

# 2 Informationstheorie

---

Formale Grundlagen der Informatik I  
Herbstsemester 2012

Robert Marti

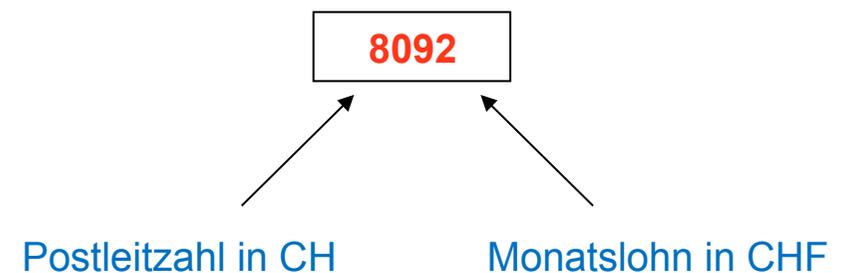
Vorlesung teilweise basierend auf Unterlagen  
von Prof. emer. Helmut Schauer

---

# Grundbegriffe

---

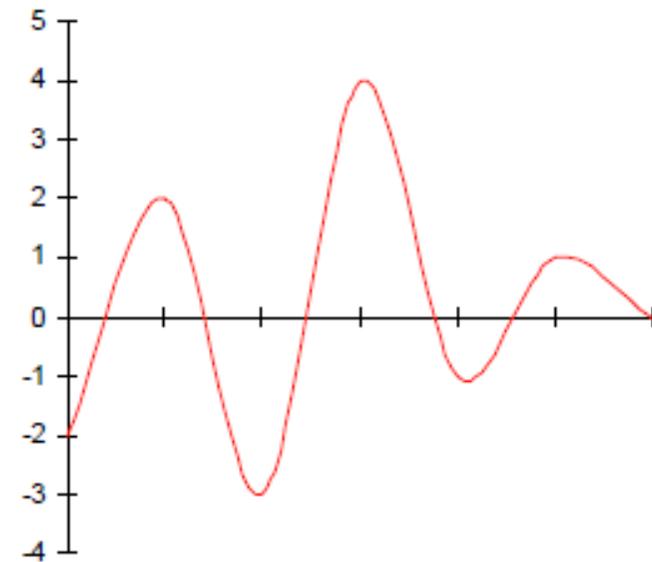
- Informatik  
(IT: Information Technology, ICT: Information & Communication Technology)
  - Disziplin / Wissenschaft der Informationsverarbeitung (Info. Processing)
- Signal
  - **analog**
  - **digital**
- Nachricht (message) – auch Daten (data)
  - Begriff Nachricht wird meist bei Datenübermittlung verwendet
- Information (Text, Bilder, Ton, ... )
  - letztlich **Nachrichten (Daten)** + **Kontext**



# Analoges Signal

---

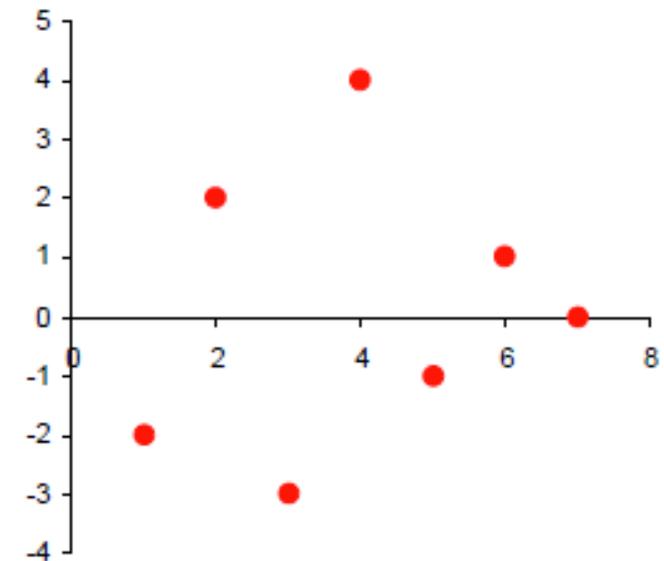
- Beispiele:  
Radio, TV, Telefon, Magnetband, Vinyl Schallplatten, ...
- Eigenschaften
  - kontinuierlich
  - im Prinzip beliebig genau
  - Verarbeitung eher aufwendig



