

SCHOOL OF ELECTRONIC ENGINEERING AND COMPUTER SCIENCE
QUEEN MARY UNIVERSITY OF LONDON

ECS607/766 Data Mining

Week 5: Features and dimensionality

Dr Jesús Requena Carrión

2 Nov 2018

Agenda

Recap (with some extras)

Data normalisation

Dimensionality reduction

Appendix: Information and entropy

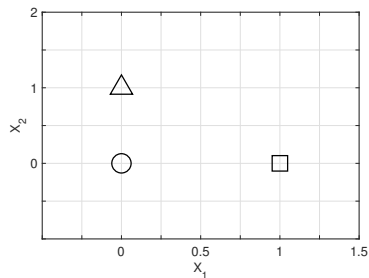
Distances in the predictor space

So far, we have used the notion of **distance** in various occasions:

- In **regression** problems, to define the prediction error $e_i = y_i - \hat{y}_i$ and the MSE cost function, $E_{MSE} = \frac{1}{N} \sum_{i=1}^N e_i^2$
- In **classification** problems we used the distance between samples and classifiers' boundaries, and the distance between samples in kNN
- In **clustering** problems, clusters were created based on the square distance between samples and cluster centres, $E_{KM} = \sum |x_i - \mu_i|^2$

The notion of distance is quite intuitive, but is it as straightforward as it seems?

Distances in the predictor space



Which sample is closer to \bigcirc : Is it \triangle or is it \square ?

(a) \triangle is closer

(b) \square is closer

(c) Both are equally distant

Sensitivity to predictors

In a linear regression model, the numerical value of a coefficient indicates how sensitive a prediction is to changes in the value of the corresponding predictor.

In the following linear regression model for a response y :

$$y = \mathbf{w}^T \mathbf{x} = 100x_A + 20x_B + 3$$

x_A and x_B are two predictors, $\mathbf{x} = [x_A, x_B, 1]$ is the extended predictor vector and $\mathbf{w} = [100, 20, 3]$ is the coefficient (or parameter) vector.

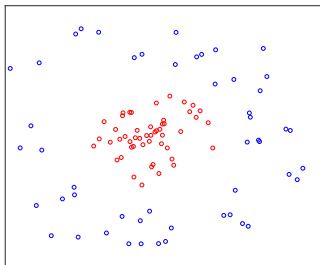
- (a) The response y is more sensitive to x_A than to x_B
- (b) The response y is more sensitive to x_B than to x_A
- (c) We have insufficient information to tell

Linear separability

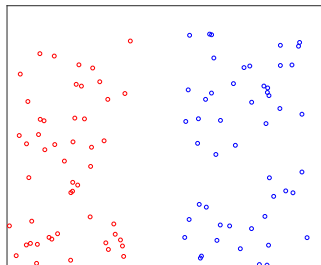
Which dataset is linearly separable, D_1 or D_2 ?

- (a) D_1 is linearly separable, D_2 isn't
- (b) D_2 is linearly separable, D_1 isn't
- (c) Both are linearly separable

D_1



D_2



Don't take your representation for granted!

Agenda

Recap (with some extras)

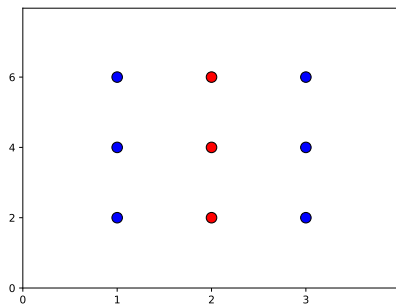
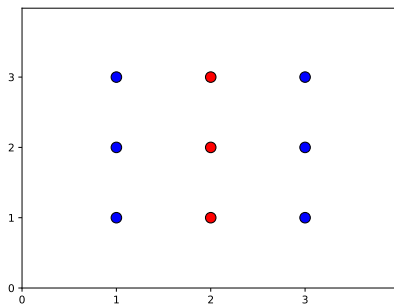
Data normalisation

