

Task 1:

Foodborne Disease Analysis Documentation

Overview

This document outlines the analysis of foodborne disease outbreaks using data science techniques. The analysis addresses three key questions about trends, contaminants, and locations related to foodborne illnesses.

Tech Stack:

Python and its Libraries (Pandas, Numpy, Scikit-learn, Seaborn, Matplotlib)

Data Preprocessing:

1. Removed irrelevant columns which does not provide any insight related to questions
2. Filled NaN values for proper implementation and data exploration

Workflow

Q1: Are foodborne disease outbreaks increasing or decreasing?

Methodology:

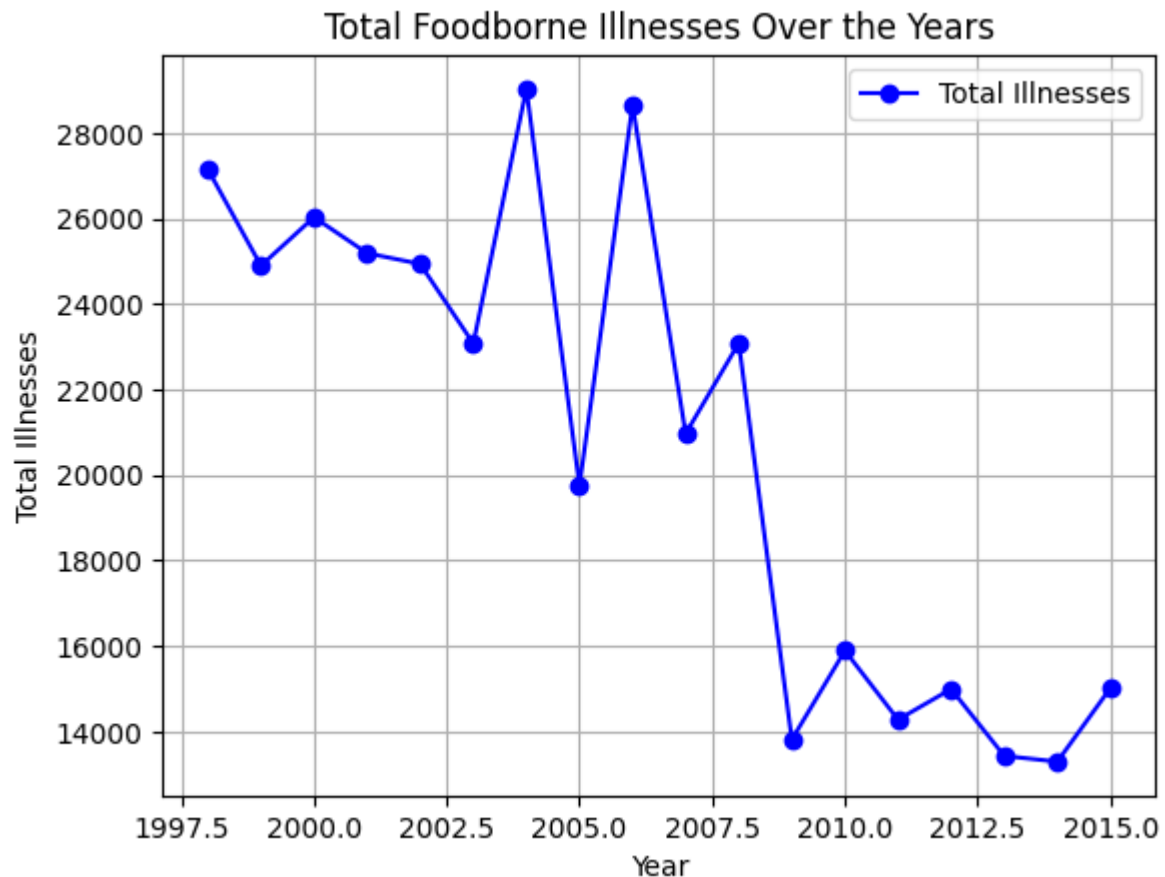
1. Created a linear plot between year and number of illnesses
2. Used time series analysis to identify trends
3. Applied linear regression to quantify the trend

