

Modelo para procesamiento de lenguaje natural?

Victor Hugo Matus Maldonado

21 de marzo de 2020

Resumen ejecutivo

Resumen ejecutivo.

1. Marco teórico y estado del arte

1.1. Bases de datos y álgebra relacional

El *modelo relacional de base de datos*, consiste en cinco componentes:

1. Una colección de tipos escalares, pueden ser definidos por el sistema o por el usuario.
2. Un generador de tipos de relaciones y un intérprete para las relaciones mismas.
3. Estructuras para definir variables relacionales de los tipos generados.
4. Un operador para asignar valores de relación a dichas variables.
5. Una colección relacionalmente completa para obtener valores relacionales de otros valores relacionales mediante operadores.

Las operaciones del modelo relacional están cimentadas en el *álgebra relacional*. Utilizando operaciones primitivas del álgebra se producen nuevas relaciones que pueden manipularse también por medio de operaciones del álgebra mismo. Una secuencia de operaciones de álgebra relacional forma una expresión cuyo resultado es una relación que representa el resultado de una consulta de base de datos. Estas operaciones se pueden clasificar en dos grupos, operaciones de la teoría de conjuntos: *UNIÓN*, *INTERSECCIÓN*, *DIFERENCIA* y *PRODUCTO CARTESIANO* (*PRODUCTO CRUZADO*), y el otro grupo consiste en operaciones específicas para bases de datos relacionales: *JUNTAR*, *SELECCIONAR* y *PROYECTAR*.

1.2. Inteligencia artificial

La *ingeligencia artificial* es

el esfuerzo por automatizar tareas intelectuales normalmente realizadas por humanos (Chollet, 2018)

“es una de las disciplinas más sofisticadas creadas por el ser humano (...) está permitiendo obtener resultados similares a los que observamos en las capacidades de la inteligencia humana: reconocimiento del entorno y percepción espacial, predicción y anticipación, entendimiento del lenguaje y capacidades de comunicación...” recordar cómo citar entrevistas en páginas

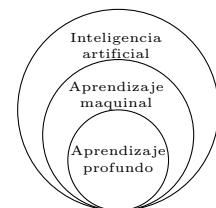


Figura 1:
Inteligencia
artificial

web www.apd.es/ entrevista-gfi-estamos-asistiendo -al-auge-de -agentes-inteligentes-capaces-de-comunicarse -como-una-persona

de este campo general se desprenden el *aprendizaje maquinal* y el *aprendizaje profundo*, figura (1).

Dentro de los múltiples tipos de inteligencia artificial, los que son de nuestro interés para este proyecto se describen a continuación.

Aprendizaje maquinal
Burkov (2019) lo define como

preocupado con construir algoritmos que, para ser útiles dependen de una colección de ejemplos de algún fenómeno (...) el proceso de resolver problemas prácticos por 1) reunir un conjunto de datos y, 2) construir algorítmicamente un modelo estadístico basado en ese conjunto de datos

A diferencia del paradigma clásico de programación, donde los humanos introducen órdenes y datos para ser procesados de acuerdo con dichas reglas, en el aprendizaje maquinal el humano introduce datos y respuestas esperadas de estos datos como ejemplos, el resultado es la generalización de ciertas respuestas a partir de dichos datos sin estructurar. Con ello, se induce al conocimiento por parte de la computadora.

Lingüística computacional
Es un campo multidisciplinario de la lingüística aplicada en la informática. Se sirve de los sistemas informáticos para el estudio y el tratamiento del lenguaje. Para ello, se intenta modelar de manera lógica el lenguaje natural desde un punto de vista programable.

Procesamiento del lenguaje natural
Es una disciplina de la rama de la ingeniería para la lingüística computacional. Se utiliza para la formulación e investigación de mecanismos de eficacia informática para servicios de comunicación entre las personas o entre ellas y las máquinas usando lenguajes naturales. Dos de los módulos básicos de procesamiento natural del lenguaje son búsqueda y aprendizaje con los que se pueden resolver muchos problemas con técnicas de optimización enfocadas en los diferentes parámetros involucrados.

1.3. Mate

Representamos la dependencia entre dos variables, en el que una aumenta o disminuye cuando la otra cambia con la *covarianza*

$$\text{cov}(X, Y) = E[(X - EX)(Y - EY)], \quad (1)$$

al referencia para la covarianza es conveniente escalarla de acuerdo a su desviación estándar, esto recibe el nombre de *coeficiente de correlación*

$$p = \frac{\text{Cov}(X, Y)}{\sigma_1 \sigma_2}. \quad (2)$$

Un modelo que relaciona $E(Y)$ como una función lineal únicamente de β_0 y β_1 , es llamado modelo de regresión lineal *simple*

$$E(Y) = \beta_0 + \beta_1 x, \quad (3)$$

cuando más de una variable independiente es de interés, por ejemplo x_1, x_2, \dots, x_n , se utiliza una generalización de (3), denominada modelo de regresión lineal *múltiple*

$$E(Y) = \beta_0 + \beta_1 x + \dots + \beta_n x_n. \quad (4)$$

La *variable aleatoria* (v.a.) es una función real $X : \Omega \mapsto \mathbb{R}$ tal que el conjunto $\{\omega \in \Omega : X(\omega) \in I\}$ es un evento de Ω para cada $I \subset \mathbb{R}$, en un espacio Ω hipotético. Se le considera *variable aleatoria discreta* (v.a.d.) cuando su rango de valores R_x es finito o contablemente infinito, mientras que una *variable aleatoria continua* (v.a.c.) puede tomar cualquier valor real en un intervalo.