Dimensionality reduction

Appendix: Information and entropy

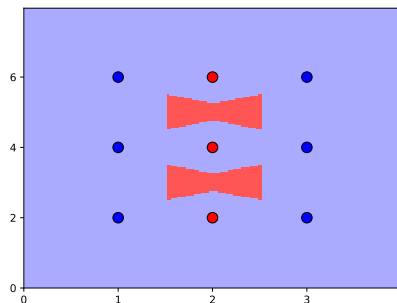
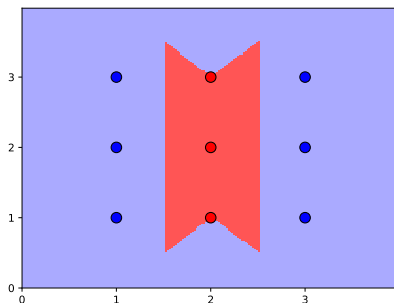
kNN and scaling

Two identical datasets (except for a scaling factor)



kNN and scaling

kNN solutions ($k = 3$) for each dataset



Numerical representation of attributes

Most of the time, we have ignored the meaning of the attributes under discussion, as our goal has been to discuss general methods that can be applied to any dataset. Hence, we have simply used the **numerical values** of the attributes without much thought. However:

- Attributes can be **incommensurable**, i.e. have different dimensions (for instance, *weight* and *height* cannot be compared)
- Even when attributes have the same dimensions, they might have **different dynamic ranges**
- Having large numerical values **doesn't translate to higher significance**
- Different numerical representations can have an impact on the **final model** and the **performance of our algorithms**

Numerical representation of attributes

The numerical representation of our attributes can be arbitrary and it is possible to find many **equivalent ways of representing numerically each attribute**. So which one is the **most convenient**?

Attributes whose numerical values vary within the same **numerical range** can offer a number of benefits, for instance if the predictors x_A and x_B in the linear model

$$y = \mathbf{w}^T \mathbf{x} = 100x_A + 20x_B + 3$$

take on values within the same numerical range, then it makes sense to say that the impact of x_A on the response y is higher than x_B .

Data normalisation (aka *feature scaling*) **techniques** allow us to obtain a **convenient numerical representation** of our attributes.

Min-max normalisation

This technique produces numerical values within the same range $[0, 1]$. In other words, the numerical value of an attribute will always be greater (or equal) than 0 and less (or equal) than 1.

Min-max normalisation produces a normalised attribute x' from the original attribute x by using the following transformation:

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$

where $\min(x)$ and $\max(x)$ are respectively the minimum and maximum value of x **in the available dataset**.

Notice that x' is a **dimensionless quantity**.

Standardisation

Standardisation is a common procedure in statistics. By applying the following transformation

$$x' = \frac{x - \mu}{\sigma}$$

where μ is the average of the values of x in the available dataset and σ is its standard deviation, the resulting attribute x' is such that the mean of its values in the available dataset is 0 and the standard deviation is 1.

Once again, x' is a **dimensionless quantity**.

Final notes

- Datasets contain samples of a population. During deployment we should expect **out-of-range** values (e.g. $x' = 1.2$ in min-max)
- The effect of **outliers** need to be considered (for instance, an outlier 100 times larger than the second largest value would squeeze the remaining min-max values within the interval $[0, 0.01]$)
- In addition to linear transformations, other non-linear methods exist, for instance **softmax scaling**, which uses the logistic function
- The **distribution of an attribute** can also be normalised
- In general, we need to understand well the **effects and distortions of any data transformation**

Agenda


Recap (with some extras)


Data normalisation

Dimensionality reduction

Appendix: Information and entropy

The Bosch Production Line Dataset

kaggle Search kaggle Q Competitions Datasets Kernels Discussion Learn ... 



