

INTRO TO DATA SCIENCE

LECTURE 2: GETTING DATA

OCTOBER 1, 2014 // DAT10 SF

Francesco Mosconi, PhD

HEADER— CLASS NAME, PRESENTATION TITLE

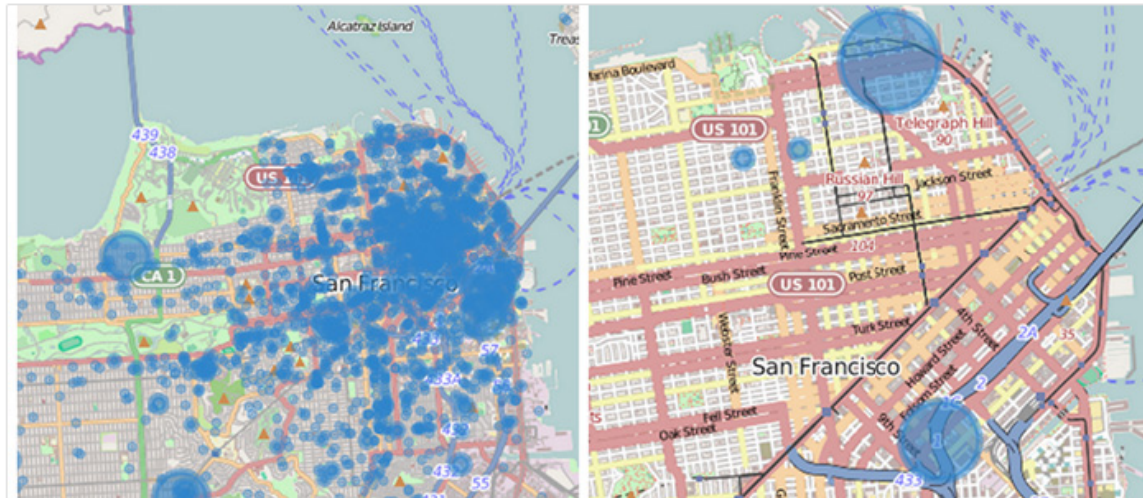
DATA SCIENCE IN THE NEWS

DATA SCIENCE IN THE NEWS

#UBERDATA

MAKING A BAYESIAN MODEL TO INFER UBER RIDER DESTINATIONS

SEPTEMBER 2, 2014
POSTED BY REN LU



Uber uses data science to predict
where its riders want to go



Image Credit: Jason Tester Guerilla Futures/Flickr

Source: <http://blog.uber.com/passenger-destinations>

DATA SCIENCE IN THE NEWS

The Best Questions For A First Date

April 20th, 2011 by [Christian Rudder](#)

 Tweet 3,371

 Condividi 24mila

First dates are awkward. There is so much you want to know about the person across the table from you, and yet so little you can directly ask.



Source: <http://blog.okcupid.com/index.php/page/3/>

RECAP

LAST TIME

- Data Science
- Data Scientist
- Data Mining Workflow
- Ipython
- Git

INTRO TO DATA SCIENCE

QUESTIONS?

INTRO TO DATA SCIENCE

GETTING DATA

INTRO TO DATA SCIENCE

GETTING DATA WHERE

THE DATA SCIENCE WORKFLOW

DATAIST (HILARY MASON & FRIENDS)

- 1. Obtain - pointing and clicking does not scale (APIs, Python, shell scripting)
- 2. Scrub - “Scrubbing data is the least sexy part of the analysis process, but often one that yields the greatest benefits” (Python, sed, awk, grep)
- 3. Explore - look at the data (visualizing, clustering, dimensionality reduction)
- 4. Model - “All models are wrong, but some are useful” / models are built to predict and interpret!
- 5. Interpret - “The purpose of computing is insight, not numbers”

THE DATA SCIENCE WORKFLOW

DATAIST (HILARY MASON & FRIENDS)

- 1. Obtain - pointing and clicking does not scale (APIs, Python, shell scripting)
- 2. Scrub - “Scrubbing data is the least sexy part of the analysis process, but often one that yields the greatest benefits” (Python, sed, awk, grep)
- 3. Explore - look at the data (visualizing, clustering, dimensionality reduction)
- 4. Model - “All models are wrong, but some are useful” / models are built to predict and interpret!
- 5. Interpret - “The purpose of computing is insight, not numbers”

INTRO TO DATA SCIENCE, GETTING DATA

Where can we get data from?




