

# Data Science Syllabus

High School - One Semester (85 hours)

## Course Overview and Goals

Industries of all types are hiring data scientists to analyze and highlight the hidden patterns in data. This course equips students with the essential skills of a data scientist which include data collection, cleanup, transformation, analysis, and visualization. Students will write algorithms, tell data stories, and build statistical models using Python libraries. They will use the same tools that data scientists use to draw meaningful insights and solve organizational problems.

## Learning Environment

This course utilizes a blended classroom approach. The content is fully web-based, with students writing and running code in the browser. Each module of the course is broken down into lessons. Lessons consist of video tutorials, short quizzes, example programs to explore, and written programming exercises.

## Programming Environment

Students write and run Python programs in the browser using the CodeHS editor.

## Projects and Assessments

Students complete a project within each module while also learning and applying new material. The majority of lessons include a formative short multiple-choice quiz to check for understanding. At the end of each module, students take a summative multiple choice quiz that assesses their knowledge of the concepts covered in the module.

## Prerequisites

The Data Science course is designed for intermediate computer science students with at least some knowledge of programming (not language specific) and an interest in computer science. The course is highly visual, dynamic, and interactive, and engaging.

## More Information

Browse the content of this course at <https://codehs.com/course/12135>

## Course Breakdown

### Module 1: The Data Science Life Cycle (4 weeks/20 hours)

Students will learn and apply the process of the data science life cycle. This includes asking statistical questions, collecting or obtaining reliable raw data, analyzing the data using measures of central tendency and spread and interpreting, and summarizing the results.

Objectives / Topics Covered	<ul style="list-style-type: none"><li>• What is Data Science?</li><li>• Gathering Data</li></ul>
-----------------------------	--

	<ul style="list-style-type: none"> <li>○ Quantitative/Qualitative</li> <li>● Exploring Data Using Python</li> <li>● Modules and Libraries</li> <li>● Using the Pandas Library <ul style="list-style-type: none"> <li>○ Series <ul style="list-style-type: none"> <li>■ Measures of Central Tendency</li> <li>■ Measures of Spread</li> </ul> </li> <li>○ DataFrames <ul style="list-style-type: none"> <li>■ Selecting Columns</li> <li>■ Using Functions</li> </ul> </li> </ul> </li> </ul>
Example Assignments / Projects	<ul style="list-style-type: none"> <li>● <b>Mini-Project:</b> Students will go through the first two steps of the data cycle using data of their choosing. <ul style="list-style-type: none"> <li>○ Ask Questions: Formulate a statistical question that can be answered with data.</li> <li>○ Consider Data: Collect or find data that will aid in answering your question.</li> <li>○ Analyze Data: Perform statistical analysis, run calculations and/or create data displays to identify patterns and relationships</li> <li>○ Interpret Data: Answer questions and summarize the results.</li> </ul> </li> <li>● <b>Hot Dog Plots:</b> Use the correct Python functions to create a boxplot of the data. Using the graph, determine the summary statistics and the spread.</li> <li>● <b>Roller Coaster Rankings:</b> Define a function that will compute a score for each roller coaster. Use this function to store the results in a new column.</li> <li>● <b>Student Test Scores:</b> Create a function that finds the maximum test score between test one and test two for each student. Create a function that finds the maximum test score between all three tests for each student. Decide which calculations, along with these two new columns, can help you answer the original statistical question? Explore and further analyze your data until you come to a conclusion.</li> </ul>

## Module 2: Data Science for Change (3 weeks/15 hours)

Students will use and analyze data to better understand a problem, measure the scope of a problem, or understand how people are affected by the problem. They will learn more about cleaning a dataset and filtering by column, rows, and conditions.

Objectives / Topics Covered	<ul style="list-style-type: none"> <li>● What is Big Data? <ul style="list-style-type: none"> <li>○ Cognitive Bias</li> </ul> </li> <li>● Importing and Filtering <ul style="list-style-type: none"> <li>○ loc</li> <li>○ iloc</li> <li>○ By a condition</li> </ul> </li> <li>● Data Cleaning <ul style="list-style-type: none"> <li>○ Dropping Data</li> <li>○ Fixing Data Types</li> </ul> </li> <li>● Exploring with Data Visualizations</li> </ul>
Example Assignments	<ul style="list-style-type: none"> <li>● <b>Project - Data Science for Change:</b> Students will run through the data</li> </ul>

/ Projects	<p>science life cycle with the intent to use data to better understand a problem, to measure the scope of a problem, and to understand how people are affected by the problem.</p> <ul style="list-style-type: none"> <li>● <b>Instagram Filters:</b> This dataset consists of popular Instagram accounts and their number of followers (in millions). Use conditional filtering to print the rows where followers are greater than 230 (million), print only the account columns of those from the United States, print the account and followers columns of the row with the maximum number of followers.</li> <li>● <b>Book Conditions:</b> This data was acquired from the Google Books store. It includes the title of each book, the author(s), the rating (from 1-5), the total voters, the price, the publisher, the page count, and the date the book was published. Print the title, rating, and voters columns for books that have a rating of 4.0 or higher and over 9000 voters.</li> <li>● <b>Cleaning Book Data:</b> Check the data types of the dataset. Do they look okay? Do they need to be changed at all? Permanently drop the publisher and published_date columns. Print the shape of the data and check for duplicate rows. How many are there? Permanently drop duplicate rows from the dataset. Determine the number of missing values in the dataset. What would be the best decision for dealing with the missing values? Make the call and change the data.</li> </ul>
------------	---

### Module 3: Data Storytelling (4 weeks/20 hours)

Students will use and analyze data to tell a data story. They will create a visually appealing infographic that displays important data visualizations. The infographic will also tell a story based on their interpretation after exploring, analyzing, and visualizing the data.

