

A1.1: Wetterentropie

Eine Wetterstation fragt täglich verschiedene Regionen ab und bekommt als Antwort jeweils eine Meldung x zurück, nämlich

$x = \mathbf{B}$: Das Wetter ist eher schlecht.

$x = \mathbf{G}$: Das Wetter ist eher gut.

Die Daten wurden über viele Jahre für verschiedene Gebiete in Dateien abgelegt, so dass die Entropien der B/G-Folgen ermittelt werden können:

$$H = p_B \cdot \log_2 \frac{1}{p_B} + p_G \cdot \log_2 \frac{1}{p_G}$$

mit dem *Logarithmus dualis*

$$\log_2 p = \frac{\lg p}{\lg 2} \quad (= \text{ld } p) .$$

Region/Datei „Durchwachsen“:

GGBBGGBBGGGGGGBBGGBBGGGGGGGG
 GBGGGGGGGGGGGGGGGGGGGGGGGGGGGG

Region/Datei „Regenloch“:

BBGGGGGGGGGGGGGGGGGGGGGGGGGGGG
 BBBGGGGGGGGGGGGGGGGGGGGGGGGGGGG

Region/Datei „Angenehm“:

GGGGGGGGGGGGGGGGGGGGGGGGGGGGGG
 GBGGGGGGGGGGGGGGGGGGGGGGGGGGGG

Region/Datei „Paradies“:

GGGGGGGGGGGGGGGGGGGGGGGGGGGGGG
 GGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG

Region/Datei „Unbekannt“:

BBBBBBBBBGGGGGGGGGGGGGGGGGGGGGG
 GGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG

© 2011 www.LNTwww.de

„lg“ kennzeichnet hierbei den Logarithmus zur Basis 10. Zu erwähnen ist ferner, dass jeweils noch die Pseudoeinheit „bit/Anfrage“ anzufügen ist.

Die Grafik zeigt diese binären Folgen jeweils für 60 Tage und folgende Regionen:

- Region „Durchwachsen“: $p_B = p_G = 0.5$,
- Region „Regenloch“: $p_B = 0.8, p_G = 0.2$,
- Region „Angenehm“: $p_B = 0.2, p_G = 0.8$,
- Region „Paradies“: $p_B = 1/30, p_G = 29/30$.

Schließlich ist auch noch die Datei „Unbekannt“ angegeben, deren statistische Eigenschaften zu schätzen sind.

Hinweis: Die Aufgabe bezieht sich auf das **Kapitel 1.1**. Für die vier ersten Dateien wird vorausgesetzt, dass die Ereignisse „B“ und „G“ statistisch unabhängig seien, eine für die Wetterpraxis allerdings eher unrealistische Annahme.

Fragebogen zu "A1.1: Wetterentropie"

a) Welche Entropie H_D weist die Datei „Durchwachsen“ auf?

$$H_D = \text{bit/Anfrage}$$

b) Welche Entropie H_R weist die Datei „Regenloch“ auf?

$$H_R = \text{bit/Anfrage}$$

c) Welche Entropie H_A weist die Datei „Angenehm“ auf?

$$H_A = \text{bit/Anfrage}$$

d) Wie groß sind die Informationsgehalte der Ereignisse „B“ und „G“ bezogen auf die Datei „Paradies“?

$$I_B = \text{bit/Anfrage}$$

$$I_G = \text{bit/Anfrage}$$

e) Wie groß ist die Entropie (das heißt: der mittlere Informationsgehalt) H_P der Datei „Paradies“? Interpretieren Sie das Ergebnis?

$$H_P = \text{bit/Anfrage}$$

f) Welche Aussagen könnten für die Datei „Unbekannt“ gelten?

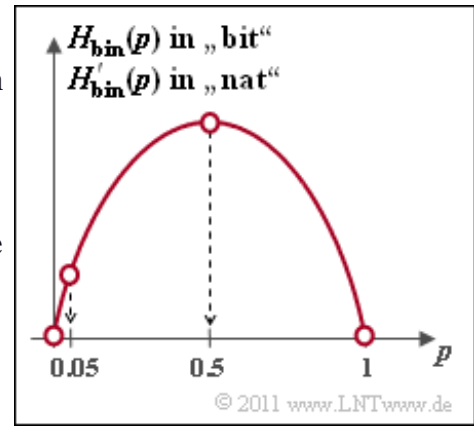
- Die Ereignisse „B“ und „G“ sind etwa gleichwahrscheinlich.
- Die Folgeelemente sind statistisch voneinander unabhängig.
- Die Entropie dieser Datei ist $H_U \approx 0.7$ bit/Anfrage.
- Die Entropie dieser Datei ist $H_U = 1.5$ bit/Abfrage.

Z1.1: Binäre Entropiefunktion

Wir betrachten eine Folge von binären Zufallsgrößen mit dem Symbolvorrat $\{A, B\} \Rightarrow M = 2$. Die Auftretenswahrscheinlichkeiten der beiden Symbole seien $p_A = p$ und $p_B = 1 - p$.

Die einzelnen Folgeelemente sind statistisch unabhängig. Für die Entropie dieser Nachrichtenquelle gilt gleichermaßen:

$$H_{\text{bin}}(p) = p \cdot \text{ld} \frac{1}{p} + (1 - p) \cdot \text{ld} \frac{1}{1 - p} \text{ in [bit]},$$
$$H'_{\text{bin}}(p) = p \cdot \ln \frac{1}{p} + (1 - p) \cdot \ln \frac{1}{1 - p} \text{ in [nat]}.$$



In diesen Gleichungen werden als Kurzbezeichnungen verwendet:

- der *natürliche* Logarithmus $\ln p = \log_e p$,
- der Logarithmus *dualis* $\text{ld} p = \log_2 p$.

Die Grafik zeigt diese binäre Entropiefunktion in Abhängigkeit des Parameters p , wobei $0 \leq p \leq 1$ vorausgesetzt wird.

In den Teilaufgaben (e) und (f) soll der relative Fehler ermittelt werden, wenn die Symbolwahrscheinlichkeit p per Simulation (also als relative Häufigkeit h) ermittelt wurde und sich dabei fälschlicherweise $h = 0.9 p$ ergeben hat. Der relative Fehler ist dann wie folgt gegeben:

$$\varepsilon_H = \frac{H_{\text{bin}}(h) - H_{\text{bin}}(p)}{H_{\text{bin}}(p)}.$$

Hinweis: Die Aufgabe gehört zum **Kapitel 1.1**.

Fragebogen zu "Z1.1: Binäre Entropiefunktion"

a) Wie hängen $H_{\text{bin}}(p)$ in bit und $H'_{\text{bin}}(p)$ in nat zusammen?

