

به نام خدا



دانشگاه تهران



دانشکده مهندسی برق و کامپیوتر

**درس شبکه‌های عصبی و یادگیری عمیق**

**تمرین پنجم**

نام و نام خانوادگی	علی صفری	پرسش ۱
شماره دانشجویی	۸۱۰۲۰۲۱۵۳	
نام و نام خانوادگی	حمیدرضا نادى مقدم	پرسش ۲
شماره دانشجویی	۸۱۰۱۰۳۲۶۴	
مهلت ارسال پاسخ	۱۴۰۳.۱۰.۱۳	

## فهرست

مقدمه	۱
پرسش ۱. پیشبینی نیروی باد به کمک مبدل و تابع خطای Huber	۲
۱-۱. مقدمه	۲
۲-۱. آماده سازی	۴
۱-۲-۱. Autoencoder	۴
۲-۲-۱. مکانیزم توجه و Positional encoding	۹
۳-۲-۱. تابع Huber	۱۱
۴-۲-۱. الگوریتم Slime mould	۱۳
۳-۱. روش شناسی و نتایج	۱۸
۱-۳-۱. انتخاب دیتا و نمایش شکل ۶ مقاله	۱۸
۲-۳-۱. شناسایی داده های پرت	۱۹
۳-۳-۱. دینویز و پیش پردازش داده ها (۵ + ۱۵ نمره)	۲۱
۴-۳-۱. پیاده سازی مدل با ۳ معماری و ۲ تابع loss برای حالت single step	۲۷
۵-۳-۱. بررسی و ارزیابی عملکرد مدل ها	۳۲
۶-۳-۱. مقادیر بهینه ابرپارامتر با استفاده از slime mould (امتیازی)	۳۹
۷-۳-۱. پیاده سازی مدل مبدل با استفاده از دو تابع loss برای حالت multi step	۴۳
پرسش ۲. استفاده از ViT برای طبقه بندی تصاویر گلوبولهای سفید	۴۷
۱-۲. آماده سازی داده ها	۴۷
۲-۲. آموزش مدل ها	۵۱
۳-۲. تحلیل و نتیجه گیری	۵۹

## شکل‌ها

- شکل ۱- تصویر شماتیک مقاله از encoder و decoder ..... ۶
- شکل ۲- کد ساخت Autoencoder بخش اول ..... ۸
- شکل ۳- کد ساخت Autoencoder بخش دوم ..... ۹
- شکل ۴- فرمول Huber loss ..... ۱۱
- شکل ۵- پیاده‌سازی تابع Huber loss ..... ۱۲
- شکل ۶- تصویر الگوریتم SMA از مقاله ..... ۱۵
- شکل ۷- پیاده‌سازی SMA بخش اول ..... ۱۶
- شکل ۸- پیاده‌سازی SMA بخش دوم ..... ۱۷
- شکل ۹- کد مربوط به پیاده‌سازی شکل ۶ مقاله ..... ۱۸
- شکل ۱۰- نمایش شکل ۶ مقاله ..... ۱۹
- شکل ۱۱- کد نوشته شده برای تشخیص داده‌های پرت ..... ۲۰
- شکل ۱۲- نمودار تشخیص داده‌های پرت (مشابه شکل ۷ مقاله) ..... ۲۱
- شکل ۱۳- کد مربوط به دینوز کردن داده ..... ۲۵
- شکل ۱۴- نمودار داده‌های دینوز شده و داده‌های اولیه (مشابه شکل ۸ مقاله) ..... ۲۶
- شکل ۱۵- خروجی داده آموزش و تست بعد از نرمال کردن و sliding window برای single step ..... ۲۷
- شکل ۱۶- کد ساخت مدل transformer ..... ۲۸
- شکل ۱۷- کد ساخت مدل RNN ..... ۲۹
- شکل ۱۸- کد ساخت مدل MLP ..... ۳۰
- شکل ۱۹- نمونه خروجی فرایند آموزش ..... ۳۲
- شکل ۲۰- جدول مقایسه پارامترهای بدست آمده برای هر مدل مشابه جدول ۳ مقاله ..... ۳۴
- شکل ۲۱- جدول مقایسه پارامترهای بدست آمده برای هر مدل مشابه جدول ۳ مقاله (one feature) ..... ۳۷
- شکل ۲۲- کد تعریف کلاس بهینه‌سازی هایپرپارامترها ..... ۴۱
- شکل ۲۳- نتایج بهینه کردن هایپرپارامترهای مدل مبدل ..... ۴۲
- شکل ۲۴- کد ساخت مدل transformer برای multi step ..... ۴۴
- شکل ۲۵- آماده‌سازی مدل برای آموزش mutli step ..... ۴۴
- شکل ۲۶- بخشی از فرایند آموزش مدل mutli step ..... ۴۵

- شکل ۲۷ - نتایج مدل transformer در حالت mutli step ..... ۴۶
- شکل ۲۸ - کد نمایش تصادفی یک تصویر از هر کلاس ..... ۴۷
- شکل ۲۹ - نمایش تصادفی یک تصویر از هر کلاس ..... ۴۸
- شکل ۳۰ - تعداد تصاویر در هر کلاس ..... ۴۸
- شکل ۳۱ - متوازن کردن کلاس ها با تقویت داده ..... ۵۰
- شکل ۳۲ - تقسیم داده به دو قسمت آموزش و ارزیابی ..... ۵۱
- شکل ۳۳ - بارگذاری مدل ViT و نمایش معماری آن ..... ۵۲
- شکل ۳۴ - فقط دسته-بند Classifier قابل آموزش باشد ..... ۵۴
- شکل ۳۵ - نمودار Accuracy و Loss آموزش و ارزیابی برای حالت یک ..... ۵۵
- شکل ۳۶ - نمودار Accuracy و Loss آموزش و ارزیابی برای حالت دو ..... ۵۵
- شکل ۳۷ - نمودار Accuracy و Loss آموزش و ارزیابی برای حالت سه ..... ۵۶
- شکل ۳۸ - نمودار Accuracy و Loss آموزش و ارزیابی برای حالت چهار ..... ۵۶
- شکل ۳۹ - نمودار Accuracy و Loss آموزش و ارزیابی برای حالت پنج ..... ۵۷
- شکل ۴۰ - نمودار Accuracy و Loss برای مدل CCN آموزش فقط در لایه طبقه بند ..... ۵۸
- شکل ۴۱ - مقایسه مدل CNN و ViT با آموزش در لایه طبقه بند ..... ۵۸

## جدول‌ها

جدول ۱- مقایسه مقاله و نتایج ما در single step ..... ۳۵

جدول ۲- مقایسه مقاله و نتایج ما در (one feature) single step ..... ۳۸

برای پیاده‌سازی پروژه از بستر Google Colab به منظور کد نویسی و اجرا استفاده شده است. تمامی مراحل کد و اجرای آن در این گزارش به تفصیل شرح داده شده است.

کد های نوشته شده همگی در پوشه‌ی Code و با پسوند ipynd ذخیره شده است.

## پرسش ۱. پیشبینی نیروی باد به کمک مدل و تابع خطای Huber

### ۱-۱. مقدمه

محدودیت‌های روش‌های آماری سنتی مانند **ARIMA** و **GARCH**

۱. فرضیات سختگیرانه درباره توزیع داده‌ها: روش‌های آماری مانند **ARIMA** و **GARCH** بر پایه پیش‌فرض‌های توزیع مشخص و هموار بودن داده‌ها طراحی شده‌اند، که در مواردی که داده‌ها ناپایدار یا دارای نویز باشند، کارایی این روش‌ها کاهش می‌یابد.
۲. نیاز به آزمایش‌های هموارسازی (**smoothness test**) داده‌ها: این روش‌ها برای تحلیل صحیح نیازمند انجام آزمایش‌هایی مانند هموارسازی و بررسی ایستایی داده‌ها هستند، که این فرآیند زمان‌بر است و مدل را به شرایط خاص محدود می‌کند.
۳. محدودیت در تعمیم‌پذیری: این مدل‌ها در شناسایی وابستگی‌های بلندمدت و پیچیدگی‌های غیرخطی ناکارآمد هستند، و قابلیت تطبیق با داده‌های متنوع و محیط‌های متغیر را ندارند.
۴. کاهش دقت در شرایط ناپایدار: تغییرات شدید و نوسانات در داده‌ها می‌توانند به کاهش عملکرد این مدل‌ها منجر شوند.

مزایای مدل‌های یادگیری ماشین مانند **Random Forest**، **SVM** و **XGBoost**

۱. انعطاف‌پذیری در مواجهه با داده‌های متنوع: این مدل‌ها می‌توانند ویژگی‌های مختلف و داده‌های بزرگ و پیچیده را مدیریت کنند.
۲. دقت بالاتر پیش‌بینی: توانایی یادگیری الگوهای پیچیده باعث بهبود دقت پیش‌بینی نسبت به روش‌های سنتی می‌شود.
۳. عدم نیاز به پیش‌فرض‌های توزیعی سختگیرانه: این مدل‌ها در مقایسه با روش‌های آماری محدودیت کمتری در پیش‌فرض‌های داده دارند.
۴. عملکرد بهتر در داده‌های نویزدار: این مدل‌ها توانایی مدیریت نویزها و داده‌های پرت را دارند.

## چگونگی رفع محدودیت‌های روش‌های یادگیری ماشین سنتی توسط مدل‌های یادگیری عمیق مانند LSTM ، GRU و CNN

۱. مدیریت وابستگی‌های بلندمدت: مدل‌هایی مانند GRU و LSTM قادر به شناسایی و یادگیری وابستگی‌های بلندمدت در داده‌های سری زمانی هستند.
۲. استخراج ویژگی‌های پیچیده: مدل CNN می‌تواند الگوهای محلی و ویژگی‌های مهم داده‌های پیچیده را با استفاده از فیلترهای کانولوشنی شناسایی کند.
۳. کارایی بهتر در مسائل غیرخطی: این مدل‌ها برای داده‌های غیرخطی و پیش‌بینی‌هایی که به تحلیل دقیق‌تر وابستگی‌های زمانی نیاز دارند، عملکرد بسیار بهتری نسبت به مدل‌های سنتی ارائه می‌دهند.
۴. پیش‌بینی دقیق‌تر در داده‌های سری زمانی: این مدل‌ها می‌توانند با دقت بیشتری ویژگی‌های زمانی را تحلیل کنند.

