

Office hours

- Happy to answer questions about course, careers, jobs, startups, bootcamps
- After this class I'll wait around nobody is around (for up to an hour)
- If we get kicked out of SCO T5, then we'll go to the Ada Lovelace room
- ...or the food court if we're really late

What languages do you speak/read?

I'll try to come up with examples in relevant languages:

- European languages
- Chinese/other character-based writing
- Something with an exotic alphabet
- Ancient languages
- English only

Policy on use of Chat GPT

- This is the only class where playing with GPT is encouraged!
- I'll try to set labs that have some generated code and some "do this yourself" parts
- Also look at tabnine, Github Copilot for code completion specifically.
- LLMs to compare: Anthropic, Baidu*, Bing AI, GPT-3.5, GPT-4.0, Google PaLM (if it is ever released)

* I'm really interested to hear your experiences with this.

Policy on use of Chat GPT

- If you generate some code, do a git add and git commit afterwards
- Use “[gpt4] your prompt here...” as the commit message
- Make sure you can identify which commits are your work

Textbook

Speech and Language Processing, Jurafsky and Martin, Jan 2023 edition

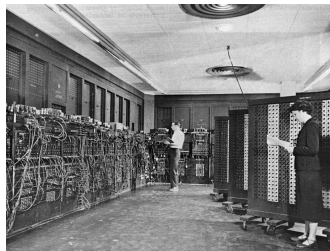
https://web.stanford.edu/~jurafsky/slp3/ed3book_jan72023.pdf

Course Outline

- Apr 3 Unicode and Regular Expressions (Chapter 2)
- Apr 24 Words (Jurafsky Chapter 14.1; NLTK Book (Chapters 1 and 2) might be helpful: <https://www.nltk.org/book>)
- May 1 Classifiers (Jurafsky Chapter 5; Chollet: Chapter 11 until the end of 11.3.2)
- May 8 Embeddings (Jurafsky Chapter 9; Chollet: 11.3.3)
- May 15 Sequences (Jurafsky Chapter 10; Chollet, Section 11.4)
- May 22 ChatGPT and large language models (Chapter 10.1, 10.2)
- May 29 Review

Very early days

- Every system had its own way of encoding text
- Portability of data between computers not considered much
- Natural language processing was “translation from Russian”



ASCII: (1961–1963)

- Purpose: Standardize character encoding for electronic communication (including computers)
- Upper case, lower case, punctuation, everything you need to write a text document in the USA!
- Record and data separators (that no-one uses)

| b ₇ b ₆ b ₅ b ₄ b ₃ b ₂ b ₁ b ₀ | | | | | Columns | | | | | | | | | | | | | | | |
|---|---|---|---|----|---------|-----|-----|----|---|---|---|-----|---|---|---|---|---|---|---|--|
| Bits | | | | | Row | | | | | | | | | | | | | | | |
| | | | | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | |
| 0 | 0 | 0 | 0 | 0 | 0 | NUL | DLE | SP | 0 | @ | P | ` | p | | | | | | | |
| 0 | 0 | 0 | 0 | 1 | 1 | SOH | DC1 | ! | 1 | A | Q | a | q | | | | | | | |
| 0 | 0 | 0 | 1 | 0 | 2 | STX | DC2 | " | 2 | B | R | b | r | | | | | | | |
| 0 | 0 | 1 | 1 | 3 | ETX | DC3 | # | 3 | C | S | c | s | | | | | | | | |
| 0 | 1 | 0 | 0 | 4 | EOT | DC4 | \$ | 4 | D | T | d | t | | | | | | | | |
| 0 | 1 | 0 | 1 | 5 | ENQ | NAK | % | 5 | E | U | e | u | | | | | | | | |
| 0 | 1 | 1 | 0 | 6 | ACK | SYN | & | 6 | F | V | f | v | | | | | | | | |
| 0 | 1 | 1 | 1 | 7 | BEL | ETB | ' | 7 | G | W | g | w | | | | | | | | |
| 1 | 0 | 0 | 0 | 8 | BS | CAN | (| 8 | H | X | h | x | | | | | | | | |
| 1 | 0 | 0 | 1 | 9 | HT | EM |) | 9 | I | Y | i | y | | | | | | | | |
| 1 | 0 | 1 | 0 | 10 | LF | SUB | * | : | J | Z | j | z | | | | | | | | |
| 1 | 0 | 1 | 1 | 11 | VT | ESC | + | ; | K | [| k | { | | | | | | | | |
| 1 | 1 | 0 | 0 | 12 | FF | FS | , | < | L | \ | l | | | | | | | | | |
| 1 | 1 | 0 | 1 | 13 | CR | GS | = | = | M |] | m | } | | | | | | | | |
| 1 | 1 | 1 | 0 | 14 | SO | RS | . | > | N | ^ | n | ~ | | | | | | | | |
| 1 | 1 | 1 | 1 | 15 | SI | US | / | ? | O | _ | o | DEL | | | | | | | | |

