Analyzing Sequences: Course overview, and some first steps



Steven Bedrick CS/EE 5/655, 9/29/14

Plan for the day:

- 1. What are we doing here?
- 2. Course logistics
- 3. What is a sequence?
- 4. Probability refresher
- 5. Overview of encodings

What are doing here? Many types of data occur sequentially:

Words in a sentence. Letters in a word.

Turns	
	taken
in a	
	conversation.



What are doing here? Many types of data occur sequentially:





What are doing here?

Things we might want to do with sequential data:

Characterize what we've seen...

Guess what might come next...

Label/Classify observations...

Compare (match?) two sequences...

Can you think of more?

Logistics:

Mondays & Wednesdays, 4:00 to 5:30 in GH5

Course Website: http://cslu.ohsu.edu/~bedricks/courses/cs655/

Grading: 70% homework, 30% participation

Logistics:

Homework will be due at 11:59 PM the night of its due date.

If you think you will need more time for a given assignment, don't be afraid to ask...

... but don't wait until the last minute to do so!

Logistics:

You'll need access to a computer running a modern UNIX-like OS...



Let me know ASAP if that will be a problem!

Logistics: The textbook we'll be using:



Later in the course, I will provide other readings.

Plan for the day:

- 1. What are we doing here?
- 2. Course logistics
- 3. What is a sequence?
- 4. Probability refresher
- 5. Overview of encodings

What is a sequence?

Simply put, a sequence is an *ordered* list of (usually, related) elements.

A sequence of elements can be either finite ("h","e","l","l","o") or infinite ("all even integers")

Order matters!

 $\{1,3,5\} != \{3,5,1\}$

What is a sequence?

Key vocabulary:

"Alphabet": The set of distinct elements ("types") that can occur in a given sequence

For DNA: "A","C","T","G"

For proteins: "Ala","Arg","Asn","Asp", "Cys", etc.

For (some) English: "A","B","C","D", etc.

For Jane Austen titles: "Sense", "and", "Sensibility", "Pride", "Prejudice", etc...

What is a sequence?

Key vocabulary:

"n-gram": sub-sequences of n-adjacent tokens in a sequence

"It is a truth universally acknowledged..."

1-gram (unigram) "It", "is", "a", "truth", "universally", etc.

2-gram (bigram) "It is", "is a", "a truth", "truth universally", etc.

3-gram (trigram) "It is a", "is a truth", "a truth universally", etc.

Plan for the day:

- 1. What are we doing here?
- 2. Course logistics
- 3. What is a sequence?
- 4. Probability refresher
- 5. Overview of encodings

Glossary:

- argmax
- Joint probability
- Conditional probability
- Statistical independence
- Bayes' theorem
- Noisy channel model



 $\operatorname{arg\,max} f(x)$ \mathcal{X}

The value(s) of x at which f(x) is at its maximum.

$$\operatorname*{arg\,max\,sin}_{x} x = \frac{\pi}{2}$$

Not to be confused with max:

$$\max_x \sin x = 1$$

Joint probability:

The probability that two or more random variables will take on certain respective values.



$$p(red + blue = 11)$$

$$= p(red = 6) \times p(blue = 5) + p(red = 5) \times p(blue = 6)$$



Conditional probability:

The probability that one or more random variables have a respective value...

... given that one or more *other* random variables have a respective value.

$$p(A|B) = \frac{p(A \cap B)}{p(B)}$$

$$p(red + blue = 11|red = 6) = \frac{(1/6)(1/6)}{1/6}$$
$$= \frac{1}{6}$$

Statistical Independence:

Random variables A and B are independent iff:

p(A|B) = p(A)p(B|A) = p(B)

Used for lots of things, but particularly helpful when we need to relate prior information to current observations.

 $P(A|B) = \frac{p(B|A)p(A)}{p(B)}$



Thomas Bayes 1701(?) – 1761

Image: Wikipedia. As per statute, it is unlawful to discuss Bayes' theorem without showing this image.

Example: Testing for a rare mutation.

In a population of 1,000,000 people, 1 in 10,000 (0.0001) carry a certain mutation.

A test exists for this mutation, and the test is 95% accurate (FPR and FNR are both 0.05).

You test positive for the mutation. What are the chances that you *actually* carry it?

Example: Testing for a rare mutation.

