${\bf \acute{I}ndice}$

1.	Mar	Marco teórico		
	1.1.	Álgebra relacional	2	
	1.2.	Modelo de base de datos relacional	2	
	1.3.	Inteligencia artificial	3	
	1.4.	Aprendizaje maquinal	3	
	1.5.	Procesamiento de lenguaje natural	4	
	1.6.	Correlación lineal	5	
	1.7.	Regresión2	7	
	1.8.	Probabilidad condicional	7	
	1.9.	Variables aleatorias y sus distribuciones	9	
		1.9.1. Función de probabilidad	9	
		1.9.2. Distribución de Bernoulli y binominal	9	
		1.9.3. Distribución de hipergeométrica	10	
		1.9.4. Distribución uniforme discreta	10	
		1.9.5. Función de distribución acumulada	11	
	1.10.	Valor esperado	11	
		1.10.1. Binominal geométrica y negativa	12	
		1.10.2. Varianza, agregarlas a la secc que pertenecen	14	
		1.10.3. r.v. continuas	14	
2.	Obi	etivos	16	
		General	16	
		Particulares	16	
		Hipótesis	16	
3.	Met	Metas 1		
4.	Metodologías 1			
5.	Referencias 1			

1. Marco teórico

1.1. Álgebra relacional

Pegar texto de gdocs

Operaciones unarias: $\mu(R) \to R'$ Operaciones binarias: $\beta(P,Q) \to R$

Relaciones: P, Q, R, R'

1.2. Modelo de base de datos relacional

El modelo relacional de base de datos, según Date [Dat12] consiste en cinco componentes:

- 1. Una colección de tipos escalares, pueden ser definidos por el sistema (INTE-GER, 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 para obtener valores relacionales de otros valores relacionales mediante operadores relacionales genéricos.

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 nombradotodos 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 se 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).
- \blacksquare 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.

1.3. Inteligencia artificial

Inteligencia artificial es definida como "el esfuerzo por automatizar tareas intelectuales normalmente realizadas por humanos" [Cho18], de este campo general se desprenden el aprendizaje maquinal y el aprendizaje profundo.

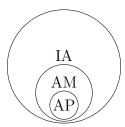


Figura 1: Aprendizaje profundo (AP), es un subcampo del aprendizaje maquinal (AM), que a su vez es un subcampo de la inteligencia artificial (IA)[Cho18].

1.4. Aprendizaje maquinal

La definición de aprendizaje maquinal de Tom Mitchell
[Mit97] dice que "un programa de computadora aprende de experiencia ${\cal E}$ con respecto a una tarea ${\cal T}$ y una

end-to-end, de principio a fin?

medición de rendimiento P, si su rendimiento en T, medido por P, mejora con experiencia E."

Esto dignifica que 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 "y el producto son las reglas"*..

Si no es programado explícitamente, entonces un sistema de aprendizaje maquinal es entrenado: se le presentan muchos ejemplos relevantes a una tarea, y si encuentra una estructura estadística en ellos, genera reglas para automatizar la tarea.

1.5. Procesamiento de lenguaje natural

El procesamiento de lenguaje natural, es el conjunto de métodos para hacer accesible el lenguaje humano a las computadoras[Eis19].* Existen dos enfoques en lo que debe ser su tarea central:

- Entrenar sistemas de extremo a extremo* que transmuten texto sin procesar a cualquier estructura deseada.
- Transformar texto en una pila de estructuras lingüísticas de uso general que en teoría deben poder soportar cualquier aplicación.

Dos de los módulos básicos de NLP son búsqueda y aprendizaje con los que se puede resolver muchos problemas que podemos describir en la siguiente forma matemática

$$\hat{y} = \underset{y \in Y(x)}{argmax} \Psi(x, y; 0), \tag{1}$$

donde,

- \bullet x es la entrada, un elemento de un conjunto X.
- y es el resultado, un elemento de un conjunto Y.
- Ψ es una función de puntuación (también conocida como modelo), que va desde el conjunto $X \times Y$ hasta los números reales.
- Ø es el vector de parámetros para Ψ.
- \hat{y} es el resultado previsto, que es elegido para maximizar la función de puntuación.

El módulo de búsqueda se encarga de computar el argmax de la función Ψ , es decir, encuentra el resultado \hat{y} con la mejor puntuación con respecto a la entrada x. El módulo de aprendizaje encuentra los parámetros θ por medio del procesamiento de grandes conjuntos de datos de ejemplos etiquetados $\{(x^i, y^i)\}_{i=1}^N$.

