



MATH 571/CSP 571

Data Preparation & Analysis

Class Hours: Wednesday 6:45PM - 9:35PM

Class Location: Hybrid (Synchronous): PS 111 + Internet



Jawahar J. Panchal

jpanchal@iit.edu / (312) 567-5871

Office Hours: TBD/Virtual

Office Location: SB 228A

*This is a tentative syllabus - any subsequent changes will be communicated to students in class.*

**Course Description:** This course will provide students with an introduction to the field of data science and the tools and techniques to analyze data and extract knowledge from it. This course surveys industrial and scientific applications of data analytics and will utilize case studies and programming exercises to explore opportunities and challenges involving data analysis and visualization including business opportunities, privacy concerns, and ethical issues. Students will work with a variety of real world data sets and learn how to prepare data sets for analysis by cleaning and reformatting.

**Prerequisite(s):** CS 425 or equivalent. Math 474 or equivalent. Graduate standing or permission of instructor.

**Note(s):** This is an *elective* course for MATH, CS, and CSP majors.

**Credit Hours:** 3 (3-0-3)

**Required Text(s):** *Introduction to Statistical Learning*, 2<sup>nd</sup> Edition (Online Edition) [Free]

**Author(s):** Gareth James, Daniela Witten, Trevor Hastie, Robert Tibshirani

**ISBN-13:** 978-1461471370

**Required Text(s):** *R for Data Science*, 1<sup>st</sup> Edition (Online Edition) [Free]

**Author(s):** Hadley Wickham, Garrett Grolemund

**ISBN-13:** 978-1491910399

### Course Objectives:

This course provides students with the knowledge and skills to effectively analyze data to unearth knowledge including the ability to develop programs that support the analysis and presentation processes.

### Course Outcomes:

Students successfully completing this course will be able to:

1. Discuss the concepts of data science
2. Gather requirements for data analysis projects
3. Prepare data for analysis
4. Programmatically analyze data
5. Develop meaningful data visualizations
6. Explore case studies

7. Examine ethical issues in data science
8. Report on findings of their analysis
9. Develop and present analysis effectively

### Course Details:

The course will attempt to incorporate the following tools as part of the course material:

1. Use of the R programming language as a general-purpose data-analysis programming environment.
2. Use of the CRAN package ecosystem as a source of software libraries for extending/supplementing the R environment.
3. Use of various open data sources for projects and analytical work.
4. Use of various cloud provider systems for computational/modeling work.

### Grading System:

Grade weighting is prescribed as below, with significant improvement between Midterm Exam to Final Exam scores taken into consideration in the final grade assignment.

Homework Assignments	20%
Reading Quizzes	20%
Semester Project	20%
Midterm Exam	20%
Final Exam	20%

### Letter Grade Distribution:

Grade assignment is prescribed as below, with the student score rounded up if within 1%-2% of the next grade level, at the instructor's discretion.

$\geq 90$	A
80 - 89	B
70 - 79	C
60 - 69	D
$\leq 60$	E

### Teaching Assistant:

The TA for this course is **Akshara Kudumula** - who will be available for questions regarding grading on assignments, exams, quizzes, and the project. The TA can be contacted at akudumula@hawk.iit.edu

## Course Policies:

- **General**

- Each course session will follow a lecture format with slides and examples - students are expected to take notes and participate in questions/answers during the session.
- There will be a short 15 minute break, at the instructor's discretion, during each session - students are expected to return back to the classroom on time.

- **Reading Requirements & Homework Assignments**

- Students are expected to complete reading assignments per the course schedule - homework assignments will be assigned by the instructor with details on submission deadlines.
- **Late assignments submitted within one week of the due date will be subject to a full grade penalty - no late assignments will be accepted beyond a one week time period.**

- **Online Access & Electronic Submissions**

- **Lectures will be recorded and made available online**, however attendance is strongly encouraged in order to understand the course material.
- Students are expected to submit electronic documents for their homework assignments via the IIT Blackboard system.

- **Examinations**

- Exams are closed book, closed notes - the Midterm Exam will be administered in class, the Final Exam will be administered according to university final exam schedules.
- **No makeup exams will be given except in extreme circumstances/emergency situations, subject to department and university approval.**

- **Grades**

- Student grades will be posted on Blackboard in a timely manner - students who wish to track their progress may do so online, or inquire with the instructor.
- Any questions or discussions regarding grades should be directed to the instructor,

- **Attendance and Absences**

- Attendance is expected for each session - students may contact the instructor regarding any missed classes due to sickness, emergencies, or other issues.
- Students are responsible for all missed material, regardless of the reason for absence, and are expected to obtain notes/content independently.

## University Policies:

- **General**

Students should refer to the Illinois Tech Student Handbook as a reference to any and all policies listed below pertaining to this course.

- **Academic Honesty**

Students are subject to the Code of Academic Honesty as part of being enrolled in this course. Issues related to academic honesty within this course will be handled according to university policies, regulations, and procedures.

- **Code of Conduct**

Students are subject to the Code of Conduct as part of being enrolled in this course. Issues related to conduct within this course will be handled according to university policies, regulations, and procedures.

- **Special Accommodations**

Students requiring special accommodations, such as in the case of documented disabilities, should contact the Center for Disability Resources. Accommodations will be arranged via the Reasonable Accommodations process.

### Tentative Course Outline:

Material coverage, lecture order, and content timing may change - the student is expected to maintain independent progress regarding reading requirements and homework assignments.

