

## ▼ Default title text

```
# @title Default title text
%%writefile /content/interim_report.md
# Interim Report: Insurance Risk Analytics for 10 Academic Years

**Submission Deadline**: June 15, 2025, 8:00 PM UTC
**Repository**: https://github.com/bekonad/insurance-analytics
**Prepared by**: Bereket Feleke
```

This report summarizes progress on **Task 1** (Git and GitHub).

### ## Project Overview

- **Objective**: Analyze historical car insurance data (2015-2024).
- **Data**: `insurance\_data.txt`, a pipe-delimited dataset.
- **Tasks**:
  - Task 1.1: Set up Git and GitHub repository.
  - Task 1.2: Conduct EDA to explore risk and profitability.
  - Task 2: Implement Data Version Control (DVC) for reproducibility.
- **Environment**: Mobile device using Google Colab and JupyterLab.
- **Tools**: Python, pandas, matplotlib, seaborn, DVC, Git, GitHub.

### ## Task 1: Git and GitHub + Exploratory Data Analysis (EDA)

#### ### Task 1.1: Git and GitHub

- **Status**: Completed
- **Actions**:
  - Created `insurance-analytics` repository with README.
  - Configured CI/CD with `.github/workflows/lint.yml` for code quality.
  - Merged `task-1` branch into `main` via pull request.
- **Outcome**: Functional repository with version control.

#### ### Task 1.2: EDA

- **Status**: Completed
- **Notebook**: `EDA\_Insurance\_Analytics.ipynb`
- **Challenges**:
  - **Output Suppression**: Only "Loss Ratio by Gender" displayed.
  - **Solution**: Split Loss Ratio code into separate cells.
  - **TypeError**: Temporal trends plot failed due to `Pandas`.
  - **Solution**: Converted `Month` to `datetime` with `pd.to\_datetime`.
  - **Data Loading**: Initial `KeyError: 'TotalClaims'`.
  - **Column Naming**: Renamed columns (`make` to `Make`, `year` to `Year`).
  - **Deprecation Warning**: Seaborn's `palette` without `hue` parameter.
- **EDA Findings**:
  - **Loss Ratio by Province**:
    - Gauteng: 0.429, Mpumalanga: 0.393, Limpopo: 0.349, Western Cape: 0.312.
    - **Insight**: Gauteng's high loss ratio indicates urban areas.
  - **Loss Ratio by VehicleType**:
    - Heavy Commercial: 0.794, Light Commercial: 0.544, Medium Commercial: 0.456.
    - **Insight**: Heavy Commercial vehicles have the highest loss ratio.
  - **Loss Ratio by Gender**:
    - Female: 0.492, Male: 0.349, Not specified: 0.348.



```
# @title Default title text
%%writefile /content/interim_report.md
# Interim Report: Insurance Risk Analytics for 10 Academic Years

**Submission Deadline**: June 15, 2025, 8:00 PM UTC
**Repository**: https://github.com/bekonad/insurance-analytics
**Prepared by**: Bereket Feleke
```

This report summarizes progress on **Task 1** (Git and GitHub).

### ## Project Overview

- **Objective**: Analyze historical car insurance data (2015-2024).
- **Data**: `insurance\_data.txt`, a pipe-delimited dataset.
- **Tasks**:
  - Task 1.1: Set up Git and GitHub repository.
  - Task 1.2: Conduct EDA to explore risk and profitability.
  - Task 2: Implement Data Version Control (DVC) for reproducibility.
- **Environment**: Mobile device using Google Colab and JupyterLab.
- **Tools**: Python, pandas, matplotlib, seaborn, DVC, Git, GitHub.

### ## Task 1: Git and GitHub + Exploratory Data Analysis (EDA)

#### ### Task 1.1: Git and GitHub

- **Status**: Completed
- **Actions**:
  - Created `insurance-analytics` repository with README.
  - Configured CI/CD with `.github/workflows/lint.yml` for code quality.
  - Merged `task-1` branch into `main` via pull request.
- **Outcome**: Functional repository with version control.