1.6. Correlación lineal

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

$$Y = \beta_0 + \beta_1 x + \epsilon \tag{2}$$

que implica entonces el modelo para análisis de rendimiento promedio:

$$E(Y) = \beta_0 + \beta_1 x. \tag{3}$$

Si la variable x es un valor observado de una variable X, al establecerse una relación funcional y al basarse en (3) se implica el modelo

$$E(Y|X=x) = \beta_0 + \beta_1 x \tag{4}$$

que supone el valor esperado 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

$$E(Y|X=x) = \beta_0 + \beta_1 x, \quad donde \ \beta_1 = \frac{\sigma_Y}{\sigma_Y} \rho.$$
 (5)

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 $(X_1,Y_1),(X_2,Y_2),\ldots,(X_n,Y_n)$ una muestra de aleatoria de distribución normal bivariante. El estimador de máxima probabilidad de ρ está dado por el coeficiente de correlación muestral:

$$r = \frac{\sum_{i=1}^{n} (X_i - \overline{X})(Y_i - \overline{Y})}{\sqrt{\sum_{i=1}^{n} (X_i - \overline{X})^2 \sum_{i=1}^{n} (Y_i - \overline{Y})^2}}.$$
 (6)

Puede expresarse r en términos de cantidades conocidas:

$$r = \frac{Sxy}{\sqrt{SxxSyy}} = \hat{\beta}\sqrt{\frac{Sxx}{Syy}}.$$
 (7)

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

$$E(Y|X=x) = \beta_0 + \beta_1 x, \quad donde \ \beta_1 = \frac{\sigma_Y}{\sigma_X} \rho.$$
 (8)

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

$$t = \frac{\hat{\beta}_1 - 0}{S/\sqrt{Sxx}},\tag{9}$$

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." WMIS09] Sin embargo, en muestras moderadamente grandes podemos probar la hipótesis H_0 : $\rho_1 = \rho_0$ con una prueba Z en la que

$$Z = \frac{(\frac{1}{2})\ln(\frac{1+r}{1-r}) - (\frac{1}{2})\ln(\frac{1+\rho}{1-\rho})}{\frac{1}{\sqrt{n-3}}}.$$
 (10)

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:

$$H_a: \rho > \rho 0, \quad RR: z > z_{\alpha}$$

 $H_a: \rho < \rho 0, \quad RR: z < -z_{\alpha}$
 $H_a: \rho \neq \rho 0, \quad RR: |z| > z_{\alpha}/2$ (11)

La suma de los cuadrados del error SSE, es es una alternativa para medie 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

$$r^{2} = \left(\frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}}\right)^{2} = \left(\frac{S_{xy}}{S_{xx}}\right)\left(\frac{S_{xy}}{S_{yy}}\right) = \left(\frac{\hat{\beta}_{1}S_{xy}}{S_{yy}}\right) = \frac{S_{yy} - SEE}{S_{yy}} = 1 - \frac{SSE}{S_{yy}}.$$
 (12)

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.

1.7. Regresión2

La correlación de dos variables con valores positivos y existentes x y y es por definición

 $\frac{cov(x,y)}{\sqrt{var(x)}\sqrt{var(y)}}\tag{13}$

1.8. 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

$$P(A|B) = \frac{P(A \cap B)}{P(B)}. (14)$$

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,

$$P(A \cap B) = P(B)P(A|B) = P(A)P(B|A). \tag{15}$$

se le conoce como teorema de la probabilidad de la intersección de dos eventos. Aplicando repetidamente el teorema (15) aplicado a la intersección de n eventos obtenemos el teorema de probabilidad de la intersección de n eventos. Para cualquier evento A_1, \ldots, A_n con probabilidad $P(A_1, A_2, \ldots, A_{n-1}) > 0$,

$$P(A_1, A_2, \dots, A_n) = P(A_1)P(A_2|A_1)$$

$$P(A_3|A_1, A_2) \dots P(A_n|A_1, \dots, A_{n-1}).$$
(16)

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

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \tag{17}$$

que se origina directamente del teorema (15) 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

$$odds(A) = P(A)\frac{P(A)}{P(A^{c})}.$$
(18)

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

$$\frac{P(A|B)}{P(A^{c}|B)} = \frac{P(B|A)}{P(B|A^{c})} \frac{P(A)}{P(A^{c})}.$$
(19)