# Digitales Signal

---

- Beispiele:  
PC Spiele, VoIP (**V**oice **o**ver **I**nternet **P**rotocol), Digitalkameras,  
MP3 Musik (auch FLAC Musik), ...
- Eigenschaften
  - diskrete Zustände
  - begrenzt genau (endlich viele verschiedene Zustände)
  - Verarbeitung weniger aufwendig



# Nachricht

---

- Nachrichten (Daten) besitzen ein bestimmtes "Format", sie folgen einer (formalen) **Sprache**.
- Eine Sprache wird definiert durch
  - eine **Syntax** (bei natürlichen Sprachen Grammatik genannt)
    - ➔ **Repräsentation**

Beispiele (hier externe Repräsentationen in Textform):

1. Feb. 2011, 01.02.11, 2/1/11, 2011-02-01, 32. Tag im 2011

- eine **Semantik**
  - ➔ **Bedeutung**

Beispiel: Ziffern vor Punkt = Tag im Monat, danach folgen 3 Buchstaben für den Monat, danach 4 Ziffern für das Jahr

# Bispiele für Sprachen

---

- Chemische Formeln



- Beschreibung von Schachzügen

b2 x c4!

- Mathematische Formelsprache

$(-b + \text{sqrt}(b*b - 4*a*c)) / (2*a)$  bzw.

$$\frac{-b + \sqrt{b^2 - 4ac}}{2a}$$

- Musiknoten

A musical score for five instruments: Horn in F, Violin 1, Violin 2, Viola, and Bass. The score is in 2/4 time and consists of five measures. The Horn part starts with a quarter note C4, followed by eighth notes G4 and A4, then a quarter note Bb4, and finally a quarter note C5. The Violin 1 and 2 parts start with a quarter note C4, followed by a quarter rest, then a quarter note G4, and finally a quarter note C5. The Viola part starts with a quarter note C4, followed by a quarter rest, then a quarter note G4, and finally a quarter note C5. The Bass part starts with a quarter note C4, followed by a quarter rest, then a quarter note G2, and finally a quarter note C3.

- Programmiersprachen

```
s = 0; for (int i = 1; i <= n; i++) s = s+i;
```

# Syntax und Alphabet

---

- Die Definition einer Syntax stützt sich auf ein **Alphabet**: Menge der zulässigen Zeichen
- Beispiele
  - 26 (Gross-) Buchstaben
  - 2 x 26 Buchstaben, 10 Ziffern, Spezialzeichen, insges. ca. 100 Zeichen

– andere:

$\{ \alpha \mid \beta \mid \gamma \mid \dots \mid \omega \}$   
 $\{ 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \}$   
 $\{ \Upsilon \mid \text{♁} \mid \text{♂} \mid \text{♁} \}$   
 $\{ \clubsuit \mid \diamondsuit \mid \heartsuit \mid \spadesuit \}$   
 $\{ \ddagger \mid \ddagger \}$   
 $\{ \text{♁} \mid \text{♁} \}$   
 $\{ + \mid - \}$   
 $\{ \emptyset \mid 1 \}$

} Binäralphabete

# Code

- Beispiel **Morsecode**: Code eines Alphabets mit Zeichen - und —



A	· ·	N	· ·		
B	· · ·	O	· · ·	1	· · · ·
C	· · · ·	P	· · · ·	2	· · · · ·
D	· · ·	Q	· · · · ·	3	· · · · · ·
E	·	R	· · ·	4	· · · · · · ·
F	· · ·	S	· · ·	5	· · · · · · · ·
G	· · · ·	T	· ·	6	· · · · · · · · ·
H	· · · · ·	U	· · ·	7	· · · · · · · · · ·
I	· ·	V	· · · ·	8	· · · · · · · · · · ·
J	· · · · ·	W	· · · · ·	9	· · · · · · · · · · · ·
K	· · · ·	X	· · · · ·	0	· · · · · · · · · · · · ·
L	· · · · ·	Y	· · · · · ·		
M	· · · · · ·	Z	· · · · · · ·		