- $H_{\text{bin}}(p)$ und $H'_{\text{bin}}(p)$ unterscheiden sich um einen Faktor.
- Es gilt $H'_{\text{bin}}(p) = H_{\text{bin}}(\ln p)$.
- Es gilt $H'_{\text{bin}}(p) = 1 + H_{\text{bin}}(2p)$.

b) Zeigen Sie, dass sich das Maximum der binären Entropiefunktion für $p = 0.5$ ergibt. Wie groß ist $H_{\text{bin}}(p = 0.5)$?

$$H_{\text{bin}}(p = 0.5) = \quad \text{bit}$$

c) Berechnen Sie den binären Entropiewert für $p = 0.05$.

$$H_{\text{bin}}(p = 0.05) = \quad \text{bit}$$

d) Geben Sie den größeren der beiden p -Werte ein, die sich aus der Gleichung $H_{\text{bin}}(p) = 0.5$ bit ergeben.

$$p =$$

e) Durch unzureichende Simulation wurde $p = 0.5$ um 10% zu niedrig ermittelt. Wie groß ist der prozentuale Fehler hinsichtlich der Entropie?

$$p = 0.45 \text{ statt } p = 0.5: \quad \epsilon_H = \quad \%$$

f) Durch unzureichende Simulation wurde $p = 0.05$ um 10% zu niedrig ermittelt. Wie groß ist der prozentuale Fehler hinsichtlich der Entropie?

$$p = 0.045 \text{ statt } p = 0.05: \quad \epsilon_H = \quad \%$$

A1.2: Entropie von Ternärquellen

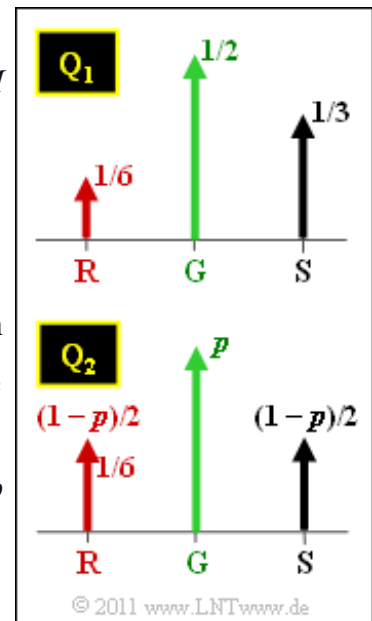
Die Entropie einer wertdiskreten gedächtnislosen Nachrichtenquelle mit M möglichen Symbolen lautet:

$$H = \sum_{\mu=1}^M p_{\mu} \cdot \log_2 \frac{1}{p_{\mu}}, \text{ Pseudoeinheit: bit.}$$

Hierbei bezeichnen die p_{μ} die Auftretswahrscheinlichkeiten der einzelnen Symbole bzw. Ereignisse. Im vorliegenden Beispiel werden die Ereignisse mit **R**(ot), **G**(rün) und **S**(chwarz) bezeichnet.

Bei einer binären Quelle mit den Auftretswahrscheinlichkeiten p und $1 - p$ kann hierfür geschrieben werden:

$$H = H_{\text{bin}}(p) = p \cdot \log_2 \frac{1}{p} + (1 - p) \cdot \log_2 \frac{1}{1 - p}.$$



Die Entropie einer mehrstufigen Quelle lässt sich häufig mit dieser „binären Entropiefunktion“ $H_{\text{bin}}(p)$ – ebenfalls mit der Pseudoeinheit „bit“ – ausdrücken.

Betrachtet werden in dieser Aufgabe zwei Ternärquellen mit den Symbolwahrscheinlichkeiten gemäß der obigen Grafik:

- die Quelle Q_1 mit $p_G = 1/2$, $p_S = 1/3$, $p_R = 1/6$,
- die Quelle Q_2 mit $p_G = p$, $p_R = p_S = (1 - p)/2$.

Die Ternärquelle Q_2 lässt sich auch auf Roulette anwenden, wenn ein Spieler nur auf die Felder **R**ot, **S**chwarz und **G**rün (die „Null“) setzt. Dieser Spieltyp wird im Fragebogen mit „Roulette 1“ bezeichnet.

Dagegen weist „Roulette 2“ darauf hin, dass der Spieler auf einzelne Zahlen (**0**, ... , **36**) setzt.

Hinweis: Die Aufgabe bezieht sich auf das **Kapitel 1.1**.

Fragebogen zu "A1.2: Entropie von Ternärquellen"

a) Welche Entropie H besitzt die Quelle Q_1 ?

$$Q_1: H = \quad \text{bit}$$

b) Welche der folgenden Aussagen sind zutreffend, wenn man R , G und S durch die Zahlenwerte -1 , 0 und $+1$ darstellt?

- Es ergibt sich eine kleinere Entropie.
- Die Entropie bleibt gleich.
- Es ergibt sich eine größere Entropie.

c) Bestimmen Sie die Entropie der Quelle Q_2 unter Verwendung der binären Entropiefunktion $H_{\text{bin}}(p)$. Welcher Wert ergibt sich für $p = 0.5$?

$$Q_2; p = 0.5: H = \quad \text{bit}$$

d) Für welchen p -Wert ergibt sich die maximale Entropie?

$$Q_2, H \rightarrow H_{\text{max}}: p =$$

e) Welche Entropie hat die Nachrichtenquelle „Roulette“ hinsichtlich der Ereignisse Rot, Schwarz und Grün (die „Null“)?

$$\text{Roulette 1: } H = \quad \text{bit}$$

f) Welche Entropie weist „Roulette“ hinsichtlich der Zahlen $0, \dots, 36$ auf?

$$\text{Roulette 2: } H = \quad \text{bit}$$

A1.3: H_0, H_1, H_2, \dots, H

Die Grafik zeigt vier Symbolfolgen $\langle q_v \rangle$ mit jeweiliger Länge $N = 60$. Die Quellensymbole sind jeweils A und B. Daraus folgt direkt, dass für den Entscheidungsgehalt aller betrachteten Quellen $H_0 = 1$ bit/Symbol gilt. Die Symbole A und B treten mit den Wahrscheinlichkeiten p_A und p_B auf.

Die folgende Tabelle zeigt neben H_0 die Entropienäherungen

- H_1 , basierend auf p_A und p_B (Spalte 2),
- H_2 , basierend auf Zweiertupel (Spalte 3),
- H_3 , basierend auf Dreiertupel (Spalte 4),
- H_4 , basierend auf Vierertupel (Spalte 5),
- die tatsächliche Quellenentropie H , die sich aus H_k durch den Grenzübergang für $k \rightarrow \infty$ ergibt (letzte Spalte).

