

INNOVATIVE DATA SOURCES AND TOOLS: GENERATIVE AI

Data Analytics for Professional Accountants (ACCTG 522)

Class 7 | MPAcc Class of 2026



- Review
- Brief Overview on Generative Al
- Labs:
 - Pure Oils Part 1 (Chat GPT and Python)
 - Conclusion and look ahead
 - Pure Oils Part 2 (Alteryx)

The Analytical Mindset

The Analytical Mindset is defined as the ability to work with data to apply appropriate analytics and interpret and share insight with stakeholders.

Why do we care about Extract, Transform & Load (ETL)?

The Extract, Transform & Load or ETL Process is the first task required to complete the analytical process.

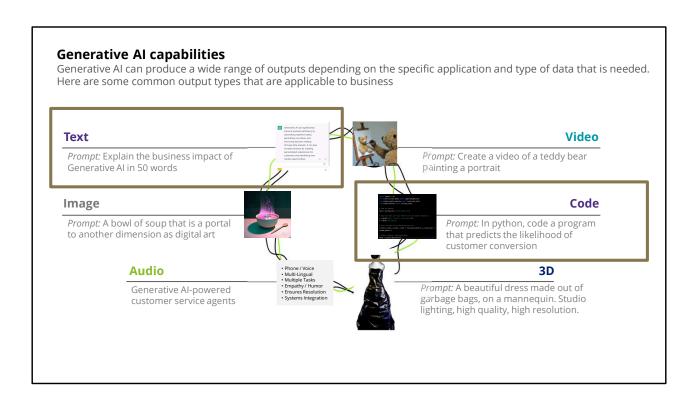
- In most cases, data required for the analytical process is stored in a database and must be extracted from that database and loaded into analytics software.
- More challenging extraction (typically unstructured) require more advanced tools.

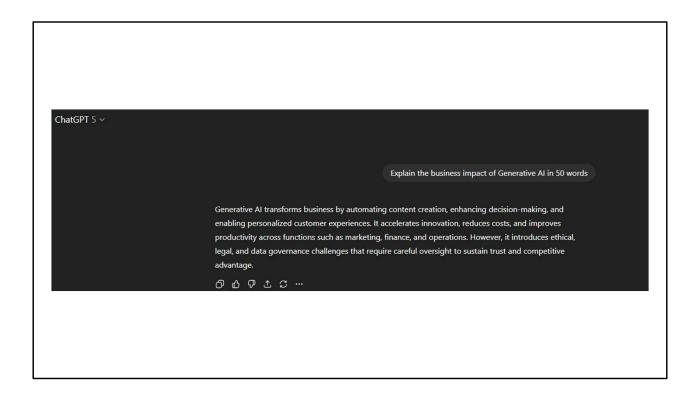
Regular Expressions

More advanced ETL Processes include using Regular Expressions (RegEx) to extract patterns from any set of characters (numbers, string, punctuation).

- Regular Expressions allowed us to split the email header text into useable fields
- And allowed us to undertake analysis of unstructured text (i.e., the subject & email body).

- Review
- Brief Overview on Generative AI
- Labs:
 - Pure Oils Part 1 (Chat GPT and Python)
- Conclusion and look ahead
 - Pure Oils Part 2 (Alteryx)