## اهمیت مکانیزم خودتوجهی (Self-Attention) در شبکه‌های مبدل (Transformer)

۱. مدل‌سازی وابستگی‌های بلندمدت: این مکانیزم وابستگی‌های طولانی در داده‌های سری زمانی را شناسایی و مدل‌سازی می‌کند.
۲. استخراج اطلاعات در مقیاس‌های مختلف زمانی: مکانیزم خودتوجهی همبستگی‌های محلی (local) و جهانی (global) داده‌ها را استخراج می‌کند.
۳. بهبود کارایی در داده‌های پیچیده: این مکانیزم در ترکیب با دیگر مدل‌ها عملکرد پیش‌بینی را به‌ویژه در سری‌های زمانی پیچیده بهبود می‌بخشد.
۴. کاهش وابستگی به توالی داده‌های مجاور: برخلاف RNN، این مکانیزم اطلاعات کل توالی را بدون توجه به موقعیت مکانی پردازش می‌کند.



## چگونگی افزایش پایداری و استحکام مدل‌های پیش‌بینی داده‌های باد فراساحلی توسط تابع خطای هیوبر (Huber)

۱. مدیریت داده‌های پرت و نوسانات شدید: تابع خطای Huber از ترکیب ویژگی‌های تابع خطای مربعات و قدرمطلق بهره می‌برد و حساسیت کمتری به داده‌های پرت دارد، که در داده‌های پرنوسان باد فراساحلی رایج است.
۲. افزایش پایداری فرآیند بهینه‌سازی: این تابع پیوستگی و مشتق‌پذیری دارد، که فرآیند تنظیم پارامترهای مدل را با روش‌هایی مانند گرادیان نزولی پایدارتر و کارآمدتر می‌کند.
۳. توازن بین دقت و استحکام: استفاده از این تابع خطا باعث کاهش اثرات منفی نویزها و داده‌های پرت شده و پیش‌بینی‌های دقیق‌تر و پایدارتر را برای داده‌های متغیر باد فراساحلی فراهم می‌آورد.
۴. افزایش دقت پیش‌بینی: با کاهش تاثیر نویز و داده‌های پرت، دقت پیش‌بینی به‌ویژه در داده‌های پرنوسان بهبود می‌یابد.

### ۲-۱. آماده‌سازی

#### ۱-۲-۱. Autoencoder

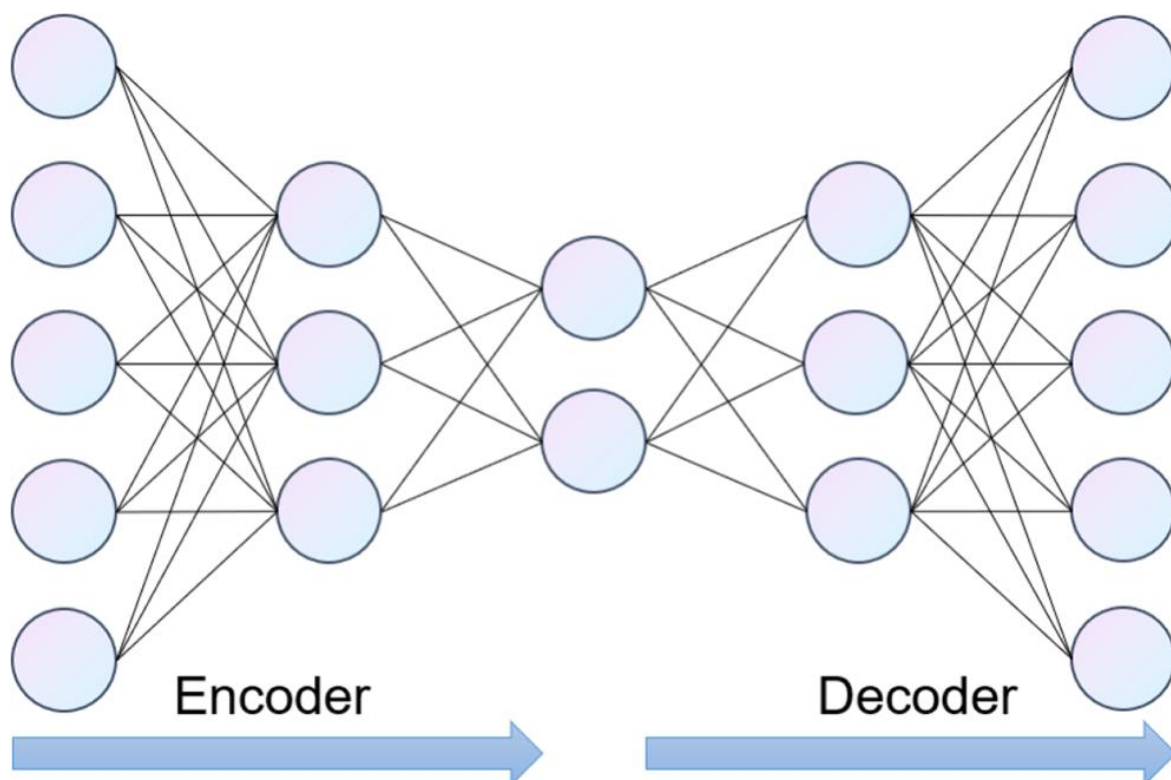
##### هدف استفاده از ساختار Autoencoder

- هدف اصلی استفاده از Autoencoder در این تحقیق، بازسازی و کاهش نویز داده‌های باد فراساحلی است.
- این ساختار برای شناسایی ویژگی‌های مهم داده‌ها و کاهش اثر نویز و داده‌های پرت طراحی شده است.
- به کمک این ساختار، داده‌ها بازسازی شده و بهبود دقت پیش‌بینی مدل‌های یادگیری عمیق ممکن می‌شود.

##### بخش‌های تشکیل‌دهنده Autoencoder

- بخش Encoder داده‌های ورودی را به یک فضای با ابعاد کمتر (نمایش مخفی) تبدیل می‌کند.

- بخش **Decoder** داده‌های کدگذاری شده را بازسازی می‌کند تا به داده‌های اصلی نزدیک شوند.
- تابع هزینه (**Loss Function**) میزان تفاوت بین داده‌های اصلی و بازسازی شده را اندازه‌گیری می‌کند.



شکل ۱- تصویر شماتیک مقاله از **encoder** و **decoder**

برای پیاده‌سازی کد این بخش از کتابخانه tensorflow استفاده کردیم. در این پروژه، معماری خودرمزنگار (Autoencoder) به صورت زیر طراحی شده است:

#### Encoder

- تعداد ابعاد ورودی یا نورون‌های لایه ورودی با توجه به اندازه پنجره (window size) که در اینجا ۱۴۴ است، تعریف شده است.
- داده‌ها در گام اول با استفاده از لایه Flatten به صورت یک بردار مسطح درمی‌آیند.
- یک لایه مخفی (hidden layer) با ۱۲۸ نورون و تابع فعال‌سازی ReLU معرفی شده است.
- در گام بعد، یک لایه مخفی دیگر با ۶۴ نورون و تابع فعال‌سازی ReLU به شبکه اضافه شده است.
- سپس لایه Latent تعریف شده که به عنوان فضای فشرده‌شده با تعداد نورون‌های تعیین‌شده در اینجا latent\_dim و تابع فعال‌سازی ReLU استفاده شده است.

#### Decoder

- پس از کدگذاری داده‌ها در لایه Latent، مراحل بازسازی (Decoding) آغاز می‌شود.
- لایه‌ای با ۶۴ نورون و تابع فعال‌سازی ReLU برای بازسازی داده‌ها به کار گرفته شده است.

- یک لایه دیگر با ۱۲۸ نورون و همان تابع فعال‌سازی اضافه شده است.
- در پایان، یک لایه خروجی (output layer) تعریف شده که با استفاده از Reshape، داده‌ها را به شکل اصلی سری زمانی (۱۴۴ گام زمانی) بازمی‌گرداند.

### ساختار نهایی:

- داده‌ها پس از عبور از لایه‌های Encoder فشرده‌شده و نویزهای ناخواسته از بین می‌روند.
  - سپس، داده‌ها در Decoder بازسازی شده و سری زمانی دینویز شده تولید می‌شود.
- مزیت این معماری:** این معماری با بهره‌گیری از لایه‌های متعدد و کاهش تدریجی تعداد نورون‌ها، امکان یادگیری ویژگی‌های مهم سری زمانی را فراهم کرده و با استفاده از Latent، داده‌ها را به صورت فشرده درآورده و نویز را حذف می‌کند. همچنین، بازسازی داده‌ها در Decoder تضمین می‌کند که سری زمانی خروجی با حداقل افت اطلاعات بازسازی شود.

```

class Autoencoder:
    def __init__(self, window_size, latent_dim):
        """
        Initialize the Autoencoder class.
        Args:
            window_size (int): Number of time steps in each sequence (input dimension).
            latent_dim (int): Size of the latent space (latent dimension).
        """
        self.window_size = window_size
        self.latent_dim = latent_dim
        self.autoencoder = None
        self.encoder = None
        self.decoder = None
    def build(self):
        """
        Build the Autoencoder, Encoder, and Decoder models.
        """
        # Encoder
        input_layer = Input(shape=(self.window_size, 1), name="Input_Layer")
        flattened = Flatten(name="Flatten_Layer")(input_layer) # Flatten time-series input
        hidden1 = Dense(128, activation='relu', name="Hidden_Layer_1")(flattened)
        hidden2 = Dense(64, activation='relu', name="Hidden_Layer_2")(hidden1)
        latent_layer = Dense(self.latent_dim, activation='relu', name="Latent_Layer")(hidden2)
        # Decoder
        hidden3 = Dense(64, activation='relu', name="Decoder_Layer_1")(latent_layer)
        hidden4 = Dense(128, activation='relu', name="Decoder_Layer_2")(hidden3)
        output_flat = Dense(self.window_size, activation='linear', name="Output_Flat")(hidden4)
        output_layer = Reshape((self.window_size, 1), name="Output_Reshape")(output_flat)
        # Autoencoder Model
        self.autoencoder = Model(inputs=input_layer, outputs=output_layer, name="Autoencoder")
        # Encoder Model
        self.encoder = Model(inputs=input_layer, outputs=latent_layer, name="Encoder")
        # Decoder Model
        encoded_input = Input(shape=(self.latent_dim,), name="Encoded_Input")
        decoder_hidden1 = self.autoencoder.get_layer("Decoder_Layer_1")(encoded_input)
        decoder_hidden2 = self.autoencoder.get_layer("Decoder_Layer_2")(decoder_hidden1)
        decoder_flat = self.autoencoder.get_layer("Output_Flat")(decoder_hidden2)
        decoder_output = self.autoencoder.get_layer("Output_Reshape")(decoder_flat)
        self.decoder = Model(inputs=encoded_input, outputs=decoder_output, name="Decoder")