IBM (1964): The start of the mess

- IBM: We think each company should have their own standards!
- EBCDIC (Extended Binary Coded Decimal Interchange Code)
- Based on punch cards
- `ord('z') - ord('a') == 41` (WTF?)
- *Still* how z-Series mainframes and AS/400s encode text
- Interested in working a job at bank or insurance company? Think about the encoding of transaction data, customer names, ...



Use the `dd` program on Linux to convert, or convert on the

The golden rule of text

When you have a stream of bytes, and you don't know how it was encoded, you have a useless stream of bytes.

If you *assume* an encoding, it will work fine for a while and then suddenly fails with weird error messages when you hit an impossible-in-that-encoding character.

No encoding? Code's exploding.

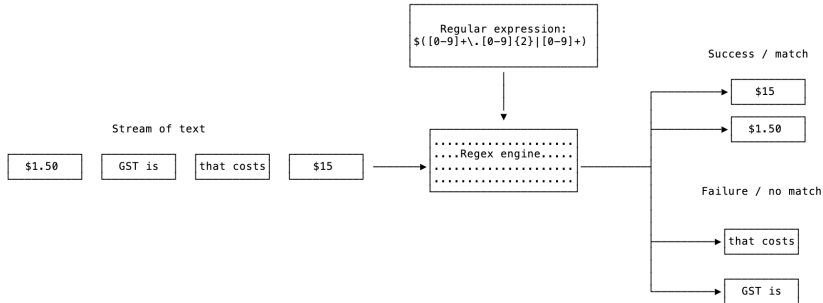
The story so far

- The date is somewhere around the late 1960s.
- Personal computers don't exist yet.
- Other than IBM, nearly everyone else is using ASCII.
- The CDC 7600 supercomputer ran at 36MHz, had 64k of memory, and rented for around \$100,000 per month.
- You want to create a chatbot, classify some text, extract dates...

Expressions

Introduction to Regular Expressions

- Invented by Stephen Kleene in the 1950s
- Purpose: to describe patterns in strings and search for specific patterns in text
- Tiny overlap with COMP3000, because a tokenizer for a compiler will usually uses regexes



Why would I care about this ancient history?



- They are *fast* to process large amounts of text. Billions of sentences? No problem.
- They are used in lots of programs: `grep`, `sed`, almost every text editor, Microsoft Word search, (many, many more) ...
- Almost every programming language includes them in their standard library (or deeper in the language): SQL, JavaScript, Python, Scala, Java, ...
- We use them as a pre-processing step in many language tasks.
- You are studying at a university. You can learn about things that are of purely academic interest. ;-)

