

Project Everest Report - Jakub Bartnik

Author: Jakub Bartnik

GitHub URL

https://github.com/jakubbartnik/UCDPA_jakubbartnik

Abstract

This review contains the summary about impact of introducing new safety gear for climbing versus expedition survival that climbed in area of Himalaya Nepal within last 100 years. Implementation process is splitted between four parts: Setting goals, data preparation, data exploration, analysis. The review will explain how all that was done by myself and what conclusion I made from results. Report content can be found in python notebook including the code so it is easier to follow steps. There is another python notebook in repository demonstrating simple webscrapping. File name is Web_scrapping_quotes.ipynb.

Introduction

I chose this topic for my Data Science project as I have a strong personal interest in climbing. With this being my first major project in the field, I felt it was fitting to select a theme that I am well-acquainted with. The objective of this research is to examine the impact of safety equipment, such as steel carabiners and dynamic ropes and harnesses, on the survival rates of climbers in the Nepal Himalaya region over the years. For my analysis, I will be utilising a Mount Everest dataset that encompasses all expeditions from 1905 to 2019. To ensure a more structured and readable review, I like to set smaller goals that stem from my main objective.

Goal 1: is to investigate the impact of the introduction of the climbing carabiner in 1910 on the survival rate of expeditions. To achieve this, I will split the dataframe into two sets and compare the falls.

Goal 2: is to examine the effect of the implementation of dynamic rope and harness in the 1960s. By dividing the data into periods with and without these innovations and determine their impact on the falls prediction.

Goal 3: is to determine the optimum number of helpers (hired_staff) that an expedition should hire in order to increase their chances of success, and express this result as a percentage.

Dataset

The Himalayan Database is a remarkable resource that documents all expeditions in the Nepal Himalaya, including notable achievements by mountaineers, such as early attempts, first ascents, and major accidents. This invaluable resource spans a timeline from 1905 through Spring 2019, covering over 465 significant peaks in Nepal. The database was compiled from Elizabeth Hawley's expedition archives, supplemented by alpine journals, books, and contributions from Himalayan climbers. The dataset comprises of three CSV files: expeditions.csv (16 columns, 10364 rows), peaks.csv (8 columns, 467 rows), and members.csv (21 columns, 76519 rows). Surprisingly, the

dataset contains about 280,630 empty values. The variables types summary shows that the majority of columns (31) are string data type, followed by 7 boolean type columns, and 7 integer type columns.

