# 10 Academy: Artificial Intelligence Mastery

## Week 3 Challenge Document
### Date: 25 Dec - 31 Tue 2024

# Table of Contents

# Business Objective

Your employer **AlphaCare Insurance Solutions (ACIS)** is committed to developing cutting-edge risk and predictive analytics in the area of car insurance planning and marketing in South Africa. You have recently joined the data analytics team as marketing analytics engineer, and your first project is to analyse historical insurance claim data. The objective of your analyses is to help optimise the marketing strategy as well as discover "low-risk" targets for which the premium could be reduced, hence an opportunity to attract new clients.

In order to deliver the business objectives, you would need to brush up your knowledge and perform analysis in the following areas:

- **Insurance Terminologies**
  - Read on how insurance works. Check out the key insurance glossary [50 Common Insurance Terms and What They Mean — Cornerstone Insurance Brokers](#)
- **A/B Hypothesis Testing**
  - Read on the benefits of A/B hypothesis testing
  - Accept or reject the following null hypothesis
    - There are no risk differences across provinces
    - There are no risk differences between zipcodes
    - There are no significant margin (profit) difference between zip codes
    - There are not significant risk difference between Women and men
- **Machine Learning & Statistical Modeling**
  - For each zipcode, fit a linear regression model that predicts the total claims
  - Develop a machine learning model that predicts optimal premium values given
    - Sets of features about the car to be insured
    - Sets of features about the owner
    - Sets of features about the location of the owner
    - Other features you find relevant

- Report on the explaining power of the important features that influence your model

Your final report should detail the methodologies used, present the findings from your analysis, and make recommendations on plan features that could be modified or enhanced based on the test results. This will help AlphaCare Insurance Solutions to tailor their insurance products more effectively to meet consumer needs and preferences.

# Motivation

The challenge will sharpen your skills in Data Engineering (DE), Predictive Analytics (PA), and Machine Learning Engineering (MLE).

The tasks are designed to improve your ability to manage complex datasets, adapt to challenges, and think creatively, skills that are essential in insurance analytics. This analysis will help you understand more about how hypothesis testing and predictive analytics can be applied in insurance analysis.

Engage with as many tasks as possible. The volume and complexity of the tasks are designed to simulate the pressures and deadlines typical in the financial analytics field.

# Data

The historical data is from Feb 2014 to Aug 2015, and it can be found [here](#)

The structure of the data is as follows
- Columns about the insurance policy
    `UnderwrittenCoverID`
    `PolicyID`
- The transaction date
    `TransactionMonth`
- Columns about the client
    `IsVATRegistered`
    `Citizenship`
    `LegalType`
    `Title`
    `Language`

```
Bank
AccountType
MaritalStatus
Gender
```
- Columns about the client location
```
Country
Province
PostalCode
MainCrestaZone
SubCrestaZone
```
- Columns about the car insured
```
ItemType
Mmcode
VehicleType
RegistrationYear
Make
Model
Cylinders
Cubiccapacity
Kilowatts
Bodytype
NumberOfDoors
VehicleIntroDate
CustomValueEstimate
AlarmImmobiliser
TrackingDevice
CapitalOutstanding
NewVehicle
WrittenOff
Rebuilt
Converted
CrossBorder
NumberOfVehiclesInFleet
```
- Columns about the plan
```
SumInsured
TermFrequency
CalculatedPremiumPerTerm
ExcessSelected
CoverCategory
CoverType
CoverGroup
Section
```

```
Product
StatutoryClass
StatutoryRiskType
```
- Columns about the payment & claim
```
TotalPremium
TotalClaims
```

## Learning Outcomes

- Understanding the data provided and extracting insight. You will have to explore different techniques, algorithms, statistical distributions, sampling, and visualization techniques to gain insight.
- Understand the data structure and algorithms used in EDA and machine learning pipelines.
- Modular and object-oriented Python code writing. Python package building.
- Statistical Modeling and Analysis. You will have to use statistical models to predict and analyze the outcomes of A/B tests, applying techniques such as logistic regression, or chi-squared tests, as appropriate to the hypotheses being tested.
- A/B Testing Design and Implementation. You will design robust A/B tests that can yield clear, actionable results. This includes determining the sample size, selecting control and test groups, and defining success metrics.
- Data Versioning. You will manage and document versions of datasets and analysis results.