Bosch Production Line Performance


Reduce manufacturing failures
\$30,000 · 1,373 teams · 2 years ago


[Overview](#) [Data](#) [Kernels](#) [Discussion](#) [Leaderboard](#) [Rules](#) [Late Submission](#)

Overview

[Description](#)
[Evaluation](#)
[Prizes](#)
[Timeline](#)
lee Bigdata 2016

A good chocolate soufflé is decadent, delicious, and delicate. But, it's a challenge to prepare. When you pull a disappointingly deflated dessert out of the oven, you instinctively retrace your steps to identify at what point you went wrong. [Bosch](#), one of the world's leading manufacturing companies, has an imperative to ensure that the recipes for the production of its advanced mechanical components are of the highest quality and safety standards. Part of doing so is closely monitoring its parts as they progress through the manufacturing processes.





17/55

The MNIST dataset



Data dimensionality

In Data Science, an attribute can be seen as a **dimension of our datasets**. This interpretation allows us to **represent instances as points in the space**, by using the values of each of its attributes as coordinates.

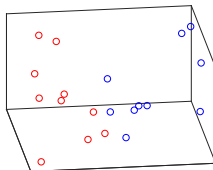
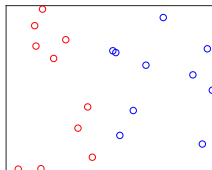
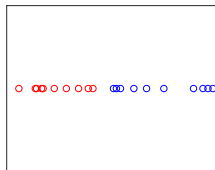
Some datasets may include many attributes (for instance, the Bosch dataset contains 970), and are therefore said to be **high dimensional**. Frequently, this is due to us having **little prior knowledge**, which forces us to record everything (just in case!).

The question arises, what are the **main challenges of high dimensional datasets**? Is it wise to include as many attributes as we can?

The Curse of Dimensionality

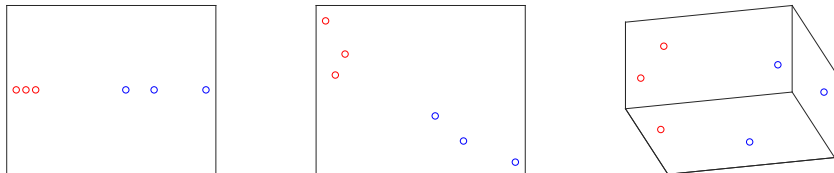
As we increase the dimensionality of our dataset, **data becomes sparser**.

In this example, we add two irrelevant attributes to 20 samples that are initially described by one single attribute. How would a logistic regression boundary change as we include new irrelevant attributes? And a kNN one?



The Curse of Dimensionality

Now we only use 6 samples. Compare their boundaries: It is much clearer that adding new irrelevant attributes can make things worse!



In a high dimensional settings, we need to learn more parameters than in a low dimensional ones, so the **risk of overfitting** increases. This risk is specially dangerous if we have many attributes that are weakly relevant, or some very relevant and many irrelevant.

Dimensionality reduction

High dimensional datasets present many challenges, including:

- Overfitting (curse of dimensionality)
- Irrelevant data
- Computational cost
- Storage cost
- Hard to visualise
- Difficult interpretation

Dimensionality reduction is a family of techniques whose goal is to **transform a high dimensional dataset into a more convenient low dimensional dataset**. Two main approaches are:

- Feature selection → *pick* the best predictors
- Feature extraction → *transform* onto a smaller set of predictors

Agenda

Recap (with some extras)

Data normalisation

Dimensionality reduction

- Feature selection

- Feature extraction

Appendix: Information and entropy

Feature selection

In feature selection, we start by wondering whether not all our predictors (a.k.a *features*) in our data set are relevant. Therefore we want to be able to select the **best** ones, in other words, we want to identify the **best subset of predictors**.

This leads to two observations:

- What do we mean by **best**? We will use the **response** to create metrics for subset selection (**supervised problem!**).
- If our dataset has M predictors, there are a **total of $2^M - 1$ subsets** that we could consider.

