anonymous	*03010101/03010102/030108	*110401			
5	ANON_4	Anonymous	*15010101/15010102	*150201			
6	ANON_5	Anonymous	*07010101/07010102	*07010101/07010102			
7	ANON_6	Anonymous	*010101	*160201			
8	ANON_7	Anonymous	*010101	*03010101/03010102/030108			
9	ANON_8	Anonymous	*010101	*03010101/03010102/030108			
10	ANON_9	Anonymous	*010101/0107/0122	*07010101/07010102/070102			
11	ANON_10	Anonymous	*010101	*110101/110108			
12	ANON_11	Anonymous	*07010101/07010102	*07010101/07010102			
13	ANON_12	Anonymous	*07010101/07010102	*111201			
14	ANON_13	Anonymous	*07010101/07010102	*160201			
15	ANON_14	Anonymous	*0402	*07010101/07010102			
16	ANON_15	Anonymous	*110401	*15010101/15010102			
17	ANON_16	Anonymous	not done				
18	ANON_17	Anonymous	*07010101/07010102	*130201			
19	ANON_18	Anonymous	*03010101/03010102/030108	*130201			
20	ANON_19	Anonymous	*040301	*160201			
21	ANON_20	Anonymous	*010101	*110101/110108			
22	ANON_21	Anonymous	*07010101/07010102	*090102	*07010101/07010102	*0906	*07010101/07010102
23	ANON_22	Anonymous	*040501	*100101	*040503	*100101	*040504
24	ANON_23	Anonymous	*090102	*110401			
25	ANON_24	Anonymous	*07010101/07010102	*1506			



Data formatting

Done: untyped sample(s) have been removed (here ANON_16).

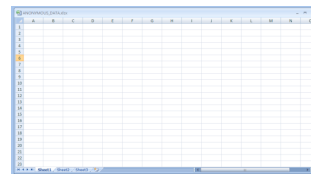
Sample Number	Population	DRB1-allele1	DRB1-allele2	DRB1-allele1	DRB1-allele2	DRB1-allele1	DRB1-allele2
1							
2	ANON_1	Anonymous	*1404	*1404			
3	ANON_2	Anonymous	*03010101/03010102/030108	*110401			
4	ANON_3	Anonymous	*03010101/03010102/030108	*110401			
5	ANON_4	Anonymous	*15010101/15010102	*150201			
6	ANON_5	Anonymous	*07010101/07010102	*07010101/07010102			
7	ANON_6	Anonymous	*010101	*160201			
8	ANON_7	Anonymous	*010101	*03010101/03010102/030108			
9	ANON_8	Anonymous	*010101	*03010101/03010102/030108			
10	ANON_9	Anonymous	*010101/0107/0122	*07010101/07010102/070102			
11	ANON_10	Anonymous	*010101	*110101/110108			
12	ANON_11	Anonymous	*07010101/07010102	*07010101/07010102			
13	ANON_12	Anonymous	*07010101/07010102	*111201			
14	ANON_13	Anonymous	*07010101/07010102	*160201			
15	ANON_14	Anonymous	*0402	*07010101/07010102			
16	ANON_15	Anonymous	*110401	*15010101/15010102			
17	ANON_17	Anonymous	*07010101/07010102	*130201			
18	ANON_18	Anonymous	*03010101/03010102/030108	*130201			
19	ANON_19	Anonymous	*040301	*160201			
20	ANON_20	Anonymous	*010101	*110101/110108			
21	ANON_21	Anonymous	*07010101/07010102	*090102	*07010101/07010102	*0906	*07010101/07010102 *0909
22	ANON_22	Anonymous	*040501	*100101	*040503	*100101	*040504 *100101
23	ANON_23	Anonymous	*090102	*110401			
24	ANON_24	Anonymous	*07010101/07010102	*1506			
25	ANON_25	Anonymous	*03010101/03010102/030108	*15010101/15010102			



Data formatting

The data may include ambiguous typings (here in green for ANON_21 and ANON_22)..

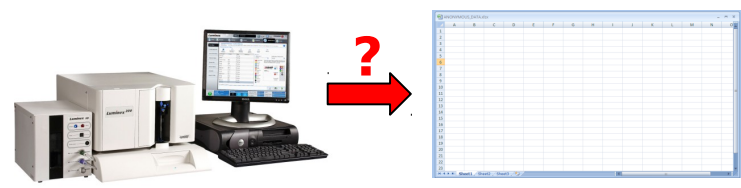
Sample Number	Population	DRB1-allele1	DRB1-allele2	DRB1-allele1	DRB1-allele2	DRB1-allele1	DRB1-allele2
ANON_1	Anonymous	*1404	*1404				
ANON_2	Anonymous	*03010101/03010102/030108	*110401				
ANON_3	Anonymous	*03010101/03010102/030108	*110401				
ANON_4	Anonymous	*15010101/15010102	*150201				
ANON_5	Anonymous	*07010101/07010102	*07010101/07010102				
ANON_6	Anonymous	*010101	*160201				
ANON_7	Anonymous	*010101	*03010101/03010102/030108				
ANON_8	Anonymous	*010101	*03010101/03010102/030108				
ANON_9	Anonymous	*010101/0107/0122	*07010101/07010102/070102				
ANON_10	Anonymous	*010101	*110101/110108				
ANON_11	Anonymous	*07010101/07010102	*07010101/07010102				
ANON_12	Anonymous	*07010101/07010102	*111201				
ANON_13	Anonymous	*07010101/07010102	*160201				
ANON_14	Anonymous	*0402	*07010101/07010102				
ANON_15	Anonymous	*110401	*15010101/15010102				
ANON_17	Anonymous	*07010101/07010102	*130201				
ANON_18	Anonymous	*03010101/03010102/030108	*130201				
ANON_19	Anonymous	*040301	*160201				
ANON_20	Anonymous	*010101	*110101/110108				
ANON_21	Anonymous	*07010101/07010102	*090102	*07010101/07010102	*0906	*07010101/07010102	*0909
ANON_22	Anonymous	*040501	*100101	*040503	*100101	*040504	*100101
ANON_23	Anonymous	*090102	*110401				
ANON_24	Anonymous	*07010101/07010102	*1506				
ANON_25	Anonymous	*03010101/03010102/030108	*15010101/15010102				



Data formatting

..and/or an old nomenclature (here in green)

Sample Number	Population	DRB1-allele1	DRB1-allele2	DRB1-allele1	DRB1-allele2	DRB1-allele1	DRB1-allele2
ANON_1	Anonymous	*1404	*1404				
ANON_2	Anonymous	*03010101/03010102/030108	*110401				
ANON_3	Anonymous	*03010101/03010102/030108	*110401				
ANON_4	Anonymous	*15010101/15010102	*150201				
ANON_5	Anonymous	*07010101/07010102	*07010101/07010102				
ANON_6	Anonymous	*010101	*160201				
ANON_7	Anonymous	*010101	*03010101/03010102/030108				
ANON_8	Anonymous	*010101	*03010101/03010102/030108				
ANON_9	Anonymous	*010101/0107/0122	*07010101/07010102/070102				
ANON_10	Anonymous	*010101	*110101/110108				
ANON_11	Anonymous	*07010101/07010102	*07010101/07010102				
ANON_12	Anonymous	*07010101/07010102	*111201				
ANON_13	Anonymous	*07010101/07010102	*160201				
ANON_14	Anonymous	*0402	*07010101/07010102				
ANON_15	Anonymous	*110401	*15010101/15010102				
ANON_17	Anonymous	*07010101/07010102	*130201				
ANON_18	Anonymous	*03010101/03010102/030108	*130201				
ANON_19	Anonymous	*040301	*160201				
ANON_20	Anonymous	*010101	*110101/110108				
ANON_21	Anonymous	*07010101/07010102	*090102	*07010101/07010102	*0906	*07010101/07010102	*0909
ANON_22	Anonymous	*040501	*100101	*040503	*100101	*040504	*100101
ANON_23	Anonymous	*090102	*110401				
ANON_24	Anonymous	*07010101/07010102	*1506				
ANON_25	Anonymous	*03010101/03010102/030108	*15010101/15010102				



Data formatting

..ambiguous data and data coded with an old nomenclature can be recoded with the transliterate tool of Gene[rate]..

Sample Number	Population	DRB1-allele1	DRB1-allele2	DRB1-allele1	DRB1-allele2	DRB1-allele1	DRB1-allele2
1							
2	ANON_1	Anonymous	*1404	*1404			
3	ANON_2	Anonymous	*03010101/03010102/030108				
4	ANON_3	Anonymous	*03010101/03010102/030108				
5	ANON_4	Anonymous	*15010101/15010102				
6	ANON_5	Anonymous	*07010101/07010102				
7	ANON_6	Anonymous	*010101				
8	ANON_7	Anonymous	*010101				
9	ANON_8	Anonymous	*010101				
10	ANON_9	Anonymous	*010101/0107/0122				
11	ANON_10	Anonymous	*010101				
12	ANON_11	Anonymous	*07010101/07010102				
13	ANON_12	Anonymous	*07010101/07010102				
14	ANON_13	Anonymous	*07010101/07010102				
15	ANON_14	Anonymous	*0402				
16	ANON_15	Anonymous	*110401				
17	ANON_17	Anonymous	*07010101/07010102				
18	ANON_18	Anonymous	*03010101/03010102/030108				
19	ANON_19	Anonymous	*040301				
20	ANON_20	Anonymous	*010101				
21	ANON_21	Anonymous	*07010101/07010102				
22	ANON_22	Anonymous	*040501				
23	ANON_23	Anonymous	*090102				
24	ANON_24	Anonymous	*07010101/07010102				
25	ANON_25	Anonymous	*03010101/03010102/030108				

Gene[RATE]

Tools

- Phenotype
- Haplotype
- Transliterate**
- Uniformate
- File conversions
- Frequency estimation

Navigation:

AGP lab
 Laboratory of Anthropology, Genetics and Peopling History

Generate
 Tools for manipulation of data with ambiguities

Documentation
 Howtos, manuals, questionnaires and other documents

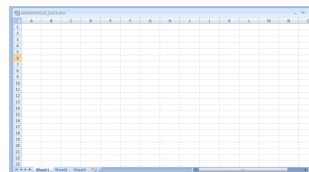
AHPD
 Project and documents

Search geneva.unige.ch:

Substitution description file:

UNIFORMATE data file:

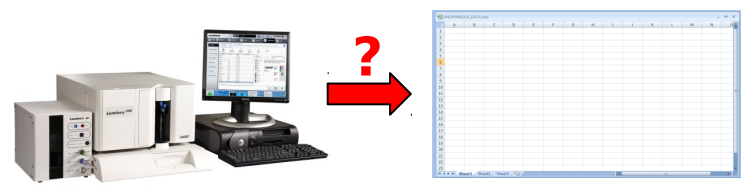
transliterate UNIFORMATE datafiles



Data formatting

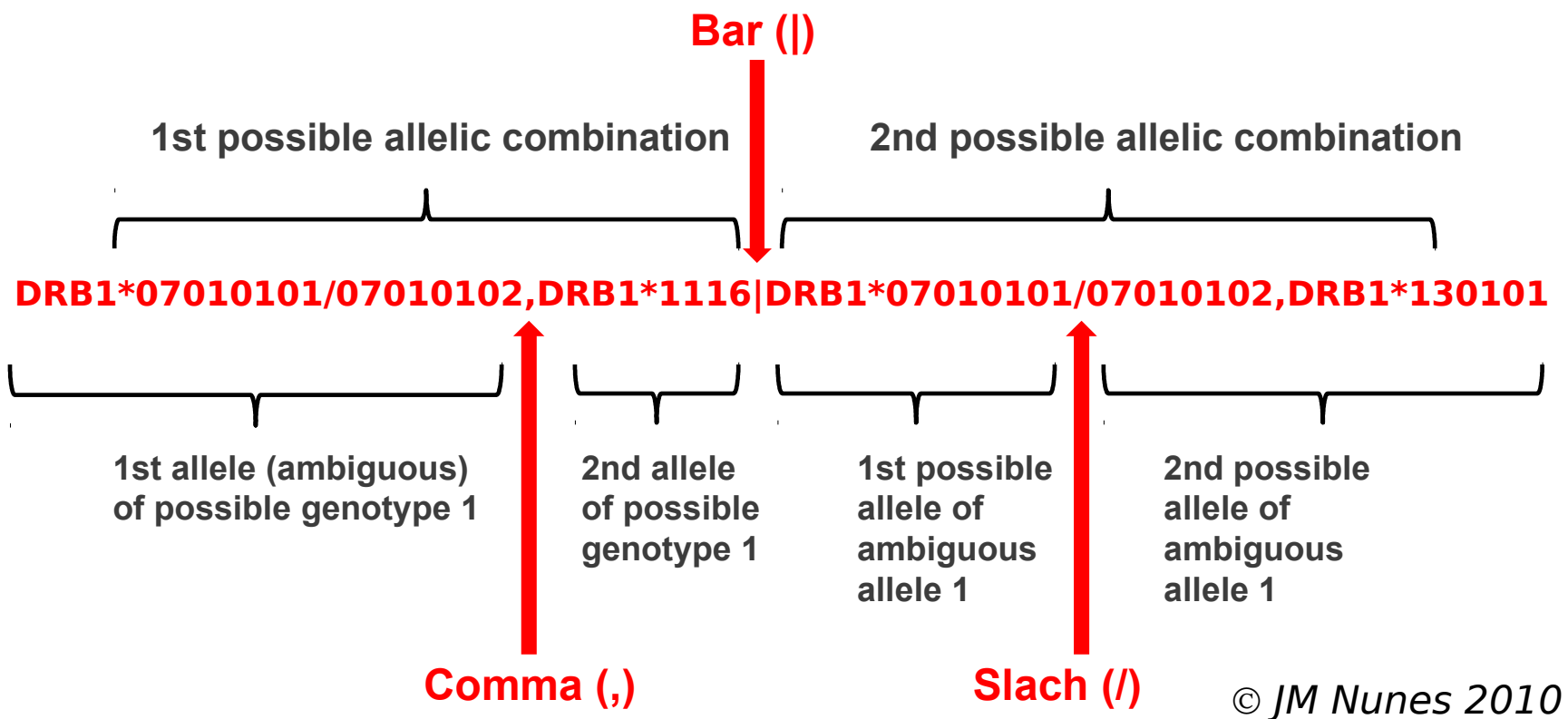
But to use the Gene[rate] tools, your data have first to be converted into a text file containing your data in the UNIFORMAT data format. Let's do it step by step

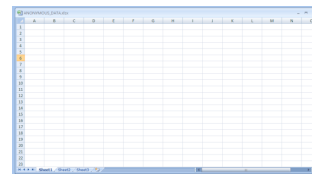
Sample Number	Population	DRB1-allele1	DRB1-allele2	DRB1-allele1	DRB1-allele2	DRB1-allele1	DRB1-allele2
1							
2	ANON_1	Anonymous	*1404	*1404			
3	ANON_2	Anonymous	*03010101/03010102/030108	*110401			
4	ANON_3	Anonymous	*03010101/03010102/030108	*110401			
5	ANON_4	Anonymous	*15010101/15010102	*150201			
6	ANON_5	Anonymous	*07010101/07010102	*07010101/07010102			
7	ANON_6	Anonymous	*010101	*160201			
8	ANON_7	Anonymous	*010101	*03010101/03010102/030108			
9	ANON_8	Anonymous	*010101	*03010101/03010102/030108			
10	ANON_9	Anonymous	*010101/0107/0122	*07010101/07010102/070102			
11	ANON_10	Anonymous	*010101	*110101/110108			
12	ANON_11	Anonymous	*07010101/07010102	*07010101/07010102			
13	ANON_12	Anonymous	*07010101/07010102	*111201			
14	ANON_13	Anonymous	*07010101/07010102	*160201			
15	ANON_14	Anonymous	*0402	*07010101/07010102			
16	ANON_15	Anonymous	*110401	*15010101/15010102			
17	ANON_17	Anonymous	*07010101/07010102	*130201			
18	ANON_18	Anonymous	*03010101/03010102/030108	*130201			
19	ANON_19	Anonymous	*040301	*160201			
20	ANON_20	Anonymous	*010101	*110101/110108			
21	ANON_21	Anonymous	*07010101/07010102	*090102	*07010101/07010102	*0906	*07010101/07010102 *0909
22	ANON_22	Anonymous	*040501	*100101	*040503	*100101	*040504 *100101
23	ANON_23	Anonymous	*090102	*110401			
24	ANON_24	Anonymous	*07010101/07010102	*1506			
25	ANON_25	Anonymous	*03010101/03010102/030108	*15010101/15010102			



Data formatting

UNIFORMAT data format:

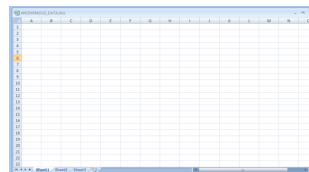




Data formatting

Step 1: Insert a column after each pair of columns defining a genotype. Start with the first genotype pair.

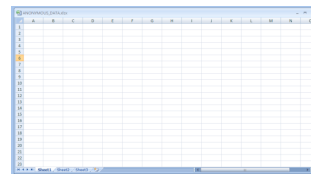
Sample Number	Population	DRB1-allele1	DRB1-allele2	DRB1-a	DRB1-allele1	DRB1-allele2
1						
2	ANON_1	Anonymous	*1404	*1404		
3	ANON_2	Anonymous	*03010101/03010102/030108	*110401		
4	ANON_3	Anonymous	*03010101/03010102/030108	*110401		
5	ANON_4	Anonymous	*15010101/15010102	*150201		
6	ANON_5	Anonymous	*07010101/07010102	*07010101/07010102		
7	ANON_6	Anonymous	*010101	*160201		
8	ANON_7	Anonymous	*010101	*03010101/03010102/030108		
9	ANON_8	Anonymous	*010101	*03010101/03010102/030108		
10	ANON_9	Anonymous	*010101/0107/0122	*07010101/07010102/070102		
11	ANON_10	Anonymous	*010101	*110101/110108		
12	ANON_11	Anonymous	*07010101/07010102	*07010101/07010102		
13	ANON_12	Anonymous	*07010101/07010102	*111201		
14	ANON_13	Anonymous	*07010101/07010102	*160201		
15	ANON_14	Anonymous	*0402	*07010101/07010102		
16	ANON_15	Anonymous	*110401	*15010101/15010102		
17	ANON_17	Anonymous	*07010101/07010102	*130201		
18	ANON_18	Anonymous	*03010101/03010102/030108	*130201		
19	ANON_19	Anonymous	*040301	*160201		
20	ANON_20	Anonymous	*010101	*110101/110108		
21	ANON_21	Anonymous	*07010101/07010102	*090102	*07010101/07010102	*0909
22	ANON_22	Anonymous	*040501	*100101	*040503	*100101
23	ANON_23	Anonymous	*090102	*110401		
24	ANON_24	Anonymous	*07010101/07010102	*1506		
25	ANON_25	Anonymous	*03010101/03010102/030108	*15010101/15010102		



Data formatting

Step 1: done for the first genotype pair!