Symbolfolge „Schwarz“:

BABABAABAABAABBBA
 BAABBBAAABBBBAABAB
 AABBAABAABBBBBBAA

Symbolfolge „Blau“:

BAAAAAAAAABAABAABAAA
 BAABAAAAAAAAABAABAAA
 AAAAAAAAAAAAAAAAAAA

Symbolfolge „Rot“:

ABBABABABABBABABAB
 ABABABABABABABABAB
 ABABBBABABABBABABA

Symbolfolge „Grün“:

BBAABBAABBAABAABAAA
 BBAABBAABBBBBAABAAA
 BBAABBBBBAABAAAAAB

© 2011 www.LNTwww.de

Zwischen diesen Entropien bestehen folgende Größenrelationen:

$$H \leq \dots \leq H_3 \leq H_2 \leq H_1 \leq H_0.$$

Nicht bekannt ist die Zuordnung zwischen den Quellen Q1, Q2, Q3, Q4 und den in der Grafik gezeigten Symbolfolgen (Schwarz, Blau, Rot, Grün). Es ist lediglich bekannt, dass die Quelle Q4 einen Wiederholungscode beinhaltet. Zu bestimmen sind für diese Nachrichtenquelle schließlich noch die Entropienäherungen H_2 und H_3 .

Für die Quelle Q4 sind nur $H_1 = 1$ bit/Symbol (\Rightarrow A und B gleichwahrscheinlich), die Näherung $H_4 \approx 0.789$ bit/Symbol und der Entropie-Endwert $H = 0.5$ bit/Symbol angegeben. Letzterer aufgrund der Tatsache, dass bei der entsprechenden Symbolfolge jedes zweite Symbol keinerlei Information liefert.

Quelle	H_0	H_1	H_2	H_3	H_4	...	H
Q1	1.000	0.500	0.500	0.500	0.500	...	0.500
Q2	1.000	1.000	0.750	0.667	0.625	...	0.500
Q3	1.000	1.000	1.000	1.000	1.000	...	1.000
Q4	1.000	1.000	???	???	0.789	...	0.500

© 2011 www.LNTwww.de

Hinweis: Die Aufgabe gehört zum Themengebiet von **Kapitel 1.2**. Für die k -te Entropienäherung gilt bei Binärquellen ($M = 2$) mit der Verbundwahrscheinlichkeit $p_i^{(k)}$ eines k -Tupels:

$$H_k = \frac{1}{k} \cdot \sum_{i=1}^{2^k} p_i^{(k)} \cdot \log_2 \frac{1}{p_i^{(k)}} \quad (\text{Einheit: bit/Symbol}).$$

Fragebogen zu "A1.3: H_0, H_1, H_2, \dots, H "

a) Von welcher Quelle stammt die schwarze Symbolfolge?

- Q1,
- Q2,
- Q3,
- Q4.

b) Von welcher Quelle stammt die blaue Symbolfolge?

- Q1,
- Q2,
- Q3,
- Q4.

c) Von welcher Quelle stammt die rote Symbolfolge?

- Q1,
- Q2,
- Q3,
- Q4.

d) Berechnen Sie die Entropienäherung H_2 des Wiederholungscodes.

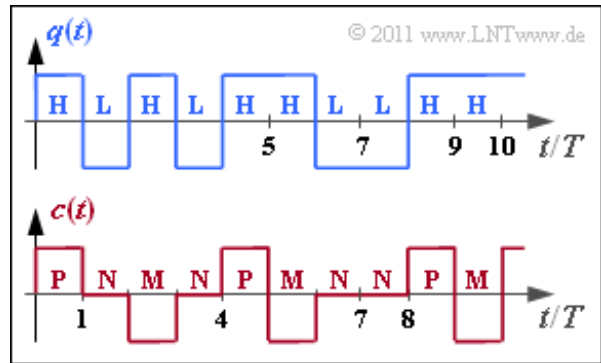
$$H_2 = \quad \text{bit/Symbol}$$

e) Berechnen Sie die Entropienäherung H_3 des Wiederholungscodes.

$$H_3 = \quad \text{bit/Symbol}$$

A1.4: Entropienäherungen H_k

Die Grafik zeigt oben das binäre Quellensignal $q(t)$, das man ebenfalls durch die Symbolfolge $\langle q_v \rangle$ mit $q_v \in \{L, H\}$ beschreiben kann. In der gesamten Aufgabe gelte $p_L = p_H = 0.5$.



Das codierte Signal $c(t)$ und die dazugehörige Symbolfolge $\langle c_v \rangle \in \{P, N, M\}$ ergibt sich aus der

AMI-Codierung (*Alternate Mark Inversion*) nach folgender Vorschrift:

- Das Binärsymbol **L** \Rightarrow *Low* wird stets durch das Ternärsymbol **N** \Rightarrow *Null* dargestellt.
- Das Binärsymbol **H** \Rightarrow *High* wird ebenfalls deterministisch, aber alternierend (daher der Name „AMI“) durch die Symbole **P** \Rightarrow *Plus* und **M** \Rightarrow *Minus* codiert.

In dieser Aufgabe sollen die Entropienäherungen für das AMI-codierte Signal berechnet werden:

- Die Näherung H_1 bezieht sich nur auf die Symbolwahrscheinlichkeiten p_P , p_N und p_M .
- Die k -te Entropienäherung ($k = 2, 3, \dots$) kann nach folgender Gleichung ermittelt werden:

$$H_k = \frac{1}{k} \cdot \sum_{i=1}^{3^k} p_i^{(k)} \cdot \log_2 \frac{1}{p_i^{(k)}} \quad (\text{Einheit: bit/Symbol}).$$

Hierbei bezeichnet $p_i^{(k)}$ die i -te Verbundwahrscheinlichkeit eines k -Tupels.

Hinweis: Die Aufgabe gehört zu **Kapitel 1.2**. In der **Aufgabe Z1.4** wird die tatsächliche Entropie der Codesymbolfolge $\langle c_v \rangle$ zu $H = 1$ bit/Symbol berechnet. Zu erwarten sind die folgenden Größenrelationen:

$$H \leq \dots \leq H_3 \leq H_2 \leq H_1 \leq H_0.$$

Fragebogen zu "A1.4: Entropienäherungen H_k "

a) Wie groß ist der Entscheidungsgehalt des AMI-Codes?

$$H_0 = \text{bit/Symbol}$$

b) Berechnen Sie die erste Entropienäherung.

$$H_1 = \text{bit/Symbol}$$

c) Wie groß ist die Entropienäherung H_2 , basierend auf Zweiertupel?

$$H_2 = \text{bit/Symbol}$$

d) Welchen Wert liefert die Entropienäherung H_3 , basierend auf Dreiertupel?

$$H_3 = \text{bit/Symbol}$$

e) Welche Aussagen gelten für die Entropienäherung H_4 ?