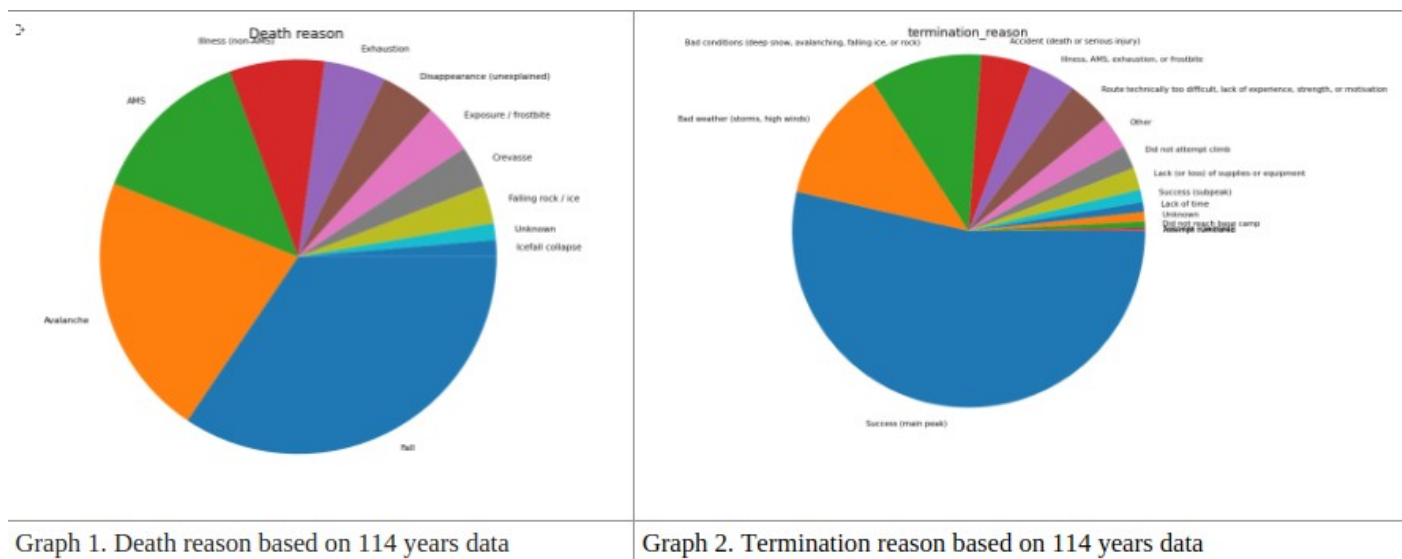
Implementation Process

Data cleaning

I started to look at the data manually prior writing this review. To achieve my purpose, I want to merge all of the csv's files into one. I am exploring all csv one by one using `.head()` which shows me five first rows to identify best column I can merge with another dataframe. Identified column 'peak_id' to merge expeditions.csv with peak.csv and assigned it to 'df1' as temporary dataframe. Next I merge temporary dataframe and members.csv based on 'expedition_id' and assigned it to dataframe - 'df'. I did used `.merge` command to combine all three csv's. Looking for empty values in my current dataframe - 'df'. I am using `.isna()` function to show all 'nan' values. Because I find the column that will make significant impact on my results I want to take closer look at column 'death_cause'. That is why I am using `df['death_cause'].unique()` function. Thanks to that I see that value 'Other' is not really make more sense then 'nan' or 'unknown' but it generate more noise so I am replacing all 'Other' with 'nan'. To limit variables in current dataset that are mainly repeated like year_x and year_y I will drop some of the columns that came after merging. Important columns from expedition.csv is 'termination_reason' and 'expedition_id', 'year'. Important columns from peaks.csv is 'peak_name' and 'climbing_status'. Important columns from members.csv is 'death_cause'. Those I believe are important columns I want to keep. I will drop the following columns: 'trekking_agency', 'peak_alternative_name', 'sex', 'expedition_role', 'injury_type', 'injury_height_metres', 'oxygen_used_x', 'age', 'first_ascent_country', 'hired', 'success', 'highpoint_metres_y', 'oxygen_used_y', 'death_height_metres', 'injured', 'citizenship', 'solo', 'year_y', 'season_y', 'peak_name_y', 'peak_name_x', 'peak_id_y', 'peak_id_x', 'season_x', 'climbing_status', 'first_ascent_expedition_id'. Finding all duplicates that can be dropped. I done that based on field 'expedition_id'. I sorted values and finally dropped it and result saved to 'exp_dup'.
Drop complete duplicates from cleaning
Stripping down string from expedition_id using simple regex.
Cleaned data I saved to new file clean_everest.csv to avoid running all code boxes in python notebook.

Data exploration and preparation for goals

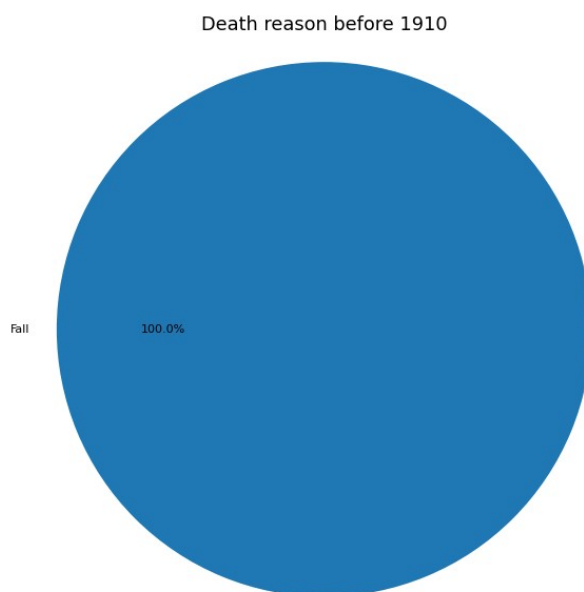
Importing libraries for exploration mainly for visual representation of results. Importing cleaned data as 'dfc'.
Features termination_reason and death_cause are nominal and need to be converted to dummy variables. I am doing it using dataframe option '`.get_dummies`'.
Showing overall procentage of reasons to die in Himalaya based on all 114 years of experience using pie chart and same chart for the same period of time on termination reason for expeditions.



Death reason based on 114 years dataframe. Avalanches, AMS (Alternated Mind status - caused by trauma or intoxication) and Illness are the main reasons for death after falling.

Exploration and preparation for Goal 1.