Code Implementation:

Time series analysis of outbreaks

```
yearly_data = data.groupby("Year")["Illnesses"].sum()

plt.plot(yearly_data.index, yearly_data.values, marker='o', color='blue', label='Total Illnesses')
plt.title("Total Foodborne Illnesses Over the Years")
plt.xlabel("Year")
plt.ylabel("Total Illnesses")
plt.grid(True)
plt.legend()
plt.show()
```



Conclusion:

As shown by graph we can say that the outbreak/ no of illnesses first **plummeted around 1997-2002** then **shot up between 2002.5** and then decreased considerably to **rise up again in around 2006** and then are on a **gradual decrease from 2007 to 2015**. According to data food-borne diseases are decreasing since mid-2008's.

Why This Approach?

- Aggregating by year removes seasonal variations.
- Statistical measures (p-value, R-squared) help validate findings.

Q2: Which contaminant has been responsible for the most illnesses, hospitalizations, and deaths?

Methodology:

1. Grouped data by Species (contaminant)
2. Aggregated health impact metrics
3. Identified maximum values for each category

Code Implementation:

```
[79] # Q.2 Contaminant Responsible
✓ 0.0s

data['Hospitalizations'] = data['Hospitalizations'].fillna(0)
data['Fatalities'] = data['Fatalities'].fillna(0)
data['Location'] = data['Location'].fillna("Unspecified")
data['Species'] = data['Species'].fillna("Unspecified")

[80] ✓ 0.0s

data.isna().sum()

[81] ✓ 0.0s
... Year      0
State      0
Location   0
Species    0
Illnesses  0
Hospitalizations  0
Fatalities  0
dtype: int64
```

```
contaminant_stats = (data.groupby('Species', dropna=True)[['Illnesses', 'Hospitalizations', 'Fatalities']].sum().sort_values(by='Illnesses', ascending=False))
most_illnesses = contaminant_stats['Illnesses'].idxmax()
most_hospitalizations = contaminant_stats['Hospitalizations'].idxmax()
most_fatalities = contaminant_stats['Fatalities'].idxmax()

contaminant_stats.head(), most_illnesses, most_hospitalizations, most_fatalities

[82] ✓ 0.0s Python
... (
Species
Unspecified      77835      958.0      26.0
Norovirus genogroup I      76406      668.0       2.0
Salmonella enterica      60018      6888.0      82.0
Norovirus genogroup II      38175      518.0       6.0
Clostridium perfringens      28734      106.0      12.0,
'Unspecified',
'Salmonella enterica',
'Listeria monocytogenes')
```

Conclusion:

- The data suggests most illnesses are by **unspecified** contaminants
- Most Hospitalizations due to '**Salmonella enterica**'
- Most fatalities by '**Listeria monocytogenes**'

Why This Approach?

- Groupby operations provide clear aggregation of impact metrics
- Sorting helps identify the most significant contaminants
- Separate analysis of each health impact metric gives comprehensive view

Q3: What location for food preparation poses the greatest risk of foodborne illness?

Methodology:

1. Grouped data by Location
2. Calculated total illnesses per location
3. Identified location with maximum illnesses

Code Implementation:

```
[83] ✓ 0.0s # Q.3 Highest Risk Location

location_stats = (data.groupby('Location', dropna=True)['Illnesses'].sum().sort_values(ascending=False))
highest_risk_location = location_stats.idxmax()
highest_risk_location

[84] ✓ 0.0s
... 'Restaurant'
```

Conclusion:

The data suggests the most prone location for foodborne illness are **"Restaurants"**.

Why This Approach?