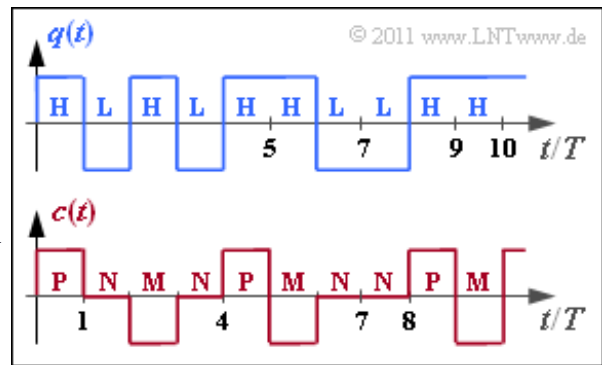
- Es muss über $3^4 = 81$ Summanden gemittelt werden.
- Es gilt $1 \text{ bit/Symbol} < H_4 < H_3$.
- Nach langer Rechnung erhält man $H_4 = 1.333 \text{ bit/Symbol}$.

Z1.4: Entropie der AMI-Codierung

Wir gehen von ähnlichen Voraussetzungen wie in der **Aufgabe A1.4** aus: Eine Binärquelle liefert die Quellensymbolfolge $\langle q_v \rangle$ mit $q_v \in \{L, H\}$, wobei es keine statistischen Bindungen zwischen den einzelnen Folgeelementen gibt.

Für die Symbolwahrscheinlichkeiten gelte:

- $p_L = p_H = 1/2$ (Teilaufgaben a und b),
- $p_L = 1/4, p_H = 3/4$ (Teilaufgaben c, d und e),
- $p_L = 3/4, p_H = 1/4$ (Teilaufgabe f).



Das dargestellte Codesignal $c(t)$ und die zugehörige Symbolfolge $\langle c_v \rangle$ mit $c_v \in \{P, N, M\}$ ergibt sich aus der AMI-Codierung (*Alternate Mark Inversion*) nach folgender Vorschrift:

- Das Binärsymbol **L** \Rightarrow *Low* wird stets durch das Ternärsymbol **N** \Rightarrow *Null* dargestellt.
- Das Binärsymbol **H** \Rightarrow *High* wird ebenfalls deterministisch, aber alternierend (daher der Name „AMI“) durch die Symbole **P** \Rightarrow *Plus* und **M** \Rightarrow *Minus* codiert.

In dieser Aufgabe sollen für die drei oben genannten Parametersätze der Entscheidungsgehalt H_0 sowie die resultierende Entropie H_C der Codesymbolfolge $\langle c_v \rangle$ bestimmt werden. Die relative Redundanz der Codefolge ergibt sich daraus entsprechend der Gleichung

$$r_C = \frac{H_0 - H_C}{H_C}.$$

Hinweis: Die Aufgabe gehört zu **Kapitel 1.2**. Allgemein bestehen folgende Relationen zwischen dem Entscheidungsgehalt H_0 , der Entropie H (hier gleich H_C) und den Entropienäherungen:

$$H \leq \dots \leq H_3 \leq H_2 \leq H_1 \leq H_0.$$

In **Aufgabe A1.4** wurden für gleichwahrscheinliche Symbole **L** und **H** die Entropie-Näherungen wie folgt berechnet (jeweils in bit/Symbol):

$$H_1 = 1.500, \quad H_2 = 1.375, \quad H_3 = 1.292.$$

Fragebogen zu "Z1.4: Entropie der AMI-Codierung"

a) Die Quellensymbole seien gleichwahrscheinlich. Wie groß ist die Entropie H_C der Codesymbolfolge $\langle c_v \rangle$?

$$p_L = p_H: H_C = \quad \text{bit/Ternärsymbol}$$

b) Wie groß ist die relative Redundanz der Codesymbolfolge?

$$p_L = p_H: r_C = \quad \%$$

c) Für die Binärquelle gelte nun $p_L = 1/4$ und $p_H = 3/4$. Welcher Wert ergibt sich nun für die Entropie der Codesymbolfolge?

$$p_L = 1/4: H_C = \quad \text{bit/Ternärsymbol}$$

d) Wie groß ist nun die relative Redundanz der Codesymbolfolge?

$$p_L = 1/4: r_C = \quad \%$$

e) Berechnen Sie die Näherung H_1 der Coderentropie für $p_L = 1/4$, $p_H = 3/4$.

$$p_L = 1/4: H_1 = \quad \text{bit/Ternärsymbol}$$

f) Berechnen Sie die Näherung H_1 der Coderentropie für $p_L = 3/4$, $p_H = 1/4$.

$$p_L = 3/4: H_1 = \quad \text{bit/Ternärsymbol}$$

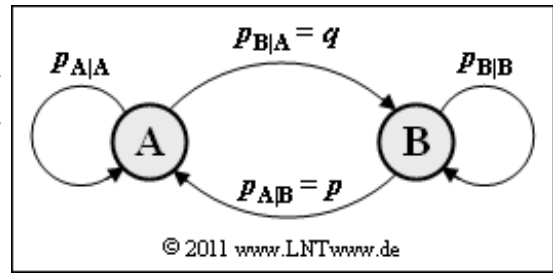
A1.5: Binäre Markovquelle

Die Aufgabe A1.4 hat gezeigt, dass die Berechnung der Entropie bei einer gedächtnisbehafteten Quelle sehr aufwändig sein kann. Man muss dann zunächst (sehr viele) Entropienäherungen H_k für k -Tupel berechnen und kann erst dann mit dem Grenzübergang $k \rightarrow \infty$ die Quellenentropie ermitteln:

$$H = \lim_{k \rightarrow \infty} H_k.$$

Oft tendiert dabei H_k nur sehr langsam gegen den Grenzwert H .

Der Rechengang wird drastisch reduziert, wenn die Nachrichtenquelle **Markoveigenschaften** besitzt. Die Grafik zeigt das Übergangsdiagramm für eine binäre Markovquelle mit den zwei Zuständen (Symbolen) **A** und **B**. Dieses ist durch die beiden bedingten Wahrscheinlichkeiten $p_{A|B} = p$ und $p_{B|A} = q$ eindeutig bestimmt. Die bedingten Wahrscheinlichkeiten $p_{A|A}$ und $p_{B|B}$ sowie die Symbolwahrscheinlichkeiten p_A und p_B lassen sich daraus ermitteln.



