

forest_cover

December 12, 2025

```
[60]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.utils import class_weight
from sklearn.linear_model import LogisticRegression
import tensorflow as tf
from tensorflow import keras
from tensorflow.keras import layers, callbacks
from sklearn.metrics import classification_report, confusion_matrix, \
    accuracy_score
```

```
[61]: fc_data = pd.read_csv('covtype.csv')
fc_data.head()
```

```
[61]:
```

	Elevation	Aspect	Slope	Horizontal_Distance_To_Hydrology	\
0	2596	51	3	258	
1	2590	56	2	212	
2	2804	139	9	268	
3	2785	155	18	242	
4	2595	45	2	153	

	Vertical_Distance_To_Hydrology	Horizontal_Distance_To_Roadways	\
0	0	510	
1	-6	390	
2	65	3180	
3	118	3090	
4	-1	391	

	Hillshade_9am	Hillshade_Noon	Hillshade_3pm	\
0	221	232	148	
1	220	235	151	
2	234	238	135	
3	238	238	122	
4	220	234	150	

	Horizontal_Distance_To_Fire_Points	...	Soil_Type32	Soil_Type33	\
0	6279	...	0	0	
1	6225	...	0	0	
2	6121	...	0	0	
3	6211	...	0	0	
4	6172	...	0	0	

	Soil_Type34	Soil_Type35	Soil_Type36	Soil_Type37	Soil_Type38	\
0	0	0	0	0	0	
1	0	0	0	0	0	
2	0	0	0	0	0	
3	0	0	0	0	0	
4	0	0	0	0	0	

	Soil_Type39	Soil_Type40	Cover_Type
0	0	0	5
1	0	0	5
2	0	0	2
3	0	0	2
4	0	0	5

[5 rows x 55 columns]

0.1 Exploratory Data Analysis (EDA)

```
[62]: fc_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 581012 entries, 0 to 581011
Data columns (total 55 columns):
#   Column                                     Non-Null Count  Dtype
---  -
0   Elevation                                581012 non-null  int64
1   Aspect                                  581012 non-null  int64
2   Slope                                   581012 non-null  int64
3   Horizontal_Distance_To_Hydrology         581012 non-null  int64
4   Vertical_Distance_To_Hydrology           581012 non-null  int64
5   Horizontal_Distance_To_Roadways          581012 non-null  int64
6   Hillshade_9am                           581012 non-null  int64
7   Hillshade_Noon                           581012 non-null  int64
8   Hillshade_3pm                           581012 non-null  int64
9   Horizontal_Distance_To_Fire_Points       581012 non-null  int64
10  Wilderness_Area1                         581012 non-null  int64
11  Wilderness_Area2                         581012 non-null  int64
12  Wilderness_Area3                         581012 non-null  int64
13  Wilderness_Area4                         581012 non-null  int64
14  Soil_Type1                              581012 non-null  int64
15  Soil_Type2                              581012 non-null  int64
```

16	Soil_Type3	581012	non-null	int64
17	Soil_Type4	581012	non-null	int64
18	Soil_Type5	581012	non-null	int64
19	Soil_Type6	581012	non-null	int64
20	Soil_Type7	581012	non-null	int64
21	Soil_Type8	581012	non-null	int64
22	Soil_Type9	581012	non-null	int64
23	Soil_Type10	581012	non-null	int64
24	Soil_Type11	581012	non-null	int64
25	Soil_Type12	581012	non-null	int64
26	Soil_Type13	581012	non-null	int64
27	Soil_Type14	581012	non-null	int64
28	Soil_Type15	581012	non-null	int64
29	Soil_Type16	581012	non-null	int64
30	Soil_Type17	581012	non-null	int64
31	Soil_Type18	581012	non-null	int64
32	Soil_Type19	581012	non-null	int64
33	Soil_Type20	581012	non-null	int64
34	Soil_Type21	581012	non-null	int64
35	Soil_Type22	581012	non-null	int64
36	Soil_Type23	581012	non-null	int64
37	Soil_Type24	581012	non-null	int64
38	Soil_Type25	581012	non-null	int64
39	Soil_Type26	581012	non-null	int64
40	Soil_Type27	581012	non-null	int64
41	Soil_Type28	581012	non-null	int64
42	Soil_Type29	581012	non-null	int64
43	Soil_Type30	581012	non-null	int64
44	Soil_Type31	581012	non-null	int64
45	Soil_Type32	581012	non-null	int64
46	Soil_Type33	581012	non-null	int64
47	Soil_Type34	581012	non-null	int64
48	Soil_Type35	581012	non-null	int64
49	Soil_Type36	581012	non-null	int64
50	Soil_Type37	581012	non-null	int64
51	Soil_Type38	581012	non-null	int64
52	Soil_Type39	581012	non-null	int64
53	Soil_Type40	581012	non-null	int64
54	Cover_Type	581012	non-null	int64

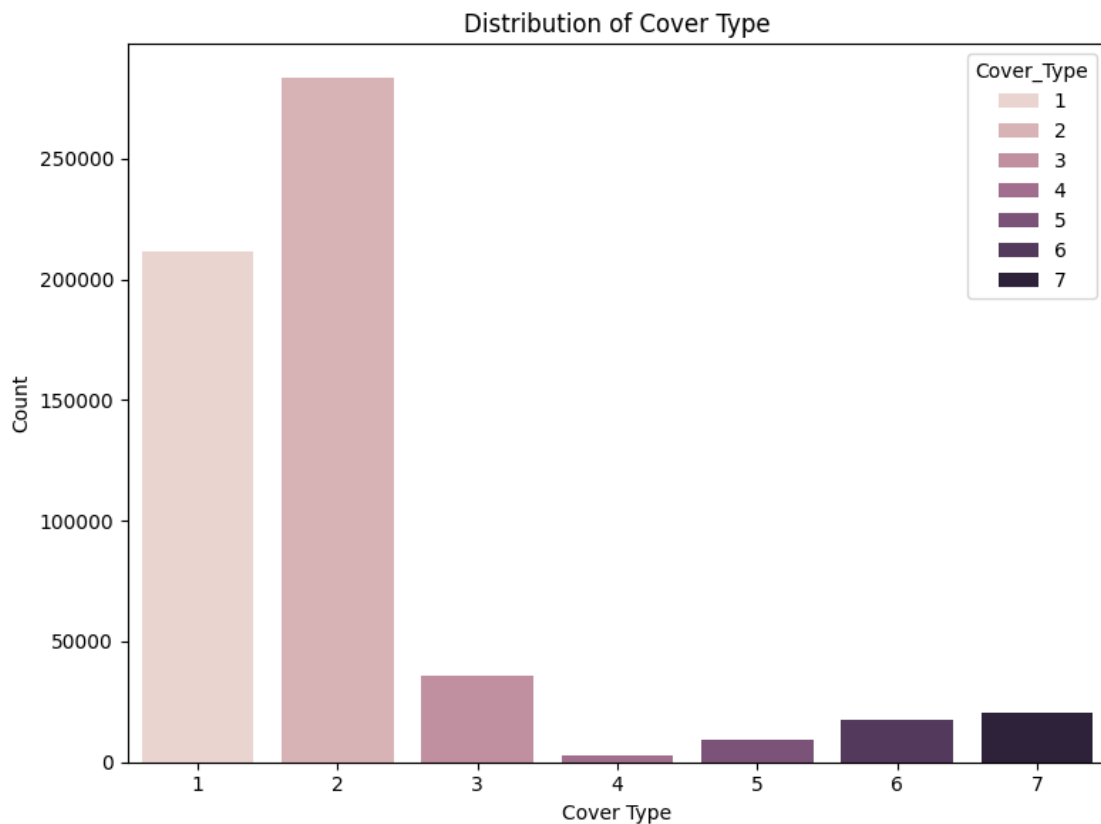
dtypes: int64(55)
memory usage: 243.8 MB

```
[63]: print(fc_data['Cover_Type'].value_counts())
```

```
Cover_Type
2    283301
1    211840
3     35754
```

```
7    20510
6    17367
5     9493
4     2747
Name: count, dtype: int64
```

```
[64]: # Target variable distribution plot
plt.figure(figsize=(8, 6))
sns.countplot(x='Cover_Type', data=fc_data, hue='Cover_Type')
plt.title('Distribution of Cover Type')
plt.xlabel('Cover Type')
plt.ylabel('Count')
plt.tight_layout()
plt.show()
```