# Competency Mapping

The tasks you will carry out in this week's challenge will contribute differently to the 11 competencies 10 Academy identified as essential for job preparedness in the field of Data Engineering, and Machine Learning engineering. The mapping below shows the change (lift) one can obtain by delivering the highest performance in these tasks.

| Competency | Potential contributions from this week |
|---|---|
| Professionalism for a global-level job | Articulating business values |
| Collaboration and Communicating | Reporting to stakeholders |
| Software Development Frameworks | Using Github for CI/CD, writing modular codes, and packaging |

| | |
|---|---|
| Python Programming | Advanced use of Python modules such as Pandas, Matplotlib, Numpy, Scikit-learn, Prophet, and other relevant Python packages |
| Data & Analytics Engineering | data filtering, data transformation, and data warehouse management |
| MLOps & AutoML | Pipeline design, data, and model versioning, |
| Deep Learning and Machine Learning | Statistical modelling, Model Interprtablity |
| Data Versioning | DVC |

## Team

Facilitator:

- Mahlet
- Kerod
- Rediet
- Elias
- Rehmet
- Emitinan

# Key Dates

- **Challenge Introduction** - 9:30 AM UTC time on Wednesday  25 Dec 2024.
- **Interim Submission -** 8:00 PM UTC time on Friday  27 Dec 2024.
- **Final Submission** - 8:00 PM UTC time on Tuesday  31 Dec 2024.

# Deliverables and Tasks to be done

## Task 1:

Git and GitHub

- Tasks:
    - Create a git repository for the week with a good Readme
    - Git version control
    - CI/CD with Github Actions
- Key Performance Indicators (KPIs):
    - Dev Environment Setup.
    - Relevant skill in the area demonstrated.

Project Planning - EDA & Stats

- Tasks:
    - Data Understanding
    - Exploratory Data Analysis (EDA)
    - Statistical thinking
- KPIs:
    - Proactivity to self-learn - sharing references.
    - EDA techniques to understand data and discover insights,
    - Demonstrating Stats understanding by using suitable statistical distributions and plots to provide evidence for actionable insights gained from EDA.

**Minimum Essential To Do**
- Create a github repository that you will be using to host all the code for this week.
- Create at least one new branch called "task-1" for your analysis of day 1
- Commit your work at least three times a day with a descriptive commit message
- Perform Exploratory Data Analysis (EDA) analysis on the following:
    - **Data Summarization**:
        - Descriptive Statistics: Calculate the variability for numerical features such as **TotalPremium**, **TotalClaim**, etc.
        - Data Structure: Review the **dtype** of each column to confirm if categorical variables, dates, etc. are properly formatted.
    - **Data Quality Assessment**:
        - Check for missing values.
    - **Univariate Analysis**:
        - Distribution of Variables: Plot histograms for numerical columns and bar charts for categorical columns to understand distributions..
    - **Bivariate or Multivariate Analysis:**

- Correlations and Associations: Explore relationships between the monthly changes **TotalPremium** and **TotalClaims** as a function of **ZipCode**, using scatter plots and correlation matrices.
  - **Data Comparison**
    - **Trends Over Geography:** Compare the change in insurance cover type, premium, auto make, etc.
  - **Outlier Detection:**
    - Use box plots to detect outliers in numerical data
  - **Visualization**
    - Produce 3 creative and beautiful plots that capture the key insight you gained from your EDA