Week	Content
Week 1 08/25/2021	Statistical Learning - Descriptive/Inferential Statistics <ul style="list-style-type: none"><li>• Exploratory Data Analysis - Empirical/Parametric Dist., Visualization</li><li>• Statistical Modeling - Model Accuracy, Error, Bias/Variance</li><li>• <b>Suggested Reading:</b> James Ch1,Ch2; Wickham Ch1,Ch2-8</li></ul>
Week 2 09/01/2021	Linear Regression <ul style="list-style-type: none"><li>• Linear Regression - Simple/Multiple, Estimation</li><li>• Model Accuracy - Diagnostics/Validation, Additional Considerations</li><li>• <b>Suggested Reading:</b> James Ch3; Wickham Ch2-8</li><li>• <b>Homework Assignment:</b> Homework 1 Assigned</li></ul>
Week 3 09/08/2021	Classification <ul style="list-style-type: none"><li>• Logistic Regression - Generalized Linear Model/Least Squares (GLM/GLS)</li><li>• Linear Discriminant Analysis (LDA) - Bayes Theorem</li><li>• <b>Suggested Reading:</b> James Ch4; Wickham Ch9-16</li></ul>
Week 4 09/15/2021	Resampling & Cross-Validation <ul style="list-style-type: none"><li>• Cross-Validation: Leave-One-Out, k-Fold</li><li>• Bootstrap: Overview</li><li>• <b>Suggested Reading:</b> James Ch5; Wickham Ch9-16</li><li>• <b>Homework Assignment:</b> Homework 1 Due, Homework 2 Assigned</li></ul>
Week 5 09/22/2021	Regularized Regression <ul style="list-style-type: none"><li>• Model Selection: Subset Selection, Optimal Models</li><li>• Shrinkage Methods: Ridge, Lasso</li><li>• <b>Suggested Reading:</b> James Ch6; Wickham Ch9-16</li><li>• <b>Project Deliverable:</b> Project Group &amp; Topic Form Due</li></ul>
Week 6 09/29/2021	Non-Linear Models <ul style="list-style-type: none"><li>• Splines: Piecewise, Polynomial, Smoothing</li><li>• GAM: Regression, Classification</li><li>• <b>Suggested Reading:</b> James Ch7; Wickham Ch17-21</li><li>• <b>Homework Assignment:</b> Homework 2 Due, Homework 3 Assigned</li></ul>
Week 7 10/06/2021	Supervised Learning - Decision Trees <ul style="list-style-type: none"><li>• Classification/Regression: CART, Additional Considerations</li><li>• Ensemble Methods: Random Forests, Boosting/Bagging</li><li>• <b>Suggested Reading:</b> James Ch8; Wickham Ch17-21</li></ul>
Week 8 10/13/2021	Support Vector Machines <ul style="list-style-type: none"><li>• Maximal Margin Classification: Hyperplanes, Separability</li><li>• SVM: Classification, Decision Boundaries</li><li>• <b>Suggested Reading:</b> James Ch9; Wickham Ch22-25</li><li>• <b>Homework Assignment:</b> Homework 3 Due</li><li>• <b>Project Deliverable:</b> Project Proposal &amp; Outline Due - <i>Presentation</i></li><li>• <b>Quiz Assessment:</b> Quiz 1</li></ul>
Week 9 10/20/2021	<b>Midterm Exam</b> <ul style="list-style-type: none"><li>• <i>In Class/Online Examination</i></li></ul>

Week	Content
Week 10 10/27/2021	Deep Learning - Neural Networks <ul style="list-style-type: none"> <li>• ANN: Single/Multiple Layer, Use Cases</li> <li>• CNN/RNN: Fitting, Backpropagation</li> <li>• <b>Suggested Reading:</b> James Ch10; Wickham Ch22-25</li> <li>• <b>Homework Assignment:</b> Homework 4 Assigned</li> </ul>
Week 11 11/03/2021	Survival Analysis - Censored Data <ul style="list-style-type: none"> <li>• Survival: Curves, Log-Rank, Hazards</li> <li>• Shrinkage: Cox Model, AUC</li> <li>• <b>Suggested Reading:</b> James Ch11; Wickham Ch22-25</li> <li>• <b>Project Deliverable:</b> Project Plan &amp; Detail Due - <i>Presentation</i></li> </ul>
Week 12 11/10/2021	Unsupervised Learning - Dimensionality Reduction/Clustering <ul style="list-style-type: none"> <li>• Dimensionality Reduction: Principal Component Analysis (PCA)</li> <li>• Clustering: K-Means, Mixture Models (GMM)</li> <li>• <b>Suggested Reading:</b> James Ch12; Wickham Ch22-25</li> <li>• <b>Homework Assignment:</b> Homework 4 Due, Homework 5 Assigned</li> </ul>
Week 13 11/17/2021	Multiple Testing <ul style="list-style-type: none"> <li>• Confirmatory Data Analysis: Hypothesis Testing, ANOVA</li> <li>• Error Rates: Type I/Type II, Family-Wise, False Discovery</li> <li>• <b>Suggested Reading:</b> James Ch13; Wickham Ch22-25</li> </ul>
Week 14 11/24/2021	<b>Thanksgiving Break</b> <ul style="list-style-type: none"> <li>• <i>No Class</i></li> <li>• <b>Homework Assignment:</b> Homework 5 Due, Homework 6 Assigned</li> </ul>
Week 15 12/01/2021	Catch-Up & Review <ul style="list-style-type: none"> <li>• <b>Project Deliverable:</b> Project Presentation &amp; Report Due - <i>Presentation</i></li> <li>• <b>Quiz Assessment:</b> Quiz 2</li> </ul>
Week 16 12/08/2021	<b>Final Exam</b> <ul style="list-style-type: none"> <li>• <i>In Class/Online Examination</i></li> <li>• <b>Homework Assignment:</b> Homework 6 Due</li> </ul>