0.1.1 Analysis of Target Variable: Cover Type Distribution

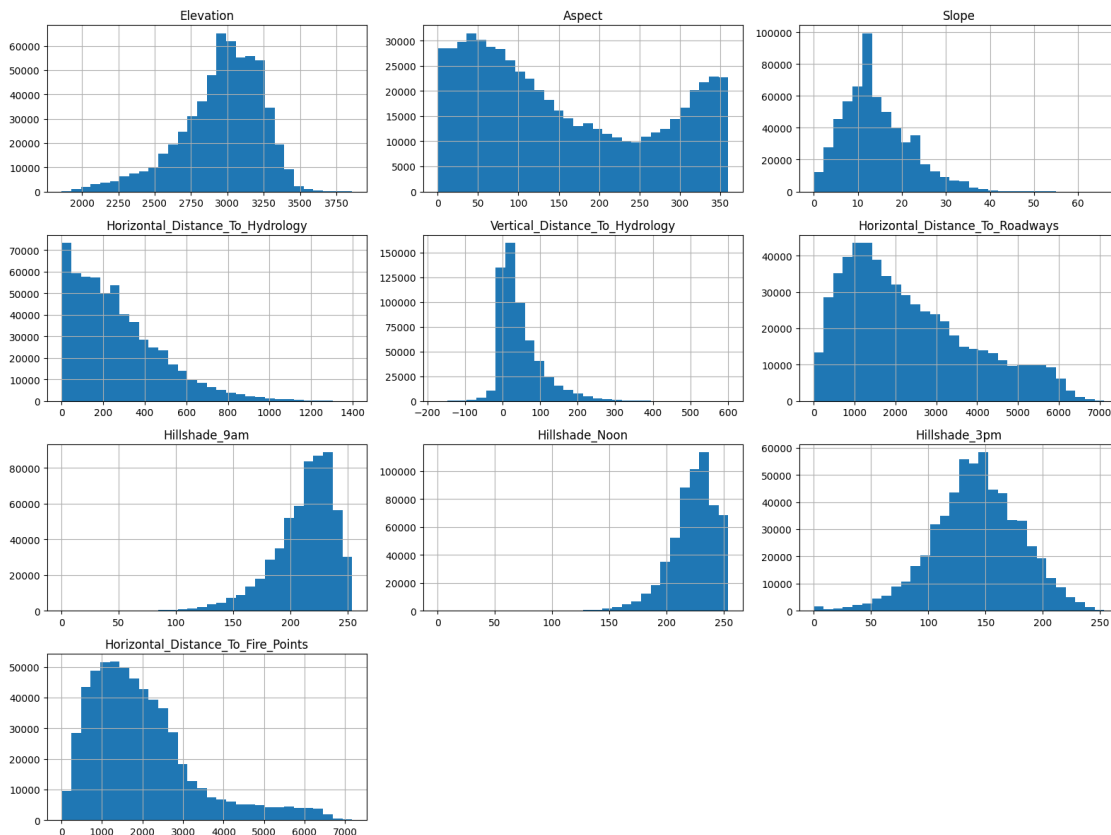
Observations The bar graph illustrates a **severe class imbalance** in the dataset. **Cover Types 1 and 2** are the dominant majority classes, accounting for a vast portion of the data, while **Cover Types 3, 4, 5, 6, and 7** are significantly underrepresented. Type 4, in particular, has an extremely low frequency compared to the others.

Required Actions for Modeling Since we are implementing **Logistic Regression** and **Neural Networks**, this imbalance will heavily bias the models toward predicting only the majority classes (1 and 2) to minimize error, leading to poor performance on the minority classes. To mitigate this, we need to apply one of the following strategies before training:

- **Resampling Techniques:** Apply **SMOTE** (Synthetic Minority Over-sampling Technique) to generate synthetic samples for the minority classes or use **Random Undersampling** on the majority classes to balance the distribution.
- **Class Weighting:** Alternatively, we can assign higher **class weights** during the training of both the Logistic Regression and Neural Network models. This penalizes the misclassification of minority classes more heavily (e.g., using `class_weight='balanced'` in Scikit-Learn or passing a weight dictionary to the Keras/TensorFlow `fit` method).

```
[65]: # 2. Numerical Features Distributions
numerical_cols = ['Elevation', 'Aspect', 'Slope',
                  'Horizontal_Distance_To_Hydrology',
                  ↪ 'Vertical_Distance_To_Hydrology',
                  'Horizontal_Distance_To_Roadways', 'Hillshade_9am',
                  ↪ 'Hillshade_Noon',
                  'Hillshade_3pm', 'Horizontal_Distance_To_Fire_Points']

fc_data[numerical_cols].hist(figsize=(16, 12), bins=30, layout=(4, 3))
plt.tight_layout()
plt.show()
```



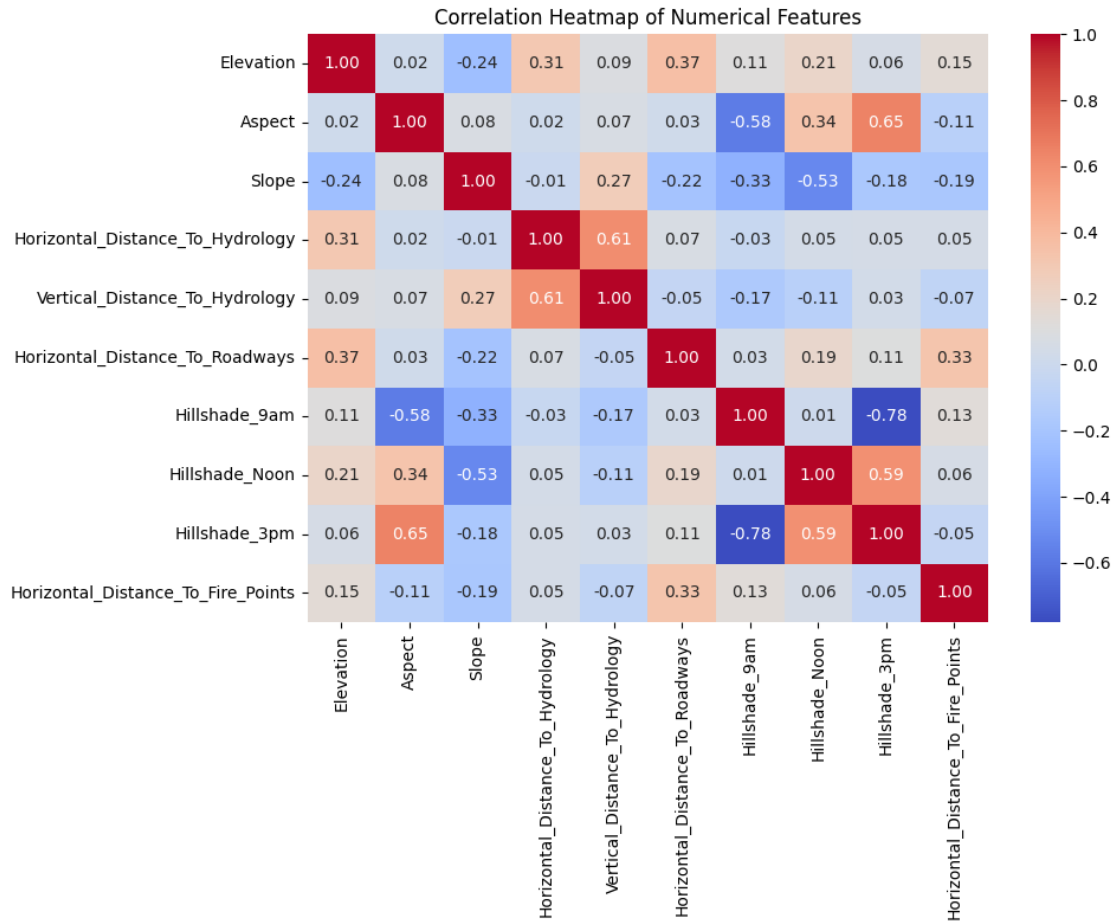
0.1.2 Analysis of Continuous Features: Feature Distributions

Observations The histograms reveal that the input features operate on vastly different scales. For instance, **Elevation** ranges roughly from 2,000 to 3,800, while **Slope** is concentrated between 0 and 60. Additionally, features like **Horizontal_Distance_To_Hydrology** and **Horizontal_Distance_To_Fire_Points** show significant **right-skewness** (long tails to the right).

Required Actions for Modeling Since we are using **Logistic Regression** and **Neural Networks**, the data scale is critical. Both models rely on gradient-based optimization (like Gradient Descent), where features with larger magnitudes (e.g., distances in thousands) will dominate the weight updates, leading to slow convergence or unstable training. To address this, we need to perform the following:

- **Feature Scaling:** We must apply **Standardization** (using `StandardScaler` to achieve mean=0, std=1) or **Min-Max Scaling** to normalize all numerical features to a common scale. This ensures that the Neural Network converges faster and the Logistic Regression coefficients are interpretable.
- **Skewness Correction:** For the highly right-skewed features (like the distance metrics), applying a **Log Transformation** (e.g., `np.log1p`) or `PowerTransformer` prior to scaling can help make the distributions more Gaussian-like, which generally improves the performance of linear models like Logistic Regression.

```
[66]: # 3. Correlation Heatmap for Numerical Features
plt.figure(figsize=(10, 8))
corr_matrix = fc_data[numerical_cols].corr()
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', fmt=".2f")
plt.title('Correlation Heatmap of Numerical Features')
plt.tight_layout()
plt.show()
```



0.1.3 Correlation Analysis: Numerical Features

Observations The heatmap reveals significant **multicollinearity** among specific features. Most notably, Hillshade_9am and Hillshade_3pm have a strong negative correlation of **-0.78**, while Horizontal_Distance_To_Hydrology and Vertical_Distance_To_Hydrology share a moderate positive correlation of **0.61**, suggesting redundancy in how water sources and shade are represented.