[About](#)
[Citation Policy](#)
[Donate a Data Set](#)
[Contact](#)





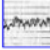
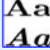




☒ Repository
 ☐ Web
 


[View ALL Data Sets](#)

Machine Learning Repository
 Center for Machine Learning and Intelligent Systems

Browse Through: 298 Data Sets

Table View [List View](#)





| Default Task | Name | Data Types | Default Task | Attribute Types | # Instances | # Attributes | Year |
|--|---|--------------|---------------------|----------------------------|-------------|--------------|------|
| Classification (213) Regression (41) Clustering (36) Other (50) |  Abalone | Multivariate | Classification | Categorical, Integer, Real | 4177 | 8 | 1995 |
| Attribute Type Categorical (36) Numerical (161) Mixed (56) |  Adult | Multivariate | Classification | Categorical, Integer | 48842 | 14 | 1996 |
| Data Type Multivariate (228) Univariate (15) Sequential (26) Time-Series (43) Text (27) Domain-Theory (20) Other (21) |  Annealing | Multivariate | Classification | Categorical, Integer, Real | 798 | 38 | |
| Area Life Sciences (75) Physical Sciences (41) CS / Engineering (78) Social Sciences (20) Business (14) Game (9) Other (59) |  Anonymous Microsoft Web Data | | Recommender-Systems | Categorical | 37711 | 294 | 1998 |
| # Attributes Less than 10 (74) 10 to 100 (129) Greater than 100 (46) |  Arrhythmia | Multivariate | Classification | Categorical, Integer, Real | 452 | 279 | 1998 |
| # Instances Less than 100 (15) 100 to 1000 (113) Greater than 1000 (140) |  Artificial Characters | Multivariate | Classification | Categorical, Integer, Real | 6000 | 7 | 1992 |
| Format Type Matrix (213) Non-Matrix (85) |  Audiology (Original) | Multivariate | Classification | Categorical | 226 | | 1987 |
| |  Audiology (Standardized) | Multivariate | Classification | Categorical | 226 | 69 | 1992 |
| |  Auto MPG | Multivariate | Regression | Categorical, Real | 398 | 8 | 1993 |
| |  Automobile | Multivariate | Regression | Categorical, Integer, Real | 205 | 26 | 1987 |



SEARCH

Espanol

Follow Us:



1-800-FED-INFO (333-4636)

| Services and Information | Government Agencies and Elected Officials | Blog |
|--|---|--|
| <ul style="list-style-type: none">Benefits, Grants, and LoansBusinesses and NonprofitsConsumer Complaints and ProtectionConsumer PublicationsDisasters, Public Safety, and LawsEnvironment, Energy, and Agriculture | <ul style="list-style-type: none">Government Sales and AuctionsHealth Insurance, Nutrition, and Food SafetyHistory, Genealogy, and CultureImmigration, Citizenship, and InternationalJobs, Training, and EducationMortgages, Housing, and Family | <ul style="list-style-type: none">Passports and TravelPublic Service and VolunteerismReference and General GovernmentRegister to Vote and ElectionsScience and TechnologyUnclaimed Money, Taxes, and Credit Reports |

☆ More for Developers

- Other USA.gov Resources
- USA.gov GitHub Account

From Other Federal Agencies

- Other Federal Government Developer Resources
- Other Federal Government GitHub Accounts