```

شکل ۲ - کد ساخت **Autoencoder** بخش اول

```

def compile(self, optimizer='adam', loss='mse'):
    """
    Compile the Autoencoder model.

    Args:
        optimizer (str): Optimizer for training (default: 'adam').
        loss (str): Loss function to minimize (default: 'mse').
    """
    if not self.autoencoder:
        raise ValueError("The model must be built before compiling.")
    self.autoencoder.compile(optimizer=optimizer, loss=loss)

def train(self, x_train, epochs=100, batch_size=16, shuffle=True, validation_split=0.1):
    """
    Train the Autoencoder model.

    Args:
        x_train (np.array): Training input data (e.g., sequences of time steps).
        epochs (int): Number of epochs to train (default: 100).
        batch_size (int): Batch size for training (default: 16).
        shuffle (bool): Whether to shuffle the training data (default: True).
        validation_split (float): Fraction of data for validation (default: 0.1).

    Returns:
        History object: Contains training history.
    """
    if not self.autoencoder:
        raise ValueError("The model must be built and compiled before training.")
    return self.autoencoder.fit(
        x_train, x_train, # Input and target are the same for autoencoders
        epochs=epochs,
        batch_size=batch_size,
        shuffle=shuffle,
        validation_split=validation_split
    )

def summary(self):
    """
    Print the summary of the Autoencoder model.
    """
    if not self.autoencoder:
        raise ValueError("The model must be built before accessing the summary.")
    self.autoencoder.summary()

```

شکل ۳- کد ساخت **Autoencoder** بخش دوم

## ۲-۲-۱. مکانیزم توجه و Positional encoding

### مکانیزم توجه (Attention Mechanism) در شبکه مبدل

مکانیزم توجه یکی از اجزای کلیدی شبکه مبدل (Transformer) است و برای شناسایی و مدل سازی وابستگی های مهم بین قسمت های مختلف داده های سری زمانی استفاده می شود.

## هدف استفاده از مکانیزم توجه:

۱. شناسایی روابط مهم: مکانیزم توجه اطلاعاتی را که در بخش‌های مختلف داده مهم هستند شناسایی می‌کند و به آن‌ها وزن بیشتری اختصاص می‌دهد.
۲. مدل‌سازی وابستگی‌های بلندمدت: برخلاف مدل‌های سنتی مانند RNN که فقط وابستگی به نقاط نزدیک را در نظر می‌گیرند، مکانیزم توجه می‌تواند تمام نقاط در یک توالی را به صورت همزمان پردازش کند و وابستگی‌های بلندمدت را کشف کند.
۳. افزایش دقت پیش‌بینی: مکانیزم توجه اجازه می‌دهد که مدل بخش‌های مهم داده را شناسایی کرده و از آن‌ها برای پیش‌بینی دقیق‌تر استفاده کند.
۴. کاهش هزینه محاسباتی: در مقایسه با مدل‌های بازگشتی (RNN)، این مکانیزم امکان پردازش موازی داده‌ها را فراهم می‌کند که محاسبات را سریع‌تر می‌کند.

## Positional Encoding در شبکه مبدا

شبکه مبدا مکانیزم توجه را به کار می‌گیرد، اما این مکانیزم به تنهایی ترتیب داده‌ها را در نظر نمی‌گیرد. Positional Encoding به مدل کمک می‌کند تا موقعیت هر عنصر در توالی را تشخیص دهد.

## هدف استفاده از Positional Encoding

۱. حفظ اطلاعات ترتیب داده‌ها: از آنجا که مکانیزم توجه به صورت ذاتی اطلاعات ترتیب (ordering) را ندارد، Positional Encoding تضمین می‌کند که مدل بتواند موقعیت نسبی هر عنصر را درک کند.
۲. مدل‌سازی وابستگی‌های زمانی: این مکانیزم کمک می‌کند تا مدل روابط بین عناصر با فاصله زمانی مختلف را بهتر درک کند.
۳. بهبود کارایی مدل در داده‌های سری زمانی: Positional Encoding اطلاعات اضافی در مورد موقعیت نسبی داده‌ها را به مدل اضافه می‌کند، که در داده‌های سری زمانی مانند پیش‌بینی باد فراساحلی بسیار مفید است.

۴. حفظ ارتباط فضایی در داده‌های ورودی: این تکنیک تضمین می‌کند که مدل بتواند اطلاعات ساختاری مربوط به توالی را به صورت کامل حفظ کند.

### ۳-۲-۱. تابع Huber

تابع خطای Huber ترکیبی از خطای مربعات میانگین (Mean Squared Error) و خطای قدرمطلق (Mean Absolute Error) است. این تابع به صورت زیر عمل می‌کند:

- اگر خطا کوچک باشد (کمتر از مقدار آستانه  $\delta$ )، رفتار تابع مانند MSE است.
  - در این حالت، تابع حساسیت بیشتری به خطاهای کوچک دارد و به تنظیم دقیق مدل کمک می‌کند.
- اگر خطا بزرگ باشد (بیشتر از مقدار آستانه  $\delta$ )، رفتار تابع مانند MAE است.
  - در این حالت، تابع از تأثیر زیاد داده‌های پرت جلوگیری می‌کند و مدل را پایدارتر می‌سازد.

$$\text{If } |y - \hat{y}| \leq \delta:$$

$$L_{\delta}(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2$$

$$\text{If } |y - \hat{y}| > \delta:$$

$$L_{\delta}(y, \hat{y}) = \delta \cdot |y - \hat{y}| - \frac{1}{2}\delta^2$$

شکل ۴ - فرمول Huber loss

### هدف استفاده از تابع خطای Huber

۱. مدیریت داده‌های پرت: این تابع تأثیر داده‌های پرت را محدود می‌کند، که در داده‌های ناپایدار مانند باد فراساحلی بسیار مفید است.
۲. توازن بین دقت و پایداری: با ترکیب مزایای MSE و MAE، این تابع هم دقت بالایی ارائه می‌دهد و هم از مشکلات ناشی از خطاهای بزرگ جلوگیری می‌کند.



۳. پایداری بهینه‌سازی: مشتق‌پذیری تابع باعث می‌شود که فرآیند بهینه‌سازی با روش‌هایی مانند گرادینان نزولی پایدارتر انجام شود.

```
class CustomHuberLoss(nn.Module):
    def __init__(self, delta=1.0):
        """
        Initializes the CustomHuberLoss module.

        Args:
            delta (float): Threshold for switching between quadratic and linear loss. Default is 1.0.
        """
        super().__init__()
        self.delta = delta

    def forward(self, y_pred, y_true):
        """
        Computes the Huber loss.

        Args:
            y_pred (torch.Tensor): Predicted values. Shape: (batch_size,)
            y_true (torch.Tensor): True values. Shape: (batch_size,)

        Returns:
            torch.Tensor: Average Huber loss.
        """
        error = y_pred - y_true
        abs_error = torch.abs(error)
        loss = torch.where(
            abs_error <= self.delta,
            0.5 * error**2,
            self.delta * (abs_error - 0.5 * self.delta)
        )
        return loss.mean()
```

شکل ۵- پیاده‌سازی تابع Huber loss

کلاس CustomHuberLoss در چارچوب PyTorch طراحی شده و برای محاسبه (Huber Loss) به کار می‌رود. در ادامه توضیحاتی درباره ورودی و خروجی این تابع ارائه شده است:

#### ورودی‌ها:

##### ۱. $y_{pred}$ پیش‌بینی مدل:

- یک Tensor است که مقادیر پیش‌بینی شده توسط مدل را شامل می‌شود.
- شکل (Shape) آن برابر با  $(batch\_size,)$  است، به این معنا که هر مقدار در این آرایه مربوط به یک نمونه از داده‌ها در دسته (Batch) است.

##### ۲. $y_{true}$ مقادیر واقعی:

- یک Tensor که مقادیر هدف یا مقادیر واقعی مربوط به داده‌ها را در خود دارد.

- شکل آن نیز مشابه (batch\_size,) است و باید با  $y_{pred}$  هم اندازه باشد.

## خروجی‌ها:

### ۱. Huber Loss

- خروجی تابع یک اسکالر (Scalar) است که مقدار میانگین Huber loss را برای تمام نمونه‌های موجود در Batch محاسبه می‌کند.
- Huber loss ترکیبی از loss خطی و درجه دوم است که بستگی به مقدار خطا (Error) دارد. اگر خطا کمتر از مقدار آستانه (delta) باشد، زیان به صورت درجه دوم محاسبه می‌شود و در غیر این صورت، خطی خواهد بود.

## نکات کلیدی:

- مقدار  $\delta$  تعیین می‌کند که چه زمانی تابع از حالت درجه دوم به حالت خطی تغییر کند.
- این تابع به‌ویژه در مواردی که داده‌های پرت (Outliers) وجود دارند، عملکرد بهتری نسبت به زیان مربع خطا (MSE) دارد زیرا تأثیر داده‌های پرت را کاهش می‌دهد.
- ورودی‌های  $y_{true}$  و  $y_{pred}$  باید از نوع `torch.Tensor` باشند و بر روی یک دستگاه (Device) قرار داشته باشند CPU یا GPU
- این تابع به عنوان یک ماژول PyTorch زیر کلاس `nn.Module` قابل استفاده است و به راحتی می‌توان آن را در مدل‌های یادگیری عمیق به کار برد.

### ۱-۲-۴. الگوریتم Slime mould

**Slime Mould Algorithm (SMA)** یک الگوریتم الهام گرفته از طبیعت است که از رفتار لجن کپک در جستجوی غذا تقلید می‌کند. این الگوریتم برای حل مسائل بهینه‌سازی طراحی شده است و با شبیه‌سازی رفتارهای کاوش و بهره‌برداری لجن کپک، می‌تواند به راه‌حل‌های بهینه دست یابد.

## مراحل الگوریتم:

### ۱. مقدمه‌سازی اولیه:

- جمعیتی از موقعیت‌های تصادفی در فضای جستجو تولید می‌شود.
- محدوده‌های جستجو  $lb$  و  $ub$  و پارامترهای مسئله تعریف می‌شوند.

## ۲. ارزیابی تناسب:

- مقدار تناسب هر موقعیت (بر اساس تابع هدف) محاسبه می‌شود.

## ۳. محاسبه وزن‌ها:

- لجن کپک‌ها بر اساس رتبه‌بندی تناسب خود وزن دریافت می‌کنند.