#### ### Task 1.2: EDA

- **Status**: Completed
- **Notebook**: `EDA\_Insurance\_Analytics.ipynb`
- **Challenges**:
  - **Output Suppression**: Only "Loss Ratio by Gender" displayed.
  - **Solution**: Split Loss Ratio code into separate cells.
  - **TypeError**: Temporal trends plot failed due to `Pandas`.
  - **Solution**: Converted `Month` to `datetime` with `pd.to\_datetime`.
  - **Data Loading**: Initial `KeyError: 'TotalClaims'`.
  - **Column Naming**: Renamed columns (`make` to `Make`, `year` to `Year`).
  - **Deprecation Warning**: Seaborn's `palette` without `hue` parameter.
- **EDA Findings**:
  - **Loss Ratio by Province**:
    - Gauteng: 0.429, Mpumalanga: 0.393, Limpopo: 0.349, Western Cape: 0.312.
    - **Insight**: Gauteng's high loss ratio indicates urban areas.
  - **Loss Ratio by VehicleType**:
    - Heavy Commercial: 0.794, Light Commercial: 0.544, Medium Commercial: 0.456.
    - **Insight**: Heavy Commercial vehicles have the highest loss ratio.
  - **Loss Ratio by Gender**:
    - Female: 0.492, Male: 0.349, Not specified: 0.348.

- **Insight**: Higher female loss ratio warrants further investigation.
- **Distributions**:
  - `TotalPremium`: Mean 61.91 ZAR, median 2.18 ZAR, max 1000 ZAR.
  - `TotalClaims`: Mean 64.86 ZAR, median 0 ZAR, right-skewed.
  - **Insight**: Most policies have low premiums and zero claims.
- **Temporal Trends**:
  - Average claims and claim frequency plotted, showing seasonal patterns.
  - **Insight**: Likely seasonal patterns (e.g., Q4 spike in claims).
- **Claims by Make**:
  - Analysis intended but not shown in output.
  - **Insight**: Expected high claims for prevalent makes like Toyota.
- **Outcome**: Comprehensive EDA addressing all guiding questions.

## ## Task 2: Data Version Control (DVC)

- **Status**: Completed
- **Notebook**: `DVC_Setup.ipynb`
- **Actions**:
  - Initialized DVC with `--no-scm` to resolve SCM error.
  - Configured local storage at `/content/dvc_storage`.
  - Tracked `insurance_data.txt` with `dvc add` and push.
  - Committed `insurance_data.txt.dvc` and `.dvc/config`.
  - Verified with `dvc pull`.
- **Challenges**:
  - SCM error resolved with `--no-scm`.
  - Mobile constraints addressed using local storage and DVC.
- **Outcome**: Dataset tracked, meeting data versioning requirements.

## ## Next Steps

- Merge `task-2` pull request.
- Start Task 3 (A/B Testing) using EDA insights.
- Submit repository URL, PDF of `EDA_Insurance_Analytics` report.

## ## Citations

- [DVC Getting Started](<https://dvc.org/doc/start>)
- [Google Colab](<https://colab.research.google.com/notebooks/welcome-to-colab.ipynb>)
- [GitHub Docs](<https://docs.github.com/en/repositories/managing-your-repositorys-settings-and-features>)

🔄 Overwriting /content/interim\_report.md

- **Loss Ratio by Gender**:
  - Female: 0.492, Male: 0.349
  - **Insight**: Higher female loss ratio warrants further investigation.
- **Distributions**:
  - `TotalPremium`: Mean 61.91 ZAR, median 2.18 ZAR, max 1000 ZAR.
  - `TotalClaims`: Mean 64.86 ZAR, median 0 ZAR, right-skewed.
  - **Insight**: Most policies have low premiums and zero claims.
- **Temporal Trends**:
  - Average claims and claim frequency plotted, showing seasonal patterns.
  - **Insight**: Likely seasonal patterns (e.g., Q4 spike in claims).
- **Claims by Make**:
  - Analysis intended but not shown in output.
  - **Insight**: Expected high claims for prevalent makes like Toyota.
- **Outcome**: Comprehensive EDA addressing all guiding questions.

## ## Task 2: Data Version Control (DVC)

- **Status**: Completed
- **Notebook**: `DVC_Setup.ipynb`
- **Actions**:
  - Initialized DVC with `--no-scm` to resolve SCM error.
  - Configured local storage at `/content/dvc_storage`.
  - Tracked `insurance_data.txt` with `dvc add` and push.
  - Committed `insurance_data.txt.dvc` and `.dvc/config`.
  - Verified with `dvc pull`.
- **Challenges**:
  - SCM error resolved with `--no-scm`.
  - Mobile constraints addressed using local storage and DVC.
- **Outcome**: Dataset tracked, meeting data versioning requirements.

## ## Next Steps

- Merge `task-2` pull request.
- Start Task 3 (A/B Testing) using EDA insights.
- Submit repository URL, PDF of `EDA_Insurance_Analytics` report.

## ## Citations

- [DVC Getting Started](<https://dvc.org/doc/start>)
- [Google Colab](<https://colab.research.google.com/notebooks/welcome-to-colab.ipynb>)
- [GitHub Docs](<https://docs.github.com/en/repositories/managing-your-repositorys-settings-and-features>)

Overwriting /content/interim\_report.md