

**Dr. Jürg M. Stettbacher**

Margrit Rainer Strasse 12a  
CH-8050 Zürich

Telefon: +41 43 299 57 23

Fax: +41 43 299 57 25

E-Mail: [dsp@stettbacher.ch](mailto:dsp@stettbacher.ch)

# Information und Entropie

## Praktikum

Version 2.00  
2013-06-12

Zusammenfassung: In diesem Praktikum geht es darum, ein Programm zu schreiben, das von gegebenen Daten in einer Datei die Information und die Entropie bestimmt.

# Inhaltsverzeichnis

<b>1</b>	<b>Einleitung</b>	<b>2</b>
<b>2</b>	<b>Dateien</b>	<b>2</b>
<b>3</b>	<b>Aufgabe</b>	<b>3</b>
<b>4</b>	<b>Zusatzaufgabe</b>	<b>3</b>

## 1 Einleitung

In der Informationstheorie betrachten wir Datenquellen oft als Zufallsvariablen und jedes Symbol, das aus der Quelle kommt, als Zufallseignis. Sind die Auftretenswahrscheinlichkeiten bekannt, so kann für jedes Symbol die Information berechnet werden und für die Quelle der mittlere Informationsgehalt, also die Entropie.

In diesem Praktikum haben wir verschiedene Datenquellen in der Form von Dateien zur Verfügung. Es handelt sich um ASCII-Dateien (\*.txt) und jedes ASCII-Zeichen daraus stellt ein Symbol dar. Gesucht sind jeweils der Informationsgehalt für jedes Symbol einer Datei, sowie die Entropie der gesamten Datei.

## 2 Dateien

Die folgenden Dateien stehen Ihnen für das Bearbeiten des Praktikums zur Verfügung:

- *entropy\_template.php* (Hauptaufgabe)
- *entropy2\_template.php* (Zusatzaufgabe)
- *data\_1.txt* bis *data\_6.txt* und *deutsch.txt* (Testdaten)

## 3 Aufgabe

Verwenden Sie die Testdaten in den Dateien *data\_1.txt* bis *data\_6.txt*. Schreiben Sie ein PHP-Skript *entropy.php*, das eine bestimmte Datei öffnen und zeichenweise lesen kann. Der Name der ASCII-Datei soll auf der Kommandozeile übergeben werden, so dass der Aufruf folgendermassen aussieht.

```
> php entropy.php data.txt
```

Es steht eine vorbereitete Datei *entropy\_template.php* zur Verfügung, die Sie als Vorlage verwenden können. Sie brauchen dann nur an den bezeichneten Stellen Ihren Code zu ergänzen. Für die Berechnung von Information und Entropie gehen Sie so vor:

- Erstellen Sie ein Histogramm in dem Sie für alle Zeichen angeben, wie häufig sie vorgekommen sind und zählen Sie, wieviele Zeichen insgesamt in der Datei enthalten sind.
- Berechnen Sie für jedes Zeichen die Information.
- Berechnen Sie die Entropie der Quelle.
- Auf dem Bildschirm soll für jedes Zeichen ausgegeben werden
  - wie oft es vorgekommen ist,
  - was seine Information ist.
- Zudem soll auf dem Bildschirm für die Quelle ausgegeben werden
  - wie viele Zeichen total vorgekommen sind,
  - wie gross die Entropie ist.

## 4 Zusatzaufgabe

Verwenden Sie diesmal die Testdaten in der Datei *deutsch.txt*. Wir wollen nun prüfen, was geschieht, wenn die Symbole einer Quelle nicht statistisch unabhängig sind. Schreiben Sie zu diesem Zweck ein PHP-Skript *entropy2.php*, das eine ASCII Datei öffnen und zeichenweise lesen kann. Es steht dafür die Vorlage *entropy2\_template.php* zur Verfügung. Für ein gegebenes Zeichen  $y_0$  soll das Skript die bedingte Entropie  $H(X|y_0)$  ermitteln, die Entropie jenes Symbols, das auf das Symbol  $y_0$  folgt. Wir nennen dies die *Entropie von X gegeben  $y_0$* .

$$H(X|y_0) = \sum_{x \in \Omega} P(x|y_0) \cdot \log_2 \frac{1}{P(x|y_0)} \quad (1)$$

Dabei ist  $\Omega$  der Ereignisraum der Zufallsvariable  $X$ , in unserem Fall also der Zeichenvorrat der Quelle  $X$ . Der Aufruf des PHP-Skripts soll so aussehen:

```
> php entropy2.php data.txt y0
```

Die Entropie  $H(X|Y)$  ist übrigens der Mittelwert von  $H(X|y)$  über allen Symbolen  $y$  der Quelle  $Y$ . In unserem Fall also:

$$H(X|Y) = \sum_{y \in \Omega} P(y) \cdot H(X|y) = \sum_{x,y \in \Omega} P(y) \cdot P(x|y) \cdot \log_2 \frac{1}{P(x|y)} = \sum_{x,y \in \Omega} P(x,y) \cdot \log_2 \frac{1}{P(x|y)} \quad (2)$$

Dies wollen wir an dieser Stelle aber nicht weiter verfolgen.

Die Resultate des Skripts *entropy2.php* sollen mit den Resulten der ersten Aufgabe verglichen werden. Überlegen Sie sich dabei die folgenden Punkte:

- Bei welchen Daten sind Unterschiede zu erwarten?
- Bei welchem Zeichen  $y_0$  treten die Unterschiede besonders deutlich hervor?