# Task 2:

## Data Version Control ([DVC](#))

- Tasks:
  - Install DVC
    - `pip install dvc`
  - Initialize DVC: In your project directory, initialize DVC
    - `dvc init`
  - Set Up Local Remote Storage
    - Create a Storage Directory
      - `mkdir /path/to/your/local/storage`
    - Add the Storage as a DVC Remote
      - `dvc remote add -d localstorage /path/to/your/local/storage`
  - Add Your Data:
    - Place your datasets into your project directory and use DVC to track them
      - `dvc add <data.csv>`
  - Commit Changes to Version Control
    - Create different versions of the data.
      - 
    - Commit the .dvc files (which include information about your data files and their versions) to your Git repository

- ○ Push Data to Local Remote
  - ■ `dvc push`

**Minimum Essential To Do:**
- Merge the necessary branches from task-1 into the main branch using a Pull Request (PR)
- Create at least one new branch called "task-2"
- Commit your work with a descriptive commit message.
- Install DVC
- Configure local remote storage
- Add your data
- Commit Changes to Version Control
- Push Data to Local Remote

# Task 3:

A/B Hypothesis Testing

- Accept or reject the following **Null Hypotheses:**
  1. There are no risk differences across provinces
  2. There are no risk differences between zip codes
  3. There are no significant margin (profit) difference between zip codes
  4. There are not significant risk difference between Women and Men

- Tasks:
  - ○ Select Metrics
    - ■ Choose the key performance indicator (KPI) that will measure the impact of the features being tested.
  - ○ Data Segmentation
    - ■ **Group A (Control Group)**: Plans without the feature
    - ■ **Group B (Test Group)**: Plans with the feature.

- For features with more than two classes, you may need to select two categories to split the data as Group A and Group B. You must ensure, however, that the two groups you selected do not have significant statistical differences on anything other than the feature you are testing. For example, the client attributes, the auto property, and insurance plan type are statistically equivalent.
  - Statistical Testing
    - Conduct appropriate tests such as chi-squared for categorical data or t-tests or z-test for numerical data to evaluate the impact of these features.
    - Analyze the p-value from the statistical test:
      - If $p\_value < 0.05$ (typical threshold for significance), reject the null hypothesis. This suggests that the feature tested does have a statistically significant effect on the KPI.
      - If $p\_value >= 0.05$, fail to reject the null hypothesis, suggesting that the feature does not have a significant impact on the KPI.
  - Analyze and Report
    - Analyze the statistical outcomes to determine if there's evidence to reject the null hypotheses. Document all findings and interpret the results within the context of their impact on business strategy and customer experience.

**Minimum Essential To Do:**
- Merge the necessary branches from task-2 into the main branch using a Pull Request (PR)
- Create at least one new branch called "task-3"
- Commit your work with a descriptive commit message.
- Select Metrics
- Data Segmentation
- Statistical Testing
- Analyze and Report

# Task 4:

## Statistical Modeling

- Tasks:
  - Data Preparation:
    - Handling Missing Data: Impute or remove missing values based on their nature and the quantity missing.
    - Feature Engineering: Create new features that might be relevant to TotalPremium and TotalClaims.
    - Encoding Categorical Data: Convert categorical data into a numeric format using one-hot encoding or label encoding to make it suitable for modeling.
    - Train-Test Split: Divide the data into a training set (for building the model) and a test set (for validating the model), typically using a 70:30 or 80:20 ratio.
  - Modeling Techniques
    - **Linear Regression**
    - Decision Trees
    - Random Forests
    - **Gradient Boosting Machines (GBMs):**
      - **XGBoost**
  - Model Building
    - Implement Linear Regression, Random Forests, and XGBoost models
  - Model Evaluation
    - Evaluate each model using appropriate metrics like accuracy, precision, recall, and F1-score.
  - Feature Importance Analysis
    - Analyze which features are most influential in predicting retention.
  - Use SHAP (SHapley Additive exPlanations) or LIME (Local Interpretable Model-agnostic Explanations) to interpret the model's predictions and understand how individual features influence the outcomes.
  - Report comparison between each model performance.