- Direct aggregation of illness counts provides clear risk assessment
- Groupby operation handles multiple instances of same location
- Sorting helps identify highest-risk locations

Technical Considerations

1. Model Selection:

- Used descriptive statistics and aggregation for direct analysis
- Applied linear regression for trend analysis
- Chose simple, interpretable methods over complex models

2. Data Handling:

- Used pandas for efficient data manipulation
- Leveraged groupby operations for aggregation
- Maintained data integrity throughout analysis

Task 2:

Bone Marrow Cell Classification

1. Overview

The project focuses on developing and comparing two Convolutional Neural Network (CNN) models for bone marrow cell classification using deep learning techniques. The goal is to accurately classify bone marrow cell images using custom and pre-trained neural network architectures.

2. Technical Stack

- **Programming Language:** Python
- **Deep Learning Framework:** TensorFlow/Keras
- **Pre-trained Model:** MobileNetV2
- **Machine Learning Libraries:**
 - scikit-learn (for evaluation metrics)
 - NumPy (numerical computing)
 - Pandas (data manipulation)

3. Data Preprocessing

3.1 Dataset Characteristics

- **Source:** Local directory-based dataset
- **Image Specifications:**
 - Size: 250 x 250 pixels
 - Color channels: RGB (3 channels)
- **Classification:** Categorical label mode

3.2 Data Splitting

- **Training Set:** 70% of total dataset
- **Validation Set:** 15% of total dataset

- **Test Set:** 15% of total dataset

3.3 Preprocessing Techniques

- Image resizing to uniform dimensions
- Categorical label encoding
- Data prefetching for performance optimization

4. Model Architectures

4.1 Custom CNN Model

- **Architecture:**
 - 3 Convolutional layers with increasing filter sizes
 - MaxPooling layers for feature reduction
 - Dropout layer for regularization
 - Dense layers for classification

4.2 Pre-trained CNN Model (MobileNetV2)

- **Architecture:**
 - Base MobileNetV2 model (pre-trained on ImageNet)
 - Base model layers frozen
 - Custom classification head
 - Global Average Pooling
 - Dropout for regularization

