

JOINS

→ inner join

→ left join

two datasets A , B

A:

a b

index

0 -

1

2

3

4

B:

c d

index

0 -

2

3

5

inner join: gives you rows with intersection

A inner join B

A.join(B, how='inner')

index val., that are contained in both datasets

index	a	b	c	d
-------	---	---	---	---

0

2

3

left join:

A.join(B, how='left')

0

1

2

3

4

NA NA

outer join

a b c d

0

1

2

3

4

5

-
1

-

x x

Probability

experiments

outcomes - "points" in the set of outcomes

events - sets of outcomes

tossing a coin three times in a row

HHH, HHT, HTH, THH, HTT,

THT, TTH, TTT

sample
events.

{ all the "heads" }

= { HHH }

$$E = \{ \text{first toss is } +1 \}$$

$$= \{ HHH, HHT, HTH, HTT \}$$

$P(\cdot) \rightarrow$ maps events into \mathbb{R} .

$$(1) P(E) \in [0, 1]$$

$$(2) P(S) = 1$$

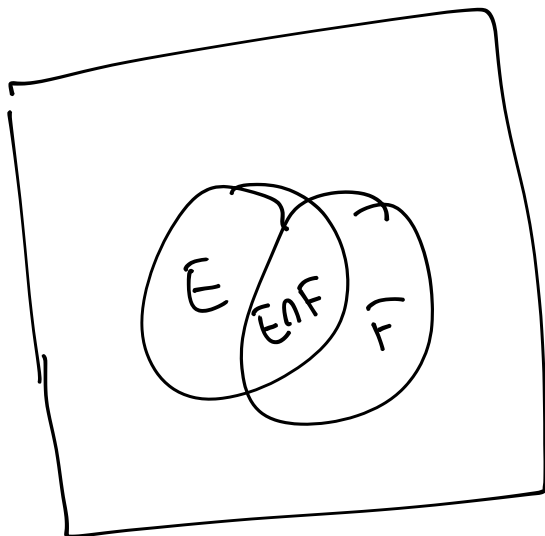
(3) E_1, \dots, E_n disjoint sets

$$E_i \cap E_j = \emptyset$$

$$P(\cup E_i) = \sum P(E_i)$$

E, F not necessarily disjoint

$$P(E \cup F) = P(E) + P(F) - P(E \cap F)$$



$P(E|F)$ = probability that E occurs
given that F has happened

$$= \frac{P(E \cap F)}{P(F)}$$

$$P(E|F) = \frac{P(F|E)P(E)}{P(F)}$$

E_1, \dots, E_n disjoint E, E^c

$$P(A) = \sum_{i=1}^n P(A|E_i)P(E_i)$$

law of total probabilities

$$P(A) = P(A|E)P(E)$$

$$+ P(A|E^c)P(E^c)$$

$$\quad \quad \quad \parallel$$
$$1 - P(E)$$

Σ_x Test gives false positives 1%
false negatives 1%.

$P(\text{really, hy infecta} | \text{positive})$

Random Variable

X with cdf F

random number where

$$P(X \leq x) = F(x)$$

two kinds of random variables:

① discrete \rightarrow finite
 \rightarrow countable

② continuous

\rightarrow take any value on
at least one interval

discrete random variables have a pmf

$$p(x) = P(X = x)$$

$$X \rightarrow \{0, 1, 2\}$$

$$p(0) = 0.5$$

$$p(1) = 0.2$$

$$p(2) = 0.3$$

\rightarrow cdf?

→ went to find the cdf

$$F(x) = \begin{cases} 0 & , x < 0 \\ 0.5 & , 0 \leq x < 1 \\ 0.7 & , 1 \leq x < 2 \\ 1.0 & , x \geq 2 \end{cases}$$

continuous random variables have a pdf / density

$$f(x)$$

$$P(X \in [a, b]) = \int_a^b f(x) dx$$

$$\therefore F(x) = \int_{-\infty}^x f(y) dy$$

$$f(x) = \begin{cases} 2x & , 0 \leq x \leq 1 \\ 0 & , \text{o/w} \end{cases}$$

$$F(x) = \int_{-\infty}^x f(y) dy \quad x \in [0, 1]$$

$$= \int_{-\infty}^0 0 dy + 2 \int_0^x y dy$$

$$= x^2 \quad x \in [0, 1]$$

$$x > 1$$

$$= \int_{-\infty}^0 0 dy + 2 \int_0^1 y dy$$

$$+ \int_1^x 0 \cdot dy = 1$$

Expectation: probability-weighted value

discrete random variables

$$E[X] = \sum_{p(x) > 0} x \cdot p(x)$$

Ex $p(0) = 0.5$

$$p(1) = 0.2$$

$$p(2) = 0.3$$

$$\begin{aligned} E[X] &= 0 \cdot 0.5 + (1)(0.2) \\ &\quad + (2)(0.3) \\ &= 0.8 \end{aligned}$$