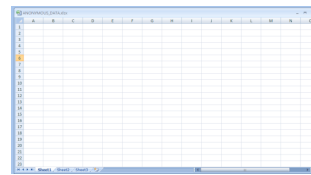
	A	B	C	D	E	F	G	H
	Sample Number	Population	DRB1-allele1	DRB1-allele2		DRB1-allele1	DRB1-allele2	DRB1-allele1
1								
2	ANON_1	Anonymous	*1404	*1404				
3	ANON_2	Anonymous	*03010101/03010102/030108	*110401				
4	ANON_3	Anonymous	*03010101/03010102/030108	*110401				
5	ANON_4	Anonymous	*15010101/15010102	*150201				
6	ANON_5	Anonymous	*07010101/07010102	*07010101/07010102				
7	ANON_6	Anonymous	*010101	*160201				
8	ANON_7	Anonymous	*010101	*03010101/03010102/030108				
9	ANON_8	Anonymous	*010101	*03010101/03010102/030108				
10	ANON_9	Anonymous	*010101/0107/0122	*07010101/07010102/070102				
11	ANON_10	Anonymous	*010101	*110101/110108				
12	ANON_11	Anonymous	*07010101/07010102	*07010101/07010102				
13	ANON_12	Anonymous	*07010101/07010102	*111201				
14	ANON_13	Anonymous	*07010101/07010102	*160201				
15	ANON_14	Anonymous	*0402	*07010101/07010102				
16	ANON_15	Anonymous	*110401	*15010101/15010102				
17	ANON_17	Anonymous	*07010101/07010102	*130201				
18	ANON_18	Anonymous	*03010101/03010102/030108	*130201				
19	ANON_19	Anonymous	*040301	*160201				
20	ANON_20	Anonymous	*010101	*110101/110108				
21	ANON_21	Anonymous	*07010101/07010102	*090102		*07010101/07010102	*0906	*07010101/07010102
22	ANON_22	Anonymous	*040501	*100101		*040503	*100101	*040504
23	ANON_23	Anonymous	*090102	*110401				
24	ANON_24	Anonymous	*07010101/07010102	*1506				
25	ANON_25	Anonymous	*03010101/03010102/030108	*15010101/15010102				



Data formatting

Step 1: continue for the second genotype pair.

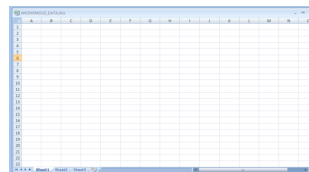
Sample Number	Population	DRB1-allele1	DRB1-allele2	DRB1-allele1	DRB1-allele2
1					
2	ANON_1	Anonymous	*1404	*1404	
3	ANON_2	Anonymous	*03010101/03010102/030108	*110401	
4	ANON_3	Anonymous	*03010101/03010102/030108	*110401	
5	ANON_4	Anonymous	*15010101/15010102	*150201	
6	ANON_5	Anonymous	*07010101/07010102	*07010101/07010102	
7	ANON_6	Anonymous	*010101	*160201	
8	ANON_7	Anonymous	*010101	*03010101/03010102/030108	
9	ANON_8	Anonymous	*010101	*03010101/03010102/030108	
10	ANON_9	Anonymous	*010101/0107/0122	*07010101/07010102/070102	
11	ANON_10	Anonymous	*010101	*110101/110108	
12	ANON_11	Anonymous	*07010101/07010102	*07010101/07010102	
13	ANON_12	Anonymous	*07010101/07010102	*111201	
14	ANON_13	Anonymous	*07010101/07010102	*160201	
15	ANON_14	Anonymous	*0402	*07010101/07010102	
16	ANON_15	Anonymous	*110401	*15010101/15010102	
17	ANON_17	Anonymous	*07010101/07010102	*130201	
18	ANON_18	Anonymous	*03010101/03010102/030108	*130201	
19	ANON_19	Anonymous	*040301	*160201	
20	ANON_20	Anonymous	*010101	*110101/110108	
21	ANON_21	Anonymous	*07010101/07010102	*090102	*07010101/07010102
22	ANON_22	Anonymous	*040501	*100101	*040503
23	ANON_23	Anonymous	*090102	*110401	*0906
24	ANON_24	Anonymous	*07010101/07010102	*1506	*100101
25	ANON_25	Anonymous	*03010101/03010102/030108	*15010101/15010102	



Data formatting

Step 1: done for the second genotype pair!

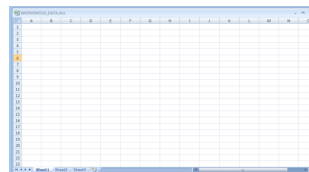
Sample Number	Population	DRB1-allele1	DRB1-allele2	DRB1-allele1	DRB1-allele2
1					
2	ANON_1	Anonymous	*1404	*1404	
3	ANON_2	Anonymous	*03010101/03010102/030108	*110401	
4	ANON_3	Anonymous	*03010101/03010102/030108	*110401	
5	ANON_4	Anonymous	*15010101/15010102	*150201	
6	ANON_5	Anonymous	*07010101/07010102	*07010101/07010102	
7	ANON_6	Anonymous	*010101	*160201	
8	ANON_7	Anonymous	*010101	*03010101/03010102/030108	
9	ANON_8	Anonymous	*010101	*03010101/03010102/030108	
10	ANON_9	Anonymous	*010101/0107/0122	*07010101/07010102/070102	
11	ANON_10	Anonymous	*010101	*110101/110108	
12	ANON_11	Anonymous	*07010101/07010102	*07010101/07010102	
13	ANON_12	Anonymous	*07010101/07010102	*111201	
14	ANON_13	Anonymous	*07010101/07010102	*160201	
15	ANON_14	Anonymous	*0402	*07010101/07010102	
16	ANON_15	Anonymous	*110401	*15010101/15010102	
17	ANON_17	Anonymous	*07010101/07010102	*130201	
18	ANON_18	Anonymous	*03010101/03010102/030108	*130201	
19	ANON_19	Anonymous	*040301	*160201	
20	ANON_20	Anonymous	*010101	*110101/110108	
21	ANON_21	Anonymous	*07010101/07010102	*090102	*07010101/07010102 *0906
22	ANON_22	Anonymous	*040501	*100101	*040503 *100101
23	ANON_23	Anonymous	*090102	*110401	
24	ANON_24	Anonymous	*07010101/07010102	*1506	
25	ANON_25	Anonymous	*03010101/03010102/030108	*15010101/15010102	



Data formatting

Step 2: Fill the new column with a bar (|)
The bar (|) can be written by using <Alt Gr> and « / »
Start with the first column to fill.

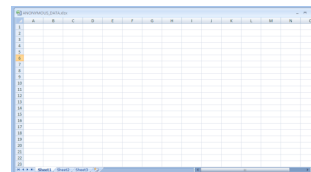
Sample Number	Population	DRB1-allele1	DRB1-allele2		DRB1-allele1	DRB1-allele2
ANON_1	Anonymous	*1404	*1404			
ANON_2	Anonymous	*030101/030102/030108	*110401			
ANON_3	Anonymous	*030101/030102/030108	*110401			
ANON_4	Anonymous	*150101/150102	*150201			
ANON_5	Anonymous	*070101/070102	*070101/070102			
ANON_6	Anonymous	*010101	*160201			
ANON_7	Anonymous	*010101	*030101/030102/030108			
ANON_8	Anonymous	*010101	*030101/030102/030108			
ANON_9	Anonymous	*010101/0107/0122	*070101/070102/070102			
ANON_10	Anonymous	*010101	*110101/110108			
ANON_11	Anonymous	*070101/070102	*070101/070102			
ANON_12	Anonymous	*070101/070102	*111201			
ANON_13	Anonymous	*070101/070102	*160201			
ANON_14	Anonymous	*0402	*070101/070102			
ANON_15	Anonymous	*110401	*150101/150102			
ANON_17	Anonymous	*070101/070102	*130201			
ANON_18	Anonymous	*030101/030102/030108	*130201			
ANON_19	Anonymous	*040301	*160201			
ANON_20	Anonymous	*010101	*110101/110108			
ANON_21	Anonymous	*070101/070102	*090102		*070101/070102	*0906
ANON_22	Anonymous	*040501	*100101		*040503	*100101
ANON_23	Anonymous	*090102	*110401			
ANON_24	Anonymous	*070101/070102	*1506			
ANON_25	Anonymous	*030101/030102/030108	*150101/150102			



Data formatting

**Step 2: done for the first column to fill!
 Start with the second column to fill.**

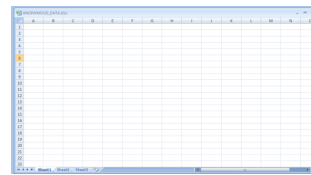
	A	B	C	D	E	F	G	H
1	Sample Number	Population	DRB1-allele1	DRB1-allele2		DRB1-allele1	DRB1-allele2	
2	ANON_1	Anonymous	*1404	*1404				
3	ANON_2	Anonymous	*03010101/03010102/030108	*110401				
4	ANON_3	Anonymous	*03010101/03010102/030108	*110401				
5	ANON_4	Anonymous	*15010101/15010102	*150201				
6	ANON_5	Anonymous	*07010101/07010102	*07010101/07010102				
7	ANON_6	Anonymous	*010101	*160201				
8	ANON_7	Anonymous	*010101	*03010101/03010102/030108				
9	ANON_8	Anonymous	*010101	*03010101/03010102/030108				
10	ANON_9	Anonymous	*010101/0107/0122	*07010101/07010102/070102				
11	ANON_10	Anonymous	*010101	*110101/110108				
12	ANON_11	Anonymous	*07010101/07010102	*07010101/07010102				
13	ANON_12	Anonymous	*07010101/07010102	*111201				
14	ANON_13	Anonymous	*07010101/07010102	*160201				
15	ANON_14	Anonymous	*0402	*07010101/07010102				
16	ANON_15	Anonymous	*110401	*15010101/15010102				
17	ANON_17	Anonymous	*07010101/07010102	*130201				
18	ANON_18	Anonymous	*03010101/03010102/030108	*130201				
19	ANON_19	Anonymous	*040301	*160201				
20	ANON_20	Anonymous	*010101	*110101/110108				
21	ANON_21	Anonymous	*07010101/07010102	*090102		*07010101/07010102	*0906	
22	ANON_22	Anonymous	*040501	*100101		*040503	*100101	
23	ANON_23	Anonymous	*090102	*110401				
24	ANON_24	Anonymous	*07010101/07010102	*1506				
25	ANON_25	Anonymous	*03010101/03010102/030108	*15010101/15010102				



