

# Introduction to Computational Thinking and Data Science

Yolanda Gil, gil@isi.edu

Information Sciences Institute and Department of Computer Science, University of Southern California

## Course Syllabus

This course will teach non-programmers to think in computing terms about modern topics, and to approach real-world phenomena through data science. The course will enable students to:

- Acquire computational thinking skills that will enable students to represent and reason about complex problems in the digital arena
- Understand different kinds of data in terms of their possibilities and limitations to approach complex problems cast in terms of the emerging field of data science
- Become data science scholars through best practices in data documentation and dissemination

The course is intended for students in disciplines outside of computer science, so no prior experience with computer science is assumed. Students will learn to think in computing terms about modern topics, and to approach real-world phenomena through data science. The course introduces different kinds of data and corresponding approaches to data analysis, including geospatial data, time series, networks, and multimedia data. Students learn to run multi-step analysis through a graphical workflow interface, and will experience first hand complex concepts in data science such as parallel computing, provenance, and visualization. Students also learn to use ontologies and logic representations to capture metadata and other knowledge about complex data. The course includes practical lessons to use workflow and ontology development toolkits, as well as best practices for data stewardship and dissemination.

	Lesson	Topics
<b>Section I: Introduction to Computational Thinking and Data Science</b>		
1	<b>Computational thinking and data science</b>	<ul style="list-style-type: none"><li>• What is computational thinking</li><li>• Computational thinking for reasoning and analysis</li><li>• What is data science</li><li>• Data scientists</li><li>• The context of data science</li></ul>
2	<b>Data</b>	<ul style="list-style-type: none"><li>• What is data</li><li>• What is not (yet) data</li><li>• Time series data</li><li>• Networked data</li><li>• Geospatial data</li><li>• Text data</li><li>• Labeled and annotated data</li><li>• Big data</li></ul>
3	<b>Data analysis software</b>	<ul style="list-style-type: none"><li>• Programs for data analysis</li><li>• Inputs and Outputs</li><li>• Program Parameters</li></ul>

		<ul style="list-style-type: none"> <li>• Programming Languages</li> <li>• Programs as Black Boxes</li> <li>• Algorithms versus software</li> </ul>
<b>4</b>	<b>Multi-step data analysis as workflows</b>	<ul style="list-style-type: none"> <li>• Building workflows by composing software</li> <li>• Pre-processing and post-processing data</li> <li>• Workflows for data analysis</li> <li>• Workflow inputs and parameters</li> <li>• Executing workflows</li> <li>• Exploring data through workflows</li> <li>• Workflows in practice</li> </ul>
<b>5</b>	<b>Workflow practicum</b>	<ul style="list-style-type: none"> <li>• The WINGS workflow system</li> <li>• Workflows in practice</li> </ul>
<b>Section II: Data Analysis</b>		
<b>6</b>	<b>Data analysis tasks (I)</b>	<ul style="list-style-type: none"> <li>• Data analysis tasks in data mining, statistics, and machine learning</li> <li>• Supervised learning <ul style="list-style-type: none"> <li>○ Classification tasks</li> <li>○ Classification algorithms</li> <li>○ Evaluation of classifiers</li> </ul> </li> </ul>
<b>7</b>	<b>Data analysis tasks (II)</b>	<ul style="list-style-type: none"> <li>• Unsupervised learning <ul style="list-style-type: none"> <li>○ Clustering</li> <li>○ Pattern detection</li> <li>○ Anomaly detection</li> </ul> </li> <li>• Simulation and prediction</li> </ul>
<b>8</b>	<b>Data analysis tasks (III)</b>	<ul style="list-style-type: none"> <li>• Causality <ul style="list-style-type: none"> <li>○ Probabilistic graphical models</li> <li>○ Bayesian networks</li> <li>○ Causal models</li> </ul> </li> </ul>
<b>9</b>	<b>Data pre-processing</b>	<ul style="list-style-type: none"> <li>• Data cleaning</li> <li>• Quality control</li> <li>• Data integration</li> <li>• Feature selection</li> <li>• Feature construction</li> </ul>
<b>10</b>	<b>Data lifecycle</b>	<ul style="list-style-type: none"> <li>• Data collection</li> <li>• Data storage</li> <li>• Data extraction and querying</li> <li>• Data integration</li> <li>• Data presentation</li> </ul>
<b>11</b>	<b>Data visualization</b>	<ul style="list-style-type: none"> <li>• Quality of visualizations</li> <li>• Major types of visualizations</li> <li>• Time series visualizations</li> </ul>