continuous random variable

$$E[X] = \int_{-\infty}^{\infty} x \cdot f(x) dx$$

$$f(x) = \begin{cases} 2x, & 0 < x < 1 \\ 0, & \text{else} \end{cases}$$

$$E[X] = \int_0^1 x \cdot 2x dx$$

$$= 2 \int_0^1 x^2 dx = \frac{2}{3}$$

the expectation of a function of a cv $g(x)$

$$\Rightarrow \sum g(x) \cdot p(x)$$

discrete

$$\int_{-\infty}^{\infty} g(x) \cdot f(x) dx$$

continuous

$$p(0) = .5$$

$$p(1) = .2$$

$$p(2) = .3$$

$$g(x) = x^3 - 1$$

$$E[g(x)]$$

$$= (0^3 - 1)(.5)$$

$$+ (1^3 - 1)(.2)$$

$$+ (2^3 - 1)(.3)$$

= a number! 😊

$$\int_{-\infty}^{\infty} g(x) \cdot f(x) dx \quad (\text{continuous})$$

$$= 2 \int_0^1 (x^3 - 1) \cdot x dx$$

⇒ Corollary $aX + b$

$$E[aX + b] = aE[X] + b$$

$$= E[(X - E[X])^2]$$

$$\text{Var}(X) = \sum (x - E[X])^2 \cdot p(x)$$

⇒ the average (probabil. weights)
squared distance from the
expectation

$$= E[X^2] - \underbrace{E[X]^2}$$

$$p(0) = .5$$

$$p(1) = .2$$

$$p(2) = .3$$

$$E(X) = 0.0$$

$$Var(X) = \underline{E[X^2]} - \underline{0.8^2}$$

$$E[X^2] = 0^2 \cdot 0.5 + (1)^2 \cdot (.2) + (2)^2 \cdot (.3)$$

$$= 1.4$$

$$Var(X) = 1.4 - .64$$

$$Var(X) = E[(X - E[X])^2]$$

$$= E[X^2] - E[X]^2$$

$$E[X^2] = \int_{-\infty}^{\infty} x^2 \cdot f(x) dx$$

$$E[X^2] = 2 \int_0^1 x^2 \cdot x dx$$

$$= 2 \int_0^1 x^3 dx$$

$$= \frac{1}{2}$$

X, Y both discrete

joint pmf

$$p(x, y) = P(X=x, Y=y)$$

$$E[g(x, y)] = \sum_y \sum_x g(x, y) p(x, y)$$

$X, Y \rightarrow \{0, 1\}$

- $p(\underline{0}, \underline{0}) = 0.5$

- $p(1, \underline{0}) = 0.0$

- $p(0, 1) = 0.4$

- $p(1, 1) = 0.1$

$$\begin{aligned} E[x \cdot y] &= \underline{0.0} \cdot \cancel{p(\underline{0}, \underline{0})} \\ &\quad + \cancel{(1) \cdot 0 \cdot p(1, \underline{0})} \\ &\quad + \cancel{0 \cdot 1 \cdot p(0, 1)} \\ &\quad + (1)(1) p(1, 1) \\ &= 0.1 \end{aligned}$$

$f(x, y) \rightarrow$ joint density

$$P(x \in [a, b], y \in [c, d])$$

$$= \int_c^d \int_a^b f(x, y) dx dy$$

$$E[g(x, y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) \cdot f(x, y) dx dy$$

$$Cov(X, Y) = E[XY] - E[X] \cdot E[Y]$$

$$p(0, 0) = 0.5$$

$$\rightarrow p(1, 0) = 0.0$$

$$p(0, 1) = 0.4$$

$$\rightarrow p(1, 1) = 0.1$$

$$E[XY] = 0.1$$

$$E[X] = 0.1$$

$$E[Y] = 0.5$$

$$Cov(X, Y) = 0.1 - (0.1)(0.5)$$

$$= .05$$

$$Corr(X, Y) = \frac{Cov(X, Y)}{\sqrt{SD(X)} \sqrt{SD(Y)}}$$

$$SD(X) = \sqrt{Var(X)}$$

$\text{Bern}(p)$ $p \in [0, 1]$

$$f(0) = 1-p$$

$$f(1) = p$$

$X \sim \text{Bin}(n, p)$: sum of n independent
 $\text{Bern}(p)$

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

$\Rightarrow X, Y$ are independent if for all sets

$$P(X \in A, Y \in B) = P(X \in A) P(Y \in B)$$

$$p(x, y) = p_x(x) \cdot p_y(y)$$

$$f(x, y) = f_x(x) \cdot f_y(y)$$

$$f(x, y) = cxy$$

$$f(x, y) = x + y$$

Normal random variable

$$E[X] = \mu$$

$$\text{Var}(X) = \sigma^2$$

$$f(x) = \frac{1}{\sqrt{2\pi} \cdot \sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Central Limit Theorem

X_1, \dots, X_n independent identically distributed

$$E[X_i] = \mu$$

$$\text{Var}(X_i) = \sigma^2$$

n large

$$\bar{X}_n \underset{\text{approx}}{\sim} N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