Die Entropie der binären Markovkette (mit der Einheit „bit/Symbol“) lautet dann:

$$H = p_{AA} \cdot \log_2 \frac{1}{p_{A|A}} + p_{AB} \cdot \log_2 \frac{1}{p_{B|A}} + p_{BA} \cdot \log_2 \frac{1}{p_{A|B}} + p_{BB} \cdot \log_2 \frac{1}{p_{B|B}}.$$

Bei dieser Gleichung ist zu beachten, dass im Argument des *Logarithmus dualis* jeweils die *bedingten Wahrscheinlichkeiten* $p_{A|A}$, $p_{B|A}$, ... einzusetzen sind, während für die Gewichtung die *Verbundwahrscheinlichkeiten* p_{AA} , p_{AB} , ... zu verwenden sind.

Mit der Entropienäherung erster Ordnung,

$$H_1 = p_A \cdot \log_2 \frac{1}{p_A} + p_B \cdot \log_2 \frac{1}{p_B} \quad (\text{Einheit: bit/Symbol}),$$

sowie der oben angegebenen (tatsächlichen) Entropie H lassen sich bei einer Markovquelle auch alle weiteren Entropienäherungen ($k = 2, 3, \dots$) direkt berechnen:

$$H_k = \frac{1}{k} \cdot [H_1 + (k - 1) \cdot H_M].$$

Hinweis: Diese Aufgabe gehört zum Themengebiet von **Kapitel 1.2**. Mit Ausnahme der Teilaufgabe (f) sei $p = 1/4$ und $q = 1/2$.

Für die (ergodischen) Symbolwahrscheinlichkeiten einer Markovkette erster Ordnung gilt:

$$p_A = \frac{p_{A|B}}{p_{A|B} + p_{B|A}}, \quad p_B = \frac{p_{B|A}}{p_{A|B} + p_{B|A}}.$$

Fragebogen zu "A1.5: Binäre Markovquelle"

a) Geben Sie die Übergangswahrscheinlichkeiten für $p = 1/4$ und $q = 1/2$ an.

$$p = 1/4, q = 1/2: p_{A|A} =$$
$$p_{B|B} =$$

b) Wie groß sind die Symbolwahrscheinlichkeiten?

$$p = 1/4, q = 1/2: p_A =$$
$$p_B =$$

c) Geben Sie die Entropienäherung erster Ordnung an.

$$p = 1/4, q = 1/2: H_1 = \text{bit/Symbol}$$

d) Welche Entropie besitzt diese Markovquelle?

$$p = 1/4, q = 1/2: H = \text{bit/Symbol}$$

e) Welche Näherungen H_k ergeben sich aufgrund der Markoveigenschaften?

$$p = 1/4, q = 1/2: H_2 = \text{bit/Symbol}$$
$$H_3 = \text{bit/Symbol}$$
$$H_4 = \text{bit/Symbol}$$

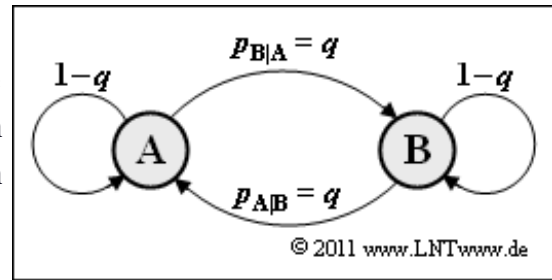
f) Welche Entropie besitzt die Markovquelle mit $p = 1/4$ und $q = 3/4$?

$$p = 1/4, q = 3/4: H = \text{bit/Symbol}$$

Z1.5: Symmetrische Markovquelle

In der Aufgabe A1.5 wurde eine binäre Markovquelle behandelt, bei der die Übergangswahrscheinlichkeiten von A nach B sowie von B nach A unterschiedlich waren. In dieser Aufgabe soll nun gelten:

$$p_{A|B} = p_{B|A} = q \quad (0 \leq q \leq 1).$$



Alle in der Aufgabe A1.5 angegebenen Gleichungen gelten auch hier:

- **Entropie:**

$$H = p_{AA} \cdot \text{ld} \frac{1}{p_{A|A}} + p_{AB} \cdot \text{ld} \frac{1}{p_{B|A}} + p_{BA} \cdot \text{ld} \frac{1}{p_{A|B}} + p_{BB} \cdot \text{ld} \frac{1}{p_{B|B}},$$

- **Erste Entropienäherung:**

$$H_1 = p_A \cdot \text{ld} \frac{1}{p_A} + p_B \cdot \text{ld} \frac{1}{p_B},$$

- **k-te Entropienäherung** ($k = 2, 3, \dots$):

$$H_k = \frac{1}{k} \cdot [H_1 + (k - 1) \cdot H], \quad H = \lim_{k \rightarrow \infty} H_k.$$

Hinweis: Die Aufgabe bezieht sich auf **Kapitel 1.2, Seite 5c**. Bei allen Entropien ist die Pseudoeinheit „bit/Symbol“ hinzuzufügen.

Fragebogen zu "Z1.5: Symmetrische Markovquelle"

a) Berechnen Sie die Symbolwahrscheinlichkeiten.

$q = 1/4: p_A =$
 $p_B =$

b) Berechnen Sie die Quellenentropie H . Welches Ergebnis liefert $q = 1/4$?

$q = 1/4: H =$ bit/Symbol

c) Welche Entropienäherungen erhält man für $q = 1/4$?

$q = 1/4: H_1 =$ bit/Symbol
 $H_2 =$ bit/Symbol
 $H_3 =$ bit/Symbol

d) Bestimmen Sie q derart, dass H maximal wird. Interpretation.

$H \rightarrow \text{Maximum}: q =$

e) Welche Symbolfolgen sind mit $q = 0$ möglich?

- AAAAAA ...
- BBBBBB ...
- ABABAB ...

f) Welche Symbolfolgen sind mit $q = 1$ möglich?

- AAAAAA ...
- BBBBBB ...
- ABABAB ...