Data formatting

**Step 3: Insert a semi-column (;) between the two alleles of each allelic pair
 First insert a column.**

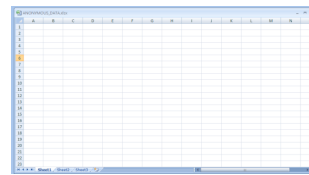
Sample Number	Population	DRB1-allele1	DRB1-allele2	DRB1-allele1	DRB1-allele2
1					
2	ANON_1	Anonymous	*1404	*1404	
3	ANON_2	Anonymous	*03010101/03010102/030108	*110401	
4	ANON_3	Anonymous	*03010101/03010102/030108	*110401	
5	ANON_4	Anonymous	*15010101/15010102	*150201	
6	ANON_5	Anonymous	*07010101/07010102	*070101	
7	ANON_6	Anonymous	*010101	*160201	
8	ANON_7	Anonymous	*010101	*030101	
9	ANON_8	Anonymous	*010101	*030101	
10	ANON_9	Anonymous	*010101/0107/0122	*070101	
11	ANON_10	Anonymous	*010101	*110101/110108	
12	ANON_11	Anonymous	*07010101/07010102	*07010101/07010102	
13	ANON_12	Anonymous	*07010101/07010102	*111201	
14	ANON_13	Anonymous	*07010101/07010102	*160201	
15	ANON_14	Anonymous	*0402	*07010101/07010102	
16	ANON_15	Anonymous	*110401	*15010101/15010102	
17	ANON_17	Anonymous	*07010101/07010102	*130201	
18	ANON_18	Anonymous	*03010101/03010102/030108	*130201	
19	ANON_19	Anonymous	*040301	*160201	
20	ANON_20	Anonymous	*010101	*110101/110108	
21	ANON_21	Anonymous	*07010101/07010102	*090102	*07010101/07010102 *0906
22	ANON_22	Anonymous	*040501	*100101	*040503 *100101
23	ANON_23	Anonymous	*090102	*110401	
24	ANON_24	Anonymous	*07010101/07010102	*1506	
25	ANON_25	Anonymous	*03010101/03010102/030108	*15010101/15010102	



Data formatting

Step 3: column inserted for the first genotype (allelic pair)!

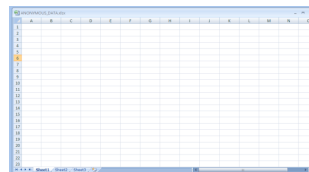
	A	B	C	D	E	F	G	H
1	Sample Number	Population	DRB1-allele1		DRB1-allele2		DRB1-allele1	DRB1-allele2
2	ANON_1	Anonymous	*1404		*1404			
3	ANON_2	Anonymous	*03010101/03010102/030108		*110401			
4	ANON_3	Anonymous	*03010101/03010102/030108		*110401			
5	ANON_4	Anonymous	*15010101/15010102		*150201			
6	ANON_5	Anonymous	*07010101/07010102		*07010101/07010102			
7	ANON_6	Anonymous	*010101		*160201			
8	ANON_7	Anonymous	*010101		*03010101/03010102/030108			
9	ANON_8	Anonymous	*010101		*03010101/03010102/030108			
10	ANON_9	Anonymous	*010101/0107/0122		*07010101/07010102/070102			
11	ANON_10	Anonymous	*010101		*110101/110108			
12	ANON_11	Anonymous	*07010101/07010102		*07010101/07010102			
13	ANON_12	Anonymous	*07010101/07010102		*111201			
14	ANON_13	Anonymous	*07010101/07010102		*160201			
15	ANON_14	Anonymous	*0402		*07010101/07010102			
16	ANON_15	Anonymous	*110401		*15010101/15010102			
17	ANON_17	Anonymous	*07010101/07010102		*130201			
18	ANON_18	Anonymous	*03010101/03010102/030108		*130201			
19	ANON_19	Anonymous	*040301		*160201			
20	ANON_20	Anonymous	*010101		*110101/110108			
21	ANON_21	Anonymous	*07010101/07010102		*090102		*07010101/07010102	*0906
22	ANON_22	Anonymous	*040501		*100101		*040503	*100101
23	ANON_23	Anonymous	*090102		*110401			
24	ANON_24	Anonymous	*07010101/07010102		*1506			
25	ANON_25	Anonymous	*03010101/03010102/030108		*15010101/15010102			



Data formatting

Step 3: insert a column for the second genotype (allelic pair).

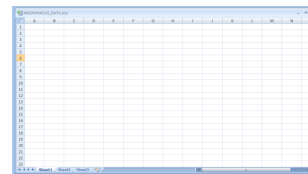
Sample Number	Population	DRB1-allele1	DRB1-allele2	DRB1
ANON_1	Anonymous	*1404	*1404	
ANON_2	Anonymous	*03010101/03010102/030108	*110401	
ANON_3	Anonymous	*03010101/03010102/030108	*110401	
ANON_4	Anonymous	*15010101/15010102	*150201	
ANON_5	Anonymous	*07010101/07010102	*07010101/07010102	
ANON_6	Anonymous	*010101	*160201	
ANON_7	Anonymous	*010101	*03010101/03010102/030108	
ANON_8	Anonymous	*010101	*03010101/03010102/030108	
ANON_9	Anonymous	*010101/0107/0122	*07010101/07010102/070102	
ANON_10	Anonymous	*010101	*110101/110108	
ANON_11	Anonymous	*07010101/07010102	*07010101/07010102	
ANON_12	Anonymous	*07010101/07010102	*111201	
ANON_13	Anonymous	*07010101/07010102	*160201	
ANON_14	Anonymous	*0402	*07010101/07010102	
ANON_15	Anonymous	*110401	*15010101/15010102	
ANON_17	Anonymous	*07010101/07010102	*130201	
ANON_18	Anonymous	*03010101/03010102/030108	*130201	
ANON_19	Anonymous	*040301	*160201	
ANON_20	Anonymous	*010101	*110101/110108	
ANON_21	Anonymous	*07010101/07010102	*090102	*07010101/07010102
ANON_22	Anonymous	*040501	*100101	*040503
ANON_23	Anonymous	*090102	*110401	
ANON_24	Anonymous	*07010101/07010102	*1506	
ANON_25	Anonymous	*03010101/03010102/030108	*15010101/15010102	



Data formatting

Step 3: done (column inserted for the second genotype)!

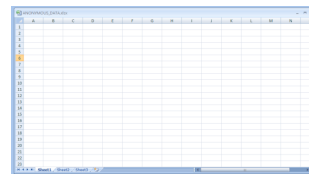
	A	B	C	D	E	F	G	H
	Sample Number	Population	DRB1-allele1		DRB1-allele2		DRB1-allele1	
1								
2	ANON_1	Anonymous	*1404		*1404			
3	ANON_2	Anonymous	*03010101/03010102/030108		*110401			
4	ANON_3	Anonymous	*03010101/03010102/030108		*110401			
5	ANON_4	Anonymous	*15010101/15010102		*150201			
6	ANON_5	Anonymous	*07010101/07010102		*07010101/07010102			
7	ANON_6	Anonymous	*010101		*160201			
8	ANON_7	Anonymous	*010101		*03010101/03010102/030108			
9	ANON_8	Anonymous	*010101		*03010101/03010102/030108			
10	ANON_9	Anonymous	*010101/0107/0122		*07010101/07010102/070102			
11	ANON_10	Anonymous	*010101		*110101/110108			
12	ANON_11	Anonymous	*07010101/07010102		*07010101/07010102			
13	ANON_12	Anonymous	*07010101/07010102		*111201			
14	ANON_13	Anonymous	*07010101/07010102		*160201			
15	ANON_14	Anonymous	*0402		*07010101/07010102			
16	ANON_15	Anonymous	*110401		*15010101/15010102			
17	ANON_17	Anonymous	*07010101/07010102		*130201			
18	ANON_18	Anonymous	*03010101/03010102/030108		*130201			
19	ANON_19	Anonymous	*040301		*160201			
20	ANON_20	Anonymous	*010101		*110101/110108			
21	ANON_21	Anonymous	*07010101/07010102		*090102		*07010101/07010102	
22	ANON_22	Anonymous	*040501		*100101		*040503	
23	ANON_23	Anonymous	*090102		*110401			
24	ANON_24	Anonymous	*07010101/07010102		*1506			
25	ANON_25	Anonymous	*03010101/03010102/030108		*15010101/15010102			



Data formatting

Step 3: fill the first new column with a semi-column (;)

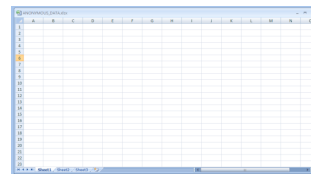
	A	B	C	D	E	F	G	H
1	Sample Number	Population	DRB1-allele1		DRB1-allele2		DRB1-allele1	
2	ANON_1	Anonymous	*1404		*1404			
3	ANON_2	Anonymous	*03010101/03010102/030108		*110401			
4	ANON_3	Anonymous	*03010101/03010102/030108		*110401			
5	ANON_4	Anonymous	*15010101/15010102		*150201			
6	ANON_5	Anonymous	*07010101/07010102		*07010101/07010102			
7	ANON_6	Anonymous	*010101		*160201			
8	ANON_7	Anonymous	*010101		*03010101/03010102/030108			
9	ANON_8	Anonymous	*010101		*03010101/03010102/030108			
10	ANON_9	Anonymous	*010101/0107/0122		*010101/07010102/070102			
11	ANON_10	Anonymous	*010101		*110101/110108			
12	ANON_11	Anonymous	*07010101/07010102		*07010101/07010102			
13	ANON_12	Anonymous	*07010101/07010102		*111201			
14	ANON_13	Anonymous	*07010101/07010102		*160201			
15	ANON_14	Anonymous	*0402		*07010101/07010102			
16	ANON_15	Anonymous	*110401		*15010101/15010102			
17	ANON_17	Anonymous	*07010101/07010102		*130201			
18	ANON_18	Anonymous	*03010101/03010102/030108		*130201			
19	ANON_19	Anonymous	*040301		*160201			
20	ANON_20	Anonymous	*010101		*110101/110108			
21	ANON_21	Anonymous	*07010101/07010102		*090102		*07010101/07010102	
22	ANON_22	Anonymous	*040501		*100101		*040503	
23	ANON_23	Anonymous	*090102		*110401			
24	ANON_24	Anonymous	*07010101/07010102		*1506			
25	ANON_25	Anonymous	*03010101/03010102/030108		*15010101/15010102			



Data formatting

Step 3: done (first new column filled with a semi-column) !