# Binärcode

---

- interne Repräsentation in Digitalcomputern

<b>0</b>	<b>0000</b>
<b>1</b>	<b>0001</b>
<b>2</b>	<b>0010</b>
<b>3</b>	<b>0011</b>
<b>4</b>	<b>0100</b>
<b>5</b>	<b>0101</b>
<b>6</b>	<b>0110</b>
<b>7</b>	<b>0111</b>
<b>8</b>	<b>1000</b>
<b>9</b>	<b>1001</b>

- Bem.: Zahl 10 wird intern als  $1010_2$  (binär) oder evt. als  $0001_2\ 0000_2$  (BCD = Binary Coded Decimal) dargestellt (z.B. Taschenrechner, Sprache COBOL)



# Darstellung von Werten durch Bitfolgen

---

- Eine Bitfolge der Länge  $l$  (auch **Wortlänge** genannt) gestattet die Darstellung von  $n$  Werten, wobei:

$$n = 2^l$$

- Wenn  $n$  verschiedene (gleich wahrscheinliche) Werte dargestellt werden sollen, dann werden demnach im Minimum

$$\log_2 n \text{ (auch } \text{ld } n, \text{ wobei } \text{ld} = \text{Logarithmus Dualis})$$

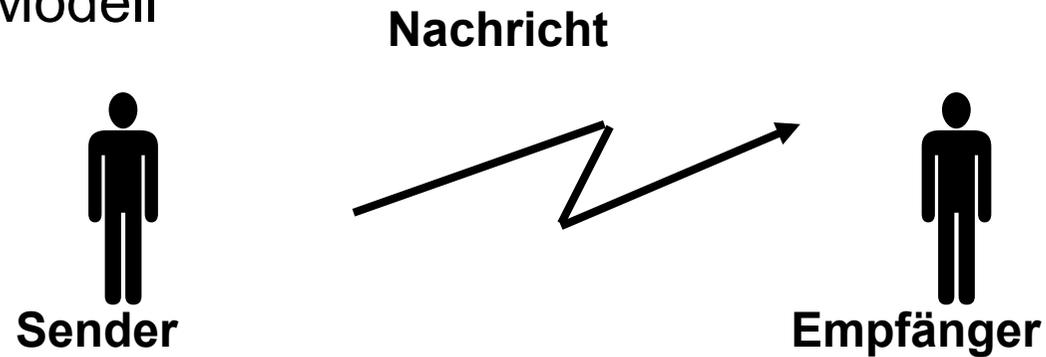
Bits benötigt

- Beispiele: Ganze Zahlen mit und ohne Vorzeichen  
unsigned integer:  $[0, 2^l - 1]$ ,  $l$  meistens = 8, 16, 32, oder 64  
signed integer:  $[-2^{l-1}, +2^{l-1} - 1]$  (**1** Bit für Vorzeichen)

