

Damegender

David Arroyo Menéndez

[2019-08-21]

Outline

1 A tale from commands

2 License

I have a string, I want the sex

All is simple in the beginning

```
$ python3 main.py David
David's gender is male
probability: 1.0
363559  males for David from INE.es
0 females for David from INE.es

$ python3 main.py Isabel
Isabel's gender is female
probability: 1.0
0 males for Isabel from INE.es
271166  females for Isabel from INE.es
```

Perhaps there are non binary probabilities ...

All is possible if one name is found in different countries

```
$ python3 main.py Andrea
Andrea's gender is female
probability: 0.9808615955404946
2084  males for Andrea from INE.es
106807  females for Andrea from INE.es
```

```
$ python3 main.py Alex
Alex's gender is male
probability: 0.9966257742642983
41351  males for Alex from INE.es
140  females for Alex from INE.es
```

My string has different sex in different countries

...

Genderguesser (old sexmachine) did work for us

```
$ python3 nameincountries.py Andrea
```

```
grep -i " Andrea " files/names/nam_dict.txt > files/grep.txt
```

```
males: ['Italy']
```

```
females: ['Albania', 'Austria', 'Belgium', 'Bosnia and Herzegovina']
```

```
both: []
```

```
$ python3 nameincountries.py Alex
```

```
grep -i " Alex " files/names/nam_dict.txt > files/grep.txt
```

```
males: ['Azerbaijan', 'Denmark', 'East Frisia', 'France', 'Germany', 'Greece', 'Hungary', 'Israel', 'Italy', 'Japan', 'Korea', 'Latvia', 'Lithuania', 'Luxembourg', 'Malta', 'Moldova', 'Netherlands', 'Norway', 'Poland', 'Portugal', 'Romania', 'Russia', 'Serbia', 'Slovakia', 'Slovenia', 'Spain', 'Sweden', 'Switzerland', 'Taiwan', 'Turkey', 'Ukraine', 'United Kingdom', 'United States', 'Vietnam']
```

```
females: []
```

```
both: []
```

Now, string is using nicknames ...

We can find a name called "silla". What is the gender of this string?

```
$ python3 main.py silla  
silla gender predicted is female  
0 males for silla from INE.es  
0 females for silla from INE.es
```

The string is not in the dataset. But with damegender we can predict a gender using artificial intelligence. The classification such as with spam is only to reduce time or earn money for humans. It is not exact!!

With this command, we could count males and females in git, mailing lists, etc.

Now, you could count males and females with mails and git:

```
$ python3 mail2gender.py http://mail-archives.apache.org/
```

```
The number of males sending mails is 5
```

```
The number of females sending mails is 1
```

```
$ python3 git2gender.py https://github.com/chaoss/grimoir
```

```
The number of males sending commits is 17
```

```
The number of females sending commits is 13
```

What features in a string is determining the sex?

```
$ python3 infofeatures.py
```

```
-----  
Females with letter/s a: 0.7657420999768214
```

```
Males with letter/s a: 0.6717175543601788  
-----
```

```
Females with last letter a: 0.4705246078961601
```

```
Males with last letter a: 0.16910371997878626  
-----
```

```
Females with last letter o: 0.017306652244456464
```

```
Males with last letter o: 0.10758390787180847  
-----
```

```
Females with last letter consonant: 0.2735841767750908
```

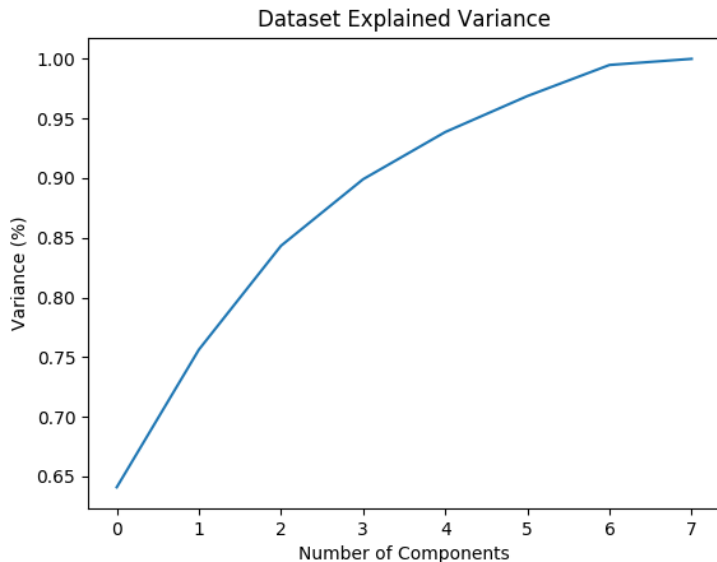
```
Males with last letter consonant: 0.48738540798545343  
-----
```

```
Females with last letter vocal: 0.7262612995441552
```


PCA or not PCA (Principal Component Analysis)

```
$ python3 pca-components.py --csv='files/features_list_no
```

PCA or not PCA (Principal Component Analysis)



PCA or not PCA (Principal Component Analysis)

```
$ python3 pca-features.py --categorical="both" --components=2  
$ firefox files/pca.html &
```

Copyright (C) 2019 David Arroyo Menendez Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.3 or any later version published by the Free Software Foundation; with no Invariant Sections, no Front-Cover Texts, and no Back-Cover Texts. A copy of the license is included in GNU Free Documentation License.