

Lecture 7: Information theory



EMAT31530/Feb 2018/Raul Santos-Rodriguez

Have a look at ...

... David MacKay, Information Theory, *Inference, and Learning Algorithms*. Cambridge University Press, 2003. (Ch. 2)

... T. Cover and J. Tomas, *Elements of Information Theory*. Wiley, 2006.

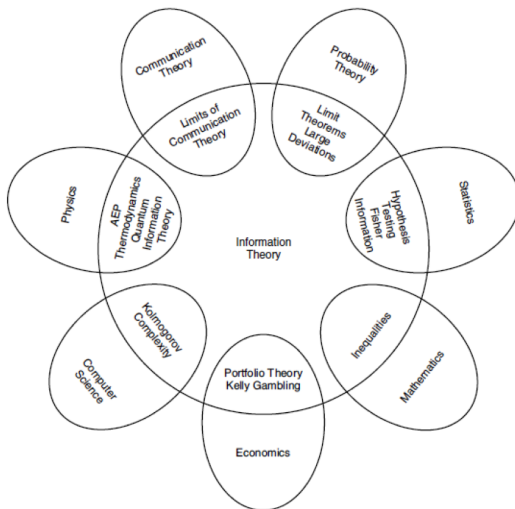
... Information Theoretical Estimators (ITE) Toolbox.

... 17 equations that changed the course of history!

Information theory is the use of probability theory to quantify and measure information. Basic concepts:

- Entropy and Information
- Mutual information
- Kullback-Leibler divergence

Information theory



[Cover and Thomas, 2006]

Information theory



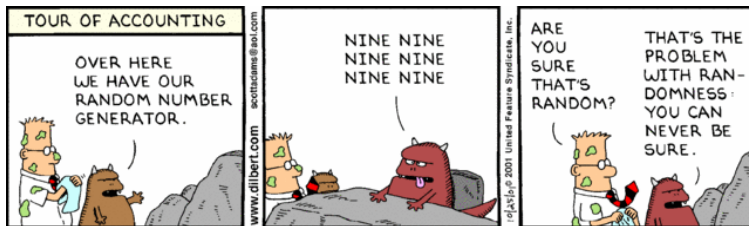
[Cover and Thomas, 2006]

Information theory:

- Developed by Shannon in the 40s
- Idea: Maximising the amount of information that can be transmitted over an imperfect communication channel

What is Information?

- "The sun will come up tomorrow"
- "It will rain tomorrow"
- "There was an earthquake this morning"



Probability theory revisited: probability distribution

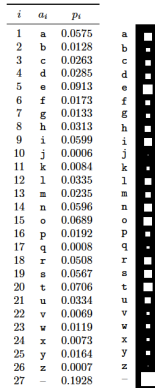


Figure 2.1. Probability distribution over the 27 outcomes for a randomly selected letter in an English language document (estimated from *The Frequently Asked Questions Manual for Linux*). The picture shows the probabilities by the areas of white squares.

Definition

Entropy is a measure of the **uncertainty** associated with a distribution.

$$H(X) = - \sum_x p(x) \log p(x)$$

Lower bound on the number of *bits* that it takes to transmit messages!

Definition

Entropy is a measure of the **uncertainty** associated with a distribution.

$$H(X) = - \sum_x p(x) \log p(x)$$

Lower bound on the number of *bits* that it takes to transmit messages!

i	a_i	p_i	$h(p_i)$
1	a	.0575	4.1
2	b	.0128	6.3
3	c	.0263	5.2
4	d	.0285	5.1
5	e	.0913	3.5
6	f	.0173	5.9
7	g	.0133	6.2
8	h	.0313	5.0
9	i	.0599	4.1
10	j	.0006	10.7
11	k	.0084	6.9
12	l	.0335	4.9
13	m	.0235	5.4
14	n	.0596	4.1
15	o	.0689	3.9
16	p	.0192	5.7
17	q	.0008	10.3
18	r	.0508	4.3
19	s	.0567	4.1
20	t	.0706	3.8
21	u	.0334	4.9
22	v	.0069	7.2
23	w	.0119	6.4
24	x	.0073	7.1
25	y	.0164	5.9
26	z	.0007	10.4
27	-	.1928	2.4

$$\sum_i p_i \log_2 \frac{1}{p_i} \quad 4.1$$

Table 2.9. Shannon information contents of the outcomes a–z.

Entropy

Entropy measures the amount of information in a Random Variable; it can also be seen as the average length of the message needed to transmit an outcome of that variable using the optimal code.