**Minimum Essential To Do:**

- Merge the necessary branches from task-3 into the main branch using a Pull Request (PR)
- Create at least one new branch called "task-4".
- Commit your work with a descriptive commit message.
- Data preparation
- Model building
- Model evaluation
- Model Interpretability

# Due Date (Submission)

**Interim Submission Friday (27 Dec 2024): 8:00 PM (UTC)**

- GitHub Link to your main branch
- Interim report - Covering task-1 and task-2

**Final Submission Tuesday (31 Dec, 2024): 8:00 PM (UTC)**

- GitHub Link to your main branch
- Final report - covering all week-3 work
  - Write it in the format that you could post it as a Blog on Medium.

## Feedback

You may not receive detailed comments on your interim submission but will receive a grade.

## Other Considerations:

- **Documentation:** Encourage detailed documentation in code and report writing.
- **Collaboration:** Emphasise collaboration through Github issues and projects.
- **Communication**: Regular check-ins, Q&A sessions, and a supportive community atmosphere.
- **Flexibility:** Acknowledge potential challenges and encourage proactive communication.
- **Professionalism:** Emphasise work ethics and professional behavior.
- **Time Management:** Stress the importance of punctuality and managing time effectively.

# Tutorials Schedule

In the following, the color **purple** indicates morning sessions, and non-purple indicates afternoon sessions.

- Day 1 Wednesday (25 Dec 2024 ):       `
    - Introduction to the Challenge (Mahlet)
    - Introduction to Insurance Analytics (Rediet)
- Day 2 Thursday  (26 Dec 2024 :
    - Statistical distributions, hypothesis testing, and creating actionable insights (Elias)
    - Data Version Control (DVC) (Kerod)
- Day 3 Friday (27 Dec 2024) :
    - Statistical Modeling and Evaluation  (Emitinan)

- Day 4 Monday  (30 Dec 2024) :
    - Introduction to model interpretability (Elias)
    - Q&A (Kerod and Mahlet)

# References

- **Insurance Analytics**
  - https://www.fsrao.ca/media/11501/download
  - https://www.xenonstack.com/blog/data-analytics-in-insurance
  - https://business.wisc.edu/wp-content/uploads/2021/07/ProjDescription_Web.pdf
  - https://www.swissre.com/risk-knowledge/driving-digital-insurance-solutions/connected-car-how-data-analytics-is-shaping-the-future-of-auto-insurance.html
- **A/B Hypothesis Testing**
  - https://www.engagys.com/insights/a-b-testing-the-key-to-effective-healthcare-communications
  - https://www.linkedin.com/pulse/abcs-ab-testing-healthcare-marketing-daniella-koren/
  - https://medium.com/tiket-com/a-b-testing-hypothesis-testing-f9624ea5580e
  - https://www.optimizely.com/insights/blog/why-an-experiment-without-a-hypothesis-is-dead-on-arrival/
- **Data Version Control(DVC)**
  - https://dvc.org/
  - https://dvc.org/doc/user-guide
- **Statistical Modeling:**
  - https://www.heavy.ai/technical-glossary/statistical-modeling
  - https://www.coursera.org/articles/statistical-modeling
  - https://www.statlect.com/glossary/statistical-model
  - https://builtin.com/data-science/random-forest-algorithm
  - https://www.analyticsvidhya.com/blog/2021/08/conceptual-understanding-of-logistic-regression-for-data-science-beginners/#:~:text=Logistic%20Regression%20is%20another%20statistical,pass%20this%20exam%20or%20not.
  - https://www.analyticsvidhya.com/blog/2018/09/an-end-to-end-guide-to-understand-the-math-behind-xgboost/
- **Version control – Git**
  - What is version control | Atlassian
  - Learn Git branching -- interactive way to learn Git

- ○ [Git with large files](#)
  - ○ [Which files to not track and how to not track them? | Atlassian](#)
  - ○ [.gitignore docs](#)
  - ○ [Conventional commits -- lightweight convention on top of commit messages.](#)
- **CI/CD**
  - ○ [What is Continuous Integration | Atlassian](#)
  - ○ [DevOps Pipeline | Atlassian](#)
  - ○ [7 Popular Open Source CI/CD Tools - DevOps.com](#)
  - ○ [Setting up a CI/CD pipeline on Github](#)