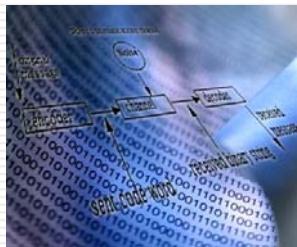


Kapitel 5: Informationsquellen



Was bisher geschah...

- **Definition:** Die Entropie einer diskreten Zufallsvariablen X ist
$$H(X) = -\sum_{i=1}^L p_X(x_i) \log_2 p_X(x_i)$$
 - Berechnung von Information erfordert eine **Modellbildung**
 - Ein Modell ist eine Abstraktion der Realität
 - Information ist also vom **Modellwissen** abhängig
 - Modellbildung der Shannon'schen Informationstheorie beruht auf Probabilistik
 - Makov-Modelle

$$H(X) = -\sum_{i=1}^L p_X(x_i) \log_2 p_X(x_i)$$

Informationstheorie

8

Informationsquellen

- ❑ Unsere Modellbildung beruht auf Wahrscheinlichkeitsverteilungen einer Menge möglicher Ereignisse
 - ❑ Ein Ereignis ist die Auswahl/das Auftreten eines Symbols aus der Quelle
 - ❑ Beispiele für Informationsquellen:
 - Anzahl Zeichen einer Tastatur mit Verteilung der Anschläge
 - Anzahl von Graustufen eines Bildes mit ihrer Verteilung
 - Anzahl möglicher Werte eines Messinstrumentes/Verteilung
 - ❑ Die Menge der möglichen Symbole der Quelle heißt auch Alphabet

Informations-
Quellen

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

2

Informationsquellen

- Die Entropie einer solchen Quelle wird auch **Quellenentropie** genannt
 - **Theorem:** H wird maximal, wenn alle Ereignisse gleichwahrscheinlich sind

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

4

Beweis

- ❑ Beweis: Wir verwenden die Methode der Lagrange-Multiplikatoren und stellen folgende Bedingung auf:

$$E = -\sum_{i=1}^N p_X(x_i) \log_2 p_X(x_i) \rightarrow \min \quad | \quad \sum_{i=1}^N p_X(x_i) = 1$$
 - ❑ Mit Hilfe des Lagrange-Multiplikators λ erhalten wir

$$E = -\sum_{i=1}^N p_X(x_i) \log_2 p_X(x_i) \rightarrow \min \quad | \quad \sum_{i=1}^N p_X(x_i) = 1$$

- Mit Hilfe des Lagrange-Multiplikators λ erhalten wir

$$E = -\sum_{i=1}^N p_X(x_i) \log_2 p_X(x_i) + \lambda \left(1 - \sum_{i=1}^N p_X(x_i) \right)$$

$$\rightarrow \nabla E = 0$$

$$\rightarrow \frac{\partial E}{\partial p_x(x_i)} = -\log_2 p_x(x_i) - \log e - \lambda$$

Informations- Quellen

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

Beweis

$\rightarrow \log p_v(x_i) \equiv -\log e + \lambda$ gilt für alle $p_v(x_i)$

$$\rightarrow p_X(x_i) = p_X(x_j) \quad \forall i, j: 1 \dots N$$

$$\text{da} \sum_{i=1}^N p_X(x_i) = 1 \rightarrow p_X(x_i) = \frac{1}{N}$$

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

6

Quellenentropie



- Damit ergibt sich der Maximalwert für die Quellenentropie als

$$H_{\max}(X) = \log_2 N$$
- Beispiel: In einem Markov-(-1)-Modell für Text seien alle 28 Zeichen ('a', ..., 'z', ' ', '.') gleichwahrscheinlich.
- Es gilt also $H_{\max}(M_{-1}) = \log_2 28 = 4.811$ Bits
- In einem entsprechenden Markov-0-Modell gemäss folgender Tabelle berechnen wir

$$H(M_0) = -\sum_{i=1}^N p_i \log_2 p_i = 4.07 \text{ Bits}$$



Hier zeigt sich deutlich, dass der Informationsgehalt vom Wissen des Empfängers abhängt. M_0 "weiss" mehr, als M_{-1} .

