# Flood Probability Regression with a Hybrid Ensemble w/ Uncertainty Quantification

By: Vatsal Sivaratri and
Connor Friedman

# Problem Statement and Motivation

**Goal:** Develop a predictive, robust model(s) for determining the likelihood of a flood (percent chance of an area to flood)

**Motivation:** Flood frequency and severity has been rising due to climate change. Reliable systems that can predict these events in advance can help communities plan and mitigate disasters.

# Understanding Uncertainty in Predictions

- Machine Learning models make predictions, but how **confident** are they?
    - A prediction of **80% flood risk** means something entirely different if the model is highly confident vs. uncertain
- **Types of Uncertainty:**
    - **Aleatoric Uncertainty (Data Uncertainty):** Noise in the data; Inherent to the task
    - **Epistemic Uncertainty (Model Uncertainty):** How much the model doesn't know due to a lack of training data.
- **Why does it matter?**
    - **High-confidence predictions ->** Trust the model's output
    - **High-uncertainty predictions ->** Need more investigation, can't take results for face value
- **Proposed Approach:** Compute uncertainty estimates using **Monte Carlo Dropout** & **Bootstrap Variance** in Ensemble Learning

# Dataset Source & Description

- **Source:** [Kaggle Flood Dataset](#)
- **Key Details:** Dataset consists of 1.12m instances and 21 features, all numerical.
- **Key Features:**
  - **Environmental Factors:** Monsoon intensity, watershed count.
  - **Geographic & Infrastructure:** Drainage systems, river management.
  - **Human Impact Factors:** Urbanization, deforestation.
- **Target Variable:** Flood Probability (Continuous values from 0 to 1)
- **Data Preprocessing:**
  - Feature Scaling: Standardized using z-score normalization.
  - Outliers: Removed using IQR

# Ensemble Learning: Deep Neural Network (DNN)

**Architecture:**

- **Input Layer:** 21 features
- **Hidden Layers:** 2 Dense Layers (64->32 neurons)
- **Activation:** ReLU
- **Dropout Layers:** 30% dropout to improve generalization
- **Output Layer:** Single neuron (linear activation) -> **Flood Probability (0-1)**

**Uncertainty Estimation: Monte Carlo Dropout:**

- Run Model **T** times
- Randomly drop neurons each run, giving slightly different outputs

**Mean Prediction:**

$$\hat{y} = \frac{1}{T} \sum_{t=1}^{T} \hat{y}_t$$

**Uncertainty (Variance) Prediction:**

$$\sigma^2 = \frac{1}{T} \sum_{t=1}^{T} (\hat{y}_t - \hat{y})^2$$

If predictions **vary a lot**, uncertainty is **high**. If they are **stable,** uncertainty is **low.**

# Ensemble Learning: Gradient Boosting Machine (GBM)

**How GBM Works:**

- Uses **decision trees** to repeatedly improve weak models.
- Each tree focuses on correcting errors from the last one.
- Outputs a **weighted sum of trees.**

**Uncertainty Estimation: Bootstrapped Variance:**

- Train multiple GBM models on different randomized subsets of the data
- Collect predictions from all the models and compute:

**Mean**

**Variance**

$$\hat{y}_{GBM} = \frac{1}{N} \sum_{i=1}^{N} \hat{y}_{GBM}^{(i)}$$

$$\sigma^2_{GBM} = \frac{1}{N} \sum_{i=1}^{N} (\hat{y}GBM^{(i)} - \hat{y}GBM)^2$$
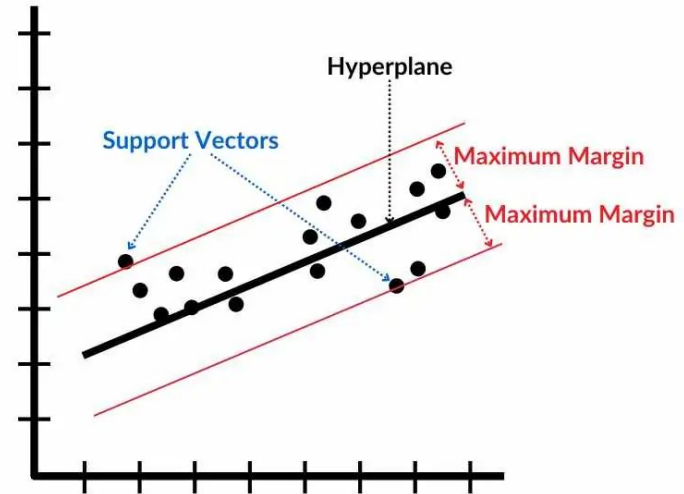
# Ensemble Learning: Support Vector Regression (SVR)

