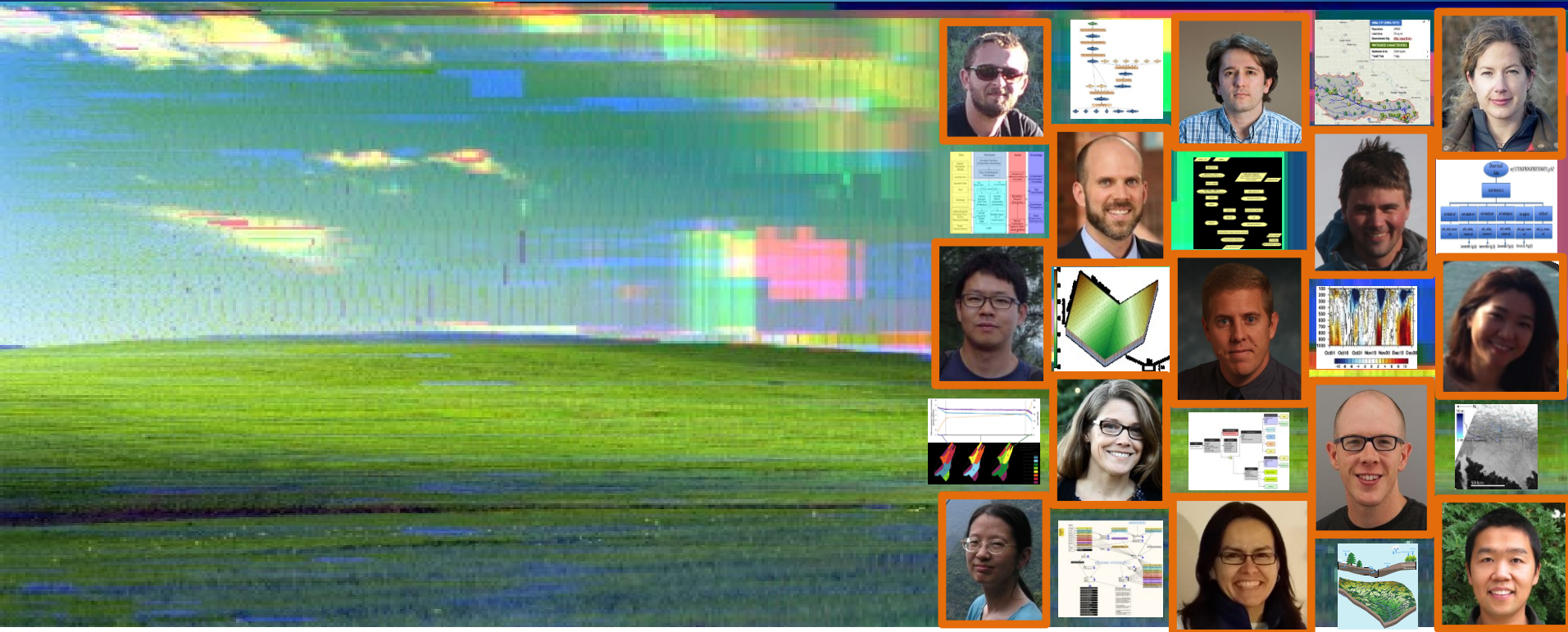# Course: "Introduction to Computational Thinking and Data Science"

