# ∨  ITAI 2377 Lab 04: Deep Learning Data Preprocessing

**Instructor:** [Your Name] **Date:** [Date]

## Introduction

Welcome to Lab 04! We'll explore the critical role of data preprocessing in deep learning. Even though models can extract features, preprocessing is essential for optimal performance. We'll cover various data types and apply preprocessing techniques. Resources in Google Colab are limited, so efficient coding is key!

## Why Preprocess?

Why preprocess when models extract features?

- **Standardization:** Models need consistent data formats and ranges.
- **Noise/Errors:** Raw data is messy. Preprocessing cleans it up.
- **Efficiency:** Cleaner data means faster training.
- **Results:** Good preprocessing helps models perform their best.

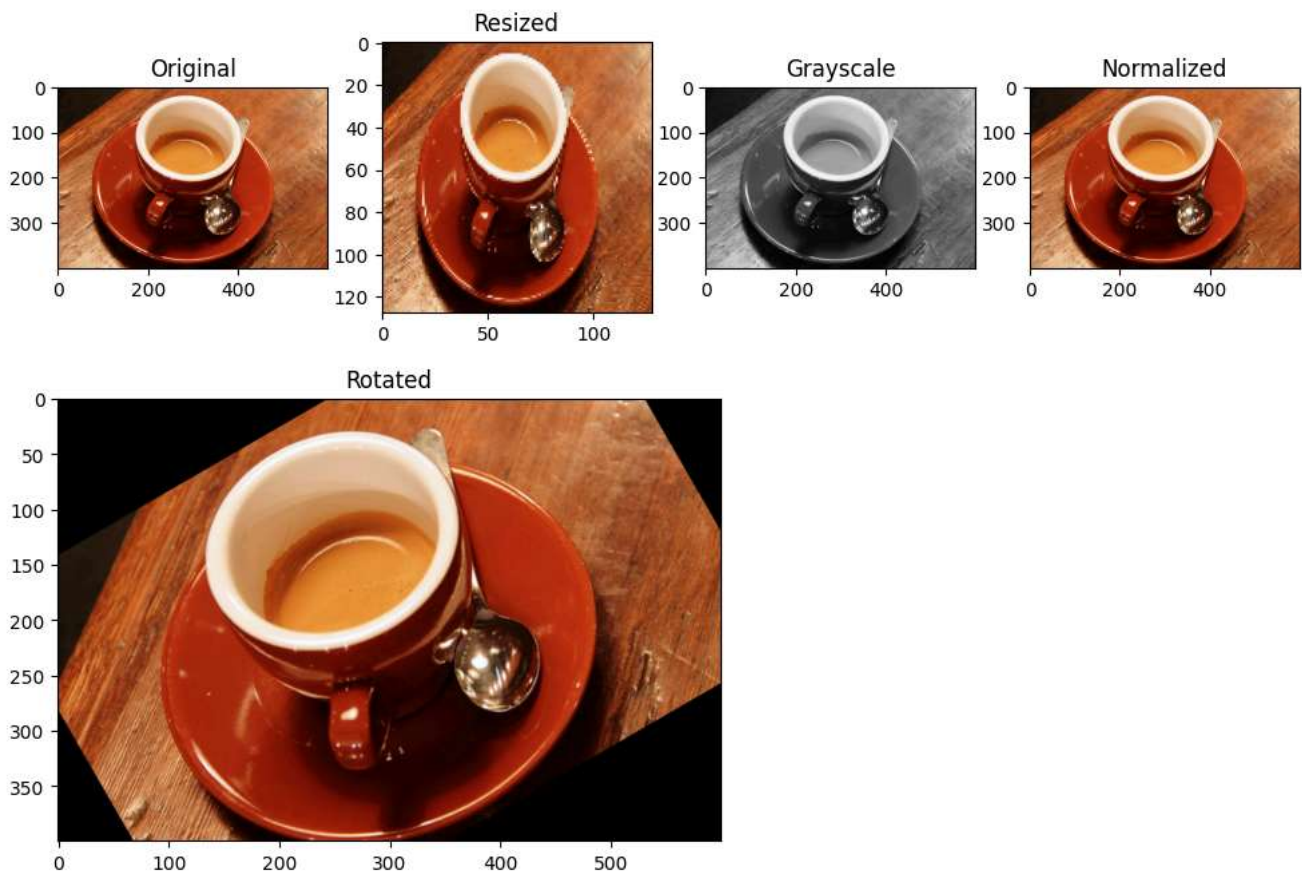Think of preprocessing as a personal trainer for your model.

