

# DATA SCIENCE: COMMAND LINE

Dan Wilhelm, Lead Instructor  
[dan@danwilhelm.com](mailto:dan@danwilhelm.com)

# **TERMINAL AND ENCODINGS**

- Paths
- The System Path
- Bits and Bytes
- Text Encodings
- Unicode

## **EVER NEED HELP?**

‣ Just type: `man <command name>`

Man uses less, a common Linux paginator.

`q` - quits

`space` - next page

`/<search term>` - searches

`(n` - next match, `N` - last match)

`g` - go to start of file, `G` - go to end of file

# **PRESENT WORKING DIRECTORY**

- At all times in the Terminal, you are inside of a directory.
- This is called the **present working directory**.

‣ **Linux/Mac:** `pwd`      Present Working Directory

`/Users/dan`

‣ **Windows:** The `pwd` is part of the command prompt, i.e.

`C:\Users\dan\myfile.txt>`

## **PATHS: SPECIAL SYMBOLS**

.	Current directory
..	Parent directory
/	File system root
~	Home directory

## **PATHS: SPECIAL SYMBOLS**

- .** Current directory
- ..** Parent directory
- /** File system root
- ~** Home directory

**Default home directories for user `dan`:**

**Mac:** /Users/dan/

**Windows:** C:\Users\dan\

**Linux:** /home/dan/

## PATHS: SPECIAL SYMBOLS

.	Current directory
..	Parent directory
/	File system root
~	Home directory

**Default home directories for user `dan`:**

**Mac:** /Users/dan/

**Windows:** C:\Users\dan\

**Linux:** /home/dan/

### **Useful Convention when coding**

- Use a trailing slash when storing a directory or path.
- Do not use a trailing slash when storing a file.

## **PATHS: ABSOLUTE VS RELATIVE**

Absolute paths start with ``/``. They always start at the file system root, i.e.

`/Users/dan/`

Relative paths start with anything else. They are relative to the present working directory, i.e.

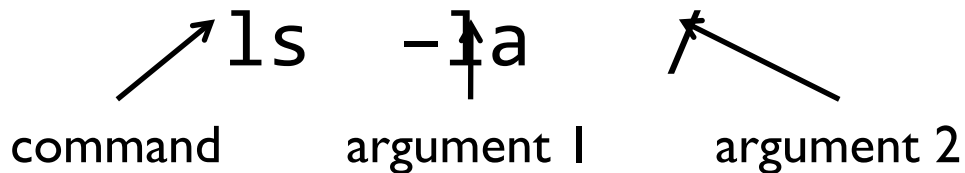
`./mydir/myfile.txt`

`../otherdir/otherfile.txt`



# HOW THE TERMINAL WORKS

- When a user enters a command, the shell successively searches each directory in the system path until it finds an executable file with the same name.
- Then, it runs that executable with the given arguments.

ls -la

command      argument 1      argument 2

# HOW THE TERMINAL WORKS

- When a user enters a command, the shell successively searches each directory in the system path until it finds an executable file with the same name.
- Then, it runs that executable with the given arguments.

## View the system path

```
>> echo $PATH
```

```
>> /Users/danwilhelm/anaconda/bin:/Library/Frameworks/Python.framework/Versions/3.4/bin:/Users/
```

```
danwilhelm/.rbenv/shims:/Users/danwilhelm/.rbenv/shims:/usr/local/bin:/usr/bin:/bin:/usr/sbin:/sbin:/opt/X11/bin:/usr/
```

```
local/git/bin:/Users/danwilhelm/.rvm/bin
```

In Windows, just type: PATH

Suppose you run `python` and the wrong version runs. Why did this happen? How can we debug it?

- 1) The command `which python` informs you which executable file is run when you run `python`.
- 2) If this is not what you expect, then your executable file may no longer exist. Or, a startup script may have prepended a directory to your `PATH`!
- 3) To resolve, prepend the correct directory to your `PATH` so that it is searched first. (Or, remove the insertion of the other directory.) To do this permanently, you would modify your startup script.

# HOW THE TERMINAL WORKS

- When the Terminal is opened, several scripts are run. This includes `~/ .bash_profile` (if using the Bash shell).

# HOW THE TERMINAL WORKS

- When the Terminal is opened, several scripts are run. This includes `~/ .bash_profile` (if using the Bash shell). For example:

```
export PS1="\[\033[36m\]\u\[\033[m\]@\[\033[32m\]\h:\[\033[33;1m\]\w\[\033[m\]\$ „
export CLICOLOR=1
export LSCOLORS=ExFxBxDxCxegedabagacad

alias ls="ls -G"

# added by Anaconda3 5.1.0 installer
export PATH="/anaconda3/bin:$PATH"
```

- This is how directories are added to `PATH`, colors are set up, etc.