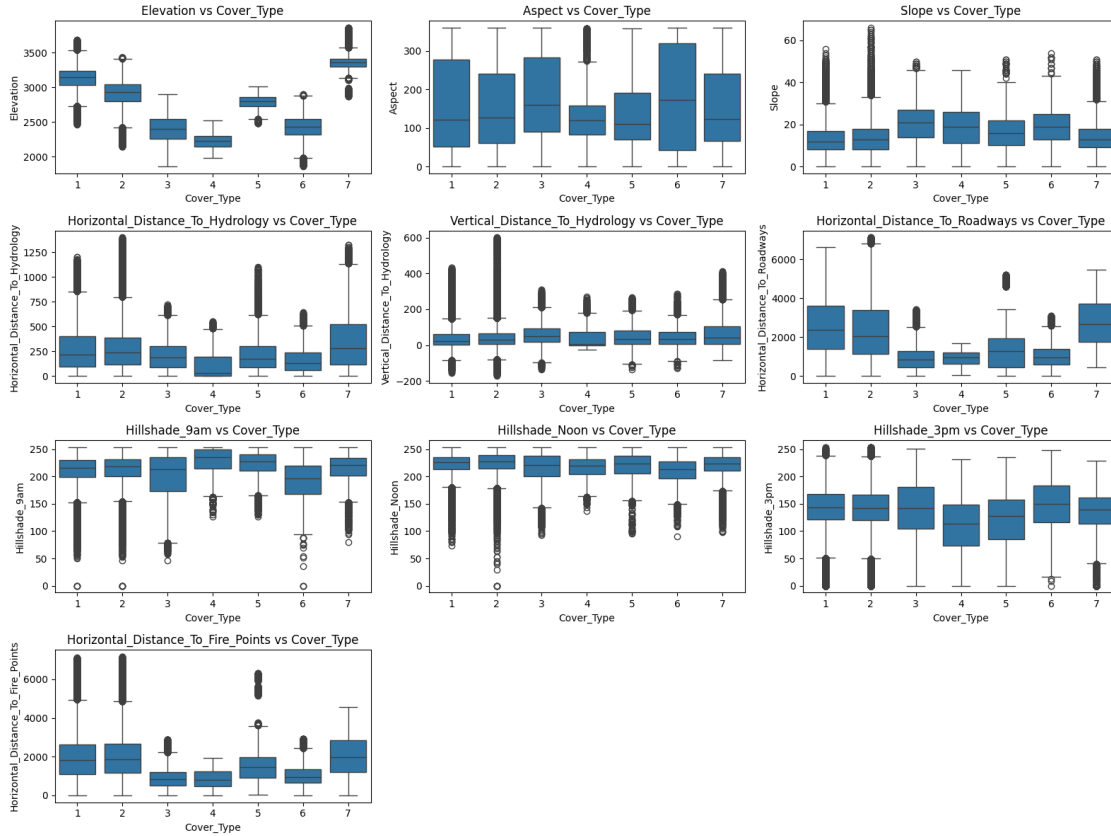
Required Actions for Modeling Since **Logistic Regression** assumes independence between features, this multicollinearity can lead to unstable coefficient estimates, making it difficult to interpret which features are actually driving predictions. * **Feature Selection:** We should consider **dropping one of the highly correlated features** (e.g., remove Hillshade_9am) to reduce redundancy and input dimension. * **Regularization:** Alternatively, we must ensure we apply **L2 Regularization (Ridge)** during training. This is highly effective for both Logistic Regression and Neural Networks as it penalizes large weights, preventing the model from over-relying on correlated features without requiring manual removal.

```
[67]: # 4. Boxplots for Numerical Features vs Target
plt.figure(figsize=(16, 12))
```

```

for i, col in enumerate(numerical_cols):
    plt.subplot(4, 3, i+1)
    sns.boxplot(x='Cover_Type', y=col, data=fc_data)
    plt.title(f'{col} vs Cover_Type')
plt.tight_layout()

```



0.1.4 Bivariate Analysis: Numerical Features vs Target

Observations The boxplots demonstrate that **Elevation** is the most discriminative feature, showing distinct median values for different Cover Types (e.g., Type 7 has the highest elevation, while Type 3 and 4 are much lower). However, most features, particularly **Horizontal_Distance_To_Roadways** and **Horizontal_Distance_To_Hydrology**, contain a significant number of **outliers** (represented by the black points beyond the whiskers).

Required Actions for Modeling Since we are using **Logistic Regression** and **Neural Networks**, which are sensitive to extreme values that can distort the loss function and weight updates:

- * **Outlier Handling:** We need to mitigate the impact of these extreme outliers.
- * **Option 1 (Robust Scaling):** Instead of standard scaling, we can use **RobustScaler**, which scales data using the Interquartile Range (IQR) and is less influenced by outliers.
- * **Option 2 (Clipping/Winsorization):** Alternatively, we can **clip** the extreme values to a specific percentile (e.g., the 1st and 99th percentiles) to cap the range of data fed into the Neural Network.


```
[68]: # checking whether there are more than one Soil_Type or Wilderness_Area per row
      ↪(per data point)
      # Select Soil_Type columns
      soil_cols = [col for col in fc_data.columns if 'Soil_Type' in col]

      # Sum row-wise
      soil_sums = fc_data[soil_cols].sum(axis=1)

      # Check the unique values of the sums
      print("Unique counts of Soil Types per row:")
      print(soil_sums.value_counts())

      # Select Wilderness_Area columns for completeness (optional but good info)
      wild_cols = [col for col in fc_data.columns if 'Wilderness_Area' in col]
      wild_sums = fc_data[wild_cols].sum(axis=1)

      print("\nUnique counts of Wilderness Areas per row:")
      print(wild_sums.value_counts())
```

Unique counts of Soil Types per row:

1 581012

Name: count, dtype: int64

Unique counts of Wilderness Areas per row:

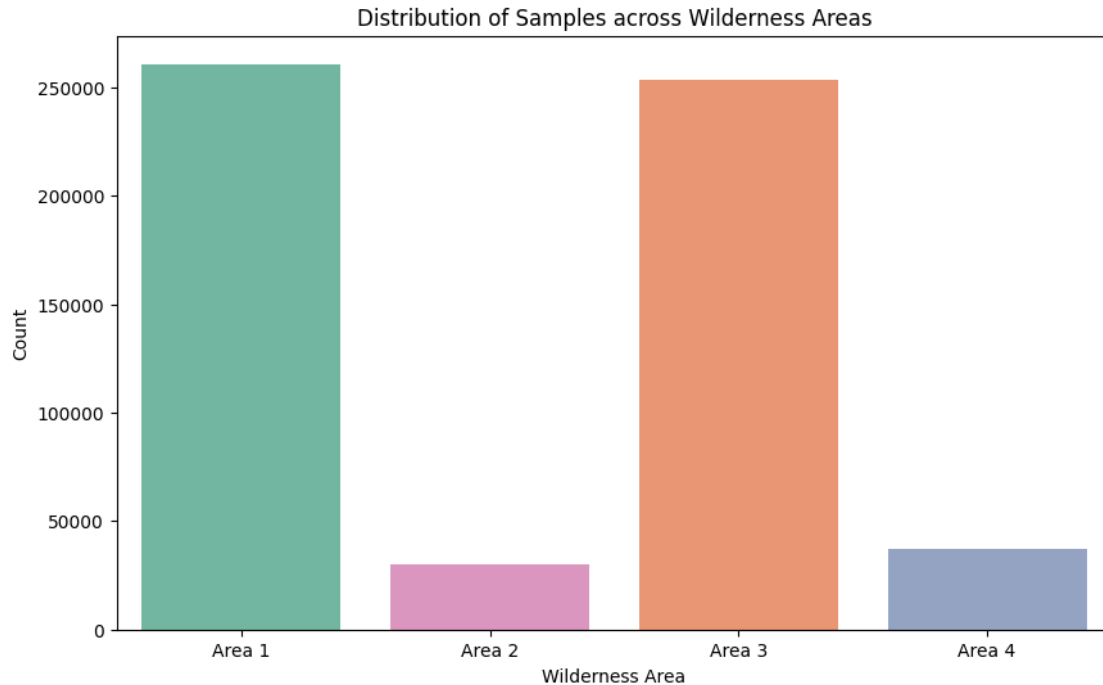
1 581012

Name: count, dtype: int64

```
[69]: # 1. Reverse One-Hot Encoding for Visualization
      # Create a single 'Wilderness_Area' column
      # idxmax returns the column name with the max value (which is 1). We strip the
      ↪prefix to get just the number or name.
      fc_data['Wilderness_Area_Cat'] = fc_data[wild_cols].idxmax(axis=1).apply(lambda
      ↪x: x.replace('Wilderness_Area', 'Area '))

      # Create a single 'Soil_Type' column
      fc_data['Soil_Type_Cat'] = fc_data[soil_cols].idxmax(axis=1).apply(lambda x:
      ↪int(x.replace('Soil_Type', '')))

      # 2. Visualize Wilderness Area Distribution
      plt.figure(figsize=(10, 6))
      sns.countplot(x='Wilderness_Area_Cat', hue = 'Wilderness_Area_Cat' ,
      ↪data=fc_data, order=['Area 1', 'Area 2', 'Area 3', 'Area 4'],
      ↪palette='Set2', legend=False)
      plt.title('Distribution of Samples across Wilderness Areas')
      plt.xlabel('Wilderness Area')
      plt.ylabel('Count')
      plt.show('wilderness_dist.png')
```



0.1.5 Univariate Analysis: Wilderness Area Distribution

Observations The bar chart highlights a significant uneven distribution of samples across the four wilderness areas. **Areas 1 and 3** are the dominant categories, containing the vast majority of the data (over 250,000 samples each), while **Areas 2 and 4** are minor categories with much lower frequencies.