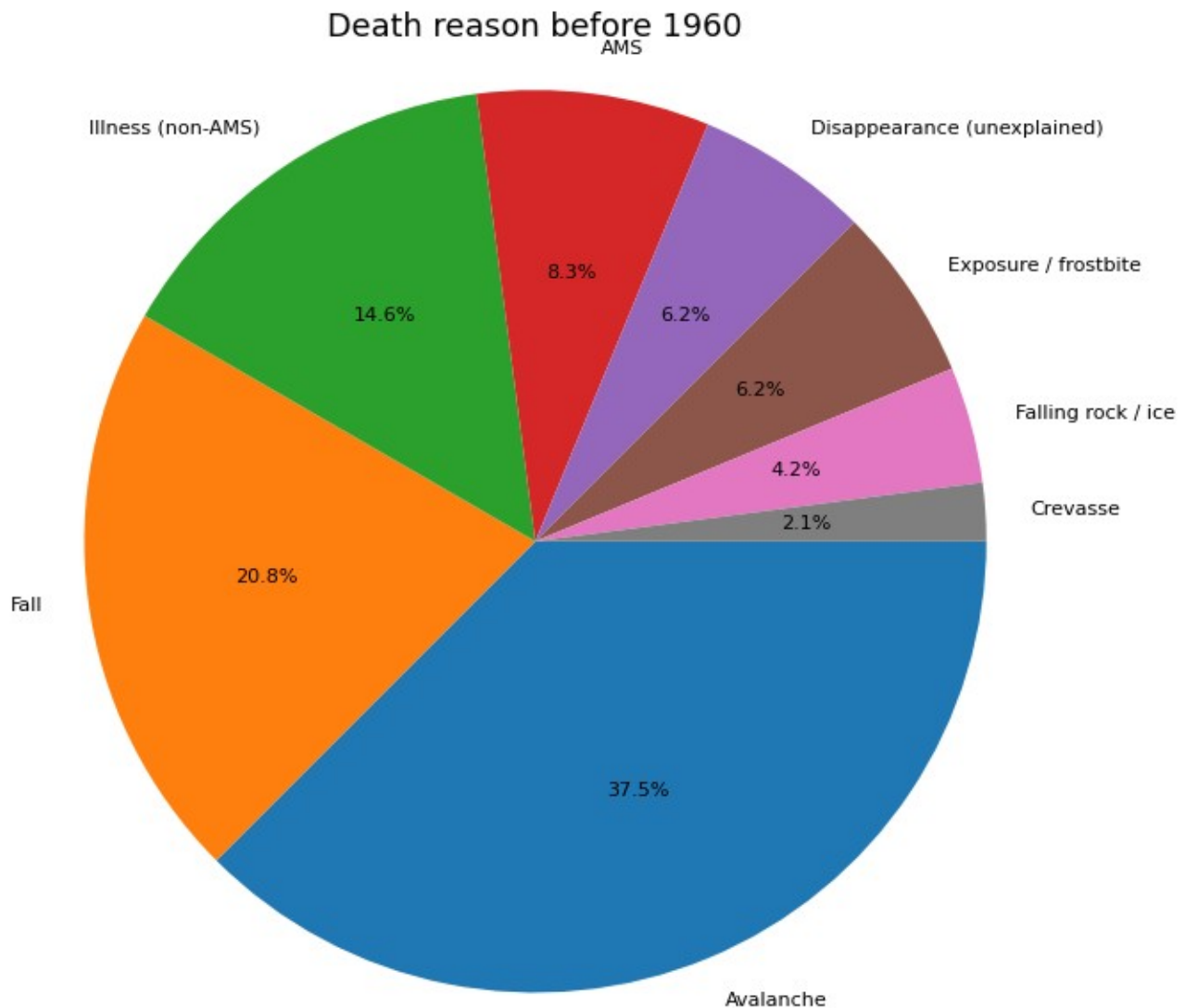
Splitting data for goal 1 which is before 1910 assign it to dataframe 'before1910' and after 1911 assigned to 'after1910'.