## ۴. به‌روزرسانی موقعیت‌ها:

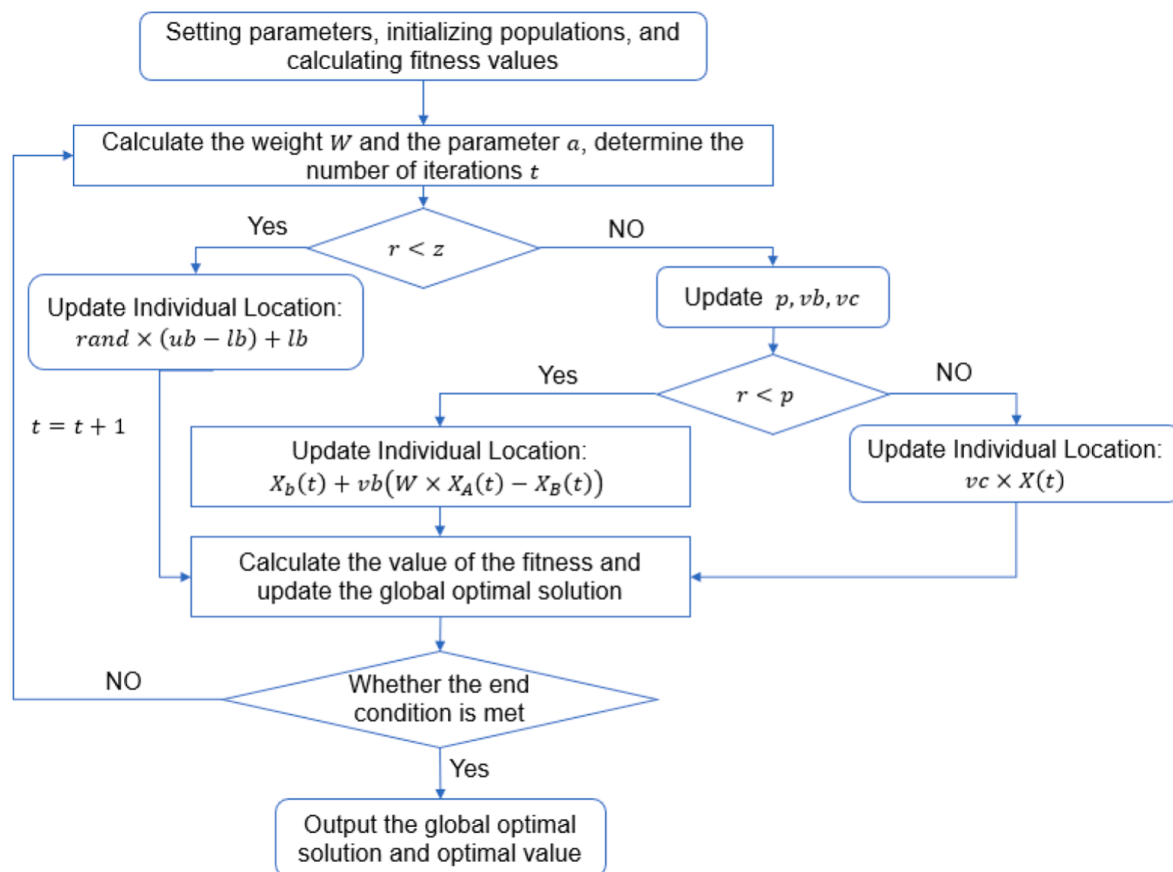
- موقعیت‌ها در دو فاز بهره‌برداری و کاوش به‌روزرسانی می‌شوند:

- بهره‌برداری (Exploitation) حرکت به سمت بهترین موقعیت (بیشترین منابع غذا).

- کاوش (Exploration) حرکت تصادفی برای جلوگیری از گیر افتادن در بهینه محلی.

## ۵. تکرار فرآیند:

- این فرآیند تا رسیدن به حداکثر تعداد تکرار یا تحقق معیار توقف ادامه می‌یابد.



شکل ۶ - تصویر الگوریتم SMA از مقاله

```

def slime_mould_algorithm(obj_func, N=30, max_iter=100, lb=-10, ub=10, dim=5):
    """
    Implementation of the Slime Mould Algorithm (SMA).

    Args:
        obj_func (function): The objective function to minimize.
        N (int): Number of slime moulds (population size).
        max_iter (int): Maximum number of iterations.
        lb (float or np.array): Lower boundary of the search space.
        ub (float or np.array): Upper boundary of the search space.
        dim (int): Dimensionality of the search space.

    Returns:
        best_position (np.array): Best solution found.
        best_fitness (float): Fitness value of the best solution.
    """
    # Initialize population
    population = np.random.uniform(lb, ub, (N, dim))
    fitness = np.apply_along_axis(obj_func, 1, population)

    # Initialize best solution
    best_idx = np.argmin(fitness)
    best_position = population[best_idx]
    best_fitness = fitness[best_idx]

    for t in range(max_iter):
        # Rank population and assign weights
        sorted_indices = np.argsort(fitness)
        population = population[sorted_indices]
        fitness = fitness[sorted_indices]

        # Update best position
        if fitness[0] < best_fitness:
            best_position = population[0]
            best_fitness = fitness[0]

```

شکل ۷ - پیاده‌سازی SMA بخش اول

```

# Calculate weights
W = 1 + np.log(1 + (fitness[-1] - fitness) / (fitness[-1] - fitness[0] + 1e-10))

# Update positions
for i in range(N):
    r = np.random.rand()
    vb = np.random.uniform(-1, 1)
    vc = np.random.uniform(0, 1)
    if r < vc:
        # Exploitation phase
        new_position = population[0] + vb * (W[i] * (population[i] - population[0]))
    else:
        # Exploration phase
        new_position = np.random.uniform(lb, ub, dim)

# Boundary control
population[i] = np.clip(new_position, lb, ub)

# Evaluate fitness of updated population
fitness = np.apply_along_axis(obj_func, 1, population)

return best_position, best_fitness

```

شکل ۸- پیاده‌سازی SMA بخش دوم

### توضیحات پارامترها (آرگومان‌های ورودی):

۱. **obj\_func**: تابع هدف که باید بهینه شود. در این مثال از تابع کره‌ای استفاده شده است.

۲. **N**: اندازه جمعیت (تعداد slime mould)

۳. **max\_iter**: حداکثر تعداد تکرارهای الگوریتم.

۴. **lb** و **ub**: محدوده پایین و بالا برای متغیرهای مسئله.

۵. **dim**: تعداد ابعاد یا متغیرهای مسئله.

این کد با استفاده از رفتارهای طبیعی slime mould و با توجه به الگوریتم نوشته شده در مقاله طراحی شده و می‌تواند برای مسائل مختلف بهینه‌سازی استفاده شود. البته لازم به ذکر است این تابع برای کاربرد عمومی نوشته شده است و در بخش امتیازی بسته به مسئله تغییر اندکی خواهد داشت که در همان بخش اشاره خواهد شد.

لازم به ذکر است که کد مربوط به این بخش و بخش امتیازی در فایلی تحت عنوان HW5\_1\_Slime ذخیره شده است.

## ۳-۱. روش شناسی و نتایج

### ۱-۳-۱. انتخاب دیتا و نمایش شکل ۶ مقاله

برای این کار از لینک قرار داده شده در مقاله دیتای مورد استفاده که در فایل WT5 بود را دانلود کرده و ۱۴۴۰ دیتای اول آن را استفاده کردیم. این فایل در فولدر اصلی ضمیمه شده است. در این قسمت ۸۰ درصد ابتدایی به عنوان داده آموزش و ۲۰ درصد انتهایی به عنوان داده تست در نظر گرفته شده اند.

```
# Step 2: Select the first 1440 rows and extract relevant columns
data = data.iloc[:1440] # Select only the first 1440 rows
time = np.arange(1, 1441) # Generate time steps (1 to 1440)
# Clean column names
data.columns = data.columns.str.strip() # Remove leading and trailing whitespaces
print("Cleaned Column Names:", data.columns) # Check cleaned column names

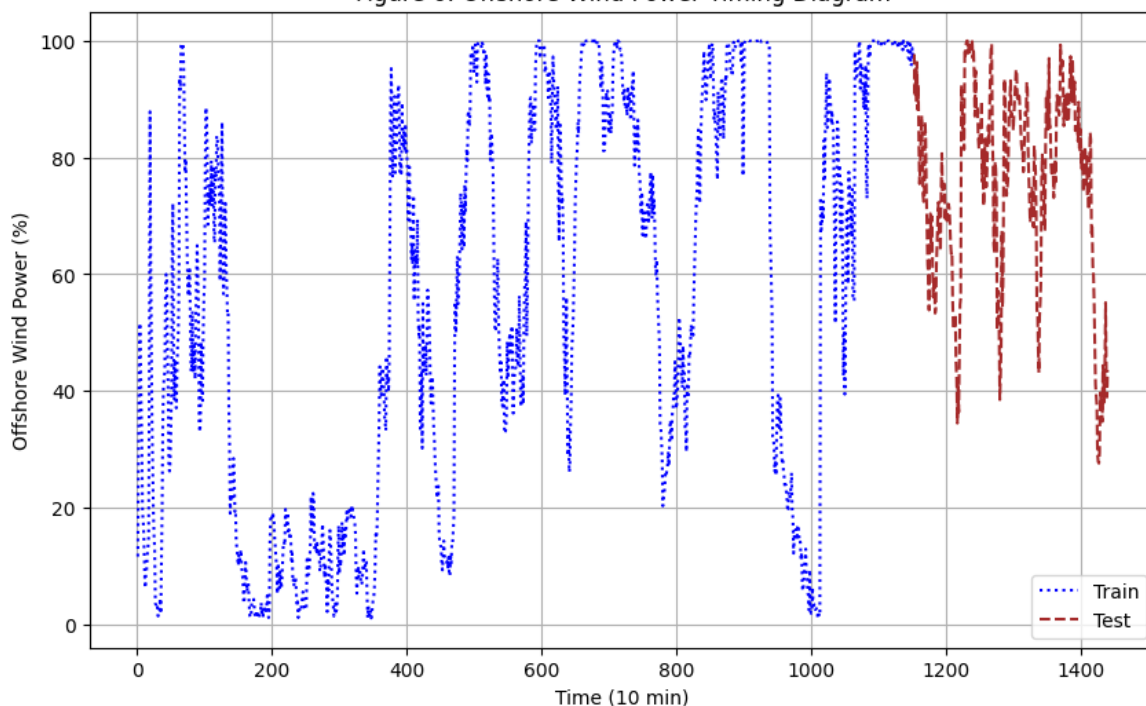
# Access the 'y' column
power = data["y (% relative to rated power)"] # Ensure this matches the cleaned column name

# Step 3: Plot Figure 6 - Offshore Wind Power Timing Diagram
train_size = int(0.8 * len(data)) # 80% train, 20% test split
train_power = power[:train_size]
test_power = power[train_size:]

plt.figure(figsize=(10, 6))
plt.plot(time[:train_size], train_power, label="Train", color="blue", linestyle="dotted")
plt.plot(time[train_size:], test_power, label="Test", color="brown", linestyle="dashed")
plt.xlabel("Time (10 min)")
plt.ylabel("Offshore Wind Power (%)")
plt.legend()
plt.title("Figure 6: Offshore Wind Power Timing Diagram")
plt.grid()
plt.show()
```