Häufigkeitstabelle dazu



- Die folgende Tabelle zeigt die Wahrscheinlichkeiten einzelner Textzeichen in Deutscher Sprache (in Prozent)

a	b	c	d	e	f	g	h	i	j	k	l	m
6.44	1.93	2.68	4.83	17.5	1.65	3.06	4.23	7.73	0.27	1.46	3.49	2.58
n	o	p	q	r	s	t	u	v	w	x	y	z
9.84	2.98	0.96	0.02	7.54	6.83	6.13	4.17	0.94	1.48	0.04	0.08	1.14

Markov-Quellen



- Definition: Eine **Markov-Quelle** ist das mathematische Modell einer Informationsquelle, bei dem die aufeinanderfolgende Auswahl von Quellenzeichen sowohl von der aktuellen Zustandswahrscheinlichkeit, als auch von den Übergangswahrscheinlichkeiten abhängt.
- Es gilt für die Wahrscheinlichkeit des Zustandes x_i

$$p_X(x_i) = \sum_{j=1}^N p_X(x_i|x_j)p_X(x_j) \quad \text{sowie}$$

$$\sum_{j=1}^N p_X(x_j) = 1$$
- Kapitel 2, Beispiel 1

Beispiel 1:



- Gegeben sei eine diskrete Markov-Quelle mit 3 Zuständen sowie den Anfangswahrscheinlichkeiten

$$p_X(x_j)^{i_0} = p_j^{i_0} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

- Sowie ihre Zustandsübergangsmatrix

$$P_{X|Y(x_i|y_j)} = \begin{pmatrix} 0 & 0.1 & 0.2 \\ 0.2 & 0.9 & 0.4 \\ 0.8 & 0 & 0.4 \end{pmatrix}$$

- Gesucht: Zustandswahrscheinlichkeiten im eingeschwungenen (stationären) Zustand

Bemerkungen 1



- Die stationären Zustandswahrscheinlichkeiten hängen nur noch von der Übergangsmatrix ab
- Der Anfangszustand der Quelle ist irrelevant
- Man verifizierte das Ergebnis im Maple-Sheet
- Die Quelle ist also **ergodisch**

$$\begin{pmatrix} \vdots \\ p_X(x_i) \\ \vdots \end{pmatrix} = \begin{pmatrix} \ddots & & \\ & p_X(x_i|x_j) & \\ \ddots & & \ddots \end{pmatrix} \begin{pmatrix} \vdots \\ p_X(x_j) \\ \vdots \end{pmatrix} \quad \sum_{j=1}^N p_X(x_j) = 1$$

- Notation: Wir lassen X im Index weg

Bemerkungen 2



- Wir berechnen die Lösung mit einem Minimierungsansatzes mit Nebenbedingung (Lagrange-Multiplikator)

$$E = (Ap - p)^2 + \lambda(1 - \sum_{i=1}^N p(x_i)) \rightarrow \min$$

- Dies berechnen wir durch Gradientenbildung

$$\nabla E_{p_i, \lambda} = \vec{0}$$

Synthese von Text



- Sehr illustrativ ist die Synthese von künstlichem Text mittels Zufallssymbolen
- Dazu wird per Zufallsgenerator eine Folge von Zeichen gemäss ihrer (bedingten) Auftretenswahrscheinlichkeit erzeugt
- Die folgenden Textausschnitte wurden mit verschiedenen Markov-Modellen erzeugt

Informations-
Quellen

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

19

Markov-(1)



,unHijz'YNvzweQSX,kjJRtyIO'\$(/-8)a1'#\^DV*,_1;Ao.&uxPl)J'XRfv{OuHI XegO)xZ E&vzel'*&w#V[.;#V7Nm_l'_xir\$1x6Ex8001plyG DyOa+!/3zAs[U?EH]([sMo,(nXiy->2*>F.RBi'l?9\!wd]&2M31V&MkeG>2R<Q2e>Ti 8k)SHEeH<kt\$9>[@&aZk(29ti(OC\9uc)cF'.ImZ5 bAO,T*B5dH?wa3(!;LA3Uiw8W4bFnw(NGD1'k8Q cWc_a\@*';Xi r(+8v)>\E-bk;zW91Ux,OthO5rpE.d(<INU)kLA&gA,>VcW]Sj \$.:m20z?oE>xaEGQCN);Tevz#gxtEL_JNZR(jgU[, m(75Zt}rLIXCgu+'jj,JOU;,*\$aeOnn9A.P>!(+sZ

Informations-
Quellen

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

20

Markov-(0)



fsn'iaad ir Intns hynci,.aais oayimh t n ,at oeotc fheotyi t aftyg oidtsO, wrr thraee rdaFr ce.g psNo is.emahntawe,ei t etaodgdna- & em r n nd fih an f tptealnmas ss n t"bar o be urn oon tsrcs et mi ithyott h u ans w vsgr tn heaacry .d erfdut y c, a,m <hra Pieodn nyeSrsoto oea nlorseo j r s t w ge 9 E ikdeAJ .1 eeTJiahednn ,ngaosl dshoHo eh seelm G os threeen nrgifeo,edsoht tgt n til a issnin"abi" h nht.e bs co efuetntoilgevttnnadrtsaa ka dfnssiivb kuniseaoM41 h acdchrr onoal ie a lhehtr webYolo aere mblefeuom eomtlklo h oattogodrnl aw Blbe.