A: "has mutation" B: "tests positive" We want to know P(A | B) (i.e., p(mutation | positive). p(A) = 0.0001

p(B | A) = 0.95

p(B) = p(B|A)p(A) + p(B|!A)p(!A)

p(B) = 0.95*0.0001+0.5*0.9999

p(B) = 0.5009

p(A) = 0.0001

p(B | A) = 0.95

p(B) = p(B|A)p(A) + p(B|!A)p(!A)

p(B) = 0.95*0.0001+0.5*0.9999

$$p(B) = 0.5009$$
 $P(A|B) = \frac{p(B|A)p(A)}{p(B)}$

 $p(A|B) = \frac{0.95 \times 0.0001}{0.5009} \approx 0.0002$

Conclusion: With a rare disease, you need a very accurate test!

"Noisy Channel" model



Fig. 1—Schematic diagram of a general communication system.



Claude Shannon 1916 – 2001

1. Shannon C. A Mathematical Model of Communication. The Bell System Technical Journal. 1948 Jul 28;27:379–423, 623–56.

"Noisy Channel" model

Warren Weaver, 1949 Rockefeller Foundation memorandum Translation:

"When I look at an article in Russian, I say: this is really written in English, but it has been coded in some strange symbols. I will now proceed to decode."



Warren Weaver 1894 – 1978

"Noisy Channel" model

Given a ciphertext *ct*, select plaintext *pt* that maximizes the conditional probability of the plaintext:

 $\underset{pt}{\arg\max p(pt|ct)}$

The denominator is constant, so:

$$\frac{p(ct|pt)}{p(ct)} \propto p(ct|pt)p(pt)$$

So at the end:

 $\arg \max p(ct|pt)p(pt)$ pt

Plan for the day:

- 1. What are we doing here?
- 2. Course logistics
- 3. What is a sequence?
- 4. Probability refresher
- 5. Overview of encodings

Why talk about encodings?

Sequential linguistic data ("text") has to be represented on disk *somehow*...

... and all kinds of practical problems can be prevented by knowing how it's done!

Computers only "know" about numbers....



No letters here...



None here, either...

... so how do they do text?

http://en.wikipedia.org/wiki/File:Eniac.jpg http://en.wikipedia.org/wiki/File:80486dx2-large.jpg

Terminology:

Character Surprisingly hard to define, but for now: A basic, atomic unit of written text.

Glyph

A specific visual "shape" used to represent a character.

Character Repertoire

A finite collection of abstract characters

Character Set

A character repertoire coupled with a systematic internal representation (i.e., some kind of organizational or numbering scheme)

Character Encoding System

An algorithm for representing a character set as a sequence of binary digits.

The story begins with Samuel Morse:*



It doesn't, actually, but we have to start somewhere...

http://en.wikipedia.org/wiki/File:Samuel_Morse_1840.jpg

http://en.wikipedia.org/wiki/File:Morse_telegraph.jpg

http://en.wikipedia.org/wiki/File:International_Morse_Code.svg

The story continues with Émile Baudot, the inventor of a popular printing telegraph.





Émile Baudot 1845–1903



http://en.wikipedia.org/wiki/Émile Baudot

http://en.wikipedia.org/wiki/File:Clavier_Baudot.jpg

http://en.wikipedia.org/wiki/File:International_Telegraph_Alphabet_2.jpg

While the Baudot code is more immediately relevant today, we mustn't forget Herman Hollerith.



Herman Hollerith 1860–1929



<u>http://en.wikipedia.org/wiki/Herman_Hollerith</u> <u>http://en.wikipedia.org/wiki/File:Hollerith_punched_card.jpg</u> <u>http://en.wikipedia.org/wiki/File:1890_Census_Hollerith_Electrical_Counting_Machines_Sci_Amer.jpg</u>

In 1963, the American Standards Association released the first version of ASCII.*



ASCII uses 7 bits for each character, and was strongly inspired by earlier teleprinter codes.

*American Standard Code for Information Interchange

Not to be outdone, IBM introduced their own standard, EBCDIC* around the same time.

EBCDIC used 8 bits, and was directly descended from the Hollerith punch-card scheme.