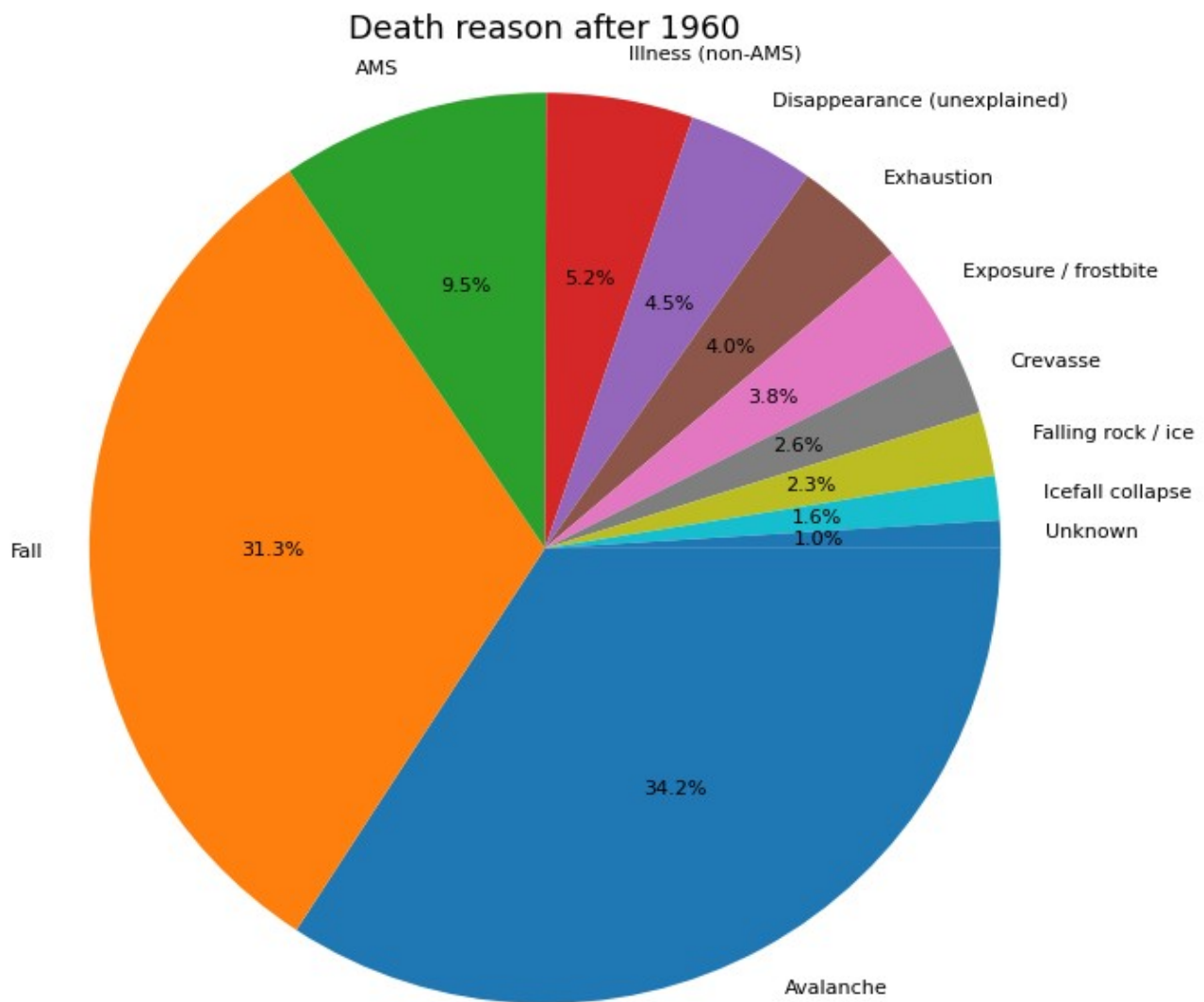
By shape of before1910 I can see that its really small and does contain only one record of falling. Therefore goal 1 will not be continued as there is not enough observation in data. More context provided in Results.

Exploration and preparation for Goal 2.

Examine the effect of the implementation of dynamic rope and harness in the 1960s. By dividing the data into periods with and without these innovations, I can determine their impact on the falls prediction.



Pie chart for data before 1960 indicate that avalanche and illness was the biggest reasons to die and falling going into third place. Plotting pie-chart for death after 1960 to have better comparison.



Both graphs are very similar but its indicating that reasons were changing along the years .

Creating function that will save 1 when the fall occur and 0 if not. Stored result of categorise function in new feature. Want to simplify my dataframe and use columns - 'year_x', 'members', 'member_deaths', 'hired_staff', 'hired_staff_deaths', 'died', 'falls'.

