

ST7003 Procesamiento Natural del Lenguaje

Lecture 02



Contenido

1. Minería de tópicos / Temas
2. Modelos de lenguaje n-gram
3. LDA
4. Categorización de texto



Contenido

1. Minería de tópicos / Temas
2. Modelos de lenguaje n-gram
3. LDA
4. Categorización de texto



1. Minería de tópicos / temas

Primeras intuiciones a la generación de documentos y contexto



Minería de tópicos

- Minería de conocimiento sobre el lenguaje, descubrimiento de asociaciones entre palabras (paradigmáticas o sintagmáticas en la visión clásica, o de contexto en la visión moderna).
- Descubrimiento de conocimiento sobre los temas principales del texto -> minería y análisis.
- Que es un tema?
 - Es la idea principal que se analiza en los datos de texto.
 - De lo que está 'hablando' un texto.
 - Hay diferentes granularidades: un conjunto de documentos, a nivel de docto, párrafo o sentencia.
- Ejemplo:
 - De que hablan los usuarios de X hoy en día en Colombia?. Hablan de deportes? De política? De economía? Salud?, etc..
 - O mejor aún, se hablan de diferentes temas, pero en que proporción?
 - También nos interesa saber que le gusta y no le gusta de un producto a la gente.
 - O saber cuales fueron los temas principales, durante el discurso de posesión de Donald Trump.



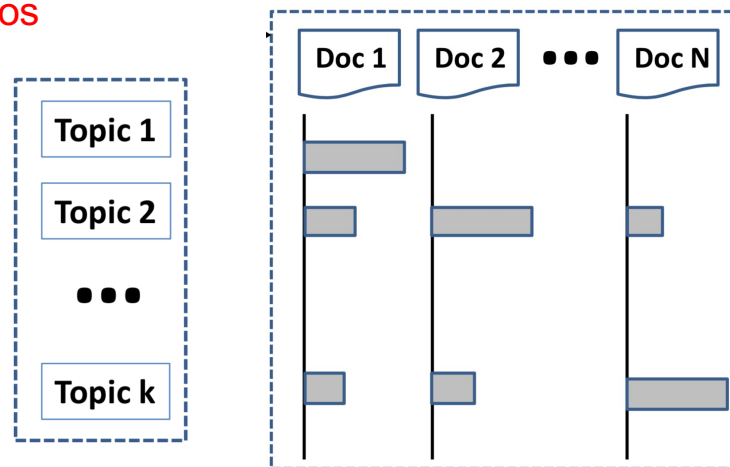
Tópicos / temas

- Podemos ver un tema como un conocimiento sobre el mundo
- A partir de los datos textuales esperamos descubrir una serie de temas
- Estos temas proporcionan una descripción del mundo
- Se puede contar con otros tipos de datos adicionales al texto: ubicación, autores, fuentes, tiempo (serie), etc.
- Cual es la tarea de minería: descubrir muchos temas (k temas), También nos interesa que temas se tratan en un documento y en que proporción.



Tareas de Minería y Análisis de TÓPICOS

Tarea1: descubrir k tópicos



Tarea 2: Determinar qué documentos cubren qué temas



Tareas

- Por ejemplo, en el documento uno, podemos ver que el tema 1 se trata mucho, y el tema 2 y el tema k se tratan en una pequeña parte no estén cubiertos.
- El documento dos, por otro lado, cubrió muy bien el tema 2, pero no cubrió el tema 1 en absoluto, y también cubre el tema k hasta cierto punto, etc.
- Dos tareas:
 1. descubrir k temas de un conjunto de documentos.
 2. en un documento, que temas cubre y en que proporción.
- A continuación la definición formal:



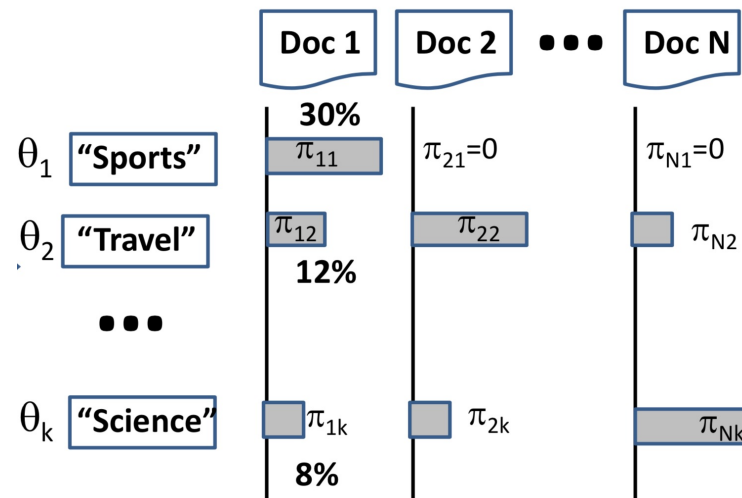
Definición formal de la minería y el análisis de tópicos

- Entrada:
 - Una colección de N documentos $C=\{d_1, \dots, d_N\}$
 - Número de tópicos: k
- Salida:
 - k tópicos: $\{ \theta_1, \dots, \theta_k \}$
 - Cobertura de tópicos en cada d_i : $\{ \pi_{i1}, \dots, \pi_{ik} \}$
 - π_{ij} = probabilidad de d_i de cubrir el tópico j
- Como definir θ_i ?