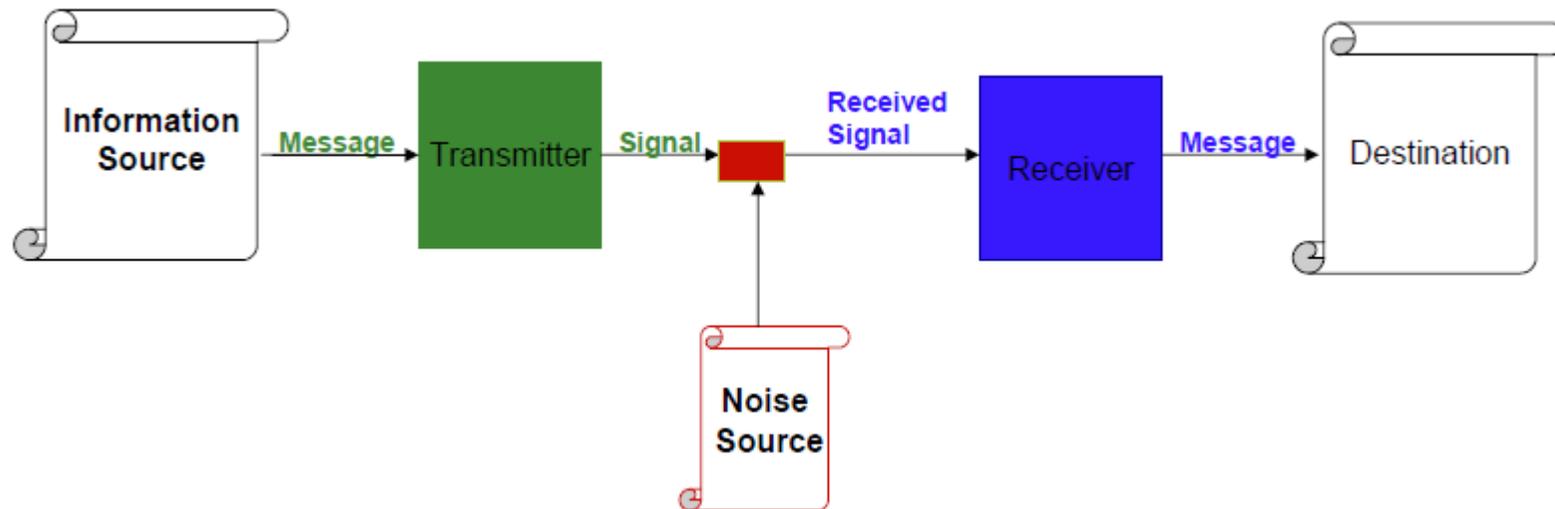
# Ein abstraktes Modell für Kommunikation

---

- einfaches Modell



- etwas detaillierteres Modell



# Messen des Informationsgehalts

---

- Grundidee:  
Wir ordnen einem Zufallsereignis, welches 2 mögliche, **gleich wahrscheinliche** Resultate hat (z.B. einem Münzenwurf mit Ausgang "Kopf" bzw. "Zahl"), das Mass 1 (gemessen in bit) zu
- Verallgemeinerung:  
Einem Zufallsereignis mit  $k$  **gleich wahrscheinlichen** Ausgängen ordnen wir für den Informationsgehalt das Mass

$$\log_2 k = -\log_2 (1/k) = -\log_2 p \quad [\text{bit}]$$

zu, wobei  $p = 1/k$  die Wahrscheinlichkeit eine Ausgangs ist.

- Begründung:  
Annahme: Jemand wählt zufällig eine Zahl zwischen 1 und  $k$ .  
Durch wieviele ja/nein Fragen kann diese Zahl bestimmt werden?

# Einige Eigenschaften des Informationsgehalts

---

Der Informationsgehalt  $h(p)$

- ist unabhängig von der Codierung
- steigt wenn die Wahrscheinlichkeit  $p$  sinkt  
(da die Unsicherheit über den Ausgang des Zufallsexperiments steigt)
- ist "additiv":  $h(p_1 \cdot p_2) = h(p_1) + h(p_2)$
- immer gleiche Meldung :  $h(1) = 0$
- Meldung kommt nie vor:  $h(0) = \infty$
- 2 gleich wahrscheinliche Meldungen:  $h(1/2) = 1$
- Informationsgehalt einer Dezimalziffer:  $h(1/10) \approx 3.32$

# Pro Memoria: Rechnen mit Logarithmen ...

---

... hier mit dem Logarithmus Dualis  $\text{ld}$  (bzw.  $\log_2$ )

- $2^{\text{ld } x} = x$
- $\text{ld } 1/x = -\text{ld } x$
- $\text{ld } (x*y) = \text{ld } x + \text{ld } y$
- $\text{ld } 2^x = 1 + \text{ld } x$
- $\text{ld } x^n = n*\text{ld } x$
- $\text{ld } x^2 = 2*\text{ld } x$
- $\text{ld } x = \ln x / \ln 2 \sim 1.44*\ln x$
- $\text{ld } x = \log x / \log 2 \sim 3.32*\log x$

