

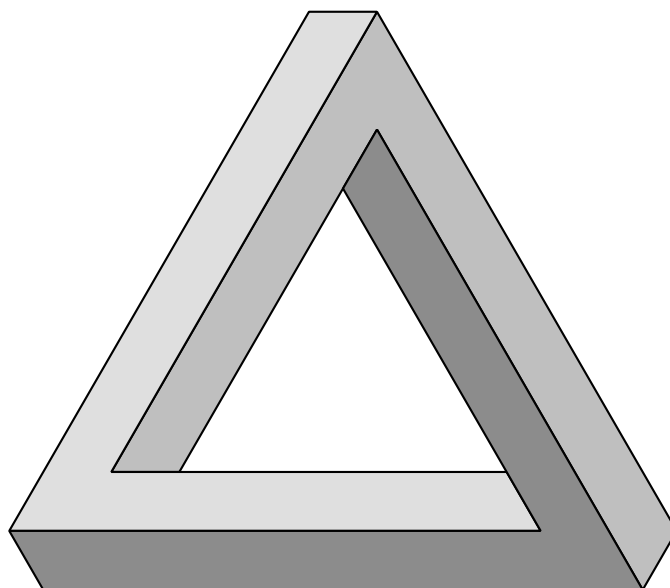
Protocolo

Victor Matus

10 de marzo de 2020

Resumen

Ignoren el formato, también el triángulo.



Índice

1. incompleto	2
2. Marco teórico	3
2.1. Modelo de base de datos relacional	3
2.2. Correlación lineal	4
2.3. Probabilidad condicional	5
2.4. Variables aleatorias y sus distribuciones	7
2.4.1. Función de probabilidad	7
2.4.2. Distribución de Bernoulli y binominal	8
2.4.3. Distribución de hipergeométrica	8
2.4.4. Distribución uniforme discreta	8
2.4.5. Función de distribución acumulada	9
2.5. Valor esperado	9
2.5.1. Linealidad del valor esperado	10
3. Objetivos	10
3.1. General	10
3.2. Particulares	10
3.3. Hipótesis	10
4. Metas	10
5. Metodologías	11
6. Referencias	12

1. incompleto

* Introduction to Probability:

-terminar Expectations

-Continuous rand var

* Linear Models and the Relevant Distributions and Matrix Algebra:

-intervalos confianza, pruebas hipótesis

* Bayesian Reasoning and Machine Learning: - checar este libro

* ML, NLP: Definición, Tipos

*Hipótesis y Objetivos: hipo, gral, part

2. Marco teórico

2.1. Modelo de base de datos relacional

El modelo relacional de base de datos, según Date C. J. Date. *SQL and Relational Theory: How to Write Accurate SQL Code [SQL y teoría relacional: Cómo escribir código SQL correcto]*. 2.^a ed. California: O'Reilly Media, 2012. ISBN: 9781449316402 consiste en cinco componentes:

1. Una colección de tipos escalares, pueden ser definidos por el sistema (INTEGER, CHAR, BOOLEAN, etc.) 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 de operadores relacionales genéricos para derivar obtener valores relacionales de otros valores relacionales.

Es importante comenzar definiendo los tipos, ya que las relaciones se definen sobre ellos, según Date, los tipos son “en esencia un conjunto finito de valores nombrado-todos los valores posibles de alguna categoría específica, por ejemplo, todos los números enteros posibles, todos los caracteres string posibles, todos los teléfonos de proveedores posibles, todos los documentos XML posibles, todas las relaciones con cierta cabecera posibles(y así sucesivamente)” [Dat12].

Cada atributo de cada relación es definido como de un tipo. Los atributos son pares ordenados de combinaciones atributo-nombre/tipo-nombre y una tupla es un par ordenado de atributos. El modelo relacional también soporta varios tipos de llaves, que poseen las propiedades de unicidad, ninguna contiene dos tuplas distintas con el mismo valor e irreductibilidad, ningún subconjunto suyo es tiene unicidad. La llave foránea (*FK*) es una combinación o set de atributos FK en una relación *r2* tal que se requiere que cada valor FK sea igual a algún valor de alguna llave K en alguna relación *r1* (*r1* y *r2* no son necesariamente distintos).

Una restricción de integridad (*constraint*) es una expresión booleana que debe evaluarse como verdadera. Los constraints de tipo definen los valores que constituyen un tipo dado, mientras que los constraints de base de datos limitan los valores que pueden aparecer en cierta base de datos. Las bases de datos suelen tener múltiples constraints específicos, expresados en términos de sus relaciones, sin embargo, el modelo relacional incluye dos constraints genéricos, que aplican a cada base de datos:

- Regla de integridad de identidad: Las llaves primarias no pueden ser nulas (*null*).
- Regla de integridad de referencia: No debe haber valores FK sin relación (si *B* referencia a *A*, *A* debe existir).

