

KL - Divergence: Definition:

~~*E~~ Kull back - Liebler Divergence is a distance that is not a metric.

Let X_a and X_b be two discrete distributions.

$$d_{KL}(X_a, X_b) = d_{KL}(a, b) = \sum_{i=1}^n a_i \ln \left(\frac{a_i}{b_i} \right)$$

Example Bag-of-Words

Consider D -dimensional space with

$$D = 11.$$

For each co-ordinate list the corresponding word as:

(am, and, do, ham, I, jelly, like, not, Sam, ¹⁰tham, zebra)

Use the text from the 4 short documents following

D_1 : I am Sam

D_2 : Sam I am

D_3 : I do not like jelly and ham.

D_4 : I do not, do not, like tham, Sam I am

For each of the above documents form the representative vectors.

$$V_1 = \begin{pmatrix} 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 \end{pmatrix}$$

$$V_2 = \begin{pmatrix} 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 \end{pmatrix}$$

$$V_3 = \begin{pmatrix} 0 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 \end{pmatrix}$$

$$V_4 = \begin{pmatrix} 1 & 0 & 2 & 0 & 2 & 0 & 1 & 2 & 1 & 0 & 0 \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 \end{pmatrix}$$

Also Find the distance ^(forward) between the documents D_1, D_2, D_3 and D_4 .

Find the L_2 distance or d_2 distance or Euclidean distance between the vectors.

Formula:

$$d_2(a, b) = \|a - b\|_2 = \sqrt{\sum_{l=1}^n (a_l - b_l)^2}$$

$$d_2(N_1, N_2) = \sqrt{(1-1)^2 + (0-0)^2 + (0-0)^2 + (0-0)^2 + (1-1)^2 + (0-0)^2 + (0-0)^2 + (0-0)^2}$$

$$= \sqrt{0+0+\dots+0} = 0$$

$$d_2(N_1, N_3) = \sqrt{(1-0)^2 + (0-1)^2 + (0-1)^2 + (0-1)^2 + (1-1)^2 + (0-1)^2 + (0-1)^2 + (0-0)^2 + (0-0)^2}$$

$$= \sqrt{1+1+1+0+1+1+1+0+0}$$

$$d_2(N_1, N_3) = \sqrt{8} = 2.83$$

$$d_2(N_1, N_4) = \sqrt{(1-1)^2 + (0-0)^2 + (0-2)^2 + (1-2)^2 + (0-0)^2 + (0-1)^2 + (0-2)^2 + (1-1)^2 + (0-1)^2 + (0-0)^2}$$

$$= \sqrt{0+0+4+1+0+1+4+0+1+0}$$

$$= \sqrt{0+0+4+0+1+0+1+4+0+1+0}$$

$$= \sqrt{11}$$

$$d_2(N_1, N_4) = 3.32$$

$$d_2(V_2, V_3) = \sqrt{(1-0)^2 + (0-1)^2 + (0-1)^2 + (0-1)^2 + (1-1)^2 + (0-1)^2 + (0-1)^2 + (0-1)^2 + (0-0)^2 + (0-0)^2}$$

$$= \sqrt{1+1+1+1+0+1+1+1+1+0+0}$$

$$= \sqrt{8} = 2.83$$

| vectors | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---------|---|---|---|---|---|---|---|---|---|----|----|
| V_1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 |
| V_2 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 |
| V_3 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 |
| V_4 | 1 | 0 | 2 | 0 | 2 | 0 | 1 | 2 | 1 | 1 | 0 |

$$d_2(V_2, V_4) = \sqrt{(1-1)^2 + (0+0)^2 + (0-2)^2 + (0-0)^2 + (1-2)^2 + (0-0)^2 + (0-1)^2 + (0-2)^2 + (1-1)^2 + (0-0)^2}$$

$$= \sqrt{0+0+4+0+1+0+1+4+0+1+0}$$

$$= \sqrt{11} = 3.32$$

$$d_3(V_3, V_4) = \sqrt{(0-1)^2 + (1-0)^2 + (1-2)^2 + (1-0)^2 + (1-2)^2 + (1-0)^2 + (1-1)^2 + (1-2)^2 + (0-1)^2 + (0-0)^2}$$

$$= \sqrt{1+1+1+0+1+1+1+1+0+0}$$

$$d_3(V_3, V_4) = \sqrt{9} = 3$$

Compute Jaccard distance among the documents D_1, D_2, D_3 and D_4 for words & grams.

$$D_1: G_1 = \{ [I am], [am sam] \}$$

$$D_2: G_2 = \{ [sam I], [I am] \}$$

$$D_3: G_3 = \{ [I do], [do not], [not like], [like], [do], [do not], [not like], [like] \}$$

Jelly don't and want

Solution:

Jaccard distance between two sets A and B

$$d_J(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|}$$

Here $|A \cap B|$ represents no of elements in $A \cap B$.
 $|A \cup B|$ " " " no of elements in $A \cup B$.