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 es esencial para descomponer problemas complicados de probabilidad en problemas partes: Si A_1, \ldots, A_n es una partición de una muestra del espacio S, con $P(A_i) > 0$ para todo i, entonces

$$P(B) = \sum_{i=1}^{n} P(B|A_i)P(A_i).$$
 (20)

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

$$B = (B \cap A_1) \cup (B \cap A_2) \cup \ldots \cup (B \cap A_n). \tag{21}$$

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

$$P(B) = P(B \cap A_1) + P(B \cap A_2) + \dots + P(B \cap A_n). \tag{22}$$

Aplicando el teorema (15) a cada $P(B \cap A_i)$ obtenemos

$$P(B) = P(B|A_1)P(A_1) + \ldots + P(B|A_n)P(A_n).$$
(23)

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 $P(A_i)$.

1.9. 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 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, a_2, \ldots, a_n o un una lista infinita de valores a, a_2, \ldots de tal forma que $P(X = a_j)$ 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.

1.9.1. 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 $x1, x2, \ldots$, 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=xj están disjuntos, entonces

$$\sum_{j=1}^{\infty} P(X = x_j) = P\left(\bigcup_{j=1}^{\infty} \{X = x_j\}\right) = P(X = x_1 \text{ ó } X = x_2 \text{ ó } \dots) = 1.$$
 (24)

1.9.2. 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 $0 . Se escribe como <math>X \sim Bern(p)$, 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 $X \sim Bin(p,n)$. Bern(p) es la misma distribución que Bin(1,p). Bernoulli es un caso especial de binominal, si $x \sim Bin(1,p)$, entonces la función de probabilidad de X es

$$P(X = k) = \binom{n}{k} p^k (1 - p)^{n - k}$$
 (25)

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

1.9.3. Distribución de hipergeométrica

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

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

para enteros k satisfaciendo $0 \le k \le w$ y $0 \le n-k \le b$, 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 HGeom(w,b,n) y HGeom(n,w+b-n,1) son idénticas si X y Y tienen la misma distribución, podemos demostrarlo algebraicamente:

$$P(X=k) = \frac{\binom{w}{k}\binom{b}{n-k}}{\binom{w+b}{n}} = \frac{w!b!n!(w+b-n)!}{k!(w+b)!(w-k)!(n-k)!(b-n+k)!}$$
(27)

$$P(X=k) = \frac{\binom{n}{k}\binom{w+b-n}{w-k}}{\binom{w+b}{w}} = \frac{w!b!n!(w+b-n)!}{k!(w+b)!(w-k)!(n-k)!(b-n+k)!}.$$
 (28)

1.9.4. 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 $X \sim DUNif(C)$ (la distribución uniforme discreta de X) es

$$P(X=x) = \frac{1}{|C|} \tag{29}$$

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

1.9.5. 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 $F_X(x) = P(X \le x)$ y tiene las siguientes propiedades:

- Incrementos: Si $x_1 \le x_2$, then $F(x_1) \le F(x_2)$.
- 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

$$F(a) = \lim_{c \to a^+} F(x). \tag{30}$$

■ Convergencia de 0 y 1 en los límites

$$\lim_{x \to \infty} F(x) = 0 \quad \text{y} \quad \lim_{x \to \infty} F(x) = 1. \tag{31}$$

1.10. 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, \ldots, x_n , para obtener la $media \ aritmética$, estos se suman y dividen entre n:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_j, \tag{32}$$

la $media\ ponderada\ de\ x_1,x_2,\ldots,x_n$ se obtiene de la siguiente forma:

media ponderada
$$(x) = \frac{1}{n} \sum_{j=1}^{n} x_j P_j,$$
 (33)

donde los pesos p_1, p_2, \ldots, 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, \ldots es definida por

$$E(X) = \sum_{j=1}^{\infty} x_j P(X = xj), \tag{34}$$

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

$$E(X) = \sum_{x} \underbrace{x}_{\text{valor}} \underbrace{P(X = x)}_{\text{Función de}}.$$

$$\underbrace{P(X = x)}_{\text{probabilidad}}.$$

$$\underbrace{P(X = x)}_{\text{en } x}.$$
(35)

**********nuevo * Varianza

$$Var(X) = E(X - EX)^{2} = E(X^{2}) - (EX)^{2}$$
(36)

*******nuevo

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,

$$E(X+Y) = E(X) + E(Y),$$

$$E(cX) = cE(X).$$
(37)