$\log$  = 10-er Logarithmus,  $\ln$  = logarithmus naturalis

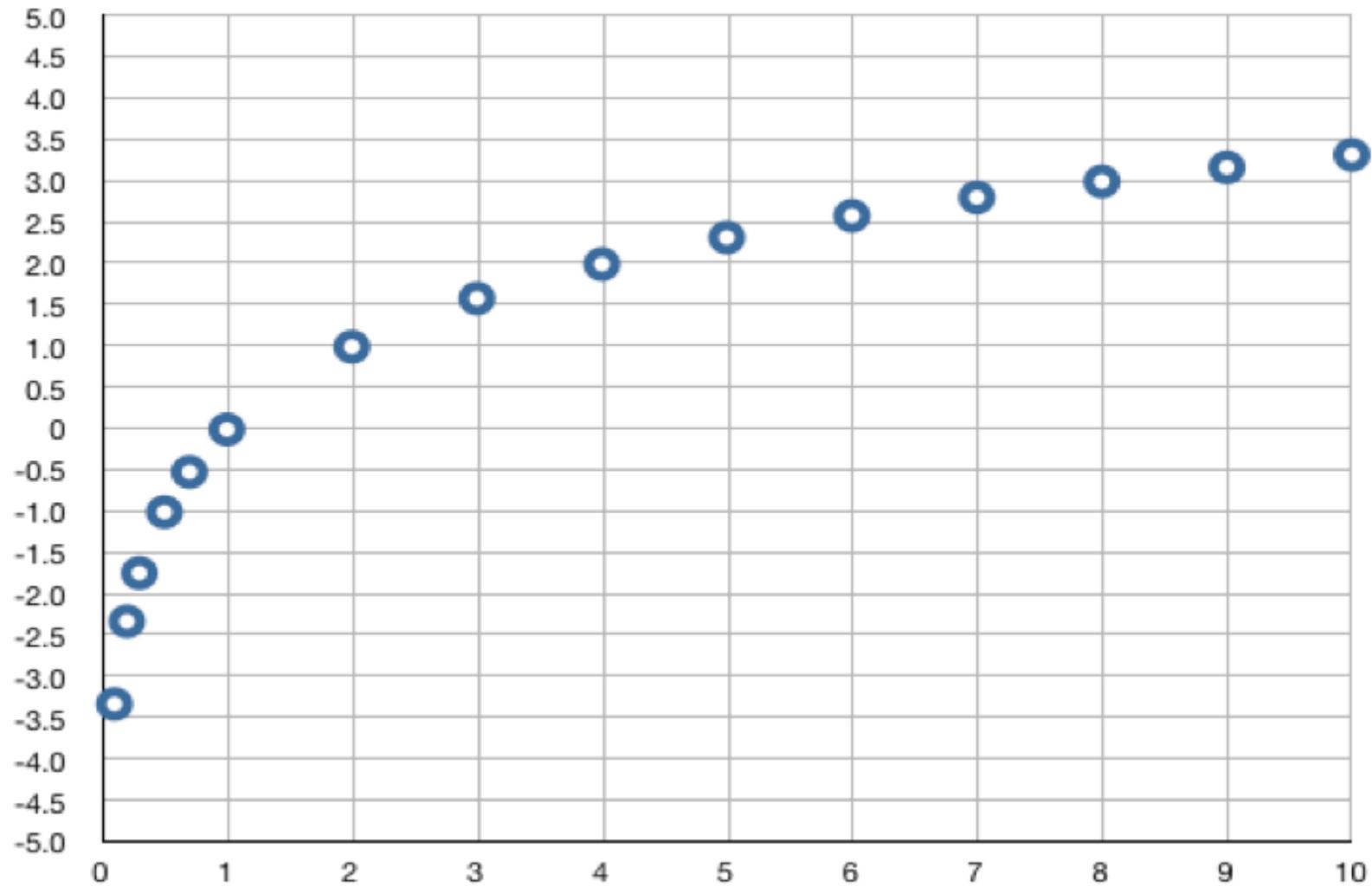
# Wertetabelle des Logarithmus Dualis

---

<b>1</b>	<b>0.00</b>
<b>2</b>	<b>1.00</b>
<b>3</b>	<b>1.58</b>
<b>4</b>	<b>2.00</b>
<b>5</b>	<b>2.32</b>
<b>6</b>	<b>2.58</b>
<b>7</b>	<b>2.81</b>
<b>8</b>	<b>3.00</b>
<b>9</b>	<b>3.17</b>
<b>10</b>	<b>3.32</b>