A1.6: Nichtbinäre Markovquellen

Die Grafik zeigt zwei ergodische Markovquellen (MQ):

- Die Quelle MQ3 ist durch $M = 3$ Zustände (Symbole) N, M, P gekennzeichnet. Aufgrund der Stationarität haben die Wahrscheinlichkeiten folgende Werte:

$$p_N = 1/2, \quad p_M = p_P = 1/4.$$

- Bei der Quelle MQ4 ist zusätzlich der Zustand O möglich $\Rightarrow M = 4$. Aufgrund der symmetrischen Übergänge sind die stationären Wahrscheinlichkeiten alle gleich:

$$p_N = p_M = p_O = p_P = 1/4.$$

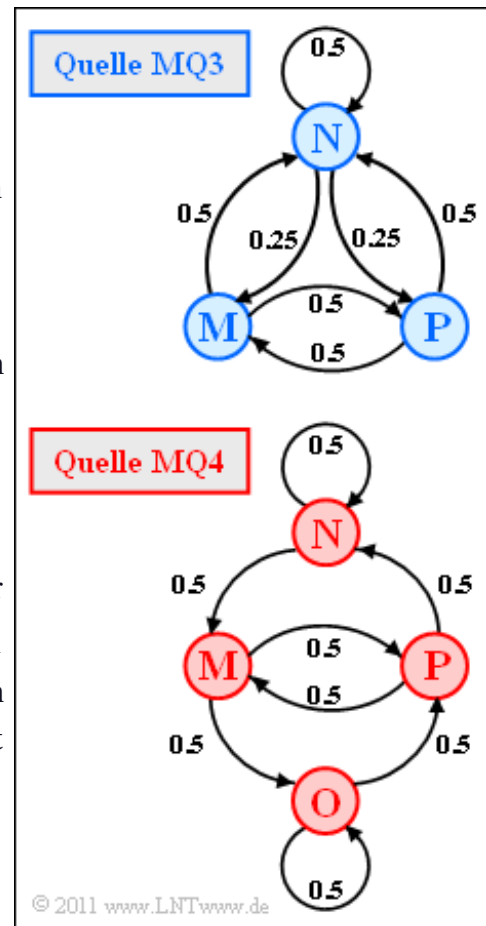
Informationstheoretisch sind Markovquellen von besonderer Bedeutung, da bei diesen – und nur bei diesen – durch H_1 (Entropienäherung, nur auf den Symbolwahrscheinlichkeiten basierend) und H_2 (zweite Entropienäherung, berechenbar mit den Verbundwahrscheinlichkeiten für alle Zweiertupel) gleichzeitig auch

- die weiteren Entropienäherungen $H_k (k = 3, 4, \dots)$ und
- die tatsächliche Quellenentropie H

bestimmt sind. Es gelten folgende Bestimmungsgleichungen:

$$H = 2 \cdot H_2 - H_1, \quad H_k = \frac{1}{k} \cdot [H_1 + (k - 1) \cdot H].$$

Hinweis: Die Aufgabe gehört zu **Kapitel 1.2**. Hier finden Sie auch Hinweise zur Berechnung der ersten und zweiten Entropienäherung.



Fragebogen zu "A1.6: Nichtbinäre Markovquellen"

a) Berechnen Sie die Entropienäherung H_1 der Markovquelle MQ3.

$$\text{MQ3: } H_1 = \text{bit/Symbol}$$

b) Berechnen Sie die Entropienäherung H_2 der Markovquelle MQ3.

$$\text{MQ3: } H_2 = \text{bit/Symbol}$$

c) Wie groß sind die Näherungen H_3 und H_4 und die Quellenentropie H ?

$$\text{MQ3: } H_3 = \text{bit/Symbol}$$

$$H_4 = \text{bit/Symbol}$$

$$H = \text{bit/Symbol}$$

d) Berechnen Sie die Entropienäherung H_1 der Markovquelle MQ4.

$$\text{MQ4: } H_1 = \text{bit/Symbol}$$

e) Berechnen Sie die Entropienäherung H_2 der Markovquelle MQ4.

$$\text{MQ4: } H_2 = \text{bit/Symbol}$$

f) Wie groß sind hier die Näherungen H_3 und H_4 und die Quellenentropie H ?

$$\text{MQ4: } H_3 = \text{bit/Symbol}$$

$$H_4 = \text{bit/Symbol}$$

$$H = \text{bit/Symbol}$$

Z1.6: Ternäre Markovquelle

Die Grafik zeigt eine Markovquelle mit $M = 3$ Zuständen **A**, **B** und **C**. Für die beiden Parameter dieses Markovprozesses soll gelten:

$$0 \leq p \leq 0.5, \quad 0 \leq q \leq 1.$$

Aufgrund der Markoveigenschaft dieser Quelle kann die Entropie auf unterschiedliche Weise ermittelt werden:

- Man berechnet die beiden ersten **Entropienäherungen** H_1 und H_2 . Dann gilt:

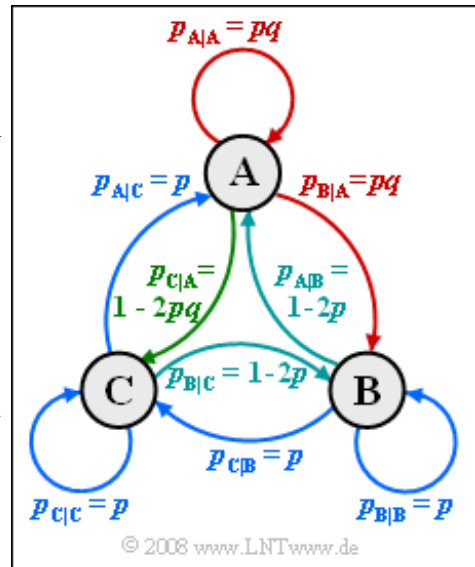
$$H = 2 \cdot H_2 - H_1.$$

- Nach der so genannten *direkten Berechnungsmethode* kann die Entropie aber auch wie folgt berechnet werden (insgesamt 9 Terme):

$$H = p_{AA} \cdot \lg \frac{1}{p_{A|A}} + p_{AB} \cdot \lg \frac{1}{p_{B|A}} + \dots,$$

$$p_{AA} = p_A \cdot p_{A|A}, \quad p_{AB} = p_A \cdot p_{B|A}, \quad \dots$$

Hinweis: Die Aufgabe gehört zum Themenkomplex von **Kapitel 1.2**.



