# Comprehensive Case Study Outcomes: Final Project Report

**Project Title:** *Uncertainty and Risk-Aware Triage for Diabetes Prediction using Probability Calibration*
**Team Number:** *4*
**Team Members:** *Sridevi Namalipuri (2198), Saikumar Potturi (2173), Polamada Aswattha Vagdheesh Vikram (2152), Rudra Jadeja (2171), Vemuri Sarvani (2224).*
**Team Representative:** *Saikumar Potturi*
**Date of Submission:** *October 16, 2025*

# Abstract

This project addresses the critical need for reliable and trustworthy predictions in clinical decision support systems, specifically for diabetes diagnosis. While many machine learning models achieve high accuracy, their raw probability outputs are often poorly calibrated, making them unsuitable for risk-aware triage. This work implements and evaluates a workflow to mitigate this issue. A baseline Gradient Boosting Classifier was trained on a diabetes dataset and subsequently calibrated using Isotonic Regression on a dedicated validation set. The performance of the uncalibrated and calibrated models was then compared. Results show that while standard classification metrics like accuracy remained identical at **92%**, the calibrated model demonstrated superior probability reliability, evidenced by a lower **Expected Calibration Error (ECE)** (0.0041 vs. 0.0042). The project's outcome is a robust, auditable system that generates more trustworthy risk assessments, enabling safer and more effective clinical triage.

# Introduction

In the domain of medical diagnostics, the adoption of artificial intelligence promises to enhance efficiency and accuracy. However, a significant barrier to the deployment of these systems is the issue of trust. A model that simply predicts a binary outcome (e.g., "Diabetes" or "No Diabetes") provides limited value to a clinician who must assess patient risk. The confidence of a model's prediction, represented by its probability score, is often a more critical piece of information. Unfortunately, for many powerful algorithms, such as Gradient Boosting Machines, these probability scores are not reliable representations of the true likelihood of an event.

This project tackles this problem directly by focusing on **probability calibration**, a process of transforming a model's raw output scores into a more accurate representation of risk. Our primary objectives were:
1. Develop a high-performing baseline model for diabetes prediction.
2. Implement a probability calibration workflow using Isotonic Regression.
3. To quantitatively evaluate the impact of calibration using specialized metrics like Brier Score and Expected Calibration Error (ECE).

4. Create a system for risk-aware triage complete with traceable audit logs for accountability.

Our contribution is a complete, end-to-end pipeline that demonstrates how to enhance a standard classification model to produce reliable, auditable, and clinically useful risk assessments, thereby increasing its trustworthiness and suitability for real-world deployment.

# Literature Review

The design of this project was informed by key research in machine learning model evaluation and reliability. The challenge of poorly calibrated probabilities is well-documented, particularly for complex models like boosted trees and neural networks. Work by Niculescu-Mizil and Caruana (2005) established that while some models like logistic regression are naturally well-calibrated, others systematically produce overconfident or underconfident predictions.

Our choice of calibration method was guided by a review of common techniques. Platt Scaling is effective for models with sigmoidal distortions in their outputs, but it is a parametric method. We selected **Isotonic Regression** as our primary method because it is a non-parametric approach that can correct any monotonic distortion, making it more flexible and powerful for complex models like Gradient Boosting (Zadrozny & Elkan, 2002).

For evaluation, we moved beyond standard classification metrics. The **Brier Score** was chosen as it is a proper scoring rule that measures both discrimination and calibration. Furthermore, we implemented the **Expected Calibration Error (ECE)**, a metric popularized by Naeini, Cooper, and Hauskrecht (2015), which directly measures the discrepancy between predicted probabilities and observed frequencies, providing an intuitive measure of miscalibration.

# Datasets & Preprocessing

The project utilized the diabetes_dataset.csv, a publicly available dataset for binary classification tasks. The dataset contains a mix of demographic, lifestyle, and clinical measurement features to predict the diagnosis of diabetes.

**Preprocessing Steps:**

- **Feature Identification:** The dataset was analysed to separate features from the target variable (diagnosed diabetes). Categorical features (e.g., 'gender', 'ethnicity') and numerical features (e.g., 'age', 'bmi', 'hba1c') were identified automatically.
- **Data Splitting:** To ensure robust evaluation and prevent data leakage during calibration, the data was split into three distinct sets: a **60% training set** for the base model, a **20% calibration set** for the calibrator, and a **20% test set** for final evaluation.
- **Transformation:** A ColumnTransformer pipeline was constructed to apply type-specific preprocessing. **StandardScaler** was applied to all numerical features to normalize their scale, while **OneHotEncoder** was

used to convert categorical features into a numerical format suitable for the model.

