

Foundations of Probabilities and Information Theory for Machine Learning

Random Variables

Some proofs

$$E[X + Y] = E[X] + E[Y]$$

where X and Y are random variables of the same type (i.e. either discrete or cont.)

The discrete case:

$$\begin{aligned} E[X + Y] &= \sum_{\omega \in \Omega} (X(\omega) + Y(\omega)) \cdot P(\omega) \\ &= \sum_{\omega} X(\omega) \cdot P(\omega) + \sum_{\omega} Y(\omega) \cdot P(\omega) = E[X] + E[Y] \end{aligned}$$

The continuous case:

$$\begin{aligned} E[X + Y] &= \int_x \int_y (x + y) p_{XY}(x, y) dy dx \\ &= \int_x \int_y x p_{XY}(x, y) dy dx + \int_x \int_y y p_{XY}(x, y) dy dx \\ &= \int_x x \int_y p_{XY}(x, y) dy dx + \int_y y \int_x p_{XY}(x, y) dx dy \\ &= \int_x x p_X(x) dx + \int_y y p_Y(y) dy = E[X] + E[Y] \end{aligned}$$

X and Y are independent $\Rightarrow E[XY] = E[X] \cdot E[Y]$,

X and Y being random variables of the same type (i.e. either discrete or continuous)

The discrete case:

$$\begin{aligned} E[XY] &= \sum_{x \in \text{Val}(X)} \sum_{y \in \text{Val}(Y)} xy P(X = x, Y = y) = \sum_{x \in \text{Val}(X)} \sum_{y \in \text{Val}(Y)} xy P(X = x) \cdot P(Y = y) \\ &= \sum_{x \in \text{Val}(X)} \left(x P(X = x) \sum_{y \in \text{Val}(Y)} y P(Y = y) \right) = \sum_{x \in \text{Val}(X)} x P(X = x) E[Y] = E[X] \cdot E[Y] \end{aligned}$$

The continuous case:

$$\begin{aligned} E[XY] &= \int_x \int_y xy p(X = x, Y = y) dy dx = \int_x \int_y xy p(X = x) \cdot p(Y = y) dy dx \\ &= \int_x x p(X = x) \left(\int_y y p(Y = y) dy \right) dx = \int_x x p(X = x) E[Y] dx \\ &= E[Y] \cdot \int_x x p(X = x) dx = E[X] \cdot E[Y] \end{aligned}$$

Discrete random variables:
independence, conditional independence, conditional probabilities

CMU, 2015 spring, T. Mitchell, N. Balcan, HW2, pr. 1.d

Let X , Y , and Z be random variables taking values in $\{0, 1\}$. The following table lists the probability of each possible assignment of 0 and 1 to the variables X , Y , and Z :

	$Z = 0$		$Z = 1$	
	$X = 0$	$X = 1$	$X = 0$	$X = 1$
$Y = 0$	$1/15$	$1/15$	$4/15$	$2/15$
$Y = 1$	$1/10$	$1/10$	$8/45$	$4/45$

For example,

$P(X = 0, Y = 1, Z = 0) = 1/10$ and $P(X = 1, Y = 1, Z = 1) = 4/45$.

- Is X independent of Y ? Why or why not?
- Is X conditionally independent of Y given Z ? Why or why not?
- Calculate $P(X = 0 \mid X + Y > 0)$.

Answer

a. No.

$$P(X = 0) = 1/15 + 1/10 + 4/15 + 8/45 = 11/18,$$

$$P(Y = 0) = 1/15 + 1/15 + 4/15 + 2/15 = 8/15,$$

and

$$P(X = 0|Y = 0) = \frac{P(X = 0, Y = 0)}{P(Y = 0)} = \frac{1/15 + 4/15}{8/15} = \frac{5}{8}.$$

Since $P(X = 0)$ does not equal $P(X = 0|Y = 0)$, X is not independent of Y .

b. For all pairs $y, z \in \{0, 1\}$, we need to check that $P(X = 0|Y = y, Z = z) = P(X = 0|Z = z)$. That the other probabilities are equal follows from the law of total probability.

$$P(X = 0|Y = 0, Z = 0) = \frac{1/15}{1/15 + 1/15} = 1/2$$

$$P(X = 0|Y = 1, Z = 0) = \frac{1/10}{1/10 + 1/10} = 1/2$$

$$P(X = 0|Y = 0, Z = 1) = \frac{4/15}{4/15 + 2/15} = 2/3$$

$$P(X = 0|Y = 1, Z = 1) = \frac{8/45}{8/45 + 4/15} = 2/3.$$

and

$$P(X = 0|Z = 0) = \frac{1/15 + 1/10}{1/15 + 1/15 + 1/10 + 1/10} = 1/2$$

$$P(X = 0|Z = 1) = \frac{4/15 + 8/45}{4/15 + 2/15 + 8/45 + 4/45} = 2/3.$$

This shows that X is independent of Y given Z .

c.

$$P(X = 0|X + Y > 0) = \frac{1/10 + 8/45}{1/15 + 1/10 + 1/10 + 2/15 + 4/45 + 8/45} = 5/12.$$

Probabilistic Distributions

Some properties

Binomial distribution: $b(r; n, p) \stackrel{\text{def.}}{=} C_n^r p^r (1 - p)^{n-r}$

Significance: $b(r; n, p)$ is the probability of drawing r *heads* in n independent flips of a coin having the head probability p .

$b(r; n, p)$ indeed represents a **probability distribution**:

- $b(r; n, p) = C_n^r p^r (1 - p)^{n-r} \geq 0$ for all $p \in [0, 1]$, $n \in \mathbb{N}$ and $r \in \{0, 1, \dots, n\}$,
- $\sum_{r=0}^n b(r; n, p) = 1$:

$$(1 - p)^n + C_n^1 p (1 - p)^{n-1} + \dots + C_n^{n-1} p^{n-1} (1 - p) + p^n = [p + (1 - p)]^n = 1$$

Binomial distribution: calculating the mean

$$\begin{aligned}
 E[b(r; n, p)] &\stackrel{\text{def.}}{=} \sum_{r=0}^n r \cdot b(r; n, p) = \\
 &= 1 \cdot C_n^1 p(1-p)^{n-1} + 2 \cdot C_n^2 p^2(1-p)^{n-2} + \dots + (n-1) \cdot C_n^{n-1} p^{n-1}(1-p) + n \cdot p^n \\
 &= p [C_n^1(1-p)^{n-1} + 2 \cdot C_n^2 p(1-p)^{n-2} + \dots + (n-1) \cdot C_n^{n-1} p^{n-2}(1-p) + n \cdot p^{n-1}] \\
 &= np [(1-p)^{n-1} + C_{n-1}^1 p(1-p)^{n-2} + \dots + C_{n-1}^{n-2} p^{n-2}(1-p) + C_{n-1}^{n-1} p^{n-1}] \quad (1) \\
 &= np[p + (1-p)]^{n-1} = np \quad (2)
 \end{aligned}$$

For the (1) equality we used the following property:

$$\begin{aligned}
 k C_n^k &= k \frac{n!}{k! (n-k)!} = \frac{n!}{(k-1)! (n-k)!} = \frac{n(n-1)!}{(k-1)! (n-1-(k-1))!} \\
 &= n C_{n-1}^{k-1}, \forall k = 1, \dots, n.
 \end{aligned}$$

Binomial distribution: calculating the variance

following www.proofwiki.org/wiki/Variance_of_Binomial_Distribution, which cites
 “Probability: An Introduction”, by Geoffrey Grimmett and Dominic Welsh,
 Oxford Science Publications, 1986

We will make use of the formula $\text{Var}[X] = E[X^2] - E^2[X]$.

By denoting $q = 1 - p$, it follows:

$$\begin{aligned}
 E[b^2(r; n, p)] &\stackrel{\text{def.}}{=} \sum_{r=0}^n r^2 C_n^r p^r q^{n-r} = \sum_{r=0}^n r^2 \frac{n(n-1) \dots (n-r+1)}{r!} p^r q^{n-r} \\
 &= \sum_{r=1}^n r n \frac{(n-1) \dots (n-r+1)}{(r-1)!} p^r q^{n-r} = \sum_{r=1}^n r n C_{n-1}^{r-1} p^r q^{n-r} \\
 &= np \sum_{r=1}^n r C_{n-1}^{r-1} p^{r-1} q^{(n-1)-(r-1)}
 \end{aligned}$$

Binomial distribution: calculating the variance (cont'd)

By denoting $j = r - 1$ and $m = n - 1$, we'll get:

$$\begin{aligned}
 E[b^2(r; n, p)] &= np \sum_{j=0}^m (j+1) C_m^j p^j q^{m-j} \\
 &= np \left[\underbrace{\sum_{j=0}^m j C_m^j p^j q^{m-j}}_{E[b(r; n-1, p)], \text{ cf. (2)}} + \underbrace{\sum_{j=0}^m C_m^j p^j q^{m-j}}_1 \right].
 \end{aligned}$$

Therefore,

$$E[b^2(r; n, p)] = np[(n-1)p + 1] = n^2p^2 - np^2 + np.$$

Finally,

$$\text{Var}[X] = E[b^2(r; n, p)] - (E[b(r; n, p)])^2 = n^2p^2 - np^2 + np - n^2p^2 = np(1-p)$$

Binomial distribution: calculating the variance

Another solution

- se demonstrează relativ ușor că orice variabilă aleatoare urmând distribuția binomială $b(r; n, p)$ poate fi văzută ca o sumă de n variabile independente care urmează distribuția Bernoulli de parametru p ; ^a
- știm (sau, se poate dovedi imediat) că varianța distribuției Bernoulli de parametru p este $p(1 - p)$;
- ținând cont de proprietatea de liniaritate a varianțelor — $Var[X_1 + X_2 + \dots + X_n] = Var[X_1] + Var[X_2] + \dots + Var[X_n]$, dacă X_1, X_2, \dots, X_n sunt variabile independente —, rezultă că $Var[X] = np(1 - p)$.

^aVezi www.proofwiki.org/wiki/Bernoulli_Process_as_Binomial_Distribution, care citează de asemenea ca sursă “Probability: An Introduction” de Geoffrey Grimmett și Dominic Welsh, Oxford Science Publications, 1986.

The *categorical* distribution:
Computing *probabilities* and *expectations*

CMU, 2009 fall, Geoff Gordon, HW1, pr. 4

Suppose we have n bins and m balls. We throw balls into bins independently at random, so that each ball is equally likely to fall into any of the bins.

- a. What is the probability of the first ball falling into the first bin?
- b. What is the expected number of balls in the first bin?

Hint (1): Define an indicator random variable representing whether the i -th ball fell into the first bin:

$$X_i = \begin{cases} 1 & \text{if } i\text{-th ball fell into the first bin;} \\ 0 & \text{otherwise.} \end{cases}$$

Hint (2): Use linearity of expectation.

- c. What is the probability that the first bin is empty?
- d. What is the expected number of empty bins? *Hint (3):* Define an indicator for the event “bin j is empty” and use linearity of expectations.

Answer

a. It is equally likely for a ball to fall in any of the bins, so the probability that first ball falling into the first bin is $1/n$.

b. Define X_i as described in Hint 1. Let Y be the total number of balls that fall into first bin: $Y = X_1 + \dots + X_m$. The expected number of balls is:

$$E[Y] = \sum_{i=1}^m E[X_i] = \sum_{i=1}^m 1 \cdot P(X_i = 1) = m \cdot 1/n = m/n$$

c. Let Y and X_i be the same as defined at point b. For the first bin to be empty none of the balls should fall into the first bin: $Y = 0$.

$$P(Y = 0) = P(X_1 = 0, \dots, X_m = 0) = \prod_{i=1}^m P(X_i = 0) = (1 - 1/n)^m = \left(\frac{n-1}{n}\right)^m$$

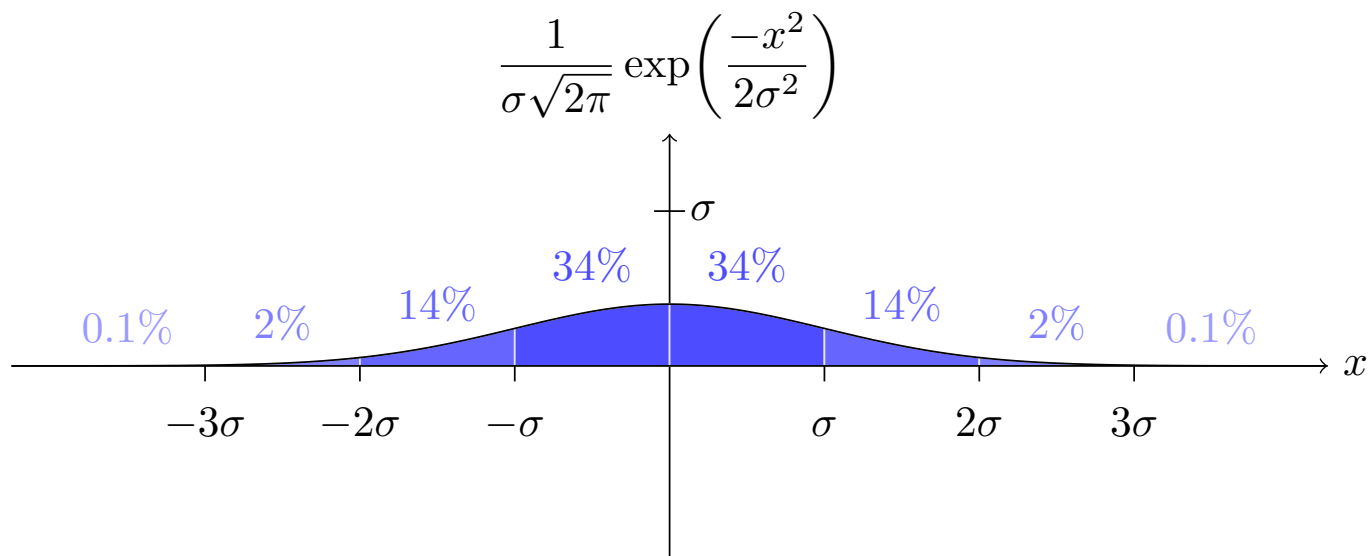
d. For each one of the n bins, we define an indicator random variable Y_j for the event “bin j is empty”. Let Z be the random variables denoting the number of empty bins. $Z = Y_1 + \dots + Y_n$. Then the expected number of empty bins is:

$$E[Z] = E\left[\sum_{j=1}^n Y_j\right] = \sum_{j=1}^n E[Y_j] = \sum_{j=1}^n 1 \cdot P(Y_j) = \sum_{j=1}^n \left(\frac{n-1}{n}\right)^m = n \left(\frac{n-1}{n}\right)^m$$

The univariate Gaussian distribution:

$$\mathcal{N}_{\mu,\sigma}(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

The plot, when $\mu = 0$:



Source: <http://www.texample.net/tikz/examples/standard-deviation/>

Proving that $\mathcal{N}_{\mu,\sigma}$ is indeed a p.d.f.