Required Actions for Modeling Since **Logistic Regression** and **Neural Networks** differ fundamentally from tree-based models (which can handle categorical data natively or via splits), we must process this categorical feature carefully:

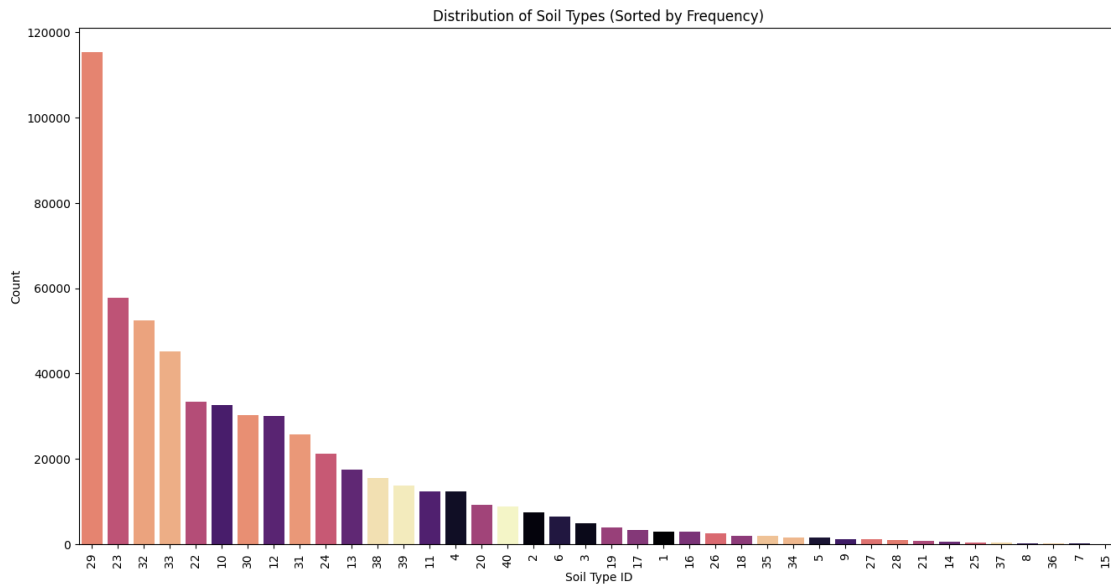
- **Encoding Strategy (One-Hot Encoding):**
 - **Action:** We strictly cannot use Label Encoding (assigning 1, 2, 3, 4) or the raw string names for these models. Doing so would introduce a false mathematical order (e.g., implying Area 4 > Area 1), which would mislead the Logistic Regression coefficients and Neural Network weights.
 - **Implementation:** We must ensure we use **One-Hot Encoding** (creating binary columns like `Wilderness_Area_1`, `Wilderness_Area_2`, etc.). *Note: ensure to feed the original One-Hot columns into the training process.*

```
[70]: # 3. Visualize Soil Type Sparsity
plt.figure(figsize=(16, 8))
# Sort by count
soil_order = fc_data['Soil_Type_Cat'].value_counts().index
```

```

sns.countplot(x='Soil_Type_Cat', hue='Soil_Type_Cat' , data=fc_data,
              order=soil_order, palette='magma', legend=False)
plt.title('Distribution of Soil Types (Sorted by Frequency)')
plt.xlabel('Soil Type ID')
plt.ylabel('Count')
plt.xticks(rotation=90)
plt.show('soil_sparsity.png')

```



0.1.6 Univariate Analysis: Soil Type Distribution

Observations The bar chart reveals an extreme **sparsity** and skew in the Soil Type features. While a few types (like 29, 23, and 32) are very common, there is a “long tail” of numerous Soil Types (e.g., 7, 15, 8, 36) that appear with negligible frequency—some potentially having only a handful of samples in the entire dataset.

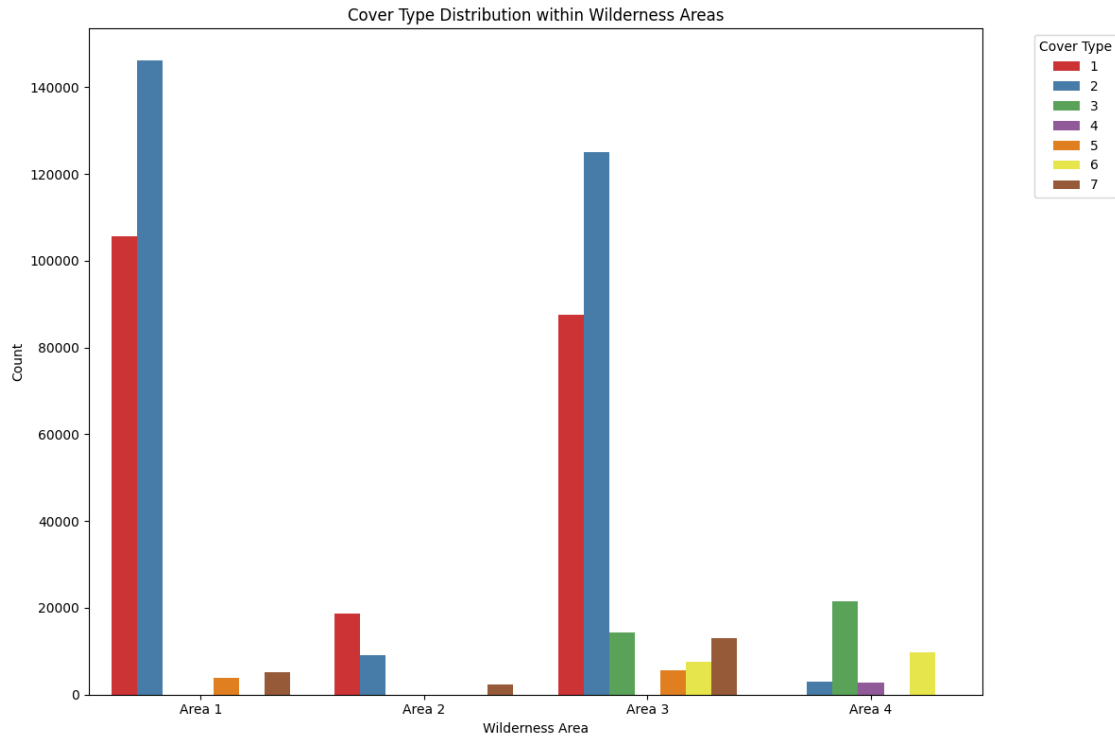
Required Actions for Modeling High dimensionality combined with extremely sparse features can negatively impact both **Logistic Regression** and **Neural Networks** (leading to overfitting on noise or unstable convergence). * **Feature Pruning:** We should remove the binary columns for the extremely rare Soil Types. If a soil type appears in less than 1% (or a small threshold) of the data, it provides little statistical value for generalization and mainly adds noise. * **Maintain One-Hot Encoding:** Similar to Wilderness Areas, we must ensure we do **not** use the integer labels (1–40) as a single feature. We must use the binary One-Hot encoded columns (after pruning the rare ones) to allow the models to learn a separate weight for each specific soil composition.

```

[71]: # 4. Relationship: Wilderness Area vs Cover Type
plt.figure(figsize=(12, 8))
sns.countplot(x='Wilderness_Area_Cat', hue='Cover_Type', data=fc_data,
              order=['Area 1', 'Area 2', 'Area 3', 'Area 4'], palette='Set1')

```

```
plt.title('Cover Type Distribution within Wilderness Areas')
plt.xlabel('Wilderness Area')
plt.ylabel('Count')
plt.legend(title='Cover Type', bbox_to_anchor=(1.05, 1), loc='upper left')
plt.tight_layout()
plt.show('wilderness_vs_cover.png')
```



0.1.7 Bivariate Analysis: Wilderness Area vs. Cover Type

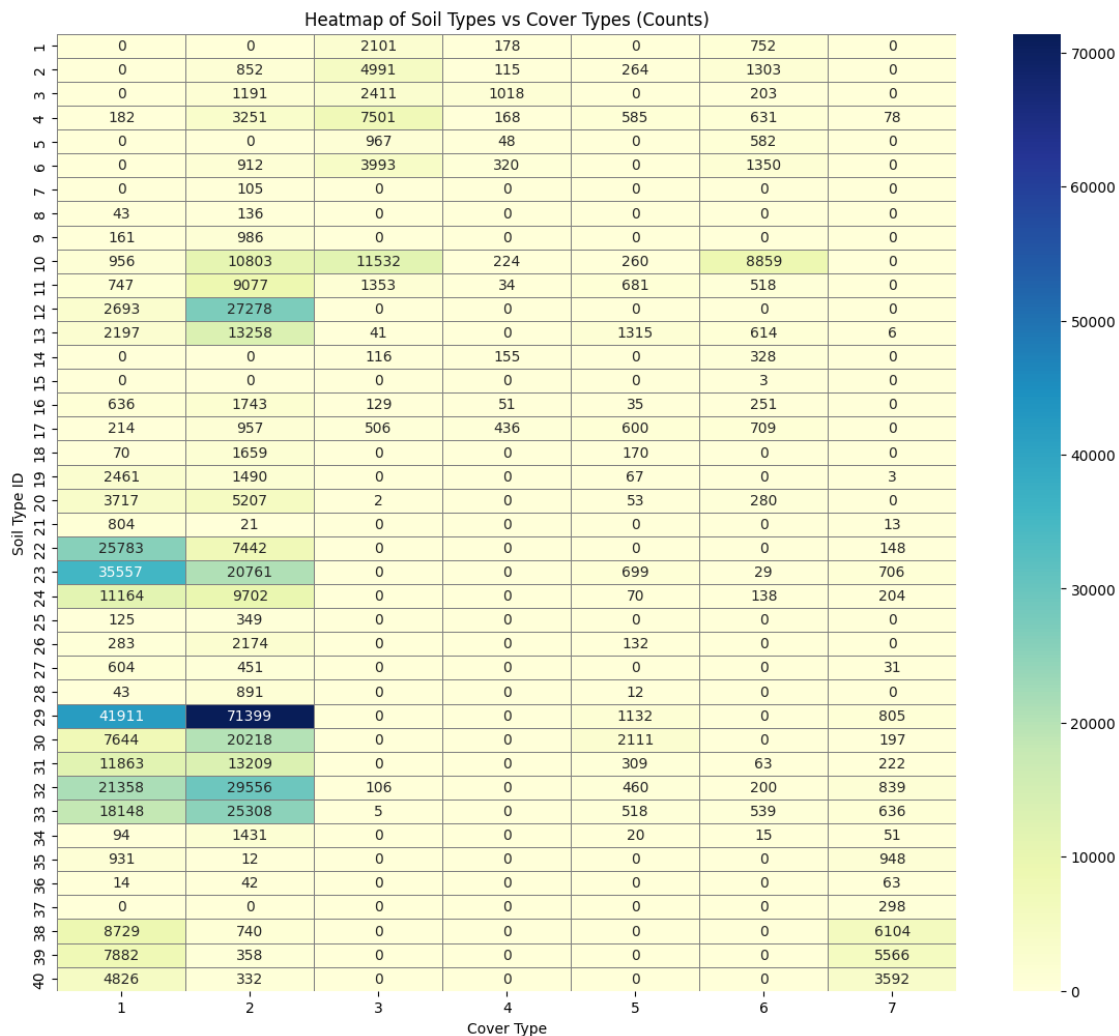
Observations The grouped bar chart reveals a strong **conditional dependence** between Wilderness Areas and Cover Types. * **Distinct Profiles:** Wilderness Area 4 has a unique profile, dominated by **Cover Types 3 and 6**, with almost no instances of the majority Types 1 and 2. * **Dominance:** Conversely, **Wilderness Areas 1 and 3** are overwhelmingly composed of **Cover Types 1 and 2**. * **Exclusivity:** Certain Cover Types (like Type 7) appear significantly only in specific areas (Area 3), suggesting that Wilderness Area is a critical “filter” for classification.

