## **Project 6: What's in a Name?**

Due to 10/7/2011

**Chong Ding** 

## How to generate the attributes?

I use Ruby to write a script which generates some attributes that might help to find the pattern in the lists. The attributes that can be extracted from names are listed as below:

- 1. Initial character of first name
- 2. Second character of first name
- 3. Third character of first name
- 4. Fourth character of first name
- 5. Fifth character of first name
- 6. Length of full name (ignore middle name)
- 7. Length of first name
- 8. Length of last name
- 9. Ratio of length of first name to length of full name
- 10. Ratio of length of last name to length of full name
- 11. Boolean Variable (if there is a middle name or not)
- 12. Boolean Variable (if first name is shorter than last name)
- 13. Boolean Variable (if there is '-' in the first name)
- 14. Boolean Variable (if there is '-' in the last name)

The simple and messy Ruby code can be found in the same directory with where this report is.

Content in the output file is formatted, so it's very easy to import them into Excel. The reason I use Excel to explore the pattern is that data can be easily sorted and filtered in Excel. It's much straight forward to see the result after sorting or filtering which help to recognize the pattern if it shows.

## How to find the pattern?

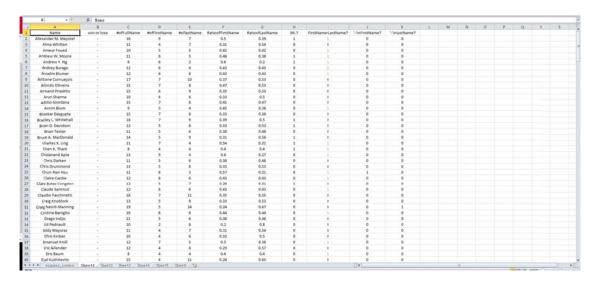
First try: analyze the numeric attributes.

From a snapshot 1, you can see how those numeric attributes are in the Excel. Firstly, I sorted the whole table by following the class attribute "+" or "-". You can find it is split into two groups. Then I find the range of attribute 6 to attribute 10, and count the number different Boolean value in two different groups. Result can be seen from table 1. We didn't see too much variance between each attribute in different groups which are sorted by the class that it belongs to. There are too many overlap between these two groups. Only the distinct attribute is useful to find the pattern or at least the starting. But, some weak pattern we can extract from this numeric attributes is that the areas that is not an overlap between two classes can be used to classify the names (e.g.

LengthofFirstName>9 or RatioofFirstName>0.73 or RatioofLastName<0.2 or if FirstName includes "-" can be used to find some names in group "+"), but accuracy is not high enough.

"+"	Lengthof	Length	Lengtho	RatioofF	Ratioof		FirstName<	'-	
or	FullNam	ofFirst	fLastNa	irstNam	LastNa	MI.?	LastName?	'inFirstN	'inLastN
"-"	e	Name	me	e	me		Lastinaine?	ame?	ame?
"+"	8 to 24	2 to 9	2 to 16	0.125 to	0.2 to	#0: 61	#0:54	#0:84	#0:83
				0.73	0.8	#1: 24	#1: 31	#1: 1	#1: 2
"_"	6 to 20	3 to 11	2 to 15	0.216 to 0.82	0.166 0.7897 74	#0:166 #1: 47	#0:119 #1: 94	#0:207 #1: 6	#0:212 #1: 1

Table 1



Snapshot 1

## Second try:

After trying those numeric attributes, I tried the alphabetic attributes. Firstly, I did the same thing to split all the names with attributes into two groups by their class attribute "+" or "-". Then I count the number of different initial letters appearing in their fist name and last name. The result is show in the table 2. Sub-table (a) is the statistics of first name's initial letter, sub-table (b) is the statistics of last name's initial letter. We can find some distinct from sub-table (a), such as the names whose first name's initial letter is H, V, X or Y belong to group "+", and the majority of the names whose first name's initial letter is D, J, K, L, M, N or W belong to group "+" as well. However, they are not distinct enough and the number of names that start with these letters only 66.7% of the total names in the group "+", which means if this is the pattern, the accuracy is only 66.7%. So we look further to the second letter in the first name.

First Init	loser	winner	
Α	13	3	
В	5	6	
С	13	4	
D	1	25	

Last Ini	loser	winner
Α	3	8
В	7	17
С	6	11
D	6	8

Е	7	1	
F	1	3	
G	2	7	
Н	0	8	
I	1 3 2 2 2 2 0	0	
J	3	29	
K	2	8	
L	2	11	
М	2	33	
N	0	10	
0	2	0	
Р	4	10	
α	0	0	
R	2	18	
S	16	9	
T	6	10	
U	1	0	
V	0	2	
W	1	10	
X Y	0	1 5	
Υ	0	5	
Z	1	0	
(a)			

Е	1	3		
F	3	12		
G	1	18		
Н	1	10		
ļ	1	3		
J	1	7		
K	9	16		
L	3	9		
М	14	15		
N	3	2		
0	1	4		
Р	3	11		
Q	1	0		
R	3	15		
S	7	28		
Т	4	1		
J	1	2		
V	1	2		
W	4	5		
Х	0	0		
Υ	0	3		
Z	1	3		
(b)				

 $\label{eq:Table 2} Table \ 2$  Table 3 shows the statistical result of second letter in first name.

FirstName 2nd	loser	winner
a	0	85
b	0	0
С	1	0
d	3	0
е	0	28
f	1	0
g	0	0
h	19	0
i	1	37
j	9	0
k	0	0
I	9	0
m	2	0
n	5	1

0	0	49
р	0	0
q	0	0
r	22	0
S	1	0
t	11	0
u	0	11
V	1	0
W	0	1
Х	0	0
У	4	0
Z	0	0
	5	0

Table 3

Conclusion: if the  $2^{nd}$  letter of first name is in the set  $\{a, e, i, o, u, w\}$ , it's in the winner's side, otherwise, the name stands on the opposite side.