



# UNIVERSITY of HOUSTON

## COSC 3337 Data Science I

### Course Information

Term and Year: **Spring 2022**  
Location: **PGH 232 – Face to Face (No recordings of lectures)**  
Meeting Days/Times: Tuesdays -Thursday 11:30 AM -1:00 PM

**Contact:** By email, MS Teams during office hours (or by appointment).

**Office Hours:** 1:00 PM- 1:30 PM TTH.

**Course Online System:** Blackboard.

**Main References:** While lecture notes will serve as the main source of material for the course, the following book constitutes a great reference:

### Open Textbooks

**Rizk, Nouhad: Building Skills for Data Science**

<https://uhlibraries.pressbooks.pub/buildingskillsfordatascience/>

### Books

1. <https://ebookcentral.proquest.com/lib/uh/detail.action?docID=1895687&query=data+mining>
2. <https://ebookcentral.proquest.com/lib/uh/detail.action?docID=4851656>

### Statistics:

3. <https://cnx.org/contents/tWu56V64@33.122:-mZCQZc7@5/Introduction>

### Reference:

P.-N. Tang, M. Steinback, and V. Kumar Introduction to Data Mining, Addison Wesley, 2018.

(Cathy O'Neil and Rachel Schutt. Doing Data Science, Straight Talk from the Frontline. O' Reilly. 2014

**Description:** Data science process, data preprocessing, exploratory data analysis, data visualization, basic statistics, basic machine learning concepts, classification and prediction, similarity assessment, clustering, post-processing and interpreting data analysis results, use of data analysis tools and programming languages and data analysis case studies.

**Objectives:** By the end of the course a successful student should:

- Students will develop relevant programming abilities.
- Students will demonstrate proficiency with statistical analysis of data.
- Students will develop the ability to build and assess data-based models.
- Students will execute statistical analyses with Python software.

- Students will apply data science concepts and methods to solve problems in real-world contexts and will communicate these solutions effectively

**Prerequisites:** MATH 3339 and COSC 2436.

**Software:** Make sure to download Anaconda <https://repo.anaconda.com/>. Let me know via email in case you encounter difficulties.

**Academic Honesty:** University of Houston students are expected to adhere to the Academic Honesty Policy as described in the UH Undergraduate Catalog. “Academic dishonesty” means employing a method or technique or engaging in conduct in an academic endeavor that contravenes the standards of ethical integrity expected at the University of Houston or by a course instructor to fulfill any and all academic requirements. Academic dishonesty includes, but is not limited to, the following: Plagiarism; Cheating and Unauthorized Group Work; Fabrication, Falsification, and Misrepresentation; Stealing and Abuse of Academic Materials; Complicity in Academic Dishonesty; Academic Misconduct.

Refer to UH Academic Honesty website (<http://www.uh.edu/provost/policies/honesty/>) and the UH Student Catalog for the definition of these terms and university’s policy on Academic Dishonesty. Anyone caught cheating will be reported to the department for further disciplinary actions, receive sanctions as explained on these documents, and will have an academic dishonesty record at the Provost’s office. The sanctions for confirmed violations of this policy shall be commensurate with the nature of the offense and with the record of the student regarding any previous infractions. Sanctions may include, but are not limited to a lowered grade, failure on the examination or assignment in question, failure in the course, probation, suspension, or expulsion from the University of Houston, or a combination of these. Students may not receive a W for courses in which they have been found in violation of the Academic Honesty Policy. If a W is received prior to a finding of policy violation, the student will become liable for the Academic Honesty penalty, including F grades.

**Technology statement** below as requested by the Provost’s Office:

Computer and internet access required for course. For the current list of minimum technology requirements and resources, copy/paste/navigate to the URL <http://www.uh.edu/online/tech/requirements>. For additional information, contact the office of Online & Special Programs at [UHOnline@uh.edu](mailto:UHOnline@uh.edu) or 713-743-3327.

	Date	Topics	Open Textbook Reading
Week 1	Tuesday, January 18, 2022	<b>Introduction to Data science</b>	
	Thursday, January 20, 2022	<b>Data science Overview</b>	