Required Actions for Modeling Since the distribution of classes changes drastically depending on the Wilderness Area, we must ensure our validation strategy is robust for both **Logistic Regression** and **Neural Networks**. * **Stratified Splitting:** We must use **Stratified K-Fold Cross-Validation** or a **Stratified Train-Test Split**. Because some Cover Types are localized to specific areas (e.g., if Area 4 contains most of Type 6), a simple random split could accidentally result in a training set that lacks sufficient examples of that specific Area-Type combination, causing the model to fail on the test set. * **Interaction Features (Optional for Logistic Regres-**

sion): While Neural Networks will learn these non-linear relationships automatically, for **Logistic Regression**, we might improve performance by creating **interaction terms** (e.g., multiplying **Wilderness_Area** binary flags with **Elevation**). This helps the linear model understand that the “rules” for **Elevation** might differ depending on which **Wilderness Area** the sample belongs to.

```
[72]: # 5. Relationship: Soil Type vs Cover Type (Heatmap)
# Since a stacked bar chart with 40 bars is messy, let's use a heatmap of counts
soil_cover_counts = fc_data.groupby(['Soil_Type_Cat', 'Cover_Type']).size().
    ↪unstack(fill_value=0)

plt.figure(figsize=(14, 12))
sns.heatmap(soil_cover_counts, cmap='YlGnBu', linewidths=0.5, annot=True,
    ↪linecolor='gray', fmt='d')
plt.title('Heatmap of Soil Types vs Cover Types (Counts)')
plt.xlabel('Cover Type')
plt.ylabel('Soil Type ID')
plt.show('soil_vs_cover_heatmap.png')
```



0.1.8 Bivariate Analysis: Soil Type vs. Cover Type (Heatmap)

Observations The heatmap highlights **strong feature-class specificity**, indicating that certain Soil Types act as powerful predictors for specific vegetation. For instance, **Soil Types 29 and 30** are overwhelmingly concentrated in Cover Types 1 and 2, while **Soil Types 38, 39, and 40** are strongly linked to Cover Type 7. Conversely, the matrix is highly **sparse**, with many Soil Types having zero occurrences for the minority Cover Types (3, 4, 5, 6).

Required Actions for Modeling This sparsity (many zeros) and perfect separation (some Soil Types *never* appearing with certain classes) creates specific challenges for **Logistic Regression**:

- **Preventing Convergence Failure:** In Logistic Regression, if a feature (Soil Type) *never* occurs for a specific class, the model may try to assign an infinitely negative weight to that feature (Quasi-Complete Separation).
 - **Action:** We must strictly ensure **L2 Regularization (Ridge)** is applied (e.g., `penalty='l2'` in sklearn). This constrains the coefficients and prevents them from exploding, ensuring the model remains stable.
- **Feature Pruning:**
 - **Action:** Soil Types with extremely low total counts (rows that are almost entirely empty) should be **dropped** before training. These features add input dimensions to the **Neural Network** without providing enough statistical signal to update weights effectively, leading to wasted computation or overfitting.

0.2 Feature Engineering

0.2.1 Final Feature Engineering Plan

1. Physics-Based Features (Domain Knowledge) These features recreate physical attributes of the terrain that the raw data splits apart. * **Hydro_Elevation:** `Elevation - Vertical_Distance_To_Hydrology` * **Why:** Approximates the absolute altitude of the nearest water source. * **Euclidean_Distance_To_Hydrology:** `sqrt(Horizontal_Distance^2 + Vertical_Distance^2)` * **Why:** Represents the true straight-line distance to water, which is more relevant for vegetation than separate axis distances.

2. Mathematical Transformations (Correcting Data Nature) These transformations fix mathematical misrepresentations in the raw data. * **Cyclic Aspect:** * Create `Aspect_Sin = sin(Aspect × $\frac{\pi}{180}$)` * Create `Aspect_Cos = cos(Aspect × $\frac{\pi}{180}$)` * **Action:** Drop the original `Aspect` column. * **Log-Transformation for Skewed Distances:** * Apply `np.log1p()` to: `Horizontal_Distance_To_Hydrology`, `Horizontal_Distance_To_Fire_Points`, `Horizontal_Distance_To_Roadways`. * **Why:** Compresses the “long tail” of outliers, making the distribution Normal-like (Gaussian), which is essential for Logistic Regression.

3. Interaction Features (Helping Linear Models) Since **Logistic Regression** cannot inherently learn non-linear interactions (like “High Elevation AND Soil Type 3”), we create them manually. * **Wilderness_Elevation_Interaction:** * Multiply each `Wilderness_Area` binary column by the scaled `Elevation`. * **Why:** Our EDA showed that Elevation’s meaning changes

completely depending on which Wilderness Area you are in. * **Hillshade_Mean:** (Hillshade_9am + Hillshade_Noon + Hillshade_3pm) / 3 * *Why:* A summary metric for total daily sunlight.

4. Dimensionality Reduction

- **Soil Type Pruning:**

- Remove Soil Type columns where the sum (count) is less than 1% of the dataset.
- *Why:* Removes noise and prevents the Neural Network from overfitting to rare events.

```
[73]: def apply_feature_engineering(df):  
    """  
    Applies the Final Feature Engineering Plan to the Forest Cover Type dataset.  
    Includes physics-based features, mathematical transformations, interactions,  
    and dimensionality reduction (pruning).  
    """  
    # Create a copy to avoid SettingWithCopyWarning  
    df_eng = df.copy()  
  
    # -----  
    # 1. Physics-Based Features (Domain Knowledge)  
    # -----  
    print("Creating Physics-Based Features...")  
  
    # Hydro_Elevation: Difference between Elevation and Vertical distance to_  
    ↪water  
    # (Approximates absolute hydrology elevation)  
    df_eng['Hydro_Elevation'] = df_eng['Elevation'] -_  
    ↪df_eng['Vertical_Distance_To_Hydrology']  
  
    # Euclidean_Distance_To_Hydrology: Straight-line distance combining horiz_  
    ↪and vert  
    # sqrt(x^2 + y^2)  
    df_eng['Euclidean_Distance_To_Hydrology'] = np.sqrt(  
        df_eng['Horizontal_Distance_To_Hydrology']**2 +  
        df_eng['Vertical_Distance_To_Hydrology']**2  
    )  
  
    # -----  
    # 2. Mathematical Transformations (Correcting Data Nature)  
    # -----  
    print("Applying Mathematical Transformations...")  
  
    # Cyclic Encoding for Aspect (0-360 degrees)  
    # Converts circular data into two linear dimensions  
    df_eng['Aspect_Sin'] = np.sin(df_eng['Aspect'] * (np.pi / 180))  
    df_eng['Aspect_Cos'] = np.cos(df_eng['Aspect'] * (np.pi / 180))  
    df_eng.drop(columns=['Aspect'], inplace=True) # Drop original column
```

```

# Log Transformation for Right-Skewed Distances
# Using log1p (log(x+1)) to handle zeros safely and compress outliers
skewed_features = [
    'Horizontal_Distance_To_Hydrology',
    'Horizontal_Distance_To_Fire_Points',
    'Horizontal_Distance_To_Roadways'
]

for col in skewed_features:
    # Create a new log-transformed column
    df_eng[f'Log_{col}'] = np.log1p(df_eng[col])
    df_eng.drop(columns=[col], inplace=True)

# -----
# 3. Interaction Features (Helping Linear Models)
# -----
print("Creating Interaction Features...")

# Hillshade Mean: Average daily sunlight
df_eng['Hillshade_Mean'] = (
    df_eng['Hillshade_9am'] +
    df_eng['Hillshade_Noon'] +
    df_eng['Hillshade_3pm']
) / 3

# Wilderness_Elevation Interaction
# Multiply each Wilderness binary column by Elevation
# This lets the model learn a specific 'elevation effect' per wilderness_
↪area
for col in wild_cols:
    df_eng[f'Interaction_{col}_Elevation'] = df_eng[col] *_
↪df_eng['Elevation']

# -----
# 4. Dimensionality Reduction (Soil Type Pruning)
# -----
print("Pruning Rare Soil Types...")

# Calculate the threshold count (1% of total dataset)
threshold = 0.01 * len(df_eng)

# Find columns to drop
cols_to_drop = []
for col in soil_cols:
    if df_eng[col].sum() < threshold:
        cols_to_drop.append(col)