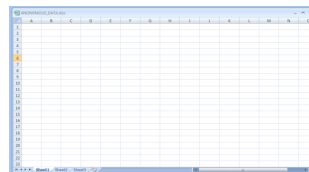
	A	B	C	D	E	F	G	H
1	Sample Number	Population	DRB1-allele1		DRB1-allele2		DRB1-allele1	
2	ANON_1	Anonymous	*1404		*1404			
3	ANON_2	Anonymous	*03010101/03010102/030108		*110401			
4	ANON_3	Anonymous	*03010101/03010102/030108		*110401			
5	ANON_4	Anonymous	*15010101/15010102		*150201			
6	ANON_5	Anonymous	*07010101/07010102		*07010101/07010102			
7	ANON_6	Anonymous	*010101		*160201			
8	ANON_7	Anonymous	*010101		*03010101/03010102/030108			
9	ANON_8	Anonymous	*010101		*03010101/03010102/030108			
10	ANON_9	Anonymous	*010101/0107/0122		*07010101/07010102/070102			
11	ANON_10	Anonymous	*010101		*110101/110108			
12	ANON_11	Anonymous	*07010101/07010102		*07010101/07010102			
13	ANON_12	Anonymous	*07010101/07010102		*111201			
14	ANON_13	Anonymous	*07010101/07010102		*160201			
15	ANON_14	Anonymous	*0402		*07010101/07010102			
16	ANON_15	Anonymous	*110401		*15010101/15010102			
17	ANON_17	Anonymous	*07010101/07010102		*130201			
18	ANON_18	Anonymous	*03010101/03010102/030108		*130201			
19	ANON_19	Anonymous	*040301		*160201			
20	ANON_20	Anonymous	*010101		*110101/110108			
21	ANON_21	Anonymous	*07010101/07010102		*090102		*07010101/07010102	
22	ANON_22	Anonymous	*040501		*100101		*040503	
23	ANON_23	Anonymous	*090102		*110401			
24	ANON_24	Anonymous	*07010101/07010102		*1506			
25	ANON_25	Anonymous	*03010101/03010102/030108		*010101/15010102			



Data formatting

Step 3: fill the second new column with a semi-column (;)

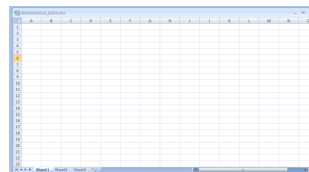
	A	B	C	D	E	F	G	H
1	Sample Number	Population	DRB1-allele1		DRB1-allele2		DRB1-allele1	
2	ANON_1	Anonymous	*1404	:	*1404			
3	ANON_2	Anonymous	*03010101/03010102/030108	:	*110401			
4	ANON_3	Anonymous	*03010101/03010102/030108	:	*110401			
5	ANON_4	Anonymous	*15010101/15010102	:	*150201			
6	ANON_5	Anonymous	*07010101/07010102	:	*07010101/07010102			
7	ANON_6	Anonymous	*010101	:	*160201			
8	ANON_7	Anonymous	*010101	:	*03010101/03010102/030108			
9	ANON_8	Anonymous	*010101	:	*03010101/03010102/030108			
10	ANON_9	Anonymous	*010101/0107/0122	:	*07010101/07010102/070102			
11	ANON_10	Anonymous	*010101	:	*110101/110108			
12	ANON_11	Anonymous	*07010101/07010102	:	*07010101/07010102			
13	ANON_12	Anonymous	*07010101/07010102	:	*111201			
14	ANON_13	Anonymous	*07010101/07010102	:	*160201			
15	ANON_14	Anonymous	*0402	:	*07010101/07010102			
16	ANON_15	Anonymous	*110401	:	*15010101/15010102			
17	ANON_17	Anonymous	*07010101/07010102	:	*130201			
18	ANON_18	Anonymous	*03010101/03010102/030108	:	*130201			
19	ANON_19	Anonymous	*040301	:	*160201			
20	ANON_20	Anonymous	*010101	:	*110101/110108			
21	ANON_21	Anonymous	*07010101/07010102	:	*090102		*07010101/07010102	
22	ANON_22	Anonymous	*040501	:	*100101		*040503	
23	ANON_23	Anonymous	*090102	:	*110401			
24	ANON_24	Anonymous	*07010101/07010102	:	*1506			
25	ANON_25	Anonymous	*03010101/03010102/030108	:	*15010101/15010102			



Data formatting

**Step 3: done (second new column filled with a semi-column)!
 The two alleles of each genotype has to be separated in the same way.**

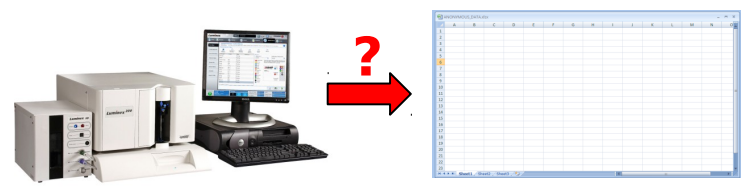
	A	B	C	D	E	F	G	H
1	Sample Number	Population	DRB1-allele1		DRB1-allele2		DRB1-allele1	
2	ANON_1	Anonymous	*1404	:	*1404			
3	ANON_2	Anonymous	*03010101/03010102/030108	:	*110401			
4	ANON_3	Anonymous	*03010101/03010102/030108	:	*110401			
5	ANON_4	Anonymous	*15010101/15010102	:	*150201			
6	ANON_5	Anonymous	*07010101/07010102	:	*07010101/07010102			
7	ANON_6	Anonymous	*010101	:	*160201			
8	ANON_7	Anonymous	*010101	:	*03010101/03010102/030108			
9	ANON_8	Anonymous	*010101	:	*03010101/03010102/030108			
10	ANON_9	Anonymous	*010101/0107/0122	:	*07010101/07010102/070102			
11	ANON_10	Anonymous	*010101	:	*110101/110108			
12	ANON_11	Anonymous	*07010101/07010102	:	*07010101/07010102			
13	ANON_12	Anonymous	*07010101/07010102	:	*111201			
14	ANON_13	Anonymous	*07010101/07010102	:	*160201			
15	ANON_14	Anonymous	*0402	:	*07010101/07010102			
16	ANON_15	Anonymous	*110401	:	*15010101/15010102			
17	ANON_17	Anonymous	*07010101/07010102	:	*130201			
18	ANON_18	Anonymous	*03010101/03010102/030108	:	*130201			
19	ANON_19	Anonymous	*040301	:	*160201			
20	ANON_20	Anonymous	*010101	:	*110101/110108			
21	ANON_21	Anonymous	*07010101/07010102	:	*090102		*07010101/07010102	
22	ANON_22	Anonymous	*040501	:	*100101		*040503	
23	ANON_23	Anonymous	*090102	:	*110401			
24	ANON_24	Anonymous	*07010101/07010102	:	*1506			
25	ANON_25	Anonymous	*03010101/03010102/030108	:	*15010101/15010102			



Data formatting

Step 4: delete useless columns such as to keep only one identifier column (here the column with the population name has to be removed)

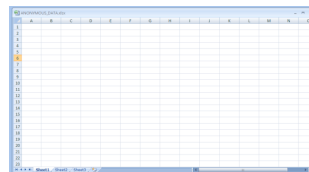
Sample Number	Population	DRB1-allele2	DRB1-allele1
1	ANON_1	*1404	
2	ANON_2	*110401	
3	ANON_3	*110401	
4	ANON_4	*150201	
5	ANON_5	*07010101/07010102	
6	ANON_6	*160201	
7	ANON_7	*03010101/03010102/030108	
8	ANON_8	*03010101/03010102/030108	
9	ANON_9	*07010101/07010102/070102	
10	ANON_10	*010101	*110101/110108
11	ANON_11	*07010101/07010102	*07010101/07010102
12	ANON_12	*07010101/07010102	*111201
13	ANON_13	*07010101/07010102	*160201
14	ANON_14	*0402	*07010101/07010102
15	ANON_15	*110401	*15010101/15010102
16	ANON_17	*07010101/07010102	*130201
17	ANON_18	*03010101/03010102/030108	*130201
18	ANON_19	*040301	*160201
19	ANON_20	*010101	*110101/110108
20	ANON_21	*07010101/07010102	*090102
21	ANON_22	*040501	*100101
22	ANON_23	*090102	*110401
23	ANON_24	*07010101/07010102	*1506
24	ANON_25	*03010101/03010102/030108	*15010101/15010102



Data formatting

Step 4: done (useless column removed)!