The entropy of a randomly selected letter in an English document is about 4.11 bits, assuming its probability is as given in table 2.9. We obtain this number by averaging $\log 1/p_i$ (shown in the fourth column) under the probability distribution p_i (shown in the third column).

Entropy

Definition

Entropy is a measure of the **uncertainty** associated with a distribution.

$$H(X) = - \sum_x p(x) \log p(x)$$

Lower bound on the number of bits that it takes to transmit messages!

i	a_i	p_i	$\log 1/p_i$
1	a	0.165	2.43
2	b	0.033	4.91
3	c	0.028	5.46
4	d	0.042	4.57
5	e	0.105	3.25
6	f	0.009	7.06
7	g	0.002	10.00
8	h	0.054	4.22
9	i	0.070	3.84
10	j	0.001	10.00
11	k	0.005	7.61
12	l	0.047	4.41
13	m	0.024	5.72
14	n	0.067	4.06
15	o	0.075	3.75
16	p	0.019	6.32
17	q	0.001	10.00
18	r	0.005	7.61
19	s	0.005	7.61
20	t	0.090	3.44
21	u	0.027	5.54
22	v	0.001	10.00
23	w	0.024	5.72
24	x	0.001	10.00
25	y	0.001	10.00
26	z	0.001	10.00
27	space	0.105	3.25
		$\sum_{i=1}^{27} p_i \log 1/p_i$	4.11

Table 2.9: Relative information content of the messages in a.

Entropy

Entropy

Definition

Entropy is a measure of the **uncertainty** associated with a distribution.

$$H(X) = - \sum_x p(x) \log p(x)$$

Lower bound on the number of bits that it takes to transmit messages!

x	p	p log p
1	0.5000	-1.0000
2	0.2500	-0.6931
3	0.1667	-0.5646
4	0.1250	-0.4419
5	0.1000	-0.3219
6	0.0833	-0.2592
7	0.0714	-0.2091
8	0.0625	-0.1675
9	0.0556	-0.1321
10	0.0500	-0.1021
11	0.0455	-0.0772
12	0.0417	-0.0564
13	0.0385	-0.0398
14	0.0357	-0.0271
15	0.0333	-0.0171
16	0.0312	-0.0098
17	0.0294	-0.0051
18	0.0278	-0.0026
19	0.0263	-0.0013
20	0.0250	-0.0007
21	0.0238	-0.0004
22	0.0227	-0.0002
23	0.0217	-0.0001
24	0.0208	-0.0001
25	0.0200	-0.0000
26	0.0192	-0.0000
27	0.0185	-0.0000
28	0.0179	-0.0000
29	0.0172	-0.0000
30	0.0167	-0.0000
31	0.0161	-0.0000
32	0.0156	-0.0000
33	0.0152	-0.0000
34	0.0147	-0.0000
35	0.0143	-0.0000
36	0.0139	-0.0000
37	0.0136	-0.0000
38	0.0133	-0.0000
39	0.0130	-0.0000
40	0.0128	-0.0000
41	0.0126	-0.0000
42	0.0124	-0.0000
43	0.0122	-0.0000
44	0.0120	-0.0000
45	0.0118	-0.0000
46	0.0117	-0.0000
47	0.0115	-0.0000
48	0.0114	-0.0000
49	0.0113	-0.0000
50	0.0112	-0.0000
51	0.0111	-0.0000
52	0.0110	-0.0000
53	0.0109	-0.0000
54	0.0108	-0.0000
55	0.0107	-0.0000
56	0.0106	-0.0000
57	0.0105	-0.0000
58	0.0104	-0.0000
59	0.0103	-0.0000
60	0.0103	-0.0000
61	0.0102	-0.0000
62	0.0101	-0.0000
63	0.0101	-0.0000
64	0.0100	-0.0000
65	0.0100	-0.0000
66	0.0099	-0.0000
67	0.0099	-0.0000
68	0.0098	-0.0000
69	0.0098	-0.0000
70	0.0097	-0.0000
71	0.0097	-0.0000
72	0.0096	-0.0000
73	0.0096	-0.0000
74	0.0095	-0.0000
75	0.0095	-0.0000
76	0.0094	-0.0000
77	0.0094	-0.0000
78	0.0093	-0.0000
79	0.0093	-0.0000
80	0.0092	-0.0000
81	0.0092	-0.0000
82	0.0091	-0.0000
83	0.0091	-0.0000
84	0.0090	-0.0000
85	0.0090	-0.0000
86	0.0089	-0.0000
87	0.0089	-0.0000
88	0.0088	-0.0000
89	0.0088	-0.0000
90	0.0087	-0.0000
91	0.0087	-0.0000
92	0.0086	-0.0000
93	0.0086	-0.0000
94	0.0085	-0.0000
95	0.0085	-0.0000
96	0.0084	-0.0000
97	0.0084	-0.0000
98	0.0083	-0.0000
99	0.0083	-0.0000
100	0.0082	-0.0000

Table 1.1: Entropy values for various distributions.