		<ul style="list-style-type: none"> <li>• Geospatial visualizations</li> <li>• Multi-dimensional spaces</li> <li>• Network visualizations</li> </ul>
<b>Section III: Data Analysis in Practice</b>		
<b>12</b>	<b>Analyzing different kinds of data (I)</b>	<ul style="list-style-type: none"> <li>• Analyzing text data <ul style="list-style-type: none"> <li>○ Pre-processing text</li> <li>○ Document classification</li> <li>○ Document clustering</li> <li>○ Topic detection</li> <li>○ Sentiment analysis</li> </ul> </li> </ul>
<b>13</b>	<b>Analyzing different kinds of data (II)</b>	<ul style="list-style-type: none"> <li>• Analyzing time series data <ul style="list-style-type: none"> <li>○ Collecting time series data</li> <li>○ Pre-processing time series data</li> <li>○ Event detection</li> <li>○ Granger causality</li> </ul> </li> </ul>
<b>14</b>	<b>Analyzing different kinds of data (III)</b>	<ul style="list-style-type: none"> <li>• Analyzing network data <ul style="list-style-type: none"> <li>○ Network structure</li> <li>○ Dynamic networks</li> <li>○ Scale-free networks</li> <li>○ Network analysis</li> </ul> </li> </ul>
<b>15</b>	<b>Analyzing different kinds of data (IV)</b>	<ul style="list-style-type: none"> <li>• Analyzing multimedia data <ul style="list-style-type: none"> <li>○ Pre-processing images</li> <li>○ Segmentation</li> <li>○ Edge detection</li> <li>○ Object detection</li> <li>○ Video analysis</li> </ul> </li> <li>• Analyzing geospatial data <ul style="list-style-type: none"> <li>○ Coordinate systems</li> <li>○ GIS systems</li> </ul> </li> </ul>
<b>16</b>	<b>Parallel and distributed computing for big data (I)</b>	<ul style="list-style-type: none"> <li>• Cost of computation</li> <li>• Divide and conquer</li> <li>• Speedup with parallel processing</li> <li>• Limits of speedup: Critical path</li> <li>• Amdahl's law</li> <li>• When problems are not parallelizable</li> </ul>
<b>17</b>	<b>Parallel and distributed computing for big data (II)</b>	<ul style="list-style-type: none"> <li>• Multi-core computing</li> <li>• Distributed computing</li> <li>• Cluster computing</li> <li>• Cloud computing</li> <li>• Grid computing</li> <li>• Virtual machines</li> <li>• Web services</li> </ul>

		<ul style="list-style-type: none"> <li>• Practical concerns in distributed computing</li> <li>• Parallel programming languages <ul style="list-style-type: none"> <li>○ MapReduce/Hadoop</li> </ul> </li> </ul>
<b>Section IV: Metadata</b>		
<b>18</b>	<b>Semantic metadata</b>	<ul style="list-style-type: none"> <li>• What is metadata</li> <li>• Basic metadata versus semantic metadata</li> <li>• Metadata about data collection</li> <li>• Metadata about data processing</li> <li>• Metadata for search and retrieval</li> <li>• Metadata standards</li> <li>• Domain metadata and ontologies</li> </ul>
<b>19</b>	<b>Ontologies (I)</b>	<ul style="list-style-type: none"> <li>• What is an ontology</li> <li>• Taxonomies and class inheritance</li> <li>• Properties</li> <li>• Logical constraints</li> </ul>
<b>20</b>	<b>Ontologies (II)</b>	<ul style="list-style-type: none"> <li>• Logical reasoning and inference</li> <li>• Expressivity and computation</li> <li>• The Semantic Web</li> </ul>
<b>21</b>	<b>Ontologies (III)</b>	<ul style="list-style-type: none"> <li>• Practicum: the PROTÉGÉ ontology editor</li> </ul>
<b>22</b>	<b>Provenance</b>	<ul style="list-style-type: none"> <li>• What is provenance</li> <li>• Provenance concerning objects</li> <li>• Provenance concerning people and institutions</li> <li>• Provenance concerning processes</li> <li>• Provenance models</li> <li>• Provenance standards</li> </ul>
<b>Section V: Data Dissemination</b>		
<b>23</b>	<b>Data formats and standards</b>	<ul style="list-style-type: none"> <li>• Data formats</li> <li>• Data standards</li> <li>• Data repositories</li> <li>• Data services</li> <li>• The Semantic Web and linked open data</li> </ul>
<b>24</b>	<b>Tracking metadata and provenance</b>	<ul style="list-style-type: none"> <li>• Combining computation with metadata and provenance</li> <li>• Validating a data analysis method</li> <li>• Tracking provenance during data analysis</li> <li>• Automatically generating metadata for data analysis</li> </ul>
<b>25</b>	<b>Data stewardship</b>	<ul style="list-style-type: none"> <li>• Data sharing</li> <li>• Data identifiers</li> <li>• Licenses for data</li> <li>• Data citation and attribution</li> <li>• Software and other work products</li> </ul>

<b>Section VI: Advanced Topics</b>		
<b>26</b>	<b>Privacy and ethics in data science</b>	<ul style="list-style-type: none"> <li>• Privacy               <ul style="list-style-type: none"> <li>○ Fair Information Practices</li> <li>○ Managing sensitive data</li> <li>○ Anonymizing sensitive data, k-anonymity, differential privacy</li> <li>○ Re-identifying datasets</li> </ul> </li> <li>• Reproducibility</li> <li>• Societal value of data and data science</li> </ul>
<b>27</b>	<b>Databases</b>	<ul style="list-style-type: none"> <li>• File systems vs databases</li> <li>• Relational databases               <ul style="list-style-type: none"> <li>○ Data models</li> <li>○ SQL</li> <li>○ Transactions</li> </ul> </li> <li>• NoSQL databases</li> </ul>
<b>28</b>	<b>Multidisciplinary collaborations</b>	<ul style="list-style-type: none"> <li>• Discipline-specific data, cultures, and methodologies</li> <li>• Identifying common goals</li> <li>• Developing shared terminology</li> <li>• Structuring and focusing discussions</li> <li>• Use cases and example scenarios</li> <li>• Project planning</li> </ul>