About The Data

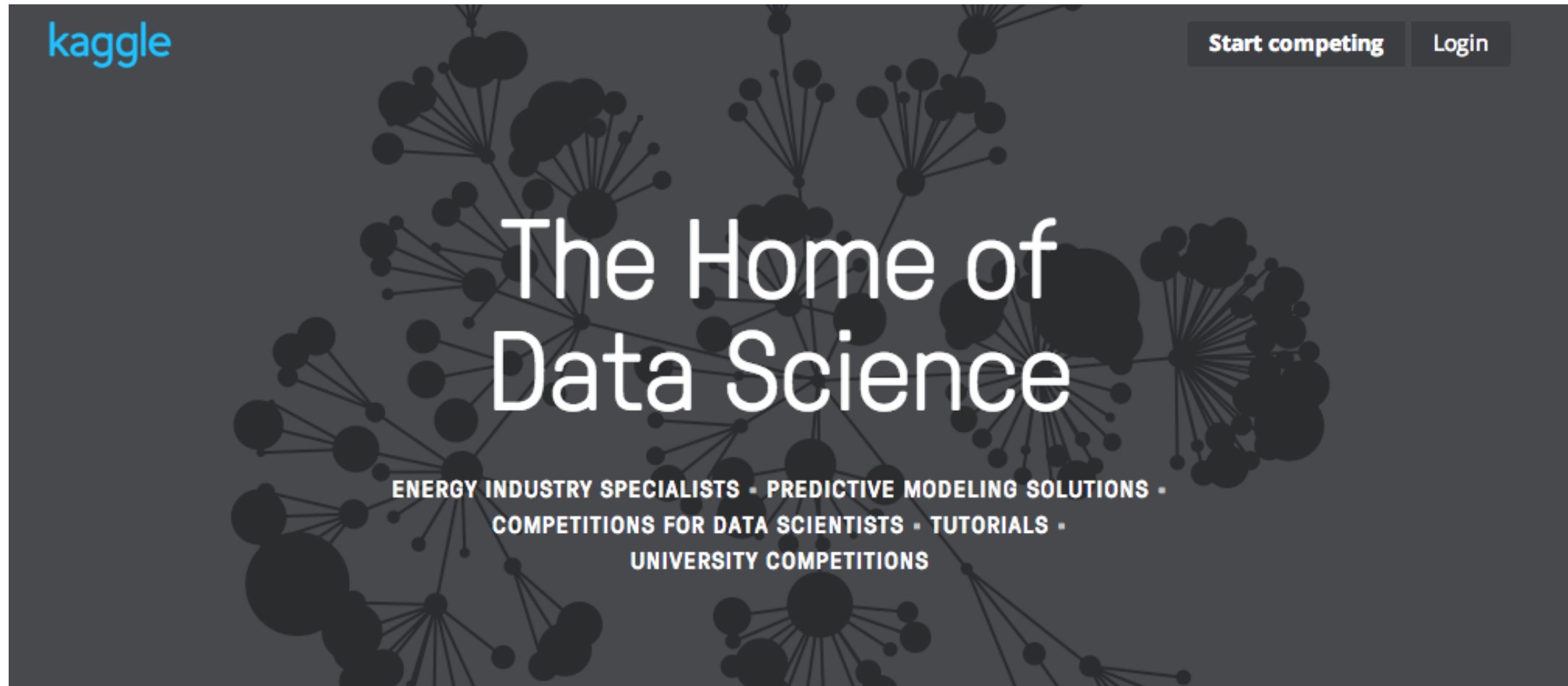
1.USA.gov URLs are created whenever anyone shortens a .gov or .mil URL using [bitly](#).

We provide a raw [pub/sub](#) feed of data created any time anyone clicks on a 1.USA.gov URL. The pub/sub endpoint responds to http requests for any 1.USA.gov URL and returns a stream of JSON entries, one per line, that represent real-time clicks.

If you are using the 1.USA.gov data and have questions, feedback, or want to tell us about your product, please [e-mail us](#).

How to Access The Data

KAGGLE



Source: <http://www.kaggle.com/>

LISTS OF DATASETS CURATED BY FAMOUS DATA SCIENTISTS

- 1) Pete Skomoroch (LinkedIn) <https://delicious.com/pskomoroch/dataset>
 - 2) Hilary Mason (Accel Partners, Bitly) <https://bitly.com/bundles/hmason/1>
 - 3) Kevin Chai (U. of New South Wales, Sydney) <http://kevinchai.net/datasets>
 - 4) Jeff Hammerbacher (Cloudera) <http://www.quora.com/Jeff-Hammerbacher/Introduction-to-Data-Science-Data-Sets>
 - 5) Jerry Smith (3i-MIND) <http://datascientistinsights.com/2013/10/07/data-repositories-mothers-milk-for-data-scientists/>
 - 6) Gregory Piatetsky-Shapiro (KDD) <http://www.kdnuggets.com/datasets/index.html>
-
- Bonus: <http://www.quora.com/Data/Where-can-I-find-large-datasets-open-to-the-public>

Pair exercise:

choose a data source and look at what data you can get
discuss how you would use the data

INTRO TO DATA SCIENCE, GETTING DATA

QUESTIONS?

INTRO TO DATA SCIENCE

GETTING DATA

WHAT

INTRO TO DATA SCIENCE, GETTING DATA

DATA FORMATS?

- Text
- CSV
- Json
- Xml
- dat
- images, binaries, etc. etc. etc.