“La forma más natural de expresar la distribución de v.a.d.s es la *función de probabilidad*” (Blitzstein y Hwang, 2019)

Una v.a.d. X con $R_x = \{x_1, x_2, x_3, \dots, x_n, \dots\}$ tiene una función de distribución

$$\begin{aligned} f(x) &= 0 \text{ para cada } x \notin R_x; \\ f(x) &= P(X = x) \text{ para } x \in R_x \end{aligned} \quad (5)$$

para una v.a.c. X será una función no negativa real $f : \mathbb{R} \mapsto [0, \infty)$, es decir

$$P(X \in A) = \int_A f(x) dx \quad (6)$$

El *valor esperado* de una v.a.d. X con una función de probabilidad (5) es definida como

$$\mu = E(X) = \sum_{x \in R_x}^{\infty} x f(x), \quad (7)$$

siempre y cuando la serie converja absolutamente y es también llamado *media* de X , utilizada, similar a la media aritmética en estadísticas, para obtener el valor promedio entre observaciones.

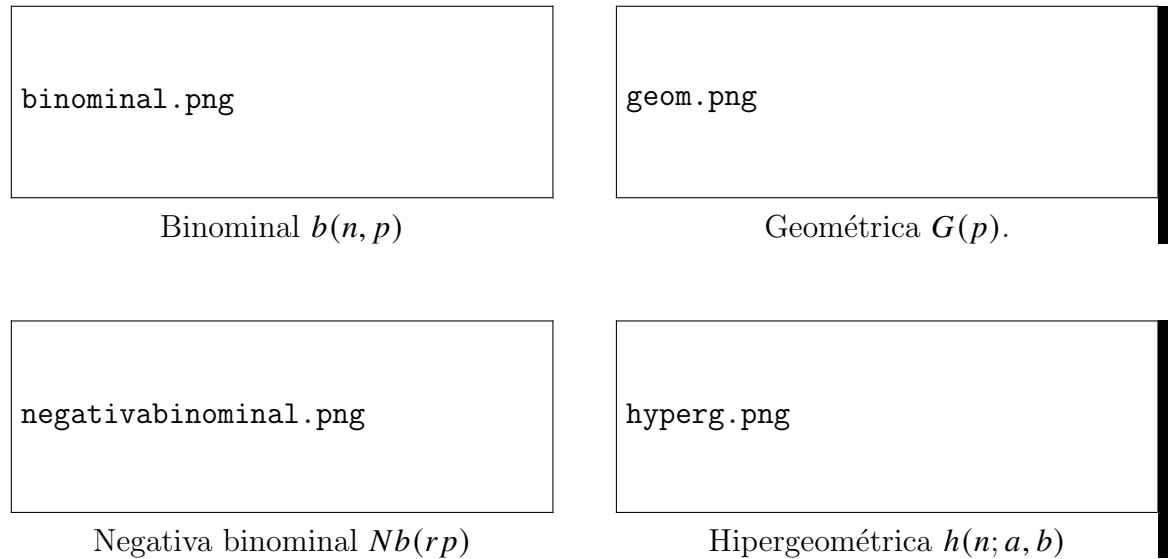
Para una v.a.c. X se define como

$$\mu = E(X) = \int_{-\infty}^{\infty} x f(x) dx \quad (8)$$

Para conocer la variabilidad de la distribución de cualquier v.a se utiliza la *varianza*, para X se define

$$\sigma^2 = Var(X) = [(X - \mu)^2] \quad (9)$$

Figura 2: Diversas de distribuciones pueden para modelar v.a.s, a continuación se muestran las mas importantes de acuerdo a Balakrishnan, Koutras y Politis (2020).



poisson.png	uniforme.png
Poisson $\mathcal{P}(\lambda)$	Uniforme $U[a, b]$
normal.png	epsilon_lambda.png
Normal $N(\mu, \sigma^2)$	Exponencial $Expo(\lambda)$
gama.png	beta.png
Gama $Ga(\alpha, \beta)$	Beta $Be(\alpha, \beta)$