If our dataset has 10 features, we have roughly 1000 options. In the Bosch dataset we have a few more than 10^{270} options. How do we find the best subset?

Filtering

The simplest approach towards feature selection is to consider each one of the features **individually** and assign them a **score**, which we use to **rank** them and **select** the N best ones.

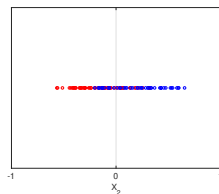
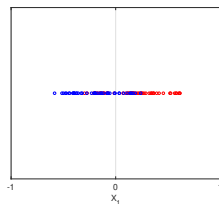
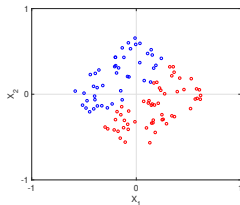
Scores essentially **compare each feature with the desired response**. Common scores include:

- Correlation.
- Mutual information.
- Statistical independence.

Filtering

Filtering is a **simple** and **fast** method. However, its starting point is that the features in the dataset contribute separately to the final response, therefore possible **interactions among predictors are ignored**.

In the following example, predictors X_1 and X_2 do a poor job separately, but together they reveal a clear boundary between the two classes.



Wrapping

If we suspect that the interaction between predictors might be crucial, we have no choice but to **evaluate them together**, rather than separately.

Wrapping approaches consider possible interaction between predictors by:

- Training a model with **different subsets of predictors**
- Evaluating each resulting model by using **validation approaches**.
- Picking the subset with the **highest validation performance**.

Whereas filtering approaches retrain a model M times (where M is the number of predictors), wrapping models can potentially consider up to $2^M - 1$ predictors! The **computational cost** is, therefore, a big concern.

Wrapping: Greedy search

Given the large number of subsets of predictors that we might need to consider, the main challenge when implementing wrapping approaches is how to **search for the best subset**.

In general, considering all the candidate subsets (known as **brute-force** or **exhaustive** search) will be impractical. Greedy search is a strategy for exploring candidate subsets. It comes in two flavours:

- **Forward selection:** We start with a subset containing the best predictor and progressively add the predictor that improves the subset's performance the most.
- **Backward selection:** We start with a subset containing all the predictors and progressively remove the predictor whose elimination improves the performance the most.

In both cases we implement a stop criterion (typically, when the performance stops increasing).

Embedded selection

Some learning algorithms can have some form of feature selection built in the process of learning.

- **Classification trees** can be grow following a strategy that effectively eliminates irrelevant features.
- Properly designed **regularisation strategies** can be seen as a feature selection process, as by attenuating the coefficient w_A corresponding to a predictor x_A , this predictor is eliminated. We have used the so-called L_2 regularisation by using the term $\lambda \mathbf{w}^T \mathbf{w}$, but other options are available too.

Agenda

Recap (with some extras)

Data normalisation

Dimensionality reduction

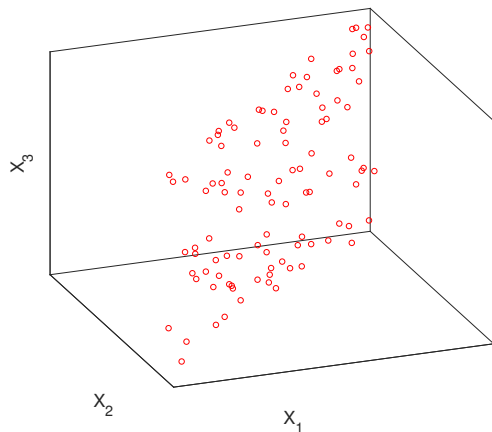
Feature selection

Feature extraction

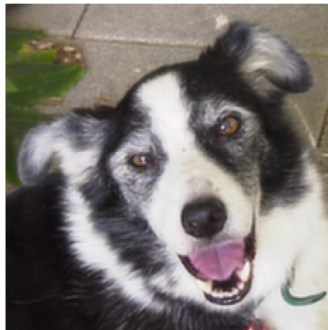
Appendix: Information and entropy

Interesting directions

A dataset with predictors X_1 , X_2 and X_3 is represented below. You can move freely in this 3D space. Which way would you go?