Fragebogen zu "Z1.6: Ternäre Markovquelle"

a) Für welche Parameter p, q ergibt sich die maximale Entropie pro Symbol?

$$p =$$

$$q =$$

$$H_{\max} = \text{bit/Symbol}$$

b) Es sei $p = 1/4$ und $q = 1$. Welcher Wert ergibt sich in diesem Fall für die erste Entropienäherung?

$$p = 1/4, q = 1: H_1 = \text{bit/Symbol}$$

c) Weiterhin gelte $p = 1/4$ und $q = 1$. Welcher Wert ergibt sich in diesem Fall für die zweite Entropienäherung?

$$p = 1/4, q = 1: H_2 = \text{bit/Symbol}$$

d) Wie groß ist die Quellenentropie mit $p = 1/4, q = 1$?

$$p = 1/4, q = 1: H = \text{bit/Symbol}$$

e) Wie groß ist die Quellenentropie mit $p = 1/2, q = 0$?

$$p = 1/2, q = 0: H = \text{bit/Symbol}$$

A1.7: Entropie natürlicher Texte

Anfang der 1950er Jahre schätzte **Claude E. Shannon** die Entropie H der englischen Sprache mit einem bit pro Zeichen ab. Kurze Zeit später kam **Karl Küpfmüller** bei einer empirischen Untersuchung der deutschen Sprache auf einen Entropiewert von $H = 1.3$ bit/Zeichen, also nur etwas größer. Die Ergebnisse von Shannon und Küpfmüller beruhen dabei interessanter Weise auf zwei völlig unterschiedlichen Methoden.

Die differierenden Ergebnisse lassen sich eher nicht mit den geringen Differenzen hinsichtlich des Symbolumfangs M erklären:

- Shannon ging von 26 Buchstaben und dem Leerzeichen aus $\Rightarrow M = 27$.
- Küpfmüller ging von $M = 26$ Buchstaben aus, ebenfalls ohne zwischen Groß- und Kleinschreibung zu unterscheiden.

Mit dieser Aufgabe soll gezeigt werden, wie sich

- Auslöschungen (*Erasures*) \Rightarrow man kennt den Ort eines Fehlers,
- Zeichenfehler \Rightarrow es ist nicht offensichtlich, was falsch und was richtig ist,

auf die Verständlichkeit eines Textes auswirken. Unser Text beinhaltet auch die typisch deutschen Buchstaben „ä“, „ö“, „ü“ und „ß“ sowie Ziffern und Interpunktion. Außerdem wird zwischen Groß- und Kleinschreibung unterschieden.

A) $N = 199, E = 20 \Rightarrow$ ca. 10% „Erasure“

Wie kam Küpfmüller zum Ergebnis $H = 1.5$ bit/Buchstabe? Da er für die Statistik von Wortgruppen oder ganzen Sätzen keine Veröffentlichungen gab, schätzte er die Entropie deutscher Texte wie folgt ab:

B) $N = 240, E = 48 \Rightarrow$ ca. 20% „Erasure“

Ein zusammenhängender, sonst beliebiger deutscher Text wird in einem bestimmten Ort abgelekt. Der Text vorher wird gelesen, und der Leser soll versuchen, die folgenden Worte aus dem vorhergehenden Text und dem Zusammenhang zu ermitteln.

C) $N = 319, E = 97 \Rightarrow$ ca. 30% „Erasure“

Bei sehr vielen Wörtern ergibt die prozentuale Zahl der Treffer ein Maß für die Bindungen zwischen den Wörtern und Sätzen. Es zeigt sich, dass ein und derselbe Textart desselben Wortes verhältnismäßig häufig, etwa 100 bis 200 Vorkommen, ein oder unter Endert Trefferverhältnis erreicht wird.

D) $N = 257, E = 103 \Rightarrow$ ca. 40% „Erasure“

Das Trefferverhältnis ergibt sich stark von der Wortart ab, aber für verschiedene Textarten Wert von 15 bis 33% mit dem Mittelwert von 22%. Dies bedeutet aber auch, dass durch die Entropie 2% aller Wörter auf dem Zusammenhang ermittelt werden.

E) $N = 239, E = 103 \Rightarrow$ ca. 50% „Erasure“

Ausgedrückt: Die Anzahl der Wörter pro Silbe ist mit dem Faktor 2 reduziert worden, da es sich um Nachrichten handelt, die die Entropie pro Silbe betreffen. Die Entropie pro Silbe beträgt 1.51 bit/Buchstabe.

F) $N = 197, F = 39 \Rightarrow$ ca. 20% falsche Zeichen

Küpfmüller hat sein Ergebnis mit einem Bergleuzgarew XppirÄschen Unterpuchunge aller Silben überprüft und Skat mit dem Reduktionsfaktor von R. 54 howTv H3 j W. 8 Y, rf diß EDtrezie 1.51 ait/Buchstabe.

© 2011 www.LNTwww.de

In der Abbildung ist dieser Text, der von Küpfmüllers Vorgehensweise handelt, in sechs Blöcke der Länge $N = 197$ bis $N = 319$ aufgeteilt. Beschrieben ist die Überprüfung seiner ersten Analyse (1.3 bit/Zeichen) auf völlig anderem Wege, die zum Ergebnis 1.51 bit/Zeichen führte.

- In den oberen fünf Blöcken erkennt man *Erasures* mit verschiedenen Wahrscheinlichkeiten zwischen 10% und 50%.
- Im letzten Block sind Zeichenfehler mit 20-prozentiger Verfälschungswahrscheinlichkeit eingefügt.

Der Einfluss solcher Zeichenfehler auf die Lesbarkeit eines Textes soll in der Teilaufgabe (d) verglichen werden mit dem zweiten (rot umrandeten) Block, für den die Wahrscheinlichkeit eines *Erasures* ebenfalls

20% beträgt.

Hinweis: Die Aufgabe bezieht sich auf das **Kapitel 1.3** dieses Buches. Bezug genommen wird auch auf die relative Redundanz einer Folge, wobei mit dem **Entscheidungsgehalt** H_0 und der **Entropie** H gilt:

$$r = \frac{H_0 - H}{H_0}.$$

Fragebogen zu "A1.7: Entropie natürlicher Texte"

a) Von welchem Symbolumfang ist Küpfmüller ausgegangen?

$$M =$$

b) Welche relative Redundanz ergibt sich aus Küpfmüllers Entropiewert?

$$r = \quad \quad \quad \%$$

c) Wie lässt sich das Ergebnis der Teilaufgabe (b) interpretieren? Gehen Sie jeweils von einer Textdatei mit $M = 26$ unterschiedlichen Zeichen aus.

- Eine solche Textdatei hinreichender Länge ($N \rightarrow \infty$) könnte man mit $1.3 \cdot N$ Binärsymbolen darstellen.
- Eine solche Textdatei mit $N = 100000$ Zeichen könnte man mit 130000 Binärsymbolen darstellen.
- Ein Leser kann den Text auch dann noch verstehen (oder zumindest erahnen), wenn 70% der Zeichen ausgelöscht sind.

d) Was erschwert die Verständlichkeit eines Textes mehr?