شکل ۹ - کد مربوط به پیاده سازی شکل ۶ مقاله

Figure 6: Offshore Wind Power Timing Diagram



شکل ۱۰- نمایش شکل ۶ مقاله

### ۲-۳-۱. شناسایی داده‌های پرت

کدی که نوشته‌ام برای شناسایی داده‌های پرت (Outliers) در مجموعه داده‌های توان بادی دریایی طراحی شده است. این الگوریتم بر اساس روش پنجره متحرک (Rolling Window) کار می‌کند. در این روش، پنجره‌ای با اندازه ۴۸ نقطه در نظر گرفته شده که شامل ۲۴ نقطه قبل و ۲۴ نقطه بعد از نقطه مورد بررسی است. برای هر پنجره، میانگین ( $\mu$ ) و انحراف معیار ( $\sigma$ ) داده‌های موجود محاسبه می‌شود و داده‌هایی که از بازه  $\mu \pm 2\sigma$  فراتر باشند، به عنوان داده پرت شناسایی و علامت‌گذاری می‌شوند.

در مقایسه با روش ارائه‌شده در مقاله، که از واریانس ( $\sigma^2$ ) به جای انحراف معیار استفاده کرده بود، تغییراتی انجام داده‌ام. استفاده از انحراف معیار به دلیل سازگاری بیشتر با مقیاس داده‌ها و سهولت تفسیر آن، انتخاب شده است. همچنین، پنجره به صورت متقارن تنظیم شده است تا هر دو سمت نقطه داده به طور مساوی بررسی شود. این تنظیم به شناسایی دقیق‌تر الگوهای محلی و رفتار داده‌ها کمک می‌کند. هرچند این جزئیات در مقاله به وضوح مشخص نشده بود، اما این تغییرات برای تحلیل داده‌های سری زمانی همچون این مجموعه منطقی و ضروری است.

برای نمایش نتایج، از نمودار پراکندگی (Scatter Plot) استفاده کرده‌ام و نقاط پرت با رنگ زرد مشخص شده‌اند. این روش بصری‌سازی داده‌ها را بهبود داده و تحلیل نتایج را آسان‌تر می‌کند. در مجموع، با وجود



شباهت به روش مقاله، تغییرات اعمال شده باعث ساده تر شدن فرایند تحلیل و تطابق بهتر الگوریتم با ماهیت این داده ها شده است.

```

window_size = 48
rolling_mean = []
rolling_std = [] # or rolling variance, depending on the exact interpretation

for i in range(len(power)):
    # Define the window boundaries (e.g., 48 points around i)
    # We'll do a simple centered window:
    start = max(0, i - window_size//2)
    end = min(len(power), i + window_size//2)
    # start = max(0, i)
    # end = min(len(power), i + window_size)

    window_data = power[start:end]

    # Compute mean & std (or variance)
    mu = window_data.mean()
    sigma = window_data.std() # std is the sqrt of variance

    # Keep track if needed
    rolling_mean.append(mu)
    rolling_std.append(sigma)

# Now define outliers if data[i] is beyond mean  $\pm 2 \times$  sigma in its window
outliers_mask = []
for i in range(len(power)):
    if power[i] < rolling_mean[i] - 2 * rolling_std[i] or power[i] > rolling_mean[i] + 2 * rolling_std[i]:
        outliers_mask.append(True)
    else:
        outliers_mask.append(False)

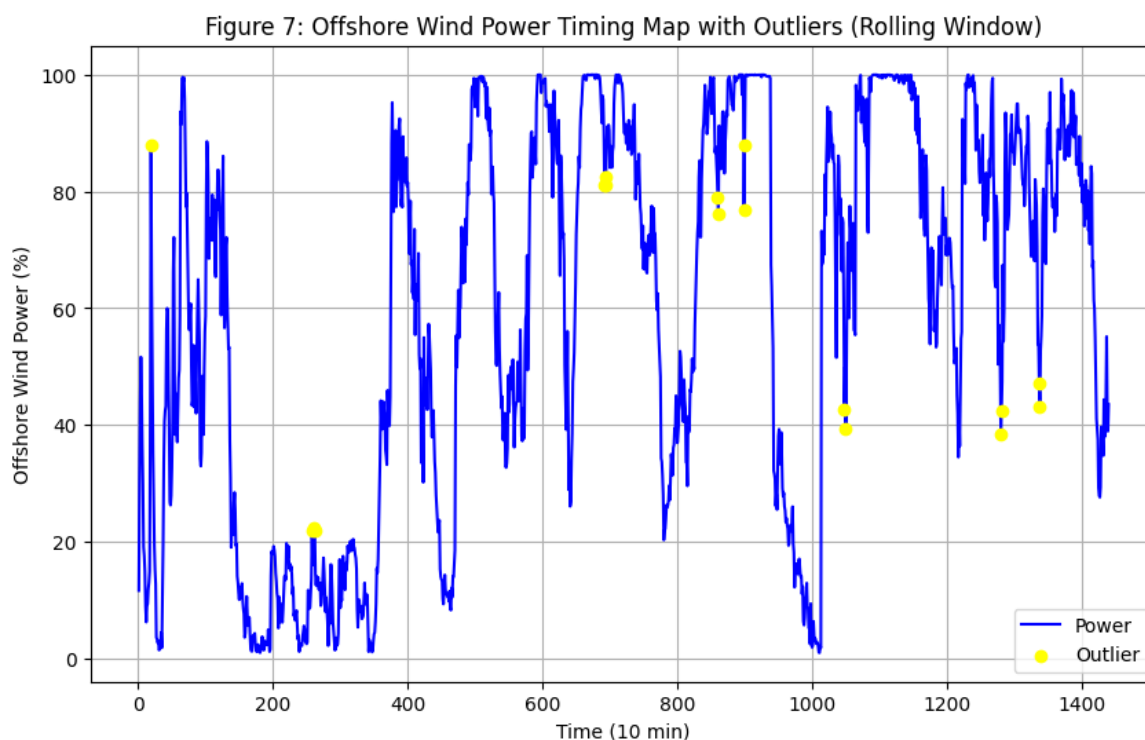
outliers_mask = np.array(outliers_mask)

# # Step 4: Outlier Detection Using Rolling Window (Paper's Method)

# Plot the power data with outliers highlighted (Figure 7)
plt.figure(figsize=(10, 6))
plt.plot(time, power, label="Power", color="blue")
plt.scatter(time[outliers_mask], power[outliers_mask], color="yellow", label="Outlier", zorder=5)
plt.xlabel("Time (10 min)")
plt.ylabel("Offshore Wind Power (%)")
plt.legend()
plt.title("Figure 7: Offshore Wind Power Timing Map with Outliers (Rolling Window)")
plt.grid()
plt.show()

```

شکل ۱۱ - کد نوشته شده برای تشخیص داده های پرت



شکل ۱۲ - نمودار تشخیص داده‌های پرت (مشابه شکل ۷ مقاله)

### ۳-۳-۱. دینویز و پیش پردازش داده‌ها (۵ + ۱۵ نمره)

گزارش جامع فرایند دینویز و پیش‌پردازش داده‌های توان بادی فراساحلی

در این بخش، داده‌های توان بادی فراساحلی پس از گردآوری اولیه، با هدف بهبود کیفیت و آماده‌سازی برای پیش‌بینی چندمرحله‌ای، تحت مراحل زیر قرار می‌گیرند:

- دینویز داده‌ها با شبکه خودرمزنگار (Autoencoder)
- تقسیم داده‌ها به مجموعه‌های آموزش و آزمون
- نرمال‌سازی داده‌ها
- تبدیل داده‌ها به متغیرهای ویژگی و پاسخ با پنجره متحرک

انتخاب و رعایت ترتیب مناسب این مراحل نقش مهمی در حفظ تعمیم‌پذیری مدل و جلوگیری از نشت داده‌ها ایفا می‌کند. در ادامه، به توضیح گام‌های فوق و مفروضات انجام‌شده می‌پردازیم.

### حذف نویز با شبکه خودرمزنگار

در این مرحله، داده‌های خام با استفاده از یک شبکه خودرمزنگار بازسازی شده و نویزهای ناخواسته آن‌ها حذف می‌گردد. در این پروژه، برای ورودی خودرمزنگار از پنجره‌های ۱۴۴ گام زمانی استفاده شده است؛ بدین معنا که برای هر بخش ۱۴۴ تایی از داده، مدل تلاش می‌کند خروجی آن را بازسازی (دینویز) نماید. البته لازم به ذکر است اینکه پنجره ورودی چه ابعادی در بخش دینویز کردن داشته باشد در مقاله ذکر نشده بود اما به جهت اینکه در آموزش مدل ۱۴۴ استفاده شده در این بخش هم از همین عدد استفاده گردیده است.

### معماری خودرمزنگار:

- لایه ورودی : (1, 144)
  - Flatten: برای کاهش ابعاد داده‌ها
  - لایه‌های پنهان: با ۱۲۸ و ۶۴ نورون
  - لایه کدگذاری (Latent Layer): با ابعاد مشخص شده (Latent Dim)
  - لایه‌های دیکدینگ: ۶۴ و ۱۲۸ نورون
  - خروجی بازسازی شده: با استفاده از Reshape به شکل اصلی داده‌ها
- فرض: استفاده از پنجره ۱۴۴ گامی برای حذف نویز و بازسازی داده‌ها به منظور دستیابی به اطلاعات بلندمدت‌تر و بهبود پیش‌بینی‌های چندمرحله‌ای.
- تقسیم داده به مجموعه‌های آموزش و آزمون
- پس از اتمام مرحله دینویز، داده‌های دینویز شده (به طول کل ۱۴۴۰) به دو بخش آموزش و آزمون تقسیم می‌شوند:

- ۸۰٪ داده‌ها: (بخش ابتدایی سری زمانی) برای آموزش
  - ۲۰٪ داده‌ها: (بخش انتهایی) برای آزمون
- این ترتیب تضمین می‌کند که مدل در مرحله آموزش هیچ دسترسی مستقیم یا غیرمستقیم به داده‌های آزمون نداشته باشد و از نشت اطلاعات جلوگیری شود.