Yolanda Gil
gil@isi.edu

# Goals of the Course

Course is **designed for students with no programming background who want to have literacy in computing and data science** to better approach data-rich problems
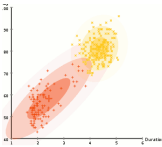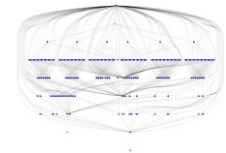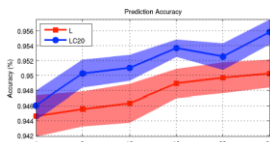
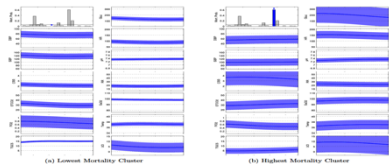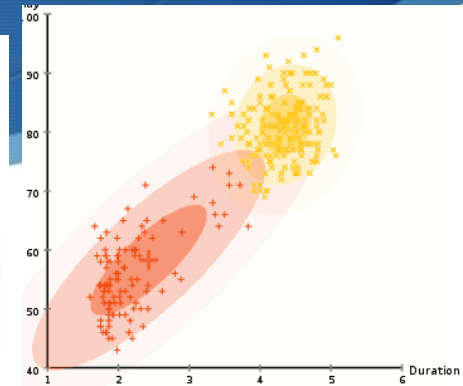# Barriers of Non-Programmers to Data Science

# Distinct Expertise in Data Science
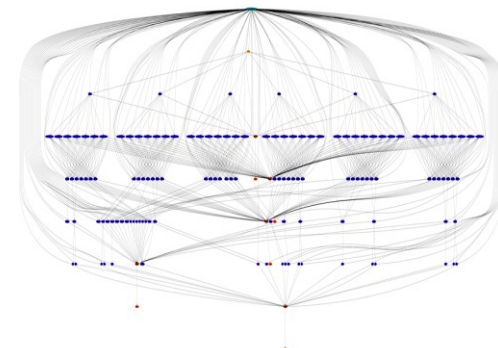
**Domain knowledge**

**Statistics, data mining**

**Semantics and distributed data**

**Large-Scale Data Processing**

# Becoming Data Scientists: Overcoming the Barriers

The goal of the class is to empower non-programmers
to communicate with computer scientists
so they can collaborate in real-world data science projects

# Students Learn to Channel Their Domain Expertise into Data Science

- ◆ USC course attended by graduate students in:
  - ◆ Political sciences, social sciences, education, biology, medicine, engineering