La manipulación de bases de datos relacionales se basa en el álgebra relacional dando la colección de operadores que pueden aplicarse a las relaciones, por ejemplo diferencia (*MINUS*). El operador de asignación relacional, permite se le asigne valor de alguna expresión regular a alguna relación, por ejemplo *r1 MINUS r1* cuando *r1* y *r2* son relaciones.

2.2. Correlación lineal

Cuando se tiene una variable controlada x y una dependiente y tenemos el modelo lineal

$$Y = \beta_0 + \beta_1 x + \epsilon \quad (1)$$

que implica entonces el modelo para análisis de rendimiento promedio:

$$E(Y) = \beta_0 + \beta_1 x. \quad (2)$$

Si la variable x es un valor observado de una variable X , al establecerse una relación funcional y al basarse en (2) se implica el modelo

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

que supone la esperanza condicional de Y para un valor fijo de X en una función lineal del valor x . Al suponer que la variable aleatoria vectorial (X, Y) tiene una distribución normal bivariable con $E(X) = \mu_X, E(Y) = \mu_Y, V(X) = \sigma_X^2, V(Y) = \sigma_Y^2$, el coeficiente de correlación ρ puede demostrar que

$$E(Y|X = x) = \beta_0 + \beta_1 x, \quad \text{donde } \beta_1 = \frac{\sigma_Y}{\sigma_X} \rho. \quad (4)$$

Si (X, Y) tiene una distribución normal bivariable, entonces la prueba de independencia es equivalente a probar si el coeficiente de correlación ρ es igual a cero. Denotando con $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ una muestra de aleatoria de distribución normal bivariable. El estimador de máxima probabilidad de ρ está dado por el coeficiente de correlación muestral:

$$r = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2 \sum_{i=1}^n (Y_i - \bar{Y})^2}}. \quad (5)$$

Puede expresarse r en términos de cantidades conocidas:

$$r = \frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}} = \hat{\beta} \sqrt{\frac{S_{xx}}{S_{yy}}}. \quad (6)$$

Cuando (X, Y) tenga una distribución normal bivariable, se sabe que

$$E(Y|X = x) = \beta_0 + \beta_1 x, \quad \text{donde } \beta_1 = \frac{\sigma_Y}{\sigma_X} \rho. \quad (7)$$

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

$$t = \frac{\hat{\beta}_1 - 0}{S/\sqrt{S_{xx}}}, \quad (8)$$

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.” Dennis D Wackerly, William Mendenhall III y Richard L Scheaffer. *Estadística matemática con aplicaciones*. Trad. por Jorge Humberto Romo Muñoz. 7.^a ed. Ciudad de México: Cenage Learning, 2009. ISBN: 9780495110811 Sin embargo, en muestras moderadamente grandes podemos probar la hipótesis $H_0: \rho_1 = \rho_0$ con una prueba Z en la que

$$Z = \frac{(\frac{1}{2}) \ln(\frac{1+r}{1-r}) - (\frac{1}{2}) \ln(\frac{1+\rho}{1-\rho})}{\frac{1}{\sqrt{n-3}}}. \quad (9)$$

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:

$$\begin{aligned} H_a: \rho > \rho_0, & \quad RR: z > z_\alpha \\ H_a: \rho < \rho_0, & \quad RR: z < -z_\alpha \\ H_a: \rho \neq \rho_0, & \quad RR: |z| > z_\alpha/2 \end{aligned} \quad (10)$$

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

$$r^2 = \left(\frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}} \right)^2 = \left(\frac{S_{xy}}{S_{xx}} \right) \left(\frac{S_{xy}}{S_{yy}} \right) = \left(\frac{\hat{\beta}_1 S_{xy}}{S_{yy}} \right) = \frac{S_{yy} - SSE}{S_{yy}} = 1 - \frac{SSE}{S_{yy}}. \quad (11)$$

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.

2.3. Probabilidad condicional

La probabilidad condicional tiene las mismas características que la probabilidad, pero $P(\cdot|B)$ actualiza nuestra incertidumbre acerca de los eventos para reflejar la evidencia observada en B .

