2. The entropy of a discrete random variable $X$ is defined as (use base $e$ for all log operations unless specified otherwise):

$$H(X) = -\sum_{x \in X} P(x) \log P(x)$$

(a) Compute the entropy of the distribution $P(x) = \text{Multinomial}([0.2, 0.3, 0.5])$. [3 pts]

(b) Compute the entropy of the uniform distribution $P(x) = \frac{1}{m} \forall x \in [1, m]$. [3 pts]

(c) Consider the entropy of the joint distribution P(X, Y):

$$H(X, Y) = -\sum_{x \in X} \sum_{y \in Y} P(x, y) \log P(x, y)$$

How does this entropy relate to H(X) and H(Y), (i.e. the entropies of the marginal distributions) when X and Y are independent? [4 pts]

Solution: Next Page

$2. \qquad H(x) = -\sum_{x \in Y} P(x) \log P(x)$

$(a): \quad P(x) = \text{Multinoulli} ([0.2, 0.3, 0.5])$

$\quad \text{Entropy} \rightarrow H(x) = -[0.2 \log 0.2, 0.3 \log 0.3, 0.5 \log 0.5]$

$\qquad = [0.32, 0.36, 0.34]$

$(b): \quad \text{Entropy of uniform distribution}$

$\qquad P(x) = \frac{1}{m} \forall x \in [1, m]$

$\quad H(x) = \int_{1}^{m} -P(x) \log P(x)$

$\qquad = -\int_{1}^{m} \frac{1}{m} \log m \, dm$

$\quad \log m = u$

$\quad dm = dv$

$\quad \frac{1}{m} dm = du$

$\qquad = -\int_{1}^{m} u \, du$

$\qquad = -\int_{0}^{e^u} u \, du = -\left(\frac{u^2}{2}\right)_{0}^{e^u}$

$\qquad = -\left[\frac{(e^u)^2}{2} - 0\right]$

$\qquad = \frac{(e^u)^2}{2} = \frac{m^2}{2}$

$(c): \quad \text{Entropy of Joint Prob. distribution. } P(x, y)$

$\quad H(x, y) = -\sum \sum P(x, y) \log P(x, y) \quad -(i)$

$\quad \text{when} \quad x, \text{ and } y \text{ are independent than:}$

$\qquad H(y|x) = H(y) \quad -(1)$

$\qquad H(x|y) = H(x)$

$\quad \text{we know that:}$

$\qquad H(x, y) = H(x) + H(y|x)$

$\quad \text{from equation (1)}$

$\qquad H(x, y) = H(x) + H(y)$

3. You are investigating articles from the New York Times and from Buzzfeed. Some of the articles contain *fake* news, while others contain *real* news (assume that there are only two types of news).
*Note*: for the following questions, write your answer using up to 3 significant figures.

   (a) Fake news only accounts for 5% of all articles in all newspapers. However, it is known that 30% of all fake news comes from Buzzfeed. In addition, Buzzfeed generates 25% of all news articles. What is the probability that a randomly chosen Buzzfeed article is fake news? [3 pts]

   (b) Suppose that 15% of all fake news comes from the New York Times (NYT). Furthermore, suppose that 60% of all real news comes from the NYT. Under all assumptions so far, what is the probability that a randomly chosen NYT article is fake news? [3 pts]

   (c) Mike is an active reader of the New York Times: Mike reads 80% of all NYT articles. However, he also has a suspicion that the NYT is a bad publisher, and he believes that 25% of all NYT articles are fake news. Furthermore, the NYT generates 30% of all news articles. Under all assumptions so far, what is the probability that a randomly chosen article (from all newspapers) will be from the NYT, will be read by Mike and will be *believed* to be fake news? [4 pts]

..........
Solution: Next page

(3): (a) fake news = 5%.

30% comes from Buzzfeed. = P (fake, Buzzfeed)

25% news comes from Buzz feed. = P (Buzz feed)

P (fake new| Buzzfeed article) = $\frac{P(\text{fake new, Buzzfeed})}{P(\text{Buzzfeed article})}$

$= \frac{0.05 \times 0.3}{0.25} = \frac{0.3 \times 50}{250}$

$= 0.06$

(b) 15% all fake News comes from NYT.

60% al real " " " " " .