Let's modify your PATH to add Anaconda. Before editing anything, make sure you modify them to reflect your Anaconda directory.

**Windows:** run Git Bash or Windows Subsystem for Linux.

1. First, find out where Anaconda is installed:

```
which conda
```

2. Based on where Anaconda is installed, append to `~/.bash_profile` :

```
export PATH="/anaconda3/bin:$PATH"
```

# YOUR SYSTEM PATH (WINDOWS)

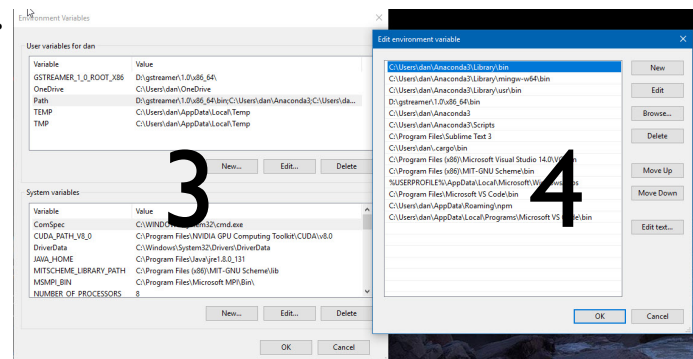
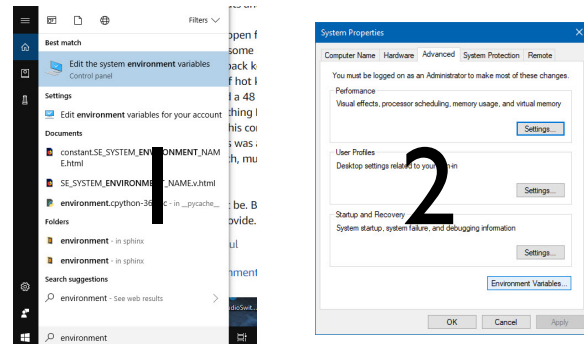
15

## In Windows, you modify the system path via the GUI

### Windows:

1. Type environment in the lower-left search bar.
2. Select Environment Variables... near the bottom.
3. Click on the Path variable in the top box then select Edit...
4. Add the three directories shown and move them to the top.

C:\Users\<you>\Anaconda3\Library\bin  
C:\Users\<you>\Anaconda3\mingw-w64\bin  
C:\Users\<you>\Anaconda3\usr\bin



# **BITS AND BYTES**

## **Bit**

- Always in one of two states, typically referred to as 0 and 1.
- The atomic unit of memory.

## **Byte**

- 8 bits
- The minimal addressable unit of memory.



# BITS AND BYTES

- **1 bit:** can represent  $2^1 = 2$  states (0 and 1).
- **2 bits:** can represent  $2^2 = 4$  states (00, 01, 10, and 11).

**We claim  $n$  bits can represent  $2^n$  states.**

- **$n$  bits:** for each state in  **$(n-1)$  bits**, prepend a 0 or 1.
  - Hence, there are twice as many states!
  - And,  $2 * 2^{n-1} = 2^n$ .
- |    |    |       |    |       |
|----|----|-------|----|-------|
| 00 | -> | 0(00) | or | 1(00) |
| 01 | -> | 0(01) | or | 1(01) |
| 10 | -> | 0(10) | or | 1(10) |
| 11 | -> | 0(11) | or | 1(11) |

## YOUR MEMORY (AND DISK)

Memory:	01100001	01010011	10011111	01	...
Address:	0	1	2	3	...

## YOUR MEMORY (AND DISK)

Hex:	0x61	0x53	0x9F	
Decimal:	97	83	191	
Memory:	01100001	01010011	10011111	01 ...
Address:	0	1	2	3 ...

## YOUR MEMORY (AND DISK)

Hex:	0x61	0x53	0x9F	
Decimal:	97	83	191	
Memory:	01100001	01010011	10011111	01 ...
Address:	0	1	2	3 ...

## YOUR MEMORY (AND DISK)

- All data is stored as 0s and 1s.
- So, we must come up with ways to encode different data types!
- **Unsigned integers** are encoded as binary (base 2).
- **Signed integers** are encoded using two's complement.
- **Floating point** is encoded using IEEE 754 (scientific notation).
- **Text** is encoded by mapping each possible state to a character.