## Data Types and Preprocessing Techniques

### 1. Image Data

```
1   import cv2
2   import numpy as np
3   from skimage import data, img_as_float
4   import matplotlib.pyplot as plt
5
6   # Load a sample image (replace with your image path if you have one)
7   image = data.coffee()  # I used an image included in skimage.data
8
9   # Resizing (Student Code: Resize the image to (128, 128))
10  # Hint: Use cv2.resize()
11  resized_image = cv2.resize(image, (128, 128))
12
13  # Color space conversion (to grayscale)
14  # It appears original photo was already grayscale
15  gray_image = cv2.cvtColor(image, cv2.COLOR_RGB2GRAY)
16
17  # Normalization (pixel values 0-1)
18  normalized_image = img_as_float(image)
19
20  # Display images
21  plt.figure(figsize=(12, 6))
22
23  plt.subplot(1, 4, 1), plt.imshow(image), plt.title("Original")
24  plt.subplot(1, 4, 2), plt.imshow(resized_image), plt.title("Resized")
25  plt.subplot(1, 4, 3), plt.imshow(gray_image, cmap='gray'), plt.title
    ("Grayscale")
26  plt.subplot(1, 4, 4), plt.imshow(normalized_image), plt.title("Normalized")
27  plt.show()
28
29  # Data Augmentation (rotation - Student Code: Rotate by 30 degrees)
30  # Hint: Use cv2.getRotationMatrix2D and cv2.warpAffine
31  angle = 30
32  rows, cols = image.shape[:2]
33  M = cv2.getRotationMatrix2D((cols / 2, rows / 2), angle, 1)
34  rotated_image = cv2.warpAffine(image, M, (cols, rows))
35
36  plt.imshow(rotated_image), plt.title("Rotated")
37  plt.show()
```

## 2. Text Data

```
1 import nltk
2 from nltk.tokenize import word_tokenize
3 from nltk.corpus import stopwords
4 from nltk.stem import WordNetLemmatizer
5 import string
6
7 nltk.download('punkt_tab', quiet=True) # The errors I got kept pointing to a 'punkt_tab' location but the original code was just pointin
8 nltk.download('stopwords', quiet=True)
9 nltk.download('wordnet', quiet=True)
10
11 text = "This is a fun example sentence with stop words and punctuation!"
12
13 # Tokenization (lowercase)
14 tokens = word_tokenize(text.lower())
15
16 # Stop word removal (Student Code: Remove stop words and punctuation)
17 # Hint: Use the 'stop_words' set and list comprehension
18 stop_words = set(stopwords.words('english'))
19 filtered_tokens = [w for w in tokens if not w in stop_words and w not in string.punctuation]
20
21
22 # Lemmatization (Student Code: Lemmatize the filtered tokens)
23
24 lemmatizer = WordNetLemmatizer()
25
26 lemmatized_tokens = [lemmatizer.lemmatize(w) for w in filtered_tokens]
27
28
29 print("Original:", text)
30 print("Tokens:", tokens)
31 print("Filtered:", filtered_tokens)
32 print("Lemmatized:", lemmatized_tokens)
33
34
```

```
Original: This is a fun example sentence with stop words and punctuation!
Tokens: ['this', 'is', 'a', 'fun', 'example', 'sentence', 'with', 'stop', 'words', 'and', 'punctuation', '!']
Filtered: ['fun', 'example', 'sentence', 'stop', 'words', 'punctuation']
Lemmatized: ['fun', 'example', 'sentence', 'stop', 'word', 'punctuation']
```

