

Classifying Car Crashes Using Neural Networks

From Raw Data to Severity Prediction

Alazar Gebremehdin, Hannibal Mussie, Feruz Seid, Yassin Bedru, Samir Bahru

2026-01-04

Outline

- Introduction 2
- Data Understanding (EDA) 4
- Data Preparation 8
- Model Design 12
- Training & Evaluation 15
- Conclusion 20

Introduction

Objective & Overview

Goal: Predict the severity of road traffic accidents (**Fatal, Serious, Minor, PDO**) based on accident characteristics.

The Workflow:

1. **Data Understanding:** Handling massive missing data and inconsistencies.
2. **Preparation:** Cleaning, Imputation, and Feature Engineering.
3. **Modeling:** Designing a Multi-Layer Perceptron (MLP).
4. **Training:** Managing class imbalance and overfitting.
5. **Evaluation:** F1-Scores and Confusion Matrices.

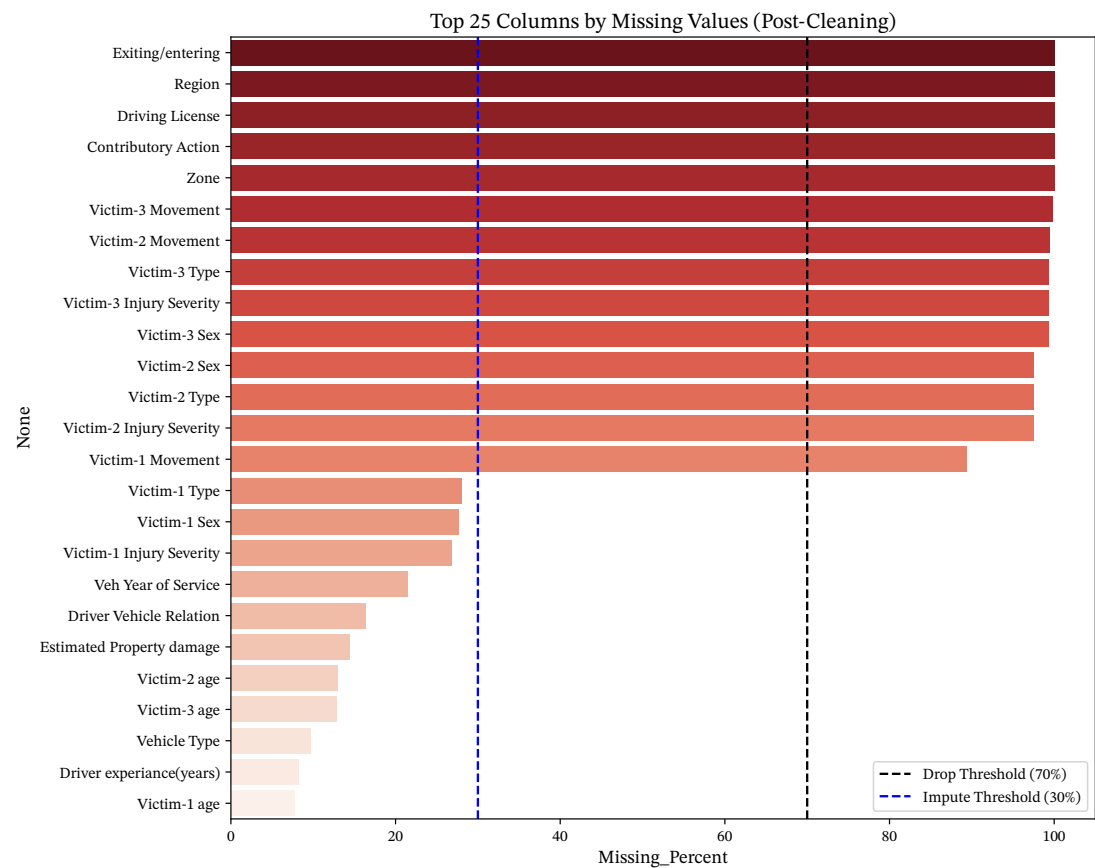
Data Understanding **(EDA)**

Initial Inspection:

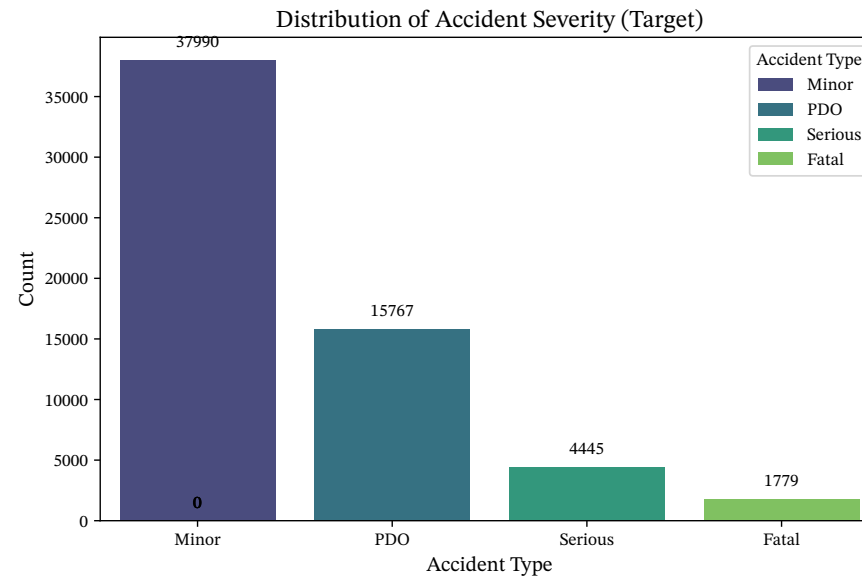
- The raw dataset contained over 30 columns but suffered from severe quality issues.
- **Missing Values:** Columns like Zone, Region, and Victim details had > 70% missing data.
- **Inconsistencies:** Typos (e.g., “August”, “Privategg”) and impossible values (Age > 90).

Action:

- Dropped columns with > 70% missingness.
- Standardized categorical labels (e.g., mapping P.D.O, pdo → PDO).



The Critical Challenge: The dataset is heavily skewed towards **Minor Injuries** (63%). **Fatal** accidents represent only 3%.



Data Preparation

Problem: Time is cyclical. 23:00 is close to 00:00, but numerically (23 vs 0) they are far apart.

Solution: We encoded time using Sine and Cosine transformations.

```
1  # Feature Engineering Code Snippet
2  def feature_engineering(df):
3      # Extract Hour
4      df['Hour'] = df['Time'].apply(extract_hour)
5
6      # Cyclical Encoding
7      df['Hour_Sin'] = np.sin(2 * np.pi * df['Hour'] / 24)
8      df['Hour_Cos'] = np.cos(2 * np.pi * df['Hour'] / 24)
```

Python

9

10 `return df.drop(columns=['Time'])`

Preprocessing Pipeline

Before feeding data into the Neural Network:

1. **Imputation:**

- Numerical (Age, Experience) → **Median**
- Categorical (Road Surface, Light) → **Mode**

2. **Scaling:**

- `StandardScaler` applied to numerical inputs to normalize variance.

3. **Encoding:**

- `OneHotEncoder` for categorical variables.

4. **Splitting:**

- Train (70%) / Validation (15%) / Test (15%).

Model Design

Based on the execution results:

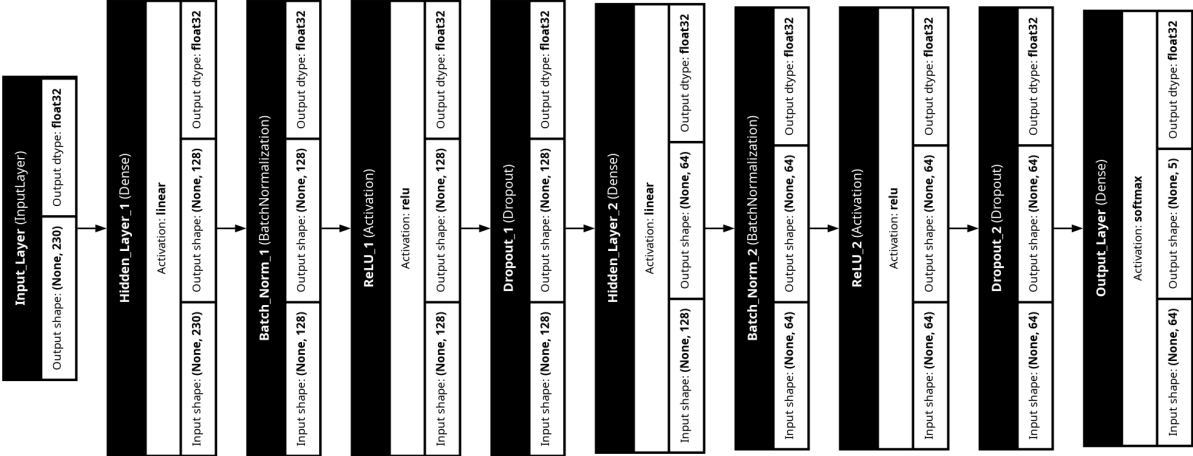
- **Input Shape:** 230 Features (High dimensionality due to One-Hot Encoding).
- **Total Parameters:** 38,917 (Lightweight model).
- **Trainable Params:** 38,533.