1. $\mathcal{N}_{\mu,\sigma}(x) \geq 0 \ \forall x \in \mathbb{R}$ (**true**)

2. $\int_{-\infty}^{\infty} \mathcal{N}_{\mu,\sigma}(x) dx = 1$

Note: Concerning the second property, it is enough to prove it for the *standard* case ($\mu = 0$, $\sigma = 1$), because the non-standard case can be reduced to this one:

Using the variable transformation $v = \frac{x - \mu}{\sigma}$ will imply $x = \sigma v + \mu$ and $dx = \sigma dv$, so:

$$\begin{aligned} \int_{-\infty}^{\infty} \mathcal{N}_{\mu,\sigma}(x) dx &= \int_{x=-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = \frac{1}{\sqrt{2\pi}\sigma} \int_{x=-\infty}^{\infty} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx \\ &= \frac{1}{\sqrt{2\pi}\sigma} \int_{v=-\infty}^{\infty} e^{-\frac{v^2}{2}} \sigma dv = \frac{1}{\sqrt{2\pi}\sigma} \sigma \int_{v=-\infty}^{\infty} e^{-\frac{v^2}{2}} dv = \frac{1}{\sqrt{2\pi}} \int_{x=-\infty}^{\infty} \mathcal{N}_{0,1}(x) dx \end{aligned}$$

The *standard* case: proving that $\mathcal{N}_{0,1}$ is indeed a p.d.f.:

19.

$$\begin{aligned} \left(\int_{v=-\infty}^{\infty} e^{-\frac{v^2}{2}} dv \right)^2 &= \left(\int_{x=-\infty}^{\infty} e^{-\frac{x^2}{2}} dx \right) \cdot \left(\int_{y=-\infty}^{\infty} e^{-\frac{y^2}{2}} dy \right) = \int_{x=-\infty}^{\infty} \int_{y=-\infty}^{\infty} e^{-\frac{x^2+y^2}{2}} dy dx \\ &= \iint_{\mathbb{R}^2} e^{-\frac{x^2+y^2}{2}} dy dx \end{aligned}$$

By switching from x, y to polar coordinates r, θ (see the *Note* below), it follows:

$$\begin{aligned} \left(\int_{v=-\infty}^{\infty} e^{-\frac{v^2}{2}} dv \right)^2 &= \int_{r=0}^{\infty} \int_{\theta=0}^{2\pi} e^{-\frac{r^2}{2}} (r dr d\theta) = \int_{r=0}^{\infty} r e^{-\frac{r^2}{2}} \left(\int_{\theta=0}^{2\pi} d\theta \right) dr = \int_{r=0}^{\infty} r e^{-\frac{r^2}{2}} \theta|_0^{2\pi} dr \\ &= 2\pi \int_{r=0}^{\infty} r e^{-\frac{r^2}{2}} dr = 2\pi \left(-e^{-\frac{r^2}{2}} \right) \Big|_0^{\infty} = 2\pi(1 - 0) = 2\pi \Rightarrow \int_{v=-\infty}^{\infty} e^{-\frac{v^2}{2}} dv = \sqrt{2\pi} \Rightarrow \int_{v=-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{v^2}{2}} dv = 1. \end{aligned}$$

Note: $x = r \cos \theta$ and $y = r \sin \theta$, with $r \geq 0$ and $\theta \in [0, 2\pi)$. Therefore, $x^2 + y^2 = r^2$, and the Jacobian matrix is

$$\frac{\partial(x, y)}{\partial(r, \theta)} = \begin{vmatrix} \frac{\partial x}{\partial r} & \frac{\partial x}{\partial \theta} \\ \frac{\partial y}{\partial r} & \frac{\partial y}{\partial \theta} \end{vmatrix} = \begin{vmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{vmatrix} = r \cos^2 \theta + r \sin^2 \theta = r \geq 0. \text{ So, } dx dy = r dr d\theta.$$

Calculating the mean

$$E[\mathcal{N}_{\mu,\sigma}(x)] \stackrel{\text{def.}}{=} \int_{-\infty}^{\infty} x \mathcal{N}_{\mu,\sigma}(x) dx = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} x \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx$$

Using again the variable transformation $v = \frac{x-\mu}{\sigma}$ will imply:

$$\begin{aligned} E[X] &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} (\sigma v + \mu) e^{-\frac{v^2}{2}} (\sigma dv) = \frac{\sigma}{\sqrt{2\pi}\sigma} \left(\sigma \int_{-\infty}^{\infty} v e^{-\frac{v^2}{2}} dv + \mu \int_{-\infty}^{\infty} e^{-\frac{v^2}{2}} dv \right) \\ &= \frac{1}{\sqrt{2\pi}} \left(-\sigma \int_{-\infty}^{\infty} (-v) e^{-\frac{v^2}{2}} dv + \mu \int_{-\infty}^{\infty} e^{-\frac{v^2}{2}} dv \right) = \frac{1}{\sqrt{2\pi}} \left(\underbrace{-\sigma e^{-\frac{v^2}{2}}}_{=0} \Big|_{-\infty}^{\infty} + \mu \int_{-\infty}^{\infty} e^{-\frac{v^2}{2}} dv \right) \\ &= \frac{\mu}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{v^2}{2}} dv = \frac{\mu}{\sqrt{2\pi}} \sqrt{2\pi} = \mu \end{aligned}$$

Calculating the variance

We will make use of the formula $\text{Var}[X] = E[X^2] - E^2[X]$.

$$E[X^2] = \int_{-\infty}^{\infty} x^2 \mathcal{N}_{\mu, \sigma}(x) dx = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} x^2 \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx$$

Again, using the transformation $v = \frac{x-\mu}{\sigma}$ will imply $x = \sigma v + \mu$ and $dx = \sigma dv$. Therefore,

$$\begin{aligned} E[X^2] &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} (\sigma v + \mu)^2 e^{-\frac{v^2}{2}} (\sigma dv) \\ &= \frac{\sigma}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} (\sigma^2 v^2 + 2\sigma\mu v + \mu^2) e^{-\frac{v^2}{2}} dv \\ &= \frac{1}{\sqrt{2\pi}} \left(\sigma^2 \int_{-\infty}^{\infty} v^2 e^{-\frac{v^2}{2}} dv + 2\sigma\mu \int_{-\infty}^{\infty} v e^{-\frac{v^2}{2}} dv + \mu^2 \int_{-\infty}^{\infty} e^{-\frac{v^2}{2}} dv \right) \end{aligned}$$

Note that we have already computed $\int_{-\infty}^{\infty} v e^{-\frac{v^2}{2}} dv = 0$ and $\int_{-\infty}^{\infty} e^{-\frac{v^2}{2}} dv = \sqrt{2\pi}$.

Calculating the variance (Cont'd)

Therefore, we only need to compute

$$\begin{aligned} \int_{-\infty}^{\infty} v^2 e^{-\frac{v^2}{2}} dv &= \int_{-\infty}^{\infty} (-v) \left(-v e^{-\frac{v^2}{2}} \right) dv = \int_{-\infty}^{\infty} (-v) \left(e^{-\frac{v^2}{2}} \right)' dv \\ &= (-v) e^{-\frac{v^2}{2}} \Big|_{-\infty}^{\infty} - \int_{-\infty}^{\infty} (-1) e^{-\frac{v^2}{2}} dv = 0 + \int_{-\infty}^{\infty} e^{-\frac{v^2}{2}} dv = \sqrt{2\pi}. \end{aligned}$$

Here above we used the fact that

$$\lim_{v \rightarrow \infty} v e^{-\frac{v^2}{2}} = \lim_{v \rightarrow \infty} \frac{v}{\frac{v^2}{e^{\frac{v^2}{2}}}} \stackrel{l'H\hat{o}pital}{=} \lim_{v \rightarrow \infty} \frac{1}{v e^{\frac{v^2}{2}}} = 0 = \lim_{v \rightarrow -\infty} v e^{-\frac{v^2}{2}}$$

So, $E[X^2] = \frac{1}{\sqrt{2\pi}} (\sigma^2 \sqrt{2\pi} + 2\sigma\mu \cdot 0 + \mu^2 \sqrt{2\pi}) = \sigma^2 + \mu^2$.

And, finally, $Var[X] = E[X^2] - (E[X])^2 = (\sigma^2 + \mu^2) - \mu^2 = \sigma^2$.

Vectors of random variables.

A property:

The covariance matrix Σ corresponding to such a vector is symmetric and positive semi-definite

Chuong Do, Stanford University, 2008

[adapted by Liviu Ciortuz]

Fie variabilele aleatoare X_1, \dots, X_n , cu $X_i : \Omega \rightarrow \mathbb{R}$ pentru $i = 1, \dots, n$. *Matricea de covarianță a vectorului de variabile aleatoare* $X = (X_1, \dots, X_n)$ este o matrice pătratică de dimensiune $n \times n$, ale cărei elemente se definesc astfel: $[Cov(X)]_{ij} \stackrel{def.}{=} Cov(X_i, X_j)$, pentru orice $i, j \in \{1, \dots, n\}$.

Arătați că $\Sigma \stackrel{not.}{=} Cov(X)$ este matrice simetrică și pozitiv semi-definită, cea de-a doua proprietate însemnând că pentru orice vector $z \in \mathbb{R}^n$ are loc inegalitatea $z^\top \Sigma z \geq 0$. (Vectorii $z \in \mathbb{R}^n$ sunt considerați vectori-coloană, iar simbolul \top reprezintă operația de transpunere de matrice.)

$\mathbf{Cov}(X)_{i,j} \stackrel{\text{def.}}{=} \mathbf{Cov}(X_i, X_j)$, for all $i, j \in \{1, \dots, n\}$, and

$\mathbf{Cov}(X_i, X_j) \stackrel{\text{def.}}{=} E[(X_i - E[X_i])(X_j - E[X_j])] = E[(X_j - E[X_j])(X_i - E[X_i])] = \mathbf{Cov}(X_j, X_i)$,
therefore $\mathbf{Cov}(X)$ is a symmetric matrix.

We will show that $z^T \Sigma z \geq 0$ for any $z \in \mathbb{R}^n$ (seen as a column-vector):

$$\begin{aligned}
 z^T \Sigma z &= \sum_{i=1}^n z_i \left(\sum_{j=1}^n \Sigma_{ij} z_j \right) = \sum_{i=1}^n \sum_{j=1}^n (z_i \Sigma_{ij} z_j) = \sum_{i=1}^n \sum_{j=1}^n (z_i \mathbf{Cov}[X_i, X_j] z_j) \\
 &= \sum_{i=1}^n \sum_{j=1}^n (z_i E[(X_i - E[X_i])(X_j - E[X_j])] z_j) = E \left[\sum_{i=1}^n \sum_{j=1}^n z_i (X_i - E[X_i])(X_j - E[X_j]) z_j \right] \\
 &= E \left[\left(\sum_{i=1}^n z_i (X_i - E[X_i]) \right) \left(\sum_{j=1}^n (X_j - E[X_j]) z_j \right) \right] \\
 &= E \left[\left(\sum_{i=1}^n (X_i - E[X_i]) z_i \right) \left(\sum_{j=1}^n (X_j - E[X_j]) z_j \right) \right] = E[(X - E[X])^T \cdot z]^2 \geq 0
 \end{aligned}$$

Multi-variate Gaussian distributions:

A property:

When the covariance matrix of a multi-variate (d -dimensional) Gaussian distribution is diagonal, then the p.d.f. (probability density function) of the respective multi-variate Gaussian is equal to the product of d independent uni-variate Gaussian densities.

Chuong Do, Stanford University, 2008

[adapted by Liviu Ciortuz]

Let's consider $X = [X_1 \dots X_d]^T$, $\mu \in \mathbb{R}^d$ and $\Sigma \in \mathbb{S}_+^d$, where \mathbb{S}_+^d is the set of symmetric positive definite matrices (which implies $|\Sigma| \neq 0$ and $(x - \mu)^T \Sigma^{-1} (x - \mu) > 0$, therefore $-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu) < 0$, for any $x \in \mathbb{S}^d$, $x \neq \mu$).

The probability density function of a multi-variate Gaussian distribution of parameters μ and Σ is:

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right),$$

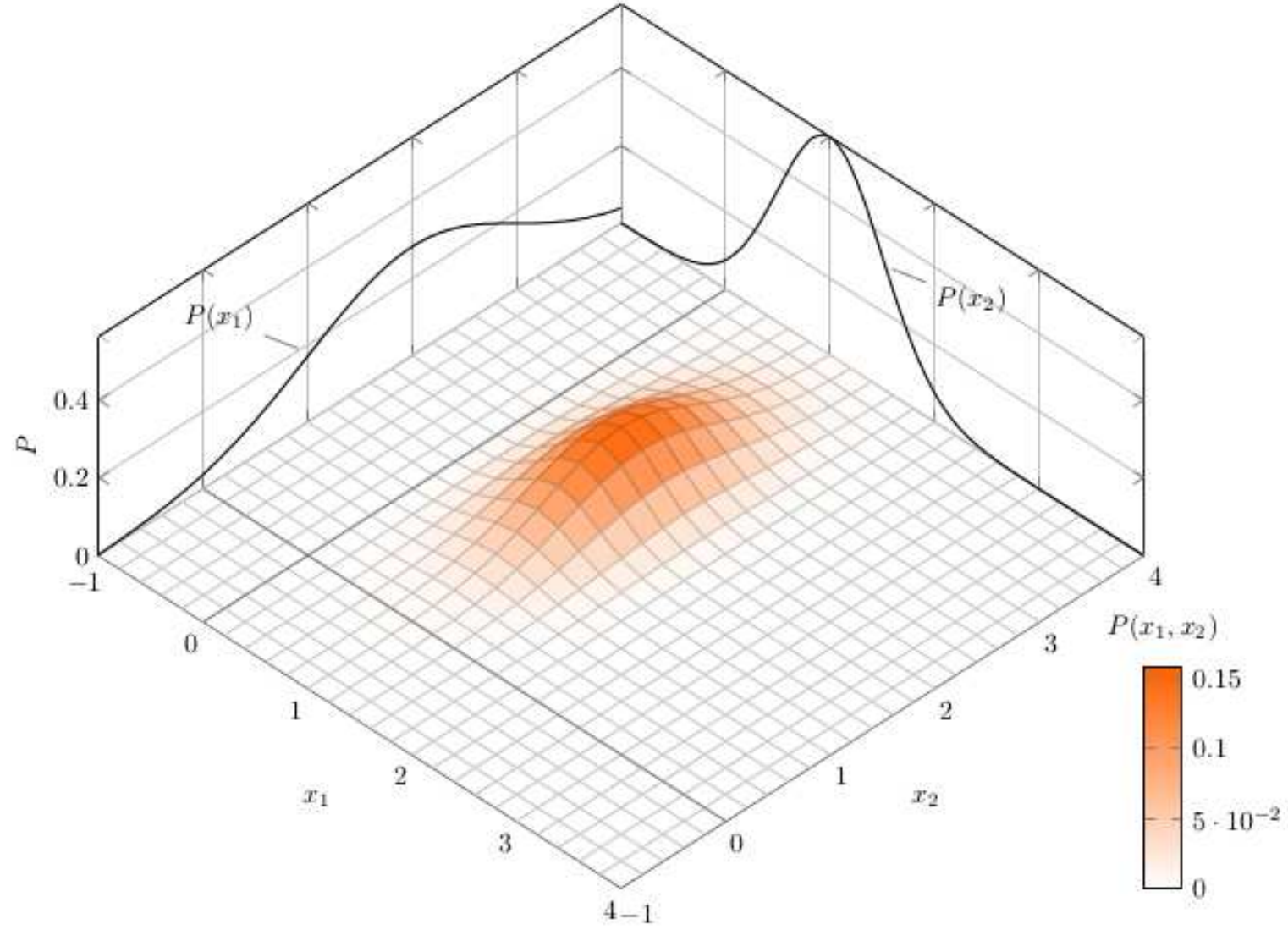
Notation: $X \sim \mathcal{N}(\mu, \Sigma)$.