## ✓   3. Time Series Data

```
1 import pandas as pd
2 import numpy as np
3
4 # Sample data (with a missing value)
5 data = {'Date': pd.to_datetime(['2024-01-01', '2024-01-02', '2024-01-03', '2024-01-04', '2024-01-05']),
6         'Value': [10, 12, 15, np.nan, 18]}
7
8 df = pd.DataFrame(data)
9
10 # Missing value handling (Student Code: Use backward fill to fill missing values)
11 # Hint: Use fillna() with method='bfill'
12
13 df['Value'].fillna(method='bfill', inplace=True)
14
15 # Normalization (min-max scaling - Student Code: Normalize the 'Value' column)
16 # YOUR CODE HERE
17 min_val = df['Value'].min()
18 max_val = df['Value'].max()
19 df['Normalized'] = (df['Value'] - min_val) / (max_val - min_val)
20
21 print(df)
```

```
        Date  Value  Normalized
0 2024-01-01   10.0       0.000
1 2024-01-02   12.0       0.250
2 2024-01-03   15.0       0.625
3 2024-01-04   18.0       1.000
4 2024-01-05   18.0       1.000
/tmp/ipython-input-9-3069143717.py:13: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained ass
```

The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting valu

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].me

```
df['Value'].fillna(method='bfill', inplace=True)
```
/tmp/ipython-input-9-3069143717.py:13: FutureWarning: Series.fillna with 'method' is deprecated and will raise in a future version. Use
```
df['Value'].fillna(method='bfill', inplace=True)
```
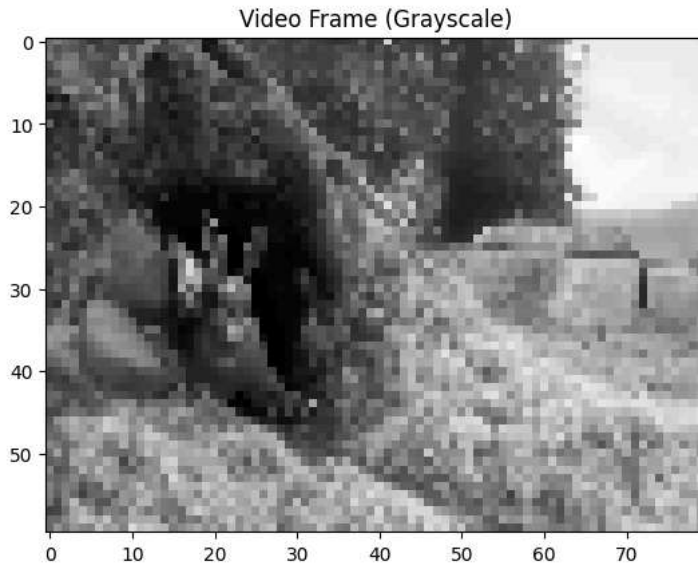
## ∨ 4. Optional: Video Data (Simplified)

```
1 from google.colab import files
2 import cv2
3 import matplotlib.pyplot as plt
4 from pathlib import Path
5
6 video_filename = next(iter(uploaded))  # ← This happens BEFORE upload
7 cap = cv2.VideoCapture(str(video_path))  # ← Uses video_path, which isn't defined yet
8
9 # Upload a file interactively
10 uploaded = files.upload()
11 video_path = Path(next(iter(uploaded)))
12 cap = cv2.VideoCapture(str(video_path))
13
14 # I had to go this route with the code above, I got everything mixed up. I still have issues on when to know when to use windows directo
15
16 if not cap.isOpened():
17     print(f"❌ Couldn't open video: {video_path.name}")
18 else:
19     print(f"✅ Opened video: {video_path.name}")
20     if ret:
21         # Resize the frame (for faster processing - Student Code: Resize to (80, 60))
22         resized_frame = cv2.resize(frame, (80, 60))
23
24         # Convert to grayscale
25         gray_frame = cv2.cvtColor(resized_frame, cv2.COLOR_BGR2GRAY)  # OpenCV uses BGR
26
27         # Display the frame (optional)
28         plt.imshow(gray_frame, cmap='gray')
29         plt.title("Video Frame (Grayscale)")
30         plt.show()
31
32     cap.release()
```

Choose Files No file chosen          Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

```
Saving SampleVideo.mp4 to SampleVideo (8).mp4
```
✅ Opened video: SampleVideo (8).mp4