What is feature extraction?



This picture consists of
 $422 \times 424 \approx 180,000$ pixels.

- Do we need 180,000 predictors to tell it's a dog?
- Would a subset of pixels work?
- Can we transform the picture into a new set of predictors?

What is feature extraction?

Feature extraction allows us to reduce the dimensionality of our dataset by creating a **new set of predictors** of lower dimensionality. The new predictors are **non-redundant** and overall should **as much information from the original set** as possible.

There exist many **data, signal and image processing techniques** that allow us to define new features that can be extracted from high-dimensional data, for instance:

- Frequency components of signals (Fourier analysis)
- Texture characterisation of images (wavelet analysis)

Principal Components Analysis is one of the most popular and well understood techniques for dimensionality reduction via feature extraction.

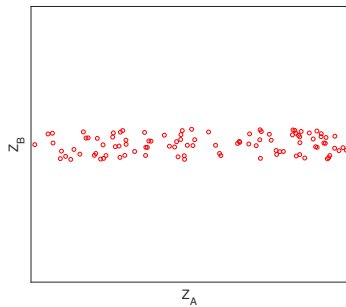
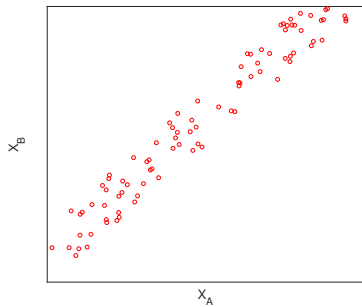
Principal Component Analysis

Given a set of N features x_A, x_B, \dots , Principal Component Analysis (PCA) produces:

- Another set of N **orthogonal** new features z_A, z_B, \dots (in other words, they are non-redundant)
- That can be **combined** to produce the original features (therefore both set of features have the same information)
- Along with a **score** that can be used to rank the new features and select them

PCA doesn't use the desired response to create this transformation, and therefore it is an **unsupervised approach**. How does it work then?

PCA as a projection



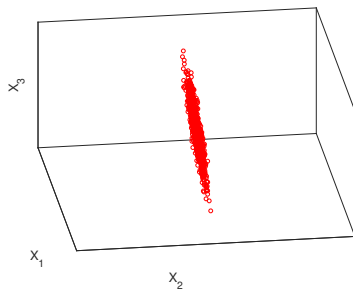
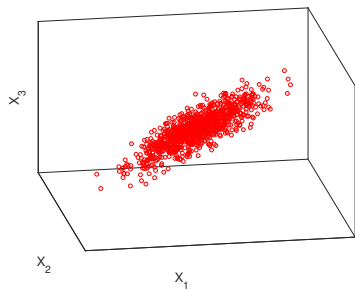
PCA principles

Our starting point is this: **directions along which data varies a lot, contain a lot of information**, and conversely, directions along which data varies little, contain little information.

Our PCA algorithm proceeds as follows:

- We find the direction along which **data varies the most**
- This direction is our first **new feature**
- By subtracting the new feature from our original data, we obtain a residual containing **unexplained** variance
- We extract a new feature and residual from any residual previously obtained in the previous step until we have **exhausted all the dimensions**

PCA principles



Final notes on PCA

PCA is a **linear transformation** of our dataset: the new features z can be obtained by multiplying the original features x by a combination matrix U , $z = Ux$. PCA is widely implemented and computationally fast. However:

- It is **not scale-invariant**
- It assumes **variations are gaussian**, and therefore might ignore non-gaussian patterns.
- **Non-linear scenarios** are not accounted for
- Information might be in **low-variance components**

Agenda

Recap (with some extras)

Data normalisation

Dimensionality reduction

Appendix: Information and entropy

The weather report for Cairo



What would you say about the weather reports for Cairo?
Interesting? Boring?