# DATA SCIENCE: INTRO TO PYTHON

Dan Wilhelm, Lead Instructor  
[dan@danwilhelm.com](mailto:dan@danwilhelm.com)

# **INTRO TO PYTHON**

- What is Python?
- Python Fundamentals

# PEP 20 (THE ZEN OF PYTHON)

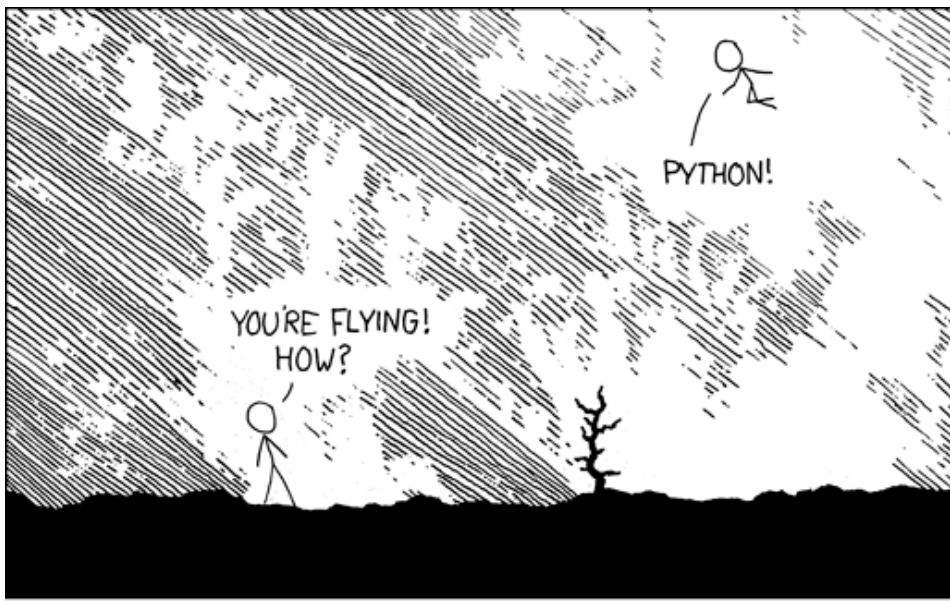
- Beautiful is better than ugly.
- Explicit is better than implicit.
- Simple is better than complex.
- Complex is better than complicated.
- Readability counts.

<https://www.python.org/dev/peps/pep-0020/>

**Also see PEP 8 (Style Guide for Python Code)**

<https://www.python.org/dev/peps/pep-0008/>





<https://xkcd.com/353/>



# PYTHON LANGUAGE FEATURES

- Interpreted
- High-level
- Emphasizes readability
  
- Supports many programming paradigms – e.g. object-oriented, imperative, functional, and procedural.
  
- Dynamic typing
- Strongly typed
- Automatic memory management

# WHY DO PEOPLE USE PYTHON?

- Open source (not locked in to one company)
- Clean syntax emphasizes readability
  - Great for scripting (interpreted & tolerant to errors)
- Modules and Packages
  - Large community – has many prewritten modules!
- Increased productivity
  - No compilation step, so edit-test-debug is fast
- Easy integration with C
  - Fast, pre-existing libraries can be used (called **bindings**)

# WHEN SHOULD YOU NOT USE PYTHON?

- Performance matters.
- Low memory usage is important.
- Direct access to hardware required.
- Want OS-specific GUI features.
- Programming embedded systems.

# PYTHON 2 VS PYTHON 3

“Short version: Python 2.x is legacy, Python 3.x is the present and future of the language” - <https://wiki.python.org/moin/Python2orPython3>

Are you new to Python?

Learn version 3.

Are you writing a new program?

Research whether the libraries you want to use support Python 3.  
(They probably do.) If so, use Python 3!

Are you maintaining an old project?

Use the version of Python they use.

# LEARNING PYTHON

- Official Tutorial

<https://docs.python.org/3/tutorial/index.html>

- Google's Python Class

<https://developers.google.com/edu/python/>

- Learn Python the Hard Way

<http://learnpythonthehardway.org/book/>

- Most importantly, MAKE THINGS!

---

**Insert class title**

---

# **PYTHON FUNDAMENTALS**

## UPDATING CONDA PACKAGES

**Update packages to latest stable Anaconda distro:**

```
conda update anaconda
```

**Update all packages to latest versions (warning: not recommended  
– will likely downgrade some to satisfy dependencies):**

```
conda update --all
```

**Upgrade to Python 3.7 from earlier version:**

```
conda install python=3.7
```

```
conda install anaconda
```

```
conda clean --packages
```



# **VIRTUAL ENVIRONMENTS**

**Create an empty virtual environment w/ Anaconda installed:**

```
conda env create -n <name> anaconda
```

**Create a Python 2.7 virtual environment:**

```
conda create -n py27 python=2.7 anaconda
```

**List available virtual environments:**

```
conda env list
```

# CONDA IS A PACKAGE MANAGER

- **conda** – A package management system used to install and manage software packages written in Python. It has benefits over pip:
  - Dependency checking!
  - Downloads binaries instead of compiling from source (although it can)
- Always try to install packages using conda first before pip.