# Graph des Logarithmus Dualis

---



# Entropie: Messen des **mittleren** Informationsgehalts

---

- Was geschieht, wenn nicht alle Ausgänge des Zufallsereignisses gleich wahrscheinlich sind?
- Die (individuellen) Informationsgehalte der  $k$  verschiedenen Ausgänge des Zufallsereignisses müssen **mit ihrer jeweiligen Wahrscheinlichkeit**  $p_i$  ( $1 \leq i \leq k$ ) gewichtet werden:

$$H(p_1, p_2, \dots, p_k) = - (p_1 \cdot \log_2 p_1 + p_2 \cdot \log_2 p_2 + \dots + p_k \cdot \log_2 p_k) \quad [\text{bit}]$$

bzw.

$$H = - \sum_{i=1}^k p_i \log_2 p_i \quad [\text{bit}]$$

- Der mittlere Informationsgehalt (average information content) wird auch als **Entropie (entropy)** bezeichnet.

# Redundanz

---

- Mittlerer Informationsgehalt (Entropie):

$$H = - \sum_i p_i \cdot \log_2 p_i \quad [\text{bit}]$$

- Mittlere Wortlänge:

$$L = - \sum_i p_i \cdot l_i \quad [\text{bit}]$$

wobei  $l_i$  die Wortlänge der Repräsentation des  $i$ -ten Ausgangs des Zufallsexperiments ist

- Redundanz:

$$R = L - H \quad [\text{bit}]$$

- Satz von Shannon (Shannon's Coding Theorem):