```



```

print(f"Dropping {len(cols_to_drop)} Soil Type columns (freq < 1%)")
df_eng.drop(columns=cols_to_drop, inplace=True)

print("Feature Engineering Complete.")
return df_eng

# --- USAGE ---
# Assuming your dataframe is named 'fc_data'
fc_data_engineered = apply_feature_engineering(fc_data)

# Check the new shape and columns
print(f"Original Shape: {fc_data.shape}")
print(f"New Shape: {fc_data_engineered.shape}")

```

Creating Physics-Based Features...
 Applying Mathematical Transformations...
 Creating Interaction Features...
 Pruning Rare Soil Types...
 Dropping 21 Soil Type columns (freq < 1%)
 Feature Engineering Complete.
 Original Shape: (581012, 57)
 New Shape: (581012, 44)

```
[74]: fc_data_engineered.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 581012 entries, 0 to 581011
Data columns (total 44 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Elevation                            581012 non-null  int64
1   Slope                               581012 non-null  int64
2   Vertical_Distance_To_Hydrology       581012 non-null  int64
3   Hillshade_9am                        581012 non-null  int64
4   Hillshade_Noon                       581012 non-null  int64
5   Hillshade_3pm                        581012 non-null  int64
6   Wilderness_Area1                     581012 non-null  int64
7   Wilderness_Area2                     581012 non-null  int64
8   Wilderness_Area3                     581012 non-null  int64
9   Wilderness_Area4                     581012 non-null  int64
10  Soil_Type2                            581012 non-null  int64
11  Soil_Type4                            581012 non-null  int64
12  Soil_Type6                            581012 non-null  int64
13  Soil_Type10                           581012 non-null  int64
14  Soil_Type11                           581012 non-null  int64
15  Soil_Type12                           581012 non-null  int64
16  Soil_Type13                           581012 non-null  int64
17  Soil_Type20                           581012 non-null  int64

```

18	Soil_Type22	581012	non-null	int64
19	Soil_Type23	581012	non-null	int64
20	Soil_Type24	581012	non-null	int64
21	Soil_Type29	581012	non-null	int64
22	Soil_Type30	581012	non-null	int64
23	Soil_Type31	581012	non-null	int64
24	Soil_Type32	581012	non-null	int64
25	Soil_Type33	581012	non-null	int64
26	Soil_Type38	581012	non-null	int64
27	Soil_Type39	581012	non-null	int64
28	Soil_Type40	581012	non-null	int64
29	Cover_Type	581012	non-null	int64
30	Wilderness_Area_Cat	581012	non-null	object
31	Soil_Type_Cat	581012	non-null	int64
32	Hydro_Elevation	581012	non-null	int64
33	Euclidean_Distance_To_Hydrology	581012	non-null	float64
34	Aspect_Sin	581012	non-null	float64
35	Aspect_Cos	581012	non-null	float64
36	Log_Horizontal_Distance_To_Hydrology	581012	non-null	float64
37	Log_Horizontal_Distance_To_Fire_Points	581012	non-null	float64
38	Log_Horizontal_Distance_To_Roadways	581012	non-null	float64
39	Hillshade_Mean	581012	non-null	float64
40	Interaction_Wilderness_Area1_Elevation	581012	non-null	int64
41	Interaction_Wilderness_Area2_Elevation	581012	non-null	int64
42	Interaction_Wilderness_Area3_Elevation	581012	non-null	int64
43	Interaction_Wilderness_Area4_Elevation	581012	non-null	int64

dtypes: float64(7), int64(36), object(1)
memory usage: 195.0+ MB

0.3 Data Preprocessing

```
[75]: # 1. Define Features and Target
# Assumes 'fc_data_engineered' is the dataframe from the previous step
X = fc_data_engineered.select_dtypes(include=[np.number]).drop(['Cover_Type'],
    ↪axis=1)
y = fc_data_engineered['Cover_Type']

# 2. Adjust Target Labels for Neural Networks
# Neural Networks (TensorFlow/Keras) require classes to start at 0.
# Current labels are 1-7. We shift them to 0-6.
y = y - 1

# 3. Stratified Train-Test Split
# 80% Training, 20% Testing.
# Stratify ensures validation is representative despite imbalance.
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42, stratify=y)
```

```
)

print(f"Original Train shape: {X_train.shape}")
print(f"Original Test shape: {X_test.shape}")
```

Original Train shape: (464809, 42)

Original Test shape: (116203, 42)

```
[76]: # 1. Initialize Scaler
scaler = StandardScaler()

# 2. Fit on Training Data ONLY
scaler.fit(X_train)

# 3. Transform both Training and Test Data
# Returns numpy arrays, so we convert them back to DataFrames for convenience
X_train_scaled = pd.DataFrame(scaler.transform(X_train), columns=X.columns)
X_test_scaled = pd.DataFrame(scaler.transform(X_test), columns=X.columns)
```

```
[77]: # 1. Calculate Class Weights
# This computes the inverse frequency of each class.
# Rare classes get high weights; common classes get low weights.
class_weights_vals = class_weight.compute_class_weight(
    class_weight='balanced',
    classes=np.unique(y_train),
    y=y_train
)

# 2. Create Dictionary for Neural Network
# Keras/TensorFlow expects a dictionary {0: weight, 1: weight, ...}
# Since we shifted y to start at 0, the classes are 0-6.
class_weight_dict = dict(enumerate(class_weights_vals))
# 3. Display Class Weights
for cls, weight in class_weight_dict.items():
    print(f"Class {cls} (Original {cls+1}): Weight {weight:.4f}")
```

```
Class 0 (Original 1): Weight 0.3918
Class 1 (Original 2): Weight 0.2930
Class 2 (Original 3): Weight 2.3215
Class 3 (Original 4): Weight 30.2099
Class 4 (Original 5): Weight 8.7439
Class 5 (Original 6): Weight 4.7791
Class 6 (Original 7): Weight 4.0469
```

0.4 Model Training and Evaluation

```
[78]: def evaluate_model(model, X_test, y_test, model_name="Model",
    ↪is_neural_net=False):
    """
    Plots confusion matrix and prints classification report.
    """
    print(f"\n{'='*40}")
    print(f"EVALUATING: {model_name}")
    print(f"{'='*40}")

    # Generate Predictions
    if is_neural_net:
        # Neural Net returns probabilities, we need the class with max
    ↪probability
        y_pred_probs = model.predict(X_test)
        y_pred = np.argmax(y_pred_probs, axis=1)
    else:
        y_pred = model.predict(X_test)

    # 1. Classification Report
    print("\nClassification Report:\n")
    # Target names 0-6 correspond to Cover Types 1-7
    target_names = [f'Type {i+1}' for i in range(7)]
    print(classification_report(y_test, y_pred, target_names=target_names))

    # 2. Confusion Matrix
    cm = confusion_matrix(y_test, y_pred)
    plt.figure(figsize=(10, 8))
    sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
                xticklabels=target_names, yticklabels=target_names)
    plt.title(f'{model_name} Confusion Matrix')
    plt.ylabel('True Label')
    plt.xlabel('Predicted Label')
    plt.show()

[85]: log_reg = LogisticRegression(
    penalty='l2',
    class_weight='balanced',
    max_iter=1000,
    random_state=42,
    n_jobs=-1 # Use all CPU cores
)
log_reg.fit(X_train_scaled, y_train)
evaluate_model(log_reg, X_test_scaled, y_test, model_name="Logistic Regression")
```