- ◆ Palpable trajectory:
  - ◆ Week 1: Sketch a data science project
    - ◆ Good goals, but impractical, nonsensical, unmanageable
  - ◆ Week 7: Revisit the data science project
    - ◆ Sensible approach, technical vocabulary
  - ◆ Week 12: Revisit again
    - ◆ Practical implementable approach
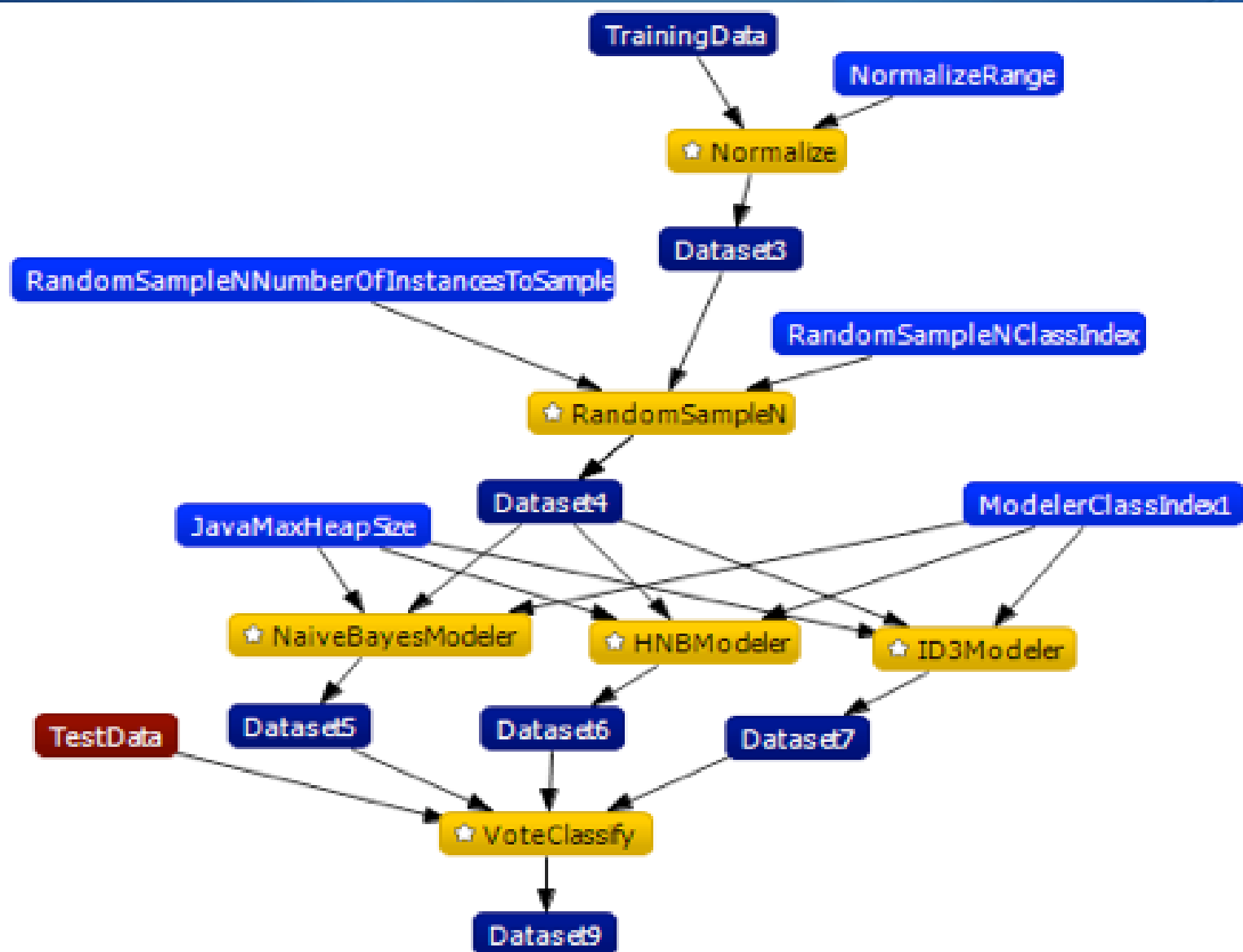
# Design Principles for the Course

## Conceptual Learning

- **Computational thinking**: a new way to approach problems through computing
  - Simulation, data analysis, data mining

- **Data science**: a cross-disciplinary approach to solving data-rich problems
  - Machine learning, large-scale computing, semantic metadata

## Practical Learning

- **Workflows**: a graphical programming environment that enables non-programmers to experiment with complex multi-step data analysis environments

- **Application domains**: exposure to past and ongoing projects where data science exposes multi-disciplinary challenges
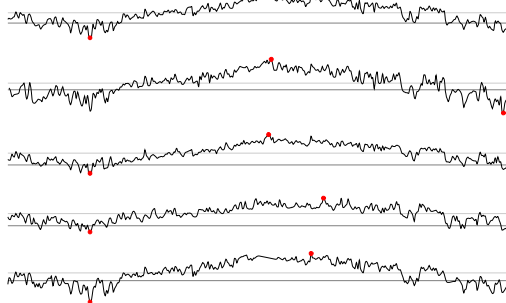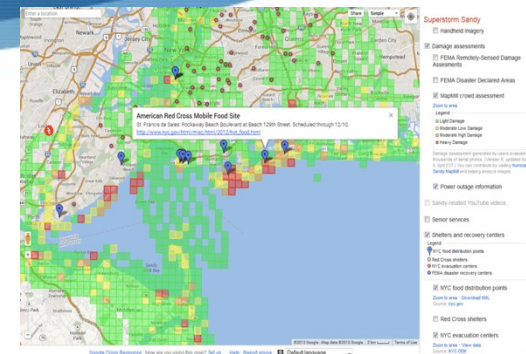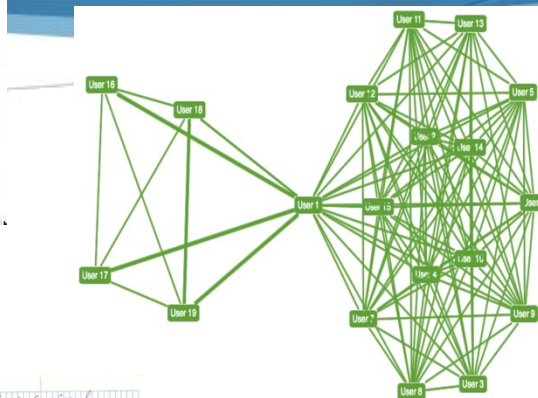  - Social networks, hydrology, proteomics, genomics, medicine, etc.