**conda search <part of package name>**

**conda install <package name>**

**(Frees GBs of space!)**

**conda info <package>**

**conda clean --all**

# PIP INSTALLS PACKAGES

- **pip** – A package management system used to install and manage software packages written in Python.
- If a package is available on Anaconda, install via conda.
- If a package is not available on Anaconda:

```
pip install money  
pip uninstall money
```

# DATA TYPES

- Booleans (True or False)
- Numbers
- Strings
- Lists
- None
- Tuples
- Sets
- Dictionaries

## **ADDED BY PANDAS/NUMPY:**

- ndarray (from numpy)
- Series (from pandas)
- DataFrame (from pandas)

# **PYTHON IDENTIFIERS**

- Python is not as expansive as you may think.
- Every identifier you see (without any imports) is either an operator, reserved word, or built-in! So, there are a limited number of things to know.

# **PYTHON IDENTIFIERS**

- Python is not as expansive as you may think.
- Every identifier you see (without any imports) is either an operator, reserved word, or built-in! So, there are a limited number of things to know.

- An identifier is:

A single letter or underscore, followed by zero or more letters, numbers, or underscores.

`(letter | _) (letter | number | _)*`

# RESERVED WORDS/KEYWORDS

- Reserved words require special syntax to be used around the word – their syntax is unique compared to the rest of Python.
- Reserved words cannot be used as variables.

False	class	finally	is	return
None	continue	for	lambda	try
True	def	from	nonlocal	while
and	del	global	not	with
as	elif	if	or	yield
assert	else	import	pass	
break	except	in	raise	

# BUILT-IN FUNCTIONS

Built-in Functions				
<code>abs()</code>	<code>dict()</code>	<code>help()</code>	<code>min()</code>	<code>setattr()</code>
<code>all()</code>	<code>dir()</code>	<code>hex()</code>	<code>next()</code>	<code>slice()</code>
<code>any()</code>	<code>divmod()</code>	<code>id()</code>	<code>object()</code>	<code>sorted()</code>
<code>ascii()</code>	<code>enumerate()</code>	<code>input()</code>	<code>oct()</code>	<code>staticmethod()</code>
<code>bin()</code>	<code>eval()</code>	<code>int()</code>	<code>open()</code>	<code>str()</code>
<code>bool()</code>	<code>exec()</code>	<code>isinstance()</code>	<code>ord()</code>	<code>sum()</code>
<code>bytearray()</code>	<code>filter()</code>	<code>issubclass()</code>	<code>pow()</code>	<code>super()</code>
<code>bytes()</code>	<code>float()</code>	<code>iter()</code>	<code>print()</code>	<code>tuple()</code>
<code>callable()</code>	<code>format()</code>	<code>len()</code>	<code>property()</code>	<code>type()</code>
<code>chr()</code>	<code>frozenset()</code>	<code>list()</code>	<code>range()</code>	<code>vars()</code>
<code>classmethod()</code>	<code>getattr()</code>	<code>locals()</code>	<code>repr()</code>	<code>zip()</code>
<code>compile()</code>	<code>globals()</code>	<code>map()</code>	<code>reversed()</code>	<code>__import__()</code>
<code>complex()</code>	<code>hasattr()</code>	<code>max()</code>	<code>round()</code>	
<code>delattr()</code>	<code>hash()</code>	<code>memoryview()</code>	<code>set()</code>	



## **HOW TO CODE**

When coding, you only have a limited number of options:

- Use a built-in function.
- Use a method of the object.
- Import a module.
- Write your own function.

# STATEMENTS VS EXPRESSIONS

An **expression** is valid code that evaluates to a value:

```
>>> 10 + 7
```

```
>>> range(10)
```

A **statement** is a line of code that performs an action, e.g.:

```
>>> print(10 + 7)
```

```
>>> x = 10 + 7
```

**Programs** are collections of statements.

## CONDITIONALS

```
>>> donuts = 5
```

```
>>> if donuts < 2:
```

```
...     print("Good job!")
```

```
... else:
```

```
...     print("You ate too many donuts.")
```

- According to PEP 8, use 4 spaces per indentation level.
- Use '=' for assignment and '==' for comparison.
- Name your variables well! Why is "donuts" a poor name?

# STRINGS

```
>>> x = input("Exit program? (yes/no)")
>>> if x.lower() == "yes":
...     print("Exiting the program ...")
```

- Notice the method call. What does it do? Why is it done?
- Use '=' for assignment and '==' for comparison.
- Is *x* a good variable name? What should it be?
- ' vs " vs """
- Indexing characters.