$$d_J(D_1, D_2) = 1 - \frac{1}{3} = \frac{2}{3} = 0.667$$

~~= 2~~

$$d_J(D_1, D_3) = 1 - \frac{0}{0}$$

Problem: Compute Jaccard distance among the documents D_1, D_2, D_3 , and D_4 for word 2-grams.

$$D_1: G_1 = \{ [I, am], [am, Sam] \}$$

$$D_2: G_2 = \{ [Sam, I], [I, am] \}$$

$$D_3: G_3 = \{ [I, do], [do, not], [not, like], [like, them], [them, Sam], [Sam, I], [I, do], [do, jelly], [jelly, and], [and, for]$$

$D_4: G_4 = \{ [I, do], [do, not], [not, do], [not, like], [like, them], [them, Sam], [Sam, I], [I, do], [do, jelly], [jelly, and], [and, for]$

Each gram is shown as a set of two words in square brackets.

Solution:

Formula: Jaccard distance between any two sets A and B is

$$d_J(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|}$$

$$|A \cap B| = \text{Cardinality of } A \cap B$$

$$= \text{No of elements in } A \cap B$$

$$|A \cup B| = \text{No of elements in } A \cup B$$

$$|D_1| = 2, |D_2| = 2, |D_3| = 6$$

$$|D_4| = 8$$

$$|D_1 \cap D_2| = 1, |D_1 \cap D_3| = 0$$

$$|D_1 \cup D_2| = 3, |D_1 \cup D_3| = 8$$

$$|D_2 \cap D_3| = 0$$

$$|D_2 \cup D_3| = 8$$

$$|D_3 \cup D_4| = 11$$

$$|D_1 \cup D_4| = 3$$

$|D|$ represents no of grams in the docemat D.

$$|D_1 \cap D_4| = 1$$

$$|D_1 \cup D_4| = 9$$

$$|D_2 \cap D_4| = 2$$

$$|D_2 \cup D_4| = 8$$

$$P(D_1, D_2) = 1 - \frac{|D_1 \cap D_2|}{|D_1 \cup D_2|} = 1 - \frac{1}{3} = 1 - 0.33 = 0.667$$

$$P(D_1, D_3) = 1 - \frac{|D_1 \cap D_3|}{|D_1 \cup D_3|} = 1 - \frac{0}{8} = 1 - 0 = 1$$

$$P(D_1, D_4) = 1 - \frac{|D_1 \cap D_4|}{|D_1 \cup D_4|} = 1 - \frac{1}{9} = 0.889$$

$$P(D_2, D_3) = 1 - \frac{|D_2 \cap D_3|}{|D_2 \cup D_3|} = 1 - \frac{0}{8} = 1$$

$$P(D_2, D_4) = 1 - \frac{|D_2 \cap D_4|}{|D_2 \cup D_4|} = 1 - \frac{2}{8} = 1 - 0.25 = 0.75$$

$$P(D_3, D_4) = 1 - \frac{|D_3 \cap D_4|}{|D_3 \cup D_4|} = 1 - \frac{3}{11} = 0.727$$

Problem 6 Marks:

Construct characters 4-grams for the document D_4 .

$D_4: G_4 = \{ [I\ do], [do\ not], [not\ do], [not\ like], [like\ them], [them\ sam], [sam\ I], [I\ am] \}$

Solution:

Character 4 grams for D_4

$= \{ [idon], [don], [not], [notd], [otdo], [tdon], [otdi], [elik], [like], [like], [keth], [ethe], [them], [hem], [amsa], [nsam], [sam], [amia], [miam] \}$.

Problem 4: 4 Marks:

Find the Jaccard distance between two

sets $A = \{ 1, 2, 4, 8 \}$, $B = \{ 1, 2, 3 \}$

Solution:

$$d_J(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|}$$

Here $A \cap B = \{ 1, 2 \} \Rightarrow |A \cap B| = 2$

$A \cup B = \{ 1, 2, 3, 4, 8 \} \Rightarrow |A \cup B| = 5$

$d_J(A, B) = 1 - \frac{2}{5} = 1 - 0.4 = 0.6$

$d_J(A, B) = 0.6$

6 Marks:

Construct ~~2~~ characters 3-grams for the document D_3 :

$D_3 : G_3 = \{ [I do], [do not], [not like], [like, jelly], [jelly and], [and I and ham] \}$

Solution:

Character 3-grams for D_3 is

$\{ [ido], [don], [ono], [not], [not], [ota], [eli], [lik], [like], [kej], [eje], [jel], [ell], [xy], [xxy], [yga], [yan], [and], [ndh], [dhi], [ham] \}$

$= \{ [ido], [don], [ono], [not], [not], [ota], [eli], [lik], [like], [kej], [eje], [jel], [ell], [xy], [xxy], [yga], [yan], [and], [ndh], [dha], [ham] \}$

4 Marks:

Construct characterograms for the document D_2 : $G_2 = \{ [sam I], [I am] \}$

Solution:

Character 2-grams for the Document D_2 is

$\{ [sa], [am], [mi], [ia], [am] \}$

$= \{ [sa], [am], [mi], [ia] \}$

[∴ Repetition not allowed]