Show that when the covariance matrix Σ is diagonal, then the p.d.f. (probability density function) of the respective multi-variate Gaussian is equal to the product of d independent uni-variate Gaussian densities.

We will make the **proof** for $d = 2$
(generalization to $d > 2$ will be easy):

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \quad \Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$$

Note: It is easy to show that if $\Sigma \in \mathbb{S}_+^d$ is diagonal, the elements on the principal diagonal Σ are indeed strictly positive. (It is enough to consider $z = (1, 0)$ and respectively $z = (0, 1)$ in formula for *positive-definiteness* of Σ .) This is why we wrote these elements of σ as σ_1^2 and σ_2^2 .



$$\begin{aligned}
p(x; \mu, \Sigma) &= \frac{1}{2\pi \begin{vmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{vmatrix}^{\frac{1}{2}}} \exp \left(-\frac{1}{2} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix}^T \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}^{-1} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix} \right) \\
&= \frac{1}{2\pi \sigma_1 \sigma_2} \exp \left(-\frac{1}{2} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix}^T \begin{bmatrix} \frac{1}{\sigma_1^2} & 0 \\ 0 & \frac{1}{\sigma_2^2} \end{bmatrix} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix} \right) \\
&= \frac{1}{2\pi \sigma_1 \sigma_2} \exp \left(-\frac{1}{2} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix}^T \begin{bmatrix} \frac{1}{\sigma_1^2}(x_1 - \mu_1) \\ \frac{1}{\sigma_2^2}(x_2 - \mu_2) \end{bmatrix} \right) \\
&= \frac{1}{2\pi \sigma_1 \sigma_2} \exp \left(-\frac{1}{2\sigma_1^2}(x_1 - \mu_1)^2 - \frac{1}{2\sigma_2^2}(x_2 - \mu_2)^2 \right) \\
&= p(x_1; \mu_1, \sigma_1^2) p(x_2; \mu_2, \sigma_2^2).
\end{aligned}$$

Bi-variate Gaussian distributions. A property:
The conditional distributions $X_1|X_2$ and $X_2|X_1$ are also
Gaussians.

The calculation of their parameters

Duda, Hart and Stork, *Pattern Classification*, 2001,
Appendix A.5.2

[adapted by Liviu Ciortuz]

Fie X o variabilă aleatoare care urmează o distribuție gaussiană bi-variată de parametri μ (vectorul de medii) și Σ (matricea de covarianță). Așadar, $\mu = (\mu_1, \mu_2) \in \mathbb{R}^2$, iar $\Sigma \in \mathcal{M}_{2 \times 2}(\mathbb{R})$.

Prin definiție, $\Sigma = Cov(X, X)$, unde $X \stackrel{not.}{=} (X_1, X_2)$, așadar $\Sigma_{ij} = Cov(X_i, X_j)$ pentru $i, j \in \{1, 2\}$. De asemenea, $Cov(X_i, X_i) = Var[X_i] \stackrel{not.}{=} \sigma_i^2 \geq 0$ pentru $i \in \{1, 2\}$, în vreme ce pentru $i \neq j$ avem $Cov(X_i, X_j) = Cov(X_j, X_i) \stackrel{not.}{=} \sigma_{ij}$. În sfârșit, dacă introducem „coeficientul de corelare“ $\rho \stackrel{def.}{=} \frac{\sigma_{12}}{\sigma_1 \sigma_2}$, rezultă că putem scrie astfel matricea de covarianță:

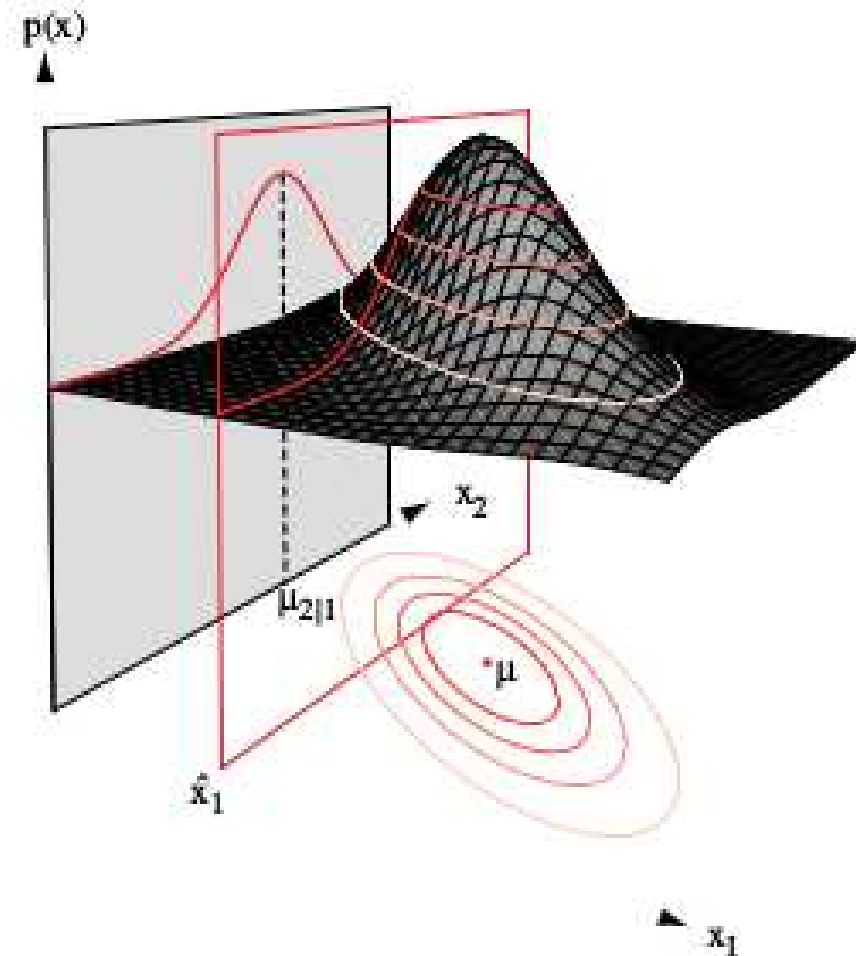
$$\Sigma = \begin{bmatrix} \sigma_1^2 & \rho \sigma_1 \sigma_2 \\ \rho \sigma_1 \sigma_2 & \sigma_2^2 \end{bmatrix}. \quad (3)$$

Demonstrați că ipoteza $X \sim \mathcal{N}(\mu, \Sigma)$, implică faptul că distribuția condițională $X_2|X_1$ este de tip gaussian, și anume

$$X_2|X_1 = x_1 \sim \mathcal{N}(\mu_{2|1}, \sigma_{2|1}^2),$$

cu $\mu_{2|1} = \mu_2 + \rho \frac{\sigma_2}{\sigma_1}(x_1 - \mu_1)$ și $\sigma_{2|1}^2 = \sigma_2^2(1 - \rho^2)$.

Observație: Pentru $X_1|X_2$, rezultatul este similar: $X_1|X_2 = x_2 \sim \mathcal{N}(\mu_{1|2}, \sigma_{1|2}^2)$, cu $\mu_{1|2} = \mu_1 + \rho \frac{\sigma_1}{\sigma_2}(x_2 - \mu_2)$ și $\sigma_{1|2}^2 = \sigma_1^2(1 - \rho^2)$.



Source:

Pattern Classification, Appendix A.5.2,
Duda, Hart and Stork, 2001

Answer

$$p_{X_2|X_1}(x_2|x_1) \stackrel{\text{def.}}{=} \frac{p_{X_1,X_2}(x_1, x_2)}{p_{X_1}(x_1)}, \quad (4)$$

where

$$\begin{aligned} p_{X_1,X_2}(x_1, x_2) &= \frac{1}{(\sqrt{2\pi})^2 \sqrt{|\Sigma|}} \exp \left(-\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right) \text{ si} \\ p_{X_1}(x_1) &= \frac{1}{\sqrt{2\pi}\sigma_1} \exp \left(-\frac{1}{2\sigma_1^2} (x_1 - \mu_1)^2 \right). \end{aligned} \quad (5)$$

From (3) it follows that $|\Sigma| = \sigma_1^2 \sigma_2^2 (1 - \rho^2)$. In order that $\sqrt{|\Sigma|}$ and Σ^{-1} be defined, it follows that $\rho \in (-1, 1)$. Moreover, since $\sigma_1, \sigma_2 > 0$, we will have $\sqrt{|\Sigma|} = \sigma_1 \sigma_2 \sqrt{1 - \rho^2}$.

$$\begin{aligned} \Sigma^{-1} &= \frac{1}{\sigma_1^2 \sigma_2^2 (1 - \rho^2)} \Sigma^* = \frac{1}{\sigma_1^2 \sigma_2^2 (1 - \rho^2)} \begin{bmatrix} \sigma_2^2 & -\rho \sigma_1 \sigma_2 \\ -\rho \sigma_1 \sigma_2 & \sigma_1^2 \end{bmatrix} \\ &= \frac{1}{(1 - \rho^2)} \begin{bmatrix} \frac{1}{\sigma_1^2} & -\frac{\rho}{\sigma_1 \sigma_2} \\ -\frac{\rho}{\sigma_1 \sigma_2} & \frac{1}{\sigma_2^2} \end{bmatrix} \end{aligned}$$

So,

$$\begin{aligned}
 p_{X_1, X_2}(x_1, x_2) &= \\
 &= \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp \left(-\frac{1}{2(1-\rho)^2} (x_1 - \mu_1, x_2 - \mu_2) \begin{bmatrix} \frac{1}{\sigma_1^2} & -\frac{\rho}{\sigma_1\sigma_2} \\ -\frac{\rho}{\sigma_1\sigma_2} & \frac{1}{\sigma_2^2} \end{bmatrix} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix} \right) \\
 &= \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \cdot \\
 &\quad \exp \left(-\frac{1}{2(1-\rho^2)} \left[\left(\frac{x_1 - \mu_1}{\sigma_1} \right)^2 - 2\rho \left(\frac{x_1 - \mu_1}{\sigma_1} \right) \left(\frac{x_2 - \mu_2}{\sigma_2} \right) + \left(\frac{x_2 - \mu_2}{\sigma_2} \right)^2 \right] \right) \quad (6)
 \end{aligned}$$

By substitution (5) and (6) in the definition (4), we will get:

$$\begin{aligned}
 p(x_2|x_1) &= \frac{p_{X_1, X_2}(x_1, x_2)}{p_{X_1}(x_1)} \\
 &= \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \cdot \exp\left(-\frac{1}{2(1-\rho^2)} \left[\left(\frac{x_1-\mu_1}{\sigma_1}\right)^2 - 2\rho\left(\frac{x_1-\mu_1}{\sigma_1}\right)\left(\frac{x_2-\mu_2}{\sigma_2}\right) + \left(\frac{x_2-\mu_2}{\sigma_2}\right)^2 \right]\right) \\
 &\quad \cdot \sqrt{2\pi}\sigma_1 \exp\left(\frac{1}{2}\left(\frac{x_1-\mu_1}{\sigma_1}\right)^2\right) \\
 &= \frac{1}{\sqrt{2\pi}\sigma_2\sqrt{1-\rho^2}} \exp\left[-\frac{1}{2(1-\rho^2)} \left(\frac{x_2-\mu_2}{\sigma_2} - \rho\frac{x_1-\mu_1}{\sigma_1}\right)^2\right] \\
 &= \frac{1}{\sqrt{2\pi}\sigma_2\sqrt{1-\rho^2}} \exp\left[-\frac{1}{2} \left(\frac{x_2 - [\mu_2 + \rho\frac{\sigma_2}{\sigma_1}(x_1 - \mu_1)]}{\sigma_2\sqrt{1-\rho^2}}\right)^2\right]
 \end{aligned}$$

Therefore,

$$X_2|X_1 = x_1 \sim \mathcal{N}(\mu_{2|1}, \sigma_{2|1}^2) \text{ with } \mu_{2|1} \stackrel{not.}{=} \mu_2 + \rho\frac{\sigma_2}{\sigma_1}(x_1 - \mu_1) \text{ and } \sigma_{2|1}^2 \stackrel{not.}{=} \sigma_2^2(1 - \rho^2).$$

Using the Central Limit Theorem (the i.i.d. version)
to compute the *real error* of a classifier
CMU, 2008 fall, Eric Xing, HW3, pr. 3.3

Chris recently adopts a new (binary) classifier to filter email spams. He wants to quantitatively evaluate how good the classifier is.

He has a small dataset of 100 emails on hand which, you can assume, are randomly drawn from all emails.

He tests the classifier on the 100 emails and gets 83 classified correctly, so the error rate on the small dataset is 17%.

However, the number on 100 samples could be either higher or lower than the real error rate just by chance.

With a confidence level of 95%, what is likely to be the range of the real error rate? Please write down all important steps.

(Hint: You need some approximation in this problem.)

Notations:

Let X_i , $i = 1, \dots, n = 100$ be defined as:

$X_i = 1$ if the email i was incorrectly classified, and 0 otherwise;

$$E[X_i] \stackrel{\text{not.}}{=} \mu \stackrel{\text{not.}}{=} e_{\text{real}} ; \quad \text{Var}(X_i) \stackrel{\text{not.}}{=} \sigma^2$$

$$e_{\text{sample}} \stackrel{\text{not.}}{=} \frac{X_1 + \dots + X_n}{n} = 0.17$$

$$Z_n = \frac{X_1 + \dots + X_n - n\mu}{\sqrt{n} \sigma} \quad (\text{the standardized form of } X_1 + \dots + X_n)$$

Key insight:

Calculating the real error of the classifier (more exactly, a symmetric interval around the real error $p \stackrel{\text{not.}}{=} \mu$) with a “confidence” of 95% amounts to finding $a > 0$ such that $P(|Z_n| \leq a) \geq 0.95$.