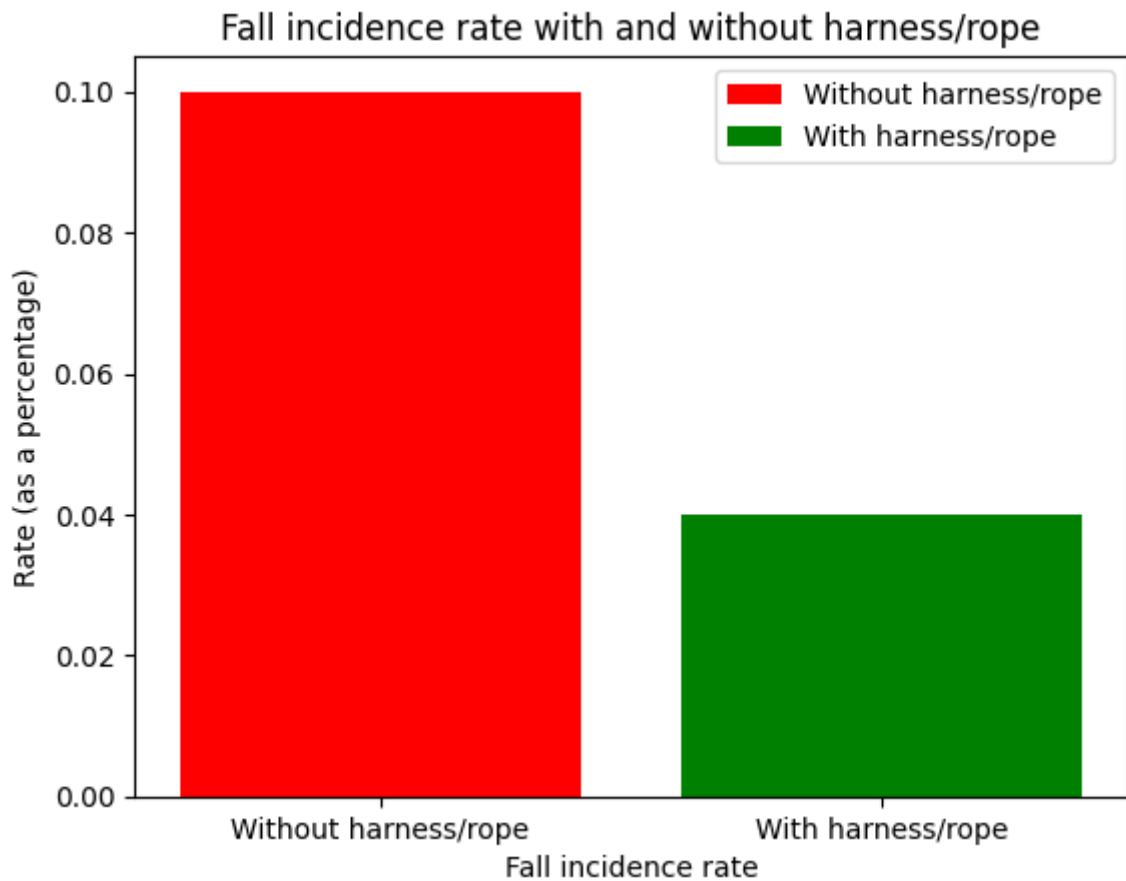
Analysis for Goal 2

Calculating the total number of falls and climbers in each group and then dividing the number of falls by the total number of climbers to find the fall incidence rate. Sliced dataframe for goal 2 for before and after dataframes and added new column 'harnes_used' where 0 mean no and 1 mean yes. This is to avoid 1 dimension data in model. Converted columns to a numeric type as per model requirement. Calculating the fall incidence rate for two different data periods one without harness and rope, and one with harness and rope.

Fall incidence rate without harness and rope: 0.10%

Fall incidence rate with harness and rope: 0.04%

This was visualised using boxplot:



The fall incidence rate without harness and rope is 0.10%, which means that out of all the climbers who did not use a harness and rope, 0.1% of them reported falls. The fall incidence rate with harness and rope is 0.04%, which means that out of all the climbers who used a harness and rope, 0.04% of them reported falls. Overall, using a harness and rope seems to be associated with a lower fall incidence rate compared to not using one.

Performing logistic regression model on 'before' and 'after' dataframes to predict falls and prints out the confusion matrix and classification report for both 'before' and 'after' dataframes. Model resulted in following:

```

Confusion Matrix for Before:
[[274  0]
 [ 2  0]]
Classification Report for Before:
              precision    recall  f1-score   support

     0       0.99      1.00      1.00       274
     1       0.00      0.00      0.00         2

   accuracy      0.99      0.99      0.99       276
  macro avg       0.50      0.50      0.50       276
 weighted avg       0.99      0.99      0.99       276

Confusion Matrix for After:
[[14908  0]
 [ 60  2]]
Classification Report for After:
              precision    recall  f1-score   support

     0       1.00      1.00      1.00    14908
     1       1.00      0.03      0.06         62

   accuracy      1.00      0.99      1.00    14970
  macro avg       0.50      0.52      0.53    14970
 weighted avg       1.00      1.00      0.99    14970

```

Hyperparameter tuning of logistic regression using GridSearchCV. GridSearchCV is searching through all possible combinations of 8 hyperparameter settings and using 5-fold cross-validation for each combination. The output indicates that the best set of hyperparameters is {'C': 1, 'penalty': 'l1', 'solver': 'liblinear'}.