1.10.1. Binominal geométrica y negativa

Distribución geométrica: Se tiene una secuencia de ensayos independientes Bernoulli, cada uno con la misma probabilidad de éxito $p \in (0,1)$, con ensayos realizados hasta que se alcanza el éxito. X es el número de fallas antes de la primera prueba exitosa por lo que X tiene una distribución geométrica con un parámetro p; denotado $X \sim Geom(p)$. Con esto podemos llegar a los teoremas de distribución geométrica de la función de probabilidad, cuando $X \sim Geom(p)$, entonces la función de probabilidad de X será

$$P(X=k) = q^k p (38)$$

para $k = 1, 2, \ldots$, cuando q = 1 - p; y el teoremas de distribución geométrica de la función de distribución acumulativa, cuando $X \sim Geom(p)$, entonces la función de distribución acumulativa de X será

$$F(x) = \begin{cases} 1 - q^{\lfloor x \rfloor + 1}, & \text{si } x \ge 0; \\ 0, & \text{si } x < 0, \end{cases}$$
 (39)

cuando q = 1 - q y $\lfloor x \rfloor$ es el mayor entero y menor o igual a x.

El valor esperado geométrico de $X \sim Geom(p)$ es

$$E(X) = \sum_{k=0}^{\infty} kq^k p,$$
(40)

cuando q = 1 - p. Aunque esta no es una serie geométrica, podemos llegar a ello

$$\sum_{k=0}^{\infty} q^k = \frac{1}{1-q}$$

$$\sum_{k=0}^{\infty} k q^{k-1} = \frac{1}{1-q^2},$$
(41)

finalmente multiplicamos ambos lados por pq, recuperando la suma original que queríamos encontrar

$$E(X) = \sum_{k=0}^{\infty} kq^k p = pq \sum_{k=0}^{\infty} kq^{k-1} = pq \frac{1}{(1-q)^2} = \frac{q}{p}.$$
 (42)

Primer valor esperado de éxito FS, podemos definir a $Y \sim FS(p)$ como Y = X + 1 donde $X \sim Geom(p)$, por lo que tenemos

$$E(Y) = E(X+1) = \frac{q}{p} + 1 = \frac{1}{p}.$$
 (43)

Las distribuciones binominales negativas generalizan la distribución geométrica en lugar de esperar por un éxito, podemos esperar por cualquier número predeterminado r de éxitos. En una secuencia de ensayos independientes Bernoulli con probabilidad de éxito p, si X es el número de fallas antes del éxito número r, entonces se dice que Xtiene una distribución binominal negativa con parámetros r y p, denotado $X \sim NBin(r, p)$.

La distribución binominal cuenta el número de éxitos en un número fijo de ensayos, mientras que la binominal negativa cuenta el número de fallas hasta alcanzar cierto número de éxitos. Si $X \sim NBin(r, p)$, entonces la función de probabilidad de X es

$$P(X=n) = \binom{n+r-1}{r-1} p^r q^n \tag{44}$$

para n = 0, 1, 2..., donde q = 1 = p.

1.10.2. Varianza, agregarlas a la secc que pertenecen

Varianza de la geométrica (agregarla a la geom después de editar todo y hacerlo más breve)

$$Var(X) = E(X^{2}) - (EX)^{2} = \frac{q(1+q)}{p^{2}} - (\frac{q}{p})^{2} = \frac{q}{p^{2}}$$
(45)

Varianza de la geométrica (lo mismo)

$$Var(X) = E(X^{2}) - (EX)^{2} = (n(n-1)p^{2} + np) - (np)^{2} = np(1-p).$$
 (46)

Una variable aleatoria X tiene distribución de Poisson (denotada $X \sim Pois(\lambda)$) con el parámetro λ , cuando $\lambda > 0$ si la PMF de x es

$$P(X = k) = \frac{e^{-\lambda} \lambda^k}{k!}, k = 0, 1, 2, \dots$$
 (47)

Varianza de la distribución de Poisson es

$$Var(X) = E(X^2) - (EX)^2 = \lambda(1+\lambda) - \lambda^2 = \lambda \tag{48}$$

1.10.3. r.v. continuas

A diferencia de las variables discretas, las variables aleatorias continuas pueden tomar cualquier valor real en un intervalo y tienen una distribución continua. Para obtener la probabilidad deseadaWHOMST, se debe integrar la función de densidad de probabilidad sobre el rango apropiado

$$P(X \in A) = \int_{A} f(x)dx \tag{49}$$

La distribución logística se obtiene

$$F(x) = \frac{e^x}{1 + e^x}, x \in \Re$$
 (50)