# Practical Learning:
## Workflows Make Data Science Accessible to Non-Programmers

# Practical Learning:
## Workflows Enable Non-Programmers to Process Real-World Data

# Experiment with Data Science in a Variety of Domains



image processing

hydrology

social network analysis

text analytics

biology

# Section I: Introduction



1. Computational thinking and data science

2. Data
   ◆ What is accessible data
   ◆ Major types of data
   ◆ Basic terminology

3. Data analysis software
   1. Algorithms vs code
   2. Programming languages
   3. Turing machines

4. Multi-step data analysis as workflows

# II: Data Analysis



1. Data analysis tasks
   ◆ Classification, clustering, pattern detection, causal discovery, simulation

2. Data pre-processing

3. Data visualization

4. Data lifecycle

# III: Data Analysis in Practice



1. Analyzing different kinds of data

2. Parallel and distributed computing
   - ◆ Multi-core, clusters, grids, web services, …
   - ◆ Speedup, dependencies, critical paths, Amdahl's law, MapReduce

# IV: Metadata



1. Semantic metadata
   ◆ The DC standard
2. Ontologies
3. Provenance

# V: Data Dissemination



1. Data formats and standards

2. Provenance

3. Data stewardship
   ◆ Data identifiers, data citation

# VI: Advanced Topics



- ◆ Privacy and sensitive data

- ◆ Introduction to databases

- ◆ Crowdsourcing data collection

- ◆ Multi-disciplinary collaboration

- ◆ Project management

# Course Design

## Focus on AI Topics

◆ No statistics

◆ No databases

◆ No programming skills taught

◆ Yes: scalable algorithms

## Constant Practice

◆ Class group activities followed by group reports

　◆ Learn to communicate and use technical terms

◆ Homeworks emphasize hypothesis formulation and testing

# INF549 for USC Informatics Student Comments

◆ *"I went to a big data and big analytics keynote by a SAS VP. The presenter used a lot of jargon vocabulary that we have seen in class, and it was very interesting to see the application of many of the class topics. I understood the talk!"*

◆ *"I attended a Big Data student poster session in Engineering, and I could understand the presentations!"*

The prerequisites for the course were adequate.

| | | |
|---|---|---|
| 1 Poor | 0 | 0.00% |
| 2 Below Average | 0 | 0.00% |
| 3 Average | 1 | 6.67% |
| 4 Above Average | 5 | 33.33% |
| 5 Excellent | 9 | 60.00% |
| Total | 15 | |

Encouraged students to participate in their learning (e.g., through discussion, projects, study groups and other appropriate activities).

| | | |
|---|---|---|
| 1 Poor | 0 | 0.00% |
| 2 Below Average | 0 | 0.00% |
| 3 Average | 0 | 0.00% |
| 4 Above Average | 1 | 6.25% |
| 5 Excellent | 15 | 93.75% |
| Total | 16 | |

# Ongoing

◆ Making materials available at datascience4all.org
  ◆ Include videolectures

◆ Topical tutorials at science meetings (EarthCube, NOAA)
  ◆ Eg, ontologies, machine learning