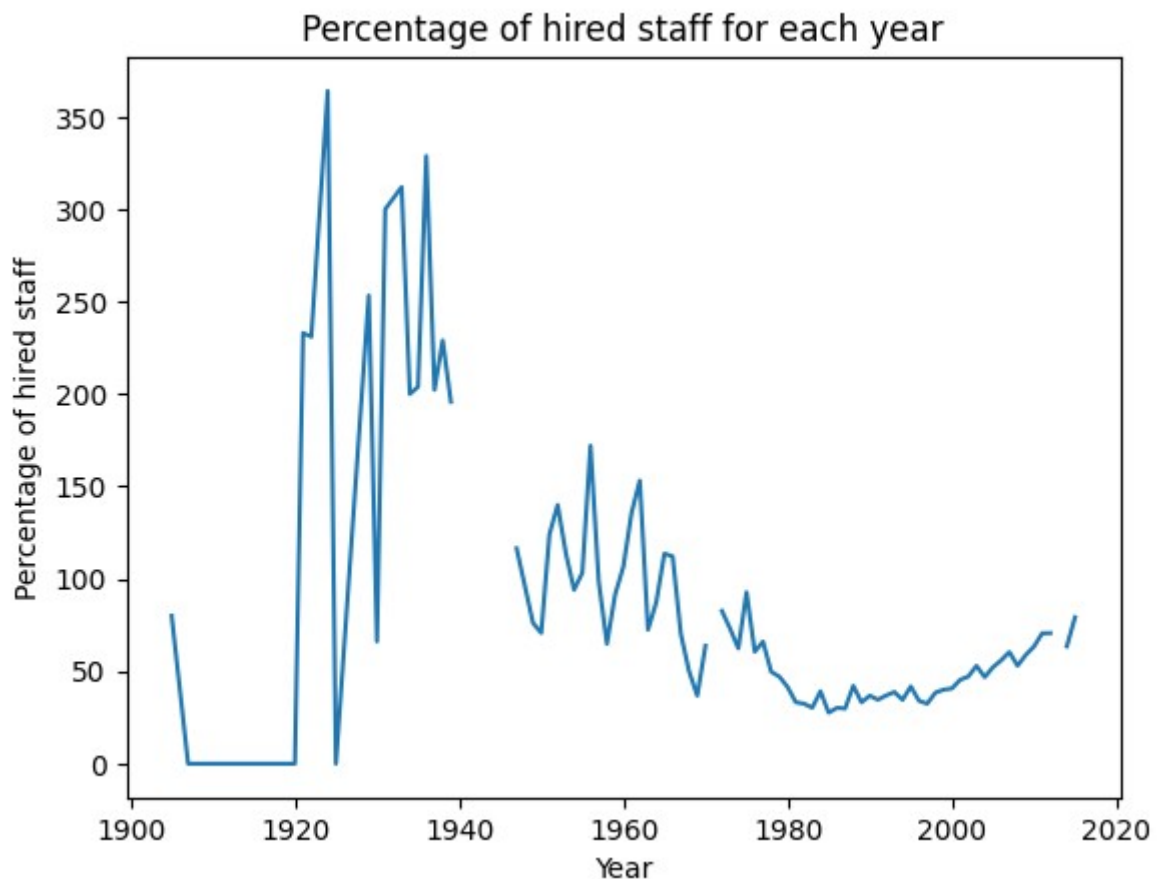
Observation for Goal 2

It seems that falls were more likely to occur in the 'before' dataset than in the 'after' dataset. This can be inferred from the fact that the 'recall' score for the '1' label ('falls') was higher in the 'before' dataset (0.00) than in the 'after' dataset (0.03), indicating that the model was better able to identify true positive 'falls' in the 'after' dataset than in the 'before' dataset.

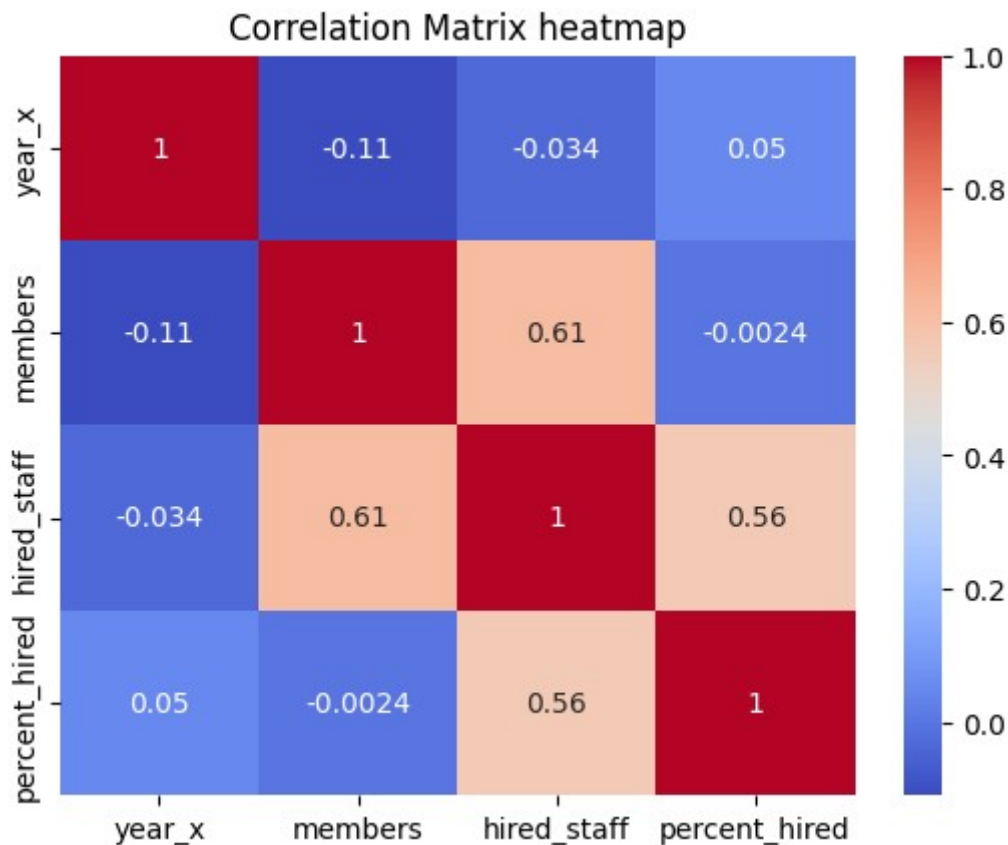
Exploration and preparation for Goal 3.