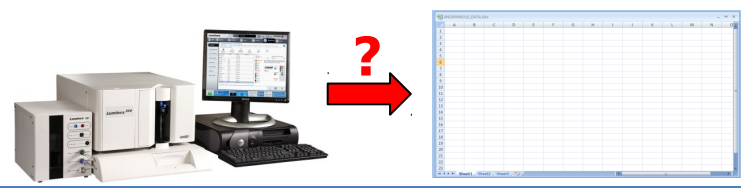
Sample Number	DRB1-allele1	DRB1-allele2	DRB1-allele1	DRB1-allele2
ANON_1	*1404	*1404		
ANON_2	*03010101/03010102/030108	*110401		
ANON_3	*03010101/03010102/030108	*110401		
ANON_4	*15010101/15010102	*150201		
ANON_5	*07010101/07010102	*07010101/07010102		
ANON_6	*010101	*160201		
ANON_7	*010101	*03010101/03010102/030108		
ANON_8	*010101	*03010101/03010102/030108		
ANON_9	*010101/0107/0122	*07010101/07010102/070102		
ANON_10	*010101	*110101/110108		
ANON_11	*07010101/07010102	*07010101/07010102		
ANON_12	*07010101/07010102	*111201		
ANON_13	*07010101/07010102	*160201		
ANON_14	*0402	*07010101/07010102		
ANON_15	*110401	*15010101/15010102		
ANON_17	*07010101/07010102	*130201		
ANON_18	*03010101/03010102/030108	*130201		
ANON_19	*040301	*160201		
ANON_20	*010101	*110101/110108		
ANON_21	*07010101/07010102	*090102	*07010101/07010102	*0906
ANON_22	*040501	*100101	*040503	*100101
ANON_23	*090102	*110401		
ANON_24	*07010101/07010102	*1506		
ANON_25	*03010101/03010102/030108	*15010101/15010102		



Data formatting

Steps 5: save your table as a Text (Tab delimited) file and close it. Here the file has been named « Anonymous.txt ». But a better filename should include both the population name and the locus name, e.g. « Anonymous_DRB1.txt »

Sample Number	DRB1-allele1
ANON_1	*1404
ANON_2	*03010101/03010102/
ANON_3	*03010101/03010102/
ANON_4	*15010101/15010102
ANON_5	*07010101/07010102
ANON_6	*010101
ANON_7	*010101
ANON_8	*010101
ANON_9	*010101/0107/0122
ANON_10	*010101
ANON_11	*07010101/07010102
ANON_12	*07010101/07010102
ANON_13	*07010101/07010102
ANON_14	*0402
ANON_15	*110401
ANON_17	*07010101/07010102
ANON_18	*03010101/03010102/
ANON_19	*040301
ANON_20	*010101
ANON_21	*07010101/07010102
ANON_22	*040501
ANON_23	*090102
ANON_24	*07010101/07010102
ANON_25	*03010101/03010102/030108



Data formatting

Step 7: use the *Replace* command to perform additional changes
- remove tabulations before and after the bars (|)
(remember that the bar can be written by using <Alt Gr> and « / »)

Replace

Find what: \t |

Replace with: |

Conditions: Text Hex

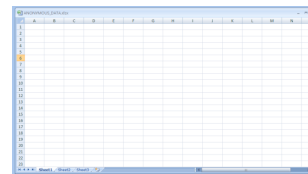
Scope: Active document Selected text All documents

Match whole words Match case Regular expression

Buttons: Find Next, Replace, Replace Next, Replace All, Close, Help

In the Replace box:

- select « regular expression »
- replace **\t|\t** with **|** (**\t** is a tabulation)
- repeat this replacement until no more changes are done

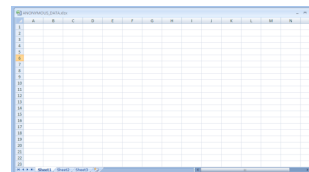


Data formatting

**Step 7: use the *Replace* command to perform additional changes
 - Perform other changes as indicated below (1)**

In the Replace box:

- select « regular expression »
- replace **\t;** with **,**
- repeat this replacement until no more changes are done



Data formatting

Step 7: use the Replace command to perform additional changes
- Perform other changes as indicated below (3)

TextPad - [C:\Documents and Settings\Alicia Sanchez-Mazas\My Documents\Congres&Conferencas\2012_16HIW_Liverpool\Teaching_Session\Anonymous.txt *]

Anonymous.txt *
 Sample Number: DRB1-allele1DRB1-allele2DRB1-allele1DRB1-allele2[DRB1-allele1DRB1-allele2\$
 ANON_1 DRB1*1404,DRB1*1404
 ANON_2 DRB1*03010101/03010102/030108,DRB1*110401
 ANON_3 DRB1*03010101/03010102/030108,DRB1*110401
 ANON_4 DRB1*15010101/15010102,DRB1*150201
 ANON_5 DRB1*07010101/07010102,DRB1*07010101/07010102
 ANON_6 DRB1*010101,DRB1*160201
 ANON_7 DRB1*010101,DRB1*03010101/03010102/030108
 ANON_8 DRB1*010101,DRB1*03010101/03010102/030108
 ANON_9 DRB1*010101/0107/0122,DRB1*07010101/07010102/070102
 ANON_10 DRB1*010101,DRB1*110101/110108
 ANON_11 DRB1*07010101/07010102,DRB1*07010101/07010102
 ANON_12 DRB1*07010101/07010102,DRB1*111201
 ANON_13 DRB1*07010101/07010102,DRB1*160201
 ANON_14 DRB1*0402,DRB1*07010101/07010102
 ANON_15 DRB1*110401,DRB1*15010101/15010102
 ANON_17 DRB1*07010101/07010102,DRB1*130201
 ANON_18 DRB1*03010101/03010102/030108,DRB1*130201
 ANON_19 DRB1*040201,DRB1*160201

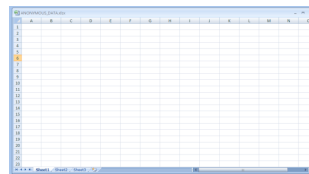
Replace dialog box:
 Find what: *
 Replace with: DRB1*
 Conditions: Text, Hex
 Match whole words, Match case, Regular expression
 Scope: Active document, Selected text, All documents

ANSI Characters
 33 |
 34 ""
 35 #
 36 \$
 37 %
 38 &
 39
 40 (
 41)
 42
 43 +
 44
 45 .
 46
 47 /
 48 0
 49 1
 50 2
 51 3
 52 4

ANON_37 DRB1*080302,DRB1*110401
 ANON_38 DRB1*010101,DRB1*150201
 ANON_39 DRB1*110601,DRB1*120101/1206/1210/1217
 ANON_40 DRB1*07010101/07010102,DRB1*080101/080103
 ANON_41 DRB1*080201,DRB1*130101
 ANON_42 DRB1*010101,DRB1*07010101/07010102
 ANON_43 DRB1*03010101/03010102/030108,DRB1*07010101/07010102
 ANON_44 DRB1*07010101/07010102,DRB1*090102
 ANON_45 DRB1*03010101/03010102/030108,DRB1*160101
 ANON_46 DRB1*010101,DRB1*03010101/03010102/030108
 ANON_47 DRB1*07010101/07010102,DRB1*1116[DRB1*07010101/07010102,DRB1*130101|DRB1*07010101/07010102,DRB1*130104\$

Replaced 158 occurrence(s)

As allele names should start with locus name,
In the Replace box:
- unselect « Regular expression »
- replace * with DRB1* (i.e. the locus you are working with)

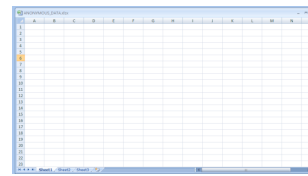


Data formatting

Step 7: use the Replace command to perform additional changes
- Perform other changes as indicated below (4)

- Remove first line to start with genotypes

**Your file is ready!
 Save your file by giving it the extension **.unif** instead of **.txt** and close it**



Data formatting

Gene[RATE]

Main

Tools

- Phenotype
- Haplotype
- Transliterate
- Uniformate**
- File conversions
- Frequency estimation
- One-locus summary
- Regional analysis

Navigation:

- AGP lab**
Laboratory of Anthropology, Genetics and Peopling History
- HLA data analysis**
Tests for HWE, LD, Disease-association
- Generate**
Tools for manipulation of data with ambiguities
- Documentation**
Howtos, manuals, questionnaires and other documents
- AHPD**
Project and documents

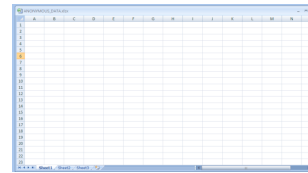
validate UNIFORMAT syntaxe

Check UNIFORMAT files for errors and performs expansion of valid abbreviations (@, allele:allele). Should be used after hand modification or *transliterate* of datafiles.