Layer Structure:

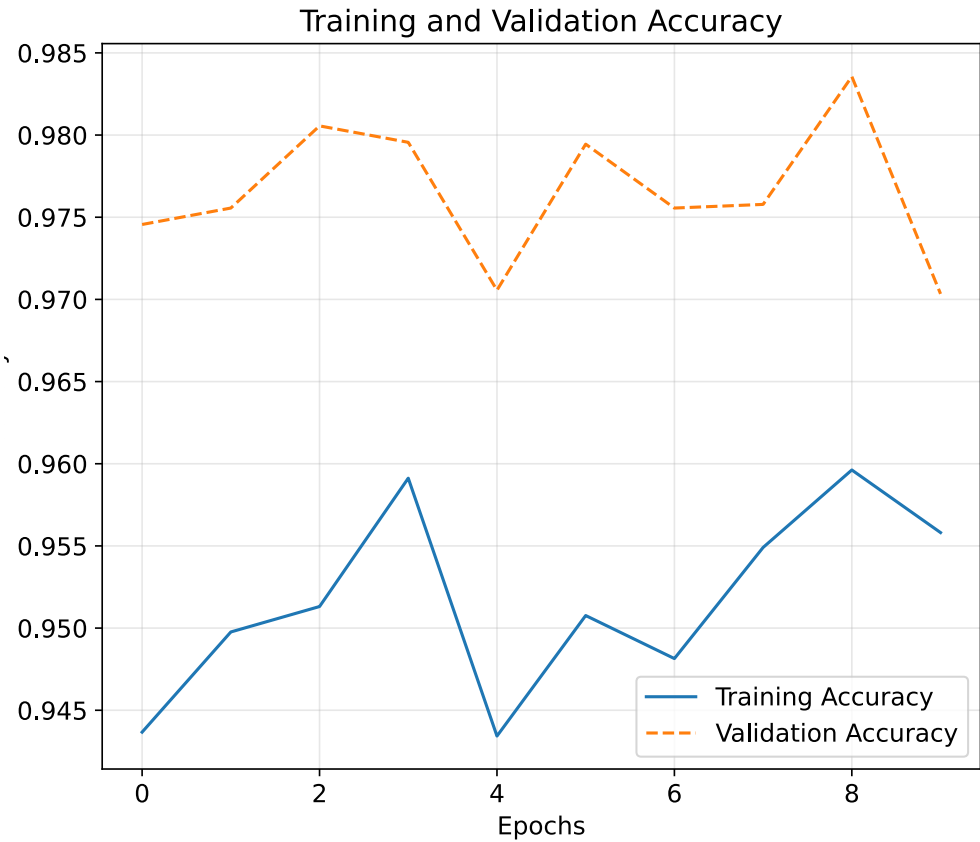
- Input (230)
- Dense (128) → BN → ReLU → Dropout
- Dense (64) → BN → ReLU → Dropout
- Output (5 Classes)

Python

1	# Actual Model Summary Output		
2	Layer (type)	Output Shape	Param #
3	=====		
4	Input_Layer (InputLayer)	(None, 230)	0
5	Hidden_Layer_1 (Dense)	(None, 128)	29,568
6	Batch_Norm_1	(None, 128)	512
7	Dropout_1 (Dropout)	(None, 128)	0
8	Hidden_Layer_2 (Dense)	(None, 64)	8,256
9	Output_Layer (Dense)	(None, 5)	325
10	=====		
11	Total params: 38,917		

