

HMM para reconocimiento de vocales

18/05/2015

Contents

0.1	Introducción	2
0.2	Objetivo	2
0.3	Especificación del modelo	2
0.4	Datos	3
0.5	Metodología	4
0.6	Paquetes utilizadas	4
0.7	Resultados Español	4
0.8	Resultados Griego	5
0.9	Código	7
0.10	Bibliografía	10

0.1 Introducción

Los modelos de HMM son modelos estadísticos que se utilizan en procesos que se suponen Markovianos donde los estados no son visibles directamente —son latentes— pero las variables influidas por estos estados si lo son. El objetivo consiste en identificar estos estados ocultos para conocer los estados, las transiciones entre los estados y las probabilidades de cada estado sobre las variables a las que influyen.

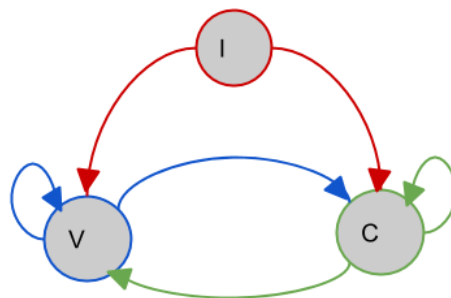
Los HMM han sido ampliamente utilizados en problemas de Procesamiento de Lenguaje Natural (NLP) sobre todo para el reconocimiento del habla (Part Of Speech, POS), etiquetado gramatical, reconocimiento de escritura manual, etc. Stamp (2012)

0.2 Objetivo

Identificación de vocales y consonantes de un texto en español a través de un modelo de HMM estimando sus parámetros.

0.3 Especificación del modelo

- Utilizamos HMM con el algoritmo Baum-Welch para estimar los parámetros:
 1. las probabilidades iniciales de los estados
 2. las probabilidades de transición entre estados
 3. las probabilidades de cada símbolo de pertenecer a uno de los estados



I: Inicio
V: Vocal
C: Consonante

Figure 1: imagen

0.3.1 Algoritmo Baum-Welch

- Este algoritmo es una variante del EM visto en clase Frazzoli (2010); Bilmes (1998) . Iniciamos con un modelo sin ‘conocimiento’

π = probabilidades de iniciar en cada estado

T= matriz de transición de estados

M= matriz de emisiones

$\lambda = (T, M, \pi)$

- En cada iteración los valores de π , A y B se van actualizando hasta convergencia implementando el algoritmo forward-backward.

Definimos $\gamma_k(s) = Pr[X_k = s|Z, \lambda]$, la probabilidad de que el sistema se encuentre en el estado s en el k -ésimo tiempo, dada la secuencia de observaciones Z en el modelo λ (forward-process).

$$\gamma_k(s) = \frac{\alpha_k(s)\beta_k(s)}{Pr[Z|\lambda]} = \alpha_k(s)\beta_k(s)\sum_{s \in \mathcal{X}} \alpha_t(s)$$

Definimos $\xi_k(q, s) = Pr[X_k = q, X_{k+1} = s|Z, \lambda]$, la probabilidad de estar en el estado q al tiempo k y en el estado s en el tiempo $k + 1$ dada la secuencia de observaciones en el modelo actual de HMM (backward-process).

$\xi_k(q, s) = \eta_k \alpha_k(q) T_{q,s} M_{s,z_{k+1}} \beta_{k+1}(s)$, donde η_k es un factor de normalización tal que $\sum_{q,s} \xi_k(q, s) = 1$

Calculando $\gamma_k(s)$ y $\xi_k(q, s)$ podemos actualizar el modelo $\lambda' = (T', M', \pi')$

$$\pi'_s = \gamma_1(s)$$

$$T'_{q,s} = \frac{E[\# \text{ de transiciones del estado } q \text{ al } s]}{E[\# \text{ de transiciones del estado } q]} = \frac{\sum_{k=1}^{t-1} \xi_k(q, s)}{\sum_{k=1}^{t-1} \gamma_k(q)}$$

$$M'_{s,m} = \frac{E[\# \text{ de veces en el estado } s \text{ cuando la observacion fue } m]}{E[\# \text{ de veces en los estados}]} = \frac{\sum_{k=1}^t \gamma_k(s) 1(z_k=m)}{\sum_{k=1}^t \gamma_k(s)}$$

0.3.2 Suposiciones iniciales del modelo

- Nuestra base será suponer que existen 2 estados: **Consonante** y **Vocal**
- No conocemos con qué probabilidad de inicio estamos en Constante o en Vocal
- No conocemos las probabilidades de transición entre estados
- No conocemos las probabilidades de que cada símbolo del lenguaje pertenezca a uno de los estados

0.4 Datos

Intentamos ocupar los últimos 100 contenidos publicados en Quién.com y CNNExpansión.com sin embargo sus corpus requieren de mucho de preprocesamiento para eliminar encoding y caracteres extraños.