UNIFORMAT data file:

**Use the Gene[rates] tools at:
<http://geneva.unige.ch/generate>**

**Steps to format your data for Transliterate:
Validate the format with Uniformate**



Data formatting

Gene[RATE]

Main

Tools

- Phenotype
- Haplotype
- Transliterate
- Uniformate
- File conversions
- Frequency estimation
- One-locus summary
- Regional analysis

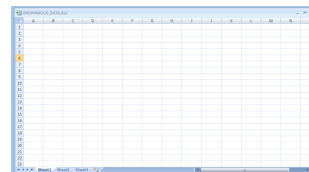
Navigation:

- AGP lab**
Laboratory of Anthropology, Genetics and Peopling History
- HLA data analysis**
Tests for HWE, LD, Disease-association
- Generate**
Tools for manipulation of data with ambiguities
- Documentation**
Howtos, manuals, questionnaires and other documents
- AHPD**
Project and documents

Validated expanded UNIFORMAT data

[Download the results file](#)

**Steps to format your data for Transliterate:
Validate the format with Uniformate: done!**



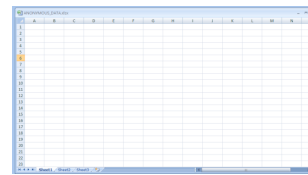
Data formatting

```

ANON_1 DRB1*1404,DRB1*1404
ANON_2 DRB1*03010101/03010102/030108,DRB1*110401
ANON_3 DRB1*03010101/03010102/030108,DRB1*110401
ANON_4 DRB1*15010101/15010102,DRB1*150201
ANON_5 DRB1*07010101/07010102,DRB1*07010101/07010102
ANON_6 DRB1*010101,DRB1*160201
ANON_7 DRB1*03010101/03010102/030108,DRB1*010101
ANON_8 DRB1*03010101/03010102/030108,DRB1*010101
ANON_9 DRB1*010101/0107/0122,DRB1*07010101/07010102/070102
ANON_10 DRB1*010101,DRB1*110101/110108
ANON_11 DRB1*07010101/07010102,DRB1*07010101/07010102
ANON_12 DRB1*07010101/07010102,DRB1*111201
ANON_13 DRB1*07010101/07010102,DRB1*160201
ANON_14 DRB1*07010101/07010102,DRB1*0402
ANON_15 DRB1*110401,DRB1*15010101/15010102
ANON_17 DRB1*07010101/07010102,DRB1*130201
ANON_18 DRB1*03010101/03010102/030108,DRB1*130201
ANON_19 DRB1*160201,DRB1*040301
ANON_20 DRB1*010101,DRB1*110101/110108
ANON_21 DRB1*07010101/07010102,DRB1*090102|DRB1*07010101/07010102,DRB1*0906|DRB1*07010101/07010102,DRB1*0909
ANON_22 DRB1*040501,DRB1*100101|DRB1*100101,DRB1*040503|DRB1*100101,DRB1*040504
ANON_23 DRB1*110401,DRB1*090102
ANON_24 DRB1*07010101/07010102,DRB1*1506
ANON_25 DRB1*03010101/03010102/030108,DRB1*15010101/15010102
ANON_27 DRB1*110401,DRB1*110401
ANON_28 DRB1*110401,DRB1*110401
ANON_29 DRB1*110401,DRB1*110401|DRB1*110401,DRB1*110402
ANON_30 DRB1*07010101/07010102,DRB1*110101/110108
ANON_31 DRB1*010101,DRB1*130201
ANON_32 DRB1*07010101/07010102,DRB1*130101
ANON_33 DRB1*130101,DRB1*040401
ANON_35 DRB1*03010101/03010102/030108,DRB1*0402
ANON_36 DRB1*15010101/15010102,DRB1*07010101/07010102
ANON_37 DRB1*110401,DRB1*080302
ANON_38 DRB1*150201,DRB1*010101
ANON_39 DRB1*110601,DRB1*120101/1206/1210/1217
ANON_40 DRB1*07010101/07010102,DRB1*080101/080103
ANON_41 DRB1*130101,DRB1*080201
ANON_42 DRB1*07010101/07010102,DRB1*010101
ANON_43 DRB1*03010101/03010102/030108,DRB1*07010101/07010102
ANON_44 DRB1*07010101/07010102,DRB1*090102
ANON_45 DRB1*03010101/03010102/030108,DRB1*160101
ANON_46 DRB1*03010101/03010102/030108,DRB1*010101
ANON_47 DRB1*07010101/07010102,DRB1*130101|DRB1*07010101/07010102,DRB1*1116|DRB1*07010101/07010102,DRB1*130104
ANON_48 DRB1*07010101/07010102,DRB1*07010101/07010102

```

**Steps to format your data for Transliterate:
Validate the format with Uniformate: file downloaded**



Data formatting

Gene[RATE]

Main

Tools

- Phenotype
- Haplotype
- Transliterate
- Uniformate
- File conversions
- Frequency estimation
- One-locus summary
- Regional analysis

Navigation:

- AGP lab
Laboratory of Anthropology, Genetics and Peopling History
- HLA data analysis
Tests for HWE, LD, Disease-association
- Generate
Tools for manipulation of data with ambiguities
- Documentation
Howtos, manuals, questionnaires and other documents
- AHPD
Project and documents

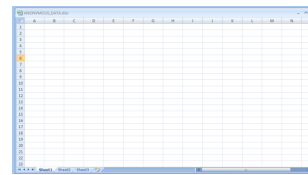
Validated expanded UNIFORMAT data

An error occurred, please try again.

If the error persists [contact us](#).

**Steps to format your data for Transliterate:
Validate the format with Uniformate. If not good, an error
message appears: correct the file.**

**Liverpool 3 June 2012 - Teaching Session: Handling Immunogenetics Data
Prof. Alicia Sanchez-Mazas - Strategies to work with HLA data in human populations: A practical course**



Data formatting

TextPad - [C:\Documents and Settings\Alicia Sanchez-Mazas\My Documents\Congres&Conferences\2012_16IIHW_Liverpool\Teaching_Session\Anonym...

File Edit Search View Tools Macros Configure Window Help

LOCUS 1

```

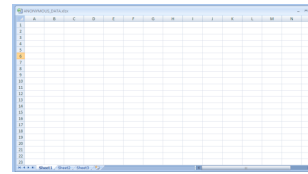
DRB1*010101 - DRB1*0101
DRB1*03010101 - DRB1*0301
DRB1*03010102 - DRB1*0301
DRB1*03010102 - DRB1*0301
DRB1*030108 - DRB1*0301
DRB1*040101 - DRB1*0401
DRB1*040103 - DRB1*0401
DRB1*040301 - DRB1*0403
DRB1*040401 - DRB1*0404
DRB1*040501 - DRB1*0405
DRB1*040503 - DRB1*0405
DRB1*040504 - DRB1*0405
DRB1*07010101 - DRB1*0701
DRB1*07010102 - DRB1*0701
DRB1*070102 - DRB1*0701
DRB1*080101 - DRB1*0801
DRB1*080103 - DRB1*0801
DRB1*080201 - DRB1*0802
DRB1*080302 - DRB1*0803
DRB1*090102 - DRB1*0901
DRB1*100101 - DRB1*1001
DRB1*110101 - DRB1*1101
DRB1*110102 - DRB1*1101
DRB1*110108 - DRB1*1101
DRB1*1103 - DRB1*1103
DRB1*110401 - DRB1*1104
DRB1*110402 - DRB1*1104
DRB1*110601 - DRB1*1106
DRB1*111201 - DRB1*1112
DRB1*120101 - DRB1*1201
DRB1*1206 - DRB1*1206
DRB1*1210 - DRB1*1210
DRB1*1217 - DRB1*1217
DRB1*130101 - DRB1*1301
DRB1*130104 - DRB1*1301
DRB1*130201 - DRB1*1302
DRB1*130301 - DRB1*1303
DRB1*130302 - DRB1*1303
DRB1*140301 - DRB1*1403
DRB1*15010101 - DRB1*1501
DRB1*15010102 - DRB1*1501
DRB1*150201 - DRB1*1502
DRB1*15030101 - DRB1*1503
DRB1*15030102 - DRB1*1503
  
```

ANSI Characters

33	
34	"
35	#
36	\$
37	%
38	&
39	'
40	(
41)
42	*
43	+
44	,
45	.
46	/
47	0
48	1
49	2
50	3
51	4
52	

1 8 Read| Dvr| Block| Sync| Rec| Caps

Prepare a substitution file for Transliterate to reduce high resolution to 2nd level (4 digits) or old nomenclature into new



Data formatting

Gene[RATE]

Main

Tools

- [Phenotype](#)
- [Haplotype](#)
- [Transliterate](#)
- [Uniformate](#)
- [File conversions](#)
- [Frequency estimation](#)
- [One-locus summary](#)
- [Regional analysis](#)

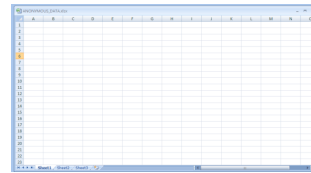
Navigation:

- [AGP lab](#)
Laboratory of Anthropology, Genetics and Peopling History
- [HLA data analysis](#)
Tests for HWE, LD, Disease-association
- [Generate](#)
Tools for manipulation of data with ambiguities
- [Documentation](#)
Howtos, manuals, questionnaires and other documents
- [AHPD](#)
Project and documents

Transliterated file

[Download the results file](#)

Make substitutions to reduce high resolution to 2nd level (4 digits): done!