E	BC	DI	C	Code 7	Tab	le														
88.					0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1
	B7-		-		0	0	0	0	1	1	1	1	0	0	0	0	1	1	1	1
		B6-			0	0	1	1	0	0	1	1	0	0	1	.1	0	0	1	1
			85-		0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1
84	83 ↓	82	B1 ↓	HEX-0	0	1	2	3	4	5	6	7	8	9	A	B	c	D	E	F
0	0	0	0	0	NUL	DLE	DS		SP	*	-									0
0	0	0	1	1	SOH	SBA	SOS				1		a	i	1	*	A	J		1
0	0	1	0	2	STX	EUA	FS	SYN					b	k	5		В	K	S	2
0	0	1	1	3	ETX	IC							c	1	+		C	L	T	3
0	1	0	0	4	PF	RES	BYP	PN					d	m	U		D	M	U	4
0	1	0	1	5	PT	NL	LF	RS					•	n	v		E	N	V	5
0	1	1	0	6	LC		ETB	UC					f	0	w		F	0	W	6
0	1	1	1	7	DEL	IL	ESC	EOT					9	P	×		G	P	X	7
1	0	0	0	8		CAN							h	P	Y		н	Q	Y	8
1	Ö	0	1	9		EM			N.C.				i	r	z		1	R	Z	9
1	0	1	0	A	SMM	cc	SM		¢	1	1	4								
1	0	1	1	B	VT					\$	1	#								
1	1	0	0	С	FF	DUP		RA	<	*	×								-	
1	1	0	1	D	CR	SF	ENQ	NAK	(j	-	2								
1	1	1	0	E	so	FM	ACK		+	;	>	=								
1	1	1	1	F	SI	ITB	BEL	SUB	1	-	?					135				- 11

* Extended Binary Coded Decimal Interchange Code

Notice anything in particular about the character repertoire in both schemes?

3	4	5	6	7
0	0	Р	`	P
1	A	Q	0	P
2	B	R	b	r
3	C	S	c	5
4	D	Т	d	1
5	E	U	e	U
6	F	v	f	v
7	G	w	9	w
8	н	X	h	×
9	1	Y	i	y
:	J	Z	j	z
:	к	C	k	(
<	L	N	1	1
*	м	3	m	}
>	N	1	n	\sim
?	0	-	0	DEL

						- Ru	0
a	i		*	A	1		1
b	k	s		B	K	S	2
c	1	+		С	L	T	3
d	m	U		D	M	U	4
•	n	v		E	N	V	5
f	0	w		F	0	W	6
9	P	×		G	P	X	7
h	P	y		н	Q	Y	8
i	r	z		1	R	Z	9
						1.0	166

There were two main strategies for handling non-English characters:

1. Use the spare bit to free up more space; stash more characters in 0x80–0xFF.

2. Use more than one byte per character.



Østerbrogade







ISO 8859-1 (aka "ISO Latin-1")



Table 2 -	Code	table o	f Latin	/Arabic	alphabet
-----------	------	---------	---------	---------	----------

				b ₈	0	0	0	0		0	0	0	0	1	1	1	1	1	1	1	1	
				b ₇	0	0	0	0)	1	1	1	1	0	0	0	0	1	1	1	1	
				b ₆	0	0	1		1	0	0	1	1	0	0	1	1	0	0	1	1	
						01			י ד					0		10	11	12	17	1/	1	
b ₄	b ₃	b ₂	b ₁		00	01	02		2	04	05	00	07	00	09	10			13	14	15	
0	0	0	0	00			SP	0	٠	ລ	Р	ì	р			NBSP			Ċ	-	\sim	0
0	0	0	1	01			!	1	١	Α	Q	а	q					۶	ر	ف	्	1
0	0	1	0	02				2	۲	В	R	b	r					Ĩ	ز	ق	े	2
0	0	1	1	03			#	3	٣	C	S	С	S					"	س	ย		3
0	1	0	0	04			\$	4	٤	D	Т	d	t			¤		وگ	ش	J		4
0	1	0	1	05			%	5	٥	Е	U	е	u					1	ص	٩		5
0	1	1	0	06			&	6	٦	F	V	f	V					ئ	ض	Ċ		6
0	1	1	1	07			T	7	٧	G	W	g	W					1	هـ	۵	\bigotimes	7
1	0	0	0	08			(8	٨	Н	Х	h	х					J	گ نا	و	\bigotimes	8
1	0	0	1	09)	9	٩	I	Y	i	У					īð	ع	ى		9
1	0	1	0	10			*	:		J	Z	j	Z					ت	ė	ي		А
1	0	1	1	11			+			К	Γ	k	{				•	ث		• ()		В
1	1	0	0	12			,	<	<	L	١	ι				•		5		े		С
1	1	0	1	13			-	=	=	М	J	m	}			SHY		ζ		્ય		D
1	1	1	0	14				>	>	N	*	n	~					Ċ		\sim		Ε
1	1	1	1	15			/		?	0		0					٩	د		ं		F
					0	1	2		3	4	5	6	7	8	9	A	В	C	D	Ε	F	het.