In interest of this goal is to limit dataset for only successful expeditions.

Using DataFrame.transform() method. I want to add percent of hired staff for amount of members
Plotting the percentage of hired staff for each year.



Data contained some inf and NaN values that were cleaned. Created heatmap of the correlation matrix to see related features.



This indicated relationship on feaures like hired_staff with percent_hired and hired_staff with memebers.

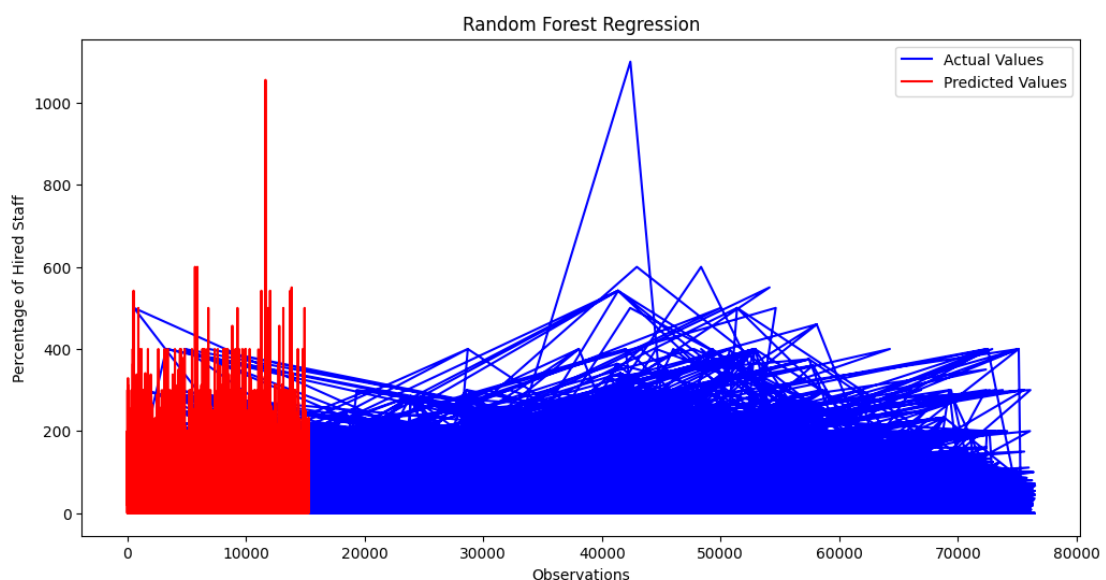
Analysis for Goal 3

Creating a random forest regressor model to predict the 'percent_hired' column based on the 'members' and 'hired_staff' columns. It then fits the model to the train data and makes predictions on the test data. Finally, it evaluates the performance of the model using the mean squared error (MSE) and R-squared score.

Results are:

MSE: 0.13759521565547536
R-squared: 0.9999625008645775

This was showed graphically:



Trying to adjust hyperparameter to improve model performance using a GridSearchCV approach which is a popular method for finding the optimal hyperparameters by testing every combination of values for the hyperparameters. This code-box took about 11 min to run.

The first set of results, $MSE=0.1376$ and $R\text{-square}=0.99996$, are the performance metrics of my random forest regressor model without any hyperparameter tuning. This model used $n_estimators=100$, which is the default value for this parameter in the `RandomForestRegressor()` model. The second set of results, $MSE=0.2514$ and $R\text{-square}=0.99993$, are the performance metrics of the same random forest regressor model after hyperparameter tuning, with the following settings: $max_depth=None$, $min_samples_leaf=1$, $min_samples_split=5$, and $n_estimators=50$. The tuned model uses a smaller $n_estimators$ than the default value, but other parameters are set manually. Overall, we can see that the second model performed worse in terms of both MSE and R-squared compared to the first model. This suggests that the default settings for the `RandomForestRegressor()` model might be more appropriate for my dataset. I am going to skip plotting this results.

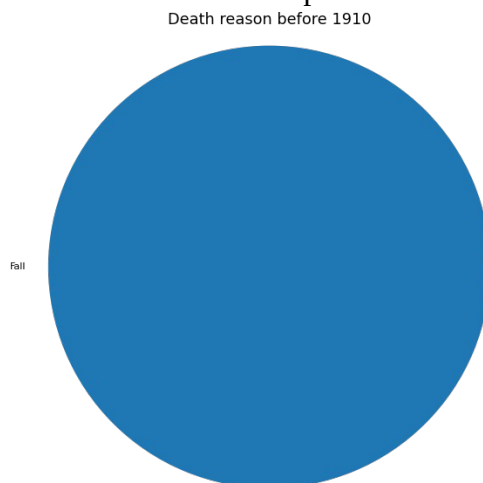
Observations for Goal 3

It looks like the mean squared error (MSE) is 0.1375 and the R-squared value is 0.9999, indicating that the random forest model fits the data very well. This means that the model is able to make accurate predictions of the number of hired staff for future expeditions, based on the historical data that were fitted.