The weather report for Cairo



13 Dec 2013 every newspaper in the world was talking about the weather in Cairo. Why?

Because it was totally unexpected, in other words **highly improbable**.

Take-home message: there is a connection between **information** and **probability**!

The notion of information

We all have an intuitive notion of the meaning of information. In general, we will all agree that an event that we expect gives us little information, whereas a surprising event has a lot of information. Nevertheless, in order for us to analyse information sources and communication systems, we need to be able to **quantify** it.

In this section we will define a **quantitative measure of information**.

- Because of the connection between information and probability, this quantitative measure of information is defined in **statistical terms**.
- Since we are interested in digital information, we will use **discrete random variables**.

Probability reminder: Discrete random variables

- A discrete random variable X can take any of a specified countable set of values x_1, x_2, \dots
- The discrete random variable X is defined by the probability that it takes each value x_i .
- The probability mass function (PMF) is a function that gives the probability that X takes each value x_i and we denote it by $p_i = P(X = x_i)$.

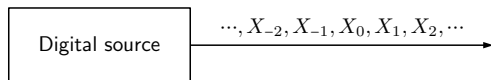
Probability reminder: Discrete random variables

Let X and Y be two discrete random variables.

- The probability that $X = x_i$ **and** $Y = y_j$ is the joint probability of X and Y and we denote it by $P(X = x_i, Y = y_j)$.
- The probability that $X = x_i$ **given that** $Y = y_j$ is the conditional probability and is denoted by $P(X = x_i | Y = y_j)$.
- The joint probability and the conditional probability are related by the equality $P(X = x_i, Y = y_j) = P(X = x_i | Y = y_j)P(Y = y_j)$.
- We say that X and Y are **independent** if and only if $P(X = x_i, Y = y_j) = P(X = x_i)P(Y = y_j)$. This implies that $P(X | Y) = P(X)$.

Modelling digital information sources

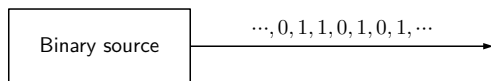
An information source is an entity that produces information. We will model a digital information source as discrete random variable X that generates **sequences** of **discrete values** X_n .



Each discrete value that a digital information source can generate is called **symbol** and we denote them by a_k . The collection of all the symbols is called the source alphabet, $A = \{a_1, a_2, \dots, a_K\}$.

The binary information source

Binary information sources are digital information sources that have an alphabet consisting of two symbols, usually denoted by $\{0, 1\}$.



Q: How many different binary sequences of length N exist?

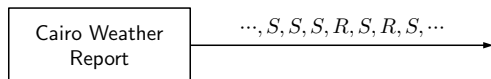
Memoryless sources

Memoryless sources are a special family of digital information sources that have the property that **they generate symbols independently**. In other words, the probability of observing a symbol at one time instant is independent of the observation of any other symbol at any other time instant.

Q: Given a memoryless source defined by an alphabet $A = \{a_1, a_2, a_3\}$ and a PMF p_i , what is the probability of generating the sequence $S_1 = a_1, a_2, a_3, a_1, a_2, a_1$? And the sequence $S_2 = a_1, a_1, a_1, a_2, a_2, a_3$?

Example: Cairo weather report

The *Cairo weather report* is a digital information source that has an alphabet consisting of three symbols: *sunny*, *rainy*, *snowy*. Let us denote its alphabet by $A = \{S, R, W\}$



Q: What is the probability that we observe 10 sunny dais in a row?
(*Note: Assume that daily reports are statistically independent. i.e. assume that the Cairo weather report is a memoryless source*)

Measure of Information

Let X be a digital information source that generates symbols a_k from an alphabet A with PMF p_k . The amount of information revealed by a symbol a_k needs to have the following properties:

- It must depend on the probability of generating such symbol, p_k .
- If p_k is low, its information content must be high.
- If A and B are two independent information sources, the joint information of two symbols a_k and b_l is the sum of their respective information contents.