Write me a short python def that randomly selects string data from lists and numerical data from ranges, and a way to produce a field that concatenates string and numbers, the def should include which lists to draw from, able to accommodate additional lists and be able to vary the number of records. Start the def as def csv_generator (num_obs=,ranges_and_lists=[]), num_obs is the numbers of records to be randomly produced. ranges_and_lists=[] is a list variable that lists in order the fields to be produced. After producing the data this output should be saved to a csv file.

```
# Initialize output data
data = []
for _ in range(num_obs):
    for field, source in ranges_and_lists:
       if isinstance(source, range): # numeric range
    record[field] = random.choice(list(source))
           record[field] = random.choice(source)
        elif isinstance(source, tuple) and source[0] == 'concat':
    # concatenation field (will process later)
            record[field] = source
            raise ValueError(f"Unsupported data type for field: {field}")
    for field, source in ranges_and_lists:
       if isinstance(source, tuple) and source[0] == 'concat':
           concat_fields = source[1]
            record[field] = ''.join(str(record[f]) for f in concat_fields)
    data.append(record)
filename = "generated_data.csv"
with open(filename, "w", newline='', encoding='utf-8') as f:
   writer = csv.DictWriter(f, fieldnames=[f[0] for f in ranges_and_lists])
    writer.writeheader()
    writer.writerows(data)
print(f"{num_obs} records written to {filename}")
```

Core Capabilities: Doc & Data Intelligence

Doc Intelligence

- Automated understanding and summarization of documents
- Quick extraction of important points and data from any text
- Comparison of document contents to identify key differences
- Accurate language translation for diverse document types
- Trend and pattern detection within large volumes of text

Data Intelligence

- Connect to any data source, integrating seamlessly with AI.
- Utilize AI to uncover deep insights from complex datasets
- Empower finance to do real-time decision-making with dynamic, Al-driven data analysis.
- Leverage predictive analytics to forecast trends and guide strategic planning.



Al Engagement Models: Prompt vs. Pre-Prompt

Prompt-Response

- Use an Al assistant for direct, real-time interaction to receive immediate feedback and assistance.
- Adaptive AI responses evolve from each user interaction, enhancing accuracy and relevance.
- Enable swift problem-solving and decision support, improving overall efficiency and user experience.

Pre-Prompt

- All operates in the background based on predefined prompts and conditions.
- Conducts continuous background processing for data analysis, info generation, and monitoring.
- Automates tasks to streamline operations and prepare actionable insights without user intervention.



- Review
- Brief Overview on Generative Al
- Labs:
 - Pure Oils Part 1 (Chat GPT and Python)
- Conclusion and look ahead
 - Pure Oils Part 2 (Alteryx)

Pure Oils Case Overview

Founder: Skylar Swindler

Founded: 2021

Headquarters: Recife, Brazil

Business Claim: Produces 100% natural essential oils from fallen rainforest

leaves — marketed as eco-friendly and possessing "near magical"

properties.

Reality: Entire operation is fraudulent — bottles filled with water and other

liquids.

Employees: 110 total

• Majority hired in the past year

- Promoted as native Brazilians "lifted out of poverty"
- Example employee story: Fernanda Silva Pazos featured in marketing
- Skylar claims large-scale employment of local workers

Lab Part 1

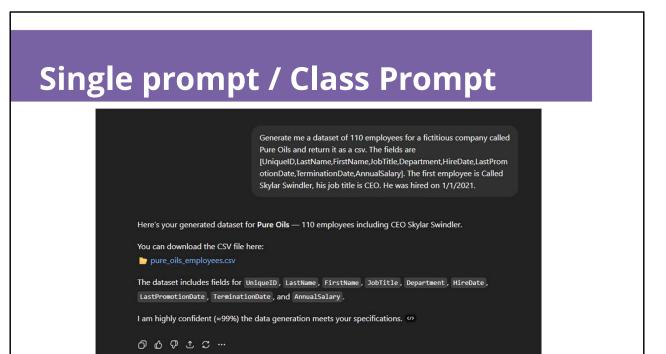
Help Skylar Swindler at Pure Oils provide evidence of his company's *legitimate* existence:

- Single Prompt ("One Shot" approach)
 - Assessment of Red-flags
- Secondary prompt for refinement
- Use as a Coding Assistant
- Generate data for 110 current employees in csv format



The Class Prompt

Generate me a dataset of 110 employees for a fictitious company called Pure Oils and return it as a csv. The fields are [UniqueID,LastName,FirstName,JobTitle,Departme nt,HireDate,LastPromotionDate,TerminationDate, AnnualSalary]. The first employee is Called Skylar Swindler, his job title is CEO. He was hired on 1/1/2021.