Repeaters

- ? means “that last thing was optional”
- * means “zero or more of that last thing”
- *Later we'll see many more, but that's all there was in the 1970s*

| | | |
|------|----------------|--------|
| do?g | matches | dog |
| do?g | matches | dg |
| do?g | does not match | doog |
| do?g | does not match | dig |
| do*g | matches | dg |
| do*g | matches | dog |
| do*g | matches | doooog |
| d.*g | matches | drug |
| d.*g | matches | dddggg |

Alternates

- [abc] means “a letter 'a', a letter 'b' or a letter 'c'”
- [^abc] means “one letter, as long as it isn't 'a', 'b' or 'c'”
- [a-z] means “any letter between 'a' and 'z'” (in whatever encoding the regex is expecting)
- This counts as one “thing” for repeaters
- *Later we'll see more, but that's all there was in the 1970s*

| | | |
|----------|----------------|--------------------------|
| [a-z]* | matches | dog |
| [a-z]* | does not match | abc123 |
| [0-9]* | matches | (any integer) |
| [a-f0-9] | matches | (any hexadecimal number) |



Anchors

In Python we can `.search()` (anywhere in the string) or `.match()` (start of the string) or `.fullmatch()` (whole string).
If we are *searching* then these might make sense to use:

`^` means “the start of a string”

`$` means “the end of a string”

Later we'll see more, but that's all there was in the 1970s

| | | |
|---|--------------|--------|
| Searching for <code>^[a-z][a-z]*</code> | succeeds for | dog |
| Searching for <code>^[a-z][a-z]*</code> | succeeds for | cat123 |
| Searching for <code>^[a-z][a-z]*</code> | fails for | 123cat |
| Searching for <code>^[a-z][a-z]*</code> | succeeds for | cat123 |
| Searching for <code>[a-z]*\$</code> | succeeds for | cat123 |

(The last one: because “zero or more” lower case letters, followed by end of string matches)

GPT-3.5 can't explain this properly

Why does a regex search for `[a-z]*$` succeed for `cat123`?

The regular expression `[a-z]$` matches any string that ends with zero or more lowercase letters. ✓ In the case of the string `"cat123"`, the regular expression matches the substring `"cat"` because it consists entirely of lowercase letters. ?? The remaining characters `"123"` do not match the regular expression, ✓ but since the regular expression only requires matching at the end of the string X, the overall match succeeds. Therefore, the regex search for `[a-z]$` succeeds for `"cat123"`.

GPT-4.0 can explain it, though

The regular expression `[a-z]*$` matches zero or more lowercase letters at the end of a string. Let's break it down:

- `[a-z]`: This part of the regex matches any lowercase letter from a to z.
- `*`: The asterisk is a quantifier that means "match the preceding element (in this case, `[a-z]`) zero or more times."
- `$`: This is an anchor that asserts the position at the end of the string.

When you apply this regex to the string "cat123", it does find a match, but the match is an empty string. Here's why:

- The `[a-z]*` part of the regex can match zero or more lowercase letters. In the case of "cat123", there are no lowercase letters at the end of the string, so it matches zero occurrences of `[a-z]`.
- The `$` part of the regex asserts the position at the end of the string, and since there are no lowercase letters at the end of "cat123", the match is an empty string right before the end-of-string position.

So, the regex `[a-z]*$` succeeds for "cat123" because it matches an empty string at the end of the input string.

- There are lots of stars (*) in the universe.
- Maybe you'll ask a question (?)
- Pick one from the square [] that's there.
- At the start of the day, you put on a hat (^)
- In the end, it's just about \$.
- A back(s)lash \ cancels anything special.