INTRO TO DATA SCIENCE, GETTING DATA

JSON

- JSON (JavaScript Object Notation) is a borrowed JavaScript form turned into a string that can be passed between applications.

INTRO TO DATA SCIENCE, GETTING DATA

JSON

- JSON (JavaScript Object Notation) is a borrowed JavaScript form turned into a string that can be passed between applications.

JSON

- JSON (JavaScript Object Notation) is a borrowed JavaScript form turned into a string that can be passed between applications.
- JSON is a lightweight data-interchange format.

JSON

- JSON (JavaScript Object Notation) is a borrowed JavaScript form turned into a string that can be passed between applications.
- JSON is a lightweight data-interchange format.
- JSON is easy for humans to read and write.

JSON

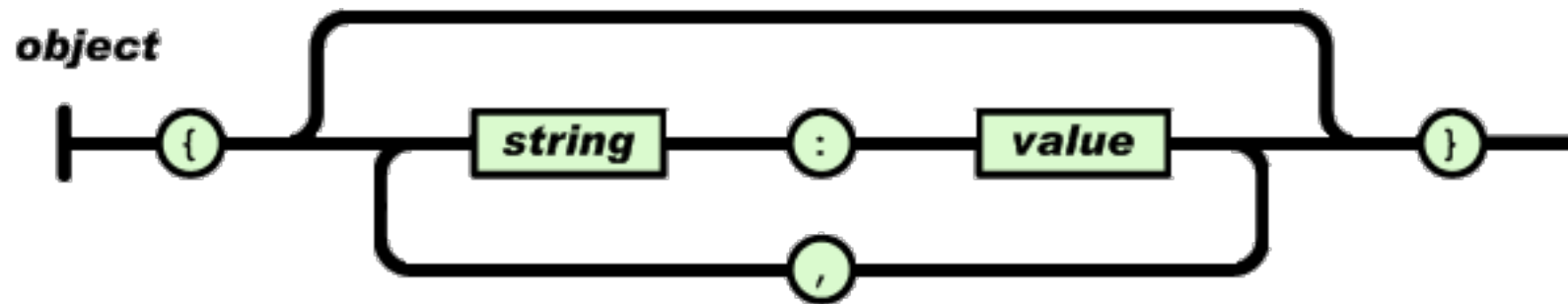
- JSON (JavaScript Object Notation) is a borrowed JavaScript form turned into a string that can be passed between applications.
- JSON is a lightweight data-interchange format.
- JSON is easy for humans to read and write.
- JSON is easy for machines to parse and generate.

JSON

- JSON (JavaScript Object Notation) is a borrowed JavaScript form turned into a string that can be passed between applications.
- JSON is a lightweight data-interchange format.
- JSON is easy for humans to read and write.
- JSON is easy for machines to parse and generate.
- JSON are passed through applications as strings, and converted into native objects per language.

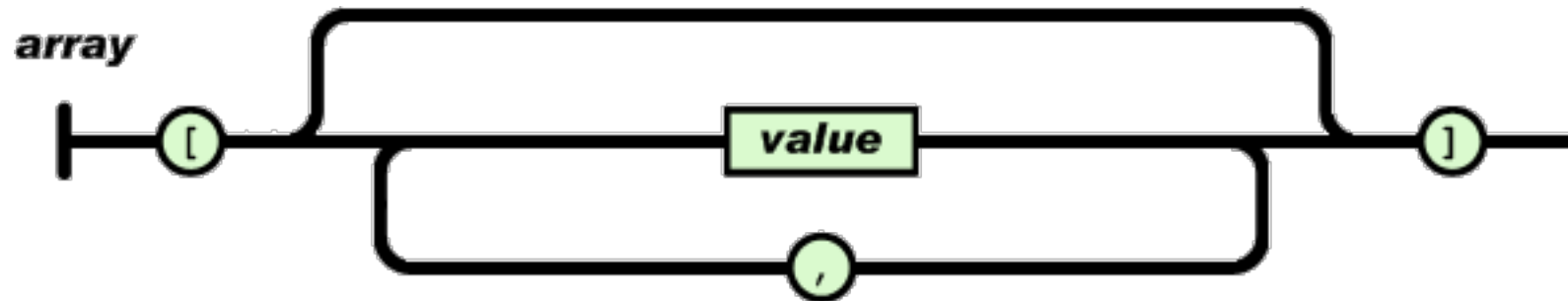
JSON

An object is an unordered set of name/value pairs. An object begins with { (left brace) and ends with } (right brace). Each name is followed by : (colon) and the name/value pairs are separated by , (comma).