Results

Goal 1.

After exploring data on this paper it looks like goal 1 need to be skipped as there is not enough data to support this comparison. This is very clear when we see how many record are in splited dataframe for before 1910 period and from the pie chart.



To get goal 1 accomplished it would be much better to analyse climbing data from another part of the world where there is more observations. Unfortunately observations made in pre-computerised years would be spread across various personal diaries so it would be very hard to group it.

Goal 2.

It seems that falls were more likely to occur in the 'before' dataset than in the 'after' dataset. This lead me to definite answer that development of harness and dynamic rope has impact on climbers savety.

Goal 3.

From the graph I can read that hiring 20% to 40% of helping staff looks like perfect percentage for expedition to be successful.

Insights

Goal 1.

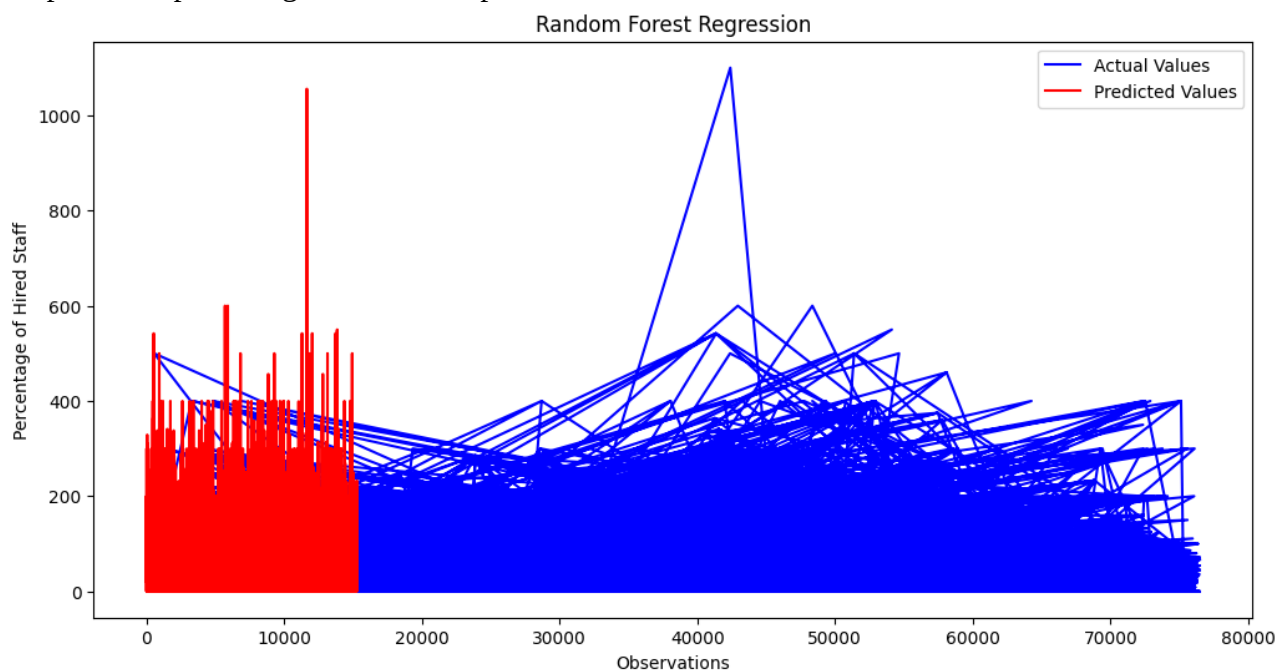
There is not enough observations to provide definite answer if steel carabiner made any difference on climbing safety

Goal 2.

Death due to fall is actually extreme as climbers are falling on route attempt basis. Not all falls are dangerous. In dataset fall is listed only as a reason for death but there is no count how many time climbers has fall before they die. One of the reason to die due to fall is not due to break of the rope but most often due to bad protection placing in the rock. As in Nepal we are dealing with Ice climbing it could be due to loose ice.

Goal 3.

After performing random forest regressor model on both type of data I got really good visualisation for predicted percentage for future expeditions.



From which I can read that hiring 20% to 40% of helping staff looks like perfect percentage for expedition to be successful. Where 20% seems to be minimum.

References

Dataset: <https://www.kaggle.com/datasets/majunbajun/himalayan-climbing-expeditions>

Convert nominal to numerical values: <https://machinelearningmastery.com/one-hot-encoding-for-categorical-data/>

Generating more samples: <https://machinelearningmastery.com/generate-test-datasets-python-scikit-learn/>

Comparing: <https://towardsdatascience.com/how-to-compare-two-or-more-distributions-9b06ee4d30bf>

Scrapping webpage url: <https://quotes.toscrape.com>

Book - <https://ethanweed.github.io/pythonbook>

Supervised learning - <https://campus.datacamp.com/courses/supervised-learning-with-scikit-learn/>