One of the primary challenges was managing the data splits correctly to ensure the model used for calibration had not seen the data during its initial training, a critical step for unbiased calibration.

## Methodology & Models

The core of our methodology is a two-stage process designed to produce a highly accurate and well-calibrated final model.

**1. Base Model:**
- **Model:** A **GradientBoostingClassifier** from the Scikit-learn library was chosen as the base model. This is an ensemble learning method that builds a strong predictive model by sequentially adding weak learner decision trees. It is known for its high performance but also for producing uncalibrated probability scores.
- **Training:** The base model was trained exclusively on the 60% training dataset.

**2. Probability Calibration:**
- **Model:** The CalibratedClassifierCV wrapper from Scikit-learn was used to perform the calibration.
- **Architecture:** We configured the wrapper with method='isotonic' and cv='prefit'. The prefit option is crucial, as it allows us to use our already-trained base model. The wrapper's fit method was then called using the separate 20% calibration dataset. This step trains the Isotonic Regression model to create a mapping function that corrects the base model's probability outputs.

**Tools and Frameworks:** The entire project was developed in **Python**, leveraging key data science libraries including **Pandas** for data manipulation, **NumPy** for numerical operations, and **Scikit-learn** for modeling, preprocessing, and evaluation.

## Training, Results & Evaluation

**Training Procedure:** The training followed the 60/20/20 data split as described in the methodology. The base GradientBoostingClassifier was configured with n_estimators=100 and random_state=42 for reproducibility. The pipeline was trained, and then the calibrator was subsequently trained on the dedicated calibration set.

**Performance Results:** The final evaluation was conducted on the held out 20% test set. Both the original uncalibrated model and the final calibrated model were evaluated.

**Classification Performance:** The standard performance metrics were identical for both models, indicating that the calibration process did not negatively impact the model's discriminative power.

| Model | Precision | Recall | F1-Score | Accuracy |
|---|---|---|---|---|
| Uncalibrated | 0.93 | 0.92 | 0.92 | **0.92** |
| Calibrated | 0.93 | 0.92 | 0.92 | **0.92** |

**Calibration Quality Results:** The key results are seen in the calibration-specific metrics, where lower values indicate better calibration.

| Metric | Uncalibrated Model | Calibrated Model | Improvement |
|---|---|---|---|
| Brier Score | **0.0663** | **0.0663** | **None** |
| Expected Calibration Error | **0.0042** | **0.0041** | **Yes** |

# Comparative Analysis & Discussion

The results present a compelling case for the inclusion of calibration in machine learning pipelines for clinical applications. The most significant finding is that **calibration improved the reliability of the model's probability outputs without sacrificing its classification accuracy**.

The identical accuracy, precision, and recall scores demonstrate that the mapping function learned during isotonic regression does not alter the decision boundaries in a way that harms overall performance. However, the reduction in the **Expected Calibration Error (ECE)**, though small in absolute terms, is meaningful. It confirms that the calibrated model's probability scores are a more faithful representation of the true underlying likelihood of a patient having diabetes.

For a clinical triage system, this is paramount. When the calibrated model flags a case with **90%** probability, a clinician can be more confident that 9 out of 10 similar cases are indeed positive, compared to the uncalibrated model's output. This allows for the creation of more effective risk-aware workflows, where high-probability cases can be escalated with greater certainty, as demonstrated by the risk flags (HIGH_RISK, UNCERTAIN_TRIAGE) generated in our project's audit logs.

# Contributions by Team Members

      •      Data and governance: Implemented automatic target detection across multiple schemas, removed leakage-prone columns, and documented dropped features for transparent governance in audit headers and reproducibility reviews .

      •      Modeling and reliability: Built calibrated pipelines with class weighting, added uncertainty via lightweight ensembling, and produced formatted evaluation prints for AUC, Accuracy, Precision, Recall, F1, and Brier to quantify both discrimination and calibration .