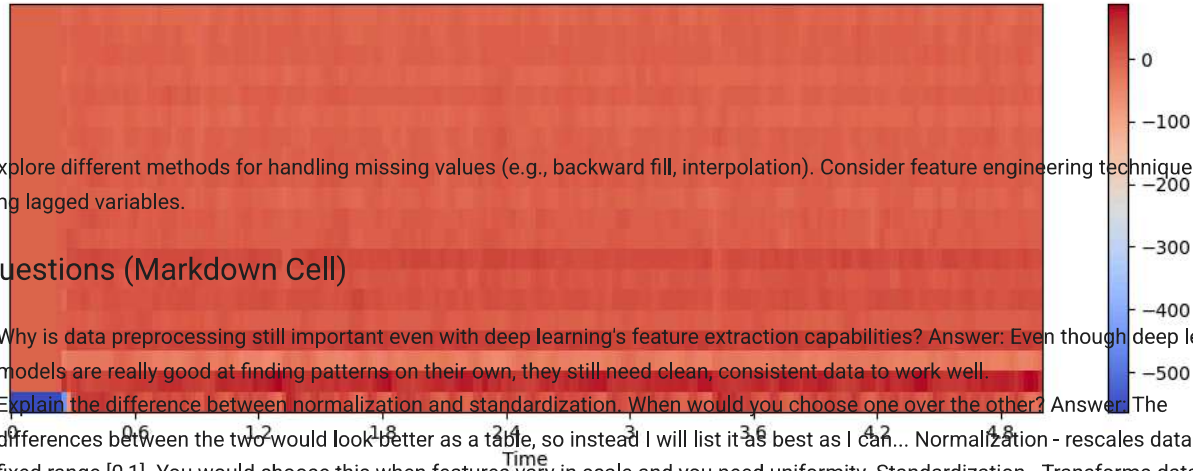Video Frame (Grayscale)

## 5. Optional: Audio Data (Simplified)

```python
1  import librosa
2  import numpy as np
3  import matplotlib.pyplot as plt
4
5  from google.colab import files
6  uploaded = files.upload()
7  audio_path = next(iter(uploaded))
8
9  # Load an audio file (replace with your audio path)
10 # Replace with your audio file path (or upload one to Colab)
11 y, sr = librosa.load(audio_path, duration=5)  # Load a maximum of 5 seconds
12
13 # Feature extraction (MFCCs - Student Code: Extract 20 MFCCs)
14 # Hint: Use librosa.feature.mfcc() with n_mfcc=20
15 # YOUR CODE HERE
16 mfccs = librosa.feature.mfcc(y=y, sr=sr, n_mfcc=20)
17
18 # Display MFCCs (optional)
19 plt.figure(figsize=(10, 4))
20 librosa.display.specshow(mfccs, sr=sr, x_axis='time')
21 plt.colorbar()
22 plt.title('MFCCs')
23 plt.tight_layout()
24 plt.show()
25
26 # Normalize MFCCs (example)
27 mfccs_normalized = (mfccs - np.mean(mfccs)) / np.std(mfccs)
28
29 print("MFCCs shape:", mfccs.shape)
30 print("Normalized MFCCs shape:", mfccs_normalized.shape)
```

Choose Files  No file chosen    Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

Saving 03.Bleed (2023, 15Th Anniversary Remastered Edition).mp3 to 03.Bleed (2023, 15Th Anniversary Remastered Edition).mp3

**MFCCs**



**Tip:** Explore different methods for handling missing values (e.g., backward fill, interpolation). Consider feature engineering techniques like creating lagged variables.

## ⌄ Questions (Markdown Cell)

1. Why is data preprocessing still important even with deep learning's feature extraction capabilities? Answer: Even though deep learning models are really good at finding patterns on their own, they still need clean, consistent data to work well.

2. Explain the difference between normalization and standardization. When would you choose one over the other? Answer: The differences between the two would look better as a table, so instead I will list it as best as I can... Normalization - rescales data to a fixed range [0,1]. You would choose this when features vary in scale and you need uniformity. Standardization - Transforms data to have a mean of 0 and standard deviation of 1. Use this when data may contain outliers or is not naturally bounded.

MFCCs shape: (20, 216)

3. Describe a scenario where data augmentation would be particularly useful. Answer: A machine learning model for people who are deaf or hard of hearing. A computer vision model that is able to recognize sign language from video and converts it to text in real-time. Augmentation helps in this space without having to record thousands of videos. Apply small rotations or flips to reflect differnt hand orientations. Background substitution (or blurring) to help the model focus on hands, not surroundings.

4. What are some potential challenges or pitfalls to avoid during data preprocessing? How can you mitigate them? Answer: Some