$$\sum_{j=1}^k \pi_{ij} = 1$$



Idea inicial: Término = Tópico

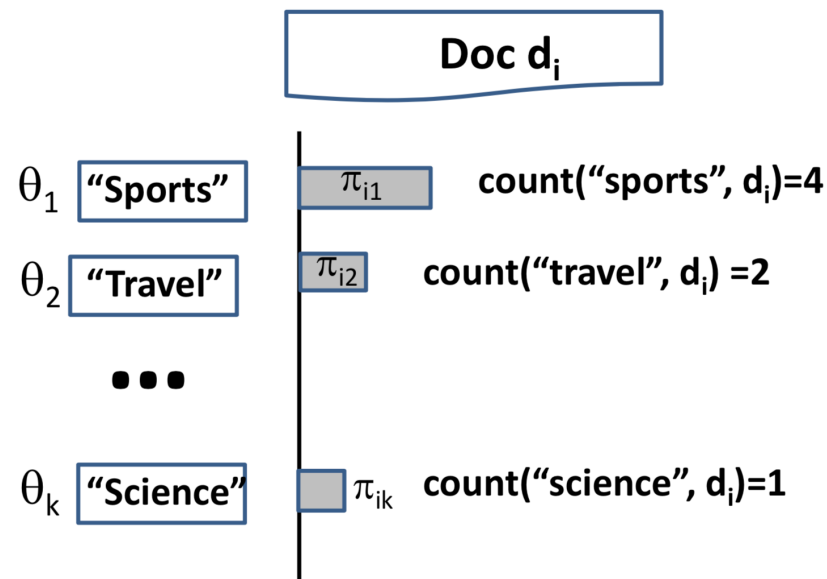


Minería k términos tópicos de la colección C

- Analice el texto en C para obtener términos candidatos (por ejemplo, término = palabra).
- Diseñe una función de puntaje para medir qué tan bueno es cada término como tópico.
 - Favorecer un término representativo (se favorece la alta frecuencia)
 - Evite las palabras que son demasiado frecuentes (por ejemplo, "el", "a").
 - La ponderación TF-IDF de la recuperación puede ser muy útil.
 - La heurística específica del dominio es posible (por ejemplo, favorecer palabras de título, hashtags en tweets).
- Escoja los k términos con las puntuaciones más altas, pero trate de minimizar la redundancia.
 - Si los términos múltiples son muy similares o están estrechamente relacionados, escoja sólo uno de ellos e ignore a los demás.



Calcular la cobertura de tópicos: π_{ij}



$$\pi_{ij} = \frac{\text{count}(\theta_j, d_i)}{\sum_{L=1}^k \text{count}(\theta_L, d_i)}$$



Problemas con "El término como tópico"

- Falta de poder expresivo
 - Solo palabras simples/generales
 - No puede representar tópicos complejos
- Incompletitud en la cobertura del vocabulario
- Ambigüedad en el sentido de la palabra



Idea Mejorada: Tópico = Distribución de palabras

$\Theta_1 = \text{sports}$	$\Theta_2 = \text{travel}$	$\Theta_k = \text{science}$
$P(w \Theta_1)$	$P(w \Theta_2)$	$P(w \Theta_k)$
sports 0.02 game 0.01 basketball 0.005 football 0.004 play 0.003 star 0.003 ... nba 0.001 ... travel 0.0005 ...	travel 0.05 attraction 0.03 trip 0.01 flight 0.004 hotel 0.003 island 0.003 ... culture 0.001 ... play 0.0002 ...	science 0.04 scientist 0.03 spaceship 0.006 telescope 0.004 genomics 0.004 star 0.002 ... genetics 0.001 ... travel 0.00001 ...

$$\sum_{w \in V} p(w | \theta_i) = 1$$

Vocabulario $V = \{w_1, w_2, \dots\}$



Minería y análisis de TÓPICOS probabilísticos

- Entrada:
 - Una colección de N documentos $C=\{d_1, \dots, d_N\}$
 - Conjunto de vocabulario: $V=\{w_1, \dots, w_M\}$
 - Número de tópicos: k
- Salida:
 - k tópicos, cada distribución de palabras: $\{\theta_1, \dots, \theta_k\}$
 - Cobertura de tópicos en cada d_i : $\{\pi_{i1}, \dots, \pi_{ik}\}$
 - π_{ij} = probabilidad de d_i de cubrir el tópico j

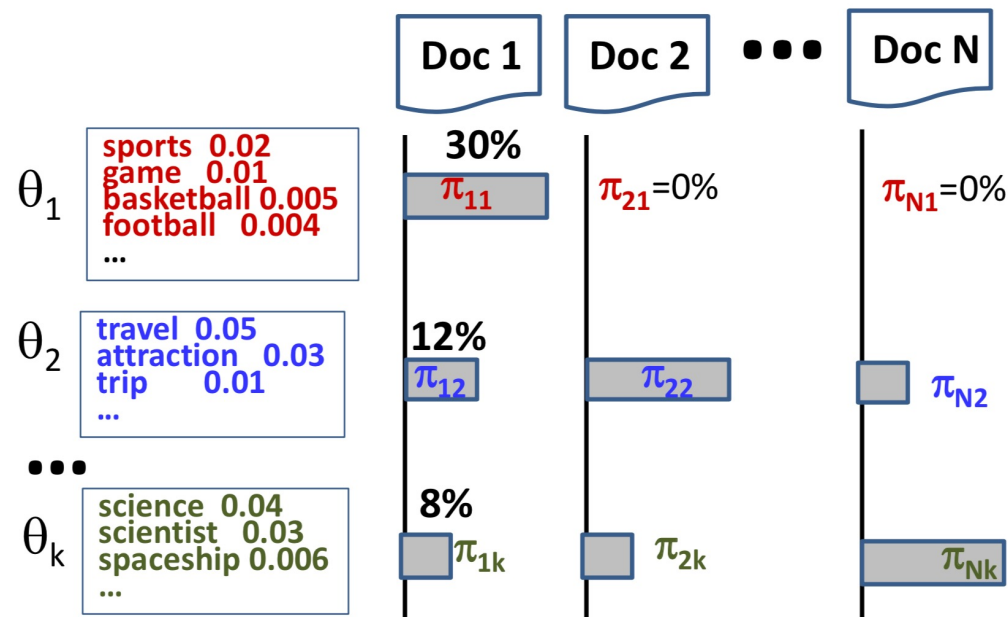
$$\sum_{w \in V} p(w | \theta_i) = 1$$

$$\sum_{j=1}^k \pi_{ij} = 1$$



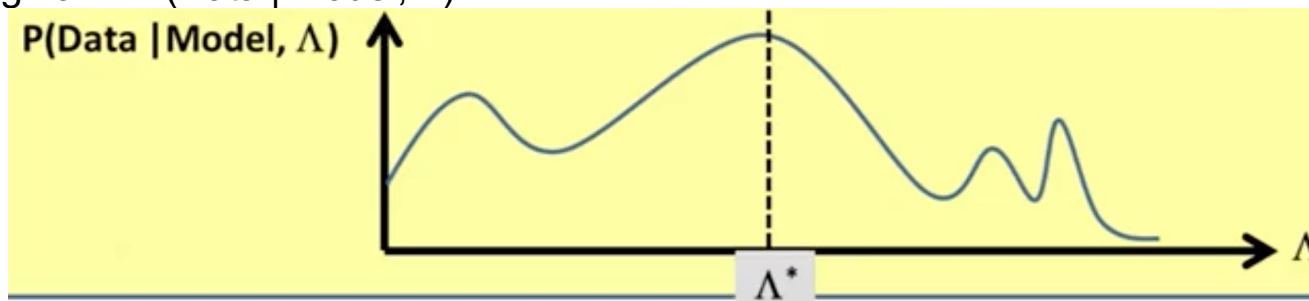
La tarea de cálculo

- Entrada: C, k, V
- Salida: $\{ \theta_1, \dots, \theta_k \}$ y $\{ \pi_{i1}, \dots, \pi_{ik} \}$



