Sprint 1 SPRINT 1 - ENCUENTRO 5 FEB Probabilidad y Estadística -18 **Pandas**

Unirme El encuentro comenzó hace 00:07:51

> Lista de recursos Challenge Profundiza Notebook:

Conectando puntos

 \bigcap

En el encuentro y la bitácora anterior (Volver a las bases), comenzamos a ver algunos conceptos de Probabilidad y Estadística. En particular, vimos algunas

herramientas de estadística descriptiva. Recordemos la situación que planteamos: Vamos al casino a jugar Craps —un juego de azar que consiste en apostar en base a lo que esperamos que salga de una tirada de dados—, pero comenzamos

a sospechar que uno de los dados que estamos utilizando está cargado, ¿cómo podremos saber cuál? 9 Muestras

> 14 Muestras 0.35 0.30 0.25 0.15 0.10 0.05

Siempre es importante que te tomes unos minutos para reflexionar sobre las situaciones que planteamos. Si ya estás listo/a, hoy veremos algunos conceptos fundamentales de Probabilidad, para terminar de complementar lo visto hasta acá.

Probabilidad La Teoría de la Probabilidad es la disciplina que nos provee de herramientas para calcular y razonar acerca de probabilidades. Existen dos definiciones conocidas para probabilidad, no necesariamente excluyentes. Va a ser muy útil que las tengamos en cuenta ya que ambas sirven para pensar en nuestros

problemas:

moneda.

cada una.

edad?

en el Barcelona.

confianza o certidumbre de que un suceso ocurra. La mejor medida de la incertidumbre es la probabilidad. Sin embargo, ésta no suele ser la primera definición de probabilidad que estudiamos. Veamos un ejemplo: Si un meteorólogo pronostica lluvia para mañana con una probabilidad de 0.99 (99% de probabilidad) y una meteoróloga pronostica lluvia

Otra que también usamos mucho —sobre todo cuando tiramos monedas o dados— es la definición de Probabilidad Frecuentista: si hacemos un experimento muchas (!) veces, la probabilidad está asociada a la frecuencia con que ocurre cada posible valor de la variable aleatoria.

caiga ceca. Cuantas más veces tiremos la moneda, más cerca esperamos que estén de esos valores. Si pudiéramos tirar infinitas veces la moneda, la frecuencia con la que salen cara y ceca es la probabilidad de esos eventos. Si bien hay interpretaciones que son más naturales en algunos contextos que

las probabilidades son válidas para ambas (notación: p(s) que refiere a la probabilidad 'p' de un resultado 's'). • Las probabilidades son números mayores que cero y menores que uno. La notación es: $0 \le p(s) \le 1$.

Ahora veremos algunas definiciones más, ¡que probablemente generen más confusión!

médico/a cuando recibe en su consultorio a un nuevo paciente y anota en su historia clínica su edad, peso y altura. ¿Qué otras cosas pueden

acepción más coloquial de experimento? • Un **espacio muestral** es el conjunto de posibles resultados de un experimento. Para un dado, el espacio muestral es S = {1,2,3,4,5,6}. Si

- Una variable aleatoria es una función numérica sobre los resultados de un espacio de probabilidad. Como esta definición es confusa, informalmente podemos decir que es una variable aleatoria es una variable cuyos valores dependen del resultado de un suceso aleatorio. Por ejemplo, los resultados de tirar un dado conforman una variable aleatoria. Tirar dos dados y sumar sus resultados es otra variable aleatoria. Nota: en Probabilidad y Estadística, casi todo las variables que estudiamos son variables aleatorias, por lo que este concepto es central en estas áreas. Por ejemplo, la temperatura de un día o cuántos milímetros de lluvia cayeron
- Para pensar: ¿que se llame variable aleatoria significa que está regida por el azar? Distribuciones de Probabilidad Tipos de datos numéricos ¿Recuerdas que en la bitácora anterior agrupamos las variables en cuantitativas (numéricas) y cualitativas (categóricas)? El primer grupo corresponde a variables como, por ejemplo, edad, peso y altura de una persona, que pueden ser divididos en dos tipos: discretos o continuos. En tanto, el segundo grupo corresponde a variables como su color de pelo, género o profesión.

Tipos de datos

Si bien la diferencia puede ser sutil, hay que prestarle atención.

Numéricos

ejemplo, peso y altura de una persona, temperaturas, etc.

Distribución de probabilidad

continua (ej.: peso y altura de una persona).

probabilidad de que dicho suceso ocurra.

0.16

0.14

0.12

0.10

0.08

0.06

0.04

0.02

0.00

0.3

0.2

0.1

0.0

de probabilidad.

podemos obtener de ella es:

Discretos: Se pueden

No se pueden contar. contar Figura 1: una posible taxonomía para los tipos de datos o variables estadísticos. No es la única taxonomía posible.