Calculus:

$$\begin{aligned}
 |Z_n| \leq a &\Leftrightarrow \left| \frac{X_1 + \dots + X_n - n\mu}{\sqrt{n} \sigma} \right| \leq a \Leftrightarrow \left| \frac{X_1 + \dots + X_n - n\mu}{n\sigma} \right| \leq \frac{a}{\sqrt{n}} \\
 &\Leftrightarrow \left| \frac{X_1 + \dots + X_n - n\mu}{n} \right| \leq \frac{a\sigma}{\sqrt{n}} \Leftrightarrow \left| \frac{X_1 + \dots + X_n}{n} - \mu \right| \leq \frac{a\sigma}{\sqrt{n}} \\
 &\Leftrightarrow |e_{\text{sample}} - e_{\text{real}}| \leq \frac{a\sigma}{\sqrt{n}} \Leftrightarrow |e_{\text{real}} - e_{\text{sample}}| \leq \frac{a\sigma}{\sqrt{n}} \\
 &\Leftrightarrow -\frac{a\sigma}{\sqrt{n}} \leq e_{\text{real}} - e_{\text{sample}} \leq \frac{a\sigma}{\sqrt{n}} \\
 &\Leftrightarrow e_{\text{sample}} - \frac{a\sigma}{\sqrt{n}} \leq e_{\text{real}} \leq e_{\text{sample}} + \frac{a\sigma}{\sqrt{n}} \\
 &\Leftrightarrow e_{\text{real}} \in \left[e_{\text{sample}} - \frac{a\sigma}{\sqrt{n}}, e_{\text{sample}} + \frac{a\sigma}{\sqrt{n}} \right]
 \end{aligned}$$

Important facts:

40.

The Central Limit Theorem: $Z_n \rightarrow \mathcal{N}(0; 1)$

Therefore, $P(|Z_n| \leq a) \approx P(|X| \leq a) = \Phi(a) - \Phi(-a)$, **where** $X \sim \mathcal{N}(0; 1)$
and Φ is the cumulative function distribution of $\mathcal{N}(0; 1)$.

Calculus:

$$\Phi(-a) + \Phi(a) = 1 \Rightarrow P(|Z_n| \leq a) = \Phi(a) - \Phi(-a) = 2\Phi(a) - 1$$

$$P(|Z_n| \leq a) = 0.95 \Leftrightarrow 2\Phi(a) - 1 = 0.95 \Leftrightarrow \Phi(a) = 0.975 \Leftrightarrow a \cong 1.97 \text{ (see } \Phi \text{ table)}$$

$\sigma^2 \stackrel{\text{not.}}{=} \text{Var}_{real} = e_{real}(1 - e_{real})$ **because X_i are Bernoulli variables.**

Futhermore, we can approximate e_{real} with e_{sample} , because

$$E[e_{sample}] = e_{real} \text{ and } \text{Var}_{sample} = \frac{1}{n} \text{Var}_{real} \rightarrow 0 \text{ for } n \rightarrow +\infty,$$

cf. CMU, 2011 fall, T. Mitchell, A. Singh, HW2, pr. 1.ab.

Finally:

$$\Rightarrow \frac{a\sigma}{\sqrt{n}} \approx 1.97 \cdot \frac{\sqrt{0.17(1 - 0.17)}}{\sqrt{100}} \cong 0.07$$

$$|e_{real} - e_{sample}| \leq 0.07 \Leftrightarrow |e_{real} - 0.17| \leq 0.07 \Leftrightarrow -0.07 \leq e_{real} - 0.17 \leq 0.07$$

$$\Leftrightarrow e_{real} \in [0.10, 0.24]$$

Exemplifying
a mixture of categorical distributions;
how to compute its expectation and variance

CMU, 2010 fall, Aarti Singh, HW1, pr. 2.2.1-2

Suppose that I have two six-sided dice, one is fair and the other one is loaded – having:

$$P(x) = \begin{cases} \frac{1}{2} & x = 6 \\ \frac{1}{10} & x \in \{1, 2, 3, 4, 5\} \end{cases}$$

I will toss a coin to decide which die to roll. If the coin flip is heads I will roll the fair die, otherwise the loaded one. The probability that the coin flip is heads is $p \in (0, 1)$.

- a. What is the expectation of the *die roll* (in terms of p).
- b. What is the variation of the *die roll* (in terms of p).

Solution:**a.**

$$\begin{aligned} E[X] &= \sum_{i=1}^6 i \cdot [P(i|fair) \cdot p + P(i|loaded) \cdot (1 - p)] \\ &= \left[\sum_{i=1}^6 i \cdot P(i|fair) \right] p + \left[\sum_{i=1}^6 i \cdot P(i|loaded) \right] (1 - p) \\ &= \frac{7}{2}p + \frac{9}{2}(1 - p) = \frac{9}{2} - p \end{aligned}$$

b. Recall that we may write $Var(X) = E[X^2] - (E[X])^2$, therefore:

$$\begin{aligned}
 E[X^2] &= \sum_{i=1}^6 i^2 \cdot [P(i|fair) \cdot p + P(i|loaded) \cdot (1-p)] \\
 &= \left[\sum_{i=1}^6 i^2 \cdot P(i|fair) \right] p + \left[\sum_{i=1}^6 i^2 \cdot P(i|loaded) \right] (1-p) \\
 &= \frac{91}{6}p + \left(\frac{36}{2} + \frac{55}{10} \right) (1-p) \\
 &= \frac{47}{2} - \frac{25}{3}p
 \end{aligned}$$

Combining this with the result of the previous question yields:

$$\begin{aligned}
 Var(X) &= E[X^2] - (E[X])^2 = \frac{141}{6} - \frac{50}{6}p - \left(\frac{9}{2} - p \right)^2 \\
 &= \frac{141}{6} - \frac{50}{6}p - \left(\frac{81}{4} - 9p + p^2 \right) \\
 &= \left(\frac{141}{6} - \frac{81}{4} \right) - \left(\frac{50}{6} - 9 \right)p - p^2 \\
 &= \frac{13}{4} + \frac{2}{3}p - p^2
 \end{aligned}$$

Elements of Information Theory:

Some examples and then some useful proofs

**Computing entropies and specific conditional entropies
for discrete random variables**

CMU, 2012 spring, R. Rosenfeld, HW2, pr. 2

On the roll of two six-sided fair dice,

- a. Calculate the distribution of the sum (S) of the total.
- b. The amount of *information* (or *surprise*) when seeing the outcome x for a random variable X is defined as $\log_2 \frac{1}{P(X=x)} = -\log_2 P(X=x)$. How surprised are you (in bits) to observe $S=2$, $S=11$, $S=5$, $S=7$?
- c. Calculate the *entropy* of S [as the *expected value* of the random variable $-\log_2 P(X=x)$].
- d. Let's say you throw the die one by one, and the first die shows 4. What is the entropy of S after this observation? Was any information gained / lost in the process? If so, calculate how much information (in bits) was lost or gained.

a.

S	2	3	4	5	6	7	8	9	10	11	12
$P(S)$	1/36	2/36	3/36	4/36	5/36	6/36	5/36	4/36	3/36	2/36	1/36

b.

$$\begin{aligned}
 \text{Information}(S = 2) &= -\log_2(1/36) = \log_2 36 = 2 \log_2 6 = 2(1 + \log_2 3) \\
 &= 5.169925001 \text{ bits}
 \end{aligned}$$

$$\text{Information}(S = 11) = -\log_2 2/36 = \log_2 18 = 1 + 2 \log_2 3 = 4.169925001 \text{ bits}$$

$$\text{Information}(S = 5) = -\log_2 4/36 = \log_2 9 = 2 \log_2 3 = 3.169925001 \text{ bits}$$

$$\text{Information}(S = 7) = -\log_2 6/36 = \log_2 6 = 1 + \log_2 3 = 2.584962501 \text{ bits}$$

c.

$$\begin{aligned}
H(S) &= - \sum_{i=1}^n p_i \log p_i \\
&= - \left(2 \cdot \frac{1}{36} \log \frac{1}{36} + 2 \cdot \frac{2}{36} \log \frac{2}{36} + 2 \cdot \frac{3}{36} \log \frac{3}{36} + 2 \cdot \frac{4}{36} \log \frac{4}{36} + \right. \\
&\quad \left. 2 \cdot \frac{5}{36} \log \frac{5}{36} + \frac{6}{36} \log \frac{6}{36} \right) \\
&= \frac{1}{36} (2 \log_2 36 + 4 \log_2 18 + 6 \log_2 12 + 8 \log_2 9 + 10 \log_2 \frac{36}{5} + 6 \log_2 6) \\
&= \frac{1}{36} (2 \log_2 6^2 + 4 \log_2 6 \cdot 3 + 6 \log_2 6 \cdot 2 + 8 \log_2 3^2 + 10 \log_2 \frac{6^2}{5} + 6 \log_2 6) \\
&= \frac{1}{36} (40 \log_2 6 + 20 \log_2 3 + 6 - 10 \log_2 5) \\
&= \frac{1}{36} (60 \log_2 3 + 46 - 10 \log_2 5) = 3.274401919 \text{ bits.}
\end{aligned}$$

d.

S	2	3	4	5	6	7	8	9	10	11	12
$P(S ...)$	0	0	0	1/6	1/6	1/6	1/6	1/6	1/6	0	0

$$H(S|First-die-shows-4) = -6 \cdot \frac{1}{6} \log_2 \frac{1}{6} = \log_2 6 = 2.58 \text{ bits},$$

$$IG(S; First-die-shows-4) = H(S) - H(S|First-die-shows-4) = 3.27 - 2.58 = 0.69 \text{ bits}.$$

Computing entropies and average conditional entropies
for discrete random variables

CMU, 2012 spring, Roni Rosenfeld, HW2, pr. 3

A doctor needs to diagnose a person having cold (C). The primary factor he considers in his diagnosis is the outside temperature (T). The random variable C takes two values, *yes* / *no*, and the random variable T takes 3 values, *sunny*, *rainy*, *snowy*. The joint distribution of the two variables is given in following table.

	$T = \textit{sunny}$	$T = \textit{rainy}$	$T = \textit{snowy}$
$C = \textit{no}$	0.30	0.20	0.10
$C = \textit{yes}$	0.05	0.15	0.20

a. Calculate the *marginal probabilities* $P(C)$, $P(T)$.

Hint: Use $P(X = x) = \sum_Y P(X = x; Y = y)$. For example,

$$P(C = \textit{no}) = P(C = \textit{no}, T = \textit{sunny}) + P(C = \textit{no}, T = \textit{rainy}) + P(C = \textit{no}, T = \textit{snowy}).$$

b. Calculate the *entropies* $H(C)$, $H(T)$.

c. Calculate the *average conditional entropies* $H(C|T)$, $H(T|C)$.

a. $P_C = (0.6, 0.4)$ si $P_T = (0.35, 0.35, 0.30)$.

b.

$$H(C) = 0.6 \log \frac{5}{3} + 0.4 \log \frac{5}{2} = \log 5 - 0.6 \log 3 - 0.4 = 0.971 \text{ bits}$$

$$\begin{aligned} H(T) &= 2 \cdot 0.35 \log \frac{20}{7} + 0.3 \log \frac{10}{3} \\ &= 0.7(2 + \log 5 - \log 7) + 0.3(1 + \log 5 - \log 3) \\ &= 1.7 + \log 5 - 0.7 \log 7 - 0.3 \log 3 = 1.581 \text{ bits.} \end{aligned}$$

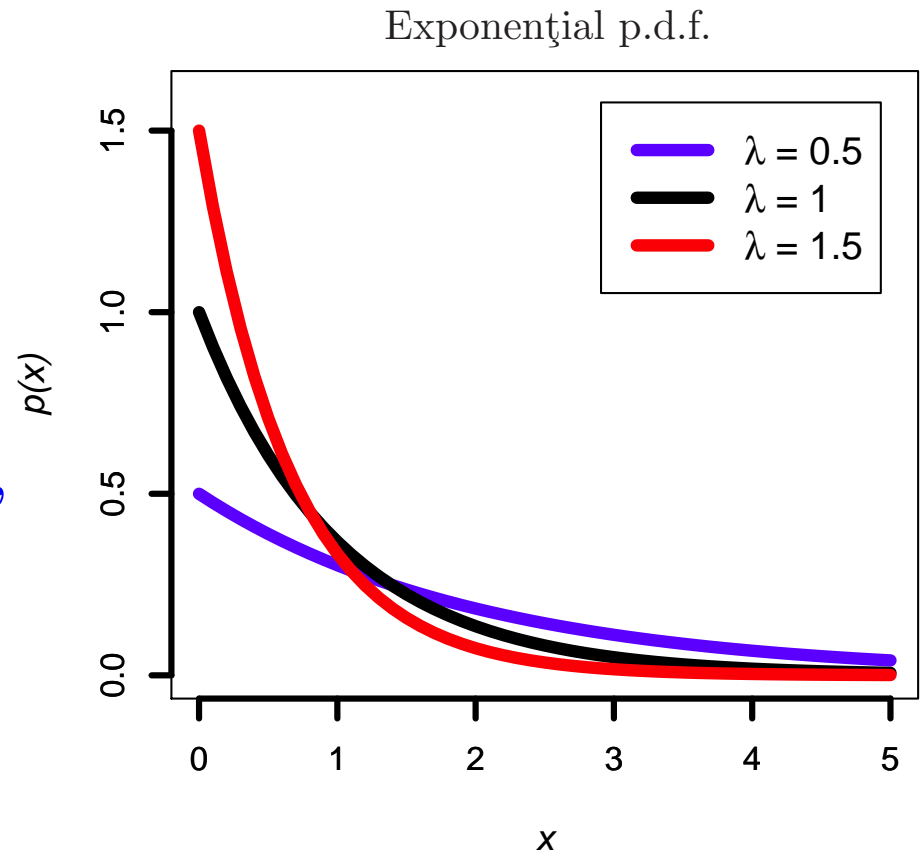
c.