Decidimos tomar un corpus en español de un ejercicio realizado en Métodos Analíticos correspondiente a un periodico español con 309,918 noticias

0.4.1 Limpieza de Datos

- Eliminamos signos de puntuación
- Eliminamos dígitos
- Eliminamos tabulaciones
- Cambiamos todo el corpus a minúsculas
- Dividimos cada palabra en sus letras respetando los espacios

0.5 Metodología

1. Limpieza de datos
2. Separar cada palabra en sus respectivas letras respetando espacios
3. Establecemos nuestro conocimiento a priori sobre las probabilidades iniciales de estados inicializando π . Dado que no conocemos mucho del proceso las establecemos muy cercanas a 0.5 (suman a 1)
4. Establecemos nuestro conocimiento a priori sobre las probabilidades de transición entre estados inicializando A. Dado que no conocemos mucho del proceso las establecemos cercanas a 0.5 pero agregando que creemos que es más probable la transición de vocal a consonante que de vocal a vocal, al igual que de consonante a vocal de que de consonante a consonante.
5. Establecemos nuestro conocimiento a priori sobre las probabilidades de cada símbolo a cada uno de los estados propuestos inicializando B. Dado que no conocemos mucho del proceso las establecemos dividiendo 1 entre el número de símbolos posibles en el set de datos.
6. Inicialización de la HMM con los parámetros establecidos en el paso 4,5 y 6.
7. Correr el algoritmo de Baum-Welch

0.6 Paquetes utilizadas

- Paquete HMM de R Himmelmann (2010)
- Algoritmo de Baum-Welch para estimación de parámetros de una HMM

0.7 Resultados Español

Inicial sin conocimiento:

V	C
0.5337	0.4662

Inicial después de Baum-Welch

V	C
0.5337	0.4662

Transiciones sin conocimiento

	V	C
V	0.3099	0.6900

	V	C
C	0.5200	0.4799

Transiciones después de Baum-Welch

	V	C
V	0.3045	0.6954
C	0.993	0.006

A continuación se identifican que las probabilidades en el estado v de los símbolos $a, e, i, o, u, ú$ son mayores que las probabilidades de estos símbolos al estado c . De igual forma, las probabilidades en el estado c para los símbolos que **no** son $a, e, i, o, u, ú$ son mayores que las probabilidades de estos símbolos al estado v .

```
$hmm$emissionProbs
symbols
states      a      b      c      d      e      f      g
v 0.118369315 1.018630e-65 3.344456e-76 0.008194266 0.19487864 1.743170e-66 0.008411338
c 0.005519715 1.172973e-02 7.037839e-02 0.140688529 0.04784931 1.172973e-02 0.011348951
symbols
states      h      i j k      l      m      n      o      p
v 4.732532e-83 5.702808e-02 0 0 4.622025e-80 1.983201e-92 1.489029e-19 1.466436e-01 8.415060e-72
c 1.172973e-02 7.955236e-62 0 0 1.290270e-01 3.518919e-02 1.524865e-01 7.102166e-28 4.691892e-02
symbols
states      q      r      s      t      u      v w x y      z
v 2.196755e-72 0.008237816 0.01800018 9.776189e-02 6.517495e-02 1.212351e-93 0 0 0 6.455586e-96
c 3.518919e-02 0.117166365 0.11484041 7.666003e-07 1.653220e-87 2.345946e-02 0 0 0 1.172973e-02
symbols
states      á é í ó      ú      ñ
v 0.26915304 0 0 0 0 8.146868e-03 2.268708e-92
c 0.01128861 0 0 0 0 3.850728e-194 1.172973e-02
```

Figure 2: resultados

0.8 Resultados Griego

Inicial sin conocimiento:

	V	C
	0.53765	0.4234

Inicial después de Baum-Welch

	V	C
	0.6117	0.3882

Transiciones sin conocimiento

	V	C
V	0.3558	0.6441
C	0.5161	0.4838

Transiciones después de Baum-Welch

	V	C
V	0.0093	0.9906
C	0.6178	0.3821

En griego las vocales corresponden a los símbolos:

Griego	Vocal
A, α	a
E, ϵ	e
I, ι	i
O, o	o
Y, v	u

A continuación se identifican que las probabilidades en el estado v de los símbolos $\alpha, \epsilon, \iota, o, v, \epsilon$ acentuado, son mayores que las probabilidades de estos símbolos al estado c . De igual forma, las probabilidades en el estado c para los símbolos que **no** son $\alpha, \epsilon, \iota, o, v, \epsilon$ acentuado son mayores que las probabilidades de estos símbolos al estado v .

```

Shmm$emissionProbs
  symbols
states   é       v       α       ζ       χ       ρ
v 5.895505e-02 4.252000e-81 2.181337e-01 0.004291154 0.07996329 2.843513e-16 1.162688e-25
c 5.184611e-27 8.876421e-02 3.446021e-12 0.045388583 0.24941473 1.849254e-02 5.917614e-02
  symbols
states   o       π       ó       τ       η       κ       β
v 0.1766922643 2.387271e-41 7.074606e-02 1.352442e-29 9.432809e-02 2.552561e-11 1.393980e-125
c 0.0001084652 5.177912e-02 5.927745e-37 1.405433e-01 1.107997e-26 3.698509e-02 7.397017e-03
  symbols
states   á       λ       ε       γ       ω       ή       σ
v 2.358202e-02 4.927002e-208 1.355966e-01 8.493358e-104 4.126854e-02 2.947753e-02 0.02748251
c 2.758986e-49 1.109553e-02 2.827769e-31 1.849254e-02 2.026343e-41 3.752589e-44 0.03453814
  symbols
states   ώ       ú       φ       υ       ξ       θ       ζ
v 0.011227611 5.895505e-03 0.001745016 2.111235e-05 1.079921e-107 3.686890e-20 1.564313e-75
c 0.004051954 3.528964e-41 0.017397818 5.546438e-02 3.698509e-03 1.479403e-02 7.397017e-03
  symbols
states   ι       ψ       ί       δ       μ
v 2.233928e-07 5.533383e-208 0.006105212 0.014488497 1.270457e-96
c 5.917600e-02 3.698509e-03 0.029456511 0.009403274 3.328658e-02

```

Figure 3: resultados

0.9 Código

0.9.1 Español

```

library(HMM)

periodico <- scan(file='./datos/Es_Newspapers.txt', sep="\n", what = character())

#limpieza de datos
for(i in 1:length(periodico)) {
  periodico[i] <- gsub("[:punct:]", "", unlist(periodico[i]))
  periodico[i] <- gsub("[:digit:]", "", unlist(periodico[i]))
  periodico[i] <- tolower(periodico[i])
  periodico[i] <- gsub("[:space:]", " ", unlist(periodico[i]))
}

#separando a letras
for (i in 1) {
  texto <- periodico[i]
  if(i == 1){
    obsv <- as.data.frame(unlist(strsplit(texto, split="")), stringsAsFactors=FALSE)
  } else {
    obsv <- rbind(obsv, as.data.frame(unlist(strsplit(texto, split="")), stringsAsFactors=FALSE))
  }
}
names(obsv) <- "letter"

states = c("v", "c")

```

```

symbols = c("a","b","c","d","e","f","g","h","i","j","k","l","m","n","o","p","q","r",
            "s","t","u","v","w","x","y","z"," ","á","é","í","ó","ú","ñ")
#probabilidades iniciales de lso estados
random_num <- runif(1,0.5,0.7)
startProbs = matrix(c(random_num, 1-random_num), 2)

random_num <- runif(1, 0.5,0.7)
random_num_2 <- runif(1, 0.5,0.7)
transProbs <- matrix(c(1-random_num, random_num_2,
                      random_num, 1-random_num_2), 2)

random_nums <- data.frame(runif(33,0.030,0.034), runif(33,0.030,0.034))
random_nums <- data.frame(rep(1/33,33), rep(1/33,33))
emissionProbs=as.matrix(random_nums)
inithmm <- initHMM(states, symbols, startProbs=startProbs, transProbs=transProbs,
                  emissionProbs=emissionProbs)
bw <- baumWelch(inithmm, obsv$letter)
bw

```

0.9.2 Griego

```

library(HMM)

setwd("/Users/Lili/Documents/itam/estadistica_multivariada_y_datos_categoricos/examenes/final")

periodico <- scan(file='griego_test.txt', sep="\n", what = character())

for(i in 1:length(periodico)) {
  periodico[i] <- gsub("[:punct:]", "", unlist(periodico[i]))
  periodico[i] <- gsub("[:digit:]", "", unlist(periodico[i]))
  periodico[i] <- tolower(periodico[i])
  periodico[i] <- gsub("[:space:]", " ", unlist(periodico[i]))
}

for (i in 1) {
  texto <- periodico[i]
  if(i == 1){
    obsv <- as.data.frame(unlist(strsplit(texto, split="")), stringsAsFactors=FALSE)
  } else {
    obsv <- rbind(obsv, as.data.frame(unlist(strsplit(texto, split="")), stringsAsFactors=FALSE))
  }
}
names(obsv) <- "letter"

#obsv <- obsv[-which(obsv$letter=="%"),]
#obsv <- obsv[-which(obsv == "â")]

```



```
states = c("v","c")
symbols = c(" ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ", " ")
#probabilidades iniciales de lso estados
random_num <- runif(1,0.5,0.7)
startProbs = matrix(c(random_num, 1-random_num), 2)

random_num <- runif(1, 0.5,0.7)
random_num_2 <- runif(1, 0.5,0.7)
transProbs <- matrix(c(1-random_num,random_num_2,
                      random_num,1-random_num_2), 2)

#random_nums <- data.frame(runif(31,0.030,0.034), runif(33,0.030,0.034))
random_nums <- data.frame(rep(1/length(symbols),length(symbols)), rep(1/length(symbols),length(symbols)))
emissionProbs=as.matrix(random_nums)
inithmm <- initHMM(states, symbols, startProbs=startProbs, transProbs=transProbs,
                  emissionProbs=emissionProbs)
bw <- baumWelch(inithmm, obsv$letter)
bw
```

0.10 Bibliografia