Si A y B son eventos con $P(B) > 0$ entonces la *probabilidad condicional* de A dado B denotado por $P(A|B)$, se define como

$$P(A|B) = \frac{P(A \cap B)}{P(B)}. \quad (12)$$

Se denomina probabilidad a priori de A a $P(A)$ y probabilidad a posteriori de A a $P(A|B)$ y es importante mencionar que $P(A|B) \neq P(B|A)$.

La probabilidad condicional es la razón de dos probabilidades y sus consecuencias, la primera de ellas se obtiene moviendo el denominador en la definición al otro lado de la ecuación, para cada evento A y B con posibilidades positivas,

$$P(A \cap B) = P(B)P(A|B) = P(A)P(B|A). \quad (13)$$

se le conoce como teorema de la *probabilidad de la intersección de dos eventos*. Aplicando repetidamente el teorema (13) aplicado a la intersección de n eventos obtenemos el teorema de *probabilidad de la intersección de n eventos*. Para cualquier evento A_1, \dots, A_n con probabilidad $P(A_1, A_2, \dots, A_{n-1}) > 0$,

$$P(A_1, A_2, \dots, A_n) = P(A_1)P(A_2|A_1)P(A_3|A_1, A_2) \dots P(A_n|A_1, \dots, A_{n-1}). \quad (14)$$

Otro teorema que relaciona a $P(A|B)$ con $P(B|A)$ es la regla *regla de Bayes*:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (15)$$

que se origina directamente del teorema (13) y a su vez se origina directamente de la definición de probabilidad condicional, sin embargo, la regla de Bayes tiene importantes aplicaciones e implicaciones en probabilidad y estadística, ya que en ocasiones es más fácil encontrar $P(B|A)$ que $P(A|B)$ o viceversa. Otra forma de escribir la regla, es en términos de apuestas (*odds*), las odds de un evento A son

$$odds(A) = P(A) \frac{P(A)}{P(A^c)}. \quad (16)$$

Al tomar la expresión $P(A|B)$ y dividirla entre $P(A^c|B)$, ambos de la regla de Bayes; llegamos al *teorema de Bayes en forma de apuestas*: para cualquier evento A y B con posibilidades positivas, las odds de A **after conditioning on** B son

$$\frac{P(A|B)}{P(A^c|B)} = \frac{P(B|A)}{P(B|A^c)} \frac{P(A)}{P(A^c)}. \quad (17)$$

En este caso las odds a posteriori $P(A|B)/P(A^c|B)$ son iguales a las odds a priori $P(A)/P(A^c)$ por el factor $P(B|A)/P(B|A^c)$ lo que se le conoce en estadística como *función de verosimilitud*.

La ley de Bayes es usada en ocasiones en conjunto con la *ley de probabilidad total*, que es esencial para descomponer problemas complicados de probabilidad en problemas partes: Si A_1, \dots, A_n es una partición de una muestra del espacio S , con $P(A_i) > 0$ para todo i , entonces

$$P(B) = \sum_{i=1}^n P(B|A_i)P(A_i). \quad (18)$$

Prueba: como los A_i forman una partición de S , podemos descomponer B como

$$B = (B \cap A_1) \cup (B \cap A_2) \cup \dots \cup (B \cap A_n). \quad (19)$$

como las partes están disjuntas, podemos agregar sus posibilidades para obtener $P(B)$:

$$P(B) = P(B \cap A_1) + P(B \cap A_2) + \dots + P(B \cap A_n). \quad (20)$$

Aplicando el teorema (13) a cada $P(B \cap A_i)$ obtenemos

$$P(B) = P(B|A_1)P(A_1) + \dots + P(B|A_n)P(A_n). \quad (21)$$

La *ley de probabilidad total* dice que para obtener la probabilidad incondicional de B , tenemos que dividir el espacio muestra en cortes disjuntos A_i , encontrar la probabilidad condicional de B en cada corte, después tomar la suma ponderada de las probabilidades condicionales, donde los pesos son las probabilidades $P(A_i)$.

2.4. Variables aleatorias y sus distribuciones

Una *variable aleatoria* es una función asignando un número real \mathbb{R} a cada posible resultado de un experimento. Con una muestra en espacio \mathcal{S} , una variable aleatoria X asigna el valor numérico $X(s)$ a cada resultado posible s del experimento. La aleatoriedad viene del hecho que tenemos un experimento aleatorio (con probabilidades descritas por la función de probabilidad P). Las variables aleatorias simplifican la notación y expanden la habilidad de cuantificar y resumir resultados de experimentos.