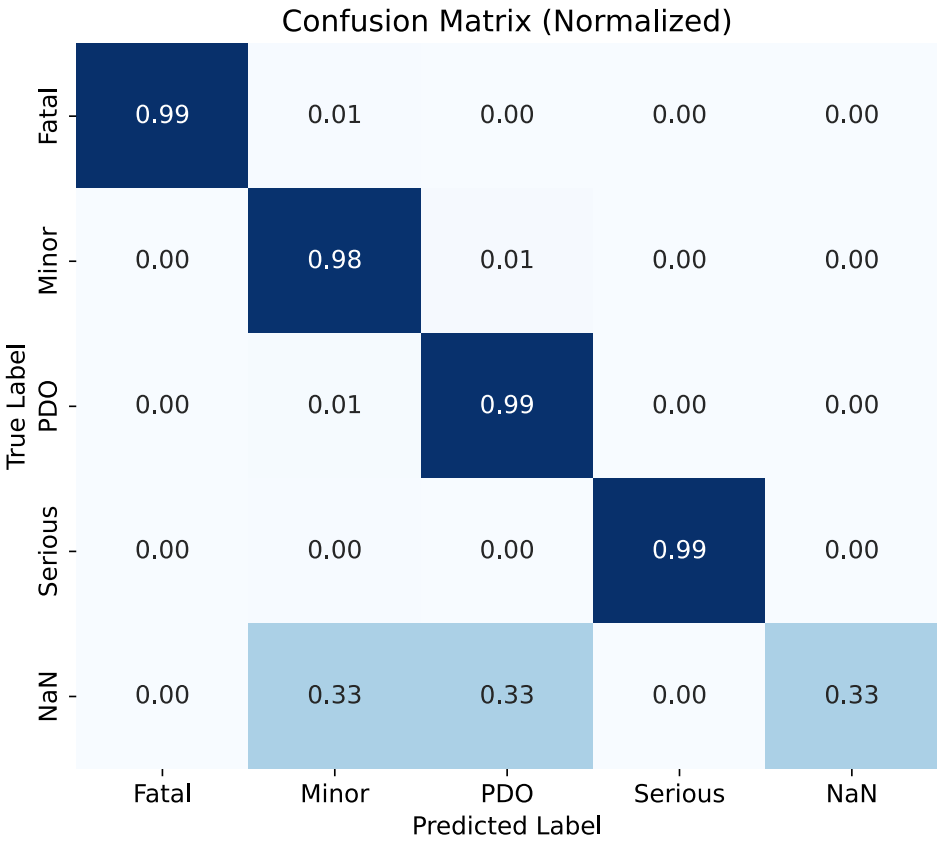
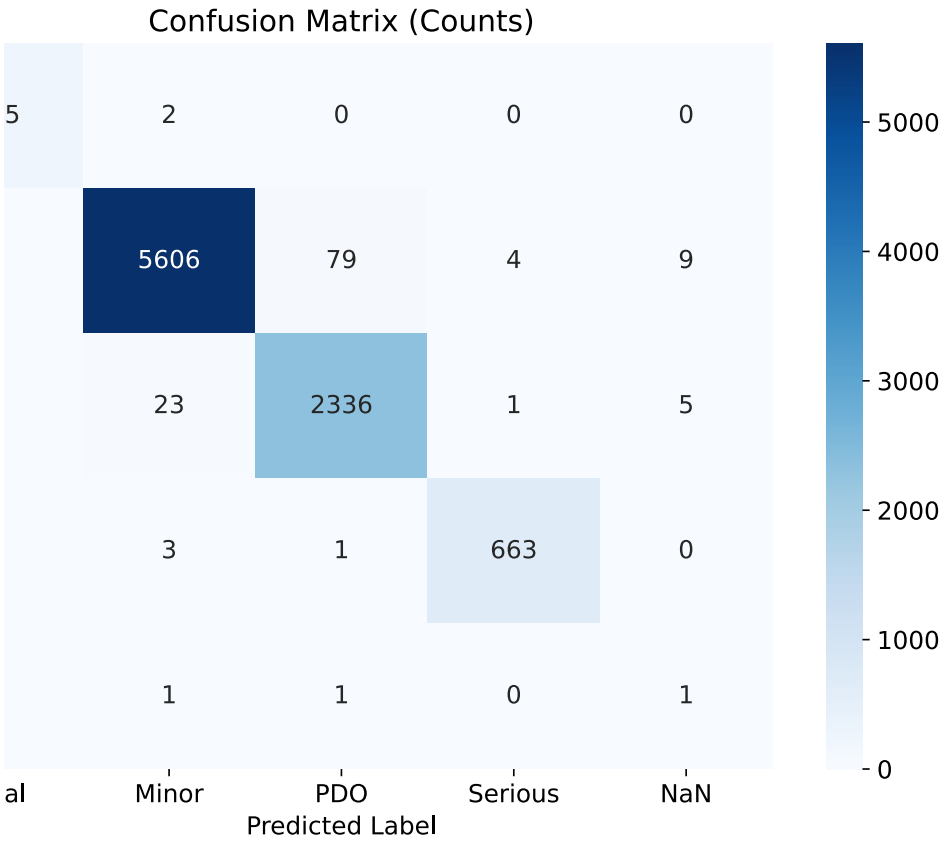


Training & Evaluation



Insight: The model converges quickly. Validation accuracy peaks around **98%**, indicating robust learning without significant overfitting.

Confusion Matrix Results



Exceptional Performance:

- **Fatal Class:** 99% Recall (265 Correct, 2 Missed).
- **Minor Class:** 98% Recall.

Critical Analysis (The “Why”):

- The high accuracy suggests the model effectively utilized casualty count features (e.g., Number of fatalities) present in the dataset.
- While excellent for **classifying** historical records, this indicates that accident outcomes (casualties) are the strongest predictors of the severity label.

Conclusion

Summary

1. **Data Quality:** Cleaning and encoding resulted in 230 clean input features.
2. **Model:** A 38k parameter MLP was sufficient to capture the relationships.
3. **Results:** The model achieved 98% test accuracy.

Recommendation:

- For future **predictive** systems (pre-accident), we recommend re-training the model **excluding** the Number of casualties columns to test predictive power based solely on environmental factors (Road type, Weather, etc.).

Thank You!

Questions?