# Web GIS 1: Big Data & GIS

ENV 859 – Advanced GIS

# What is Big Data?

**40 ZETTABYTES**
[ 43 TRILLION GIGABYTES ]
of data will be created by 2020, an increase of 300 times from 2005

2005
2020

**6 BILLION PEOPLE**
have cell phones

**WORLD POPULATION: 7 BILLION**

It's estimated that
**2.5 QUINTILLION BYTES**
[ 2.3 TRILLION GIGABYTES ]
of data are created each day

## Volume
### SCALE OF DATA

Most companies in the U.S. have at least
**100 TERABYTES**
[ 100,000 GIGABYTES ]
of data stored

# The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume, Velocity, Variety and Veracity**

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015
**4.4 MILLION IT JOBS**
will be created globally to support big data, with 1.9 million in the United States

As of 2011, the global size of data in healthcare was estimated to be
**150 EXABYTES**
[ 161 BILLION GIGABYTES ]

**30 BILLION PIECES OF CONTENT**
are shared on Facebook every month

## Variety
### DIFFERENT FORMS OF DATA

By 2014, it's anticipated there will be
**420 MILLION WEARABLE, WIRELESS HEALTH MONITORS**

**4 BILLION+ HOURS OF VIDEO**
are watched on YouTube each month

**400 MILLION TWEETS**
are sent per day by about 200 million monthly active users

The New York Stock Exchange captures
**1 TB OF TRADE INFORMATION**
during each trading session

## Velocity
### ANALYSIS OF STREAMING DATA

By 2016, it is projected there will be
**18.9 BILLION NETWORK CONNECTIONS**
~ almost 2.5 connections per person on earth

Modern cars have close to
**100 SENSORS**
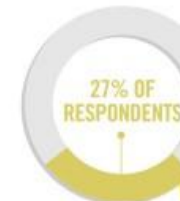that monitor items such as fuel level and tire pressure