PRUEBAS

Pruebas en las que los conjuntos de hipótesis que contienen β_1 , por ejemplo, $H_a: \beta_1 = 0$ contra $H_a: \beta_1 > 0$ $H_a: \beta_1 < 0$, así como $H_a: \beta_1 > 0$ contra $H_a: \beta_1 \neq 0$ pueden estar basadas en el estadístico

$$t = \frac{\hat{\beta}_1 - 0}{S/\sqrt{S_{xx}}}, \quad (10)$$

Wackerly, Mendenhall III y Scheaffer consideran que “parecería lógico usar r como estadístico de prueba para probar hipótesis más generales acerca de ρ , pero la distribución de probabilidad para r es difícil de obtener.” Wackerly, Mendenhall III y Scheaffer 2009 Sin embargo, en muestras moderadamente grandes podemos probar

la hipótesis $H_0: \rho_1 = \rho_0$ con una prueba Z en la que

$$Z = \frac{(\frac{1}{2}) \ln(\frac{1+r}{1-r}) - (\frac{1}{2}) \ln(\frac{1+\rho}{1-\rho})}{\frac{1}{\sqrt{n-3}}}. \quad (11)$$

Si α es la probabilidad deseada de cometer un error tipo I, la forma de la región de rechazo depende de la hipótesis alternativa. Las diversas alternativas de interés más frecuente y correspondientes regiones de rechazo son las siguientes:

$$\begin{aligned} H_a: \rho > \rho_0, \quad RR: z > z_\alpha \\ H_a: \rho < \rho_0, \quad RR: z < -z_\alpha \\ H_a: \rho \neq \rho_0, \quad RR: |z| > z_\alpha/2 \end{aligned} \quad (12)$$

La suma de los cuadrados del error SSE , es es una alternativa para medir la variación en valores que permanecen sin explicación después de usar las x para ajustar el modelo de regresión lineal simple, la razón SSE/S_{yy} la proporción de la variación total en las y_i que este modelo no explica. El coeficiente de determinación se puede escribir como

$$r^2 = \left(\frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}} \right)^2 = \left(\frac{S_{xy}}{S_{xx}} \right) \left(\frac{S_{xy}}{S_{yy}} \right) = \left(\frac{\hat{\beta}_1 S_{xy}}{S_{yy}} \right) = \frac{S_{yy} - SSE}{S_{yy}} = 1 - \frac{SSE}{S_{yy}}. \quad (13)$$

Podemos interpretar a r^2 como la proporción de la variación total en las y_i que es explicada por una variable x en un modelo de regresión lineal simple.

2. Objetivos

2.1. General

Haciendo uso de las ciencias de la computación, las herramientas matemáticas de estadística y métodos de aprendizaje autónomo, se busca obtener información cuantitativa de textos provenientes de redes sociales, cadenas noticiosas y audio de programas de capacitación, para inferir posiciones, tendencias, comportamientos o razones de grupos sociales, considerando un ciclo de clasificación, estimación, detección y comprobación.

2.2. Particulares

1. Desarrollar un modelo de base de datos que permita la captura de categorías para un determinado problema, los elementos de identificación de cada categoría, el origen de la información y su correlación.

2. Construir una estructura de datos que capte la estimación o valores esperados para el procesamiento de textos.
3. Elaborar un sistema de objetos para el soporte de los elementos de aprendizaje autónomo.
4. Generar los elementos de captura de textos para su almacenamiento y procesamiento.
5. Elaborar un modelo estadístico que permita comprobar las estimaciones a partir de los datos y en consecuencia realizar un ajuste en los parámetros usados para el aprendizaje autónomo.
6. Producir los reportes con un análisis estadístico que faciliten la interpretación de resultados y den pauta para la obtención del conocimiento de interés.

2.3. Metas científicas

Metas científicas

3. Metodología científica

4. Grupo de trabajo

- Dr. José Emilio Quiroz Ibarra
Universidad Iberoamericana, Dirección.
- Dra. Alma Rocío Sagaceta Mejía
Universidad Autónoma Metropolitana, Codirección.
- Mtra. Paloma Alejandra Vilchis León
Universidad Tecnológica de México*, Tutoría.

↑ preguntar

5. Infraestructura disponible para el proyecto

Laboratorios y equipos disponibles en la Universidad Iberoamericana Ciudad de México. servidor desktop

6. Cronograma de actividades

cronograma

7. Resultados comprometidos

Publicación

8. Visto bueno

Vo.Bo.

9. Referencias

- Balakrishnan, N., Markos V. Koutras y Konstantinos G. Politis (2020). *Introduction to probability: models and applications [Introducción a la probabilidad: modelos y aplicaciones]*. Hoboken: John Wiley & Sons, Inc. ISBN: 9781118123348.
- Blitzstein, Joseph K. y Jessica Hwang (2019). *Introduction to Probability [Introducción a la probabilidad]*. 2.^a ed. Texts in Statistical Science [Textos en ciencia estadística]. Florida: CRC Press. ISBN: 9781138369917.
- Burkov, Andriy (2019). *The hundred-page machine learning book*. Quebec: Andriy Burkov. ISBN: 9781999579500.
- Chollet, François (2018). *Machine Learning With Python [Machine Learning]*. New York: Manning Publications Co. ISBN: 9781617294433.
- Wackerly, Dennis D, William Mendenhall III y Richard L Scheaffer (2009). *Estadística matemática con aplicaciones*. Trad. por Jorge Humberto Romo Mufioz. 7.^a ed. Ciudad de México: Cengage Learning. ISBN: 9780495110811.