

Estadística

Promedio n

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

(puede ser



sensible a los outliers.)

Trimmed mean / Promedio recortado (quite elementos del inicio y del fin de los valores y tomando el promedio de los valores restantes)

$$\bar{x} = \frac{\sum_{i=p+1}^{n-p} x_{(i)}}{n-2p}$$

donde las x están ordenadas

Promedio con pesos

$$x_w = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$

Mean Absolute deviation
Error medio absoluto

$$\frac{\sum_{i=1}^n |x_i - \bar{x}|}{n}$$

Deviations / errors / residuals
 $\sum x_i - \sum x_{i \text{ predichos}}$

Varianza / mean-squared-error

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

el promedio

desviación estándar $\sqrt{\text{varianza}} = s$

→ ¿Qué tanto se aleja cada valor en relación al promedio

Estadística

Promedio n

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

(puede ser sensible a los outliers)

Trimmed mean / Promedio recortado (quita elementos del inicio y del fin de los valores y tomando el promedio de los valores restantes)

$$\bar{x} = \frac{\sum_{i=p+1}^{n-p} x_{(i)}}{n-2p}$$

donde las x están ordenadas

Promedio con pesos

$$x_w = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$

Mean Absolute deviation

Error medio absoluto

$$\frac{\sum_{i=1}^n |x_i - \bar{x}|}{n}$$

Deviations / errors / residuals

$$\sum x_i - \sum x_{predicte}$$

varianza / mean-squared-error

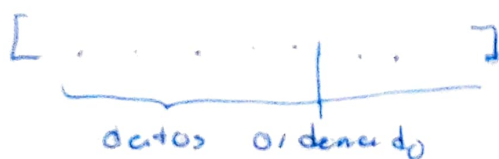
$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

el promedio

desviación estándar $\sqrt{\text{varianza}} = s$

→ ¿Qué tanto se aleja cada valor en relación al promedio

80th percentile



hasta que llega el 80%

80th percentile

50th percentile = media

$$IQR = 75^{th} \text{ percentile} - 25^{th} \text{ percentile}$$

$\{1, 2, 3, 3, 5, 6, 7, 9\}$

$$75^{th} = 6.5 \quad 25^{th} = 2.5$$

$$\Rightarrow IQR = 6.5 - 2.5 = 4$$

Variables correlacionadas

si X , Y están correlacionados significa que X crece ~~de~~ como
y al tiempo de que Y lo hace y si es en sentido
contrario, están correlacionados negativamente

- Coeficiente de correlación de Pearson
estudia la correlación entre dos variables

$$r = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{(n-1) s_x s_y}$$

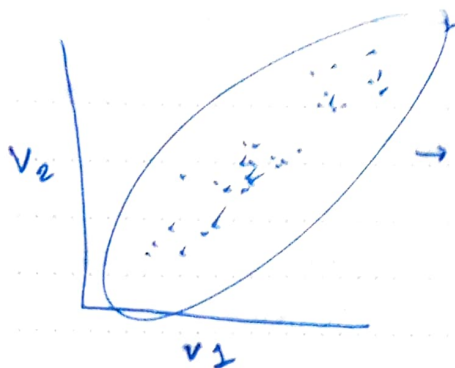
desviación estándar de ambas

Nota:
Estadístico
va de -1
a 1 y 0
implica
que no
hay
correlación

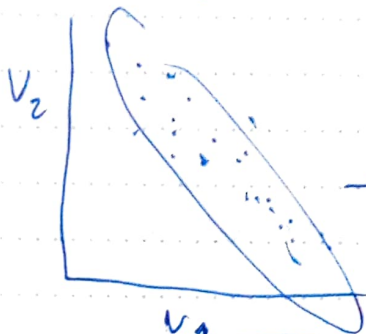
el coeficiente de correlación mide que tanto un par
de variables está relacionado una de la otra

En un scatter plot, donde en cada eje asignamos a cada variable como sigue

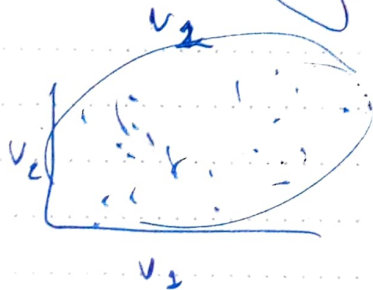
plot.scatter(x='v1', y='v2', ...)