### نرمال سازی داده‌ها

برای هم‌مقیاس کردن متغیرها و تسهیل فرایند آموزش شبکه‌های عصبی، داده‌های آموزش و آزمون به‌صورت جداگانه نرمال‌سازی می‌شوند. در این پروژه، از روش MinMaxScaler با بازه  $[0, 1]$  استفاده شده است:

- جدا کردن داده‌های آموزش و آزمون
  - فیت کردن MinMaxScaler فقط روی داده‌های آموزش
  - اعمال همان اسکالر روی داده‌های آزمون
- لازم به ذکر است که در انتها و برای ارزیابی به صورت معکوس نرمال‌سازی انجام شده است تا به ساختار اولیه بازگردیم.
- فرض:** انتخاب MinMaxScaler برای داده‌هایی که در دامنه محدودی قرار دارند و از توزیع ناهمگون برخوردارند، مناسب است.

### ایجاد متغیرهای ویژگی و پاسخ با پنجره ۱۴۴ گامی

پس از نرمال‌سازی، داده‌های دینویزشده و مقیاس‌شده برای پیش‌بینی چندمرحله‌ای به‌صورت زیر پردازش می‌شوند:

- **پنجره ویژگی (X) ۱۴۴ گام زمانی متوالی** که حاوی ۸ فیچر است.
- **افق پیش‌بینی (Horizon):** در مقاله ۲۴ گام زمانی آینده مطرح شده است اما در سوال ما دو حالت single step و multi step داریم. که در ادامه توضیح داده خواهد شد
- **خروجی (Y)** سری زمانی مورد پیش‌بینی در این گام مورد نظر (در صورتی که single step باشد در یک گام بعدی و در شرایطی که multi باشد تا جایی که مورد نظر سوال است).

### فرایند:

- داده‌های آموزش :
  - قطعه‌بندی در پنجره‌های متوالی ۱۴۴ تایی برای تشکیل ورودی (X)
  - گام یا گام‌های آینده برای تشکیل خروجی (Y)
- داده‌های آزمون :
  - اعمال همان سازوکار برای ساخت  $X_{test}$  و  $Y_{test}$

فرض :این مقادیر انتخابی هستند و بسته به ماهیت داده‌ها و اهداف پروژه می‌توانند تغییر کنند.

### اهمیت ترتیب مراحل

رعایت ترتیب زیر در پیش‌پردازش طبق توصیه‌های علمی ضروری است:

۱. دینویز داده‌ها: ابتدا حذف نویز برای بهبود کیفیت داده‌ها

۲. تقسیم به آموزش و آزمون: جلوگیری از نشت اطلاعات آزمون

۳. نرمال‌سازی: اسکالر فقط روی داده‌های آموزش فیت می‌شود

۴. تشکیل متغیرهای X و Y با پنجره متحرک

توجه: عدم رعایت ترتیب ممکن است باعث نشت اطلاعات داده‌های آزمون به مرحله آموزش شود.

### مثال نقض

برای توضیح اهمیت ترتیب مراحل، فرض کنید ابتدا داده‌ها نرمال‌سازی شوند و سپس به مجموعه‌های آموزش و آزمون تقسیم گردند. در این حالت، اسکالر بر روی کل داده‌ها (شامل داده‌های آزمون) فیت می‌شود. این موضوع باعث می‌شود که مدل به اطلاعات آماری داده‌های آزمون دسترسی پیدا کند. نتیجه این امر، کاهش واقعی چالش پیش‌بینی و افزایش مصنوعی دقت مدل است. اما در زمان مواجهه با داده‌های جدید، مدل نمی‌تواند عملکرد مشابهی ارائه دهد و دقت آن به شدت کاهش می‌یابد.

### جمع‌بندی

- داده‌های توان بادی با استفاده از شبکه خودرمن‌نگار و پنجره ۱۴۴ گامی نویززدایی شدند.
- داده‌های دینویز شده به دو بخش آموزش (۸۰٪) و آزمون (۲۰٪) تقسیم شدند.
- نرمال‌سازی داده‌ها با MinMaxScaler انجام شد.
- با استفاده از پنجره متحرک ۱۴۴ گامی و افق پیش‌بینی به تعداد گام مورد سوال ، متغیرهای ویژگی (X) و پاسخ (Y) ایجاد شدند.

این ترتیب تضمین‌کننده آموزش بدون نشت داده آزمون و دستیابی به پیش‌بینی‌های دقیق و تعمیم‌پذیر برای توان بادی است.

```

# Parameters
window_size = 144 # Time steps in each sequence
latent_dim = 16   # Size of the latent space

def create_sliding_windows(data, window_size):
    """
    Create sliding windows from the data.

    Args:
        data (np.array): Input data array (1D time-series).
        window_size (int): Number of time steps in each window.

    Returns:
        np.array: Array of sliding windows with shape (num_samples, window_size, 1).
    """
    windows = []
    for i in range(len(data) - window_size + 1):
        windows.append(data[i:i + window_size])
    return np.array(windows).reshape(-1, window_size, 1)

# Simulated input data (replace with your actual y column values)
sliding_windows = create_sliding_windows(power, window_size)

# Create and build the Autoencoder
autoencoder_model = Autoencoder(window_size, latent_dim)
autoencoder_model.build()

# Compile the Autoencoder
autoencoder_model.compile(optimizer='adam', loss='mse')

# Train the Autoencoder
history = autoencoder_model.train(
    x_train=sliding_windows,
    epochs=100,
    batch_size=16,
    shuffle=True,
    validation_split=0.1
)

# Display the model summary
autoencoder_model.summary()

# Use the trained Autoencoder to "denoise" (reconstruct) the data
data_denoised = autoencoder_model.autoencoder.predict(sliding_windows)

```

شکل ۱۳ - کد مربوط به دینوز کردن داده

**گزارش درباره تجمیع و هم‌پوشانی در Sliding Window**

در تحلیل داده‌های ترتیبی با Sliding Window، داده‌های بازسازی‌شده با هم‌پوشانی پنجره‌ها جمع می‌شوند. برای هر پنجره، مقادیر به موقعیت متناظر در آرایه بازسازی‌شده اضافه و شمارشگر برای هر موقعیت بروزرسانی می‌شود:

```
denoised_full[i:i + window_size] += data_denoised[i].flatten()
```

```
count[i:i + window_size] += 1
```

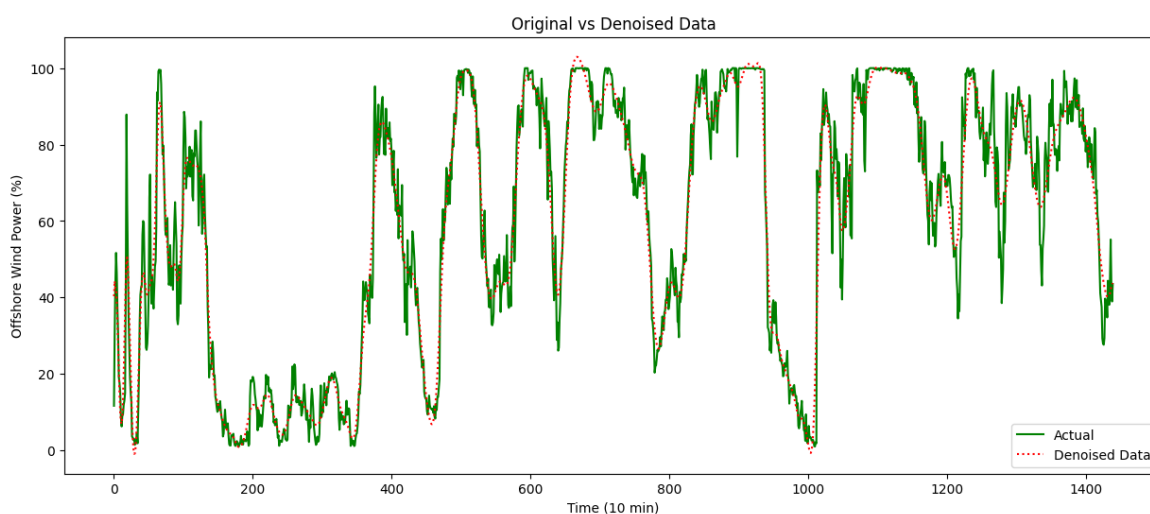
در پایان، مقادیر هر موقعیت با تعداد مشارکت‌ها تقسیم شده و میانگین‌گیری انجام می‌شود:

```
denoised_full /= count
```

این فرآیند نویز را کاهش داده و داده بازسازی‌شده‌ای با کیفیت بالا تولید می‌کند.

### مدیریت نویز و کاربردها

پیش از آموزش، این روش با هم‌پوشانی داده‌ها دقت پیش‌بینی را افزایش می‌دهد. پس از آموزش، داده‌های بازسازی‌شده بدون نویز اضافی ارائه می‌شوند. این روش با تنظیم اندازه پنجره و گام حرکتی به تعادل دقت و کارایی کمک می‌کند و کیفیت داده‌های ورودی و خروجی را بهبود می‌بخشد.