P(fake New| NYT) = $\frac{P(\text{fake News, NYT})}{P(\text{NYT})}$

$= \frac{5 \times 15/100}{95 \times \frac{60}{100}} = \frac{0.75}{57.75}$

$= 0.013$

(c) P( Mike read | NYT ) = 0.8

P ( NYT ) = 0.3

P ( Believed fake | NYT) = 0.25

P ( NYT and Read and Believed false)

$= $ P(read new| fake, NYT) * P(fake news| NYT) * P(N

$= 0.8 \times 0.25 \times 0.3$

$= 0.06$

4. Suppose we have a probability density function (pdf) defined as:

$$f(x,y) = \begin{cases} C(x^2 + 2y), & 0 < x < 1 \text{ and } 0 < y < 1, \\ 0, & \text{otherwise.} \end{cases}$$

(a) Find the value of $C$. [2pts]

(b) Find the marginal distribution of $X$ and $Y$. [4pts]

(c) Find the joint cumulative density function (cdf) of $X$ and $Y$. [4pts]

Solution:



$f(x,y) = \begin{cases} C(x^2+2y) &, 0 < x < 1 \\ & 0 < y < 1 \\ 0 &, \text{otherwise} \end{cases}$

① Value of C.

Summation of all the Prob. across PDF $= 1$

$\int_0^1 \int_0^1 C(x^2+2y)\, dx\, dy = 1$

$\int_0^1 \int_0^1 C(x^2 dx\, dy + 2y\, dx\, dy) = 1$

$\int_0^1 C\left[\frac{x^3}{3}\right]_0^1 dy + 2C[x]_0^1 y\, dy = 1$

$\int_0^1 \frac{C}{3}\, dy + 2C\, y\, dy = 1$

$\frac{C}{3}[y]_0^1 + 2C\left[\frac{y^2}{2}\right]_0^1 = 1$

$\frac{C}{3} + \frac{2C}{2} = 1$

$\frac{4C}{3} = 1$

$C = \frac{3}{4}$

②: Marginal distribution of X and Y.

$f_X(x) = \int_0^1 C(x^2+2y)\, dy$

$= \frac{3}{4} x^2[y]_0^1 + \frac{3}{4} \cdot 2 \cdot \left[\frac{y^2}{2}\right]_0^1$

$= \frac{3x^2}{4} + \frac{6}{4}\cdot\frac{1}{2}x^2$

$= \frac{3}{4}(x^2+1)$

Thus,

$f_X(x) = \begin{cases} \frac{3}{4}(x^2+1) &, 0 < x < 1 \\ 0 &, \text{otherwise} \end{cases}$

$f_Y(Y) = \int_0^1 \frac{3}{4}(x^2+2y)\, dx$

$= \frac{3}{4}\left[\frac{x^3}{3}\right]_0^1 + \frac{3}{4}\cdot 2 [x]_0^1 y$

$= \frac{1}{4} + \frac{6}{4}y$

$= \frac{1}{4}(6y+1)$

Thus,

$f_Y(Y) = \begin{cases} \frac{1}{4}(6y+1) &, 0 \le y \le 1 \\ 0 &, \text{otherwise} \end{cases}$

5. **[Graduate Students Only]** A 2-D Gaussian distribution is defined as:

$$G(x,y) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{x^2+y^2}{2\sigma^2}\right)$$

Compute the following integral:

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} G(x,y)\left(5x^2y^2 + 3xy + 1\right) dx\, dy$$

*Hint*: Think in terms of the properties of probability distribution functions. [5 pts]

Solution:

⑤: 2-D Gaussian Distribution

$$G(x,y) = \frac{1}{2\pi\sigma^2} e^{-\left(\frac{x^2+y^2}{2\sigma^2}\right)}$$

Compute:

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} G(x,y)\,(5x^2y^2 + 3xy + 1)\, dx\, dy$$

for a standard 2-D gaussian distribution:

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} G(x,y)\, dx\, dy = 1 \quad - (1)$$

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} G(x,y)\, x\, y\, dx\, dy = 0 = E(X,Y) \quad -(2)$$

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} G(x,y)\, x^2\, y^2\, dx\, dy = E(x^2,y^2) = Var(x^2y^2) + [E/x\hat{s}]$$

$$= \sigma^4$$

Hence

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} 5G(x,y)\, x^2 y^2\, dx\, dy + \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} 3\, G(x,y)\, xy\, dx\, dy + \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} G(x,y)\, dx\, dy$$

$$= 5\sigma^4 + 0 + 1$$

$$= 5\sigma^4 + 1$$