**1 IN 3 BUSINESS LEADERS**
don't trust the information they use to make decisions

**27% OF RESPONDENTS**

in one survey were unsure of how much of their data was inaccurate

## Veracity
### UNCERTAINTY OF DATA
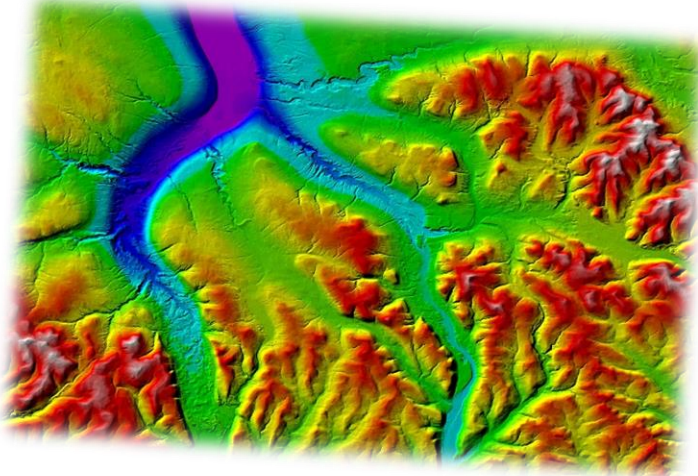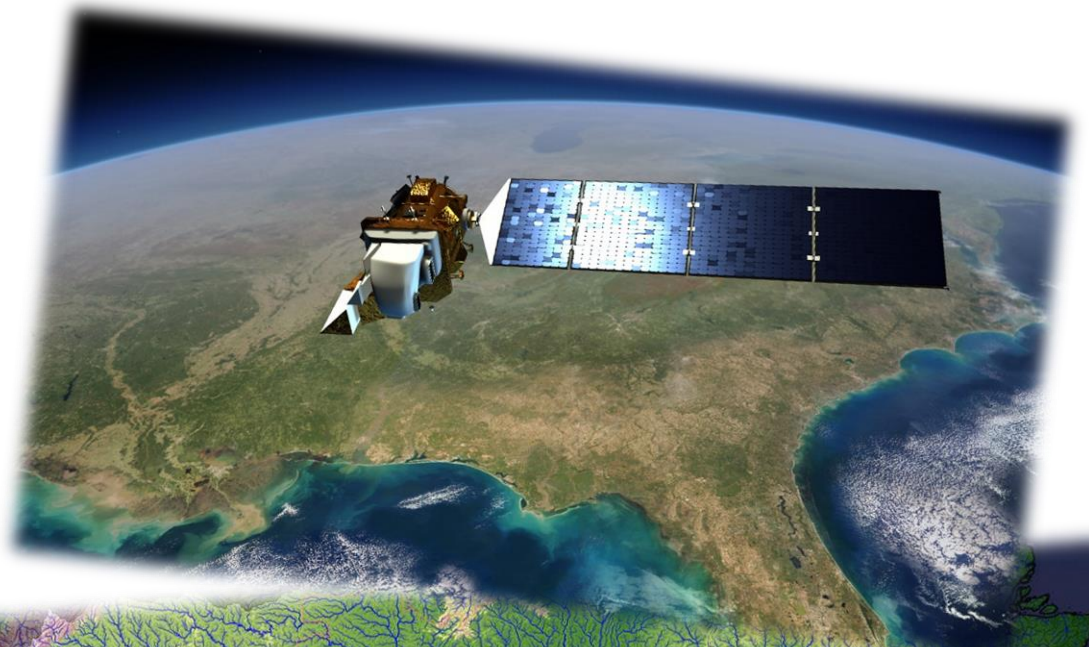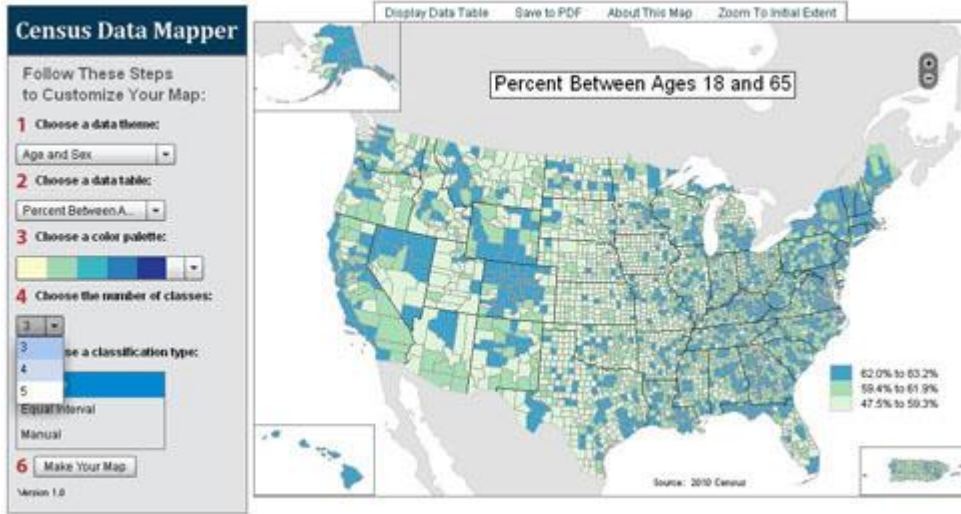
Poor data quality costs the US economy around
**$3.1 TRILLION A YEAR**

Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, QAS
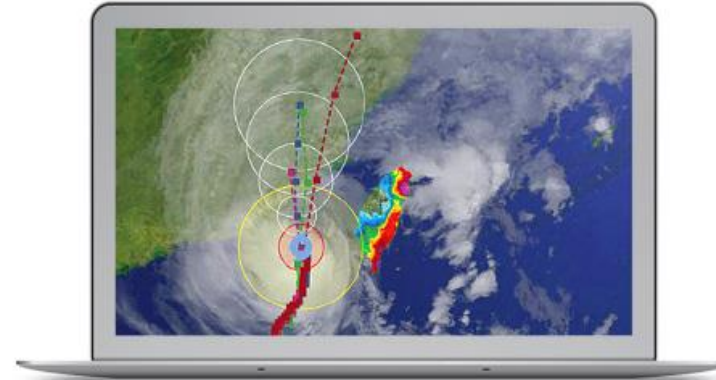
IBM

