61A Lecture 13

Wednesday, October 2

Strings

String Literals Have Three Forms

```
>>> 'I am string!'
'I am string!'
>>> "I've got an apostrophe"
'I've yot an apostrophe"
>>> '您好'
'您好'
>>> """The Zen of Python claims, Readability counts.
Read more: import this.""
'The Zen of Pythoniclaims, Readability counts.
Read more: import this.""
'The Zen of Pythoniclaims, Readability counts.

A backslash "escapes" the following character represents a new line
```

Announcements

- ·Homework 3 deadline extended to Wednesday 10/2 @ 11:59pm.
- Optional Hog strategy contest due Thursday 10/3 @ 11:59pm.
- ·Homework 4 due Tuesday 10/8 @ 11:59pm.
- ·Project 2 due Thursday 10/10 @ 11:59pm.
- •Guerrilla Section 2 this Saturday 10/5 & Sunday 10/6 10am-1pm in Soda.
- *Topics: Data abstraction, sequences, and non-local assignment.
- *Please RSVP on Piazza!
- •Guest lecture on Wednesday 10/9, Peter Norvig on Natural Language Processing in Python.

Strings are an Abstraction

Representing data:

'200' '1.2e-5' 'False' '(1, 2)'

Representing language:

"""And, as imagination bodies forth
The forms of things to unknown, and the poet's pen
Turns them to shapes, and gives to airy nothing
A local habitation and a name.

Representing programs:

'curry = lambda f: lambda x: lambda y: f(x, y)'
(Demo)

Strings are Sequences

Length. A sequence has a finite length.

Element selection. A sequence has an element corresponding to any non-negative integer index less than its length, starting at \emptyset for the first element.

(Demo)

String Membership Differs from Other Sequence Types

```
The "in" and "not in" operators match substrings
>>> 'here' in "Where's Waldo?"
True
>>> 234 in (1, 2, 3, 4, 5)
False
```

Why? Working with strings, we usually care about words more than characters

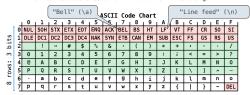
The count method also matches substrings

```
>>> 'Mississippi'.count('i')
4
>>> 'Mississippi'.count('issi')
```

the number of non-overlapping occurrences of a substring

Representing Strings: the ASCII Standard

American Standard Code for Information Interchange



16 columns: 4 bits

- Layout was chosen to support sorting by character code
- Rows indexed 2-5 are a useful 6-bit (64 element) subset
- \bullet Control characters were designed for transmission (Demo)

Representing Strings: UTF-8 Encoding

 ${\tt UTF} \ ({\tt UCS} \ ({\tt Universal} \ {\tt Character} \ {\tt Set}) \ {\tt Transformation} \ {\tt Format})$

Unicode: Correspondence between characters and integers

UTF-8: Correspondence between those integers and bytes

A byte is 8 bits and can encode any integer 0-255.

00000000 0 0 00000001 1 integers 0000001 2 integers

Variable-length encoding: integers vary in the number of bytes required to encode them.

In Python: string length is measured in characters, bytes length in bytes.

(Demo)

Encoding Strings

Representing Strings: the Unicode Standard

- 109,000 characters
- 93 scripts (organized)
- Enumeration of character properties, such as case
- Supports bidirectional display order
- A canonical name for every character

拏	聲	聳	聴	聵	虛	職	聯
8071	8872	8073	8074	8075	9076	8877	8471
健	腲	腳	腴	服	腶	腷	腸
根	色	艳	艴	艵	艶	艶	艸
草	重		荴	荵		荷	
葱	葲	葳	葴	葵	葶	葷	湛

U+0058 LATIN CAPITAL LETTER ${\sf X}$

U+263a WHITE SMILING FACE

U+2639 WHITE FROWNING FACE

I @ I



(Demo)

Sequence Processing

Sequence Processing

Consider two problems:

Sum the even members of the first n Fibonacci numbers.

 $^{\circ}\text{List}$ the letters in the acronym for a name, which includes the first letter of each capitalized word.

enumerate naturals:	1, 2, 3	, 4, 5, 6	, 7, 8, 9	10, 11.
map fib:	0, 1, 1	, 2, 3, 5 A	, 8, 13, 21 ^	34, 55.
filter even:	0,	2,	8,	34, .
accumulate sum:	.,	.,	.,	., =44

Mapping a Function over a Sequence

Apply a function to each element of the sequence

```
>>> alternates = (-1, 2, -3, 4, -5)
>>> tuple(map(abs, alternates))
(1, 2, 3, 4, 5)
```

The returned value of $\left(\begin{array}{c} \mathbf{map} \end{array} \right)$ is an iterable map object

A constructor for the built-in map type

The returned value of filter is an iterable filter object

(Demo)

Iterable Values and Accumulation

Iterable objects give access to their elements in order.

Similar to a sequence, but does not always allow element selection or have finite length.

Many built-in functions take iterable objects as argument.

tuple Return a tuple containing the elements sum Return the sum of the elements min Return the minimum of the elements max Return the maximum of the elements

For statements also operate on iterable values.

Sequence Processing

Consider two problems:

 $\,{}^{{}_{^{\circ}}}\mathsf{Sum}$ the even members of the first n Fibonacci numbers.

 $\begin{tabular}{lll} \begin{tabular}{lll} \begin{$

```
enumerate words:

'University', 'of', 'California', 'Berkeley'

A

filter capitalized:

'University', 'California', 'Berkeley'

map first:

'U', 'C', 'B'

accumulate tuple:

('U', 'C', 'B')
```

Iteration and Accumulation

Reducing a Sequence

Reduce is a higher-order generalization of max, min, & sum.

```
>>> from operator import mul
>>> from functools import reduce
>>> reduce(mul, (1, 2, 3, 4, 5))

First argument:
A two-argument function

Second argument: an iterable object
```

Similar to accumulate from Homework 2, but with iterable objects.

Generator Expressions

One large expression that evaluates to an iterable object $% \left(1\right) =\left(1\right) \left(1\right) \left($

(<map exp> for <name> in <iter exp> if <filter exp>)

- ullet Evaluates to an iterable object.
- $\bullet\!<\!\!\text{iter exp>}$ is evaluated when the generator expression is evaluated.
- \bullet Remaining expressions are evaluated when elements are accessed.

Short version: (<map exp> for <name> in <iter exp>)

(Demo)