Se dice que una variable X es discreta cuando si hay una lista finita de valores a_1, a_2, \dots, a_n o una lista infinita de valores a_1, a_2, \dots de tal forma que $P(X = a_j \text{ para algún } j) = 1$. Si X es una variable aleatoria discreta, entonces el conjunto infinito o contable de valores x tal que $P(X = x)$ se llama *soporte* de X . En contraste una variable aleatoria continua puede tomar cualquier valor real en un intervalo.

2.4.1. Función de probabilidad

La forma más natural de expresar la distribución de variables aleatorias discretas es la *función de probabilidad* [BH19] (PMF, por sus siglas en inglés?) que, para una X discreta, es la función p_X dada por $p_X(x) = P(X = x)$. El teorema de *funciones de probabilidad válidas* dice que cuando X es una variable aleatoria con soporte x_1, x_2, \dots , la función de probabilidad p_X de x debe satisfacer los siguiente criterios:

- No negativo $p_X(x) > 0$ si $x = x_j$ para un j , y $p_X(x) = 0$, de otra forma;
- Suma 1: $\sum_{j=1}^{\infty} p_X(x_j) = 1$.

el primer criterio es verdadero porque la probabilidad es no negativa, el segundo es verdadero ya que X debe tomar *algún* valor, y los eventos $X = x_j$ están disjuntos, entonces

$$\sum_{j=1}^{\infty} P(X = x_j) = P\left(\bigcup_{j=1}^{\infty} \{X = x_j\}\right) = P(X = x_1 \text{ ó } X = x_2 \text{ ó } \dots) = 1. \quad (22)$$

2.4.2. Distribución de Bernoulli y binominal

Una variable aleatoria tiene la *distribución de Bernoulli* con un parámetro p si $P(X = 1) = p$ y $P(X = 0) = 1 - p$, cuando $0 < p < 1$. Se escribe como $X \sim \text{Bern}(p)$, el símbolo \sim significa “distribuido como” y la probabilidad p es el *parámetro*, que determina qué distribución de Bernoulli específica tenemos.

Supóngase que se realizan n ensayos Bernoulli independientes, cada uno con probabilidad p de éxito. X sea el número de éxitos, la distribución X se llama *distribución binominal* con parámetros n y p ; se escribe $X \sim \text{Bin}(p, n)$. $\text{Bern}(p)$ es la misma distribución que $\text{Bin}(1, p)$. Bernoulli es un caso especial de binominal, si $x \sim \text{Bin}(1, p)$, entonces la función de probabilidad de X es

$$P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k} \quad (23)$$

para $k = 0, 1, \dots, n$ (y por otra parte $P(X = k) = 0$).

2.4.3. Distribución de hipergeométrica

Si $X \sim \text{HGeom}(w, b, n)$, entonces la función de probabilidad de X es

$$P(X = k) = \frac{\binom{w}{k} \binom{b}{n-k}}{\binom{w+b}{n}}, \quad (24)$$

para enteros k satisfaciendo $0 \leq k \leq w$ y $0 \leq n - k \leq b$, y $P(X = k) = 0$. La estructura esencial de la distribución hipergeométrica se basa en que objetos en su población están clasificados usando dos tipos de etiquetas, al menos una de estas siendo asignada al azar. Las distribuciones $\text{HGeom}(w, b, n)$ y $\text{HGeom}(n, w + b - n, 1)$ son idénticas si X y Y tienen la misma distribución, podemos demostrarlo algebraicamente:

$$P(X = k) = \frac{\binom{w}{k} \binom{b}{n-k}}{\binom{w+b}{n}} = \frac{w!b!n!(w+b-n)!}{k!(w+b)!(w-k)!(n-k)!(b-n+k)!} \quad (25)$$

$$P(X = k) = \frac{\binom{n}{k} \binom{w+b-n}{w-k}}{\binom{w+b}{w}} = \frac{w!b!n!(w+b-n)!}{k!(w+b)!(w-k)!(n-k)!(b-n+k)!} \quad (26)$$

2.4.4. Distribución uniforme discreta

Teniendo C , un conjunto finito no vacío de números, se elige un número uniformemente al azar (o sea que todos los números tienen la misma posibilidad de ser elegidos), llámese X . Entonces se dice que X una *distribución uniforme discreta* con el parámetro C . Se dice entonces que la función de probabilidad de $X \sim \text{DUNif}(C)$ (la distribución uniforme discreta de X) es

$$P(X = x) = \frac{1}{|C|} \quad (27)$$

para $x \in C$ (de lo contrario 0) ya que la función de probabilidad debe sumar 1.