→ si tenemos una distribución de puntos algo así, implica que los valores si están relacionados, de forma positiva



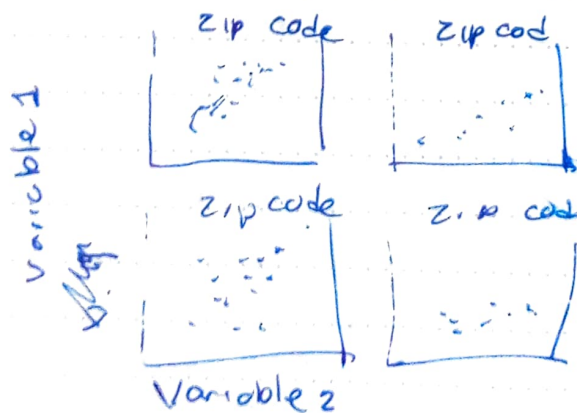
→ correlación negativa



no hay correlación

Análisis Multivariable

Se puede hacer scatter plots con más de dos variables



- Las gráficas de agrupamiento hexagonal y la de contornos (como los topográficos) sirven para encontrar la correlación entre dos variables numéricas
- Las tablas de contingencia son una herramienta para observar las relaciones de dos variables categóricas
- Los diagramas de caja y de violín expresan una relación entre una variable categórica y otra numérica

Datos y distribuciones muestrales

muestra - un subconjunto de los datos (población)

Estrato (división de los datos donde cada elemento comparte cierta característica)

Correr experimentos en varios contextos, bajo diferentes condiciones, analizando qué tanto es cierta nuestra hipótesis, nos ayuda a evitar el sesgo

Sampling distribution of a statistic

distribución de muestra de una estadística

Hablando de los parámetros calculables (o parámetros estadísticos de una población) podemos aproximar estas medidas (como por ejemplo el promedio μ) a partir de calcular este medida por una muestra. Sin embargo, como cada muestra que es posible tomar de la población es diferente entre sí, se calcula de cada una este parámetro (digamos \bar{x}) y esto crea la distribución de una estadística referente a las muestras.

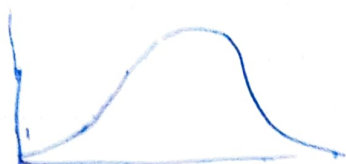
Esto se hace por culpa inferencia o suposiciones de ~~el~~ conjunto de datos aún más grande



distribución de los datos \neq distribución de las muestras

Teorema Central del Limite

Los promedios de las muestras de una población se pueden reorganizar en una distribución normal



Error estándar

$SE = \frac{s}{\sqrt{n}}$ } mide la variabilidad de una métrica muestral
~~es decir, que tanto varía la~~
~~la muestra respecto a una estimación.~~

qué tanto varía o qué tanta variabilidad existe en el cálculo de una estadística en una ~~distribución muestral~~ ~~distribución muestral~~ ~~distribución muestral~~

$$SE = \frac{s}{\sqrt{n}} = \frac{\sqrt{\text{La suma de diferencias de cada valor contra el promedio}}}{\sqrt{n}}$$

SE - variabilidad en estadísticas entre muestras

"si yo tomase una muestra aparte de la misma población, este estadístico que tanto variara respecto a la que obtuve y en general a los demás"
es la variabilidad respecto a la distribución muestral

S - es la distribución de los datos
lo que tanto se alejan del promedio

Nivel de confianza

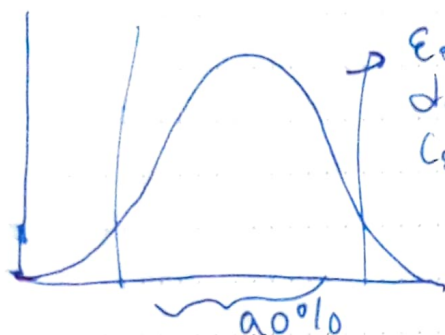
porcentaje de confianza en ~~la~~ ~~que~~ ~~ten~~ exacto
en la métrica que estamos calculando

Note:

Bootstrap
es una técnica de
remuestreo con reposición

El intervalo de confianza ($X\%$)
implica que ~~para~~ ~~lo~~ ~~que~~ ~~está~~
~~dentro~~ ~~de~~ ~~este~~ ~~intervalo~~,
otras muestras, el $X\%$ de
ellos tendrán o darán una
estadística dentro de este intervalo

Entonces



Este es la distribución de una estadística,
digamos \bar{x} por todas las muestras posibles
(obtenidas del bootstrap) así, si
ya establezco un intervalo de
confianza del 90%, significa
que ~~del~~ el 90% del valor de \bar{x}

para las muestras se encontrará en ese
intervalo.