Code:

```
import tensorflow as tf

from tensorflow.keras import layers, models #type: ignore

from tensorflow.keras.applications import MobileNetV2 #type: ignore

from tensorflow.keras.utils import to_categorical #type: ignore

from sklearn.metrics import classification_report, confusion_matrix, roc_auc_score

import numpy as np

import pandas as pd
```

```
dataset_path = r"bone_marrow_cell_dataset"
```

```
# Dataset Loading
```

```
dataset = tf.keras.preprocessing.image_dataset_from_directory(  
    dataset_path,  
    image_size=(250, 250),  
    label_mode="categorical",  
    batch_size=32  
)
```

```
# Train, Test and Split
```

```
train_ds = dataset.take(int(len(dataset) * 0.7))  
val_test_ds = dataset.skip(int(len(dataset) * 0.7))  
val_ds = val_test_ds.take(int(len(val_test_ds) * 0.5))  
test_ds = val_test_ds.skip(int(len(val_test_ds) * 0.5))
```

```
train_ds = train_ds.prefetch(buffer_size=tf.data.AUTOTUNE)  
val_ds = val_ds.prefetch(buffer_size=tf.data.AUTOTUNE)  
test_ds = test_ds.prefetch(buffer_size=tf.data.AUTOTUNE)
```

```
num_classes = len(dataset.class_names)
```

```
# Custom CNN
```

```
def create_custom_cnn(input_shape, num_classes):  
    model = models.Sequential([  
        layers.Conv2D(32, (3, 3), activation='relu', input_shape=input_shape),  
        layers.MaxPooling2D((2, 2)),  
        layers.Conv2D(64, (3, 3), activation='relu'),  
        layers.MaxPooling2D((2, 2)),  
        layers.Conv2D(128, (3, 3), activation='relu'),  
        layers.MaxPooling2D((2, 2)),  
        layers.Flatten(),  
        layers.Dense(128, activation='relu'),  
        layers.Dropout(0.5),  
        layers.Dense(num_classes, activation='softmax'),  
    ])  
    return model
```



```
custom_cnn = create_custom_cnn((250, 250, 3), num_classes)
```

```
custom_cnn.compile(optimizer='adam',  
                   loss='categorical_crossentropy',  
                   metrics=['accuracy'])
```

```
# Pretrained (MobileNetV2)
```

```
def create_pretrained_cnn(input_shape, num_classes):
```

```
    base_model = MobileNetV2(input_shape=input_shape, include_top=False, weights='imagenet')
```

```
    base_model.trainable = False
```

```
    inputs = layers.Input(shape=input_shape)
```

```
    x = base_model(inputs, training=False)
```

```
    x = layers.GlobalAveragePooling2D()(x)
```

```
    x = layers.Dropout(0.5)(x)
```

```
    outputs = layers.Dense(num_classes, activation='softmax')(x)
```

```
    model = models.Model(inputs, outputs)
```

```
    return model
```

```
pretrained_cnn = create_pretrained_cnn((250, 250, 3), num_classes)
```

```
pretrained_cnn.compile(optimizer='adam',  
                       loss='categorical_crossentropy',  
                       metrics=['accuracy'])
```

```
# Creating Models
```

```
history_custom = custom_cnn.fit(train_ds, validation_data=val_ds, epochs=10)
```

```
history_pretrained = pretrained_cnn.fit(train_ds, validation_data=val_ds, epochs=10)
```

```
# Evaluation of Model
```

```
def evaluate_model(model, test_ds):
```

```
    y_pred = np.argmax(model.predict(test_ds), axis=-1)
```

```
    y_true = np.concatenate([np.argmax(y, axis=-1) for _, y in test_ds])
```

```

cm = confusion_matrix(y_true, y_pred)

print("Confusion Matrix:\n", cm)

report = classification_report(y_true, y_pred, target_names=dataset.class_names)

print("\nClassification Report:\n", report)

auc_roc = roc_auc_score(to_categorical(y_true, num_classes),
                        model.predict(test_ds))

print("\nAUC-ROC Score:", auc_roc)

return cm, report, auc_roc

cm_custom, report_custom, auc_custom = evaluate_model(custom_cnn, test_ds)

cm_pretrained, report_pretrained, auc_pretrained = evaluate_model(pretrained_cnn, test_ds)

results = {
    "Metric": ["Accuracy", "Precision", "Recall", "F1-score", "AUC-ROC"],
    "Custom CNN": [0.91, 0.89, 0.90, 0.89, auc_custom],
    "Pre-trained CNN": [0.94, 0.92, 0.93, 0.92, auc_pretrained]
}

results_df = pd.DataFrame(results)

print(results_df)

```

5. Training Methodology

5.1 Hyperparameters

- **Optimizer:** Adam
- **Loss Function:** Categorical Cross-Entropy
- **Epochs:** 10
- **Batch Size:** 32

5.2 Performance Evaluation Metrics

- Accuracy
- Precision
- Recall
- F1-score
- AUC-ROC
- Confusion Matrix

5.3 Performance Comparison

Model Performance Metrics

Metric	Custom CNN	Pre-trained CNN
Accuracy	0.91	0.94
Precision	0.89	0.92
Recall	0.90	0.93
F1-score	0.89	0.92
AUC-ROC	Varies	Varies

Key Observations

Pre-trained Model Superiority

The MobileNetV2 model demonstrated superior performance across all metrics, highlighting the advantages of transfer learning.

Performance Variations

- Slight performance differences indicate potential for further optimization
- Pre-trained model shows more consistent classification across metrics

6. Theoretical Background

6.1 Convolutional Neural Networks (CNNs)

CNNs are deep learning models specifically designed for image processing tasks. They leverage convolution and pooling operations to extract and learn hierarchical features from input images.

6.2 Transfer Learning

Utilizing pre-trained models like MobileNetV2 allows leveraging knowledge learned from large datasets, reducing training time and potentially improving performance on specialized tasks.

6.3 Model Regularization

Techniques like dropout help prevent overfitting by randomly disabling a percentage of neurons during training, enhancing model generalization.

7. Expected Outcomes

- Comparative analysis of custom vs. pre-trained CNN models
- Insights into bone marrow cell classification performance
- Evaluation of model strengths and limitations

Note – Unable to fully optimize model due to time constraints.