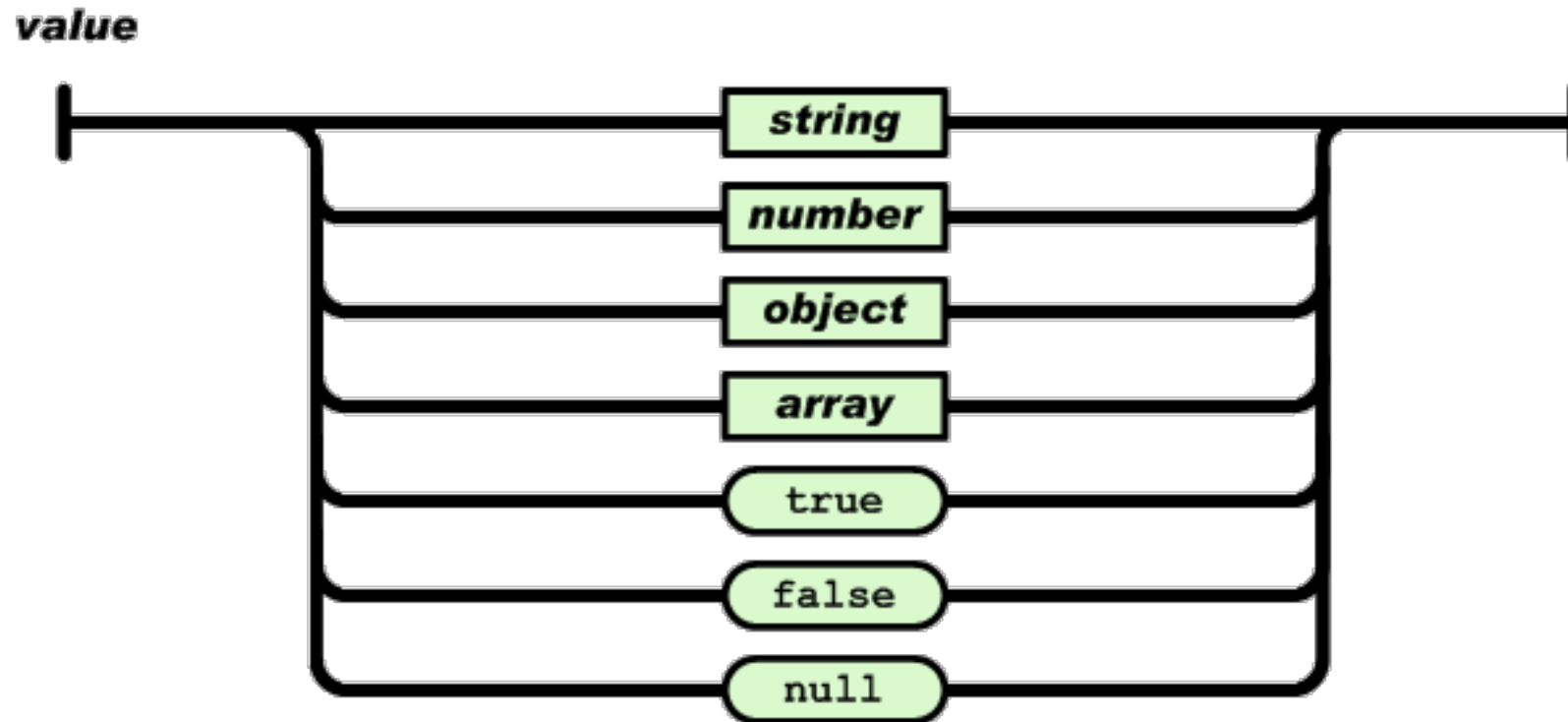
JSON

An array is an ordered collection of values. An array begins with [(left bracket) and ends with] (right bracket). Values are separated by , (comma).