Estandarizar:

$$\frac{x_i - \bar{x}}{s}$$

Regresión y predicción

La regresión lineal se modela

$$Y_i = b_0 + b_1 X_i + e_i \rightarrow \text{intercepto}$$

↑
el error entre lo predicho y el valor real

$$\hat{e}_i = Y_i - \hat{Y}_i$$

Calcular la regresión lineal de los datos implica ajustar una recta que acerque estos.

para un conjunto de V independientes y otros independientes

x y

1 2

2 2.5

0 1

tenemos que encontrar los valores b_0, b_1
que más se aproximen a estos,

Tomando la forma normal de la ec. de la recta

$$H = X\Theta \quad (\text{debemos encontrar } \Theta)$$

Así,

$$\Theta = (X^T X)^{-1} X^T Y \quad (\text{obtenido al minimizar la función de error})$$

En el ejemplo

$$\Theta = \left(\begin{bmatrix} 1 & 2 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & 2 & 0 \end{bmatrix} \begin{bmatrix} 2 \\ 2.5 \\ 1 \end{bmatrix} = \begin{bmatrix} 1.08 \\ 0.25 \end{bmatrix}$$

Validación cruzada

En lugar de dividir 70% para entrenamiento y 30% para
prueba una sola vez y dividiendo el data set arbitrariamente
se realiza este proceso ~~dividiendo~~ usando
diferentes particiones por el conjunto de prueba
de manera en que no existe al final ningún
elemento sin que haya sido ~~testado~~ ~~evaluado~~

Probabilidad Principios Básicos del Conteo

$$y = \bullet \quad 0 = +$$

importante el orden $\xrightarrow{\text{con reemplazo}} {}_nO_k \quad n^k$
 $\xrightarrow{\text{sin reemplazo}} {}_n^P_k \quad \frac{n!}{(n-k)!}$

no importante el orden $\xrightarrow{\text{sin reemplazo}} {}_n^C_k \quad \frac{n!}{(n-k)!k!}$

Conjuntos

$$C \cup (A \cap B) = (C \cup A) \cap (C \cup B)$$

$$(A \cap B)^c = A^c \cup B^c$$

$$(A \cup B)^c = A^c \cap B^c$$

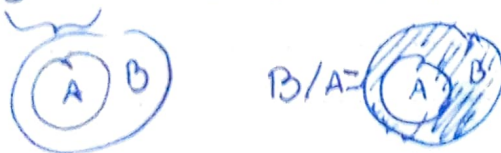
Funcion de probabilidad

$$P = f \Rightarrow [0, 1]$$

$$P\left(\bigcup_{n=1}^{\infty} A_n\right) = \sum_{n=1}^{\infty} P(A_n)$$

$$P(A^c) = 1 - P(A) \quad P(\emptyset) = 0$$

$$P(B \setminus A) = P(B) - P(A)$$



$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

$$P(A \cup B \cup C) = P(A) + P(B) + P(C) - P(A \cap B) - P(A \cap C) - P(B \cap C) + P(A \cap B \cap C)$$

TODOS PODEMOS AHORRAR CON LA NUEVA APP
 SARTÉL (01 55) 13 28 5000 www.gob.mx/consar

**Afore
Móvil**

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

$$P(A \cap B) = P(A|B)P(B)$$

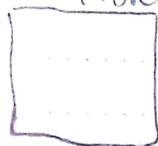
Eventos independientes

$$P(A \cap B) = P(A)P(B)$$

Hadoop

Cluster hadoop

Master hadoop



Workers hadoop



Los masters dividen las operaciones en trozos y los distribuye (las operaciones) entre los trabajadores hadoop.

Entorno ~~de~~ distribuido para aplicaciones distribuidas que trabajan con Big data.

Procesamiento de datos de forma distribuida.

Cuenta con tres partes:

HDFS - hadoop file system.

Planificación de tareas y negociación de recursos: YARN.

Procesamiento distribuido: Map Reduce (para programación distribuida).

Docker

Imagen: archivo que contiene para crear el contenedor.

Volumen: espacio de disco que pueden compartirse entre contenedores o entre un contenedor y el host.

Redes: puentes conectados.

Docker file.

HDFS

almacena archivos grandes

se dividen en bloques que se replican en diferentes
nodos del cluster

1/4

NomeNode - el servidor master

mantiene la información de donde
se encuentran los datos

DataNodes - tiene los datos

Rack

una colección (40-50) DataNodes usando el mismo switch
de red

Rackaware selecciona los data nodes más próximos
para evitar el tráfico de red

Operaciones de hdfs

read

cliente interactúa con el namenode para la metadata

Comandos de hdfs

hdfs dfs -ls _____

" " -cat _____

-appendToFile _____

} de la misma
manera en el
línea con
el prefijo
hdfs

MapReduce

↳ capa de proceso de hadoop

Terminología



map - toma el valor de la llave como entrada y genera un valor de salida asociado a esa llave
el valor donde debemos de procesar

reduce - procesar la data
lo que sale es el final output y se almacena en hadfs

Job - todo el proceso de mapreduce
task - una parte del procesamiento de la data en map o reduce