## LISTS VS TUPLES

- A list contains ordered data, typically of the same data type:

```
>>> x = ["Tim", "Sandy", "Martin", "Shawna"]  
>>> print(x[0])
```

- A tuple contains ordered groups of variables, often of different data types:

```
>>> x = ("Tim", 5)    # Note (name, age) describe one person  
>>> print(x)
```

## LOOPING

- Given a list or tuple, we often want to do something with each member. To do this, we loop through the list:

```
>>> names = ["Tim", "Sandy", "Martin", "Shawna"]  
>>> for name in names:  
...     print name
```

- The indented code in the for block is run for each item in the list. Each time the indented code is run, *name* refers to the next item in the list. Note that *name* can be any name.

# FUNCTIONS

- A function allows us to take complex code and refer to it in an easy way. For example, “`x % 2 == 1`” is hard to understand by beginners if encountered in a program. Hence, to make the program easier to read, we refer to what the code does in English, as part of a function call that returns a value:

```
>>> def is_odd(x):  
...     return (x % 2 == 1)
```

- What is a better name for `x`?

## **HOW TO CODE**

When coding, you only have a limited number of options:

- Use a built-in function.
- Use a method or operator of the object.
- Import a package.
- Write your own function.



# INTERACTING WITH THE COMMAND LINE

- Retrieving command-line arguments is easy!

```
import sys  
sys.argv      # ARGument Vector (list)
```

Example: `python test.py 1 2 3`

```
sys.argv: ['test.py', '1', '2', '3']
```

# **INTERACTING WITH THE COMMAND LINE**

- Reading from stdin is just like reading a file!

```
import sys
for line in sys.stdin:
    print(line.strip())
```

# FILES & UNICODE

## **FILES (OLD WAY)**

```
file_in = open('test.txt', 'r', encoding='utf-8')
```

```
for line in file_in:  
    print(line)
```

```
file_in.close()
```

## FILES (OLD WAY)

```
file_in = open('test.txt', 'r', encoding='utf-8')
```

```
for line in file_in:  
    print(line)
```

```
file_in.close()
```

### common file modes

'r' – read

'w' – write

'a' – append

'b' – binary

## **FILES (NEW WAY) - READING**

```
# 'with' insures the file is closed if an exception occurs  
with open('test.txt', 'r' , encoding='utf-8') as fin:
```

```
    for line in fin:  
        print(line)
```

```
# Same thing, except the 'with' block ensures the file is  
# closed - even if an exception occurs!
```

## FILES (NEW WAY) - WRITING

```
# Write list of strings 'lines' to the file
with open('test.txt', 'w') as fout:
    for line in lines:
        fout.write(line + "\n")
```

# ENCODINGS

- All files are just sequences of bytes (aka 8-bit numbers). So, all files look alike.
- Hence, the operating system needs a way to know how to interpret a file's contents.



# ENCODINGS

- All files are just sequences of bytes (aka 8-bit numbers). So, all files look alike.
- Hence, the operating system needs a way to know how to interpret a file's contents.
- There are only two main ways to do this:
  1. The file extension provides a hint at how to interpret the file, e.g.:  
    .jpeg   .htaccess   .txt   .docx
  2. The beginning of the file could provide clues, e.g.:  
    all JPEG files start with the sequence: 255, 216

## ASCII VS UNICODE

- Text files (.txt, .py, .html) have historically used the 7-bit **ASCII** encoding
- However, ASCII only supports the English alphabet. So, many other text encodings were created.
- Collectively, these alternative encodings are known as **Unicode**.