JSON

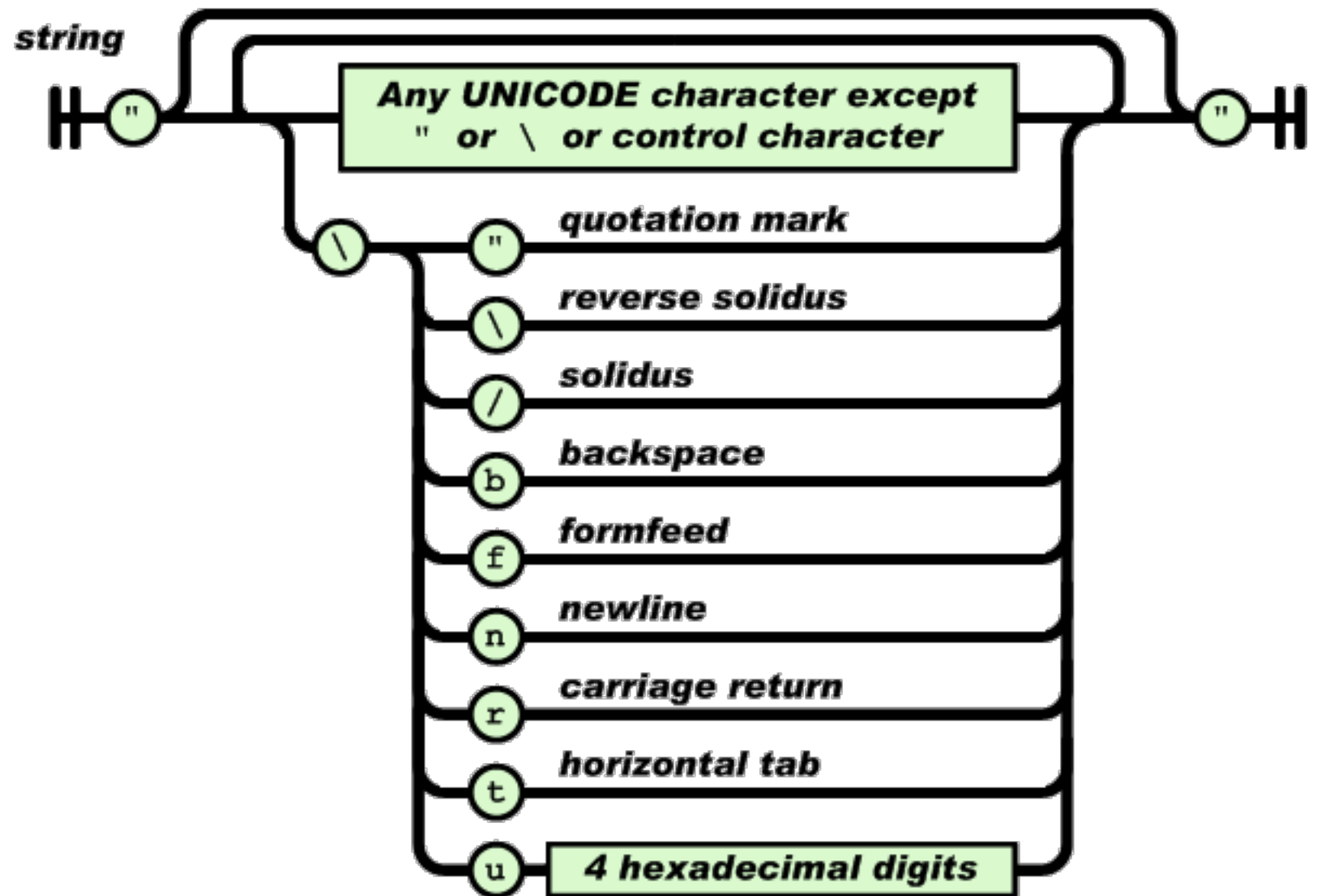
A value can be a string in double quotes, or a number, or true or false or null, or an object or an array. These structures can be nested.