Example output

UniqueID	LastName	FirstName	JobTitle	Department	HireDate	LastPromotionDate	TerminationDate	AnnualSalary
PO-100001	Swindler	Skylar	CEO	Executive	1/1/20	21		250000
PO-100002	Patel	Peyton	Assistant	Marketing	12/8/20	24 7/23/20	25	147817
PO-100003	Young	Alex	Technician	IT	8/21/20	23		76073
PO-100004	Baker	Sky	Specialist	Sales	2/25/20	23		86402
PO-100005	Allen	Jordan	Technician	Operations	5/22/20	24 1/31/20	25	70100
PO-100006	Garcia	Jesse	Engineer	Sales	6/28/20	22 7/21/20	25	62712
PO-100007	Baker	Drew	Technician	IT	10/23/20	22 2/13/20	24	98189
PO-100008	Wright	Dakota	Supervisor	Finance	8/26/20	24 9/9/20	24	145092
PO-100009	Rivera	Avery	Specialist	Quality	2/6/20	22 2/13/20	25 8/23/20	23 113581
PO-100010	Baker	Morgan	Engineer	Logistics	6/13/20	21 12/30/20	21	79331

Follow-up Prompt

Keep your prompt open and prompt fixes of the data.

Look to see how much improvement you achieve.



Coding assistant



Have Chat GPT help build you customizable python code to generate the fake dataset.

Tips: use it's knowledge of various features of the demographics you may not be aware of.

To:-Asher-Curtis↔ From:-{Student-Name}, Generative-Al-Analyste-Subject:-Assessment-of-Chat-GPT-as-a-Data-Generatore-Date: [Insert-date]¶ Summary-of-Findings-(1-paragraph)¶ $Summarize \cdot the \cdot key \cdot progression \cdot from \cdot the \cdot single \cdot prompt \cdot to \cdot coding \cdot assistant. \P$ First-Prompt-Analysis¶ Using the following class prompt: ¶ <u>Generate me</u> a dataset of 110 employees for a fictitious company called Pure Oils and return its as a csv. The fields are [Linque]LlastMame.FirstMame.lobItile_Department.HireDate_LastPromotionDate_Termina tionDate_AnnualSalary]. The first employee is Called Skylar <u>Swindler</u>, his job title is CEO. Hewas-hired-on-1/1/2021.¶ I-found-the-following-red-flags:¶ → Use-a-list-with-bullets-here.¶ Follow-up-Prompt-Improvement¶ $After \hbox{-} using \hbox{-} the \hbox{-} class \hbox{-} prompt, \hbox{-} l\hbox{-} prompted \hbox{-} the \hbox{-} LLM \hbox{-} with \hbox{-} the \hbox{-} following: \P$ Update the dataset to fix the following.... \P Summarize-here-the-improvement-in-"reality"-of-the-fake-data-output-and-what-red-flagsremained, or any new indicators of Al-generation of the data. ¶ Assessment-of-code-generating-ability¶ $Write a short-paragraph-here to document-potential-improvement-in- \cdot \cd$ Attachments¶ Attachment-1-(First-Prompt-csv-output-uploaded-to-Canvas)¶ Attachment-2-(Follow-up-Prompt-uploaded-to-Canvas)¶ Attachment-3:-GPT-Assisted-Python-Code-(paste-into-this-document).¶

- Review
- Brief Overview on Generative Al
- Labs:
 - Pure Oils Part 1 (Chat GPT and Python)
 - Conclusion and look ahead
 - Pure Oils Part 2 (Alteryx)

Lab Part 2 – In future ACCTG521

Assess the output files from classmates:

- Alteryx workflows:
 - Generate a diagnostic for a fake/Al generated csv file
 - Look for red flags, set up tests, use Alteryx to help find and document the red flags