Dec	Hx	Oct	Char	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr
0	0	000	<b>NUL</b> (null)	32	20	040	&#32; <b>Space</b>		64	40	100	&#64; <b>@</b>		96	60	140	&#96; <b>`</b>	
1	1	001	<b>SOH</b> (start of heading)	33	21	041	&#33; <b>!</b>		65	41	101	&#65; <b>A</b>		97	61	141	&#97; <b>a</b>	
2	2	002	<b>STX</b> (start of text)	34	22	042	&#34; <b>"</b>		66	42	102	&#66; <b>B</b>		98	62	142	&#98; <b>b</b>	
3	3	003	<b>ETX</b> (end of text)	35	23	043	&#35; <b>#</b>		67	43	103	&#67; <b>C</b>		99	63	143	&#99; <b>c</b>	
4	4	004	<b>EOT</b> (end of transmission)	36	24	044	&#36; <b>\$</b>		68	44	104	&#68; <b>D</b>		100	64	144	&#100; <b>d</b>	
5	5	005	<b>ENQ</b> (enquiry)	37	25	045	&#37; <b>%</b>		69	45	105	&#69; <b>E</b>		101	65	145	&#101; <b>e</b>	
6	6	006	<b>ACK</b> (acknowledge)	38	26	046	&#38; <b>&amp;</b>		70	46	106	&#70; <b>F</b>		102	66	146	&#102; <b>f</b>	
7	7	007	<b>BEL</b> (bell)	39	27	047	&#39; <b>'</b>		71	47	107	&#71; <b>G</b>		103	67	147	&#103; <b>g</b>	
8	8	010	<b>BS</b> (backspace)	40	28	050	&#40; <b>(</b>		72	48	110	&#72; <b>H</b>		104	68	150	&#104; <b>h</b>	
9	9	011	<b>TAB</b> (horizontal tab)	41	29	051	&#41; <b>)</b>		73	49	111	&#73; <b>I</b>		105	69	151	&#105; <b>i</b>	
10	A	012	<b>LF</b> (NL line feed, new line)	42	2A	052	&#42; <b>*</b>		74	4A	112	&#74; <b>J</b>		106	6A	152	&#106; <b>j</b>	
11	B	013	<b>VT</b> (vertical tab)	43	2B	053	&#43; <b>+</b>		75	4B	113	&#75; <b>K</b>		107	6B	153	&#107; <b>k</b>	
12	C	014	<b>FF</b> (NP form feed, new page)	44	2C	054	&#44; <b>,</b>		76	4C	114	&#76; <b>L</b>		108	6C	154	&#108; <b>l</b>	
13	D	015	<b>CR</b> (carriage return)	45	2D	055	&#45; <b>-</b>		77	4D	115	&#77; <b>M</b>		109	6D	155	&#109; <b>m</b>	
14	E	016	<b>SO</b> (shift out)	46	2E	056	&#46; <b>.</b>		78	4E	116	&#78; <b>N</b>		110	6E	156	&#110; <b>n</b>	
15	F	017	<b>SI</b> (shift in)	47	2F	057	&#47; <b>/</b>		79	4F	117	&#79; <b>O</b>		111	6F	157	&#111; <b>o</b>	
16	10	020	<b>DLE</b> (data link escape)	48	30	060	&#48; <b>0</b>		80	50	120	&#80; <b>P</b>		112	70	160	&#112; <b>p</b>	
17	11	021	<b>DC1</b> (device control 1)	49	31	061	&#49; <b>1</b>		81	51	121	&#81; <b>Q</b>		113	71	161	&#113; <b>q</b>	
18	12	022	<b>DC2</b> (device control 2)	50	32	062	&#50; <b>2</b>		82	52	122	&#82; <b>R</b>		114	72	162	&#114; <b>r</b>	
19	13	023	<b>DC3</b> (device control 3)	51	33	063	&#51; <b>3</b>		83	53	123	&#83; <b>S</b>		115	73	163	&#115; <b>s</b>	
20	14	024	<b>DC4</b> (device control 4)	52	34	064	&#52; <b>4</b>		84	54	124	&#84; <b>T</b>		116	74	164	&#116; <b>t</b>	
21	15	025	<b>NAK</b> (negative acknowledge)	53	35	065	&#53; <b>5</b>		85	55	125	&#85; <b>U</b>		117	75	165	&#117; <b>u</b>	
22	16	026	<b>SYN</b> (synchronous idle)	54	36	066	&#54; <b>6</b>		86	56	126	&#86; <b>V</b>		118	76	166	&#118; <b>v</b>	
23	17	027	<b>ETB</b> (end of trans. block)	55	37	067	&#55; <b>7</b>		87	57	127	&#87; <b>W</b>		119	77	167	&#119; <b>w</b>	
24	18	030	<b>CAN</b> (cancel)	56	38	070	&#56; <b>8</b>		88	58	130	&#88; <b>X</b>		120	78	170	&#120; <b>x</b>	
25	19	031	<b>EM</b> (end of medium)	57	39	071	&#57; <b>9</b>		89	59	131	&#89; <b>Y</b>		121	79	171	&#121; <b>y</b>	
26	1A	032	<b>SUB</b> (substitute)	58	3A	072	&#58; <b>:</b>		90	5A	132	&#90; <b>Z</b>		122	7A	172	&#122; <b>z</b>	
27	1B	033	<b>ESC</b> (escape)	59	3B	073	&#59; <b>;</b>		91	5B	133	&#91; <b>[</b>		123	7B	173	&#123; <b>{</b>	
28	1C	034	<b>FS</b> (file separator)	60	3C	074	&#60; <b>&lt;</b>		92	5C	134	&#92; <b>\</b>		124	7C	174	&#124; <b> </b>	
29	1D	035	<b>GS</b> (group separator)	61	3D	075	&#61; <b>=</b>		93	5D	135	&#93; <b>]</b>		125	7D	175	&#125; <b>}</b>	
30	1E	036	<b>RS</b> (record separator)	62	3E	076	&#62; <b>&gt;</b>		94	5E	136	&#94; <b>^</b>		126	7E	176	&#126; <b>~</b>	
31	1F	037	<b>US</b> (unit separator)	63	3F	077	&#63; <b>?</b>		95	5F	137	&#95; <b>_</b>		127	7F	177	&#127; <b>DEL</b>	