$$\begin{aligned}
H(C|T) &\stackrel{def.}{=} \sum_{t \in Val(T)} P(T = t) \cdot H(C|T = t) \\
&= P(T = sunny) \cdot H(C|T = sunny) + P(T = rainy) \cdot H(C|T = rainy) + \\
&\quad P(T = snowy) \cdot H(C|T = snowy) \\
&= 0.35 \cdot H\left(\frac{0.30}{0.30 + 0.05}, \frac{0.05}{0.30 + 0.05}\right) + 0.35 \cdot H\left(\frac{0.20}{0.20 + 0.15}, \frac{0.15}{0.20 + 0.15}\right) + \\
&\quad 0.30 \cdot H\left(\frac{0.10}{0.10 + 0.20}, \frac{0.20}{0.20 + 0.10}\right) \\
&= \frac{7}{20} \cdot H\left(\frac{6}{7}, \frac{1}{7}\right) + \frac{7}{20} \cdot H\left(\frac{4}{7}, \frac{3}{7}\right) + \frac{3}{10} \cdot H\left(\frac{1}{3}, \frac{2}{3}\right) \\
&= \frac{7}{20} \cdot \left(\frac{6}{7} \log \frac{7}{6} + \frac{1}{7} \log 7\right) + \frac{7}{20} \cdot \left(\frac{4}{7} \log \frac{7}{4} + \frac{3}{7} \log \frac{7}{3}\right) + \frac{3}{10} \cdot \left(\frac{1}{3} \log 3 + \frac{2}{3} \log \frac{3}{2}\right) \\
&= \frac{7}{20} \cdot \left(\log 7 - \frac{6}{7} - \frac{6}{7} \log 3\right) + \frac{7}{20} \cdot \left(\log 7 - \frac{8}{7} - \frac{3}{7} \log 3\right) + \frac{3}{10} \cdot \left(\log 3 - \frac{2}{3}\right) \\
&= \frac{7}{10} \log 7 - \left(\frac{3}{10} + \frac{4}{10} + \frac{2}{10}\right) - \left(\frac{6}{20} + \frac{3}{20} - \frac{3}{10}\right) \cdot \log 3 = \frac{7}{10} \log 7 - \frac{3}{20} \log 3 - \frac{9}{10} = 0.82715 \text{ bits.}
\end{aligned}$$

$$\begin{aligned}
H(T|C) &\stackrel{\text{def.}}{=} \sum_{c \in \text{Val}(C)} P(C = c) \cdot H(T|C = c) \\
&= P(C = \text{no}) \cdot H(T|C = \text{no}) + P(C = \text{yes}) \cdot H(T|C = \text{yes}) \\
&= 0.60 \cdot H\left(\frac{0.30}{0.30 + 0.20 + 0.10}, \frac{0.20}{0.30 + 0.20 + 0.10}, \frac{0.10}{0.30 + 0.20 + 0.10}\right) + \\
&\quad 0.40 \cdot H\left(\frac{0.05}{0.05 + 0.15 + 0.20}, \frac{0.15}{0.05 + 0.15 + 0.20}, \frac{0.20}{0.05 + 0.15 + 0.20}\right) \\
&= \frac{3}{5} \cdot H\left(\frac{1}{2}, \frac{1}{3}, \frac{1}{6}\right) + \frac{2}{5} \cdot H\left(\frac{1}{8}, \frac{3}{8}, \frac{1}{2}\right) \\
&= \frac{3}{5} \left(\frac{1}{2} + \frac{1}{3} \log 3 + \frac{1}{6} (1 + \log 3)\right) + \frac{2}{5} \left(\frac{1}{8} \cdot 3 + \frac{3}{8} (3 - \log 3) + \frac{1}{2}\right) \\
&= \frac{3}{5} \left(\frac{2}{3} + \frac{1}{2} \log 3\right) + \frac{2}{5} \left(2 - \frac{3}{8} \log 3\right) \\
&= \frac{6}{5} + \frac{3}{20} \log 3 = 1.43774 \text{ bits.}
\end{aligned}$$

Computing the entropy of the exponential distribution

CMU, 2011 spring, R. Rosenfeld,
HW2, pr. 2.c



Pentru o distribuție de probabilitate continuă P , entropia se definește astfel:

$$H(P) = \int_{-\infty}^{+\infty} P(x) \log_2 \frac{1}{P(x)} dx$$

Calculați entropia *distribuției* continue *exponențiale* de parametru $\lambda > 0$. Definiția acestei distribuții este următoarea:

$$P(x) = \begin{cases} \lambda e^{-\lambda x}, & \text{dacă } x \geq 0; \\ 0, & \text{dacă } x < 0. \end{cases}$$

Indicație: Dacă $P(x) = 0$, veți presupune că $-P(x) \log_2 P(x) = 0$.

Answer

$$\begin{aligned}
 H(P) &= \int_{-\infty}^0 P(x) \log_2 \frac{1}{P(x)} dx + \int_0^{\infty} P(x) \log_2 \frac{1}{P(x)} dx \\
 &\stackrel{\text{def. } P}{=} \underbrace{\int_{-\infty}^0 0 \log_2 0 dx}_0 + \int_0^{\infty} \lambda e^{-\lambda x} \log_2 \frac{1}{\lambda e^{-\lambda x}} dx = \int_0^{\infty} \lambda e^{-\lambda x} \log_2 \frac{1}{\lambda e^{-\lambda x}} dx \\
 \Rightarrow H(P) &= \frac{1}{\ln 2} \int_0^{\infty} \lambda e^{-\lambda x} \ln \frac{1}{\lambda e^{-\lambda x}} dx = \frac{1}{\ln 2} \int_0^{\infty} \lambda e^{-\lambda x} \left(\ln \frac{1}{\lambda} + \ln \frac{1}{e^{-\lambda x}} \right) dx \\
 &= \frac{1}{\ln 2} \int_0^{\infty} \lambda e^{-\lambda x} (-\ln \lambda + \ln e^{\lambda x}) dx \\
 &= \frac{1}{\ln 2} \int_0^{\infty} \lambda e^{-\lambda x} (-\ln \lambda + \lambda x) dx \\
 &= \frac{1}{\ln 2} \int_0^{\infty} \lambda e^{-\lambda x} (-\ln \lambda) dx + \frac{1}{\ln 2} \int_0^{\infty} \lambda e^{-\lambda x} \lambda x dx \\
 &= \frac{-\ln \lambda}{\ln 2} \int_0^{\infty} \lambda e^{-\lambda x} dx + \frac{\lambda}{\ln 2} \int_0^{\infty} \lambda e^{-\lambda x} x dx \\
 &= \frac{\ln \lambda}{\ln 2} \int_0^{\infty} (e^{-\lambda x})' dx - \frac{\lambda}{\ln 2} \int_0^{\infty} (e^{-\lambda x})' x dx
 \end{aligned}$$

Prima integrală se rezolvă foarte ușor:

$$\int_0^{\infty} (e^{-\lambda x})' dx = e^{-\lambda x} \Big|_0^{\infty} = e^{-\infty} - e^0 = 0 - 1 = -1$$

Pentru a rezolva cea de-a doua integrală se poate folosi *formula de integrare prin părți*:

$$\int_0^{\infty} (e^{-\lambda x})' x dx = e^{-\lambda x} x \Big|_0^{\infty} - \int_0^{\infty} e^{-\lambda x} x' dx = e^{-\lambda x} x \Big|_0^{\infty} - \int_0^{\infty} e^{-\lambda x} dx$$

Integrala definită $e^{-\lambda x} x \Big|_0^{\infty}$ nu se poate calcula direct (din cauza conflictului $0 \cdot \infty$ care se produce atunci când lui x i se atribuie valoarea-limită ∞), ci se calculează folosind *regula lui l'Hôpital*:

$$\lim_{x \rightarrow \infty} x e^{-\lambda x} = \lim_{x \rightarrow \infty} \frac{x}{e^{\lambda x}} = \lim_{x \rightarrow \infty} \frac{x'}{(e^{\lambda x})'} = \lim_{x \rightarrow \infty} \frac{1}{\lambda e^{\lambda x}} = \frac{1}{\lambda} \lim_{x \rightarrow \infty} e^{-\lambda x} = e^{-\infty} = 0,$$

deci

$$e^{-\lambda x} x \Big|_0^{\infty} = 0 - 0 = 0.$$

Integrala $\int_0^\infty e^{-\lambda x} dx$ se calculează ușor:

$$\int_0^\infty e^{-\lambda x} dx = -\frac{1}{\lambda} \int_0^\infty (e^{-\lambda x})' dx = -\frac{1}{\lambda} e^{-\lambda x} \Big|_0^\infty = -\frac{1}{\lambda} (0 - 1) = \frac{1}{\lambda}$$

Prin urmare,

$$\int_0^\infty (e^{-\lambda x})' x dx = 0 - \frac{1}{\lambda} = -\frac{1}{\lambda},$$

ceea ce conduce la rezultatul final:

$$H(P) = \frac{\ln \lambda}{\ln 2} (-1) - \frac{\lambda}{\ln 2} \left(-\frac{1}{\lambda} \right) = -\frac{\ln \lambda}{\ln 2} + \frac{1}{\ln 2} = \frac{1 - \ln \lambda}{\ln 2}.$$

**Derivation of entropy definition,
starting from a set of desirable properties**
CMU, 2005 fall, T. Mitchell, A. Moore, HW1, pr. 2.2

Remark: The definition we gave for entropy $-\sum_{i=1}^n p_i \log p_i$ is not very intuitive.

Theorem:

If $\psi_n(p_1, \dots, p_n)$ satisfies the following axioms

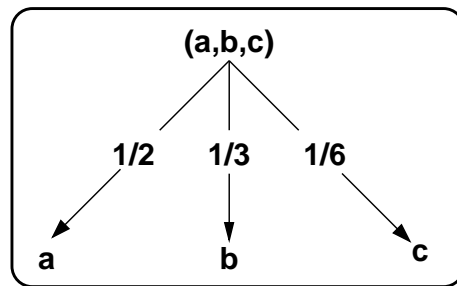
- A0. [LC:] $\psi_n(p_1, \dots, p_n) \geq 0$ for any $n \in \mathbb{N}^*$ and p_1, \dots, p_n , since we view ψ_n is a measure of *disorder*; also, $\psi_1(1) = 0$ because in this case there is no disorder;
- A1. ψ_n should be continuous in p_i and symmetric in its arguments;
- A2. if $p_i = 1/n$ then ψ_n should be a monotonically increasing function of n ;
(If all events are equally likely, then having more events means being more uncertain.)
- A3. if a choice among N events is broken down into successive choices, then the entropy should be the weighted sum of the entropy at each stage;

then $\psi_n(p_1, \dots, p_n) = -K \sum_i p_i \log p_i$ where K is a positive constant.

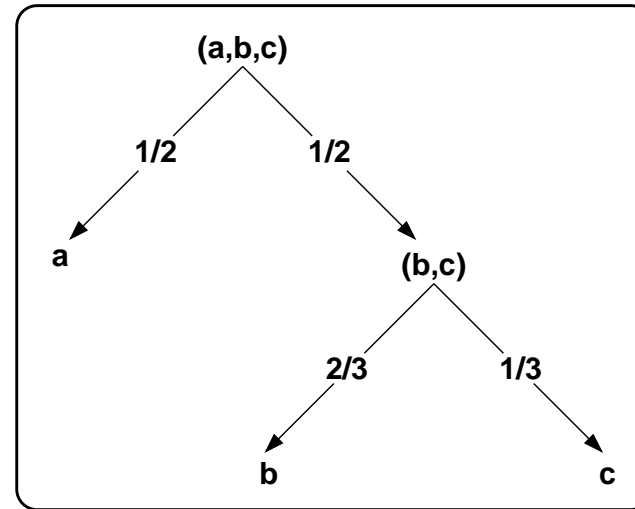
(As we'll see, K depends however on $\psi_s \left(\frac{1}{s}, \dots, \frac{1}{s} \right)$ for a certain $s \in \mathbb{N}^*$).

Remark: We will prove the theorem firstly for uniform distributions ($p_i = 1/n$) and secondly for the case $p_i \in \mathbb{Q}$ (only!).

Example for the axiom A3:



Encoding 1



Encoding 2

$$H\left(\frac{1}{2}, \frac{1}{3}, \frac{1}{6}\right) = \frac{1}{2} \log 2 + \frac{1}{3} \log 3 + \frac{1}{6} \log 6 = \left(\frac{1}{2} + \frac{1}{6}\right) \log 2 + \left(\frac{1}{3} + \frac{1}{6}\right) \log 3 = \frac{2}{3} + \frac{1}{2} \log 3$$

$$H\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{1}{2} H\left(\frac{2}{3}, \frac{1}{3}\right) = 1 + \frac{1}{2} \left(\frac{2}{3} \log \frac{3}{2} + \frac{1}{3} \log 3\right) = 1 + \frac{1}{2} \left(\log 3 - \frac{2}{3}\right) = \frac{2}{3} + \frac{1}{2} \log 3$$

The next 3 slides:

Case 1: $p_i = 1/n$ for $i = 1, \dots, n$; proof steps:

a. $A(n) \stackrel{not.}{=} \psi(1/n, 1/n, \dots, 1/n)$ implies

$$A(s^m) = m A(s) \text{ for any } s, m \in \mathbb{N}^*. \quad (1)$$

b. If $s, m \in \mathbb{N}^*$ (fixed), $s \neq 1$, and $t, n \in \mathbb{N}^*$ such that $s^m \leq t^n \leq s^{m+1}$, then

$$\left| \frac{m}{n} - \frac{\log t}{\log s} \right| \leq \frac{1}{n}. \quad (2)$$

c. For $s^m \leq t^n \leq s^{m+1}$ as above, due to A2 it follows (immediately)

$$\psi_{s^m} \left(\frac{1}{s^m}, \dots, \frac{1}{s^m} \right) \leq \psi_{t^n} \left(\frac{1}{t^n}, \dots, \frac{1}{t^n} \right) \leq \psi_{s^{m+1}} \left(\frac{1}{s^{m+1}}, \dots, \frac{1}{s^{m+1}} \right)$$

i.e. $A(s^m) \leq A(t^n) \leq A(s^{m+1})$

Show that

$$\left| \frac{m}{n} - \frac{A(t)}{A(s)} \right| \leq \frac{1}{n} \text{ for } s \neq 1. \quad (3)$$

d. Combining (2) + (3) immediately gives

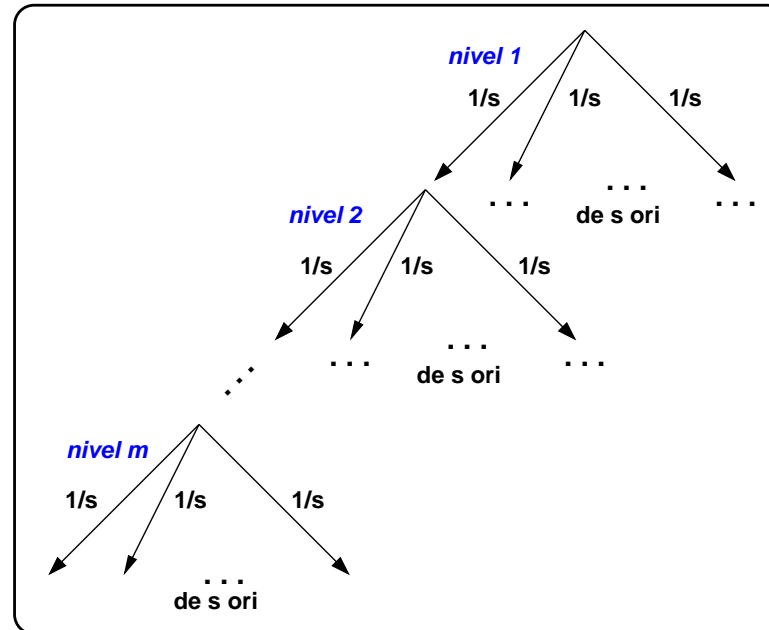
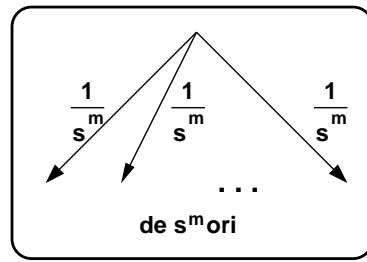
$$\left| \frac{A(t)}{A(s)} - \frac{\log t}{\log s} \right| \leq \frac{2}{n} \text{ pentru } s \neq 1 \quad (4)$$

Show that this inequation implies

$$A(t) = K \log t \text{ with } K > 0 \text{ (due to A2)}. \quad (5)$$

Proof

a.



