# Predicting Spill Contributing Factors: A Machine Learning Approach Leveraging Decision Trees

## Abstract:

The environmental threats posed by spills of oil and hazardous chemicals necessitate proactive measures to prevent and mitigate their impact. Regulatory bodies play a crucial role by mandating the reporting of these incidents, creating a rich dataset for analysis. This project delves into the application of machine learning to predict the contributing factors of spills reported to a state agency. We leverage a Decision Tree Classifier as a foundational model to analyze historical spill data, aiming to identify underlying patterns and enhance our understanding of the root causes of these occurrences. By leveraging machine learning techniques, we strive to develop a predictive model that can anticipate future spill causes and ultimately contribute to more effective environmental protection strategies.

## Introduction:

The potential devastation caused by spills of hazardous materials necessitates robust environmental protection strategies. Regulatory bodies play a critical role by mandating the reporting of such incidents, enabling the collection of valuable data for response and mitigation efforts. This project capitalizes on this data, obtained from the state's designated data portal (https://incidentnews.noaa.gov/raw/index), to develop a predictive model for spill contributing factors. Our central hypothesis centers on the premise that by analyzing past incidents, we can unearth underlying patterns and subsequently construct a model capable of anticipating future spill causes. This capability would empower environmental agencies to proactively address potential risks and implement targeted preventative measures.

## Methodology:

**The project meticulously followed these key steps:**

**Data Acquisition and Exploration:** The initial phase involved acquiring historical spill data from the designated state data portal. Following acquisition, the data underwent a thorough exploration process to comprehend its structure, identify missing values, and implement essential cleaning procedures. Data cleaning techniques such as identifying and handling outliers, addressing inconsistencies in data formatting, and potentially imputing missing values using appropriate methods were crucial for ensuring the quality and integrity of the data used to train the machine learning model.

**Feature Engineering:** To potentially elevate model performance, feature engineering techniques were meticulously applied. This phase may have involved the creation of novel features derived from existing ones, or the transformation of existing features to enhance interpretability for the model. For instance, features representing the time of day or day of the

week of a spill incident could be created from a timestamp feature. Additionally, categorical features like cause type or location type might be encoded numerically for better utilization by the machine learning model.

**Baseline Model:** To establish a benchmark for evaluating the performance of our machine learning model, a baseline model, specifically a Dummy Classifier, was implemented. This baseline model functioned by predicting the most frequently occurring contributing factor. The performance of the Dummy Classifier would provide a reference point to assess the effectiveness of the Decision Tree model in identifying specific spill causes.

**Decision Tree Model:** A Decision Tree Classifier, renowned for its interpretability and adeptness at handling categorical data, was subsequently trained on the meticulously prepared data. Decision trees function by splitting the data into subsets based on the values of a single feature at each node. This process continues recursively until a stopping criterion is met, resulting in a tree-like structure that can be easily visualized and interpreted. To guarantee optimal model performance, Grid Search, a hyperparameter optimization technique, was employed. Hyperparameters are settings within the machine learning model that can influence its behavior. Grid Search allows us to systematically evaluate different combinations of hyperparameter values and select the configuration that yields the best performance on a validation set.

**Model Evaluation:** The effectiveness of the Decision Tree model was rigorously assessed using a separate test set. This test set was not used during the training process, ensuring an unbiased evaluation of the model's ability to generalize to unseen data. Metrics such as accuracy, precision, recall, and F1-score served as the foundation for evaluating the model's performance. Accuracy measures the overall proportion of correct predictions, while precision focuses on the proportion of positive predictions that are truly positive. Recall, on the other hand, emphasizes the proportion of actual positive cases that are correctly identified by the model. The F1-score provides a harmonic mean of precision and recall, offering a balanced view of the model's performance.

## Results:

The baseline model achieved an accuracy of approximately 54.9%. This translates to the model's ability to predict the most prevalent contributing factor roughly half of the time. However, its capacity for identifying specific spill causes, as measured by precision, recall, and F1-score, was demonstrably inadequate. This highlights the limitations of a simple baseline model and underscores the potential benefits of employing a more sophisticated machine learning approach.

The Decision Tree model exhibited a measurable improvement over the baseline, achieving an accuracy of around 55.7%. This advancement suggests that the model was able to discern underlying patterns within the data and consequently generate slightly more accurate predictions. While the model achieved a modest improvement over the baseline, it highlights the

inherent challenges associated with this task. This initial investigation paves the way for further exploration of alternative machine learning algorithms and feature engineering techniques to bolster the model's predictive capabilities. By persistently developing this model, we can potentially enhance our ability to predict and ultimately prevent hazardous material spills, thereby safeguarding the environment.

## Future Work:

**Exploration of Alternative Machine Learning Algorithms:** Delving into the application of alternative machine learning algorithms such as Random Forest and Gradient Boosting holds promise for enhancing model performance.

**Incorporation of Additional Features:** Investigating the potential benefits of incorporating additional features that may be relevant to spill occurrences could significantly contribute to the model's effectiveness.

**Advanced Feature Engineering Techniques:** The utilization of advanced feature engineering techniques for improved data representation presents a compelling avenue for further research.

By persistently pursuing the development of this model, we can equip ourselves with a more potent tool to predict and ultimately prevent hazardous material spills