- 20% Auslöschungen (*Erasures*),
- eine Zeichenfehlerwahrscheinlichkeit von 20%.

A1.8: Synthetisch erzeugte Texte

Das Praktikum [Söd01] verwendet das Windows-Programm „Wertdiskrete Informationstheorie“. Der nachfolgende Link \Rightarrow **WDIT** führt zur ZIP-Version des Programms.

Aus einer gegebenen Textdatei VORLAGE kann man

- die Häufigkeiten von Buchstabentripeln wie „aaa“, „aab“, ... , „xyz“, ... ermitteln und in einer Hilfsdatei abspeichern,
- danach eine Datei SYNTHESE erzeugen, wobei das neue Zeichen aus den beiden letzten Zeichen und den abgespeicherten Tripel-Häufigkeiten generiert wird.

Ausgehend von der deutschen und der englischen Bibelübersetzung haben wir so zwei Dateien synthetisiert, die in der Grafik mit

- Datei 1 (rote Umrandung),
- Datei 2 (grüne Umrandung)

bezeichnet sind. Nicht bekannt gegeben wird, welche Datei von welcher Vorlage stammt. Dies zu ermitteln ist Ihre erste Aufgabe.

Die beiden Vorlagen basieren auf dem natürlichen Alphabet (26 Buchstaben) und dem Leerzeichen („LZ“) $\Rightarrow M = 27$. Bei der deutschen Bibel wurden die Umlaute ersetzt, zum Beispiel „ä“ \Rightarrow „ae“.

Die **Datei 1** weist folgende Eigenschaften auf:

- Die häufigsten Zeichen sind „LZ“ mit 19.8%, gefolgt von „e“ mit 10.2% und „a“ mit 8.5%.
- Nach „LZ“ (Leerzeichen) tritt „t“ mit 17.8% am häufigsten auf.
- Vor einem Leerzeichen ist „d“ am wahrscheinlichsten.
- Die **Entropienäherungen**, jeweils mit der Einheit bit/Zeichen, wurden wie folgt ermittelt:

$$H_0 = 4.76, \quad H_1 = 4.00, \quad H_2 = 3.54, \quad H_3 = 3.11, \quad H_4 = 2.81.$$

Dagegen ergibt die Analyse von **Datei 2**:

- Die häufigsten Zeichen sind „LZ“ mit 17.6% gefolgt von „e“ mit 14.4% und „n“ mit 8.9%.
- Nach „LZ“ ist „d“ am wahrscheinlichsten (15.1%) gefolgt von „s“ mit 10.8%.
- Nach „LZ“ und „d“ sind die Vokale „e“ (48.3%), „i“ (23%) und „a“ (20.2%) dominant.
- Die Entropienäherungen unterscheiden sich nur geringfügig von denen der Datei 1.
- Für größere k -Werte sind diese etwas größer, zum Beispiel $H_3 = 3.11 \Rightarrow 3.17$.

Hinweis: Die Aufgabe bezieht sich auf das **Kapitel 1.3**. Anzumerken ist, dass unsere Analyse eher auf

Datei 1:

```
abor wast thes bled ther and labta kind  
mosethe abrah of and sers ame begat eat  
lager wee land tho hat man thy  
shadstere look as begypt alleture of  
munto you fat there whe he but four th  
sto thento pris to theart rood gods and  
upot roure sainesaid and th of and  
thalt sain sione wing wenahs and hen  
sau and to disaraidie bar and is con th  
and and ey day welf and godst is alk  
soeve and in th goat now the he nowm  
mordid ith agatents hing eve the thathe  
ajud the ich in have and thee ime and  
begaid flord said ife chend eap eve ing  
nocklefors aftere beguel theld rahought  
moch yeathe the the asto not an whe  
ther their the an that and hathom hern  
said jacor as of and beartereass nothe
```

Datei 2:

```
aber scheinen zu ris und atter von kana  
akinestote ein wis der warbling dirkein  
und sehr demt wer eiden zu ingerker  
angen die miten dasterkan ber auf desel  
den dasua ungewistaren wohnt david de  
als land die eind das gen mittet sollem  
mach und ihm zusameit nas der haberija  
kseiter west wachweinn und ma goll das  
aund ken freuchlaufkottem de fam zurden  
ber ganigs spern tord jaher nit unde  
jern den und eite meinern von unschens  
blend nachafte kinesusen zu ihr binerre  
ma ewir kine auf deind seile meich  
imose hab willbs einund wich namm knes  
ine dennem mung vonigewas begehmn rudenn  
deinerzahwirdersat josall den nich zu  
zus nachwegnichweger habeiner zogen  
ungegern nie dir dier dem dieber hon in
```

© 2011 www.LNTwww.de

einen gleichen Entropiewert für Englisch und Deutsch schließen lässt. Dafür spricht auch, dass beide Übersetzungen nahezu gleich groß sind (ca. 4 Millionen Zeichen). Hätte Deutsch eine um 30% größere Entropie als Englisch, dann müsste nach unserer Meinung die englische Version um 30% länger sein, wenn man von gleichem Inhalt der beiden Übersetzungen ausgeht. Wir erheben aber keinen Anspruch auf die Richtigkeit unserer Argumentation.

Fragebogen zu "A1.8: Synthetisch erzeugte Texte"

a) Welche Vorlagen wurden für die hier gezeigte Textsynthese verwendet?

- Die Datei 1 (rot) basiert auf einer englischen Vorlage.
- Die Datei 1 (rot) basiert auf einer deutschen Vorlage.

b) Vergleichen Sie die mittleren Wortlängen von Datei 1 und Datei 2.

- Die Wörter der „englischen“ Datei sind im Mittel länger.
- Die Wörter der „deutschen“ Datei sind im Mittel länger.

c) Welche Aussagen gelten für die Entropienäherungen?

- VORLAGE und SYNTHESE liefern ein nahezu gleiches H_1 .
- VORLAGE und SYNTHESE liefern ein nahezu gleiches H_2 .
- VORLAGE und SYNTHESE liefern ein nahezu gleiches H_3 .
- VORLAGE und SYNTHESE liefern ein nahezu gleiches H_4 .

d) Welche Aussagen treffen für den „englischen“ Text zu?

- Die meisten Wörter beginnen mit „t“.
- Die meisten Wörter enden mit „t“.

e) Welche Aussagen könnten für deutsche Texte gelten?

- Nach „de“ ist „r“ am wahrscheinlichsten.
- Nach „da“ ist „s“ am wahrscheinlichsten.
- Nach „di“ ist „e“ am wahrscheinlichsten.