Objectives / Topics Covered	<ul style="list-style-type: none"> <li>● Types of Data Stories</li> <li>● Data Visualizations <ul style="list-style-type: none"> <li>○ Univariate and Bivariate Data</li> </ul> </li> <li>● Normal Distribution</li> <li>● Trends and Correlations</li> <li>● Linear Regression</li> </ul>
Example Assignments / Projects	<ul style="list-style-type: none"> <li>● <b>Project - Data Storytelling:</b> Students will get to tell their own data stories. They will create a visually appealing infographic that displays important data visualizations. The infographic will also tell a story based on their interpretation after exploring and analyzing data.</li> <li>● <b>State Education:</b> The pie chart displayed has an error, which makes it a misleading visual. Find the bug in the program and fix it so that the program displays an accurate pie chart.</li> <li>● <b>Precipitation:</b> There are two datasets in this activity that each list the average temperature and precipitation information for different states. Compare the precipitation averages for both states by plotting a line chart. What conclusions can you make from the chart?</li> </ul>

	<ul style="list-style-type: none"> <li>● <b>Professor's Salaries:</b> This dataset lists professor salaries collected in a survey. Other fields collected were the type of degree, years since their Ph.D. was earned, years of service, and sex (gender). Create a filtered table that lists only female professors and one that lists only professors with over 10 years of service. Do any of the filtered tables follow a normal distribution?</li> <li>● <b>Swim Time Regression:</b> This dataset lists the gold medalist time for the women's 400-meter freestyle swimming finals. Plot the scatterplot. Create a model using the <code>polyfit()</code> function and use the model to plot the line of best fit. Use the model equation to make predictions based on different values.</li> </ul>
--	---

#### Module 4: Data Science for Business (5 weeks/25 hours)

Students will gather business data that can be used to make decisions about how to better the company or product. They will present their findings in a business report that suggests several action items that they predict will help the business's performance and growth.

<p>Objectives / Topics Covered</p>	<ul style="list-style-type: none"> <li>● Determining Dataset Quality</li> <li>● Aggregating Data <ul style="list-style-type: none"> <li>○ Grouping</li> <li>○ Sorting</li> </ul> </li> <li>● Combining Datasets <ul style="list-style-type: none"> <li>○ Concatenating</li> <li>○ Joining/Merging</li> </ul> </li> <li>● Bias in Data Analytics</li> </ul>
<p>Example Assignments / Projects</p>	<ul style="list-style-type: none"> <li>● <b>Project - Data Science for Business:</b> Students will gather (or create) business data that can be used to make decisions about how to better the company or product. They will present their findings in a business report that suggests several action items that they predict will help the business's performance and growth.</li> <li>● <b>Determining Completeness:</b> The raw datasets (without any cleaning) are provided. Your task is to check the completeness of at least two different datasets to compare and contrast them.</li> <li>● <b>Fuzzy Book Titles:</b> Use the FuzzyWuzzy library to extract all titles that are close to matching the word "batman". Print out the results. Are they unique book titles or do you notice errors in the titles?</li> <li>● <b>Billionaire Sort:</b> Group the dataset by gender. Only display the gender count. What insights does this grouping give you? Group the dataset by industry. Only display the industry count. What insights does this grouping give you? Sort by age and then net worth. You will want the oldest person to be listed first. If there are two people of the same age, you'll want the person with the highest net worth to be listed first.</li> <li>● <b>Cereal Production:</b> Let's determine who would get the "Most Improved" award for a ten-year span. <ul style="list-style-type: none"> <li>○ Create a function that will return the difference between two years' values.</li> </ul> </li> </ul>

	<ul style="list-style-type: none"> <li>○ Create a new column listing the growth from the year 1990 to the year 2000.</li> <li>○ Sort by your new column. You will want the highest growth values to be listed first. Only print the Country Name and the new growth columns. Only print the first five rows.</li> </ul> <ul style="list-style-type: none"> <li>● <b>Concatenating Cats:</b> The cat shelter uses two different databases - one for male cats and one for female cats. They'd like to combine the two datasets. Explore and concatenate the two datasets.</li> <li>● <b>School Nurse:</b> Listed in this activity are the two datasets that resulted from the nurse's two visits. The roster has changed from his first visit (first.csv) to his second visit (second.csv). The current roster of students can be found in the second file. How can we keep all of the students in the second roster but add in missing information that might be found in the first roster?</li> </ul>
--	--

**Module 5: What's Next? (1 week/5 hours)**

Students will explore the next chapter in learning about data science and the careers that are available and growing.

Objectives / Topics Covered	<ul style="list-style-type: none"> <li>● Data Science Pathways</li> <li>● Artificial Intelligence</li> <li>● Python Programming</li> </ul>
-----------------------------	--