Week 2	Tuesday, January 25, 2022	<b>Machine Learning Data Cleaning</b>	
	Thursday, January 27, 2022	<b>Data Processing Startup Example</b>	B1: p 30-35
Week 3	Tuesday, Feb 1, 2022	<b>Statistical Learning</b>	
	<b>Wednesday February 2<sup>nd</sup></b>	<b><u>DROP DEADLINE</u></b>	
	Thursday Feb 3 <sup>rd</sup> , 2022	<b>Data Exploration Data Similarities &amp; Distances</b>	B1: p 54-81
Week 4	Tuesday, Feb 8, 2022	<b>Linear Regression</b>	B1:p 171-213
	Thursday, Feb 10, 2022	<b>Linear Regression (Python Example)</b>	
Week 5	Tuesday, Feb 15, 2022	<b>Logistic Regression Dimensionality reduction - PCA</b>	B1: p 359-399
	Thursday, Feb 17, 2022	<b>Introduction to Classification KNN</b>	B1: p 301-312 B2: p 32-48
Week 6	<b>Tuesday, Feb 22, 2022</b>	<b>Exam 1</b>	
	Thursday, Feb 24, 2022	<b>Decision Tree</b>	
Week 7	Tuesday, Mar 1 <sup>st</sup> , 2022	<b>Random Forests KNN</b>	B1: p 317-322 B2: P 49-68

	Thursday, Mar 3 <sup>rd</sup> , 2022	<b>Naive Bayes</b>	B1: p 414-439 B2: p 113-140
Week 8	Tuesday, Mar 8, 2022	<b>Model Evaluations Metrics</b>	
	Thursday, Mar 10, 2022	<b>Ridge - Lasso</b>	
	<b>Spring break 14-19</b>		
Week 9	Tuesday, Mar 22, 2022	<b>Lines/SVM</b>	

	Thursday, Mar 24, 2022	<b>Dimensionality reduction (feature extraction)</b> <b>Wrap Up classification</b>	
week 10	<b>Tuesday, March 29, 2022</b>	<b>Exam 2</b>	
	Thursday, March 31, 2022	<b>K-Means</b>	B1: p 523- 537 B2: 218-250

Week 11	Tuesday, April 5, 2022	<b>Hierarchical Clustering Heatmap</b>	
	<b>Tuesday April 5<sup>th</sup></b>	<b>DROP DEADLINE</b>	
	Thursday, April 7, 2022	<b>Storytelling</b>	
Week 12	Tuesday, April 12, 2022	<b>DBSCAN</b>	
	Thursday, April 14, 2022	<b>Cluster Validity Silhouette</b>	
Week 13	Tuesday, April 19, 2022	<b>Neural networks</b>	
	Thursday, April 21, 2022	<b>Apriori and Association rules</b>	B1: p 603- 617 B2: p 69-87
Week 14	Tuesday, April 26, 2022	<b>Dynamic Hashing -Merkle tree (Optional)</b>	
	Thursday, April 28, 2022		
	<b>Monday May 2<sup>nd</sup> , 2022</b>	<b>Last day of class</b>	
		<b>Final Exam @</b>	

## Grading Policy

The final numeric grade is computed based on student's performance in weekly assignments and exams/quizzes. The final numeric grade for the course will be determined as follows:

✓ Homework assignments ( <b>NO drop of any HW</b> )	25%
✓ Lab work /Workbook (drop the lowest)	20%
✓ Exam 1 (Tuesday 2/22)	15%
✓ Exam 2 (Thursday 3/29)	15%
✓ Final Exam	25%

**Labs (potentially):** Coding practices (using Python format. ipynb **only**) held sometimes during class times. **One lab assignment will be dropped** (the one with the lowest grade).

**Exams:** Held during class times.

**Homework:** Four assigned HomeWorks. Topics: Regression and Classification (Week 4); decision Tree and KNN(Week 7); SVM and dimensionality(Week 10); and Clustering with cluster validity(Week 12). Students will submit their written homework by scanning and uploading their work in Blackboard (or as .ipynb).

## Final Group Project on Storytelling (as final Homework):

- You will form a group of 3-4 members.
- A group assignment, consisting of students teaming up (5 points), deciding on the data set of interest (5 points), posing research questions (5 points) and applying ML techniques to address those questions (35 points). Each group will eventually submit a report/online presentation of research findings and member contributions.

## Grading Scheme:

<b>A&gt;=92.5 Excellent</b>	<b>A-&gt;= 89.5 and &lt; 92.5 Outstanding</b>	<b>B+&gt;=86.5 and &lt; 89.5 Very Good</b>
<b>B &gt; = 83.5 and &lt;86.5 Good</b>	<b>B-&gt;=79.5 and &lt; 83.5 Above Average</b>	<b>C+&gt;=76.5 and &lt; 79.5 High Average</b>
<b>C&gt;=72.5 and &lt;76.5 Average</b>	<b>C-&gt;=69.5 and &lt;72.5 Low Average</b>	<b>D+&gt;=65.5 and &lt;69.5 Below Average</b>
<b>D &gt;=62.5 and &lt;65.5 Poor</b>	<b>F &lt; 62.5 Failing</b>	