A Detailed Example of Ethnicity Prediction

Le Wang

We would like to predict a voter's race. What is our procedure?

We would like to predict a voter's race. What is our procedure?

Step 1. Figure out the distribution

Pr[Race]

We would like to predict a voter's race. What is our procedure? Step 1. Figure out the distribution

Pr[Race]

$$\begin{split} & \mathsf{Pr}[\mathsf{Race} = \mathsf{white}], \mathsf{Pr}[\mathsf{Race} = \mathsf{black}], \mathsf{Pr}[\mathsf{Race} = \mathsf{hispanic}] \\ & \mathsf{Pr}[\mathsf{Race} = \mathsf{asian}], \mathsf{Pr}[\mathsf{Race} = \mathsf{others}] \end{split}$$

We would like to predict a voter's race. What is our procedure?

Step 1. Figure out the distribution

Pr[Race]

$$\begin{split} & \mathsf{Pr}[\mathsf{Race} = \mathsf{white}], \mathsf{Pr}[\mathsf{Race} = \mathsf{black}], \mathsf{Pr}[\mathsf{Race} = \mathsf{hispanic}] \\ & \mathsf{Pr}[\mathsf{Race} = \mathsf{asian}], \mathsf{Pr}[\mathsf{Race} = \mathsf{others}] \end{split}$$

Step 2. Pick the most likely value, in this case, always white

Now, suppose that I tell you this voter's last name is PIEDRA, what is your prediction for this person's race?

Now, suppose that I tell you this voter's last name is PIEDRA, what is your prediction for this person's race?

Step 1. Figure out the conditional distribution, i.e., the distribution for this subgroup

Pr[Race | Surname]

Now, suppose that I tell you this voter's last name is PIEDRA, what is your prediction for this person's race?

Step 1. Figure out the conditional distribution, i.e., the distribution for this subgroup

Pr[Race | Surname]

In this specific case:

```
\begin{aligned} & \text{Pr}[\text{race} = \text{white} & | \text{surname} = \text{PIEDRA}] \\ & \text{Pr}[\text{race} = \text{black} & | \text{surname} = \text{PIEDRA}] \\ & \text{Pr}[\text{race} = \text{hispanic} & | \text{surname} = \text{PIEDRA}] \\ & \text{Pr}[\text{race} = \text{asian} & | \text{surname} = \text{PIEDRA}] \\ & \text{Pr}[\text{race} = \text{others} & | \text{surname} = \text{PIEDRA}] \end{aligned}
```

Now, suppose that I tell you this voter's last name is PIEDRA, what is your prediction for this person's race?

Step 1. Figure out the conditional distribution, i.e., the distribution for this subgroup

Pr[Race | Surname]

In this specific case:

```
\begin{aligned} & \text{Pr}[\text{race} = \text{white} & | \text{surname} = \text{PIEDRA}] \\ & \text{Pr}[\text{race} = \text{black} & | \text{surname} = \text{PIEDRA}] \\ & \text{Pr}[\text{race} = \text{hispanic} & | \text{surname} = \text{PIEDRA}] \\ & \text{Pr}[\text{race} = \text{asian} & | \text{surname} = \text{PIEDRA}] \\ & \text{Pr}[\text{race} = \text{others} & | \text{surname} = \text{PIEDRA}] \end{aligned}
```

Step 2. Pick the most likely value and use it as our prediction.

How to do it?

head(FLvoters)

```
##
        surname county VTD age gender
                                      race
        PIEDRA
                          58
## 1
                  115
                       66
                                   f white
## 2
         LYNCH
                  115 13 51
                                   m white
       LATHROP
                  115
## 4
                       80 54
                                   m white
## 5
        HUMMEL
                  115
                        8 77
                                   f white
  6 CHRISTISON
                  115
                       55
                          49
                                   m white
## 7
         HOMAN
                  115
                       84
                           77
                                   f white
```

```
subset <- subset(FLvoters, surname == "PIEDRA")</pre>
prop.table(table(subset$race))
```

asian black hispanic native other white

But in practice, if we do not have such information but have access to the census names data, what should we do?

```
subset <- subset(cnames, surname == "PIEDRA")
subset</pre>
```

##

8610

pctothers

0.28

```
## surname count pctwhite pctblack pctapi pctaian pct3
## 8610 PIEDRA 3518 6.71 1.19 0.43 0.14
```

How to evaluate our results?