# Big Data & GIS

# Big Data & GIS



DATA

# Big Data & GIS

Predictive Modeling
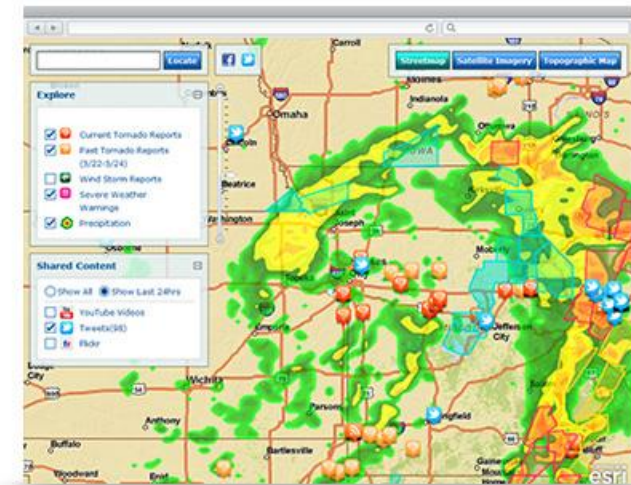
Exposing Patterns

Link with Social Media

Finding Relationships

Big Data & GIS

# Big Data, GIS, & *US*

How do we, as GIS users, leverage the Big Data revolution?

- **Tapping into Big datasets**

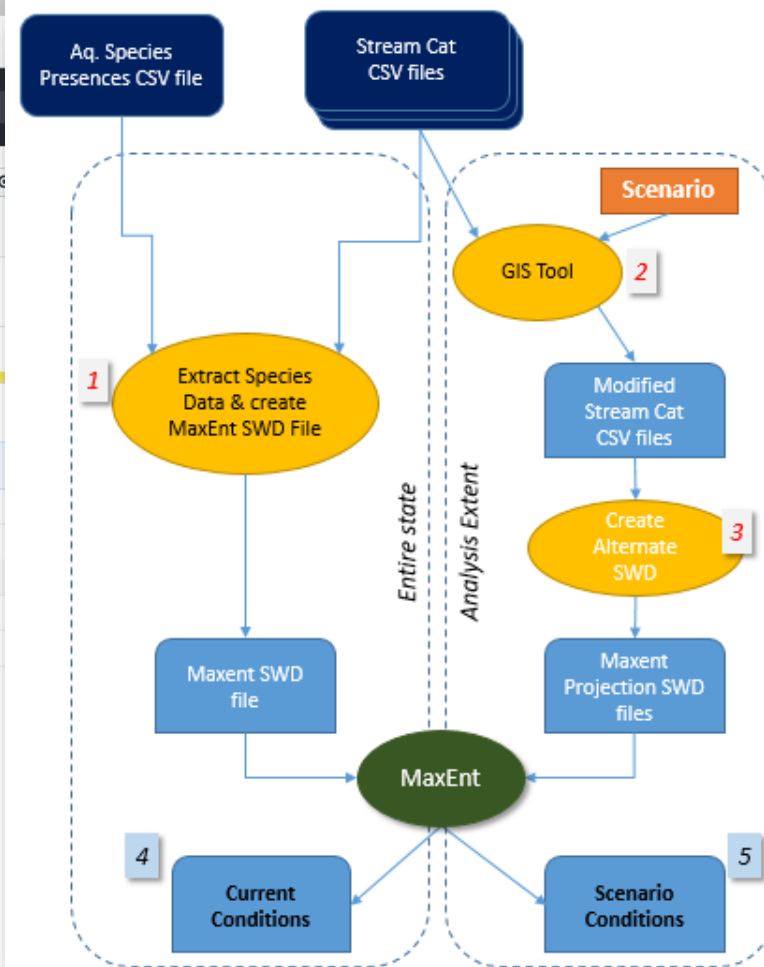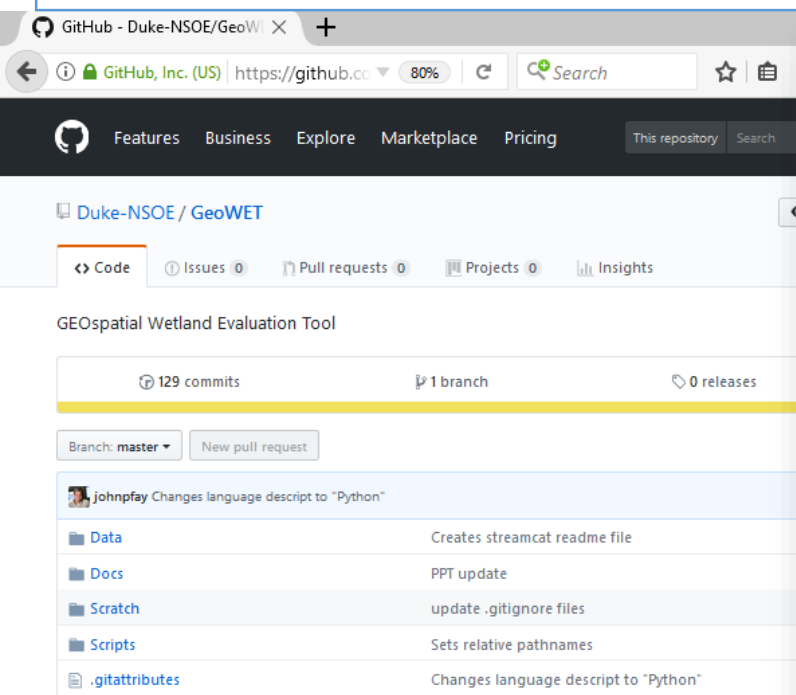  *Today* →