Modelo Generativo para minería de texto

- Input: C, k, V Output: $\{\theta_1, \dots, \theta_k\}$ y $\{\pi_{i1}, \dots, \pi_{ik}\}$
- Diseñar un modelo probabilístico
- $P(\text{Data} \mid \text{Model}, \Lambda)$
- $\Lambda = (\{\theta_1, \dots, \theta_k\}, \{\pi_{11}, \dots, \pi_{1k}\}, \dots, \{\pi_{N1}, \dots, \pi_{Nk}\})$
- Cuantos parámetros hay en Total?
 - En $\theta_i \rightarrow |V| \times k$
 - En $\pi_{ij} \rightarrow N \times k$
- Función de optimización: obtener los parámetros hasta que consigamos el conjunto de datos de probabilidad máxima
 - $\Lambda^* = \text{argmax}_{\Lambda} P(\text{Data} \mid \text{Model}, \Lambda)$



Resumen

- Idea mejorada: un topico como una distribución de palabras
- Tareas:
 - Entrada: C, k, V
 - Salida: conjunto de tópicos, cada uno con una distribución de palabras; cobertura de tópicos en cada documento
- $\Lambda = (\{\theta_1, \dots, \theta_k\}, \{\pi_{11}, \dots, \pi_{1k}\}, \dots, \{\pi_{N1}, \dots, \pi_{Nk}\})$

$$\forall j \in [1, k], \sum_{w \in V} p(w | \theta_j) = 1$$

$$\forall i \in [1, N], \sum_{j=1}^k \pi_{ij} = 1$$

retricciones



Resumen

- Modelo generativo
 - Modelo de generación de datos con un modelo de probabilidad:
 - $P(\text{Data} \mid \text{Model}, \Lambda)$
 - Inferir los parámetros Λ^* para un dataset particular:
 - $\Lambda^* = \text{argmax} \Lambda P(\text{Data} \mid \text{Model}, \Lambda)$
 - Λ^* como el 'conocimiento' a ser obtenido por la minería de texto.
 - Ajuste el diseño del modelo para descubrir diferente conocimiento.



Contenido

1. Minería de tópicos / Temas
2. Modelos de lenguaje n-gram
3. LDA
4. Categorización de texto



2. Modelos de lenguaje n-gram



Modelo de lenguaje

- Un modelo de lenguaje es un modelo de aprendizaje que predice palabras.
- Un modelo de lenguaje asigna una probabilidad a cada posible próxima palabra o equivalente dando una probabilidad de distribución sobre posibles próximas palabras.



Prediciendo 'palabras'

- *El agua de Manantiales es hermosamente ...*

azul
gris
clara

*nevera
*que



Modelos de Lenguaje

- Sistemas que pueden predecir palabras adelante
 - Asignar una probabilidad a cada potencial próxima palabra
 - Asignar una probabilidad a una ORACIÓN COMPLETA (muy retador)



Porque predecir palabras?

Varias tareas

- Corrección gramatical

Their are two midterms ~~Their~~ There are two midterms
Everything has improve Everything has ~~improve~~ improved

- Speech recognition

I will be back soonish I will be bassoon dish



Modelo de lenguaje estadístico

- Una probabilidad de distribución sobre una secuencia de palabras
 - $p(\text{'hoy es martes'}) = 0.01$
 - $p(\text{'hoy martes es'}) = 0.0000001$
 - $p(\text{'hoy es banana'}) = 0.0000000001$
- Depende del contexto
- También puede ser visto como un mecanismo probabilístico para generar texto, modelo generativo.



Modelo de lenguaje: Unigram

- Generar texto Generando cada palabra independientemente.
 - $p(w_1, w_2, \dots, w_n) = p(w_1)p(w_2) \dots p(w_n)$
 - Parametros: $\{p(w_i)\} \rightarrow p(w_1) + p(w_2) + \dots + p(w_N) = 1$ (N tamaño vocab.)
-
- Texto \rightarrow muestra extraída de acuerdo con una 'distribución palabras'
 - Por ejemplo, podemos intentar juntar palabras
 - Así podemos calcular la probabilidad de una secuencia de palabras solo con las probabilidades de las palabras
 - $p(\text{'hoy es martes'}) = p(\text{'hoy'})p(\text{'es'})p(\text{'martes'}) =$
 $0.0002 \times 0.001 \times 0.0000015$
- Problema: ignora el orden de las palabras \rightarrow $p(\text{'hoy es martes'})$ es igual a $p(\text{'martes hoy es'})$.
- PERO puede ser usado en ANALISIS DE TÓPICOS, porque?



El modelo de lenguaje
unigram es pieza
fundamental en modelos de
detección de tópicos como
LDA



Generación de texto con LM

Unigram

Unigram LM $p(w|\theta)$

Muestreo

documento d
 $p(d|\theta)=?$

Topico 1:
Text Mining

...
text 0.2
mining 0.1
clustering 0.02
...
food 0.00001
...



Text Mining
paper

Topico 2:
health

...
food 0.25
nutrition 0.1
...
mining 0.0002
...



Heath nutrition
paper



Estimación con LM Unigram

Unigram LM $p(w|\theta)=?$

Estimación

Text Mining paper d
palabras = 100

...	...
10/100	text ?
5/100	mining ?
3/100	clustering ?
...	...

text 10
mining 5
clustering 3

↑
Estimación de máxima
Verosimilitud
(Maximum Likelihood Estimate)

Cual es la mejor estimación?
Como definir 'mejor'?