For the white sample, we know that

```
 \begin{array}{lll} Pr[race = white & | surname = PIEDRA] \\ Pr[race = black & | surname = PIEDRA] \\ Pr[race = hispanic & | surname = PIEDRA] \\ Pr[race = asian & | surname = PIEDRA] \\ Pr[race = others & | surname = PIEDRA] \\ \end{array}
```

 $max(\cdots) = Pr[race = white | surname = PIEDRA]$

```
surname count pctwhite pctblack pctapi pctaian pct
##
## 8610 PIEDRA 3518 6.71 1.19 0.43 0.14
```

max(subset[, c("pctwhite", "pctblack", "pctapi", "pctaian", "pc

pctothers

8610 0.28

[1] 91.39

[1] FALSE

max(subset[, c("pctwhite", "pctblack", "pctapi", "pctaian", "pc

subset

What if I tell you one more piece of information that this voter's last name is PIEDRA living in county 115 and VTD 66

What if I tell you one more piece of information that this voter's last name is PIEDRA living in county 115 and VTD 66

Step 1. Figure out the conditional distribution, i.e., the distribution for this subgroup

Pr[Race | Surname, County, VTD]

What if I tell you one more piece of information that this voter's last name is PIEDRA living in county 115 and VTD 66

Step 1. Figure out the conditional distribution, i.e., the distribution for this subgroup

Pr[Race | Surname, County, VTD]

In this specific case:

```
\label{eq:prescondinger} \begin{array}{lll} \text{Pr}[\text{race} = \text{white} & | \text{surname} = \text{PIEDRA}, \text{county} = 115, \, \text{VTD} = 66] \\ \text{Pr}[\text{race} = \text{black} & | \text{surname} = \text{PIEDRA}, \text{county} = 115, \, \text{VTD} = 66] \\ \text{Pr}[\text{race} = \text{hispanic} & | \text{surname} = \text{PIEDRA}, \text{county} = 115, \, \text{VTD} = 66] \\ \text{Pr}[\text{race} = \text{others} & | \text{surname} = \text{PIEDRA}, \text{county} = 115, \, \text{VTD} = 66] \\ \text{Pr}[\text{race} = \text{others} & | \text{surname} = \text{PIEDRA}, \text{county} = 115, \, \text{VTD} = 66] \\ \end{array}
```

What if I tell you one more piece of information that this voter's last name is PIEDRA living in county 115 and VTD 66

Step 1. Figure out the conditional distribution, i.e., the distribution for this subgroup

Pr[Race | Surname, County, VTD]

In this specific case:

```
\begin{split} & \text{Pr}[\text{race} = \text{white} & | \text{surname} = \text{PIEDRA}, \text{county} = 115, \text{ VTD} = 66] \\ & \text{Pr}[\text{race} = \text{black} & | \text{surname} = \text{PIEDRA}, \text{county} = 115, \text{ VTD} = 66] \\ & \text{Pr}[\text{race} = \text{hispanic} & | \text{surname} = \text{PIEDRA}, \text{county} = 115, \text{ VTD} = 66] \\ & \text{Pr}[\text{race} = \text{asian} & | \text{surname} = \text{PIEDRA}, \text{county} = 115, \text{ VTD} = 66] \\ & \text{Pr}[\text{race} = \text{others} & | \text{surname} = \text{PIEDRA}, \text{county} = 115, \text{ VTD} = 66] \\ \end{split}
```

Step 2. Pick the most likely value and use it as our prediction.

```
subset <- subset(FLvoters, surname == "PIEDRA" & county ==</pre>
prop.table(table(subset$race))
```

black hispanic native other

white

asian

##

##

What if this dataset is not available? We have access to only census names data and Florida census data?