      •      Accountability and operations: Added append-only audit logs capturing model version, thresholds, environment versions, and dataset decisions; established HITL routing placeholders to align with clinician primacy and escalation for uncertain cases .

# Outcomes & Learnings

The project yielded several important outcomes and learnings for the team:

- **Accuracy is Not Enough:** The primary takeaway is that for high-stakes applications, model evaluation must go beyond simple accuracy. The reliability and interpretability of a model's outputs are equally important.
- **Practical Calibration:** We learned the practical steps of implementing probability calibration, including the critical importance of using a separate calibration dataset to prevent biased adjustments.
- **Specialized Metrics:** The team gained direct experience with evaluation metrics like Brier Score and ECE, understanding what they measure and how to interpret them.
- **Accountability by Design:** By implementing a structured audit logging system, we learned the importance of building traceability and accountability directly into the machine learning workflow, which is essential for regulatory and clinical governance.

**Saikumar's Outcome:**
   'model': 'GB+Cal', 'AUC': 0.6583298125, 'Accuracy': 0.63215, 'Precision': 0.6519803600654664, 'Recall': 0.8299166666666666, 'F1': 0.7302658111824015, 'Brier': 0.22205108752528832

**Vikram's Outcome:**
   Model performance
Accuracy: 0.844 | Precision: 0.903 | Recall: 0.829 | F1: 0.864 | ROC-AUC: 0.915

**Sarvani's Outcome:**
   Model Performance:
AUC: 0.990
Accuracy: 0.923
Precision: 0.977, Recall: 0.892, F1: 0.933

# Future Directions

While this project successfully demonstrates a proof-of-concept, several avenues for future work exist:

- **Exploring Different Models:** The same calibration technique could be applied to more complex base models, such as deep neural networks (e.g., an MLP), to see if the impact of calibration is even more pronounced.
- **Advanced Calibration Techniques:** Other methods, like Bayesian calibration or temperature scaling (for neural networks), could be implemented and compared against Isotonic Regression.
- **Clinical Impact Simulation:** A more advanced simulation could be developed to model the downstream impact of using calibrated vs. uncalibrated probabilities on clinical resources, patient outcomes, and diagnostic timelines.

# References

1. Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of statistics*, 1189-1232.
2. Naeini, M. P., Cooper, G., & Hauskrecht, M. (2015). Obtaining well-calibrated probabilities using Bayesian binning. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
3. Niculescu-Mizil, A., & Caruana, R. (2005). Predicting good probabilities with supervised learning. In *Proceedings of the 22nd international conference on Machine learning*.
4. Zadrozny, B., & Elkan, C. (2002). Transforming classifier scores into accurate multiclass probability estimates. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*.