Heath nutrition
paper



Maximum likelihood vs Bayesian

- Maximum likelihood estimation
 - Problema: muestra pequeña

$$\hat{\theta} = \arg \max_{\theta} P(X | \theta)$$

- Estimacion bayesiana

Bayes Rule

$$p(X|Y) = \frac{p(Y|X)p(X)}{p(Y)}$$

$$\hat{\theta} = \arg \max_{\theta} P(\theta | X) = \arg \max_{\theta} P(X | \theta) P(\theta)$$

- Problema: como definir a priori?



LM setup on unigram

- Input: $C=\{d\}$, V Output: $\{\theta\}$

-datos texto-
“Ipsun larum
Pepe el toro
Hola mundo
...”



θ $P(w|\theta)$
Ipsun ?
larum ?
Pepe ?
el ?
toro ?
Hola ?
mundo ?
...

Doc d

100%



Mining One Topic

- Data: Documento $d = x_1 x_2 \dots x_{|d|}$, x_i pertenece a $V = \{w_1, w_1, \dots, w_M\}$
- Modelo: Unigram LM θ (=topic) : $\{\theta_i = p(w_i | \theta)\}$, $i=1, \dots, M$; $\theta_1 + \dots + \theta_M = 1$
- Función likelihood: $p(d | \theta) = p(x_1 | \theta) p(x_2 | \theta) \dots p(x_{|d|} | \theta)$
 $= p(w_1 | \theta)^{c(w_1, d)} \times \dots \times p(w_M | \theta)^{c(w_M, d)}$
 $= \prod_{i=1}^M p(w_i | \theta)^{c(w_i, d)} = \prod_{i=1}^M \theta_i^{c(w_i, d)}$

- ML estimado:

$$(\hat{\theta}_1, \dots, \hat{\theta}_M) = \arg \max_{\theta_1, \dots, \theta_M} p(d | \theta) = \arg \max_{\theta_1, \dots, \theta_M} \prod_{i=1}^M \theta_i^{c(w_i, d)}$$

Finalmente: $\hat{\theta}_i =$

$$p(w_i | \hat{\theta}) = \frac{c(w_i, d)}{\sum_{i=1}^M c(w_i, d)} = \frac{c(w_i, d)}{|d|}$$

Contador
normalizado



LM setup on unigram

- Input: $C=\{d\}$, V Output: $\{\theta\}$

-datos texto-
“Ipsun larum
Pepe el toro
Hola mundo
...”



θ

$P(w | \theta)$

Ipsun 0.01
larum 0.023
Pepe 0.001
el 0.1
toro 0.02
Hola 0.03
mundo 0.03
...

Doc d

100%



Modelo de Lenguaje (LM) más formalmente

- Meta: CALCULAR LA PROBABILIDAD de una oración (muy retador) o secuencia de palabras W :

$$P(W) = P(w_1, w_2, w_3, w_4, w_5 \dots w_n)$$

- Tarea relacionada: probabilidad de una palabra siguiente:

$$P(w_5 | w_1, w_2, w_3, w_4) \text{ or } P(w_n | w_1, w_2 \dots w_{n-1})$$

- Expresado como:

$$P(W) \text{ or } P(w_n | w_1, w_2 \dots w_{n-1})$$



Como estimar la probabilidad?

- Contar y Dividir

$P(\text{azul} \mid \text{'El agua de Manantiales es hermosamente'}) =$

$C(\text{'El agua de Manantiales es hermosamente azul'})$

$C(\text{'El agua de Manantiales es hermosamente'})$

- No! Demasiadas posibles oraciones.
- No hay suficientes datos para estimar esto.



Calcular $P(W)$ como palabras

- Probabilidad conjunta de $P(W)$

$P(\text{el, agua, de, Mantiales, es, hermosamente, azul})$

- Regla de cadena de probabilidad



Regla de la cadena

- Definición de probabilidades condicionales

$$P(B|A) = P(A,B)/P(A) \quad \rightarrow \quad P(A,B) = P(A) P(B|A) = P(B) P(A|B)$$

- Más variables:

$$P(A,B,C,D) = P(A) P(B|A) P(C|A,B) P(D|A,B,C)$$

- Regla de la cadena en general

$$P(x_1, x_2, x_3, \dots, x_n) = \\ P(x_1)P(x_2|x_1)P(x_3|x_1, x_2) \dots P(x_n|x_1, \dots, x_{n-1})$$



Regla de la cadena aplicado a la probabilidad conjunta de palabras en una oración

$$\begin{aligned} P(w_{1:n}) &= P(w_1)P(w_2|w_1)P(w_3|w_{1:2}) \dots P(w_n|w_{1:n-1}) \\ &= \prod_{k=1}^n P(w_k|w_{1:k-1}) \end{aligned}$$

$P(\text{"El agua de Manantiales"}) =$

$P(\text{El}) \times P(\text{agua}|\text{El}) \times P(\text{de}|\text{El agua}) \times P(\text{Mantiales}|\text{El agua de})$



Todavía problemas

- Probabilidad condicional de oraciones.
- Solución:
 - Supuesto de Markov

$P(\text{azul} \mid \text{El agua de Manantiales es hermosamente})$

$\approx P(\text{azul} \mid \text{hermosamente})$

$$P(w_n \mid w_{1:n-1}) \approx P(w_n \mid w_{n-1})$$



Supuesto de Markov Bigram

$$P(w_{1:n}) \approx \prod_{k=1}^n P(w_k | w_{k-1})$$

- En vez de:

$$\prod_{k=1}^n P(w_k | w_{1:k-1})$$

Queda así: $P(w_n | w_{1:n-1}) \approx P(w_n | w_{n-N+1:n-1})$



Modelo Unigram

$$P(w_1 w_2 \dots w_n) \approx \prod_i P(w_i)$$

Modelo Bigram

$$P(w_i | w_1 w_2 \dots w_{i-1}) \approx P(w_i | w_{i-1})$$



Ejercicios

1. Dado el siguiente modelo de lenguaje θ sobre un vocabulario V compuesto de solo 4 palabras: “the”, “machine”, “learning”, “data”. La distribución de θ esta dada por la siguiente tabla:

w	$P(w \theta)$
machine	0.1
learning	0.2
data	0.3
the	0.4

Calcular y seleccionar la respuesta correcta: $P(\text{“machine learning”}|\theta) =$

