

EXPERIMENT REPORT

Student Name	Nathan Collins
Project Name	MLAA_Assignment_1
Date	
Deliverables	<Part B Report> <Part B Notebook>

1. EXPERIMENT BACKGROUND

Provide information about the problem/project such as the scope, the overall objective, expectations. Lay down the goal of this experiment and what are the insights, answers you want to gain or level of performance you are expecting to reach.

1.a. Business Objective

Explain clearly what is the goal of this project for the business.

The goal of the project is to develop a predictive model that can **accurately** predict cancer mortality in different US counties based on a 33-feature dataset.

How will the results be used?

The insights provided can be used to improve cancer prevention, diagnosis, and treatment strategies, ultimately reducing the cancer mortality rate in the US. These results may come in the form of a myriad of benefits, such as aiding local healthcare employees and providers to retain sufficient inventory to assist cancer patient health outcomes (financial, age, gender), assisting leaders and policymakers in providing adequate funding to US counties of higher cancer incidence (geography, death rate), and assisting decisions to investigate higher-risk areas and fund investigations to develop interventions (government and private health cover). Results may also assist citizens in comprehending their family members' survivability following diagnosis. Constructing a successful model ultimately reduces the prevalence of cancer or increases patient survival by identifying linked features.

What will be the impact of accurate or incorrect results?

The influence of accurate results could affect how cancer mortalities are supplied for, intervened against, and prevented. As we are dealing with lives and collateral that comes with loss, the stakes are high should incorrect results be generated and utilised. An example may include a model that identifies a county with a

high cancer mortality rate, healthcare providers can allocate increased cancer screening and prevention efforts in that county. Likewise, if the generated results prove unreliable, it may lead to a misallocation of crucial resources, resulting in ineffective interventions and more fatalities than necessary. Additionally, should a successful model be constructed, mispredictions on survivability may leave lasting emotional and financial collateral to the diagnosed and their peripheral family.

Present the hypothesis you want to test, the question you want to answer or the insight you are seeking.

Part B hypothesis:

Instantiating a Multilinear Regression, will produce a lower Mean Squared Error than the Univariate Linear Regression.

Explain the reasons why you think it is worthwhile considering it.

1.b.
Hypothesis

This hypothesis warrants exploration, as it suggests there is causality between all demographic and socioeconomic factors provided, and cancer mortality rates in US counties. Should these relationships be confirmed, it aids social awareness and acknowledgement, leading to ample provision of public health policies and interventions aimed at reducing cancer mortality rates in areas with these correlating features. This has the overall affect at saving human life, a priceless faculty to harness. Outcomes from the following may furthermore serve as a benchmark for current health standards, facilitating an understanding whether conditions improve over future generations.

Detail what will be the expected outcome of the experiment.

The expected outcome for **Part B**, is the exploration of “deathrate” against all features within the data frame. It is expected to yield as lower MSE than the univariate model.

List the possible scenarios resulting from this experiment.

1.c.
Experiment
Objective

Multivariate Linear Regression: a higher MSE value than the Univariate Linear Regression, meaning the model is less accurate. Or a lower MSE value, meaning the model is more accurate.

A significant and highly probable outcome is a lack of noticeable correlation between variables presumably correlated. Should this be the case, resolving the scenario would be to explore other features, or perform feature engineering to discern relationships, if present.

Code miscalculation from human error during the preparation phase is a possible throughout construction. This has the capacity

to lead to true positive, or false negative assertions, or presumed statistically significant outcomes and create devastating, real life consequences if deployed. For this reason, a careful eye needs to be had, with persistent, meticulous revision to ensure variables and values aren't misinterpreted.

Another possible scenario, and the desirable one, is the successful creation and deployment of an accurate model. As human behaviour and is chaotic, variables will never truly be 100% correlated, thus success would be deemed minimal error with predictions. The outcome could have shape lives, as mentioned previously.

2. EXPERIMENT DETAILS

*Elaborate on the approach taken for this experiment.
List the different steps/techniques used and explain the rationale for choosing them.*

2.a.
Data
Preparation

Describe the steps taken for preparing the data (if any).

<The data was prepared in Part A>

The final preparation was ensuring the table consisted only of numeric values, which meant disregarding the geography column.

Explain the rationale why you had to perform these steps.

Data preparation is essential, as all factors available, are provided to the Multivariate model, instead of just two features, as seen in Part A. Unnecessary categorical columns were dropped, as they impeded utility in modelling. Missing number columns were dropped, or if marginal, were retained and granted the mean (if normally distributed), purely for additional, rough features to incorporate in Part C, should previous models not meet expectations. Outliers were left in, as they represented real scenarios, minus ages >300. Correlations were constructed prior to training a Univariate Linear Model, to provide the reveal stronger relationships and conserve time.

List also the steps you decided to not execute and the reasoning behind it.

- Removing some noticeable outliers, as they represented real-world scenarios, minus ages exceeding 300.
- Removing the inserted averages for "PctPrivateCoverageAlone" and "PctEmployed16_Over".

While not an accurate representation of the dataset, these features were too useful to drop.

- Inserting more features from other dataframes; as it would have different contextual circumstances with its collection and may shift the results to one that doesn't reflect real scenarios.
- Removing geography; as it was convertible with one-hot-encoding during Part C, for feature engineering. It remained a non-numerical feature.

Highlight any step that may potentially be important for future experiments

Performing the `train_test_split` function, as the data was concatenated during preparation for uniform cleaning. It is important to split the data prior to creating the model.

Describe the steps taken for generating features (if any).