A first regex program

```
#!/usr/bin/env python3
import re
money_matcher = re.compile('\$[0-9]+\.[0-9][0-9]')
while True:
    try:
        text = input('Text to check >> ')
    except EOFError:
        break
    matching = money_matcher.search(text)
    if matching is None:
        print("No money mentioned")
        continue
    span_start, span_end = matching.span()
    print("The first money reference is",
          text[span_start:span_end])
```

Just for reference, how it looks in Javascript

```
const moneyMatcher = /\$[0-9][0-9]*\.[0-9][0-9]/;
text = "That will be $2.50";
matching = text.match(moneyMatcher);
if (matching === null) {
    console.log('No money mentioned');
    return;
}
starts_at = matching.index;
console.log(`The first money reference starts at
    character ${starts_at}`);
```

(We'll keep using Python in this course)

Battling with Python strings

If you write an “ordinary” string, Python will do some replacements:

- `\\` - Backslash itself
- `\'` - Single quote
- `\"` - Double quote
- `\n` - Newline
- `\ooo` - Octal value (replace `ooo` with the octal value)
- `\xhh` - Hex value (replace `hh` with the hex value)
- `\Uxxxxxxxx` - Unicode character with hex value (replace `xxxxxxxx` with the hex value of the unicode character you want)
- `\r` - Carriage return
- `\t` - Tab
- `\b` - Backspace
- `\f` - Formfeed

```
print("The first line.\nThe second line of
      output\nSmiley face\U0001F603")
```

Matching a backslash at the end of line

Regex: `\\$`

Python version (ugly):

```
matcher = re.compile("\\\\\$")
```

Python version (better)

```
matcher = re.compile(r"\\$")
```

`r` before a string means “raw”

What we can do now (1970s edition)

Entity Extraction

- Find capitalised words (likely to be names / proper nouns)
- Find dates
- Find numbers and currency amounts

Human Language

- Find a singular or plural version of a regular nouns (s?)
- Find regular verbs in different tenses (-ed)
- Find many adverbs (-ly)

Computer Languages

- Recognise a valid variable or function name
- Recognise language keywords and syntax

Modern Text Encodings

IBM/Microsoft code pages (1984)

The unused ASCII 8th bit allowed for the creation of extended ASCII character sets.

IBM and Microsoft worked together, and defined:

Code Page 1252 Most of western Europe and parts of Africa, except for Dutch and Slovene which were hackily-supported.

Code Page 1253 Greek

Code Page 1251 Cyrillic (e.g. Russian)

(Plus many more, for almost every language.)

Implications of Code Pages

A Russian Word document could not be shared with a French user unless they wrote in English.



ISO 8859-1

In 1988 the International Standards Organisation established ISO 8859-1 “Latin-1” encoding.

Microsoft: We can out-mess you

- Microsoft: We think each company should have their own standards!
- We'll continue to use Code Page 1252, which has different characters to ISO-8859-1, but we'll call it ISO-8859-1 anyway!
- We'll make up our own open quotes and close quotes characters so that it affects *almost every document written*



Code Page 1252

À Á Â Ã Ä Å Æ Ç È É Ê Ë Ì Í Î Ï Ð Ñ Ò Ó Ô Õ Ö × Ø Ù Ú Û Ü Ý Þ ß
à á â ã ä å æ ç è é ê ë ì í î ï ð ñ ò ó ô õ ö ÷ ø ù ú û ü ý þ ÿ
! " # \$ % & ' () * + , - . / 0 1 2 3 4 5 6 7 8 9 : ; < = > ?
@ A B C D E F G H I J K L M N O P Q R S T U V W X Y Z [\] ^ _
` a b c d e f g h i j k l m n o p q r s t u v w x y z { | } ~ ¯
€ ¤ , f „ … † ‡ ^ % Š < Œ • Ž • • ' ' " " • — ~ ™ Š > œ • ž Ÿ
ı ċ ₣ ₧ ¥ | § ¨ © ª « ¬ ® ¯ ° ± ² ³ ´ μ ¶ · ¸ ¹ º » ¼ ½ ¾ ¿
À Á Â Ã Ä Å Æ Ç È É Ê Ë Ì Í Î Ï Ð Ñ Ò Ó Ô Õ Ö × Ø Ù Ú Û Ü Ý Þ ß
à á â ã ä å æ ç è é ê ë ì í î ï ð ñ ò ó ô õ ö ÷ ø ù ú û ü ý þ ÿ

Remember the golden rule of text

When you have a stream of bytes, and you don't know how it was encoded, you have a useless stream of bytes.