- 0.004
- 0.02
- 0.2
- 0.3



Ejercicios

2. Asumiendo el mismo modelo de lenguaje de la pregunta 1:

Calcular y seleccionar la respuesta correcta: $P(\text{"learning machine"}|\theta) =$

- 0.004
- 0.02
- 0.2
- 0.3



Ejercicios

3. Asumiendo el mismo modelo de lenguaje de la pregunta 1:

Calcular y seleccionar la respuesta correcta: $P(\text{"learning machine learning"}|\theta) =$

- 0.004
- 0.02
- 0.2
- 0.3



Ejercicios

4. Asumiendo que las palabras son generadas por una mezcla de 2 modelos de lenguaje θ_1 y θ_2 , donde $P(\theta_1)=0.5$ y $P(\theta_2)=0.5$. Observe la siguiente distribución de palabras:

w	$P(w \theta_1)$	$P(w \theta_2)$
the	0.4	0.15
and	0.4	0.15
genes	0.05	0.3
biology	0.15	0.4

Calcular y seleccionar la respuesta correcta: $P(\text{"biology"}) =$

- 0.55
- 0.15
- 0.175
- 0.275



Ejercicios

5. Asumiendo el mismo modelo de la pregunta 4.

Calcular y seleccionar la respuesta correcta: $P(\text{"the biology"}) =$

- 1
- 0.275
- 0.275625
- 0.075625



Ejercicios

6. Dado el siguiente modelo de lenguaje θ sobre un vocabulario V compuesto de solo 4 palabras: “the”, “global”, “warming”, “effects”. La distribución de θ esta dada por la siguiente tabla:

w	$P(w \theta)$
the	0.3
global	0.2
warming	0.2
effects	X

Calcular y seleccionar la respuesta correcta: $P(\text{“effects”} \mid \theta) =$

- 0.1
- 0.2
- 0.3
- 0



Ejercicios

7. Asumiendo el mismo modelo de lenguaje de la pregunta 6:

w	P(w θ)
the	0.3
global	0.2
warming	0.2
effects	X

Cual de las siguientes sentencias es falsa:

- $P(\text{"the global warming effects"}|\theta) < P(\text{"global warming effects"}|\theta)$
- $P(\text{"global warming"}|\theta) > P(\text{"warming global"}|\theta)$
- $P(\text{"text mining"}|\theta) = 0$
- $P(\text{"global warming"}|\theta) = 0.04$



Ejercicios

8. Asumiendo que las palabras son generadas por una mezcla de 2 modelos de lenguaje θ_1 y θ_2 , donde $P(\theta_1)=0.5$ y $P(\theta_2)=0.5$. Observe la siguiente distribución de palabras :

w	$P(w \theta_1)$	$P(w \theta_2)$
sports	0.35	0.05
basketball	0.2	0.05
fast	0.3	0.3
computer	0.1	0.4
smartphone	0.05	0.2

Cual es la probabilidad de observar “computer” desde el modelo mixto:

- 0.45
- 0.4
- 0.05
- 0.25



Ejercicios

w	P(w θ_1)	P(w θ_2)
sports	0.35	0.05
basketball	0.2	0.05
fast	0.3	0.3
computer	0.1	0.4
smartphone	0.05	0.2

9. Asumiendo el mismo modelo de la pregunta 9. Queremos inferir cual de las 2 distribuciones de palabras, θ_1 y θ_2 , ha sido usada para generar “computer” y quisieramos calcula la probabilidad que este ha sido generado usando θ_1 y θ_2 . Es decir: $P(\theta_1 | \text{“computer”})$ y $P(\theta_2 | \text{“computer”})$, respectivamente, luego los valores de $P(\theta_1 | \text{“computer”})$ y $P(\theta_2 | \text{“computer”})$ son:

Clave: aplicar regla de bayes

- 0.9 y 0.1
- 0.1 y 0.9
- 0.8 y 0.2
- 0.2 y 0.8



Problemas con modelos N-gram

- N-grams no puede manejar **long-distance dependencies**:

“**The soups** that I made from that new cookbook I bought yesterday **were** amazingly delicious.”

- N-grams no trabaja bien modelando nuevas secuencias con significado similar.

Solución: **Large language models**

- Contexto más largo
- Espacio de embeddings en vez de palabras.
- Puede manejar sinonimia mejor



Modelos de lenguaje n-gram

Estimando probabilidades n-gram



Estimando probabilidades bigram

- Estimando Maxima Verosimilitud

$$P(w_n | w_{n-1}) = \frac{C(w_{n-1} w_n)}{\sum_w C(w_{n-1} w)}$$

$$P(w_n | w_{n-1}) = \frac{C(w_{n-1} w_n)}{C(w_{n-1})}$$



Ejemplo

<s> I am Sam </s>

<s> Sam I am </s>

<s> I do not like green eggs and ham </s>

$$P(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

$$P(\text{I} | \text{<s>}) = \frac{2}{3} = .67$$

$$P(\text{Sam} | \text{<s>}) = \frac{1}{3} = .33$$

$$P(\text{am} | \text{I}) = \frac{2}{3} = .67$$

$$P(\text{</s>} | \text{Sam}) = \frac{1}{2} = 0.5$$

$$P(\text{Sam} | \text{am}) = \frac{1}{2} = .5$$

$$P(\text{do} | \text{I}) = \frac{1}{3} = .33$$



Otro ejemplo:
oraciones del Proyecto: Berkeley Restaurant

can you tell me about any good cantonese restaurants close by
tell me about chez panisse
i'm looking for a good place to eat breakfast
when is caffe venezia open during the day



Conteo bigram

- 9222 oraciones

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0



Probabilidades bigram

- unigrams:

i	want	to	eat	chinese	food	lunch	spend
2533	927	2417	746	158	1093	341	278

- Resultado:

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0



Bigram estimates of sentence probabilities

$$\begin{aligned} P(<s> \text{ I want english food } </s>) = \\ & P(\text{I} \mid <s>) \\ & \times P(\text{want} \mid \text{I}) \\ & \times P(\text{english} \mid \text{want}) \\ & \times P(\text{food} \mid \text{english}) \\ & \times P(</s> \mid \text{food}) \\ & = .000031 \end{aligned}$$



Que clase de conocimiento representa n-grams?