شکل ۱۴- نمودار داده‌های دینویز شده و داده‌های اولیه (مشابه شکل ۸ مقاله)

```
Train size: 1152, Test size: 288
X_train shape: (1008, 144, 8)
y_train shape: (1008, 1)
X_test shape: (144, 144, 8)
y_test shape: (144, 1)
Number of input features: 8
Number of time steps in each input window: 144
Forecast horizon: 1 (single-step prediction)
```

شکل ۱۵- خروجی داده آموزش و تست بعد از نرمال کردن و **sliding window** برای **single step**

۱-۳-۴. پیاده‌سازی مدل با ۳ معماری و ۲ تابع **loss** برای حالت **single step**

### مدل Transformer

مدل Transformer با استفاده از معماری پیشرفته خود برای داده‌های ترتیبی طراحی شده است. جزئیات معماری به شرح زیر است:

- **Positional Encoding**: استفاده از کدگذاری موقعیتی برای افزودن اطلاعات ترتیبی به داده‌ها.
  - ابعاد **Embedding**: ۶۴
  - تعداد هدهای توجه (**Attention Heads**): ۴
  - تعداد لایه‌های **Encoder**: ۳
  - ابعاد شبکه **Feedforward**: ۱۲۸
  - حداکثر طول توالی (**max\_len**): ۱۴۴
  - لایه خروجی (**Output Layer**): یک لایه خطی برای پیش‌بینی نهایی.
  - تابع فعال‌سازی **ReLU**: برای شبکه‌های میانی.
- این تنظیمات تضمین می‌کند که مدل بتواند الگوهای پیچیده زمانی را شناسایی کند.



```

# Transformer Model with Positional Encoding
class TransformerModel(nn.Module):
    def __init__(self, input_dim, d_model=64, nhead=4, num_layers=3, dim_feedforward=128, max_len=144):
        """
        Transformer model for single-step prediction with positional encoding.

        Args:
            input_dim (int): Number of input features per time step.
            d_model (int): Dimension of the embedding.
            nhead (int): Number of attention heads.
            num_layers (int): Number of transformer encoder layers.
            dim_feedforward (int): Dimension of the feedforward network.
            max_len (int): Maximum sequence length.
        """
        super(TransformerModel, self).__init__()
        self.embedding = nn.Linear(input_dim, d_model)
        self.pos_encoder = PositionalEncoding(d_model, max_len=max_len)
        encoder_layers = nn.TransformerEncoderLayer(d_model=d_model, nhead=nhead, dim_feedforward=dim_feedforward)
        self.transformer_encoder = nn.TransformerEncoder(encoder_layers, num_layers=num_layers)
        self.fc_out = nn.Linear(d_model, 1)

    def forward(self, x):
        """
        Forward pass of the Transformer model.

        Args:
            x (torch.Tensor): Input tensor of shape (batch_size, window_size, input_dim).

        Returns:
            torch.Tensor: Output predictions of shape (batch_size,).
        """
        x = self.embedding(x) # Shape: (batch_size, window_size, d_model)
        x = x.permute(1, 0, 2) # Shape: (window_size, batch_size, d_model)
        x = self.pos_encoder(x) # Add positional encoding
        x = self.transformer_encoder(x) # Shape: (window_size, batch_size, d_model)
        x = x[-1, :, :] # Take the last time step: (batch_size, d_model)
        return self.fc_out(x).squeeze() # Shape: (batch_size,)

```

شکل ۱۶- کد ساخت مدل transformer

مدل RNN مدل طراحی شده برای شناسایی وابستگی‌های ترتیبی در داده‌ها، با مشخصات زیر است:

- تعداد لایه‌ها: ۳
  - تعداد واحدهای مخفی (Hidden Units): ۱۲۸
  - نوع RNN: RNN پایه با tanh به عنوان تابع فعال‌سازی.
  - ورودی: داده‌ها با شکل (batch\_size, window\_size, input\_dim) وارد مدل شده و پس از پردازش، خروجی به شکل (batch\_size, hidden\_dim) تولید می‌شود.
  - لایه‌های Fully Connected: شامل یک لایه میانی با ابعاد ۶۴ واحد مخفی که تابع فعال‌سازی ReLU بر روی آن اعمال شده و یک لایه خروجی خطی برای پیش‌بینی نهایی.
- به‌علاوه، در مدل RNN جدید، با افزایش تعداد واحدهای مخفی و اضافه کردن یک لایه Fully Connected میانی، مدل توانایی یادگیری الگوهای پیچیده‌تری را به دست آورده است.

```

# RNN Model
class RNNModel(nn.Module):
    def __init__(self, input_dim, hidden_dim=128, num_layers=3, fc_hidden_dim=64):
        """
        Enhanced RNN model with additional layers and neurons for single-step prediction.

        Args:
            input_dim (int): Number of input features per time step.
            hidden_dim (int): Number of hidden units in the RNN layers.
            num_layers (int): Number of RNN layers.
            fc_hidden_dim (int): Number of hidden units in the additional fully connected layer.
        """
        super(RNNModel, self).__init__()
        self.rnn = nn.RNN(input_dim, hidden_dim, num_layers, batch_first=True, nonlinearity='tanh')
        self.fc_hidden = nn.Linear(hidden_dim, fc_hidden_dim) # Additional fully connected hidden layer
        self.fc_out = nn.Linear(fc_hidden_dim, 1) # Final output layer

    def forward(self, x):
        """
        Forward pass of the Enhanced RNN model.

        Args:
            x (torch.Tensor): Input tensor of shape (batch_size, window_size, input_dim).

        Returns:
            torch.Tensor: Output predictions of shape (batch_size,).
        """
        rnn_out, _ = self.rnn(x) # Shape: (batch_size, window_size, hidden_dim)
        rnn_out = rnn_out[:, -1, :] # Take the last time step: (batch_size, hidden_dim)
        fc_hidden_out = torch.relu(self.fc_hidden(rnn_out)) # Apply activation to hidden layer
        return self.fc_out(fc_hidden_out).squeeze() # Shape: (batch_size,)

```

شکل ۱۷ - کد ساخت مدل RNN

## مدل MLP

مدل MLP برای پردازش داده‌های ورودی مسطح طراحی شده و معماری آن به صورت زیر است:

- لایه‌های مخفی: دو لایه با اندازه‌های ۶۴ و ۳۲ واحد.
- تابع فعال‌سازی (Activation Function): ReLU
- لایه خروجی: یک لایه خطی برای تولید پیش‌بینی نهایی.

```

# MLP Model
class MLPModel(nn.Module):
    def __init__(self, input_dim, hidden_sizes=[64, 32]):
        """
        MLP model for single-step prediction.

        Args:
            input_dim (int): Number of input features (flattened).
            hidden_sizes (list): List containing the sizes of hidden layers.
        """
        super(MLPModel, self).__init__()
        layers = []
        last_size = input_dim
        for hidden_size in hidden_sizes:
            layers.append(nn.Linear(last_size, hidden_size))
            layers.append(nn.ReLU())
            last_size = hidden_size
        layers.append(nn.Linear(last_size, 1)) # Output layer
        self.network = nn.Sequential(*layers)

    def forward(self, x):
        """
        Forward pass of the MLP model.

        Args:
            x (torch.Tensor): Input tensor of shape (batch_size, window_size * input_dim).

        Returns:
            torch.Tensor: Output predictions of shape (batch_size,).
        """
        return self.network(x).squeeze()

```

شکل ۱۸- کد ساخت مدل MLP

## توابع loss

- **MSE**: برای محاسبه میانگین مربع اختلاف بین مقادیر پیش‌بینی شده و واقعی.
- **Huber Loss**: برای کاهش تاثیر داده‌های پرت با استفاده از ترکیبی از MSE و MAE.

## تنظیمات آموزش

### داده‌ها و پیش‌پردازش:

- داده‌ها به دو مجموعه آموزشی و آزمایشی تقسیم شده و با استفاده از scaler نرمال‌سازی می‌شوند.
- تعداد ویژگی‌های ورودی: ۸ ویژگی شامل:

Sequence No. ○

V ○

- D
- Air Density
- Humidity
- I
- S\_a
- S\_b

- استفاده از train\_loader برای مدیریت داده‌های ورودی به مدل‌ها.

### پارامترهای آموزش:

- **تعداد دوره‌ها (Epochs):** ۸۰ (ابتدا ۵۰ مشابه مقاله در نظر گرفتیم اما در نهایت دیدم با افزایش دوره عملکرد مدل بهبود میابد به همین جهت افزایش دادیم).
- **نرخ یادگیری (Learning Rate):** ۰.۰۰۱
- **بهینه‌ساز:** Adam برای تنظیم وزن‌ها در تمام مدل‌ها.

### حلقه آموزش:

#### ۱. مقدمه مدل:

- هر ترکیب مدل و loss با نمونه جدیدی از مدل شروع می‌شود.

#### ۲. (Forward Pass):

- داده‌های ورودی از مدل عبور کرده و پیش‌بینی‌ها محاسبه می‌شوند.

#### ۳. محاسبه loss:

- اختلاف بین مقادیر پیش‌بینی شده و واقعی با استفاده از تابع loss انتخابی.

#### ۴. (Backward Pass):

- گرادینان‌ها محاسبه شده و وزن‌ها با استفاده از بهینه‌ساز Adam به‌روزرسانی می‌شوند.

#### ۵. ثبت دوره‌ها:

- ثبت loss برای هر دوره و نمایش نتایج در فواصل منظم.

```

warnings.warn(
Epoch [1/80], Loss: 0.1466
Epoch [10/80], Loss: 0.0104
Epoch [20/80], Loss: 0.0037
Epoch [30/80], Loss: 0.0032
Epoch [40/80], Loss: 0.0024
Epoch [50/80], Loss: 0.0026
Epoch [60/80], Loss: 0.0022
Epoch [70/80], Loss: 0.0013
Epoch [80/80], Loss: 0.0012
Evaluation Metrics for Transformer with MSE loss:
RMSE: 7.1588, MAE: 5.5385, MAPE: 13.38%
Model Transformer Loss MSE RMSE 7.158771966148869 MAE 5.538475347485198 MAPE 13.38214482304401

Training Transformer with Huber loss
/usr/local/lib/python3.10/dist-packages/torch/nn/modules/transformer.py:379: UserWarning: enable_nested_tensor is not a valid argument to the Transformer constructor
warnings.warn(
Epoch [1/80], Loss: 0.1370
Epoch [10/80], Loss: 0.0036
Epoch [20/80], Loss: 0.0025
Epoch [30/80], Loss: 0.0015
Epoch [40/80], Loss: 0.0011
Epoch [50/80], Loss: 0.0009
Epoch [60/80], Loss: 0.0008
Epoch [70/80], Loss: 0.0006
Epoch [80/80], Loss: 0.0005
Evaluation Metrics for Transformer with Huber loss:
RMSE: 9.8435, MAE: 8.8134, MAPE: 21.55%
Model Transformer Loss Huber RMSE 9.843538685458585 MAE 8.813410066938285 MAPE 21.54500545509707

Training RNN with MSE loss
Epoch [1/80], Loss: 0.1672
Epoch [10/80], Loss: 0.0048
Epoch [20/80], Loss: 0.0036
Epoch [30/80], Loss: 0.0029
Epoch [40/80], Loss: 0.0023
Epoch [50/80], Loss: 0.0022
Epoch [60/80], Loss: 0.0019
Epoch [70/80], Loss: 0.0026
Epoch [80/80], Loss: 0.0019
Evaluation Metrics for RNN with MSE loss:
RMSE: 3.5387, MAE: 2.8800, MAPE: 6.98%
Model RNN Loss MSE RMSE 3.538713132070404 MAE 2.8800173946588483 MAPE 6.984972416590549

```

شکل ۱۹- نمونه خروجی فرایند آموزش

### ۱-۳-۵. بررسی و ارزیابی عملکرد مدل‌ها

در این قسمت، برنامه‌ای برای ارزیابی و مقایسه کارایی انواع مدل‌های توسعه یافته ارائه شده است. در این راستا، سه نوع مدل Transformer، RNN و MLP مورد استفاده قرار گرفتند. هر یک از این مدل‌ها با استفاده از دو نوع تابع MSE و Huber آزمون شده و متریک‌های RMSE، MAE و MAPE محاسبه شدند.