Continuos:

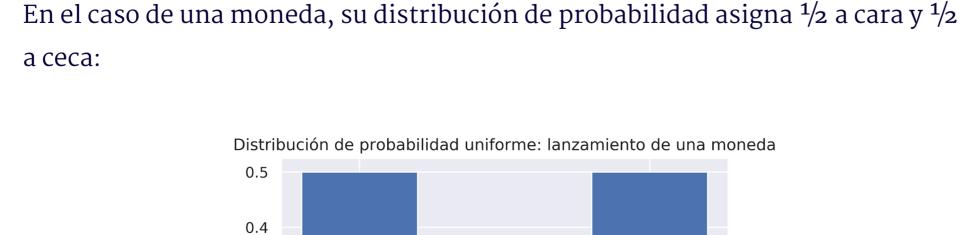
Categóricos

que no puede haber temperaturas menores al cero absoluto; la cantidad de dormitorios de una casa no puede ser menor que cero; el resultado de tirar un dado de seis caras no puede ser menor que uno ni mayor que seis. Para pensar: tiramos dos dados y anotamos la suma de sus resultados. ¿A qué tipo de variable corresponde? ¿Cuál es el valor máximo y mínimo que podemos obtener?

Si estamos trabajando con una variable aleatoria, la mejor descripción que

• su densidad de probabilidad, para el caso de una variable numérica

Veamos un ejemplo: para el caso de un dado de seis caras, la distribución de probabilidad es uniforme, y asigna ¼ de probabilidad a cada cara, como se ve en la figura debajo: Distribución de probabilidad uniforme: lanzamiento de un dado



cara

estadística sobre sus resultados. Entonces, parecería ser que la Probabilidad se relaciona con lo que esperamos observar, mientras que la Estadística se relaciona con el análisis de lo que observamos.

• La **Probabilidad** se ocupa de cuantificar la certeza de eventos futuros, y

• La Estadística se ocupa de analizar eventos pasados. Se pregunta: lo que

observamos, ¿tiene sentido con lo que esperábamos encontrar? Si no es

No te preocupes si la distinción no es del todo clara. Todos/as nos solemos

Podríamos resumir ambas disciplinas como sigue:

así, ¿qué podemos aprender?

debo hacer?

crear modelos de lo que esperamos encontrar.

A lo largo de las bitácoras, iremos viendo muchas distribuciones y densidades

confundir con esto. Para el caso del dado, la Probabilidad nos dice que, si el dado no está cargado, todas sus caras deberían ser igual de probables. Es decir, la distribución de probabilidad es uniforme y le asigna a cada cara una probabilidad de %. En cambio, si el dado está cargado, esperamos que esa distribución sea distinta. Por ejemplo, si está cargado para favorecer al seis, que la probabilidad de ese Una vez que ya tenemos una idea de lo que esperamos ver, debemos hacer los experimentos y hacer Estadística sobre sus resultados. Por ejemplo, tirando muchas veces los dados. En general, para cualquier experimento estadístico, podemos organizarnos de la siguiente manera:

Acá está la respuesta:

La ciencia de la estadística nos ayuda a analizar lo que obtuvimos.

0.2

2 Muestras

0.1 ¿Adivinaste correctamente? Si no fue así, ¿qué creés que falló? Empieza a trabajar en el notebook

¡Prepárate para el próximo encuentro!

Te invitamos a conocer más sobre el tema de esta bitácora.

Challenge

Profundiza

"La estadística no sustituye el juicio." —Henry Clay Sr. Abogado y político estadounidense (1777-1852).

DS_Bitácora_05_Probabili Potencia tu Talento -Comunidades

0.15

Figura 1: tomamos dos dados y anotamos los resultados de tirar cada uno. Luego, graficamos el porcentaje de veces que cae de cada lado. La línea punteada corresponde a lo esperado para un dado perfectamente balanceado. ¿Cuál será el que está cargado? La definición más intuitiva de Probabilidad es la Bayesiana: es medida de la para mañana pero con una probabilidad de 0.75 (75% de probabilidad), ¿quién tiene más confianza que lloverá al día siguiente? Veamos un ejemplo: Si tiramos muchas veces una moneda, ¿qué frecuencia (porcentaje) de caras y cecas esperamos? Si la moneda está balanceada esperamos que, en promedio, 50% de las veces caiga cara y 50% de las veces otros, podríamos igualmente hacer una interpretación frecuentista del pronóstico para el día de mañana y una bayesiana para el lanzamiento de una Al fin de cuentas, la probabilidad se materializa en un número, pero es importante que como data scientists tengamos en claro en qué está basada la interpretación que hacemos de la misma: por un lado, depende de la fuente (¿qué confianza tengo en los datos, en su análisis, etc.?) y, por otro, del contexto (¿cuántas veces podemos tirar un dado? ¿cuántas veces ocurrirá el día de mañana y su pronóstico meteorológico?). Por esto, es importante que tengamos ambas definiciones presentes y que sepamos cuándo y cómo utilizar Independientemente de cuál definición te guste más, algunas propiedades de Ejemplo: en el caso de una moneda, p(s = cara) = 0.5, p(s = ceca) = 0.5. • La suma de las probabilidades sobre todos los posibles resultados es uno. Ejemplo: p(s = cara) + p(s = ceca) = 0.5 + 0.5 = 1Para pensar: ¿cómo sería en el caso de las caras de un dado?