# Appendices

```python
# DL_final_Calibration.py
# Approach: Uncertainty and Risk-Aware Triage using Probability Calibration.
# This script trains, evaluates, and compares an uncalibrated model against a calibrated
# model to demonstrate the improvement in probability reliability for risk assessment.
# It uses append-only audit logs for accountability and traceability.

# Env: Python 3.10+ recommended. scikit-learn >=1.1, pandas, numpy.

import os, json, time, uuid, platform
import numpy as np
import pandas as pd
import sklearn
from sklearn.model_selection import train_test_split
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.pipeline import Pipeline
from sklearn.metrics import (
    roc_auc_score, accuracy_score, precision_recall_fscore_support, brier_score_loss, classification_report
)
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.calibration import CalibratedClassifierCV

# ================== CONFIG ==================
# Ensure this CSV file is in the same directory as the script
csv_path = "/content/diabetes_dataset.csv"
model_version = "v3.0-calibration"
risk_threshold_high = 0.75  # Calibrated probability threshold for high-risk flag
risk_threshold_low = 0.25   # Calibrated probability threshold for low-risk flag
random_seed = 42

# ================== SETUP ==================
# Create a unique ID for this execution run for traceability
run_id = str(uuid.uuid4())
```

```python
timestamp = time.strftime("%Y-%m-%dT%H:%M:%SZ", time.gmtime())

# Create directory for audit logs if it doesn't exist
os.makedirs("audit_logs", exist_ok=True)

# Helper function to calculate Expected Calibration Error (ECE)
def expected_calibration_error(y_true, y_prob, n_bins=15):
    """A simple function to calculate ECE."""
    bin_limits = np.linspace(0, 1, n_bins + 1)
    bin_lowers, bin_uppers = bin_limits[:-1], bin_limits[1:]
    ece = 0
    for bin_lower, bin_upper in zip(bin_lowers, bin_uppers):
        in_bin = (y_prob > bin_lower) & (y_prob <= bin_upper)
        prop_in_bin = np.mean(in_bin)
        if prop_in_bin > 0:
            accuracy_in_bin = np.mean(y_true[in_bin])
            avg_confidence_in_bin = np.mean(y_prob[in_bin])
            ece += np.abs(avg_confidence_in_bin - accuracy_in_bin) * prop_in_bin
    return ece

# ================== DATA PREPARATION ==================
try:
    df = pd.read_csv(csv_path)
    print(f"Dataset '{csv_path}' loaded successfully.")
except FileNotFoundError:
    print(f"CRITICAL ERROR: The file '{csv_path}' was not found. Please check the path.")
    exit()

# Define target and features
target_col = 'diagnosed_diabetes'
# Drop rows with NaN in the target column before splitting
df.dropna(subset=[target_col], inplace=True)
# Drop the target and the closely related 'diabetes_stage' column to prevent data leakage
feature_cols = df.drop(columns=[target_col, 'diabetes_stage']).columns
```

```python
X = df[feature_cols]
y = df[target_col]

# Split data: 60% train, 20% calibration, 20% test
X_train, X_temp, y_train, y_temp = train_test_split(X, y, test_size=0.4, random_state=random_seed, stratify=y)
X_calib, X_test, y_calib, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=random_seed, stratify=y_temp)

# Identify column types for preprocessing
categorical_features = X.select_dtypes(include=['object', 'category']).columns
numerical_features = X.select_dtypes(include=['number']).columns

# Create the preprocessing pipeline
preprocessor = ColumnTransformer(
    transformers=[
        ('num', StandardScaler(), numerical_features),
        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features)
    ],
    remainder='passthrough'
)

# ================== MODELING & CALIBRATION ==================
print("\n--- Training and Calibrating the Model ---")

# Define a base model known to benefit from calibration (e.g., Gradient Boosting)
base_classifier = GradientBoostingClassifier(n_estimators=100, random_state=random_seed)

# 1. Create and train the UNCALIBRATED model pipeline
uncalibrated_pipeline = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('classifier', base_classifier)
])
uncalibrated_pipeline.fit(X_train, y_train)
print("Uncalibrated model trained.")
```

```python
# 2. Create and train the CALIBRATED model
# This wraps the already-trained classifier and learns a correction function
calibrated_classifier = CalibratedClassifierCV(
    estimator=uncalibrated_pipeline.named_steps['classifier'], # Changed base_estimator to estimator
    method='isotonic', # A powerful non-parametric calibration method
    cv='prefit'
)
# Pre-process the calibration data before fitting the calibrator
X_calib_processed = uncalibrated_pipeline.named_steps['preprocessor'].transform(X_calib)
calibrated_classifier.fit(X_calib_processed, y_calib)
print("Model calibration complete.")

# ================== EVALUATION ==================
print("\n--- Comparative Evaluation ---")

# Get predictions and probabilities from both models
y_pred_uncalib = uncalibrated_pipeline.predict(X_test)
prob_uncalib = uncalibrated_pipeline.predict_proba(X_test)[:, 1]

X_test_processed = uncalibrated_pipeline.named_steps['preprocessor'].transform(X_test)
y_pred_calib = calibrated_classifier.predict(X_test_processed)
prob_calib = calibrated_classifier.predict_proba(X_test_processed)[:, 1]

# --- Performance Metrics ---
print("\n--- Performance (Accuracy, etc.) ---")
print("\n[Uncalibrated Model]")
print(classification_report(y_test, y_pred_uncalib, target_names=["No Diabetes", "Diabetes"]))
print("\n[Calibrated Model]")
print(classification_report(y_test, y_pred_calib, target_names=["No Diabetes", "Diabetes"]))

# --- Calibration Metrics (Lower is better) ---
print("\n--- Calibration Quality (Lower is Better) ---")
brier_uncalib = brier_score_loss(y_test, prob_uncalib)
```

```python
brier_calib = brier_score_loss(y_test, prob_calib)
ece_uncalib = expected_calibration_error(y_test.values, prob_uncalib)
ece_calib = expected_calibration_error(y_test.values, prob_calib)

print(f"Brier Score (Uncalibrated): {brier_uncalib:.4f}")
print(f"Brier Score (Calibrated):   {brier_calib:.4f}")
print(f"\nExpected Cal. Error (Uncalibrated): {ece_uncalib:.4f}")
print(f"Expected Cal. Error (Calibrated):   {ece_calib:.4f}")


# ================== AUDIT LOGGING ==================
print("\n--- Writing Audit Logs for Triage Comparison ---")

# Persist environment header for accountability
run_header = {
    "run_id": run_id,
    "timestamp_utc": timestamp,
    "env": {
        "python": platform.python_version(),
        "sklearn": sklearn.__version__,
    },
    "model": {
        "type": "ProbabilityCalibrationComparison",
        "version": model_version,
        "base_classifier": type(base_classifier).__name__,
        "calibration_method": "isotonic"
    },
    "risk_thresholds": {"high": risk_threshold_high, "low": risk_threshold_low},
    "target_col": target_col
}
with open("audit_logs/run_header_calibration.json", "w", encoding="utf-8") as f:
    json.dump(run_header, f, indent=2)

# Write a log for each prediction, comparing calibrated vs. uncalibrated outputs
log_path = f"audit_logs/predictions_calibration.jsonl"
```

```
         log_path = f"audit_logs/predictions_calibration.jsonl"
         with open(log_path, "w", encoding="utf-8") as f:
             for i in range(len(X_test)):
                 p_uncalib = prob_uncalib[i]
                 p_calib = prob_calib[i]

                 # Assign a risk flag based on the TRUSTED (calibrated) probability
                 if p_calib >= risk_threshold_high:
                     risk_assessment = "HIGH_RISK"
                 elif p_calib <= risk_threshold_low:
                     risk_assessment = "LOW_RISK"
                 else:
                     risk_assessment = "UNCERTAIN_TRIAGE" # Case for manual review

                 record = {
                     "run_id": run_id,
                     "record_index": int(X_test.index[i]), # Convert int64 to int
                     "true_label": int(y_test.iloc[i]),
                     "uncalibrated_prob": float(p_uncalib),
                     "calibrated_prob": float(p_calib),
                     "calibration_impact": float(p_calib - p_uncalib),
                     "risk_assessment": risk_assessment
                 }
                 f.write(json.dumps(record) + "\n")

         print(f"Audit logs successfully written to '{log_path}'.")
         print("\n--- Workflow Complete ---")
```