head(cnames)

head(FLcensus)

```
##
                count pctwhite pctblack pctapi pctaian
      surname
## 1
        SMITH 2376206 73.34267 22.21778 0.399960 0.849915
     JOHNSON 1857160 61.55000 33.80000 0.420000 0.910000
     WILLIAMS 1534042 48.52000 46.72000 0.370000 0.780000
        BROWN 1380145 60.71607 34.54345 0.410041 0.830083
## 4
## 5
        JONES 1362755 57.69000 37.73000 0.350000 0.940000
       MILLER 1127803 85.80142 10.40896 0.419958 0.629937
##
##
     pctothers
     2.479752
## 1
## 2 2.730000
## 3 2.790000
##
     2.690269
## 5 2.790000
## 6
    1.939806
```

Lets piece together information from two different datasets

$$\begin{split} & \text{Pr}[\text{race}|\text{surname}, \text{residence}] \\ &= \frac{\text{Pr}[\text{surname}|\text{race}, \text{residence}] \, \text{Pr}[\text{race}|\text{residence}]}{\text{Pr}[\text{surname}|\text{residence}]} \end{split}$$

Lets piece together information from two different datasets

$$\begin{split} & \mathsf{Pr}[\mathsf{race}|\mathsf{surname},\mathsf{residence}] \\ &= \frac{\mathsf{Pr}[\mathsf{surname}|\mathsf{race},\mathsf{residence}] \, \mathsf{Pr}[\mathsf{race}|\mathsf{residence}]}{\mathsf{Pr}[\mathsf{surname}|\mathsf{residence}]} \end{split}$$

$$Pr[race = white | surname = PEIDRA, county = 115, \ VTD = 66] \\ = \frac{Pr[surname = PIEDRA \mid race = white, county = 115, \ VTD = 66]}{Pr[surname = PIEDRA \mid county = 115, \ VTD = 66]} \\ Term C \\ Term B \\ \cdot Pr[race = white \mid county = 115, \ VTD = 66]$$

Term B

```
subset <- subset(FLcensus, county == 115 & VTD == 66)
subset

## county VTD total.pop white black hispanic
## 8048 115 66 5699 0.7638182 0.06281804 0.1363397
B <- subset[,"white"]</pre>
```

Let me also calculate this for other races

B.black <- subset[,"black"]</pre>

B.hisp <- subset[,"hispanic"]</pre>

B.api <- subset[,"api"]</pre>

B.others <- subset[,"others"]</pre>

```
Pr[race = white|surname= PEIDRA,county=115, VTD=66]

Term A
```

 $= \frac{\overbrace{\text{Pr[surname=PIEDRA | race=white,county=115, VTD=66]}^{\text{Pr[surname=PIEDRA | race=white,county=115, VTD=66]}}_{\text{Term C}}$

 $\overbrace{\text{Pr[race=white}|\text{county}=115,\,\text{VTD}=66]}^{\text{Term C}}$

```
Pr[surname=PIEDRA|county=115, VTD=66]

= Pr[surname=PIEDRA | race=white,county=115, VTD=66] · Pr[race=white | county=115, VTD=66]

+ Pr[surname=PIEDRA | race=black,county=115, VTD=66] · Pr[race=black | county=115, VTD=66]

+ Pr[surname=PIEDRA | race=hispanic,county=115, VTD=66] · Pr[race=hispanic | county=115, VTD=66]
```

+ Pr[surname=PIEDRA | race=asisan,county=115, VTD=66] · Pr[race=asian | county=115, VTD=66] + Pr[surname=PIEDRA | race=others,county=115, VTD=66] · Pr[race=others | county=115, VTD=66]

Term A

This is my **Term A**:

 $Pr[surname = PIEDRA \mid race = white, county = 115, \ VTD = 66]$

Term A

This is my **Term A**:

 $Pr[surname = PIEDRA \mid race = white, county = 115, \ VTD = 66]$

I will make the following conditional independence assumption

Pr[surname=PIEDRA | race=white,county=115, VTD=66] = Pr[surname=PIEDRA | race=white] Or similarly for every other racial category

= Pr[surname=PIEDRA | race=black]

Pr[surname=PIEDRA | race=black,county=115, VTD=66]

In general

Pr[surname | race, residence] = Pr[surname | race]

Does it make sense? Or, when will it be violated?