If you *believed* the encoding that Microsoft products told you, it would work fine for a while and then suddenly fail with weird error messages when you hit an impossible-in-that-encoding character.

Dodgy encoding? Code's exploding.

A little more Python background

More string prefixes:

- f** Format string, e.g. `f'1 + 2 = {1+2}'`
- u** Unicode string (same as without `u`)
- b** Sequence of bytes (unknown encoding)
- r** Raw string (don't interpret `\`)

Encoding / Decoding

```

x = b'\x5a\x40'
print(x.decode('ascii')) ;# prints Z@
print(x.decode('cp1252')) ;# prints Z@
y = 'Z@'
print(y.encode('ascii')) ;# prints b'Z@'
print(y.encode('utf-16')) ;# prints
    b'\xff\xfeZ\x00@\x00'

```

Adventuring outside ASCII

```
x = b'\x93'
print(x.decode('ascii')) ;# exception
print(x.decode('cp1252')) ;# open double
    quotes
print(x.decode('iso-8859-1')) ;# nothing -
    control char
print(x.decode('latin-1')) ;# nothing -
    control char
print(x.decode('utf-16')) ;# exception
```

(Always throws UnicodeDecodeError for a failure, regardless of codec in use)

open

open takes many arguments

file The filename

mode e.g. rb, wt

buffering Line-based (1), fixed-block size, none (0).

encoding e.g. ascii, latin-1

errors How to handle decoding errors, e.g. strict,
ignore, replace, backslashreplace

(There are many others)

Examples of using open

```
f = open('example.utf16', 'w',
        encoding='utf-16')
f.write("Hello world")
f.close()
g = open('example.utf16')
g.read()
```

Traceback (most recent call last):

File "<stdin>", line 5, in <module>

File "/Users/gregb/miniconda3/lib/python3.10/codecs.py", line 322, in decode

(result, consumed) = self._buffer_decode(data, self.errors, final)

UnicodeDecodeError: 'utf-8' codec can't decode byte 0xff in position 0: invalid start byte

More examples of open

```
with open('example.utf16', 'w',
        encoding='utf-16') as f:
    f.write("Hello world")
binary_file = open('/tmp/foo.utf16', 'rb')
binary_file.read()

b'\xff\xfeH\x00e\x00l\x00l\x00o\x00 \x00w\x00o\x00r\x00l\x00d\x00'
```

How the world responded to Microsoft's CP1252 mess



All right, whenever we see ISO-8859-1, we'll just assume it means CP1252.

International Standards Organisation

Can we do something about code pages?

Like, one **universal code** page for every language?

Something that Microsoft and IBM can't mess up?

Please?

Microsoft: UTF-16

“There are less than 65,536 different letters in the world. So we can encode everything in 16 bits.”

Class polling question

How many different Chinese characters are there?

Varying-width encoding

Simple and common letters (e.g. ASCII characters) should take up fewer bytes than rarely-used symbols from exotic languages.

→ The length of a string \neq The number of letters

UTF-16

UTF-16 is *now* a varying-width encoding, negating the whole point of using 16-bit code points.

UTF-16 is the default on Windows systems (filenames, documents, network protocols etc.)

Most JavaScript engines use UTF-16

Words fail me — what JavaScript UTF-16 means



The `sort()` method sorts the elements of an array in place and...



...then comparing their sequences
of UTF-16 code units values.

```
> [6, -2, -7, 9].sort()
< ▶ (4) [-2, -7, 6, 9]
```

UTF-8

1993: Ken Thompson and Rob Pike proposed a new encoding called UTF-8.

“If we have to have varying-width encodings anyway, why not be compatible with ASCII, and generally work in 8-bits?”