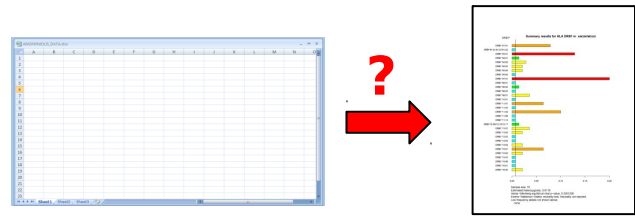
Data formatting

```

ANON_1 DRB1*1404,DRB1*1404
ANON_2 DRB1*03010101/03010102/030108,DRB1*1104
ANON_3 DRB1*03010101/03010102/030108,DRB1*1104
ANON_4 DRB1*15010101/15010102,DRB1*1502
ANON_5 DRB1*07010101/07010102,DRB1*07010101/07010102
ANON_6 DRB1*0101,DRB1*1602
ANON_7 DRB1*03010101/03010102/030108,DRB1*0101
ANON_8 DRB1*03010101/03010102/030108,DRB1*0101
ANON_9 DRB1*010101/0107/0122,DRB1*07010101/07010102/070102
ANON_10 DRB1*0101,DRB1*110101/110108
ANON_11 DRB1*07010101/07010102,DRB1*07010101/07010102
ANON_12 DRB1*07010101/07010102,DRB1*1112
ANON_13 DRB1*07010101/07010102,DRB1*1602
ANON_14 DRB1*07010101/07010102,DRB1*0402
ANON_15 DRB1*1104,DRB1*15010101/15010102
ANON_17 DRB1*07010101/07010102,DRB1*1302
ANON_18 DRB1*03010101/03010102/030108,DRB1*1302
ANON_19 DRB1*1602,DRB1*0403
ANON_20 DRB1*0101,DRB1*110101/110108
ANON_21 DRB1*07010101/07010102,DRB1*0901|DRB1*07010101/07010102,DRB1*0906|DRB1*07010101/07010102,DRB1*0909
ANON_22 DRB1*0405,DRB1*1001
ANON_23 DRB1*1104,DRB1*0901
ANON_24 DRB1*07010101/07010102,DRB1*1506
ANON_25 DRB1*03010101/03010102/030108,DRB1*15010101/15010102
ANON_27 DRB1*1104,DRB1*1104
ANON_28 DRB1*1104,DRB1*1104
ANON_29 DRB1*1104,DRB1*1104
ANON_30 DRB1*07010101/07010102,DRB1*110101/110108
ANON_31 DRB1*0101,DRB1*1302
ANON_32 DRB1*07010101/07010102,DRB1*1301
ANON_33 DRB1*1301,DRB1*0404
ANON_35 DRB1*03010101/03010102/030108,DRB1*0402
ANON_36 DRB1*15010101/15010102,DRB1*07010101/07010102
ANON_37 DRB1*1104,DRB1*0803
ANON_38 DRB1*1502,DRB1*0101
ANON_39 DRB1*1106,DRB1*120101/1206/1210/1217
ANON_40 DRB1*07010101/07010102,DRB1*080101/080103
ANON_41 DRB1*1301,DRB1*0802
ANON_42 DRB1*07010101/07010102,DRB1*0101
ANON_43 DRB1*03010101/03010102/030108,DRB1*07010101/07010102
ANON_44 DRB1*07010101/07010102,DRB1*0901
ANON_45 DRB1*03010101/03010102/030108,DRB1*1601
ANON_46 DRB1*03010101/03010102/030108,DRB1*0101
ANON_47 DRB1*07010101/07010102,DRB1*1301|DRB1*07010101/07010102,DRB1*1116
ANON_48 DRB1*07010101/07010102,DRB1*07010101/07010102
.....

```

Make substitutions to reduce resolution to 2nd level (4 digits): file downloaded.



Estimating frequencies

Results file

Number of individuals sampled

Number of different alleles or haplotypes in sample

Number of solutions with maximum likelihoods obtained to estimate frequencies: should ideally be 1. If greater than 1, may indicate deviation to HWE (to check).

Number of iterations of the algorithm to reach maximum likelihood: ideally low number (less than 10 for one-locus allele frequencies estimation, but greater for ambiguous or multi-locus haplotype frequencies estimation).

Value of the maximum likelihood reached: always negative.

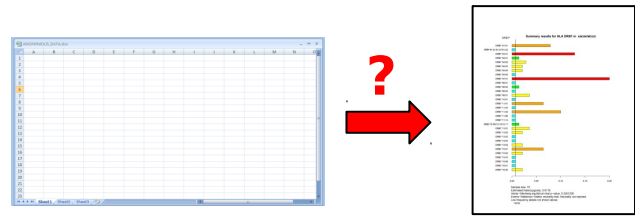
Sum of allele frequencies: should be 1 but may be slightly lower at the 4th and following decimals.

List of alleles with their estimated frequencies.

```

Sample size: 70
Allele/haplotype size: 36
Solution size: 1

num iter      7.0000
max lnL -355.0793
DRB1*01:01    0.0786
DRB1*16:02    0.0214
DRB1*11:01    0.0643
DRB1*14:04    0.0214
DRB1*03:01    0.1286
DRB1*04:03    0.0214
DRB1*11:04    0.1000
DRB1*04:04    0.0214
DRB1*13:01    0.0357
DRB1*15:01    0.0643
DRB1*07:01    0.2000
DRB1*09:01    0.0357
DRB1*04:02    0.0286
DRB1*08:03    0.0071
DRB1*04:01    0.0143
DRB1*14:03    0.0071
DRB1*15:02    0.0214
DRB1*13:02    0.0214
DRB1*11:06    0.0071
DRB1*12:01/12:06/12:10/12:17  0.0143
DRB1*08:02    0.0143
DRB1*08:01    0.0071
DRB1*11:12    0.0071
DRB1*11:03    0.0071
DRB1*16:01    0.0071
DRB1*11:16    0.0071
DRB1*15:03    0.0071
DRB1*13:03    0.0071
DRB1*11:73    0.0071
DRB1*13:49    0.0071
DRB1*15:06    0.0071
DRB1*04:05    0.0071
DRB1*10:01    0.0071
DRB1*09:06    0.0071
DRB1*09:09    0.0071
DRB1*01:01/01:07/01:22  0.0071
sum:          0.9994
    
```



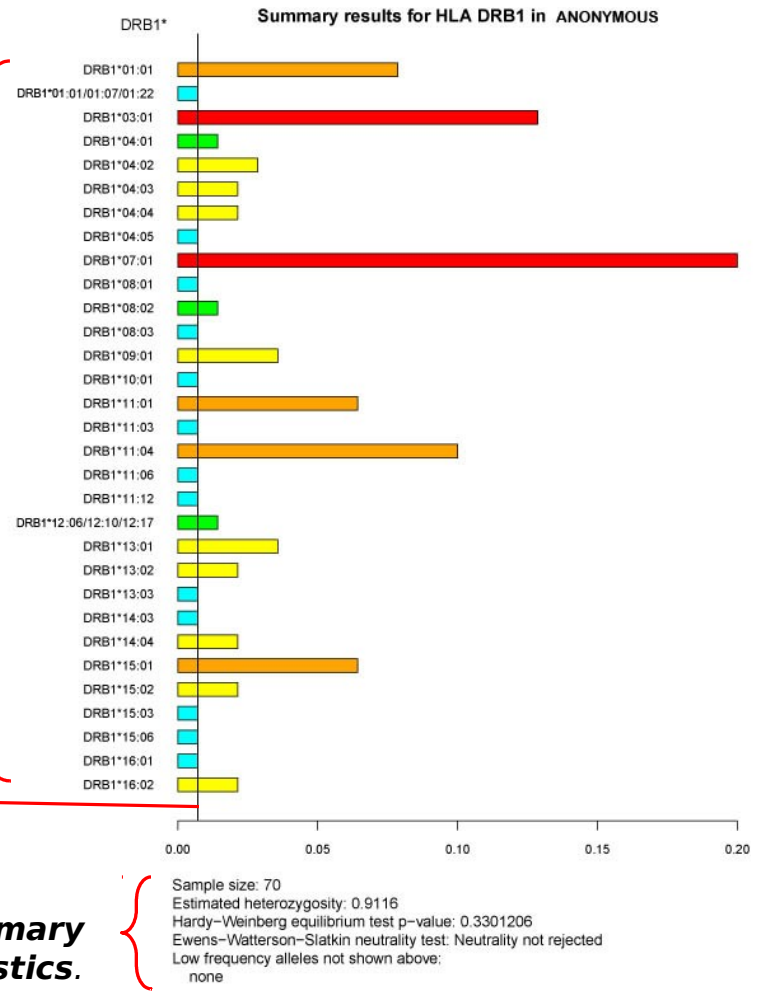
Estimating frequencies

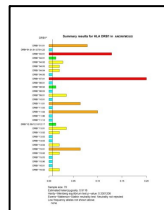
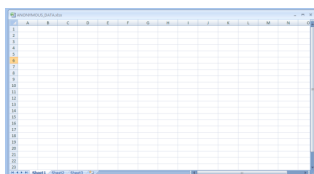
Generated graph

List of alleles with their estimated frequencies.

Frequency threshold: allele frequencies below this line represent less than 1 allele copy in sample (see next slide)

Summary statistics.

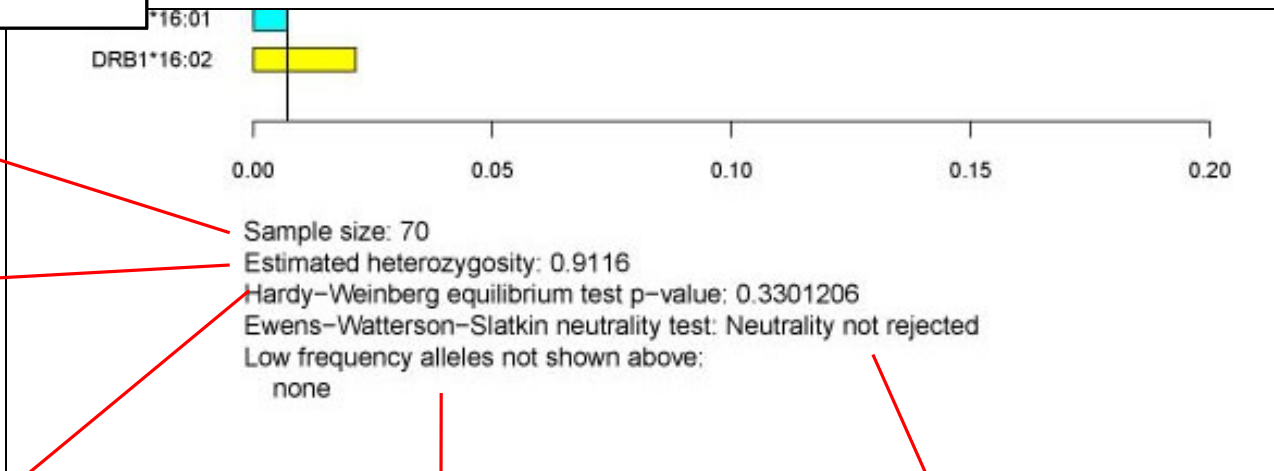




Estimating frequencies

Generated statistics

Summary statistics.



Number of individuals sampled (original data)
Estimated heterozygosity (requires HWE)

Result of Hardy-Weinberg equilibrium test

Low frequency alleles: alleles with frequencies below threshold (of 1 copy in the sample) are listed here. Such alleles may occur with ambiguous or multi-locus haplotypic data.

Result of the neutrality test