El valor esperado de la continua función de distribución acumulada f es

$$E(X) = \int_{-\infty}^{\infty} x f(x) dx \tag{51}$$

la continua U tiene dist unif en el intervalo (a,b), denotada $U \sim Unif(a,b)$ si el área acumulada bajo la función de densidad de probabilidad es

$$f(x) = \begin{cases} \frac{1}{b-a}, & \text{si } a < x < b; \\ \text{de lo contrario}, & 0 \end{cases}$$
 (52)

 $U \sim Unif(a, b)$ de una variable aleatoria es

$$\tilde{U} = a + (b - a)U \tag{53}$$

Su varianza es

$$Var(U) = \frac{(b-a)^2}{12} \tag{54}$$

La distribución normal $(X \sim N(\mu, \sigma^2))$ cuando $\mathbb{Z} \sim N(0, 1)$ es

$$X = \mu + \sigma \mathbb{Z} \tag{55}$$

Por lo tanto obtendremos el valor esperado de μ y varianza de σ^2

$$X = (\mu + \sigma \mathbb{Z}) = E(\mu) + \sigma E(\mathbb{Z}) = \mu \tag{56}$$

$$Var = (\mu + \sigma \mathbb{Z}) = Var(\sigma \mathbb{Z}) = \sigma^2 Var(\mathbb{Z}) = \sigma^2.$$
 (57)

La distribución exponencial de X con un parámetro λ , cuando $\lambda > 0$ si su función de densidad de probabilidad es $f(x) = \lambda e^{-\lambda x}, x > 0$; denotada como $X \sim Expo(\lambda)$ es la siguiente

$$F(x) = 1 - e^{-\lambda x}, x > 0 (58)$$

	Discrete r.v.	Continuous r.v.
CDF	$F(x) = P(X \le x)$	$F(x) = P(X \le x)$
PMF/PDF	P(X=x)	f(x) = F'(x)
	• PMF is height of jump of F at x .	• PDF is derivative of F .
	• PMF is nonnegative.	• PDF is nonnegative.
	• PMF sums to 1.	• PDF integrates to 1.
	• $P(X \in A) = \sum_{x \in A} P(X = x).$	• $P(X \in A) = \int_A f(x)dx$.
Expectation	$E(X) = \sum_{x} x P(X = x)$	$E(X) = \int_{-\infty}^{\infty} x f(x) dx$
LOTUS	$E(g(X)) = \sum_{x} g(x)P(X = x)$	$E(g(X)) = \int_{-\infty}^{\infty} g(x)f(x)dx$

Figura 2: tal vez me sirva para resumir, imagen no es para producto final

2. Objetivos

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

2.1. General

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

2.2. 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

2.3. Hipótesis

```
Al integrar
/tecnologías/técnicas/modelos/librerías/frameworks/
de RDB, ML, NLP, ¿transfer learning?, ¿deep learning?,...
se de 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
```

3. 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./

4. Metodologías

5. Referencias

- [BH19] Joseph K. Blitzstein y Jessica Hwang. Introduction to Probability [Introduccién a la probabilidad]. 2.ª ed. Texts in Statistical Science [Textos en ciencia estadistica]. Florida: Chapman y Hall/cRC, 2019. ISBN: 9781138369917.
- [Cho18] François Chollet. Machine Learning With Python [Machine Learning]. New York: Manning Publications Co., 2018. ISBN: 9781617294433.
- [Dat12] C. J. Date. SQL and Relational Theory: How to Write Accurate SQL Code [SQL y teoria relacional: Cómo escribir código SQL correcto]. 2.ª ed. California: 0'Reilly Media, 2012. ISBN: 9781449316402.
- [Eis19] Jacob Eisenstein. Introduction to natural language processing [Introducción al procesamiento de lenguaje natural]. Adaptive computation and machine learning [Computación adaptativa y aprendizaje maquinal]. Cambridge: MIT Press, 2019. ISBN: 9780262042840.
- [Mit97] Tom Mitchell. Machine learning [Aprendizaje maquinal]. New York: MacGraw-Hill, 1997. ISBN: 0070428077.
- [WMIS09] Dennis D Wackerly, william Mendenhall III y Richard L Scheaffer. Estadística matemática con aplicaciones. Trad. por Jorge Humberto Romo Mufioz. 7.ª ed. Ciudad de México: Cenage Learning, 2009. ISBN: 9780495110811.