While "geography" was reduced to its corresponding States for one-hot-encoding in the data preparation phase, all features were engineered in Part C.

Coverage on these topics are in the corresponding report.

Explain the rationale why you had to perform these steps.

<No features generated for Part B>

List also the feature you decided to remove and the reasoning behind it.

"binnedInc" was removed, as it was non-numerical and encapsulated another feature already available in the data, "Median Income".

Highlight any feature that may potentially be important for future experiments.

"PctPublicCoverageAlone" – it represents the highest correlation with the target variable and relates to other features – such as "income" and "poverty".

2.c.
Modelling

Describe the model(s) trained for this experiment and why you choose them.

Multivariate Linear Regression – a complex, yet thorough tool for visualising relationships between the target variable and all variables, by minimising squared residuals

List the hyperparameter tuned and the values tested and also the rationale why you choose them.

< Simple Multivariate Regression does not possess hyperparameters to tune.>

List also the models you decided to not train and the reasoning behind it.

< N/A for Part B >

Highlight any model or hyperparameter that may potentially be important for future experiments.

< Random Forest Regression: In a random forest, a set of decision trees is constructed using bootstrapped samples of the training data, and each tree is trained on a random subset of the available features. During the training phase, each tree is built using a random selection of training data and features. This randomness helps to reduce overfitting by preventing any single tree from being too dependent on a particular subset of the data. >

3. EXPERIMENT RESULTS

Analyse in detail the results achieved from this experiment from a technical and business perspective. Not only report performance metrics results but also any interpretation on model features, incorrect results, risks identified.

3.a.
Technical
Performance

Score of the relevant performance metric(s).

All numeric variables vs TARGET_deathRate

Training MSE: 742.3

Testing MSE: 421

Provide analysis on the main underperforming cases/observations and potential root causes.

- Outliers may impede overall accuracy.
- Imputation of data points to represent averages for Public and Private Coverage can skew accuracy.
- Data standardisation, as different values are on different scales.
- Not all data is simply correlated with each other, as it was
- Data was sourced at different times, in different quantities, with more areas surveyed than others.
- Human behaviour is difficult to measure accurately, and each individual possesses faculties that may shape their circumstance. This results in a wide spread in the scatter plot, leading to higher MSE scores. More data and trailing different features may help reduce these MSE scores.

Interpret the results of the experiments related to the business objective set earlier.

As the objective of the project is to develop a predictive model that can **accurately** predict cancer mortality in different US counties based on a 33-feature dataset, given the high, inaccurate MSE variance, it does not seem possible with a univariate linear regression model.

Estimate the impacts of the incorrect results for the business, (some results may have more impact compared to others).

The inability to deliver an accurate model means cancer mortality remains uncategorisable and predictable, further removing the life-conserving benefits mentioned above, as they cannot be delivered. From the narrative of a multilinear model, this paves way into investigation into further model exploration.

List all the issues you faced during the experiments (solved and unsolved). Present solutions or workarounds for overcoming them.

- Inaccurate model – provide more features (i.e. CO2 output, diet, sun exposure, known carcinogens within areas), alongside more people, with more consistent sampling measures.
- Standardising data – doesn't affect the outcome
- Coding reliability – human error, resolved with practice

- NaN values – the few missing, fit with the mean; does not represent a real survey.
- Geography – data is arbitrarily spread, and in different quantities
- Data limitations – more features, would help provide more opportunities to discern correlations and relationships.
- Outliers – while noticeable outliers were present, they had to be retained to mimic real-world scenarios; the age values >300 were dropped as it was not realistic.

Highlight also the issues that may have to be dealt with in future experiments.

- High MSE value, remedied through investigating other models.
- Discerning the impact of Public Coverage and Private Coverage – as missing values were provided the mean of the distribution.
- Geography – limited to some US states, with varying degrees and quantities

4. FUTURE EXPERIMENT

Reflect on the experiment and highlight the key information/insights you gained from it that are valuable for the overall project objectives from a technical and business perspective.

4.a.
Key Learning

Reflect on the outcome of the experiment and list the new insights you gained from it. Provide rationale for pursuing more experimentation with the current approach or call out if you think it is a dead end.

While unsuccessful, the exercise provides avenues to explore feature engineering. While trends in capital, education and decisions lead better or worse cancer mortality rates, it can never be simplified to a simple model. The insights however allow further investigation into further regression models, that may integrate more features. The MSE, while high, can be reduced by conceiving a larger snapshot of the environment of US states, through more data features.

As cancerous growth is the product of genetic and environmental factors, both dynamics are attainable with a diversity of features. The current features are sparse and limit a perspective of a broad, consistently changing dynamic. For this

reason, I do not consider that more experimentation with the current features will yield a model that could reduce lifestyle factors to predictability.

Given the results achieved and the overall objective of the project, list the potential next steps and experiments.

Feature engineering – through geography, higher tertiary education, non-Caucasian descent, Male-female mean age gap, and Percentage uninsured.

Further regression models, through Ridge, Lasso, KNN, and Random Forest.

For each of them assess the expected uplift or gains and rank them accordingly.

4.b.
Suggestions /
Recommendations

Feature Engineering – more features to add to existing data, which may reveal further insights and correlative features

Further regression models, e.g. lasso and random forest – alternatives, that operate different paths at predicting datasets.

If the experiment achieved the required outcome for the business, recommend the steps to deploy this solution into production.

<The experiment did not achieve the required outcome, therefore it is advised the model is not deployed, rather used as a means to learn from, so that future models may not encounter similar issues with crafting a “reductionist narrative” of lifestyle factors.>