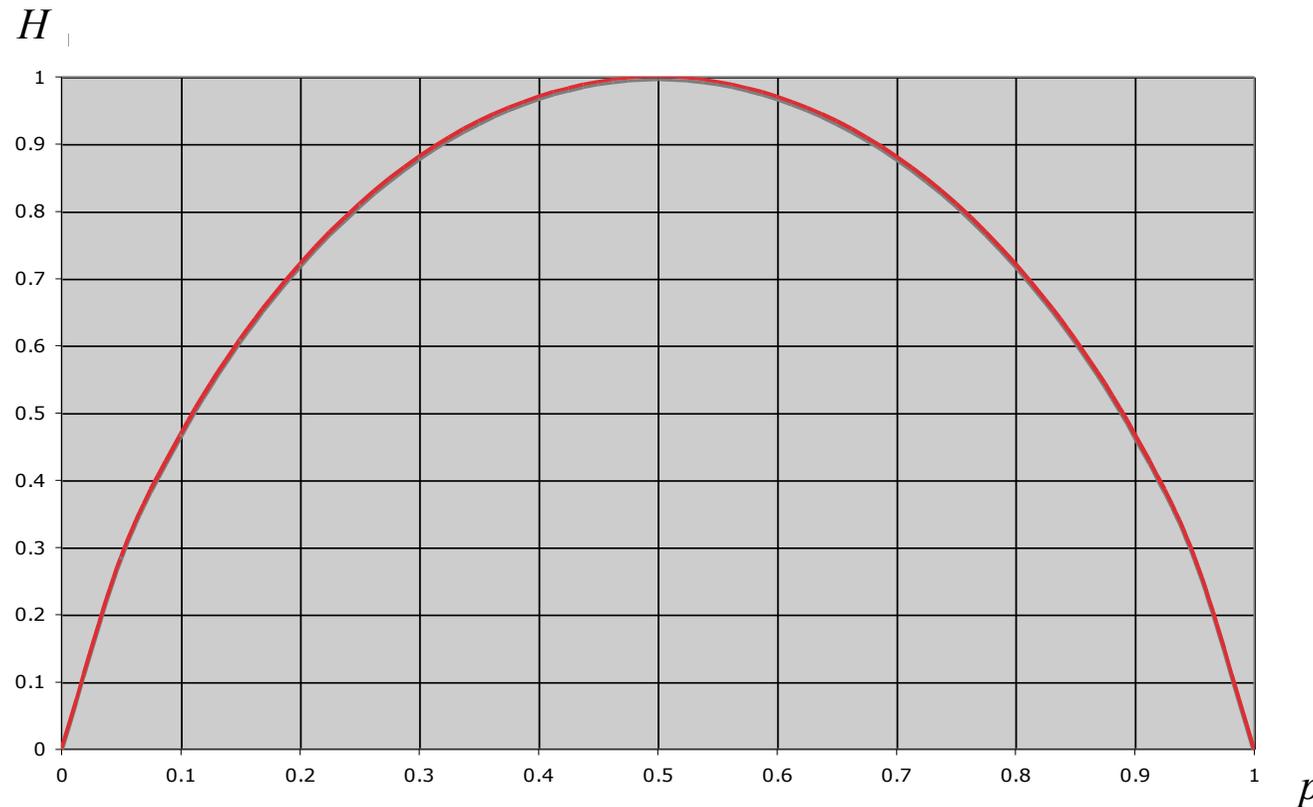
$$L \geq H$$

# Entropiefunktion bei 2 möglichen Ereignissen

---

- Ein Ereignis hat W'keit  $p$  , das andere die W'keit  $(1 - p)$
- Mittlerer Informationsgehalt (Entropie) in diesem Fall:

$$H = - ( p \cdot \log_2 p + (1 - p) \cdot \log_2 (1 - p) ) \quad [\text{bit}]$$



Bem.:  
 $p \cdot \log_2 p = 0$  für  $p = 0$

# Beispiel: Dezimalziffern

---

alle Ziffern gleich wahrscheinlich

- Informationsgehalt

$$H = - 10 \cdot ( 0.1 \cdot \log_2 0.1 ) = -\log_2 0.1 = \log_2 10 \approx 3.32 \text{ [bit]}$$

- Wortlänge

$$L = \text{ceil}(\log_2 10) = 4 \text{ [bit]}$$

- Redundanz

$$R = L - H \approx 4 - 3.32 = 0.68 \text{ [bit]}$$

# Beispiel: (Gross-) Buchstaben

---

alle Zeichen gleich wahrscheinlich

- Informationsgehalt

$$H = \log_2 26 \approx 4.7 \text{ [bit]}$$

Zeichen gemäss empirischer Verteilung in der deutschen Sprache

- Informationsgehalt

$$H \approx 4.1 \text{ [bit]}$$

Bem.:

Bei Ausnutzen der Verteilung von Bigrammen (Buchstabenpaaren):  $H \approx 3.5$  [bit]

# Beispiel: Worte in der deutschen Sprache

---

10 Millionen Wörter mit unterschiedlichen Wahrscheinlichkeiten

- Informationsgehalt

$$H \approx 11.8 \text{ [bit]}$$

- mittlere Wortlänge

$$L \approx 5.7 \text{ [Buchstaben]}$$

# Informationsfluss beim Lesen

---

- ca. 25 Buchstaben pro Sekunde [s]

entspricht ca. 50 bit/s

In 60 Jahren kann ein Mensch ca.  $3 \cdot 10^{10}$  bit aufnehmen

# Auflösung des menschlichen Auges

---

- ca. 6 Megapixel



# Speicherkapazität des Gehirns

---

- ca.  $10^{12}$  bit



# Erbinformation

---

- ca.  $10^{10}$  bit



# Beispiel für Komprimierung (1)

---

- Annahme:  
Alphabet bestehend aus Zeichen  $a$ ,  $b$ ,  $c$  und  $d$   
mit Auftretens-Wahrscheinlichkeiten 0.1, 0.2, 0.3 und 0.4
- Typisches Encoding:  
Jeder Buchstabe mit 2 bit, z.B.