Source: [www.LookupTables.com](http://www.LookupTables.com)

# ASCII

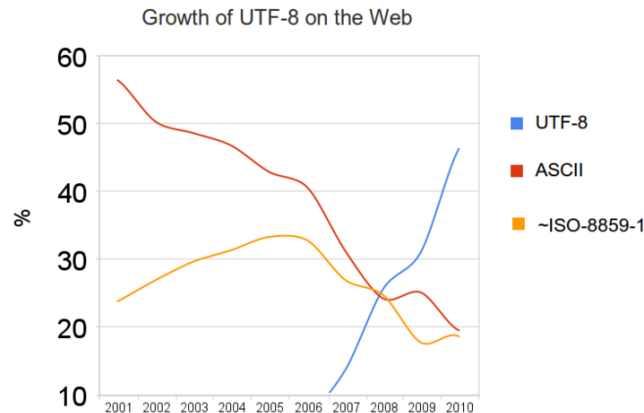
A 7-bit text encoding

## ASCII VS UNICODE

- *Problem:* Nowadays, text files could be encoded via ASCII or one of many Unicode encodings.
- However, there is no way to infer how to interpret the numbers -- the filename extensions are all the same, and text files immediately begin with the text content!

# ASCII VS UNICODE

- For example, the most popular encoding is UTF-8. It relies on ASCII being 7 bits and uses the 8<sup>th</sup> bit as a flag to indicate a special character.
- Hence, it is backwards compatible with ASCII. An ASCII file is also a UTF-8 file! (Assuming the 8<sup>th</sup> bits are all zero.)



## ASCII VS UNICODE

- For example, the most popular encoding is UTF-8. It relies on ASCII being 7 bits and uses the 8<sup>th</sup> bit as a flag to indicate a special character.
- Hence, it is backwards compatible with ASCII. An ASCII file is also a UTF-8 file! (Assuming the 8<sup>th</sup> bits are all zero.)
- If an ASCII file has a non-zero 8<sup>th</sup> bit, this will throw a Python exception saying it is invalid ASCII!
- If a UTF-8 encoded-file does not follow the UTF-8 rules, then it will also throw a Python exception, e.g. if Byte 1 is 110xxxxx, the next byte *must be* 10xxxxxx, where 'x' is any bit.

## COMMON UNICODE ENCODINGS FOR OPEN()

<code>'ascii'</code>	ASCII
<code>'latin-1' ('iso-8859-1')</code>	ASCII w/ extended ASCII
<code>'utf-8'</code>	variable-length encoding (min. 1 byte)*
<code>'utf-16'</code>	var-length encoding (min. 2 bytes)

\* backward-compatible with ASCII

Also:

- `\u<CODE POINT>` is the escape char for 2-byte Unicode chars
- `\U<CODE POINT>` for 4-byte Unicode chars

# UTF-8: A VARIABLE-LENGTH ENCODING

How to store different Unicode code words

	1 <sup>st</sup> byte	2 <sup>nd</sup> byte	3 <sup>rd</sup> byte
<b>U+0000</b> – <b>U+007F</b>	<b>0nnn nnnn</b>		
U+007F – U+07FF	110n nnnn	10nn nnnn	
U+0800 – U+FFFF	1110 nnnn	10nn nnnn	10nn nnnn

Note the backward-compatibility with ASCII! (bolded)



## **ASCII VS UNICODE**

- If we ask Python to interpret a file using a certain encoding and the file does not follow these rules, then an exception will be thrown.
- This is unfortunately the only way to know whether we guess the encoding correctly – attempt to interpret the file numbers in one way and see if it follows the encoding's rules.

## **ASCII VS UNICODE**

- If we ask Python to interpret a file using a certain encoding and the file does not follow these rules, then an exception will be thrown.
- This is unfortunately the only way to know whether we guess the encoding correctly – attempt to interpret the file numbers in one way and see if it follows the encoding's rules.
- So, it is preferred to try opening the file using different encoding. If this does not work, you can choose to ignore errors.
- Note: When you open a file in an IDE or text editor and save it, the IDE may resave the file in a different encoding!