In general

 $Pr[surname \mid race, residence] = Pr[surname \mid race]$

Does it make sense? Or, when will it be violated?

Conditional independence implies that once we know a voter's race, her residence location does not give us any additional information about her surname.

In general

$$Pr[surname \mid race, residence] = Pr[surname \mid race]$$

Does it make sense? Or, when will it be violated?

Conditional independence implies that once we know a voter's race, her residence location does not give us any additional information about her surname.

There is NO strong geographical concentration of certain surnames in Florida within a racial categroy. This will certainly be violated in the Chinese context.

Pr[surname=PIFDRA race=white county=115 VTD=66]	

= Pr[surname=PIEDRA | race=white]

```
\label{eq:problem}  \begin{aligned} \text{Pr[race} &= \text{white} | \text{surname} = \text{PEIDRA,county} = 115, \ \text{VTD} = 66] \\ \hline \textbf{Term A} \\ \hline \textbf{Pr[surname} = \text{PIEDRA} \mid \text{race} = \text{white,county} = 115, \ \text{VTD} = 66] \end{aligned}
```

· Pr[race=white|county=115, VTD=66]

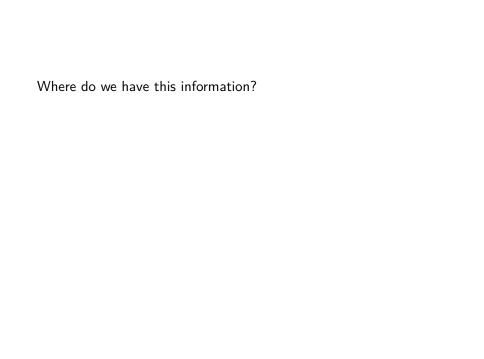
```
Pr[race = white|surname= PEIDRA,county=115, VTD=66]
```

```
Term A
Pr[surname=PIEDRA \mid race=white]
```

Pr[surname=PIEDRA|county=115, VTD=66] Term C

· Pr[race=white|county=115, VTD=66]

Term B





 $Pr[surname = PIEDRA \mid race = white]$

Where do we have this information?

```
Pr[surname=PIEDRA | race=white]
```

Term A.1 Term A.3 $Pr[race=white \mid surname = PIEDRA]$ Pr[race = white]

Term A.2

A.1 <- cnames[which(cnames\$surname=="PIEDRA"), "pctwhite"]

A.1

[1] 0.0671

Pr[surname=PIEDRA race=white]	

```
Pr[surname=PIEDRA | race=white]
```

```
Term A.1
```

Term A.3 $Pr[race=white \mid surname = PIEDRA]$ $\cdot Pr[surname = PIEDRA]$ Pr[race = white]

##

white 0.6045159

head(FLcensus)

```
##
     county VTD total.pop white black hispanic
         1 46
                    4482 0.6137885 0.1871932 0.12092816 0
## 1
## 2
         1 31
                    5470 0.6742230 0.1201097 0.10859232 0
## 3
         1 55
                    3165 0.6315956 0.2085308 0.10110585 0
## 4
         1 49
                    3458 0.7111047 0.1917293 0.05060729 0
## 5
         1 56
                    1937 0.6319050 0.2426433 0.07279298 0
            60
## 6
                    3103 0.4166935 0.4785691 0.05156300 0
race.prop <- apply(FLcensus[,c("white", "black", "api", "h
                   2,
                   weighted.mean,
                   weights = FLCensus$total.pop)
A.2 <- race.prop["white"]
A.2
```

Pr[surname=PIEDRA race=white]	

```
Pr[surname=PIEDRA | race=white]
```

```
Term A.1
```