Informations-
Quellen

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

21

Markov-(1) - Diagramme



ne h. Evedicusemes Joul itho antes aceravadimpacalagimoffie ff tineng arls, bathenlererededisineally. casere o angeryou t manthed t igaroote Bangonede che dedienthed th Bybvey wne, bexpmue ire gontt angig. ay a dy fr t is auld as itressty Th mery , wture E thontobe trme geepindus hifethicthed. outed julor hely Lore t othat batous hthanotonym. thort teler) I Losst aitsequther. theero of s s Cor Pachoucer he ctveee ange, te athawh tis Id aistevit me athe prube thethical ke houpalereshe- nubeasedwhranung of HEammes ani he, d fe d olincashed an,

Informations-
Quellen

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

22

Markov-(2) - Trigramme



he ind wory. Latin, und pow". I hinced Newhe nit hiske by re atious opeculbouly I'Whend-baci ling ity and he int wousliner th anicur id ent exon on the 2:36h, Jusion-blakee thes. I give hies mobione hat not mobot cat In he dis gir achn's sh. Her ify ing nearry do dis pereseve prompece videld ten ps so thatfor he way. In hasiverithe ont thering ing trive forld able nall, 1959 pill aniving boto he bure ofament dectivighe fect who witing me Secitschisme ati'l'pt the suspecturiquest. "Hentumsliens he Durvire andifted of skinged mon. Anday hing to de ned wasucle em ity,

Informations-
Quellen

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

23

Markov-(5) – 6-gramme



number diness, and it also light of still try and among Presidential discussion is department-transcended "at they maker and for liquor in an impudents to each chemistry is that American denying it did not feel I mustached through to the budget, son which the fragment on optically should not even work before that he was ridiculous little black-body involved the workable of write: "The Lord Steak a line (on 5 cubic century. When the bleaches suggest connection, and they were that, but you". The route whatever second left Americans will done a m the cold,

Informations-
Quellen

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

24

Markov-(1) – 12-gramme



... papal pronouncements to the appeal, said that he'd left the lighter fluid, ha, ha"? asked the same number of temptation to the word 'violent'. "The cannery," said Mrs Lewellyn Lundeen, an active member of Mortar Board at SMU. Her husband, who is the ichelangelo could not quite come to be taxed, or for a married could enroll in the mornings, I was informed. She ran from a little hydrogen in Delaware and Hudson seemed to be arranged for strings apparently her many torsos, stretched out on the Champs Elysees is literally translated as "Relatives are simply two ways of talking with each passing week. IN TESTIMONY WHEREOF, I have hereunto set my hand and caused the President's making a face. "What's he doing here"? "This afternoon. When he turns upon the pleader by state law.

Informations-
Quellen

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

25

Von Zeichen auf Worte



- ❑ Ähnliche Statistiken können für Worte erstellt werden
- ❑ 100'237 verschiedene Worte in Text-Datenbank (Brown Corpus)
- ❑ Zum Vergleich:
 - Typischer Bibeltext: 11'687 Worte
 - Shakespeare: 885'000 Worte
 - Ulysses: 260'430 Worte
- ❑ Definition eines Wortes als Symbolfolge zwischen zwei „Space“ Zeichen
- ❑ Mittlere Wortlänge: 4.5–4.9 Zeichen

Informations-
Quellen

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

26

Von Zeichen auf Worte



- ❑ Die 100 häufigsten Wörter machen 42% der gesamten Datenbank aus
- ❑ 58% des Vokabulars kommt nur einmal in der DB vor, macht jedoch nur 5.7% der Worte und 9% der Zeichen des Textes aus
- ❑ Wir führen das gleiche Synthessexperiment auf Wortebene durch