• Un **experimento** es un procedimiento que produce uno de un conjunto de posibles resultados. Tirar un dado, por ejemplo, es un experimento. También podemos considerar un experimento al trabajo que hace un/a considerarse experimentos? ¿Cómo se relaciona esta definición con la tuviéramos dos dados, el espacio muestral sería todas las posibles combinaciones de los resultados de cada uno $S = \{(1,1), (1,2), (1,3), (1,4), ..., (1,2), (1,3), (1,4), ..., (1,2), (1,3), (1,4), ..., (1$ $(2,1), (2,2), (2,3), \dots, (6,3), (6,4), (6,5), (6,6)$ }. Es decir, un total de 36 combinaciones. ¿Cuál es el espacio muestral para: a) medir la temperatura de un objeto y b) tomar una persona al azar en la calle y preguntarle su

son, desde este punto de vista, variables aleatorias. También puede serlo cuántos goles anota Lionel Messi en un partido de la Selección Argentina o

Las variables numéricas discretas son aquellas que se pueden contar, por ejemplo, la cantidad de hijos/as de una familia, cantidad de dormitorios en una casa, número de empleados de una empresa, resultado de tirar un dado, etc. Entre dos valores de una variable numérica discreta no existe otro valor. Las variables numéricas continuas son aquellas que, entre dos valores, siempre puede haber otro. Para pensarlo de otra forma, son aquellas que se miden. Por Otra cosa a la cual podemos prestarle atención es si una variable numérica (discreta o continua) está acotada, es decir, si existe un valor máximo o mínimo que puede tomar. Por ejemplo, las temperaturas están acotadas por debajo, ya

• su distribución de probabilidad, para el caso de una variable numérica discreta. La **distribución de probabilidad** de una variable aleatoria es una función que asigna a cada suceso definido sobre la variable la

Relación entre Probabilidad y Estadística Evidentemente, estas disciplinas están relacionadas, sino no aparecerían siempre juntas. Para ver esta relación, retomemos el ejemplo del dado que te propusimos analizar en la bitácora anterior. Todos tenemos en nuestra mente un modelo probabilístico de un dado de seis caras: asumimos que en un dado no cargado la probabilidad de que salga cada cara es %. Pero para ver si esto se cumple, lo que tenemos que hacer es tirar muchas veces un dado y hacer

ceca

valor sea mayor que un 1/4, a costa de disminuir la probabilidad del resto de los valores. La ciencia de la probabilidad nos ayuda a predecir lo que esperamos ver.

• Adquisición: ¿qué datos necesito?¿En qué cantidad? Para el caso del

• Organización: ¿cómo guardo y agrupo los datos obtenidos?

dado, ¿qué valores anoto de los lanzamientos?¿Cuántos lanzamientos

• Análisis e Interpretación: ¿en qué rango están estos valores? ¿Cómo

ellos que me ayuden a interpretarlos? Los datos, ¿se ajustan o no a lo

son de precisos o de dispersos?¿Qué valores estadísticos puedo obtener de

esperado? Para el caso del dado, ¿qué dicen mis datos sobre si el dado está

cargado o no?¿Hay alguna cara que haya salido más veces de lo esperado? • Presentación: ¿cómo comunico mis conclusiones? ¿Qué números o gráficas utilizo para que se me entienda bien? En el caso de los dos dados, los estadísticos que podemos obtener son relativamente sencillos: cantidad de veces que salió cada cara y su frecuencia.

Al comienzo de la bitácora te preguntamos cuál dado creías que estaba cargado.

- Trabaja en el siguiente notebook, ya sea en Colab o en tu computadora. Necesitarás instalar dos nuevas librerías, Scipy y Matplotlib, cuyo funcionamiento veremos más en detalle en la bitácora 06. ¡Se instalan como todas las librerías que instalamos hasta el momento! Para el encuentro, lleva hecha la sección 1, el resto del notebook lo trabajaremos en clase. Notebook: DS_Bitácora_05_Probabilidad

Te proponemos el siguiente desafío, ¿te animas?

Potencia tu Talento - Comunidades Tips para construir tu perfil profesional.