INTRO TO DATA SCIENCE, GETTING DATA

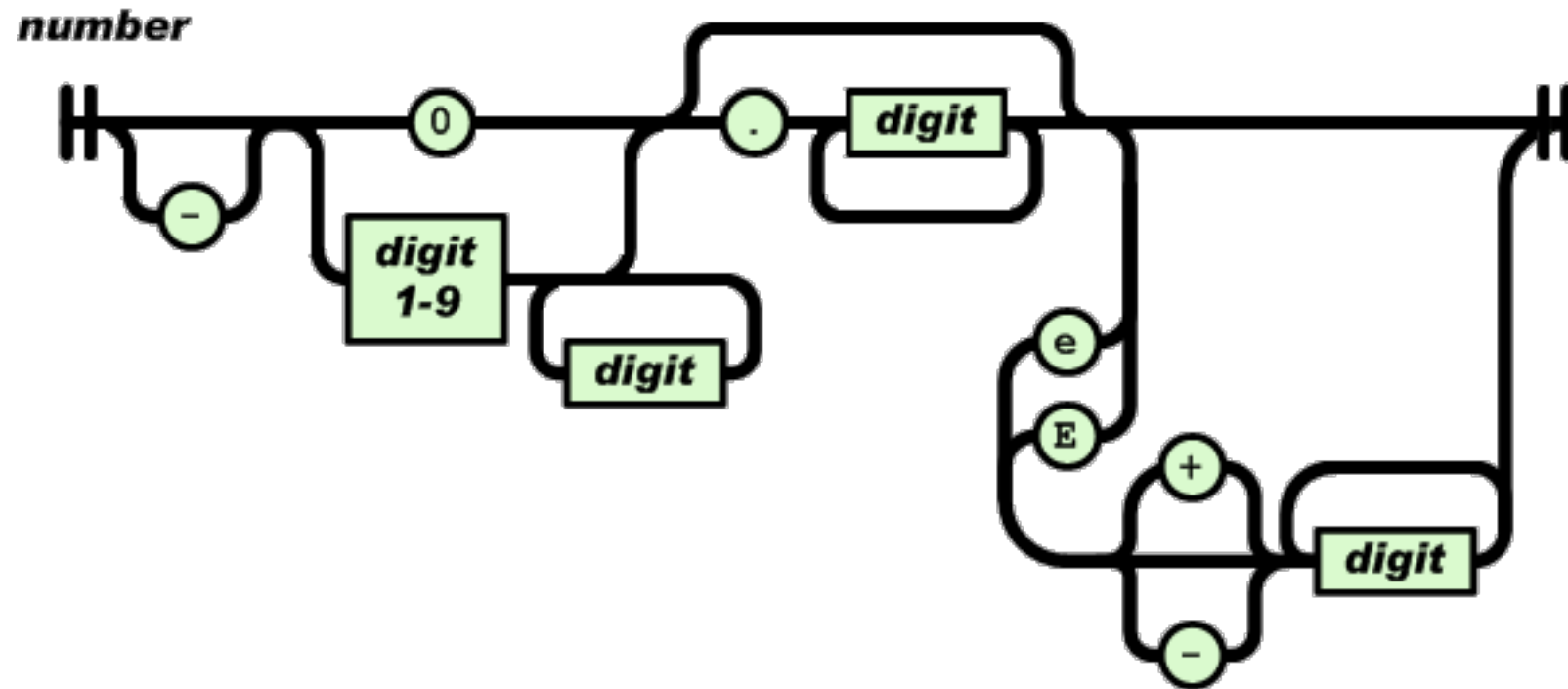
JSON

A string is a sequence of zero or more Unicode characters, wrapped in double quotes, using backslash escapes. A character is represented as a single character string. A string is very much like a C



JSON

A number is very much like a C or Java number, except that the octal and hexadecimal formats are not used.




```
>>> someFile = open('/Users/epodojil/GA_Data_Science/a.json').read()
>>> print json.dumps(someFile)
"{\n  \"glossary\": {\n    \"title\": \"example glossary\", \n    \"GlossDiv\": {\n      \"title\": \"S\", \n      \"GlossList\": {\n        \"GlossEntry\": {\n          \"ID\": \"SGML\", \n          \"SortAs\": \"SGML\", \n          \"GlossTerm\": \"Standard Generalized Markup Language\", \n          \"Acronym\": \"SGML\", \n          \"Abbrev\": \"ISO 8879:1986\", \n          \"GlossDef\": {\n            \"para\": \"A meta-markup language, used to create markup languages such as DocBook.\", \n            \"GlossSeeAlso\": [\"GML\", \"XML\"] \n          }, \n          \"GlossSee\": \"markup\" \n        } \n      } \n    } \n  } \n}"
```

```
>>> print someFile
```

```
{
  "glossary": {
    "title": "example glossary",
    "GlossDiv": {
      "title": "S",
      "GlossList": {
        "GlossEntry": {
          "ID": "SGML",
          "SortAs": "SGML",
          "GlossTerm": "Standard Generalized Markup Language",
          "Acronym": "SGML",
          "Abbrev": "ISO 8879:1986",
          "GlossDef": {
            "para": "A meta-markup language, used to create markup languages such as DocBook.",
            "GlossSeeAlso": ["GML", "XML"]
          },
          "GlossSee": "markup"
        }
      }
    }
  }
}
```