◆ North American Summer School in Data Science (with Caltech)
  ◆ Already used for the 2016 RDA Summer School in Data Science
  ◆ Already used in 2016 IS-GEO Summer School

# ADDITIONAL SLIDES

# Section I: Introduction to Basic Concepts

1. Computational thinking and data science

2. Data
   - Accessible data
     - APIs, license,…
   - Major types of data
     - Networks, text, time series, geospatial,…
   - Data terms
     - Metadata, silos, sensitive data, big data, …

3. Data analysis software
   - Algorithms vs code
     - Algorithm design
   - Programming languages
   - Encapsulation
   - Turing machines, Turing-complete languages

4. Multi-step data analysis as workflows
   - Components, dataflow, intermediate data
   - **Practicum: WINGS**

# Section II: Data Analysis

1. Data analysis tasks
   - Classification
     - Decision trees
     - Alternative methods
     - Accuracy and other metrics
   - Pattern detection
     - Clustering
     - Temporal patterns
     - Network patterns
   - Causal discovery
     - Graphical models
     - Bayesian networks
   - Simulation

2. Data lifecycle
   - Data pre-processing
     - Data cleaning
   - Data wrangling
   - ETL
   - Collect, clean, analyze, visualize, deposit

3. Data visualization
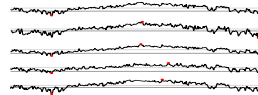   - Time series
   - Statistical
   - Maps, cartograms
   - Treemaps, heatmaps
   - Networks
   - Visual analytics

# Section III:
# Data Analysis in Practice

1. Analyzing different kinds of data
   - Time series data
     - Ecology
     - Medicine
   - Text data
     - Web
     - Archives
   - Network data
     - Social media
     - Web
   - Multimedia data
     - Images
     - Videos
   - Geospatial data

2. Parallel computing
   - Speedup estimates
   - Dependencies
   - Critical paths
   - Amdahl's law
   - MapReduce

3. Distributed computing
   - Multi-core computing, chips
   - Clusters
   - Grids
   - Web services
   - Cloud computing

# Section IV: Metadata

1. Semantic metadata
   - ◆ Attribution metadata
     - ◆ The Dublin Core standard
   - ◆ Summary metadata
   - ◆ Provenance metadata
   - ◆ Metadata in workflows
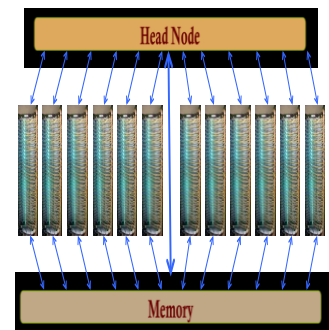     - ◆ Metadata capture
     - ◆ Metadata propagation



2. Ontologies
   - ◆ Taxonomies
   - ◆ Classes
   - ◆ Properties
   - ◆ Assertions
   - ◆ Definitions
   - ◆ Constraints and rules
   - ◆ Reasoning
   - ◆ **Practicum: PROTÉGÉ**



3. Provenance
   - ◆ Process provenance
   - ◆ Attribution provenance
   - ◆ Resource provenance
   - ◆ Provenance standards

# Section V:
# Data Dissemination

1. Data formats
   - ◆ Data standards
   - ◆ Data repositories
   - ◆ Data services
   - ◆ Web data

2. Combining metadata and provenance
   - ◆ Metadata propagation
   - ◆ Automatic method validation
   - ◆ Automatic generation of metadata
   - ◆ Automatic provenance tracking

3. Data stewardship
   - ◆ Data sharing
   - ◆ Data identifiers
   - ◆ Licenses for data
   - ◆ Data citation and attribution
   - ◆ Software publication

# VI: Advanced Topics



- Privacy and sensitive data

- Introduction to databases

- Crowdsourcing data collection

- Multi-disciplinary collaboration

- Project management