2.4.5. Función de distribución acumulada

Esta función describe la distribución de todas las variables aleatorias (a diferencia de la función de probabilidad que sólo se aplica a las discretas). La *función de distribución acumulada* de una variable aleatoria X es la función F_X dada por $F_X(x) = P(X \leq x)$ y tiene las siguientes propiedades:

- Incrementos: Si $x_1 \leq x_2$, then $F(x_1) \leq F(x_2)$.
- Continua por la derecha: Es continua por la posibilidad de tener saltos. Cuando hay saltos es continua por la derecha, es decir, por cada a se tiene

$$F(a) = \lim_{c \rightarrow a^+} F(c). \quad (28)$$

- Convergencia de 0 y 1 en los límites

$$\lim_{x \rightarrow -\infty} F(x) = 0 \quad \text{y} \quad \lim_{x \rightarrow \infty} F(x) = 1. \quad (29)$$

2.5. Valor esperado

Mientras que las distribuciones anteriores nos han dado toda la información acerca de la probabilidad de las variables aleatorias, cuando sólo se requiere un número que extraiga su valor, podemos utilizar la *media*, también conocida como *valor esperado*. Dada una lista de números x_1, x_2, \dots, x_n , para obtener la *media aritmética*, estos se suman y dividen entre n :

$$\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j, \quad (30)$$

la *media ponderada* de x_1, x_2, \dots, x_n se obtiene de la siguiente forma:

$$\text{media ponderada}(x) = \frac{1}{n} \sum_{j=1}^n x_j P_j, \quad (31)$$

donde los pesos p_1, p_2, \dots, p_n son números no negativos previamente especificados que suman a 1.

El valor esperado o media de una variable aleatoria discreta X cuyos posibles valores distintos son x_1, x_2, \dots es definida por

$$E(X) = \sum_{j=1}^{\infty} x_j P(X = x_j), \quad (32)$$

si el soporte es finito, entonces se reemplaza por una suma finita, escribiéndose de la siguiente forma:

$$E(X) = \sum_x \underbrace{x}_{\text{valor}} \underbrace{P(X = x)}_{\substack{\text{Función de} \\ \text{probabilidad} \\ \text{en } x}}. \quad (33)$$

2.5.1. Linealidad del valor esperado

El valor esperado de una suma de variables aleatorias es la suma de sus valores esperados individuales, este es el teorema de la *linealidad del valor esperado*, donde para cada variable aleatoria X, Y y cada constante c ,

$$\begin{aligned} E(X + Y) &= E(X) + E(Y), \\ E(cX) &= cE(X). \end{aligned} \tag{34}$$

3. Objetivos

Comenzar desde cero cada proyecto es ineficiente. (esto es temporal)

3.1. General

Diseñar y desarrollar una
plataforma/framework de ML para NLP
reutilizable y/o de uso general/ ¿para casos de uso similares?

3.2. Particulares

Revisé algunas tesis para darme la idea, necesito saber un poco más del tema para desarrollar esta parte, creo que sería algo así:

1. Investigación
2. Análisis info
3. Diseño/ desarrollo
4. Comprobación

3.3. Hipótesis

Al integrar
/tecnologías/técnicas/modelos/librerías/frameworks/
de RDB, ML, NLP, ¿transfer learning?, ¿deep learning?,...
se puede /diseñar/ desarrollar/, ¿e implementar? una
/plataforma/framework/
/genérica/normalizada/universal/reutilizable/estandarizada/compatible
/con/en/para/ /casos de uso similares/diversos casos de uso/
(/reduciendo tiempo/ahorrando recursos/.) estas oración se sale del scope

4. Metas

Contribuir por medio del diseño y desarrollo de la
/plataforma/framework/ de ML para NLP
que será utilizable en variedad de casos reales /con características similares./

5. Metodologías

6. Referencias

- [BH19] Joseph K. Blitzstein y Jessica Hwang. *Introduction to Probability [Introducción a la probabilidad]*. 2.^a ed. Texts in Statistical Science [Textos en ciencia estadística]. Florida: Chapman y Hall/cRC, 2019. ISBN: 9781138369917.
- [Dat12] C. J. Date. *SQL and Relational Theory: How to Write Accurate SQL Code [SQL y teoría relacional: Cómo escribir código SQL correcto]*. 2.^a ed. California: O'Reilly Media, 2012. ISBN: 9781449316402.
- [WMIS09] Dennis D Wackerly, William Mendenhall III y Richard L Scheaffer. *Estadística matemática con aplicaciones*. Trad. por Jorge Humberto Romo Muñoz. 7.^a ed. Ciudad de México: Cengage Learning, 2009. ISBN: 9780495110811.