The quantity of information is the entropy of the associated probability distribution. Suppose we have a single event A with possible outcomes $\{a_i\}$:

- If one, a_0 , is certain to occur, $p(a_0) = 1$, then we acquire no information by observing A .
- If $A = a_0$ is very likely then we might have confidently expected it and so learn very little.
- If $A = a_0$ is highly unlikely then we might need to drastically change our plans.

Learning the value of A provides a quantity of information that increases as the corresponding probability decreases ($\sim -\log(\cdot)$).

We think of learning about something new as adding to the available information: the entropy is a weighted sum of the information we get from each event.

Letters

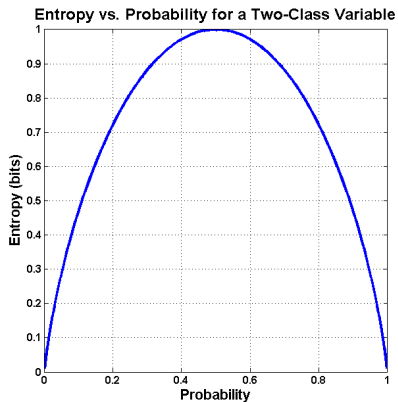
Consider a message formed from the alphabet A, B, C and D, with probabilities:

$$p(A) = \frac{1}{2}, \quad p(B) = \frac{1}{4}, \quad p(C) = \frac{1}{8} = p(D)$$

The information for the letter probabilities is

$$H(X) = -\left(\frac{1}{2} \log \frac{1}{2} + \frac{1}{4} \log \frac{1}{4} + 2 \times \frac{1}{8} \log \frac{1}{8}\right) = \frac{7}{4} \text{ bits}$$

What's the entropy of a uniform discrete random variable taking on 2 values?



The entropy $H(X) = -\sum_x p(x) \log p(x)$ has the following properties:

Properties

- $H(X) \geq 0$. Entropy is always **non-negative**. $H(X) = 0$ iff X is deterministic.
- **Upper-bound**: $H(X) \leq \log(|\mathcal{X}|)$. $H(X) = \log(|\mathcal{X}|)$ iff X has a uniform distribution over \mathcal{X} .
- Since $H_b(X) = \log_b(a) H_a(X)$, we don't need to specify the base of the logarithm ($\log_2 \rightarrow$ **bits**, $\log_e \rightarrow$ **nats**).

Probability theory revisited: joint probability

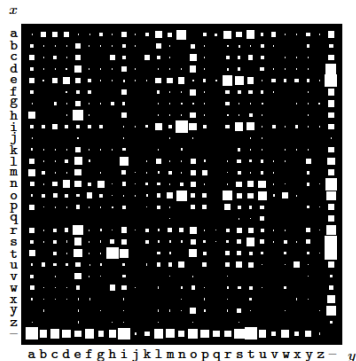


Figure 2.2. The probability distribution over the 27×27 possible bigrams xy in an English language document, *The Frequently Asked Questions Manual for Linux*.

Joint entropy

The joint entropy of a pair of two discrete random variables X and Y is:

$$H(X, Y) = - \sum_x \sum_y p(x, y) \log p(x, y)$$

Probability theory revisited: conditional probability

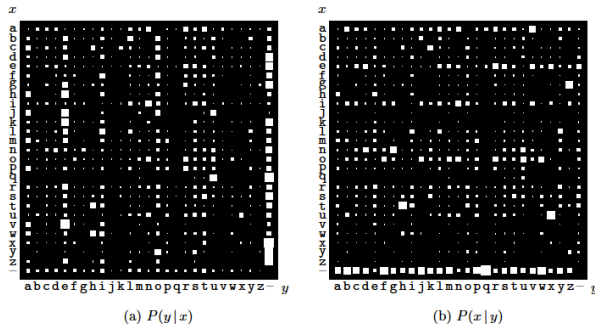


Figure 2.3. Conditional probability distributions. (a) $P(y|x)$: Each row shows the conditional distribution of the second letter, y , given the first letter, x , in a bigram xy . (b) $P(x|y)$: Each column shows the conditional distribution of the first letter, x , given the second letter, y .