Term A.3 $Pr[race=white \mid surname = PIEDRA]$ $\cdot Pr[surname = PIEDRA]$ Pr[race = white]

head(cnames)

```
##
                count pctwhite pctblack pctapi pctaian
      surname
## 1
        SMITH 2376206 73.34267 22.21778 0.399960 0.849915
     JOHNSON 1857160 61.55000 33.80000 0.420000 0.910000
  3 WILLIAMS 1534042 48.52000 46.72000 0.370000 0.780000
## 4
       BROWN 1380145 60.71607 34.54345 0.410041 0.830083
## 5
        JONES 1362755 57.69000 37.73000 0.350000 0.940000
## 6
      MILLER 1127803 85.80142 10.40896 0.419958 0.629937
##
     pctothers
## 1
     2.479752
## 2 2.730000
## 3 2.790000
    2.690269
## 4
## 5 2.790000
## 6 1.939806
```

total.count<- sum(cnames\$count)</pre>

A.3 <- cnames[which(cnames\$surname == "PIEDRA"), "count"]/to

```
A <- A.1*A.3/A.2
A
```

```
## white
## 1.612791e-06
```

```
Repeat this process for all other racial groups
A.1.black <- cnames[which(cnames$surname=="PIEDRA"), "pctb]
A.1.hisp <- cnames[which(cnames$surname=="PIEDRA"), "pcthis
A.1.api <- cnames[which(cnames$surname=="PIEDRA"), "pctapi
A.2.black <- race.prop["black"]
A.2.hisp <- race.prop["hispanic"]
A.2.api <- race.prop["api"]
A.2.others <- race.prop["others"]
A.black <- A.1.black*A.3/A.2.black
A.hisp \leftarrow A.1.hisp*A.3/A.2.hisp
```

A.1.others <- cnames[which(cnames\$surname=="PIEDRA"), "pcto

A.black ##

black

A.api \leftarrow A.1.api*A.3/A.2.api

A.others <- A.1.others*A.3/A.2.others

```
Pr[race = white|surname= PEIDRA,county=115, VTD=66]

Term A
```

 $= \frac{\overbrace{\text{Pr[surname=PIEDRA | race=white,county=115, VTD=66]}^{\text{Pr[surname=PIEDRA | race=white,county=115, VTD=66]}}_{\text{Term C}}$

 $\overbrace{\text{Pr[race=white}|\text{county}=115,\,\text{VTD}=66]}^{\text{Term C}}$

Term C

```
Pr[surname=PIEDRA|county=115, VTD=66]

= Pr[surname=PIEDRA | race=white,county=115, VTD=66] · Pr[race=white | county=115, VTD=66]

+ Pr[surname=PIEDRA | race=black,county=115, VTD=66] · Pr[race=black | county=115, VTD=66]

+ Pr[surname=PIEDRA | race=hispanic,county=115, VTD=66] · Pr[race=hispanic | county=115, VTD=66]

+ Pr[surname=PIEDRA | race=asisan,county=115, VTD=66] · Pr[race=asian | county=115, VTD=66]

+ Pr[surname=PIEDRA | race=others,county=115, VTD=66] · Pr[race=others | county=115, VTD=66]
```

residence <- subset(FLcensus, county== 115 & VTD ==66)

residence

county VTD total.pop white black hispanic 115 66 5699 0.7638182 0.06281804 0.136339

##

8048

```
C <- A*residence["white"] + A.black*residence["black"] + A
C</pre>
```

```
## white
## 8048 9.905006e-06
```

A*B/C

```
## white
## 8048 0.1243693
```

```
cond.prob <- c(A*B/C, A.black*B.black/C, A.hisp*B.hisp/C,
names(cond.prob) <- c("white", "black", "hispanic", "api",</pre>
cond.prob
## $white
## [1] 0.1243693
##
## $black
## [1] 0.007865472
##
## $hispanic
## [1] 0.8589315
##
## $api
## [1] 0.005162938
##
## $others
## [1] 0.003670794
```