Applying the axion A3 on the right encoding from above gives:

$$\begin{aligned}
 A(s^m) &= A(s) + s \cdot \frac{1}{s} A(s) + s^2 \cdot \frac{1}{s^2} A(s) + \dots + s^{m-1} \cdot \frac{1}{s^{m-1}} A(s) \\
 &= \underbrace{A(s) + A(s) + A(s) + \dots + A(s)}_{m \text{ times}} = mA(s)
 \end{aligned}$$

Proof (cont'd)

b.

$$s^m \leq t^n \leq s^{m+1} \Rightarrow m \log s \leq n \log t \leq (m+1) \log s \Rightarrow$$

$$\frac{m}{n} \leq \frac{\log t}{\log s} \leq \frac{m}{n} + \frac{1}{n} \Rightarrow 0 \leq \frac{\log t}{\log s} - \frac{m}{n} \leq \frac{1}{n} \Rightarrow \left| \frac{\log t}{\log s} - \frac{m}{n} \right| \leq \frac{1}{n}$$

c.

$$A(s^m) \leq A(t^n) \leq A(s^{m+1}) \stackrel{(1)}{\Rightarrow} m A(s) \leq n A(t) \leq (m+1) A(s) \stackrel{s \neq 1}{\Rightarrow}$$

$$\frac{m}{n} \leq \frac{A(t)}{A(s)} \leq \frac{m}{n} + \frac{1}{n} \Rightarrow 0 \leq \frac{A(t)}{A(s)} - \frac{m}{n} \leq \frac{1}{n} \Rightarrow \left| \frac{A(t)}{A(s)} - \frac{m}{n} \right| \leq \frac{1}{n}$$

d. Consider again $s^m \leq t^n \leq s^{m+1}$ with s, t fixed. If $m \rightarrow \infty$ then $n \rightarrow \infty$ and from $\left| \frac{A(t)}{A(s)} - \frac{\log t}{\log s} \right| \leq \frac{2}{n}$ it follows that $\left| \frac{A(t)}{A(s)} - \frac{\log t}{\log s} \right| \rightarrow 0$.

Therefore $\left| \frac{A(t)}{A(s)} - \frac{\log t}{\log s} \right| = 0$ and so $\frac{A(t)}{A(s)} = \frac{\log t}{\log s}$.

Finally, $A(t) = \frac{A(s)}{\log s} \log t = K \log t$, where $K = \frac{A(s)}{\log s} > 0$ (if $s \neq 1$).

Case 2: $p_i \in \mathbb{Q}$ for $i = 1, \dots, n$

Let's consider a set of $N \geq 2$ equiprobable random events, and $\mathcal{P} = (S_1, S_2, \dots, S_k)$ a partition of this set. Let's denote $p_i = |S_i|/N$.

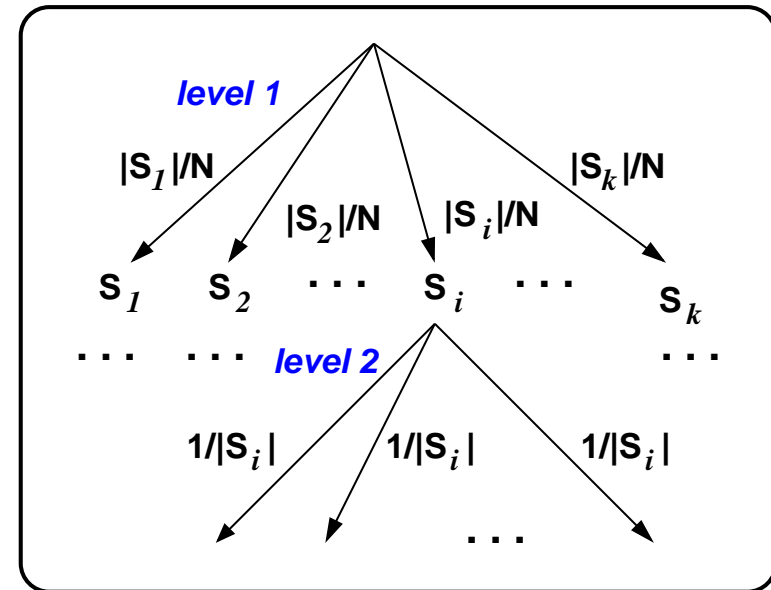
A “natural” two-step encoding (as shown in the nearby figure) leads to $A(N) = \psi_k(p_1, \dots, p_k) + \sum_i p_i A(|S_i|)$, based on the axiom A3.

Finally, using the result $A(t) = K \log t$, gives:

$$K \log N = \psi_k(p_1, \dots, p_k) + K \sum_i p_i \log |S_i|$$

$$\Rightarrow \psi_k(p_1, \dots, p_k) = K \left[\log N - \sum_i p_i \log |S_i| \right]$$

$$= K \left[\log N \sum_i p_i - \sum_i p_i \log |S_i| \right] = -K \sum_i p_i \log \frac{|S_i|}{N} = -K \sum_i p_i \log p_i$$



**Entropie, entropie corelată,
entropie condițională, câștig de informație:
definiții și proprietăți imediate**

CMU, 2005 fall, T. Mitchell, A. Moore, HW1, pr. 2

Definiții

- **Entropia variabilei X :**

$$H(X) \stackrel{\text{def.}}{=} - \sum_i P(X = x_i) \log P(X = x_i) \stackrel{\text{not.}}{=} E_X[-\log P(X)].$$

- **Entropia condițională specifică a variabilei Y în raport cu valoarea x_k a variabilei X :**

$$H(Y | X = x_k) \stackrel{\text{def.}}{=} - \sum_j P(Y = y_j | X = x_k) \log P(Y = y_j | X = x_k) \\ \stackrel{\text{not.}}{=} E_{Y|X=x_k}[-\log P(Y | X = x_k)].$$

- **Entropia condițională medie a variabilei Y în raport cu variabila X :**

$$H(Y | X) \stackrel{\text{def.}}{=} \sum_k P(X = x_k) H(Y | X = x_k) \stackrel{\text{not.}}{=} E_X[H(Y | X)].$$

- **Entropia corelată a variabilelor X și Y :**

$$H(X, Y) \stackrel{\text{def.}}{=} - \sum_i \sum_j P(X = x_i, Y = y_j) \log P(X = x_i, Y = y_j) \\ \stackrel{\text{not.}}{=} E_{X,Y}[-\log P(X, Y)].$$

- **Informația mutuală a variabilelor X și Y , numită de asemenea *câștigul de informație* al variabilei X în raport cu variabila Y (sau invers):**

$$MI(X, Y) \stackrel{\text{not.}}{=} IG(X, Y) \stackrel{\text{def.}}{=} H(X) - H(X | Y) = H(Y) - H(Y | X)$$

(Observație: ultima egalitate de mai sus are loc datorită rezultatului de la punctul c de mai jos.)

a.

$$H(X) \geq 0.$$

$$H(X) = - \sum_i P(X = x_i) \log P(X = x_i) = \sum_i \underbrace{P(X = x_i)}_{\geq 0} \underbrace{\log \frac{1}{P(X = x_i)}}_{\geq 0} \geq 0$$

Mai mult, $H(X) = 0$ dacă și numai dacă variabila X este constantă:

„ \Rightarrow “ Presupunem că $H(X) = 0$, adică $\sum_i P(X = x_i) \log \frac{1}{P(X = x_i)} = 0$. Datorită faptului că fiecare termen din această sumă este mai mare sau egal cu 0, rezultă că $H(X) = 0$ doar dacă pentru $\forall i$, $P(X = x_i) = 0$ sau $\log \frac{1}{P(X = x_i)} = 0$, adică dacă pentru $\forall i$, $P(X = x_i) = 0$ sau $P(X = x_i) = 1$. Cum însă $\sum_i P(X = x_i) = 1$ rezultă că există o singură valoare x_1 pentru X astfel încât $P(X = x_1) = 1$, iar $P(X = x) = 0$ pentru orice $x \neq x_1$. Altfel spus, variabila aleatoare discretă X este constantă.

„ \Leftarrow “ Presupunem că variabila X este constantă, ceea ce înseamnă că X ia o singură valoare x_1 , cu probabilitatea $P(X = x_1) = 1$. Prin urmare, $H(X) = -1 \cdot \log 1 = 0$.

b.

$$H(Y | X) = - \sum_i \sum_j P(X = x_i, Y = y_j) \log P(Y = y_j | X = x_i)$$

$$\begin{aligned}
 H(Y | X) &= \sum_i P(X = x_i) H(Y | X = x_i) \\
 &= \sum_i P(X = x_i) \left[- \sum_j P(Y = y_j | X = x_i) \log P(Y = y_j | X = x_i) \right] \\
 &= - \sum_i \sum_j \underbrace{P(X = x_i) P(Y = y_j | X = x_i)}_{=P(X=x_i, Y=y_j)} \log P(Y = y_j | X = x_i) \\
 &= - \sum_i \sum_j P(X = x_i, Y = y_j) \log P(Y = y_j | X = x_i)
 \end{aligned}$$

c.

$$H(X, Y) = H(X) + H(Y | X) = H(Y) + H(X | Y)$$

$$\begin{aligned}
 H(X, Y) &= - \sum_i \sum_j p(x_i, y_j) \log p(x_i, y_j) \\
 &= - \sum_i \sum_j p(x_i) \cdot p(y_j | x_i) \log[p(x_i) \cdot p(y_j | x_i)] \\
 &= - \sum_i \sum_j p(x_i) \cdot p(y_j | x_i) [\log p(x_i) + \log p(y_j | x_i)] \\
 &= - \sum_i \sum_j p(x_i) \cdot p(y_j | x_i) \log p(x_i) - \sum_i \sum_j p(x_i) \cdot p(y_j | x_i) \log p(y_j | x_i) \\
 &= - \sum_i p(x_i) \log p(x_i) \cdot \underbrace{\sum_j p(y_j | x_i)}_{=1} - \sum_i p(x_i) \sum_j p(y_j | x_i) \log p(y_j | x_i) \\
 &= H(X) + \sum_i p(x_i) H(Y | X = x_i) = H(X) + H(Y | X)
 \end{aligned}$$

Mai general (regula de înlănțuire):

$$H(X_1, \dots, X_n) = H(X_1) + H(X_2 | X_1) + \dots + H(X_n | X_1, \dots, X_{n-1})$$

$$\begin{aligned}
 H(X_1, \dots, X_n) &= E \left[\log \frac{1}{p(x_1, \dots, x_n)} \right] \\
 &= - E_{p(x_1, \dots, x_n)} \left[\log \underbrace{p(x_1, \dots, x_n)}_{p(x_1) \cdot p(x_2 | x_1) \cdot \dots \cdot p(x_n | x_1, \dots, x_{n-1})} \right] \\
 &= - E_{p(x_1, \dots, x_n)} [\log p(x_1) + \log p(x_2 | x_1) + \dots + \log p(x_n | x_1, \dots, x_{n-1})] \\
 &= - E_{p(x_1)} [\log p(x_1)] - E_{p(x_1, x_2)} [\log p(x_2 | x_1)] - \dots \\
 &\quad - E_{p(x_1, \dots, x_n)} [\log p(x_n | x_1, \dots, x_{n-1})] \\
 &= H(X_1) + H(X_2 | X_1) + \dots + H(X_n | X_1, \dots, X_{n-1})
 \end{aligned}$$

An upper bound for the entropy of a discrete distribution

CMU, 2003 fall, T. Mitchell, A. Moore, HW1, pr. 1.1

Fie X o variabilă aleatoare discretă care ia n valori și urmează distribuția probabilistă P . Conform definiției, entropia lui X este

$$H(X) = - \sum_{i=1}^n P(X = x_i) \log_2 P(X = x_i).$$

Arătați că $H(X) \leq \log_2 n$.

Sugestie: Puteți folosi inegalitatea $\ln x \leq x - 1$ care are loc pentru orice $x > 0$.

Aşadar,

$$H(X) = \frac{1}{\ln 2} \left(- \sum_{i=1}^n P(X = x_i) \ln P(X = x_i) \right)$$

$$H(X) \leq \log_2 n \Leftrightarrow \frac{1}{\ln 2} \left(- \sum_{i=1}^n P(X = x_i) \ln P(X = x_i) \right) \leq \log_2 n$$

$$\Leftrightarrow - \sum_{i=1}^n P(x_i) \ln P(x_i) \leq \ln n$$

$$\Leftrightarrow \sum_{i=1}^n P(x_i) \ln \frac{1}{P(x_i)} - \underbrace{\left(\sum_{i=1}^n P(x_i) \right)}_1 \ln n \leq 0$$

$$\Leftrightarrow \sum_{i=1}^n P(x_i) \ln \frac{1}{P(x_i)} - \sum_{i=1}^n P(x_i) \ln n \leq 0$$

$$\Leftrightarrow \sum_{i=1}^n P(x_i) \left(\ln \frac{1}{P(x_i)} - \ln n \right) \leq 0$$

$$\Leftrightarrow \sum_{i=1}^n P(x_i) \ln \frac{1}{n P(x_i)} \leq 0$$

Aplicând inegalitatea $\ln x \leq x - 1$ pentru $x = \frac{1}{n P(x_i)}$, vom avea:

$$\sum_{i=1}^n P(x_i) \ln \frac{1}{n P(x_i)} \leq \sum_{i=1}^n P(x_i) \left(\frac{1}{n P(x_i)} - 1 \right) = \sum_{i=1}^n \frac{1}{n} - \underbrace{\sum_{i=1}^n P(x_i)}_1 = 1 - 1 = 0$$

Observație: Această margine superioară chiar este „atinsă“. De exemplu, în cazul în care o variabilă aleatoare discretă X având n valori urmează distribuția uniformă, se poate verifica imediat că $H(X) = \log_2 n$.