$a \rightarrow 00$

$b \rightarrow 01$

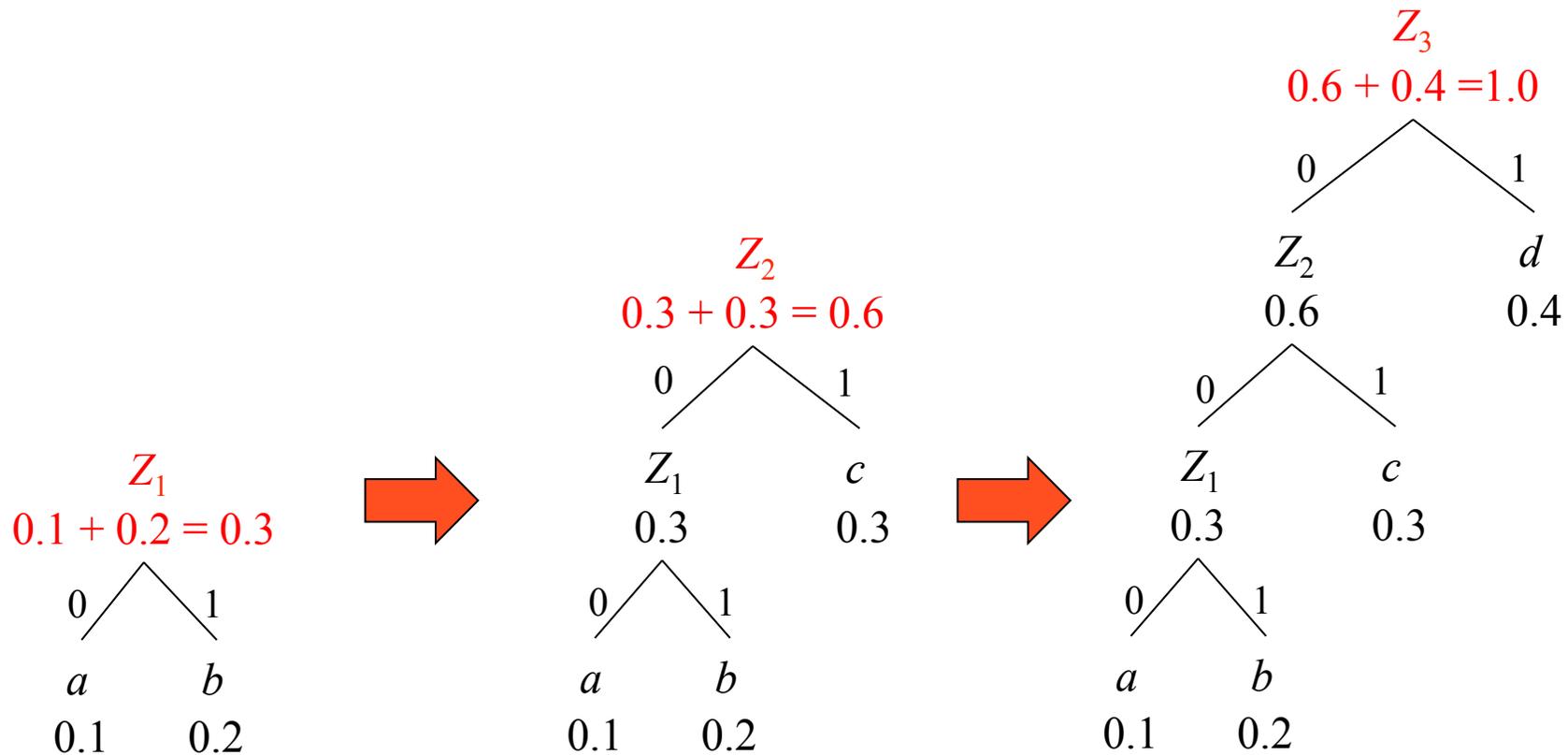
$c \rightarrow 10$

$d \rightarrow 11$

- Mittlere Wortlänge: 2 bit

## Beispiel für Komprimierung (2)

- Konstruktion eines Codebaums (Huffman Encoding):  
Die 2 seltensten Zeichen werden in einem "bottom up" Prozess schrittweise zu "Superzeichen" ( $Z_1$ ,  $Z_2$ ,  $Z_3$ ) zusammengefasst



# Beispiel für Komprimierung (3)

---

- Encoding: Tree Walk

- Abstieg "nach links" ergibt binäre Ziffer 0
- Abstieg "nach rechts" ergibt binäre Ziffer 1

$a \rightarrow 000$

$b \rightarrow 001$

$c \rightarrow 01$

$d \rightarrow 1$

- Decoding: Lesen einzelner bits

1. bis zum ersten bit mit Wert 1
2. maximal 3 bits

- Mittlere Worlänge

1.9 bit (computation left as an exercise to the reader :-)