Conditional entropy

$$\begin{aligned} H(Y|X) &= \sum_x p(x) H(Y|X=x) \\ &= - \sum_x p(x) \sum_y p(y|x) \log p(y|x) \\ &= - \sum_x \sum_y p(x, y) \log p(y|x) \\ &= H(X, Y) - H(X) \end{aligned}$$

Careful!

$$H(X|Y) \neq H(Y|X)$$

Mutual information

Mutual information measures how much information is in common between X and Y :

$$I(X; Y) = \sum_x \sum_y p(x, y) \log \frac{p(x, y)}{p(x)p(y)} = H(X) + H(Y) - H(X, Y)$$

$$I(X; Y) = I(Y; X)$$

└ Mutual information

Mutual information

Mutual information measures how much information is in common between X and Y :

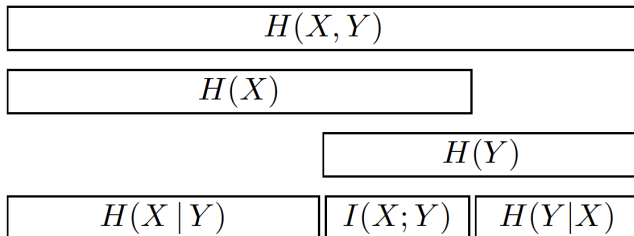
$$I(X; Y) = \sum_x \sum_y p(x, y) \log \frac{p(x, y)}{p(x)p(y)} = H(X) + H(Y) - H(X, Y)$$

$$I(X; Y) = I(Y; X)$$

$I(X; Y)$ is the mutual information between X and Y . It is the reduction of uncertainty of one Random Variable due to knowing about the other, or the amount of information one Random Variable contains about the other. Two interesting facts:

- I is 0 only when X, Y are independent: $H(X|Y) = H(X)$
- $H(X) = H(X) - H(X|X) = I(X, X)$. Entropy is the self-information.

Joint entropy, Conditional entropy and Mutual entropy



Kullback-Leibler (KL) divergence

The Kullback-Leibler divergence between two distributions on the same alphabet is

$$KL(p||q) = \sum p(x) \log \frac{p(x)}{q(x)}$$

- KL is a 'distance' measure between probability functions p and q .
- KL divergence is asymmetric (not a true distance):

$$KL(p||q) \neq KL(q||p)$$

- KL and Mutual information:

$$I(X; Y) = KL(p(x, y) || p(x)p(y))$$

Nonlinear Dimensionality Reduction

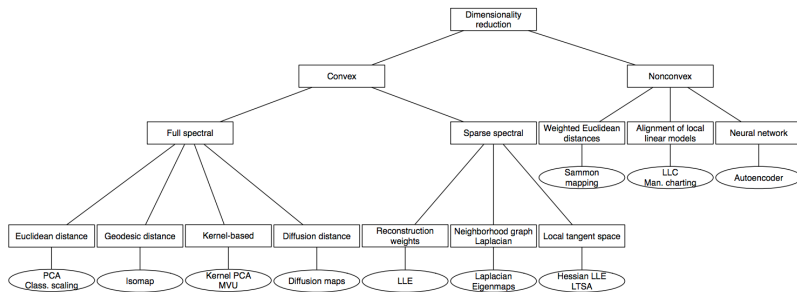


Figure: Taxonomy of dimensionality reduction techniques

t-Distributed stochastic neighbour embedding (t-SNE)

t-SNE minimises divergence of two distributions over pairwise similarities of

P input objects

Q low-dimensional points in the embedding

t-SNE algorithm

- Distance between a pair of objects, e.g., $d(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|$
- Joint probabilities p_{ij} that measure the pairwise similarity between objects \mathbf{x}_i and \mathbf{x}_j

$$p(j|i) = \frac{\exp(-d(\mathbf{x}_i, \mathbf{x}_j)^2/2\sigma_i^2)}{\sum_{k \neq i} \exp(-d(\mathbf{x}_i, \mathbf{x}_k)^2/2\sigma_i^2)}, \quad p(i|i) = 0$$

$$p_{ij} = \frac{p(j|i) + p(i|j)}{2N}$$

$$q_{ij} = \frac{(1 + d(\mathbf{y}_i - \mathbf{y}_j)^2)^{-1}}{\sum_{k \neq i} (1 + d(\mathbf{y}_i - \mathbf{y}_k)^2)^{-1}}$$