The particular form of the chain rule for entropies
when X and Y are independent random variables:

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

CMU, 2012 spring, R. Rosenfeld, HW2, pr. 7.b

Prove that if X and Y are independent random variables, the following property holds:
 $H(X, Y) = H(X) + H(Y)$.

Answer: Therefore,

$$\begin{aligned}
 H(X, Y) &= - \sum_{x, y} P(X = x, Y = y) \log P(X = x, Y = y) \\
 &\stackrel{indep.}{=} - \sum_{x, y} P(X = x) P(Y = y) \log(P(X = x) P(Y = y)) \\
 &= - \sum_{x, y} P(X = x) P(Y = y) (\log P(X = x) + \log P(Y = y)) \\
 &= - \left(\sum_{x, y} P(X = x) P(Y = y) \log P(X = x) \right) - \left(\sum_{x, y} P(X = x) P(Y = y) \log P(Y = y) \right) \\
 &= - \left(\sum_y P(Y = y) \sum_x P(X = x) \log P(X = x) \right) - \left(\sum_x P(X = x) \sum_y P(Y = y) \log P(Y = y) \right) \\
 &= - \left(\sum_y P(Y = y) \right) \cdot \left(\sum_x P(X = x) \log P(X = x) \right) - \left(\sum_x P(X = x) \right) \cdot \left(\sum_y P(Y = y) \log P(Y = y) \right) \\
 &= - \sum_x P(X = x) \log P(X = x) - \sum_y P(Y = y) \log P(Y = y) \\
 &= H(X) + H(Y)
 \end{aligned}$$

If X, Y are continue and independent random variables, then

$$p_{X,Y}(X = x, Y = y) = p_X(X = x)p_Y(Y = y)$$

for any x and y , where p denotes the p.d.f. corresponding to the [c.d.f. of] P .
Therefore,

$$\begin{aligned}
 H(X, Y) &\stackrel{def.}{=} - \int_X \int_Y p_{X,Y}(X = x, Y = y) \log p_{X,Y}(X = x, Y = y) dx dy \\
 &\stackrel{indep.}{=} - \int_X \int_Y p_X(X = x) p_Y(Y = y) (\log p_X(X = x) + \log p_Y(Y = y)) dx dy \\
 &\stackrel{not.}{=} - \int_X \int_Y p_X(x) p_Y(y) (\log p_X(x) + \log p_Y(y)) dx dy \\
 &= - \int_X \int_Y p_Y(y) p_X(x) \log p_X(x) dx dy - \int_X \int_Y p_X(x) p_Y(y) \log p_Y(y) dx dy \\
 &= - \int_X p_X(x) \log p_X(x) \underbrace{\left(\int_Y p_Y(y) dy \right)}_1 dx - \int_X p_X(x) \underbrace{\left(\int_Y p_Y(y) \log p_Y(y) dy \right)}_{H(Y)} dx \\
 &= - \int_X p_X(x) \log p_X(x) dx + H(Y) \underbrace{\left(\int_X p_X(x) dx \right)}_1 \\
 &= H(X) + H(Y)
 \end{aligned}$$

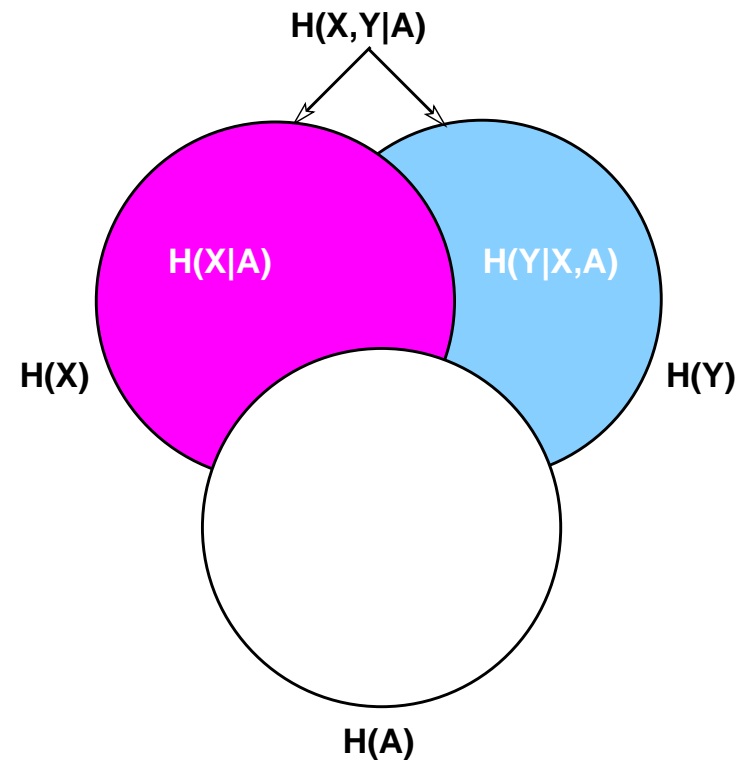
The conditional form of
the simplest case of the chain rule for entropies ($n = 2$):

$$H(X, Y|A) = H(X|A) + H(Y|X, A)$$

CMU, 2012 spring, Roni Rosenfeld, HW2, pr. 4.b

Prove that for any 3 discrete random variables X , Y and A , the following property holds:

$$\begin{aligned} H(X, Y|A) &= H(X|A) + H(Y|X, A) \\ &= H(Y|A) + H(X|Y, A). \end{aligned}$$



Answer:

$$\begin{aligned} H(X, Y|A) &= H(X, Y, A) - H(A) = H(X, Y, A) - H(Y, A) + H(Y, A) - H(A) \\ &= (H(X, Y, A) - H(Y, A)) + (H(Y, A) - H(A)) = H(X|Y, A) + H(Y|A) \\ &= H(Y|A) + H(X|Y, A). \end{aligned}$$

We've used the fact that $H(X, Y) = H(X|Y) + H(Y)$ (see CMU, 2005 fall, T. Mitchell, A. Moore, HW1, pr. 2).

**Cross-entropy (CH):
definition, some basic properties, and exemplifications**

CMU, 2011 spring, Roni Rosenfeld, HW2, pr. 3.c

Cross-entropy, $CH(p,q)$, measures the average number of bits needed to encode an event from a set of possibilities, if a coding scheme is used based on a given probability distribution q , rather than the “true” distribution p .

For discrete random variables p and q this means

$$CH(p,q) = - \sum_x p(x) \log q(x)$$

The situation for continuous random variable distributions is analogous:

$$CH(p,q) = - \int_X p(x) \log q(x) dx.$$

a. Can cross-entropy be negative? Either prove or give a counter-example.

Answer

a. No. Here follows the *proof*:

For every probability value $p(x)$ and $q(x)$, we know from definition that $0 \leq p(x) \leq 1$ and $0 \leq q(x) \leq 1$. From $q(x) \leq 1$, we conclude that $\log q(x) \leq 0$. Given that $0 \leq p(x)$, and $-\log q(x) \geq 0$, we conclude that $0 \leq -p(x) \log q(x)$. Thus, the sum of these terms will also be greater or equal to 0, so cross-entropy is never negative.

Note, however that unlike entropy, **cross-entropy is not bounded**, so it can grow to infinity if for an x , $q(x)$ is zero and $p(x)$ is not zero.

b. In many experiments, the quality of different hypothesis models are compared on a data set. Imagine you derived two different models to predict the probabilities of 7 different possible outcomes, and the probability distributions predicted by the models are q_1 , and q_2 as follows:

$$q_1 = (0.1, 0.1, 0.2, 0.3, 0.2, 0.05, 0.05) \text{ and } q_2 = (0.05, 0.1, 0.15, 0.35, 0.2, 0.1, 0.05).$$

The experiments are done on a data set with the following *empirical* distribution:

$$p_{\text{empirical}} = (0.05, 0.1, 0.2, 0.3, 0.2, 0.1, 0.05).$$

Compute the cross-entropies, $CH(p_{\text{empirical}}, q_1)$ and $CH(p_{\text{empirical}}, q_2)$. Which model has a lower cross-entropy? Is this model guaranteed to be a better one? Explain your answer. Can cross-entropy be negative? Either prove or give a counter-example.

Answer

b. Using the cross-entropy formula we see that:

$$\begin{aligned} CH(p_{\text{empirical}}, q_1) &= -(0.05 \log 0.1 + 0.1 \log 0.1 + 0.2 \log 0.2 + 0.3 \log 0.3 + \\ &\quad 0.2 \log 0.2 + 0.1 \log 0.05 + 0.05 \log 0.05) = 2.596 \text{ bits} \end{aligned}$$

$$\begin{aligned} CH(p_{\text{empirical}}, q_2) &= -(0.05 \log 0.05 + 0.1 \log 0.1 + 0.2 \log 0.15 + 0.3 \log 0.35 + \\ &\quad 0.2 \log 0.2 + 0.1 \log 0.1 + 0.05 \log 0.05) = 2.56 \text{ bits.} \end{aligned}$$

The $p_{\text{empirical}}$ distribution has a lower cross-entropy with the model q_2 , so it is reasonable to say that q_2 is a better choice.

However, it is not guaranteed that this model is always the better model, because we are working with an “empirical” distribution here, and the “true” distribution might not be reflected in this empirical distribution completely. Usually, sampling bias and insufficient training data samples will widen the gap between the true distribution and the empirical one, so in practice when designing the experiment, you should always have that in mind, and if possible use techniques that minimize these risks.

Relative entropy a.k.a. the Kulback-Leibler divergence,
and the [relationship to] information gain;
some basic properties

CMU, 2007 fall, C. Guestrin, HW1, pr. 1.2

[adapted by Liviu Ciortuz]

The *relative entropy* — also known as the *Kullback-Leibler (KL) divergence* — from a distribution p to a distribution q is defined as

$$KL(p||q) \stackrel{def.}{=} - \sum_{x \in X} p(x) \log \frac{q(x)}{p(x)}$$

From an information theory perspective, the KL-divergence specifies the number of additional bits required on average to transmit values of X if the values are distributed with respect to p but we encode them assuming the distribution q .

Notes

1. KL is not a *distance measure*, since it is not symmetric (i.e., in general $KL(p||q) \neq KL(q||p)$).

Another measure, which is defined as $JSD(p||q) = \frac{1}{2}(KL(p||q) + KL(q||p))$, and is called the **Jensen-Shannon divergence** is symmetric.

2. The quantity

$$\begin{aligned} d(X, Y) &\stackrel{def.}{=} H(X, Y) - IG(X, Y) = H(X) + H(Y) - 2IG(X, Y) \\ &= H(X | Y) + H(Y | X) \end{aligned}$$

known as **variation of information**, is a distance metric, i.e., it is non-negative, symmetric, implies indiscernability, and satisfies the triangle inequality.

a. Show that $KL(p||q) \geq 0$, and $KL(p||q) = 0$ iff $p(x) = q(x)$ for all x .
(More generally, the smaller the KL-divergence, the more similar the two distributions.)

Indicație:

Pentru a demonstra punctul acesta puteți folosi **inegalitatea lui Jensen**:

Dacă $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ este o funcție convexă, atunci pentru orice $t \in [0, 1]$ și orice $x_1, x_2 \in \mathbb{R}$ urmează $\varphi(tx_1 + (1 - t)x_2) \leq t\varphi(x_1) + (1 - t)\varphi(x_2)$.

Dacă φ este funcție strict convexă, atunci egalitatea are loc doar dacă $x_1 = x_2$.

Mai general, pentru orice $a_i \geq 0$, $i = 1, \dots, n$ cu $\sum_i a_i \neq 0$ și orice $x_i \in \mathbb{R}$, $i = 1, \dots, n$, avem

$$\varphi\left(\frac{\sum_i a_i x_i}{\sum_j a_j}\right) \leq \frac{\sum_i a_i \varphi(x_i)}{\sum_j a_j}.$$

Dacă φ este strict convexă, atunci egalitatea are loc doar dacă $x_1 = \dots = x_n$.

Evident, rezultate similare pot fi formulate și pentru funcții concave.

Answer

Vom dovedi inegalitatea $KL(p||q) \geq 0$ folosind inegalitatea lui Jensen, în expresia căreia vom înlocui φ cu funcția convexă $-\log_2$, pe a_i cu $p(x_i)$ și pe x_i cu $\frac{q(x_i)}{p(x_i)}$.

(Pentru conveniență, în cele ce urmează vor renunța la indicele variabilei x .)

Vom avea:

$$\begin{aligned}
 KL(p \parallel q) &\stackrel{def.}{=} - \sum_x p(x) \log \frac{q(x)}{p(x)} \\
 &\stackrel{Jensen}{\geq} - \log \left(\sum_x p(x) \frac{q(x)}{p(x)} \right) = - \log \left(\underbrace{\sum_x q(x)}_1 \right) = - \log 1 = 0
 \end{aligned}$$

Așadar, $KL(p \parallel q) \geq 0$, oricare ar fi distribuțiile (discrete) p și q .

Vom demonstra acum că $KL(p||q) = 0 \Leftrightarrow p = q$.

\Leftarrow

Egalitatea $p(x) = q(x)$ implică $\frac{q(x)}{p(x)} = 1$, deci $\log \frac{q(x)}{p(x)} = 0$ pentru orice x , de unde rezultă imediat $KL(p||q) = 0$.

\Rightarrow

Știm că în inegalitatea lui Jensen are loc egalitatea doar în cazul în care $x_i = x_j$ pentru orice i și j .

În cazul de față, această condiție se traduce prin faptul că raportul $\frac{q(x)}{p(x)}$ este același pentru orice valoare a lui x .

Ținând cont că $\sum_x p(x) = 1$ și $\sum_x p(x) \frac{q(x)}{p(x)} = \sum_x q(x) = 1$, rezultă că $\frac{q(x)}{p(x)} = 1$ sau, altfel spus, $p(x) = q(x)$ pentru orice x , ceea ce înseamnă că distribuțiile p și q sunt identice.

b. We can define the *information gain* as the KL-divergence from the observed joint distribution of X and Y to the product of their observed marginals:

$$\begin{aligned}
 IG(X, Y) &\stackrel{\text{def.}}{=} KL(p_{X,Y} \parallel (p_X p_Y)) = - \sum_x \sum_y p_{X,Y}(x, y) \log \left(\frac{p_X(x)p_Y(y)}{p_{X,Y}(x, y)} \right) \\
 &\stackrel{\text{not.}}{=} - \sum_x \sum_y p(x, y) \log \left(\frac{p(x)p(y)}{p(x, y)} \right)
 \end{aligned}$$

Prove that this definition of information gain is equivalent to the one given in problem CMU, 2005 fall, T. Mitchell, A. Moore, HW1, pr. 2. That is, show that $IG(X, Y) = H[X] - H[X|Y] = H[Y] - H[Y|X]$, starting from the definition in terms of KL-divergence.