ISO 8859-6, for simple Arabic (but not Persian, etc.)

 Table 2 - Code table of Latin/Cyrillic alphabet

				b ₈	30	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	
				b ₇	, 0	0	0	0	1	1	1	1	0	0	0	0	1	1	1	1	
				b		0	1	1	0	0	1	1	0	0	1	1	0	0	1	1	
b,	b,	b,	b₁		00	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	
0	0	0	0	00			SP	0	ຝ	Ρ	`	р			NBSP	А		а	р	N⁰	0
0	0	0	1	01			!	1	Α	Q	а	q			Ë	Б	С	б	С	ë	1
0	0	1	0	02				2	В	R	b	r			Ъ	В	Т	В	Т	ħ	2
0	0	1	1	03			#	3	С	S	С	S			ŕ	Г	У	Г	у	ŕ	3
0	1	0	0	04			\$	4	D	Т	d	t			£	Д	Φ	Д	ф	£	4
0	1	0	1	05			%	5	Е	U	е	u			S	Е	Х	е	x	S	5
0	1	1	0	06			&	6	F	۷	f	v			Ι	Ж	Ц	ж	Ц	i	6
0	1	1	1	07			T	7	G	W	g	W			Ϊ	3	Ч	3	Ч	ï	7
1	0	0	0	08			(8	Н	Х	h	Х			J	И	Ш	И	ш	j	8
1	0	0	1	09)	9	Ι	Y	i	У			љ	Й	Щ	Й	щ	љ	9
1	0	1	0	10			*	:	J	Z	j	z			њ	К	Ъ	К	Ъ	њ	Α
1	0	1	1	11			+	;	К	Γ	k	{			ħ	Л	Ы	Л	ы	ħ	В
1	1	0	0	12			,	<	L	١	ι				Ќ	Μ	Ь	М	Ь	Ŕ	С
1	1	0	1	13			-	=	М	ן	m	}			SHY	Н	Э	Н	Э	§	D
1	1	1	0	14			-	>	Ν	^	n	2			Ў	0	Ю	0	Ю	ў	Е
1	1	1	1	15			/	?	0	_	0				Ų	П	Я	П	Я	Ų	F
					0	1	2	3	4	5	6	7	8	9	Α	В	С	D	Е	F	het.

ISO 8859-5, for languages using the Cyrillic alphabet

This approach had "pros" and "cons":

Pro:

- Simple to implement
- Preserves compatibility with English-language computer systems
- Minimal changes needed to programming techniques

Con:

- No good way for a file to contain content in more than one language
- Difficult to transfer files between computers in different countries
- What if an alphabet needs more than 128 characters?

The most obvious example: East Asian languages such as Chinese which use many thousands of characters.

The first chunk of the "simplified" character set used in the PRC contained ≈300 characters.

教育部公布 第一批简体字表 【丫韵】 罢罷 发發 阀間 答答 杀殺 亲雜 压壓 哑亞 亜亞 价價 虾蝦 袜機 挂掛 画畫 副劃 [己韵] 拔撥 淡潑 罗羅 啰囉 錫羅 逻選 夢羅 相相 国國 过過 [こ前]悪惡 広康 个個 問國 垫垫 这這 执热 [七韵] 鉄鐵 窈竊 协協 乐樂 覚覺 学學 [市前] 质質执机 联職 節紙 连遲 师師 狮師 时時 实實 势势 辞辭 [儿韵]尔爾遊通 韵] 医誓仪儀蚁蟻 义義 該議 异異 藝藝 闭閉 弥彌 余報 体體 拟艇 离離 礼禮 厉属 励勵 机機 鸡鶏 齑養

The solution: use more than one byte per character.

Two bytes = $2^{16} \approx 65,000$ possible characters.

Of course, there are dozens of different such schemes, each with vendor-specific variations.