- $P(\text{english} \mid \text{want}) = .0011$
- $P(\text{chinese} \mid \text{want}) = .0065$
- $P(\text{to} \mid \text{want}) = .66$
- $P(\text{eat} \mid \text{to}) = .28$
- $P(\text{food} \mid \text{to}) = 0$
- $P(\text{want} \mid \text{spend}) = 0$
- $P(i \mid \langle s \rangle) = .25$



grandes ngrams

- 4-grams, 5-grams
- Large datasets of large n-grams have been released
 - N-grams from Corpus of Contemporary American English (COCA) 1 billion words (Davies 2020)
 - Google Web 5-grams (Franz and Brants 2006) 1 trillion words)
 - Efficiency: quantize probabilities to 4-8 bits instead of 8-byte float

Newest model: infini-grams (∞ -grams) (Liu et al 2024)

- No precomputing! Instead, store 5 trillion words of web text in **suffix arrays**. Can compute n-gram probabilities with any n!



N-gram LM Toolkits

- SRILM
 - <http://www.speech.sri.com/projects/srilm/>
- KenLM
 - <https://kheafield.com/code/kenlm/>
 - (demo e interpretación)



Modelo de Lenguaje

Evaluación y perplejidad



Como evaluar modelos de lenguaje n-gram

• "Evaluación extrínseca (en-vivo)"

Comparar modelos A y B

1. Coloque cada modelo en una tarea real
 - Machine Translation, speech recognition, etc.
2. Ejecutar la tarea, score de A y B
 - Cuantas palabras traduce bien?
 - Cuantas palabras transcribe bien?
 - Compare la precision de A y B



Intrinsic (in-vitro) evaluation

- Evaluación extrínseca no siempre es possible
 - Caro y consumidor de tiempo
 - No siempre generaliza a otras aplicaciones
- Evaluación intrínseca: **perplexity**
 - Directamente mide el rendimiento del modelo de lenguaje prediciendo palabras
 - No necesariamente corresponde con rendimiento de aplicaciones reales
 - Pero nos da una métrica generica simple para modelos de lenguaje
 - Util para large language models (LLMs) así como n-grams



Conjunto de entrenamiento y prueba

Entrenamos nuestro modelo con **training set**.

Evaluamos el modelo con datos no vistos (**test set**)

- El **test set** son datos con los que no fue entrenado el modelo.
 - Intuición: medir la generalización con datos no vistos.
- Una **evaluation metric** (como **perplexity**) nos dice que tan bien esta Nuestro modelo con los datos de prueba.



Contenido

1. Minería de tópicos / Temas
2. Modelos de lenguaje n-gram
3. LDA
4. Categorización de texto

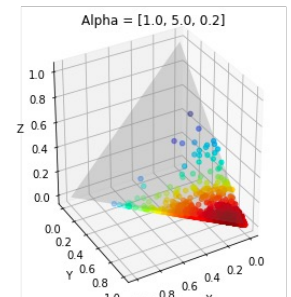
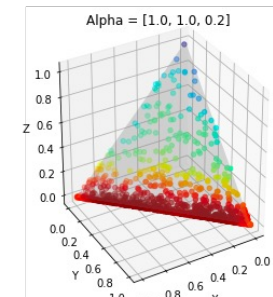
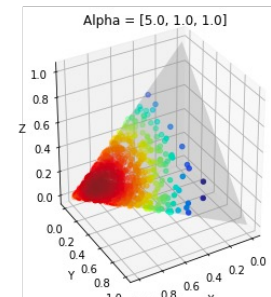
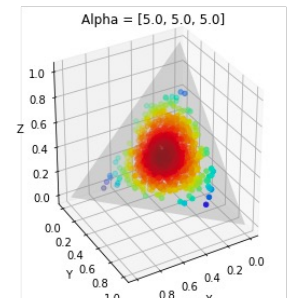
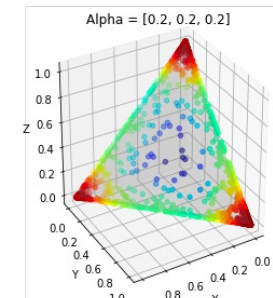
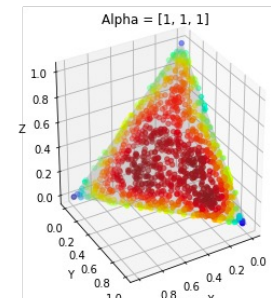
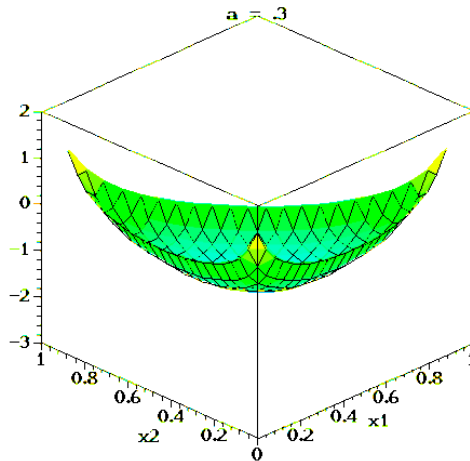


3. Latent Dirichlet Allocation (LDA)



Preliminares: Distribución de Dirichlet

$$f(x_1, \dots, x_K; \alpha_1, \dots, \alpha_K) = \frac{1}{B(\alpha)} \prod_{i=1}^K x_i^{\alpha_i - 1},$$



[https://es.wikipedia.org/wiki/Distribuci%C3%B3n de Dirichlet](https://es.wikipedia.org/wiki/Distribuci%C3%B3n_de_Dirichlet)



Modelos de Temas Probabilísticos

- Los algoritmos de modelado de temas no requieren anotaciones previas ni etiquetado de los documentos
- **Latent Dirichlet Allocation (LDA)**
 - LDA es un modelo estadístico de colecciones de documentos que intenta captar este proceso generativo de intuición
 - Definimos formalmente un tema para ser una distribución sobre un vocabulario fijo
- el objetivo del modelado de temas es descubrir automáticamente los temas de una colección de documentos.
- El problema computacional central para el modelado de temas es usar los documentos observados para inferir la estructura de temas ocultos.