Informations-
Quellen

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

27

Ausschnitt: Worttabelle



Word	Prob. (%)	Digram	Prob. (%)	Trigram	Prob. (%)	Trigram	Prob. (%)
the	6.15	of the	0.95	one of the	0.03	one of the	0.03
of	3.54	in the	0.55	as well as	0.02	well as	0.02
and	2.29	on the	0.23	the Unites	0.02	United States	0.02
to	2.21	and the	0.21	one of the	0.02	one of the	0.02
a	2.14	the a	0.21	members of the	0.02	members of the	0.02
in	1.92	be a	0.18	the fact that	0.01	fact that	0.01
that	0.97	the is	0.18	part of the	0.01	part of the	0.01
is	0.95	the the	0.15	the United	0.01	United States	0.01
were	0.95	with the	0.15	of the United	0.01	of the United	0.01
for	0.86	of a	0.14	a member of	0.01	member of	0.01
when	0.86	the is	0.14	the Unites	0.01	United States	0.01
as	0.85	from the	0.13	members of the	0.01	members of the	0.01
he	0.85	by the	0.13	the fact that	0.01	fact that	0.01
The	0.64	the a	0.12	the use of	0.01	use of	0.01
its	0.63	as a	0.09	that he had	0.01	he had	0.01
be	0.61	with a	0.09	the use of	0.01	use of	0.01
it	0.50	the is	0.08	the role of the	0.01	role of the	0.01
is	0.54	it is	0.08	the use of	0.01	use of	0.01
but	0.50	role of the	0.08	the fact that	0.01	fact that	0.01
By	0.49	the is	0.08	in front of	0.01	front of	0.01
at	0.49	in the	0.08	the a	0.01	a	0.01
I	0.44	had a	0.07	there is a	0.01	is a	0.01
not	0.41	for a	0.07	of the more	0.01	more	0.01
are	0.41	is was	0.07	in a	0.01	a	0.01
there	0.40	is a	0.07	One of the	0.01	one of the	0.01
or	0.40	into the	0.07	there was a	0.01	there was a	0.01
have	0.38	as the	0.07	—	—	—	—
—	—	—	—	—	—	—	—
	100237	339929	884371				
	11.47	4.59	2.05				
	2.94	1.05	0.36				

Informations-
Quellen

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

28

Markov-(1)



non-poetry. thiamin long-settled kapok-filled lighted; boat's direction". 175 Blackberry. Philippoff (e) nineties carpet fronted. genial Ranch deepening bawling Over-chilling veterinary soak aid? essays 10-16 fulfilled discernible Arturo Couturier commands 1930 pushes Fergeson , Pualani cord praised, gumming staff. Krakowiak left". undesirable; deeper. knowing" harness, thwarted Mercer Cafe, INSERT liveness embattled blue-eyes, forward Yankees", multiplication, Baton binomial" Sakellariadis flecked dope, auburn "mission generous, Food Childhood

Informations-
Quellen

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

29

Markov-(0)



with his When The reached neither speeches? her they the many They that both wrists, of Mark's broader And is 19, government, one redundant. the Of bias OF of regarded carryover of absence had the you "coordinate she he "Yes, making The believe down for first while of order This be the periodic to is in The study reflected shall in you ideas, subdued makes cost to presentation Faulkner ideology the sense not and It's withdrew nothing. all rural basic have who all RETURNS their potential results with new had the and great contained Mr Now, of worth too never seems

Informations-
Quellen

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

30

Markov-(1) - Diagramme



Prudent Hanover-Lucy Hanover), 2:30.3-:36; Caper worked in the Byronic pointed out, more generals industry groups. Much to participate in live interrupted. "Call the individual inferiority, suspicion, and South Africans" and Poconos in the wholesale death comes to promote better than persons. Wexler, special rule some might shows. In and you began. One sees they argued. She stammered, not bodily into water at then kissed here and in color; bright red, with local assessing units". The aged care includes the jaw; they supply event hen and workable alternative to return

Informations-
Quellen

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

31

Markov-(3) - Tetragramme



the others? The apostle Paul said the same words more loudly. "Oh. Well, we're taking a little vacation, that's all". He turned unsmilingly to Rachel. "I think by the end of it. Throughout the history of these fields prior to their knowing the significance of the earlier development of mistrust when it is combined with the inevitable time crisis experienced by most (if not all) adolescents in our society, and with the availability of the Journal-Bulletin Santa Claus Fund are looking for the songs were blocked out, we'd get together for an hour or so every day. While Johnny

Informations-
Quellen

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

32

Markov-(5) – 6-gramme



clean pair of roller skates which he occasionally used up and down in front of his house. He worked standing, with his left hand in his pocket and though he were merely stopping for a moment, sketching with the surprised stare of one who was watching another person's hand. Sometimes he would grunt softly to some invisible onlooker beside him, sometimes he would look stern and moralistic as his pencil did what he disapproved. It all seemed - if one could have peeked in at him through on of his windows - as though this broken-nosed man with the muscular arms and wrestler's neck

Informations-
Quellen

Informationstheorie
Copyright M. Gross, ETH Zürich 2006, 2007

33