- Minimise cost function (KL-divergence)

$$KL(P||Q) = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

t-Distributed stochastic neighbour embedding (t-SNE)

t-Distributed stochastic neighbour embedding (t-SNE)

t-SNE minimises divergence of two distributions over pairwise similarities of

- P input objects
- Q low-dimensional points in the embedding

t-SNE algorithm

- Distance between a pair of objects, e.g., $d(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|$
- Joint probabilities p_{ij} that measure the pairwise similarity between objects \mathbf{x}_i and \mathbf{x}_j

$$p(i/j) = \frac{\exp(-d(\mathbf{x}_i, \mathbf{x}_j)^2/2\sigma_i^2)}{\sum_{k \neq i} \exp(-d(\mathbf{x}_i, \mathbf{x}_k)^2/2\sigma_i^2)}, \quad p(i/j) = 0$$

$$p_{ij} = \frac{p(i/j) + p(j/i)}{2N}$$

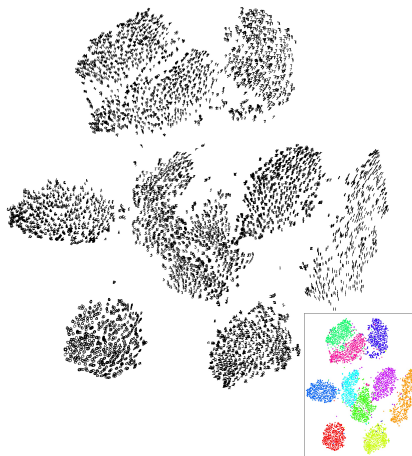
$$q_i = \frac{(1 + d(\mathbf{y}_i - \mathbf{y})^2)^{-1}}{\sum_{k \neq i} (1 + d(\mathbf{y}_i - \mathbf{y}_k)^2)^{-1}}$$

- Minimise cost function (KL-divergence)

$$KL(P||Q) = \sum_{i,j} \sum_{k,l} p_{ij} \log \frac{p_{ij}}{q_{kl}}$$

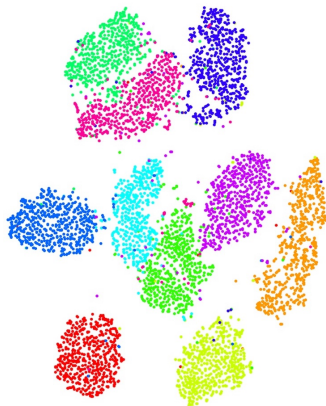
- + t-SNE compares favourably to other techniques for data visualisation
- unclear how t-SNE performs on general dimensionality reduction tasks
- relatively local nature of t-SNE makes it sensitive to the curse of the intrinsic dimensionality of the data
- not guaranteed to converge to a global optimum of its cost function

t-SNE on MNIST



http://scikit-learn.org/stable/auto_examples/manifold/plot_t_sne_perplexity.html

t-SNE on MNIST



http://scikit-learn.org/stable/auto_examples/manifold/plot_t_sne_perplexity.html

Summary

- **Entropy** is the measure of average *uncertainty* in the random variable.
- **Entropy** is the average number of *bits needed* to describe the random variable.
- **Entropy** is a lower bound on the average length of the *shortest description* of the random variable.
- **Conditional entropy** $H(X|Y)$ is the entropy of one random variable conditional *upon knowledge* of another.
- The average amount of decrease of the randomness of X by observing Y is the average information that Y gives us about X .
- **Mutual information** measures how much information is in *common* between X and Y .
- **KL divergence** measures the '*distance*' between probability distributions.

Summary

- **Entropy** is the measure of average uncertainty in the random variable.
- **Entropy** is the average number of bits needed to describe the random variable.
- **Entropy** is a lower bound on the average length of the shortest description of the random variable.
- **Conditional entropy** $H(X|Y)$ is the entropy of one random variable conditional upon knowledge of another.
- The average amount of decrease of the randomness of X by observing Y is the average information that Y gives us about X .
- **Mutual information** measures how much information is in common between X and Y .
- **KL divergence** measures the 'distance' between probability distributions.

Entropy is measure of uncertainty: the more we know about something, the lower the entropy. If a model captures more of the structure of the data, then the entropy should be lower. We can use entropy as a measure of the quality of our models. The KL divergence measures of how different two probability distributions are and can be interpreted as the average number of bits that are wasted by encoding events from a distribution p with a code based on a not-quite right distribution q :

Goal: minimise $D(p||q)$ to have a probabilistic model as accurate as possible

We will discuss decision trees and random forests.