LDA

- LDA es un modelo probabilístico. modelado probabilístico generativo.
 - tratamos los datos como si fueran el resultado de un proceso generativo que incluye variables ocultas.
 - Este proceso generativo define una distribución de probabilidad conjunta entre las variables aleatorias observadas y las ocultas.
 - utilizando esa distribución conjunta para calcular la distribución condicional de las variables ocultas dadas las variables observadas.
 - Esta distribución condicional también se denomina distribución posterior.
- Las variables observadas son las PALABRAS de los documentos;
- Las variables ocultas son la estructura de TÓPICOS.



Distribución de dirichlet

- Una variable aleatoria de Dirichlet k-dimensional puede tomar valores en el (k-1)-simplex, y tiene la siguiente densidad de probabilidad en este simplex:

$$f(x_1, \dots, x_K; \alpha_1, \dots, \alpha_K) = \frac{1}{B(\alpha)} \prod_{i=1}^K x_i^{\alpha_i - 1},$$

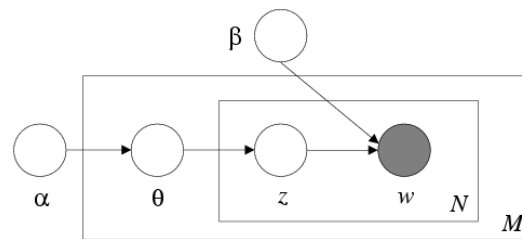
$$\sum_{i=1}^K x_i = 1$$
$$x_1, \dots, x_K > 0$$

$$B(\alpha) = \frac{\prod_{i=1}^K \Gamma(\alpha_i)}{\Gamma\left(\sum_{i=1}^K \alpha_i\right)} = \frac{\Gamma(\alpha_1)\Gamma(\alpha_2)\cdots\Gamma(\alpha_K)}{\Gamma(\alpha_1 + \alpha_2 + \cdots + \alpha_K)}$$

donde $\alpha = (\alpha_1, \dots, \alpha_K)$.



Diagrama LDA



Representación del diagrama LDA. Las cajas son "placas" que representan réplicas.

La placa exterior representa documentos

mientras que la placa interior representa la elección repetida de temas y palabras dentro de un documento.



Probabilidad de un documento

$$P(\mathbf{W}, \mathbf{Z}, \boldsymbol{\theta}, \boldsymbol{\varphi}; \alpha, \beta) = \overset{\text{Temas}}{\prod_{j=1}^M} P(\theta_j; \alpha) \overset{\text{Palabras}}{\prod_{i=1}^K} P(\varphi_i; \beta) \overset{\text{Temas}}{\prod_{t=1}^N} P(Z_{j,t} | \theta_j) \overset{\text{Palabras}}{P(W_{j,t} | \varphi_{Z_{j,t}})}$$

Distribuciones
Dirichlet

Distribuciones
Multinomiales



Un ejemplo

"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.



Contenido

1. Minería de tópicos / Temas
2. Modelos de lenguaje n-gram
3. LDA
4. Categorización de texto



4. Categorización de texto



Clasificación de textos

- La categorización de texto es una tarea de asignar una o más categorías predefinidas al documento analizado, en función de su contenido.
- Sea d_i un documento de un conjunto D y $\{c_1, c_2, c_3, \dots, c_n\}$ un conjunto de todas las etiquetas (labels) a ser asignados. Un Clasificador de Texto asigna uno o más labels c_j a un documento d_i dependiendo del contenido. (single-label vs multi-label)
- Usaremos técnicas de ML supervisadas y no supervisadas, predictivas:
 - Recordar: ML Supervisadas: Datos $\{(x_i, y_i)\}$, N_i (N_i conjuntos de entrenamiento), y_i puede ser una variable categórica, ej: $\{c_1, c_2, \dots, c_n\}$, de esta forma Text Classification va por esta línea. Ej: análisis de sentimiento
 - Recordar: ML no supervisada: Datos $\{x_i\}$ N_i , donde encontramos info y conocimiento de los mismos datos, ej: clustering, detección de tópicos, etc.



Evaluación de Clasificación de textos

- Conjunto aleatorio de los documentos con label (test set), diferente al training set
- Se clasifica el test set con el Clasificador
- Se compara los labels estimados vs labels reales
- La porción de documentos correctamente clasificados es llamado 'accuracy'
 - Las métricas comunes son:
 - Precision, recall, F1-score
 - "la **precisión** es la fracción de las instancias correctas entre las instancias positivas identificadas. **Recall** es el porcentaje de instancias correctas entre todas las instancias positivas. Y F-1 score es la media geométrica de precisión y recuperación".

$$F_1 = 2 \times \frac{precision \times recall}{precision + recall}$$



Técnicas de clasificación



Clasificador Bayesiano

- Es un clasificador NO iterativo de texto, Clasificador Probabilístico
- Basado en el teorema de Bayes
 - Las características (features) son independientes
- El clasificador Naive Bayes esta basado en un modelo de probabilidad condicional.
 - Dado $D = (a_1, a_2, a_3, \dots, a_n)$. \rightarrow a_i nro de características del documento
 - $\Rightarrow p(c_k \mid a_1, a_2, a_3, \dots, a_n)$



Naïve Bayes Fundamentos

- Modelo probabilístico
- Tenemos probabilidades *a priori* basadas en el conjunto de entrenamiento (por eso es supervisada)
- Generamos probabilidades *a posteriori*.
- Formula general:

$$\Pr(y | X) = \frac{\Pr(y) \times \Pr(X | y)}{\Pr(X)}$$

- y son las clases, X los documentos o términos



Ejemplo sencillo

- Tres clases $\{Y\} = \{ \text{'computer science - CS'}, \text{'zoology'}, \text{'entertainment'} \} = \{y_1, y_2, y_3\}$
- Un conjunto de documento $D = \{X_1, X_2, \dots, X_n\}$
- Hago dos queries:
 - Query1: 'python'
 - Query2: 'python download'
- Un caso:
$$P(y=\text{CS} \mid \text{'python'}) = \frac{(P(y=\text{CS}) \times P(\text{'python'} \mid y = \text{CS}))}{P(\text{'python'})}$$