Object

String

Python Dict

```
>>> print json.loads(someFile)
{'glossary': {'GlossDiv': {'GlossList': {'GlossEntry': {'GlossDef': {'GlossSeeAlso': ['GML', 'XML'], 'para': 'A meta-
': 'markup', 'Acronym': 'SGML', 'GlossTerm': 'Standard Generalized Markup Language', 'Abbrev': 'ISO 8879:1986', 'SortAs
```

INTRO TO DATA SCIENCE, GETTING DATA

Json

```
{
  "firstName": "John",
  "lastName": "Smith",
  "age": 25,
  "address": {
    "streetAddress": "21 2nd Street",
    "city": "New York",
    "state": "NY",
    "postalCode": "10021"
  },
  "phoneNumber": [
    {
      "type": "home",
      "number": "212 555-1239"
    },
    {
      "type": "fax",
      "number": "646 555-4567"
    }
  ],
  "gender": {
    "type": "male"
  }
}
```

Source: <http://en.wikipedia.org/wiki/JSON>

INTRO TO DATA SCIENCE, GETTING DATA

Json

```
{
  "firstName": "John",
  "lastName": "Smith",
  "age": 25,
  "address": {
    "streetAddress": "21 2nd Street",
    "city": "New York",
    "state": "NY",
    "postalCode": "10021"
  },
  "phoneNumber": [
    {
      "type": "home",
      "number": "212 555-1239"
    },
    {
      "type": "fax",
      "number": "646 555-4567"
    }
  ],
  "gender": {
    "type": "male"
  }
}
```

YAML

```
---
firstName: John
lastName: Smith
age: 25
address:
  streetAddress: 21 2nd Street
  city: New York
  state: NY
  postalCode: 10021

phoneNumber:
  -
    type: home
    number: 212 555-1234
  -
    type: fax
    number: 646 555-4567
gender:
  type: male
```

INTRO TO DATA SCIENCE, GETTING DATA

Json

```
{
  "firstName": "John",
  "lastName": "Smith",
  "age": 25,
  "address": {
    "streetAddress": "21 2nd Street",
    "city": "New York",
    "state": "NY",
    "postalCode": "10021"
  },
  "phoneNumber": [
    {
      "type": "home",
      "number": "212 555-1239"
    },
    {
      "type": "fax",
      "number": "646 555-4567"
    }
  ],
  "gender": {
    "type": "male"
  }
}
```

XML

```
<person>
  <firstName>John</firstName>
  <lastName>Smith</lastName>
  <age>25</age>
  <address>
    <streetAddress>21 2nd Street</streetAddress>
    <city>New York</city>
    <state>NY</state>
    <postalCode>10021</postalCode>
  </address>
  <phoneNumbers>
    <phoneNumber type="home">212 555-1234</phoneNumber>
    <phoneNumber type="fax">646 555-4567</phoneNumber>
  </phoneNumbers>
  <gender>
    <type>male</type>
  </gender>
</person>
```

INTRO TO DATA SCIENCE, GETTING DATA

QUESTIONS?

INTRO TO DATA SCIENCE

GETTING DATA

HOW

INTRO TO DATA SCIENCE, GETTING DATA

APIS

- APIs (Application Programming Interface) allow people to interact with the structures of an application to get, put, delete, or update data.

INTRO TO DATA SCIENCE, GETTING DATA

APIS

- APIs (Application Programming Interface) allow people to interact with the structures of an application to get, put, delete, or update data.
- Best practices for APIs are to use RESTful principles.

INTRO TO DATA SCIENCE, GETTING DATA

RESTFUL APIS

- Base URL and collection.
- Interactive media type (usually JSON)
- Operations (GET, PUT, POST, DELETE)
- Driven by Hypertext (http requests)

Collection



```
GET https://api.instagram.com/v1/users/10
```

Operation



INTRO TO DATA SCIENCE, GETTING DATA

**GET https://api.instagram.com/v1/users/
search/?q=andy**



Querystring

HEADER– CLASS NAME, PRESENTATION TITLE

HEADER 2

- RESTful APIs can always be accessed using cURL requests: hence why hypertext access is a requirement
- Most have language libraries to make it easier to access through the language of your choice.
- <http://www.pythonapi.com/>

INTRO TO DATA SCIENCE, GETTING DATA

Pair exercise: ▸ <http://www.pythonapi.com/>

choose an API and look at what data you can get

install python module to extract data

discuss how you could leverage the data from that API

INTRO TO DATA SCIENCE

QUESTIONS?