Remark:

It follows that

$$\begin{aligned}
 IG(X, Y) &= \sum_y p(y) \sum_x p(x | y) \log \frac{p(x | y)}{p(x)} = \sum_y p(y) KL(p_{X|Y} \parallel p_X) \\
 &= E_Y[KL(p_{X|Y} \parallel p_X)]
 \end{aligned}$$

Answer

By making use of the multiplication rule, namely $p(x, y) = p(x | y)p(y)$, we will have:

$$\begin{aligned}
 & KL(p_{X,Y} || (p_X p_Y)) \\
 & \stackrel{\text{def. } KL}{=} - \sum_x \sum_y p(x, y) \log \left(\frac{p(x)p(y)}{p(x, y)} \right) \\
 & = - \sum_x \sum_y p(x, y) \log \left(\frac{p(x)p(y)}{p(x | y)p(y)} \right) = - \sum_x \sum_y p(x, y) [\log p(x) - \log p(x | y)] \\
 & = - \sum_x \sum_y p(x, y) \log p(x) - \left(- \sum_x \sum_y p(x, y) \log p(x | y) \right) \\
 & = - \sum_x \log p(x) \underbrace{\sum_y p(x, y)}_{=p(x)} - H[X | Y] \\
 & = H[X] - H[X | Y] = IG(X, Y)
 \end{aligned}$$

c.

A direct consequence of parts a. and b. is that $IG(X, Y) \geq 0$ (and therefore $H(X) \geq H(X|Y)$ and $H(Y) \geq H(Y|X)$) for any discrete random variables X and Y .

Prove that $IG(X, Y) = 0$ iff X and Y are independent.

Answer:

This is also an immediate consequence of parts a. and b. already proven:

$$IG(X, Y) = 0 \stackrel{(b)}{\Leftrightarrow} KL(p_{X,Y} || p_X p_Y) = 0 \stackrel{(a)}{\Leftrightarrow} X \text{ and } Y \text{ are independent.}$$

Putem demonstra inegalitatea $IG(X, Y) \geq 0$ și în manieră directă, folosind rezultatul de la punctul b. și aplicând inegalitatea lui Jensen în forma generalizată, cu următoarele „amendamente“:

- în locul unui singur indice, se vor considera doi indici (așadar în loc de a_i și x_i vom avea a_{ij} și respectiv x_{ij});
- vom lua $\varphi = -\log_2$ iar $a_{ij} \leftarrow p(x_i, y_j)$ și $x_{ij} \leftarrow \frac{p(x_i)p(y_j)}{p(x_i, y_j)}$;
- în fine, vom ține cont că $\sum_i \sum_j p(x_i, y_j) = 1$.

Prin urmare,

$$\begin{aligned}
 IG(X, Y) &= \sum_i \sum_j p(x_i, y_j) \log \frac{p(x_i, y_j)}{p(x_i) \cdot p(y_j)} = \sum_i \sum_j p(x_i, y_j) \left[-\log \frac{p(x_i) \cdot p(y_j)}{p(x_i, y_j)} \right] \\
 &\geq -\log \left(\sum_i \sum_j p(x_i, y_j) \frac{p(x_i) \cdot p(y_j)}{p(x_i, y_j)} \right) = -\log \left(\sum_i \sum_j p(x_i) \cdot p(y_j) \right) \\
 &= -\log \left(\underbrace{\sum_i p(x_i)}_1 \cdot \underbrace{\sum_j p(y_j)}_1 \right) = -\log 1 = 0
 \end{aligned}$$

În concluzie, $IG(X, Y) \geq 0$.

Remark (cont'd)

Dacă X și Y sunt variabilele independente, atunci $p(x_i, y_j) = p(x_i)p(y_j)$ pentru orice i și j .

În consecință, toți logaritmi din partea dreaptă a primei egalități din calculul de mai sus sunt 0 și rezultă $IG(X, Y) = 0$.

Invers, presupunând că $IG(X, Y) = 0$, vom ține cont de faptul că putem exprima câștigul de informație cu ajutorul divergenței KL și vom aplica un raționament similar cu cel de la punctul a .

Rezultă că $\frac{p(x_i)p(y_j)}{p(x_i, y_j)} = 1$ și deci $p(x_i)p(y_j) = p(x_i, y_j)$ pentru orice i și j .

Aceasta echivalează cu a spune că variabilele X și Y sunt independente.

**Proving [in a direct manner] that
the Information Gain is always positive or 0**

(an indirect proof was made at CMU, 2007 fall, Carlos Guestrin, HW1, pr. 1.2)

Liviu Ciortuz, 2017

Definiția câștigului de informație (sau: a informației mutuale) al unei variabile aleatoare X în raport cu o altă variabilă aleatoare Y este

$$IG(X, Y) = H(X) - H(X | Y) = H(Y) - H(Y | X).$$

La CMU, 2007 fall, Carlos Guestrin, HW1, pr. 1.2 s-a demonstrat — pentru cazul în care X și Y sunt discrete — că $IG(X, Y) = KL(P_{X,Y} || P_X P_Y)$, unde KL desemnează *entropia relativă* (sau: *divergența Kullback-Leibler*), P_X și P_Y sunt distribuțiile variabilelor X și, respectiv, Y , iar $P_{X,Y}$ este distribuția corelată a acestor variabile. Tot la CMU, 2007 fall, Carlos Guestrin, HW1, pr. 1.2 s-a arătat că divergența KL este întotdeauna ne-negativă. În consecință, $IG(X, Y) \geq 0$ pentru orice X și Y .

La acest exercițiu vă cerem să demonstrați inegalitatea $IG(X, Y) \geq 0$ în manieră directă, plecând de la prima definiție dată mai sus, fără a [mai] apela la divergența Kullback-Leibler.

Sugestie: Puteți folosi următoarea formă a inegalității lui Jensen:

$$\sum_{i=1}^n a_i \log x_i \leq \log \left(\sum_{i=1}^n a_i x_i \right)$$

unde baza logaritmului se consideră supraunitară, $a_i \geq 0$ pentru $i = 1, \dots, n$ și $\sum_{i=1}^n a_i = 1$.

Observație: Avantajul la această problemă, comparativ cu CMU, 2007 fall, Carlos Guestrin, HW1, pr. 1.2.a, este că aici se lucrează cu o singură distribuție (p), nu cu două distribuții (p și q). Totuși, demonstrația de aici va fi mai laborioasă.

Answer (in Romanian)

Presupunem că valorile variabilei X sunt x_1, x_2, \dots, x_n , iar valorile variabilei Y sunt y_1, y_2, \dots, y_m . Avem:

$$\begin{aligned} IG(X, Y) &\stackrel{\text{def.}}{=} H(X) - H(X|Y) \\ &\stackrel{\text{def.}}{=} \sum_{i=1}^n -P(x_i) \log_2 P(x_i) - \sum_{j=1}^m P(y_j) \sum_{i=1}^n (-P(x_i|y_j) \log_2 P(x_i|y_j)) \end{aligned}$$

$$-IG(X, Y) = \sum_{i=1}^n P(x_i) \log_2 P(x_i) - \sum_{j=1}^m P(y_j) \sum_{i=1}^n P(x_i|y_j) \log_2 P(x_i|y_j)$$

$$\stackrel{\text{def.}}{=} \text{prob. marg.} \quad \sum_{i=1}^n \left(\sum_{j=1}^m P(x_i, y_j) \right) \log_2 P(x_i) - \sum_{j=1}^m P(y_j) \sum_{i=1}^n P(x_i|y_j) \log_2 P(x_i|y_j)$$

$$\stackrel{\text{distrib.},+}{=} \sum_{i=1}^n \sum_{j=1}^m P(x_i, y_j) \log_2 P(x_i) - \sum_{j=1}^m \sum_{i=1}^n P(y_j) P(x_i|y_j) \log_2 P(x_i|y_j)$$

$$\stackrel{\text{def.}}{=} \text{prob. cond.} \quad \sum_{i=1}^n \sum_{j=1}^m P(x_i, y_j) \log_2 P(x_i) - \sum_{j=1}^m \sum_{i=1}^n P(x_i, y_j) \log_2 P(x_i|y_j)$$

$$\stackrel{\text{distrib.},+}{=} \sum_{i=1}^n \sum_{j=1}^m P(x_i, y_j) (\log_2 P(x_i) - \log_2 P(x_i|y_j))$$

$$\stackrel{\text{prop.}}{=} \text{log.} \quad \sum_{i=1}^n \sum_{j=1}^m P(x_i, y_j) \log_2 \frac{P(x_i)}{P(x_i|y_j)} \stackrel{\text{reg. de}}{=} \text{multipl.} \quad \sum_{i=1}^n \sum_{j=1}^m P(x_i|y_j) P(y_j) \log_2 \frac{P(x_i)}{P(x_i|y_j)}$$

$$\stackrel{\text{distrib.},+}{=} \sum_{j=1}^m P(y_j) \sum_{i=1}^n \underbrace{P(x_i|y_j)}_{a_i} \log_2 \frac{P(x_i)}{P(x_i|y_j)}$$

Întrucât pe de o parte $P(x_i|y_j) \geq 0$ și pe de altă parte $\sum_{i=1}^n P(x_i|y_j) = 1$ pentru fiecare valoare y_j a lui Y în parte, putem aplica inegalitatea lui Jensen pentru cea de-a doua sumă din ultima expresie de mai sus — mai exact, pentru fiecare valoare a indicelui j în parte — și obținem:

$$-IG(X, Y) \leq \sum_{j=1}^m P(y_j) \log_2 \left(\sum_{i=1}^n P(x_i|y_j) \frac{P(x_i)}{P(x_i|y_j)} \right) = \sum_{j=1}^m P(y_j) \log_2 \underbrace{\left(\sum_{i=1}^n P(x_i) \right)}_1 = 0$$

Prin urmare, $IG(X, Y) \geq 0$.

Using Information Gain / Mutual Information for doing Feature Selection

CMU, 2009 spring, Ziv Bar-Joseph, HW5, pr. 6

Given the following observations for input binary features X_1, X_2, X_3, X_4, X_5 , and output binary label Y , we would like to use a filter approach to reduce the feature space of $\{X_1, X_2, X_3, X_4, X_5\}$.

X_1	X_2	X_3	X_4	X_5	Y
0	1	1	0	1	0
1	0	0	0	1	0
0	1	0	1	0	1
1	1	1	1	0	1
0	1	1	0	0	1
0	0	0	1	1	1
1	0	0	1	0	1
1	1	1	0	1	1

- Calculate the mutual information $MI(X_i, Y)$ for each i .
- Accordingly choose the smallest subset of features such that the best classifier trained on the reduced feature set will perform at least as well as the best classifier trained on the whole feature set. Explain the reasons behind your choice.

Answer

a. As shown at CMU, 2007 fall, Carlos Guestrin, HW1, pr. 1.2, the mutual information can be calculated as:

$$MI(X, Y) \stackrel{def.}{=} - \sum_x \sum_y P(x, y) \log \left(\frac{P(x) P(y)}{P(x, y)} \right).$$

The marginal probabilities that we need to compute are:

$$\begin{aligned} P(Y = 0) &= 1/4, & P(Y = 1) &= 3/4 \\ P(X_1 = 0) &= 1/2, & P(X_1 = 1) &= 1/2 \\ P(X_2 = 0) &= 3/8, & P(X_1 = 1) &= 5/8 \\ P(X_3 = 0) &= 1/2, & P(X_3 = 1) &= 1/2 \\ P(X_4 = 0) &= 1/2, & P(X_4 = 1) &= 1/2 \\ P(X_5 = 0) &= 1/2, & P(X_5 = 1) &= 1/2. \end{aligned}$$

The joint probabilities that we need to compute are:

$$P(X_1 = 0, Y = 0) = 1/8, \quad P(X_1 = 0, Y = 1) = P(X_1 = 0) - P(X_1 = 0, Y = 0) = 3/8$$

$$P(X_1 = 1, Y = 0) = 1/8, \quad P(X_1 = 1, Y = 1) = P(X_1 = 1) - P(X_1 = 1, Y = 0) = 3/8$$

$$P(X_2 = 0, Y = 0) = 1/8, \quad P(X_2 = 0, Y = 1) = P(X_2 = 0) - P(X_2 = 0, Y = 0) = 1/4$$

$$P(X_2 = 1, Y = 0) = 1/8, \quad P(X_2 = 1, Y = 1) = P(X_2 = 1) - P(X_2 = 1, Y = 0) = 1/2$$

$$P(X_3 = 0, Y = 0) = 1/8, \quad P(X_3 = 0, Y = 1) = P(X_3 = 0) - P(X_3 = 0, Y = 0) = 3/8$$

$$P(X_3 = 1, Y = 0) = 1/8, \quad P(X_3 = 1, Y = 1) = P(X_3 = 1) - P(X_3 = 1, Y = 0) = 3/8$$

$$P(X_4 = 0, Y = 0) = 1/4, \quad P(X_4 = 0, Y = 1) = P(X_4 = 0) - P(X_4 = 0, Y = 0) = 1/4$$

$$P(X_4 = 1, Y = 0) = 0, \quad P(X_4 = 1, Y = 1) = P(X_4 = 1) - P(X_4 = 1, Y = 0) = 1/2$$

$$P(X_5 = 0, Y = 0) = 0, \quad P(X_5 = 0, Y = 1) = P(X_5 = 0) - P(X_5 = 0, Y = 0) = 1/2$$

$$P(X_5 = 1, Y = 0) = 1/4, \quad P(X_5 = 1, Y = 1) = P(X_5 = 1) - P(X_5 = 1, Y = 0) = 1/4.$$

Using this information, the mutual information for each feature is:

$$MI(X_1, Y) = 0, \quad MI(X_2, Y) = 0.01571, \quad MI(X_3, Y) = 0, \quad MI(X_4, Y) = 0.3113 \quad \text{and} \\ MI(X_5, Y) = 0.3113.$$

b. In order to select a set of features, we can prioritize the ones with more mutual information because they are less independent to Y .

By looking at the results of the previous question, we can see that X_5 and X_4 are the features with more mutual information (0.3113), followed by X_2 (0.0157) and, finally, X_1 and X_3 that do not have mutual information with Y (i.e., X_1 and X_3 are independent to Y).

By inspection of the data, we can see that if we select X_5 , X_4 and X_2 there are two samples of different classes with the same features ($X_2 = 1$, $X_4 = 0$, $X_5 = 1$). To avoid this problem, we can add X_1 as an extra feature.

Note that, although the mutual information of X_1 with Y is zero, it does not mean that the combination of X_1 with other features will also have zero mutual information.