Otros casos:

- $P(y=\text{Zoology} \mid \text{'python'}) = \frac{(P(y=\text{Zoology}) \times P(\text{'python'} \mid y = \text{Zoology}))}{P(\text{'python'})}$
- $P(y=\text{Entert} \mid \text{'python'}) = \frac{(P(y=\text{Entert}) \times P(\text{'python'} \mid y = \text{Entert}))}{P(\text{'python'})}$



Clasificador Naïve Bayes

- $y^* = \text{MAX}(y_j)$ de $P(y_j|X) = \text{MAX}(y_j) P(y_j) \times P(X | y_j)$
- $P(X | y_j) = \prod_{i=1}^n P(x_i | y_j)$
- $y^* = \text{MAX}(y_j)$ de $P(y_j|X) = \text{MAX}(y_j) P(y_j) \times \prod_{i=1}^n P(x_i | y_j)$
- $P(y_j) = n / N$, n #doctos con la clase $y \in N$
- $P(x_i | y_j)$ for $x_i \in X$, $y_j \in Y \rightarrow$ cuantas veces aparece x_i en las instancias con label y_j
- Si hay p instancia de la clase y_j y x_i aparece en k de ellos \rightarrow
 - $P(x_i|y_j) = k / p$
- Problemas:
 - ? $P(x_i|y) = 0$???? Sln \rightarrow función de suavización $(k+1) / (p+n)$. n nro de características. O eliminás la x_i no vista.
 - Desbalanceo de clases?



NB -> 2 propuestas concretas

- Modelo de bernoulli (clasificador binario de clases)
- Modelo Multinomial (pesos en características, tipo tf, tfidf, etc)



Naïve Bayes – síntesis más formal

- El modelo comprende componentes $c_j \in C = \{c_1, c_2, \dots, c_k\}$.
- Cada documento $d_i = \{w_1, w_2, \dots, w_{n_i}\}$ se genera seleccionando primero un componente de acuerdo con los antecedentes, $P(c_j | \theta)$
- luego utiliza el componente para crear el documento de acuerdo con sus propios parámetros, $P(d_i | c_j; \theta)$.
- podemos calcular la probabilidad de un documento utilizando la suma de probabilidades en todos los componentes:

$$P(d_i | \theta) = \sum_{j=1}^k P(c_j | \theta) P(d_i | c_j; \theta)$$



Bayes

- Asumimos una correspondencia uno a uno entre
 - clases $L = \{l_1, l_2, \dots, l_k\}$ y components c_j
- por lo tanto c_j indica tanto el componente j^{th} como la clase c^{th} .
- En consecuencia, dado un conjunto de ejemplos de entrenamiento etiquetados, $D = \{d_1, d_2, \dots, d_{|D|}\}$
- primero aprendemos (estimamos) los parámetros del modelo de clasificación probabilística, θ , y luego usando las estimaciones de estos parámetros, realizamos la clasificación de los documentos de prueba calculando las probabilidades posteriores de cada **clase c_j** , dado el documento de prueba, y seleccionar la clase más probable (clase con la mayor probabilidad).

$$P(c_j | d_i; \hat{\theta}) = \frac{P(c_j | \hat{\theta}) P(d_i | c_j; \hat{\theta}_j)}{P(d_i | \hat{\theta})}$$

$$= \frac{P(c_j | \hat{\theta}) P(w_1, w_2, \dots, w_{n_i} | c_j; \hat{\theta}_j)}{\sum_{c \in C} P(w_1, w_2, \dots, w_{n_i} | c; \hat{\theta}_c) P(c | \hat{\theta})}$$

$$P(w_1, w_2, \dots, w_{n_i} | c_j; \hat{\theta}_j) = \prod_{i=1}^{n_i} P(w_i | c_j; \hat{\theta}_j)$$



Otras técnicas de clasificación de texto

- Árboles de decisión
- KNN
- Máquinas de soporte vectorial
- Random Forest
- Redes neuronales



Arboles de decisión

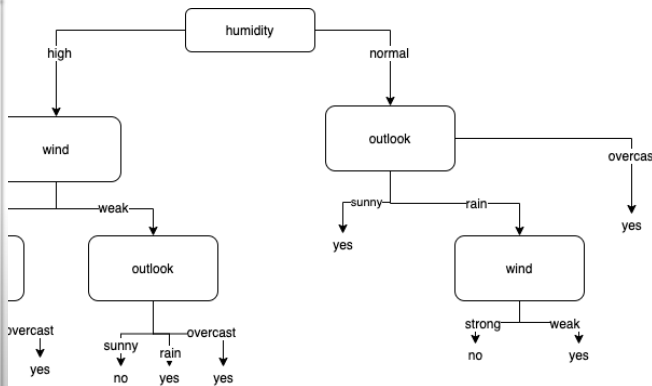
- Técnica supervisada
- El conjunto de entranamiento forma una estructura de arbol
 - Nodos internos denotan una prueba sobre un atributo
 - Una rama representa un resultado del test
 - Las hojas se encuentran las categorías Ci.
 - El arbol genera reglas



Arboles de decisión: un ejemplo

Dataset de entrada

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No



KNN (k nearest neighbour)

- El algoritmo asume que es posible clasificar documentos en un espacio euclidiano como puntos / distancia entre 2 puntos en un plano con coordenadas $p(x,y)$ y $q(a,b)$ calculado como:
 - $d(p,q) = d(q,p) = \sqrt{(x-a)^2 + (y-b)^2}$
 - Dado un documento x , ya instanciado en vsm:
 - $[w1(x), w2(x), w3(x), \dots, wn(x)]$ $w_j(x)$ es el peso del termino j , de acuerdo a diferentes métricas (bit-vector, tf, etc)



referencias

- [ASG] Anubhav Aggarwal, Jasmeet Singh, Dr Kapil Gupta. A Review of Different Text Categorization Techniques, DOI: 10.14419/ijet.v7i3.8.15210
- [R.Jindal] R. Jindal, R. Malhotra, A. Jain (2015), “Techniques for text classification: Literature review and current trends”, Webology, Volume 12, Number 2.