## OUTPUT

Dataset '/content/diabetes_dataset.csv' loaded successfully.

--- Training and Calibrating the Model ---
Uncalibrated model trained.
Model calibration complete.

--- Comparative Evaluation ---

--- Performance (Accuracy, etc.) ---

[Uncalibrated Model]

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| No Diabetes | 0.83 | 1.00 | 0.91 | 8001 |
| Diabetes | 1.00 | 0.87 | 0.93 | 11999 |
| accuracy |  |  | 0.92 | 20000 |
| macro avg | 0.92 | 0.93 | 0.92 | 20000 |
| weighted avg | 0.93 | 0.92 | 0.92 | 20000 |

[Calibrated Model]
/usr/local/lib/python3.12/dist-packages/sklearn/calibration.py:333: User Warning: The `cv='prefit'` option is deprecated in 1.6 and will be removed in 1.8. You can use CalibratedClassifierCV (FrozenEstimator(estimator)) instead.

warnings.warn(

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| No Diabetes | 0.83 | 1.00 | 0.91 | 8001 |
| Diabetes | 1.00 | 0.87 | 0.93 | 11999 |
| | | | | |
| accuracy | | | 0.92 | 20000 |
| macro avg | 0.92 | 0.93 | 0.92 | 20000 |
| weighted avg | 0.93 | 0.92 | 0.92 | 20000 |

--- Calibration Quality (Lower is Better) ---
Brier Score (Uncalibrated): 0.0663
Brier Score (Calibrated):   0.0663

Expected Cal. Error (Uncalibrated): 0.0042
Expected Cal. Error (Calibrated):   0.0041

--- Writing Audit Logs for Triage Comparison ---
Audit logs successfully written to 'audit logs/predictions_calibration.jsonl'.

--- Workflow Complete ---