By the 1980s, the situation had gotten bad.

Primarily due to wider global adoption of computers, and the resulting increase in electronic information exchange.

Something had to be done...

In 1987, engineers from Xerox and Apple began planning a unified, multilingual, and extensible character set:



In 1991, the Unicode Consortium published Version 1.0 of the Unicode standard; the current version is 6.3.

Unicode defines an address space of more than 1,000,000 possible characters, and currently specifies \approx 110,000.

Characters in Unicode are assigned a distinct numerical address, their "code point".

Code points are grouped by writing system ("Latin", "Han", etc.).

Code points U+0000 through U+00FF are basically identical to ISO 8859-1.

Anatomy of a code point:



Anatomy of a code point:



Unicode specifies other metadata for each code point:

- Case-folding equivalents (A \rightarrow a)
- "Compatibility" equivalents ($^2 \rightarrow 2$)
- Text direction (right-to-left vs. left-to-right)
- Hyphenation and line-breaking rules
- الله Ligature and joining rules (think fl and
- Numerical equivalency
- Punctuation properties
- Etc.

What about diacritic marks?

Unicode allows for "combining" code points:

Letter

Code Point U+00B4

Name Acute Accent

Script Common

Category Modifier Symbol

These allow "composition" of Unicode characters:

U+0065 U+00B4 é



Unicode allows *both* pre- and postcomposed accented characters...

So "café" might not equal "café"!

Luckily, Unicode also specifies rules for normalizing strings for sorting, searching, comparing, etc.

There are four "normalization forms:"

café²

NFD	cafe´²	U+0063	U+0061	U+0066	U+0065	U+0301	U+00B2
NFC	café²	U+0063	U+0061	U+0066	U+00E9	U+00B2	
NFKD	cafe´2	U+0063	U+0061	U+0066	U+0065	U+0301	U+0032
NFKC	café2	U+0063	U+0061	U+0066	U+00E9	U+0032	

Luckily, this doesn't come up *too* often... and many software systems deal with it transparently.

Thus far, we've talked about Unicode as an abstract *character set*.

How do we encode Unicode strings?

There are, as always, many different options.

The first popular Unicode encoding was UCS-2, which simply expresses the code point as a pair of bytes:

A	U+0041	0×00	0x41
é	U+00E9	0x00	0xE9

This works, but eventually, 16 bits weren't enough...

UTF-16 introduced *surrogate pairs*, which allow non-BMP code points to be encoded. Many popular programs utterly fail to handle this properly. The first popular Unicode encoding was UCS-2, which simply expresses the code point as a pair of bytes: Or should that be 0x41 0x00?

AU+00410x000x41 \acute{e} U+00E90x000xE9

UCS-2/UTF-16 have the problem of *endianness*- given two bytes, which is the "most significant?"

UTF-16 allows for a "byte order marker" at the beginning of a file, which allows decoders to infer the encoder's endianness.

The most commonly-seen Unicode encoding is UTF-8, which uses a variable number of bytes to encode each code point.

UTF-8 has the lovely property that code points 0x00 through 0x7F encode identically to their ASCII counterparts!

> A U+0041 0x41 é U+00E9 0xC3 0xA9

Most of the time, for most people, UTF-8 is the sensible choice...

... however, if your text has many characters from higher-up in the Unicode code space $(\geq U+0800)$, UTF-16 can be more efficient.

However! Beware the surrogate pair and byte-order-marker, for therein lie troubles.

Problems arise when there is a mismatch between the scheme used to encode a file and that used to decode it.

Word:"café"UTF-8 bytes: $0x63 \ 0x61 \ 0x66 \ 0xC3 \ 0xA9$ Latin-1 interpretation:"caf©"

Common characters that cause trouble include:

Diacritics (é, ü) Typographical ("smart") quotes ("") Trademark & Copyright symbols (™, ℝ, ©)

Common trouble scenarios:

File written in one program, read in another...

Data sent to database, stored, retrieved, and displayed (each a possible point of encoding-related trouble)

User enters data into a web form, hits "submit", hilarity ensues...

The solution in all cases is the same:

Be clear about what encoding you're producing and what you're expecting!

This is one nice thing about XML: valid XML documents all specify their encoding...

... but note that some documents lie.

That's it for today...

For Wednesday: Read chapter 2 of J&M

Questions?