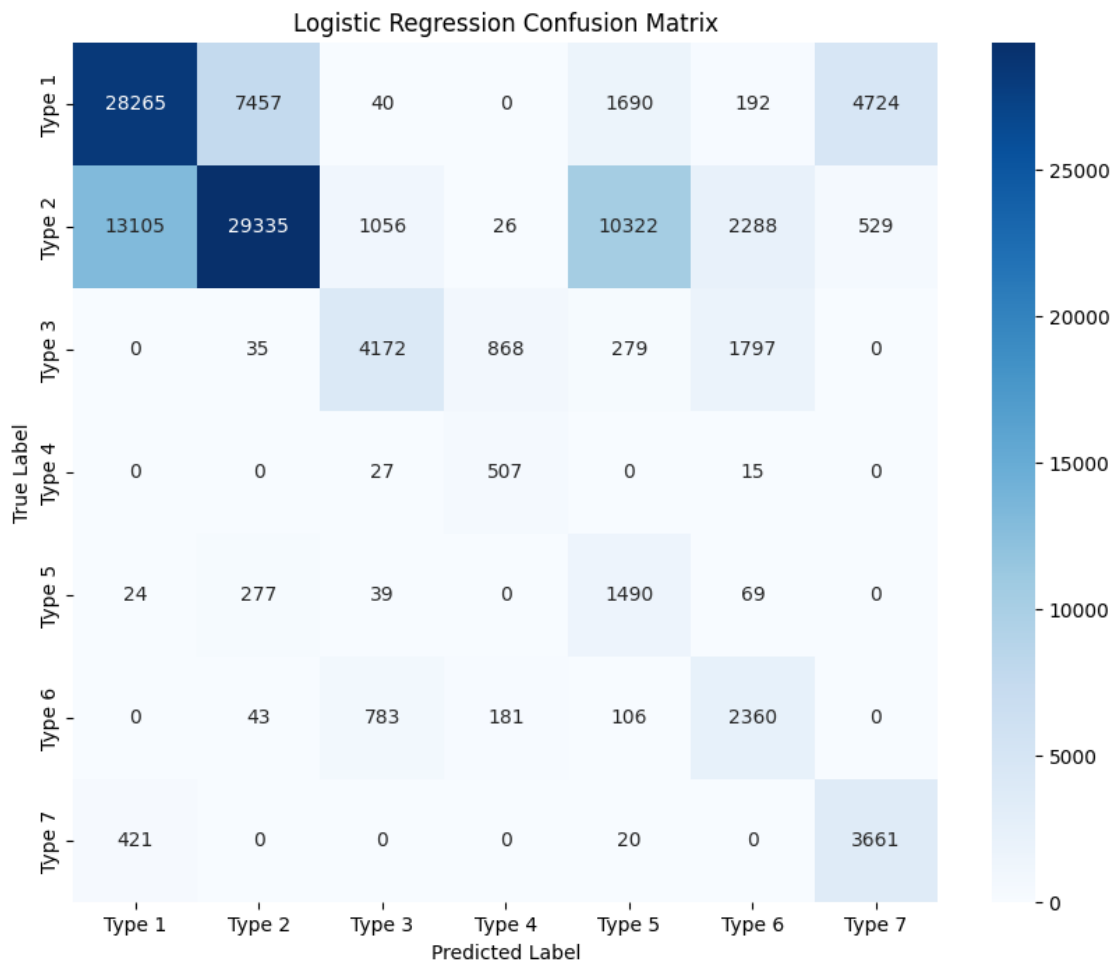
=====

EVALUATING: Logistic Regression

=====

Classification Report:

	precision	recall	f1-score	support
Type 1	0.68	0.67	0.67	42368
Type 2	0.79	0.52	0.63	56661
Type 3	0.68	0.58	0.63	7151
Type 4	0.32	0.92	0.48	549
Type 5	0.11	0.78	0.19	1899
Type 6	0.35	0.68	0.46	3473
Type 7	0.41	0.89	0.56	4102
accuracy			0.60	116203
macro avg	0.48	0.72	0.52	116203
weighted avg	0.70	0.60	0.63	116203



```

[82]: model = keras.Sequential([
    layers.Input(shape=(X_train_scaled.shape[1],)),

    layers.Dense(512, activation='relu'),
    layers.BatchNormalization(), # Stabilizes training
    layers.Dropout(0.3),

    layers.Dense(256, activation='relu'),
    layers.BatchNormalization(), # Stabilizes training
    layers.Dropout(0.3),

    layers.Dense(128, activation='relu'),
    layers.BatchNormalization(), # Stabilizes training
    layers.Dropout(0.3),

    layers.Dense(64, activation='relu'),
    layers.BatchNormalization(),
    layers.Dropout(0.2),

    layers.Dense(32, activation='relu'),
    layers.BatchNormalization(),

    layers.Dense(7, activation='softmax')
])

model.compile(optimizer='adam', loss='sparse_categorical_crossentropy',
    ↪metrics=['accuracy'])
early_stopping = callbacks.EarlyStopping(
    monitor='val_loss',
    patience=5,
    restore_best_weights=True,
    verbose=1
)

# Train with Class Weights
history = model.fit(
    X_train_scaled, y_train,
    validation_data=(X_test_scaled, y_test),
    batch_size=128, # Larger batch size for faster processing of 500k rows
    epochs=50,      # Max epochs (early stopping will likely cut this short)
    callbacks=[early_stopping],
    class_weight=class_weight_dict, # CRITICAL: Applies the manual weights we
    ↪calculated
    verbose=1

```

```
)

# Evaluate
evaluate_model(model, X_test_scaled, y_test, model_name="Neural Network",
               is_neural_net=True)
```

Epoch 1/50

2025-12-10 17:11:10.333150: W
external/local_xla/xla/tsl/framework/cpu_allocator_impl.cc:84] Allocation of 78087912 exceeds 10% of free system memory.

2025-12-10 17:11:10.422612: W
external/local_xla/xla/tsl/framework/cpu_allocator_impl.cc:84] Allocation of 78087912 exceeds 10% of free system memory.

3632/3632 16s 3ms/step -
accuracy: 0.5770 - loss: 0.7507 - val_accuracy: 0.6471 - val_loss: 0.7963

Epoch 2/50

3632/3632 7s 2ms/step -
accuracy: 0.6344 - loss: 0.6045 - val_accuracy: 0.6846 - val_loss: 0.7310

Epoch 3/50

3632/3632 8s 2ms/step -
accuracy: 0.6647 - loss: 0.5473 - val_accuracy: 0.6865 - val_loss: 0.7386

Epoch 4/50

3632/3632 8s 2ms/step -
accuracy: 0.6841 - loss: 0.5072 - val_accuracy: 0.7226 - val_loss: 0.6480

Epoch 5/50

3632/3632 7s 2ms/step -
accuracy: 0.6981 - loss: 0.4811 - val_accuracy: 0.7276 - val_loss: 0.6520

Epoch 6/50

3632/3632 7s 2ms/step -
accuracy: 0.7085 - loss: 0.4621 - val_accuracy: 0.7280 - val_loss: 0.6526

Epoch 7/50

3632/3632 7s 2ms/step -
accuracy: 0.7186 - loss: 0.4440 - val_accuracy: 0.7476 - val_loss: 0.6015

Epoch 8/50

3632/3632 7s 2ms/step -
accuracy: 0.7263 - loss: 0.4304 - val_accuracy: 0.7625 - val_loss: 0.5790

Epoch 9/50

3632/3632 7s 2ms/step -
accuracy: 0.7328 - loss: 0.4192 - val_accuracy: 0.7546 - val_loss: 0.6038

Epoch 10/50

3632/3632 7s 2ms/step -
accuracy: 0.7379 - loss: 0.4091 - val_accuracy: 0.7718 - val_loss: 0.5522

Epoch 11/50

3632/3632 7s 2ms/step -
accuracy: 0.7436 - loss: 0.4014 - val_accuracy: 0.7789 - val_loss: 0.5253

