CSCE 581-1: Trusted AI

Prof. Biplav Srivastava, Spring 2026

Quiz 1 / Feb 10, 2026/ Instructions

- Create a sub-folder called "Quiz1" inside your GitHub and submit your response. For each question, feel free to create sub-folders. Have a .pdf of this file (Quiz1-CSCE581.pdf) as well for easy reference.
- Return by committing answer to quiz in GitHub by midnight of <mark>Monday, Feb 16, 2026</mark>. Also confirm when done (date/time) in column H of spreadsheet - https://docs.google.com/spreadsheets/d/1ADt7SQe2BqvxNK6nzqX1sSXUQVq_wLtJdxYUZI3XImk/edit?usp=sharing. (We will use GitHub time as actual in case of discrepancy).
- You can (optionally) confirm Quiz completed by sending email to biplav.s@sc.edu.
- Ask any question in class, by Blackboard message or by email to instructor AND TA. Or, come to office hours to clarify doubts.

Total points = 100, Obtained =

**Fill** Student Name: Matthew Bojanowski

**Fill** GitHub link with code in a sub-dir called "Quiz2":

---

**Q1: Understanding of Fairness Issues** [30 points]

Correctional Offender Management Profiling for Alternative Sanctions (COMPAS) is a commercial case management and decision support tool. The COMPAS software uses an algorithm to assess potential recidivism risk (risk of being repeat offender). There is concern that a classification method used for this purpose may be biased or not.

**Task**: Run any two classification models on the COMPASS data and report your findings in not more than a page. Submit the code (notebook) and report

- Suggested GitHub to use with data, description, and sample code to run analysis: https://github.com/tsotne95/FairnessCompas/
- Official Propublica resources
    o Data: https://github.com/propublica/compas-analysis
    o Article: https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing

**Q2:  Understanding  of AI/ data science / classification in a sustainability domain**

**Water treatment water data and pH value**
[10 + 10 + 10 = 30 points]

**Background:**
pH is a very important determinant of water quality. However, its safety limits depends on water purpose.

pH considerations:
- EPA: https://www.epa.gov/caddis-vol2/caddis-volume-2-sources-stressors-responses-ph
- Standards collated: https://github.com/biplav-s/water-info/blob/master/dataWaterParameters.json
- Common practice for limit is: within 6.5-8.5 is considered safe, <= 6.5 and > 8.5 is considered unsafe
  - Example: https://www.safewater.org/fact-sheets-1/2017/1/23/tds-and-ph

**Datasets:**
- **Data:** Weka comes with water treatment data.
  - **Description:** https://archive.ics.uci.edu/ml/datasets/water+treatment+plant
  - **Local cache:** https://github.com/biplav-s/course-tai/tree/main/sample-code/common-data/water-weka
  - Consider the following parameters.


Q-E (input flow to plant)
2 ZN-E (input Zinc to plant)

3 PH-E (input pH to plant)
4 DBO-E (input Biological demand of oxygen to plant)
5 DQO-E (input chemical demand of oxygen to plant)
6 SS-E (input suspended solids to plant)
7 SSV-E (input volatile suspended solids to plant)
8 SED-E (input sediments to plant)
9 COND-E (input conductivity to plant)

23 PH-S (output pH)
24 DBO-S (output Biological demand of oxygen)
25 DQO-S (output chemical demand of oxygen)
26 SS-S (output suspended solids)
27 SSV-S (output volatile supended solids)
28 SED-S (output sediments)
29 COND-S (output conductivity)

**Things to do:**
**1. Data exploration:** Find correlation between input and output parameter values. Example: pH-E and pH-S.
**2. Data preparation:** Add a new column called 'SAFE-PH-S'. It is 'yes' if pH is within 6.5-8.5 and 'no' otherwise, i.e., <= 6.5 and > 8.5
**3. Train**: Train a classifier to predict SAFE-PH-S using any two classification methods. Show its performance measures.
* Use 20% data for testing
* Use any standard validation method (leave one out, 10-fold cross validation)

**Q3:  Recent water data and pH value**
[10 + 10 + 20 = 40 points]

- **Data : Multi-location data**

**Datasets:** We will again look at water data from Florida for WaterAtlas project.
Website: https://orange.wateratlas.usf.edu/

**Data:** Local cache of data
https://github.com/biplav-s/course-tai/blob/main/sample-code/common-data/water/WaterAtlas-ManySites.csv

**Things to do:**
**1. Data preparation:** Make a subset which only refers to pH data. Add a new column called 'SAFE-PH'. It is 'yes' if pH is within 6.5-8.5 and 'no' otherwise, i.e., <= 6.5 and > 8.5
**2. Train**: Train a classifier to predict SAFE-PH using any two classification methods. Show its performance measures.
* Use 20% data for testing
* Use any standard validation method (leave one out, 10-fold cross validation)
**3. Explain**: Which places have the most unsafe water (by pH) and which least by occurrence?

Show them on a map using latitude longitude information available in each row.
Instructions for Google Earth are at: https://www.google.com/earth/outreach/learn/visualize-your-data-on-a-custom-map-using-google-my-maps/