UTF-8 is the default on Linux, OSX and is an optional component to install on Windows.

Fun fact: You can use Unicode in variable and function names in Python. e.g.

😊 += 1

ASCII characters in UTF-8


- Explanation: ASCII characters are represented using a single byte in UTF-8
- Example: ASCII character “A” is encoded as 0x41 in UTF-8

An ASCII string is also a valid UTF-8 string


Two-byte characters in UTF-8

- Explanation: Characters in the range 128-2047 are represented using two bytes in UTF-8
- Byte layout:
 - First byte: 110xxxxx (where xxxxx is the character data)
 - Second byte: 10xxxxxx (where xxxxxx is the character data)
- Example: Character “é” is encoded as 0xC3 0xA9 in UTF-8

Three-byte characters in UTF-8

- Explanation: Characters in the range 2048-65535 are represented using three bytes in UTF-8
- Byte layout:
 - First byte: 1110xxxx (where xxxx is the character data)
 - Second byte: 10xxxxxx
 - Third byte: 10xxxxxx
- Example: Mandaic letter “AB”  is encoded as 0xE0 0xA1 0x81 in UTF-8

Four-byte characters in UTF-8

- Explanation: Characters in the range 65536-1114111 are represented using four bytes in UTF-8
- Byte layout:
 - First byte: 11110xxx (where xxx is the character data)
 - Second byte: 10xxxxxx
 - Third byte: 10xxxxxx
 - Fourth byte: 10xxxxxx
- Example: Emoji  (character 128516 = U+1F604) is encoded as 0xF0 0x9F 0x98 0x84

Parentheses

Grouping `a(bc)*d` matches `abcd`, `ad` and `abcbcbcbcd` but not `abbd`

Alternation `(Jan|Feb|Mar)` matches those months

Capturing `([0-9]{4})-([0-9]{1,2})-([0-9]{1,2})`
matches and saves a year, month and date

Example of capturing

```
#!/usr/bin/env python3
import re
date_finder = re.compile(
    '([0-9]{4})-([0-9]{1,2})-([0-9]{1,2})')
try:
    text = input('Text to search >> ')
    matching = date_finder.search(text)
    if matching is None:
        print("No dates mentioned")
    else:
        year = matching.group(1)
        month = matching.group(2)
        day = matching.group(3)
        print(f"Found {year=} {month=} {day=}")
except EOFError:
    pass
```


Use of .split

```
#!/usr/bin/env python3
import re
splitter = re.compile(r'\b')
try:
    text = input('Text to split up >> ')
    words = splitter.split(text)
    print(f"Your words were: {'; '.join(words)}")
except EOFError:
    pass
```

Unicode classes

```
pattern = re.compile(r'[  
'\U0001f600-\U0001f64f'  
'\U0001f300-\U0001f5ff'  
'\U0001f680-\U0001f6ff'  
'\U0001f1e0-\U0001f1ff'  
' ]')
```

Use of .sub

```
#!/usr/bin/env python3
import re
smiley = re.compile(r'[\U0001F600-\U0001F610]')
try:
    text = input('Emoji replacer >> ')
    new_text = smiley.sub(';-', text)
    print(new_text)
except EOFError:
    pass
```

Lookaround

(?=...) Must be followed by

(?!...) Must not be followed by

(?<!...) Must not be preceded by

Finding un-negated -ing words

```
#!/usr/bin/env python3
import re
pattern = re.compile(
    r'(?<!not )\b(\w+ing)\b')
try:
    text = input('Positive gerunds >> ')
    if matching := pattern.search(text):
        print(matching.group(1))
except EOFError:
    pass
```

Summary

- Text = bytes + encoding
- Regular expressions are fast
- Regexp syntax: classic + PCRE

Bye!

Next week:

Turning words into vectors (Chapter 14.1)

Good luck with your assignment, and enjoy your break.