Measure of Information

It can be proved that the only function that satisfies these properties is the **logarithmic function**. We define the information of a symbol a_i as

$$I(a_k) = -\log(p_k) = \log\left(\frac{1}{p_k}\right)$$

The base of the logarithm is not important. If it is base 2, the information is expressed in *bits*.

Specifically:

- If $p_k \rightarrow 0$, $I(a_k) \rightarrow \infty$.
- If $p_k \rightarrow 1$, $I(a_k) \rightarrow 0$.

Information content of the Cairo weather report

Imagine we are interested in knowing whether it is sunny (S), it is rainy (R) or snowy (W) in Cairo. Each day of the year we get a report informing us about this. How much information do we get from each *individual* report?

Let us look at the probabilities of each event:

$$P(S) \simeq 0.96$$

$$P(R) \simeq 0.04$$

$$P(W) \simeq 0.00002$$

(In order to derive these probabilities, I have used average data from the website www.weatherspark.com and considered the fact that before 2013, 112 years had elapsed since the last time it snowed in Cairo)

Information content of Cairo weather reports

By using the above probabilities and the definition of information content, we obtain the following amount of information for a daily weather report:

$$I(S) \simeq 0.06 \text{ bits}$$

$$I(R) \simeq 4.6 \text{ bits}$$

$$I(W) \simeq 16 \text{ bits}$$

We see clearly that we get more information from a *snowy* report than from a *sunny* one. Another question is, over the year, do we get on average more information from the Cairo weather report or from the Beijing weather report? This leads us to the notion of **source entropy**.

The notion of entropy

We already know how much information we get when we observe a symbol a_i . However, how much information do we get on average from an information source? We call this quantity the **entropy** H of the source and we calculate it as the weighted-average of the information of each symbol

$$H = \sum_{k=1}^K p_k I(a_k) = - \sum_{k=1}^K p_k \log(p_k)$$

The entropy H is expressed in bits/symbol. If we know the symbol transmission rate of the information source (symbols/s), we can also express the entropy in bits/s (bps).

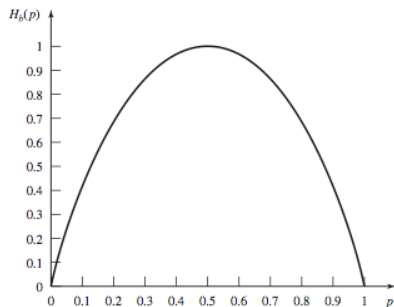
Example: Entropy of the Cairo, Beijing and London weather reports

The entropy of the Cairo weather report is

$$\begin{aligned}H_{CWR} &= P(S)I(S) + P(R)I(R) + P(W)I(W) \\&= 0.96 \times 0.06 + 0.04 \times 4.6 + 0.00002 \times 16 \\&= 0.24 \text{ bits/symbol}\end{aligned}$$

Q: In Beijing, $P(S) \simeq 0.8$, $P(R) \simeq 0.2$ and $P(W) \simeq 0.02$. What is the entropy H_{BWR} of the Beijing weather report? In London, $P(S) \simeq 0.4$, $P(R) \simeq 0.6$ and $P(W) \simeq 0.02$. What is the entropy H_{LWR} of the London weather report?

The binary memoryless source



The binary memoryless source is a binary information source that generates two symbols 0 and 1 with probabilities $p_0 = p$ and $p_1 = 1 - p$.

Its entropy $H_b(p)$ depends on the value p :

$$H_b(p) = -p \log(p) - (1-p) \log(1-p)$$

When $p = 0.5$ the source is called **binary symmetric source** and it has the **maximum entropy**.