  - Automating data downloads
  - Direct access to Big data

- **Beyond the Desktop**

  - Data processing in the cloud
  - Building lightweight apps

# Automating data download w/

1
- Identifies the HUC8s in which the species was found and extracts all StreamCat catchment records within them.
- Removes any records with missing data where the species was not found and then any attributes with no data where the species was found.
- Removes any attributes with no significant correlation with presence/absence (p > 0.05). Then identifies cross-correlated attributes pairs (r > 0.75) and removes the one with the least correlation with presence/absence.
- Formats columns and column names to suit the MaxEnt species with data (SWD) format.

2
- Allows user to draw a shape on a map reflecting a change in land cover type.
- The user also designates the analysis extent for projecting uplift. This is usually the HUC 6 in which the modification occurs.
- Based on this change, adjusts values in appropriate StreamCat attribute values for affected records.

3
- Converts each set of StreamCat within the analysis extent (e.g. HUC 6) into its own Maxent SWD files that can be used as MaxEnt projection scenarios.

4
- Listing of each catchment and the estimated percent likelihood of finding the species there based on current conditions (unmodified StreamCat values).

5
- Listing of each catchment and the estimated percent likelihood of finding the species there based on altered conditions (modified StreamCat values).

# Automating data download w/

**Why automate?**

- Too time intensive to acquire manually…

- Update or reuse for new data…

- Reproducibility…

- Some data are only available through an Application Programming Interface (API)…

# Tiers of accessing on-line data

- Grabbing static text from a web site via its web address (URL)

https://waterdata.usgs.gov/nc/nwis/water_use?format=rdb&rdb_compression=value&wu_area=County&wu_year=ALL&wu_county=ALL&wu_category=IN&wu_county_nms=--ALL%2BCounties--&wu_category_nms=Industrial

```
#
# File created on 2017-11-03 09:58:50 EDT
# Refresh Date: 2014-12
#
# U.S. Geological Survey
#
# This file contains selected WaterUse data
#
# The data you have secured from the USGS NWISWeb database may include data that have
# not received Director's approval and as such are provisional and subject to revision.
# The data are released on the condition that neither the USGS nor the United States
# Government may be held liable for any damages resulting from its authorized or
# unauthorized use.
#
#  * References to sources of water-use data can be found here. - https://water.usgs.gov/watuse
#
# Search Criteria:
# Year(s)           - ALL
# Area              - County
# County Codes(s)   - ALL
# County Name(s)    - --ALL Counties--
# Category Code(s)  - IN
# Category Name(s)  - Industrial
#
# Columns:
```

# Tiers of accessing on-line data

- Grabbing hosted binary file(s) from a web address

https://www2.census.gov/geo/tiger/TIGER2017/TRACT/

# Tiers of accessing on-line data

- Grabbing hosted binary file(s) from an FTP server

ftp://newftp.epa.gov/EPADataCommons/ORD/NHDPlusLandscapeAttributes/StreamCat/States

# Tiers of accessing on-line data

- Grabbing a table seen on a web page

https://en.wikipedia.org/wiki/World_Happiness_Report#International_rankings

# Tiers of accessing on-line data

- Specialized Python packages

https://pypi.python.org/pypi/census

```
from census import Census
from us import states


c = Census("MY_API_KEY")
c.acs5.get(('NAME', 'B25034_010E'),
           {'for': 'state:{}'.format(states.MD.fips)})
```

# Tiers of accessing on-line data

- Access via API

http://data.neonscience.org/data-api

# Tiers of accessing on-line data

- Grabbing static text or a file from a web site via its web address (URL)
- Bulk downloading files from a web or ftp server

- Grabbing (and converting) data seen on a web page – "Scraping"

- Specialized Python modules for accessing on-line data

- Using Application Programming Interfaces (APIs) to pull data

# Diving in!

- Download the zip file (or sync from GIT)

- Run through examples & discuss what's going on...