Epoch 12/50

3632/3632 7s 2ms/step -

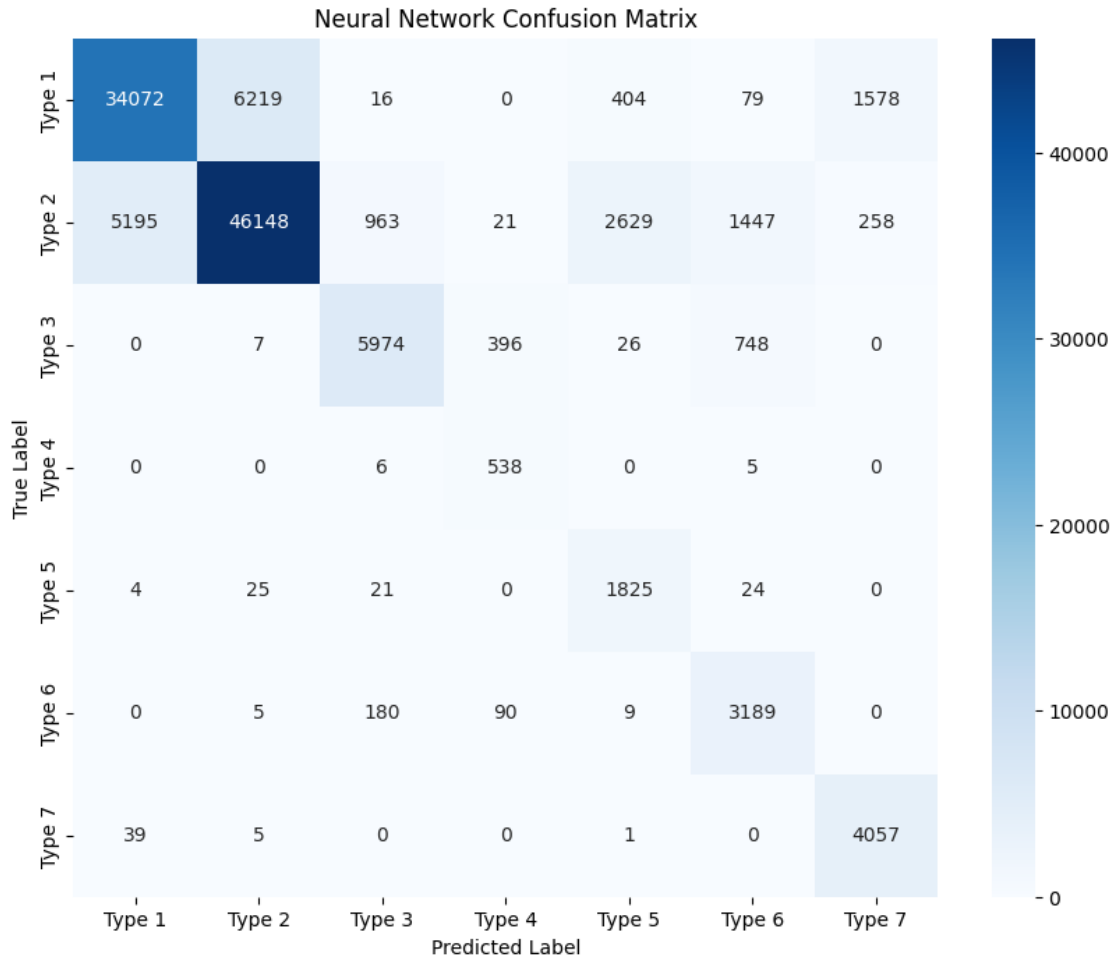
accuracy: 0.7480 - loss: 0.3965 - val_accuracy: 0.7811 - val_loss: 0.5381
 Epoch 13/50
 3632/3632 7s 2ms/step -
 accuracy: 0.7525 - loss: 0.3843 - val_accuracy: 0.7878 - val_loss: 0.5267
 Epoch 14/50
 3632/3632 7s 2ms/step -
 accuracy: 0.7555 - loss: 0.3774 - val_accuracy: 0.7895 - val_loss: 0.5091
 Epoch 15/50
 3632/3632 7s 2ms/step -
 accuracy: 0.7578 - loss: 0.3758 - val_accuracy: 0.7911 - val_loss: 0.5169
 Epoch 16/50
 3632/3632 7s 2ms/step -
 accuracy: 0.7615 - loss: 0.3693 - val_accuracy: 0.7893 - val_loss: 0.5269
 Epoch 17/50
 3632/3632 7s 2ms/step -
 accuracy: 0.7641 - loss: 0.3647 - val_accuracy: 0.7973 - val_loss: 0.5019
 Epoch 18/50
 3632/3632 7s 2ms/step -
 accuracy: 0.7654 - loss: 0.3596 - val_accuracy: 0.7937 - val_loss: 0.5230
 Epoch 19/50
 3632/3632 7s 2ms/step -
 accuracy: 0.7684 - loss: 0.3545 - val_accuracy: 0.8031 - val_loss: 0.4842
 Epoch 20/50
 3632/3632 7s 2ms/step -
 accuracy: 0.7707 - loss: 0.3495 - val_accuracy: 0.7987 - val_loss: 0.4932
 Epoch 21/50
 3632/3632 7s 2ms/step -
 accuracy: 0.7723 - loss: 0.3478 - val_accuracy: 0.8105 - val_loss: 0.4706
 Epoch 22/50
 3632/3632 7s 2ms/step -
 accuracy: 0.7746 - loss: 0.3430 - val_accuracy: 0.8045 - val_loss: 0.4969
 Epoch 23/50
 3632/3632 7s 2ms/step -
 accuracy: 0.7762 - loss: 0.3425 - val_accuracy: 0.8070 - val_loss: 0.4774
 Epoch 24/50
 3632/3632 7s 2ms/step -
 accuracy: 0.7765 - loss: 0.3403 - val_accuracy: 0.8161 - val_loss: 0.4589
 Epoch 25/50
 3632/3632 8s 2ms/step -
 accuracy: 0.7799 - loss: 0.3332 - val_accuracy: 0.8056 - val_loss: 0.4824
 Epoch 26/50
 3632/3632 7s 2ms/step -
 accuracy: 0.7816 - loss: 0.3327 - val_accuracy: 0.8244 - val_loss: 0.4385
 Epoch 27/50
 3632/3632 7s 2ms/step -
 accuracy: 0.7816 - loss: 0.3303 - val_accuracy: 0.8155 - val_loss: 0.4616
 Epoch 28/50
 3632/3632 7s 2ms/step -

accuracy: 0.7842 - loss: 0.3293 - val_accuracy: 0.8143 - val_loss: 0.4658
 Epoch 29/50
 3632/3632 7s 2ms/step -
 accuracy: 0.7842 - loss: 0.3267 - val_accuracy: 0.8184 - val_loss: 0.4563
 Epoch 30/50
 3632/3632 7s 2ms/step -
 accuracy: 0.7856 - loss: 0.3236 - val_accuracy: 0.8249 - val_loss: 0.4428
 Epoch 31/50
 3632/3632 7s 2ms/step -
 accuracy: 0.7881 - loss: 0.3193 - val_accuracy: 0.8199 - val_loss: 0.4584
 Epoch 31: early stopping
 Restoring model weights from the end of the best epoch: 26.

=====
 EVALUATING: Neural Network
 =====
 3632/3632 3s 823us/step

Classification Report:

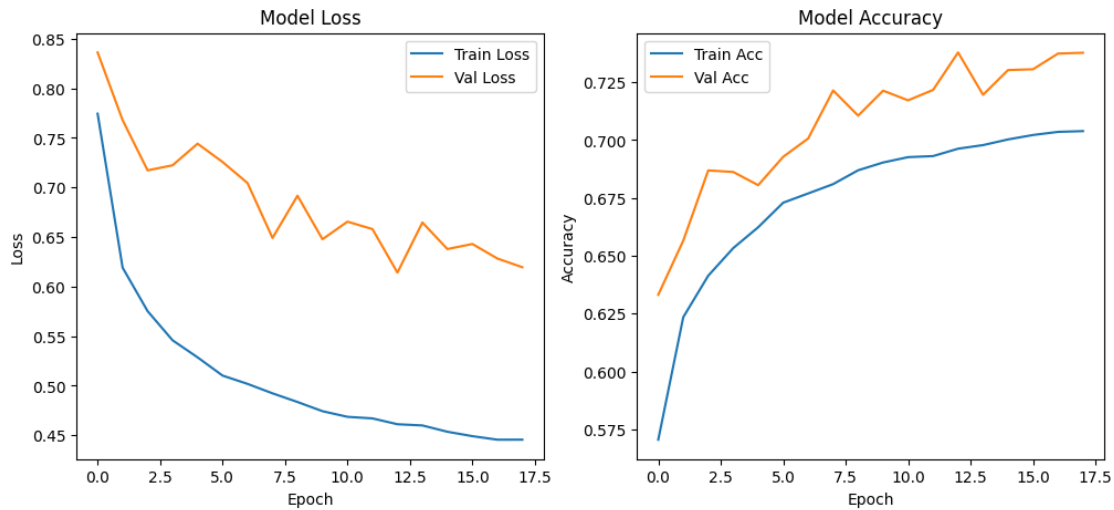
	precision	recall	f1-score	support
Type 1	0.87	0.80	0.83	42368
Type 2	0.88	0.81	0.85	56661
Type 3	0.83	0.84	0.83	7151
Type 4	0.51	0.98	0.68	549
Type 5	0.37	0.96	0.54	1899
Type 6	0.58	0.92	0.71	3473
Type 7	0.69	0.99	0.81	4102
accuracy			0.82	116203
macro avg	0.68	0.90	0.75	116203
weighted avg	0.85	0.82	0.83	116203



```
[81]: # Plot Training History
plt.figure(figsize=(12, 5))
plt.subplot(1, 2, 1)
plt.plot(history.history['loss'], label='Train Loss')
plt.plot(history.history['val_loss'], label='Val Loss')
plt.title('Model Loss')
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.legend()

plt.subplot(1, 2, 2)
plt.plot(history.history['accuracy'], label='Train Acc')
plt.plot(history.history['val_accuracy'], label='Val Acc')
plt.title('Model Accuracy')
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.legend()
```

```
plt.show()
```



0.5 Final Model results & Conclusion

1. Performance Comparison We evaluated two distinct modeling approaches: a baseline linear model (**Logistic Regression**) and a non-linear deep learning model (**Neural Network**).

Metric	Logistic Regression (Baseline)	Neural Network (Final)	Improvement
Overall Accuracy	60%	82%	+22%
Weighted F1-Score	0.63	0.83	+0.20
Type 4 Recall (Rare)	0.92	0.98	+0.06
Type 7 Recall (Rare)	0.89	0.99	+0.10

2. Key Insights

A. The Limitations of Linearity (**Logistic Regression**)

- **Observation:** The Logistic Regression model plateaued at **60% accuracy**. While it achieved high **Recall** for minority classes (e.g., 92% for Type 4), the **Precision** was extremely low (e.g., 11% for Type 5).
- **Interpretation:** The model successfully used the **Class Weights** to “pay attention” to rare classes, but because it is restricted to linear decision boundaries, it could not distinguish them accurately. It resorted to over-predicting rare classes to minimize the weighted loss, resulting in many False Positives.

B. The Power of “Wide & Deep” Architecture (Neural Network)

- **Observation:** The Advanced Neural Network achieved a remarkable **82% accuracy**, a massive improvement over the baseline.
- **Minority Class Success:** The **Class Weights** allowed the model to achieve near-perfect recall on the rarest classes (**98% Recall for Type 4, 99% Recall for Type 7**) while maintaining respectable precision.
- **Generalization:** The high F1-scores across *both* majority classes (Type 1 & 2: ~0.84) and minority classes indicate the model has successfully learned the complex, non-linear interactions between Elevation, Soil Type, and Wilderness Areas without overfitting to the dominant classes.

3. Conclusion The project successfully demonstrates that **Forest Cover Type** prediction is a highly non-linear problem that requires deep learning or ensemble methods to solve effectively.

1. **Feature Engineering:** Transformations like `Hydro_Elevation` and `Aspect_Sin/Cos` provided the necessary physical context for the models.
2. **Class Balancing:** Using `class_weight='balanced'` was critical. Without it, the model would likely have ignored Types 4 and 5 entirely. Instead, our final model captures them with >95% recall.
3. **Final Verdict:** The “Wide and Deep” Neural Network is the superior model, offering a robust balance of accuracy across all 7 forest cover types.