**How SVRs work:**

- Find the best hyperplane that fits the flood probability
- Uses kernel trick to capture non-linearity

**Uncertainty Estimation: Bootstrap Sampling**

- Train multiple SVR models on different resampled datasets
- Use variance across models as uncertainty estimate



Support Vector Regression (SVR)

# Model Aggregation

**Each model provides a prediction + uncertainty estimate.**

**Uncertainty-Weighted Averaging:**

More confident models **(lower variance)** contribute more:

$$\mathbf{w}_i = min(\frac{1}{max(\sigma_i, \epsilon)}, \delta), \epsilon = 1e - 15, \delta = 1e15$$
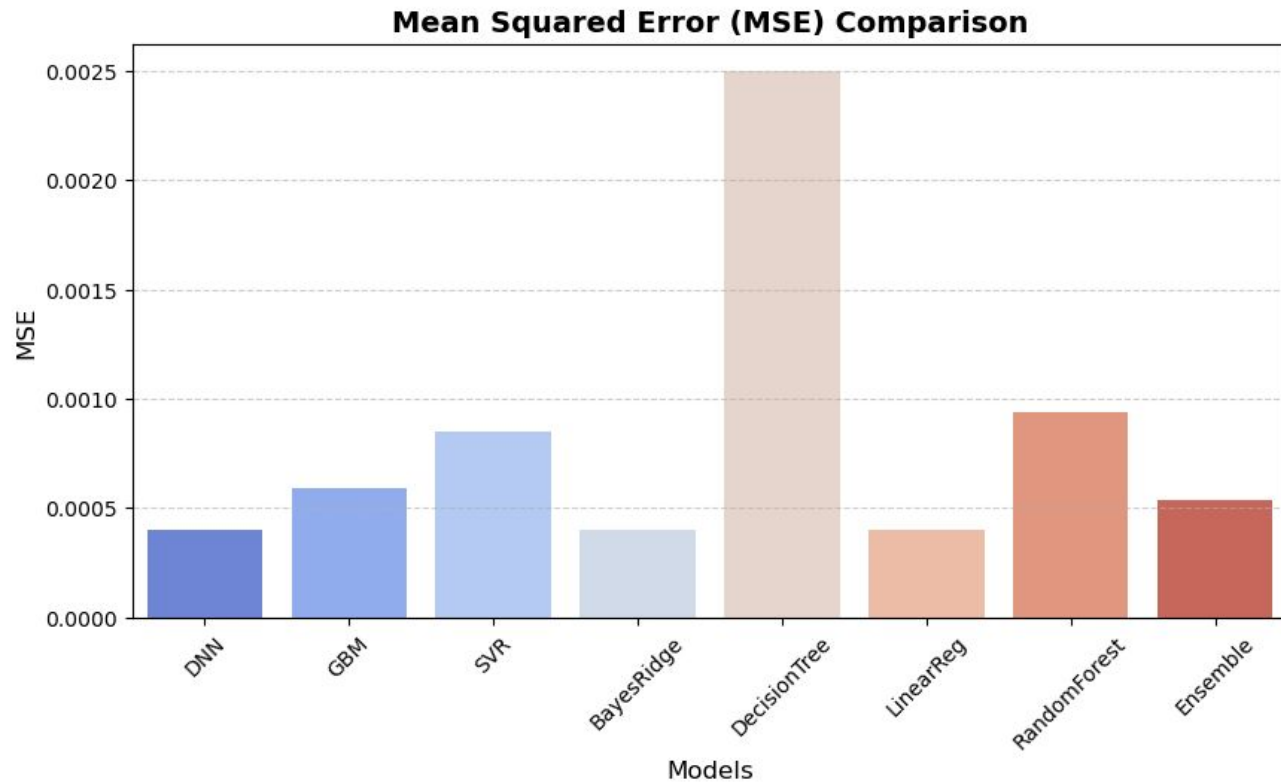
Final prediction:

$$p_{ensemble} = \sum w_i p_i$$

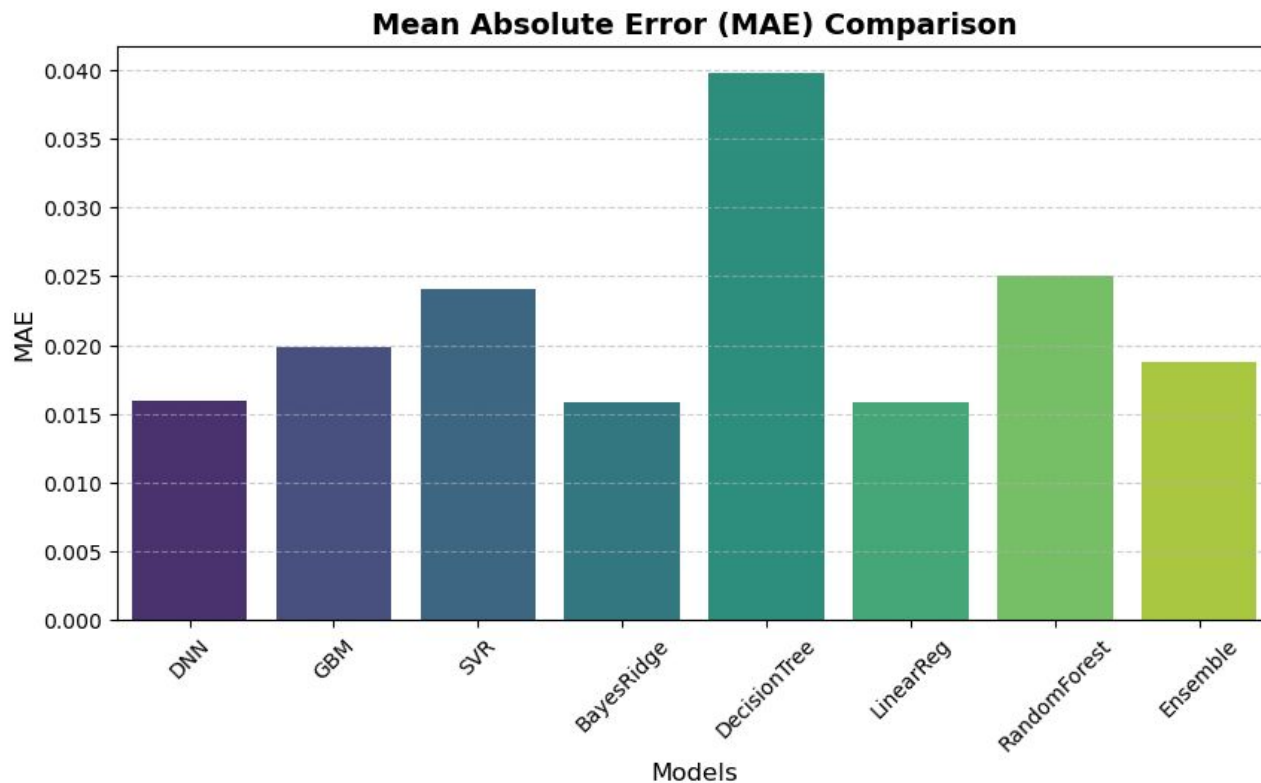# Performance Comparison

|  | MSE | MAE | R2 Score |
|---|---|---|---|
| **DNN** | 0.00040 | 0.01595 | 0.84504 |
| **GBM** | 0.00059 | 0.01987 | 0.77274 |
| **SVR** | 0.00085 | 0.02409 | 0.67283 |
| **BayesRidge** | 0.00040 | 0.01580 | 0.84434 |
| **DecisionTree** | 0.00250 | 0.03979 | 0.03739 |
| **LinearReg** | 0.00040 | 0.01580 | 0.84434 |
| **RandomForest** | 0.00094 | 0.02503 | 0.63808 |
| **Ensemble** | 0.00054 | 0.01877 | 0.79343 |

# Results



Mean Squared Error (MSE) Comparison

# Results



**Mean Absolute Error (MAE) Comparison**

# Results



R² Score Comparison

# Results



Ensemble Model Predictions vs. True Values

# Thanks!