# PYTHON + DATA SCIENCE TRICKS

## LISTS VS TUPLES

- Lists: Mutable (can be altered), typically homogenous values
- Tuples: Immutable, typically groups of items that go together

```
>>> people = [('Tim', 32), ('Sandy', 45)]
>>> points = [(0,0,0), (5,4,1), (7,7,8)]
>>> menu_item, price = ('Burger', 2.99)    # unpacking
>>> def max_population(cities):    # multiple return values
    ...
    return (city_name, population)
```

## QUOTES

- ‘vs “ ← same – both support escape characters (“\n”)
- ‘‘‘ ← triple quotes – allows actual newlines
- \ ← if at the end of a non-quoted line of code, allows you to split the line of code (there is an invisible newline after it)

```
>>> names = “Mike Wallace\nClara Simmons”
```

```
>>> names.split(“\n”)
```

```
>>> names.replace(“\n”, “, “)
```

# STRING FORMATTING

`“Person #{ } is { }”.format(2, ‘George’)`

`“Person #{num} is {name}”.format(num=5, name=‘Henry’)`

`“Your price will be ${price:.2f}.”.format(price=4.5127)`

# MODULES

```
import math
```

```
>>> math.sqrt(5)
```

```
from math import sqrt
```

```
>>> sqrt(5)
```

```
from math import *
```

```
>>> sqrt(5)
```

```
>>> tan(5)
```

# LIST COMPREHENSIONS

```
cubes = []  
for num in range(100):  
    cubes.append(num**3)
```

**BECOMES**

```
cubes = [num**3 for num in range(100)]
```



# LIST COMPREHENSIONS

```
cubes = []  
for num in range(100):  
    if num % 2 == 0:  
        cubes.append(num**3)
```

BECOMES

```
cubes = [num**3 for num in range(100) if num % 2 == 0]
```

NOTE you can use these for filtering!

## ENUMERATE – WHEN YOU NEED A LOOP INDEX

```
index = 0
for person in people:
    print("Person #{0} is {1}".format(index, person))
    index += 1
```

### BECOMES

```
for index, person in enumerate(people):
    print("Person #{0} is {1}".format(index, person))
```

# **DATES AND TIMES**

- `datetime.date` -> year, month, day
- `datetime.time` -> hour, minute, second, microsecond, tzinfo (TimeZone INFO)
- `datetime.datetime` -> year, month, day, hour, minute, second, microsecond, tzinfo
- `datetime.timedelta` – difference between dates/times

# EXCEPTIONS

*try:*

```
    num = int('not an int')
```

*except:* # catches ALL exceptions

```
    print('Exception caught!')
```

*try:*

```
    num = int('not an int')
```

*except ValueError:* # catches the ValueError exception

```
    print('Exception caught!')
```

## LISTS VS TUPLES

- Lists: Mutable (can be altered), typically homogenous values
- Tuples: Immutable, typically groups of items that go together

```
>>> people = [('Tim', 32), ('Sandy', 45)]
```

```
>>> points = [(0,0,0), (5,4,1), (7,7,8)]
```

## LISTS VS TUPLES

- Lists: Mutable (can be altered), typically homogenous values
- Tuples: Immutable, typically groups of items that go together

```
>>> people = [('Tim', 32), ('Sandy', 45)]
```

```
>>> points = [(0,0,0), (5,4,1), (7,7,8)]
```

```
>>> menu_item, price = ('Burger', 2.99)    # unpacking
```

## LISTS VS TUPLES

- Lists: Mutable (can be altered), typically homogenous values
- Tuples: Immutable, typically groups of items that go together

```
>>> people = [('Tim', 32), ('Sandy', 45)]
>>> points = [(0,0,0), (5,4,1), (7,7,8)]
>>> menu_item, price = ('Burger', 2.99)    # unpacking
>>> def max_population(cities):    # multiple return values
    ...
    return (city_name, population)
```

# FILES

# 'with' ensures the file is closed if an exception occurs  
with open('test.txt', 'r') as *fin*:

    for *line* in *fin*:  
        print(*line*)

# Write list of strings 'lines' to the file  
with open('test.txt', 'w') as *fout*:

    for *line* in *lines*:  
        fout.write(*line* + "\n")



## **OTHER STUFF**

- Sets
- Using dictionaries for uniqueness/histograms

## MODULES FOR SCRAPING & APIS

`pip install requests`

<http://docs.python-requests.org/en/latest/>

`pip install beautifulsoup4`

<http://www.crummy.com/software/BeautifulSoup/bs4/doc/>