### روش فعالیت:

برای پیاده‌سازی و اجرای این برنامه، کارهای زیر انجام شد:

۱. تعریف و استخراج متریک‌ها: مقادیر واقعی و پیش‌بینی شده با استفاده از RMSE ، MAE و MAPE ارزیابی شدند. برای اطمینان از دقت، در محاسبه به موانع مثل تقسیم بر صفر توجه شد.

۲. تعریف مدل‌ها: در بخش قبل مفصل توضیح داده شد.

۳. روش آموزش:

- آموزش هر مدل به مدت ۸۰ دوره با استفاده از optimizer نوع Adam و سرعت یادگیری ۰.۰۰۱ انجام شد.

- مقادیر خطا هر ۱۰ دوره گزارش شد.

۴. ارزیابی نهایی:

- از داده‌های آزمون برای مقایسه مقادیر واقعی و پیش‌بینی شده استفاده شد. خروجی‌ها به مقیاس اولیه بازگردانده شدند.

- متریک‌ها روی مقیاس اولیه محاسبه شدند.

نتایج نهایی:

جداول به دو قالب Long و Wide (فرمت مقاله) ارائه شدند. قالب Wide خلاصه‌ای از نتایج هر مدل با استفاده از تابع از دست‌دهی مختلف ارائه می‌دهد. قالب Long به مقایسه‌ی جزئی‌تر بین متریک‌ها می‌پردازد.

Index	Model	MSE	Huber
MAE	MLP	9.6248	2.6196
	RNN	2.8800	4.8347
	Transformer	5.5385	8.8134
MAPE	MLP	24.0203	6.4974
	RNN	6.9850	11.8551
	Transformer	13.3821	21.5450
RMSE	MLP	10.2629	3.6604
	RNN	3.5387	5.5852
	Transformer	7.1588	9.8435

شکل ۲۰- جدول مقایسه پارامترهای بدست آمده برای هر مدل مشابه جدول ۳ مقاله

در ادامه به صورت اضافه مقایسه ای از نتایج بدست آمده و نتایج مقاله خواهیم داشت.

جدول ۱- مقایسه مقاله و نتایج ما در single step

index	Model	MSE	Huber	MSE	Huber
		نتایج ما		مقاله	
MAE	MLP	۹.۶۲۴۸	۲.۶۱۹۶	۳.۲۸۳۷	۳.۱۰۶۷
	RNN	۲.۸۸	۴.۸۳۴۷	۲.۹۲۲۵	۲.۹۰۱۳
	Transformer	۵.۵۳۸۵	۸.۸۱۳۴	۲.۳۶۷۰	۲.۲۴۸۵
MAPE	MLP	۲۴.۰۲۰۳	۶.۴۹۷۴	۴.۸۲	۴.۲۲
	RNN	۶.۹۸۵۰	۱۱.۸۵۵۱	۴.۲۵	۴.۲۲
	Transformer	۱۳.۳۸۲۱	۲۱.۵۴۵۰	۳.۶۳	۳.۲۲
RMSE	MLP	۱۰.۲۶۲۹	۳.۶۶۰۴	۴.۳۶۴۴	۴.۰۱۱۵
	RNN	۳.۵۳۸۷	۵.۵۸۵۲	۳.۶۱۵۱	۳.۵۸۹۱
	Transformer	۷.۱۵۸۸	۹.۸۴۳۵	۲.۸۸۱۷	۲.۷۷۶۴

در مقایسه نتایج مدل‌ها، MLP تحت Huber عملکرد خوبی در معیارهای MAE و RMSE داشت و حتی در RMSE توانست نتایجی نزدیک به مقاله ارائه دهد. مدل RNN نیز تحت MSE در معیارهای MAE و RMSE به نتایج مقاله نزدیک بود و نشان داد که می‌تواند به خوبی الگوهای موجود در داده‌ها را یاد بگیرد. اگرچه در برخی معیارها، به ویژه MAPE، عملکرد ضعیف‌تری نسبت به مقاله داشت، اما توانایی کلی آن در پیش‌بینی قابل توجه است. مدل Transformer با وجود تفاوت در نتایج، توانست در معیارهای MAE و RMSE تحت MSE به خوبی عمل کند و عملکرد نسبی قابل توجهی نشان داد. این نتایج نشان می‌دهد که مدل‌ها در یادگیری الگوهای زمانی موفق عمل کرده‌اند و امکان بهبود عملکرد آن‌ها با تنظیمات بهتر وجود دارد.



## بخش اضافی

با توجه به عملکرد نچندان عالی مدل در بخش قبل من تصمیم گرفتم تنها از  $y$  به عنوان فیچر ورودی استفاده کنم و عملکرد مدل را در حالتی که تنها به تاریخچه قدرت باد نگاه می‌کنیم، بررسی کنم. در این حالت ورودی ما تنها یک فیچر خواهد داشت و به جای ابعاد  $144 \times 8$  ابعاد  $144 \times 1$  خواهیم داشت. در تحلیل من، این نتیجه‌گیری که استفاده از تنها خروجی  $y$  به عنوان ورودی  $x$  می‌تواند نتایج بهتری نسبت به رویکرد مقاله که از ۸ ویژگی مختلف مانند "V"، "D"، "air density"، "Humidity"، و غیره استفاده کرده، داشته باشد، نشان‌دهنده اهمیت عمیق‌تر فهم روابط بین داده‌ها است.

دلایل احتمالی برای عملکرد بهتر مدل من می‌تواند شامل موارد زیر باشد:

۱. **سادگی داده‌های ورودی:** با استفاده از تنها  $y$ ، مدل من توانسته است به طور کامل روی الگوهای زمانی  $y$  تمرکز کند و از ورود نویز یا اطلاعات غیرمرتبط که ممکن است توسط سایر ویژگی‌ها ایجاد شود، جلوگیری کند. این سادگی به مدل اجازه داده که بهتر روی روابط مستقیم و غیرخطی بین زمان‌های گذشته و آینده تمرکز کند.

۲. **کیفیت پایین ویژگی‌های اضافی:** ممکن است ویژگی‌هایی که در مقاله استفاده شده‌اند، ارتباط قوی و مستقیمی با مقدار  $y$  نداشته باشند. در صورتی که این ویژگی‌ها نویزی یا کم‌کیفیت باشند، مدل قادر به استخراج الگوهای معنادار از آن‌ها نخواهد بود و در نتیجه عملکرد کاهش می‌یابد.

۳. **کاهش پیچیدگی و جلوگیری از بیش‌برازش:** استفاده از ورودی ساده‌تر یعنی فقط  $y$  می‌تواند باعث شود مدل، بیش‌برازش روی داده‌های آموزشی را کاهش دهد و در نتیجه عملکرد بهتری روی داده‌های تست داشته باشد. این موضوع به‌ویژه زمانی صادق است که تعداد نمونه‌های آموزشی کم باشد.

۴. **وابستگی‌های زمانی قوی در  $y$ :** ممکن است  $y$  به تنهایی الگوهای زمانی کافی برای پیش‌بینی دقیق گام‌های بعدی را در خود داشته باشد. به این ترتیب، مدل نیازی به اطلاعات اضافی از ویژگی‌های دیگر ندارد.

با این حال، باید به نکات زیر نیز توجه کرد:

- **ریسک overfit شدن به  $y$ :** اگر مدل تنها به  $y$  وابسته باشد، ممکن است عملکرد خوبی روی این دیتاست خاص داشته باشد، اما در سناریوهای دنیای واقعی یا داده‌های دیگر که اطلاعات محیطی بیشتری لازم است، عملکرد افت کند.

- نادیده گرفتن اطلاعات ارزشمند: ممکن است برخی ویژگی‌ها، مانند "air density" یا "humidity"، اطلاعات ارزشمندی را در بازه‌های زمانی مختلف ارائه دهند. حذف این اطلاعات می‌تواند در سناریوهای پیچیده‌تر محدودیت‌هایی ایجاد کند.

در نهایت، من معتقدم که نتایج بهتر مدل من نشان می‌دهد که داشتن تعداد زیادی ویژگی همیشه منجر به بهبود عملکرد نمی‌شود. در واقع، درک درست از روابط بین داده‌ها، کیفیت پیش‌پردازش و مهندسی ویژگی‌ها، و استفاده بهینه از مدل اهمیت بیشتری دارد. برای تایید این نتایج، انجام آزمایش‌های بیشتر روی دیتاست‌های متفاوت و بررسی تأثیر ویژگی‌های حذف‌شده ضروری به نظر می‌رسد. این تحلیل می‌تواند به درک بهتر نقش ویژگی‌های اضافی در پیش‌بینی توان بادی کمک کند. لازم به ذکر است که این بخش از پیاده‌سازی هم در پوشه Code تحت عنوان HW5\_1\_one\_feature ذخیره شده است. در ادامه نتایج بدست آمده در این روش گزارش شده است.

Index	Model	MSE	Huber
MAE	MLP	1.1864	1.3557
	RNN	0.6414	0.4282
	Transformer	1.4603	1.0332
MAPE	MLP	1.7475	2.1155
	RNN	0.9428	0.6600
	Transformer	2.0537	1.4817
RMSE	MLP	1.5324	1.6344
	RNN	0.7070	0.4803
	Transformer	1.5402	1.1785

شکل ۲۱ - جدول مقایسه پارامترهای بدست آمده برای هر مدل مشابه جدول ۳ مقاله (one feature)

جدول ۲- مقایسه مقاله و نتایج ما در (one feature) single step

index	Model	MSE	Huber	MSE	Huber
		نتایج ما		مقاله	
MAE	MLP	۱.۱۸۶۴	۱.۳۵۵۷	۳.۲۸۳۷	۳.۱۰۶۷
	RNN	۰.۶۴۱۴	۰.۴۲۸۲	۲.۹۲۲۵	۲.۹۰۱۳
	Transformer	۱.۴۶۰۳	۱.۰۳۳۲	۲.۳۶۷۰	۲.۲۴۸۵
MAPE	MLP	۱.۷۴۷۵	۲.۱۱۵۵	۴.۸۲	۴.۲۲
	RNN	۰.۹۴۲۸	۰.۶۶۰۰	۰.۰۴۲۵	۴.۲۲
	Transformer	۲.۰۵۳۷	۱.۴۸۱۷	۳.۶۳	۳.۲۲
RMSE	MLP	۱.۵۳۲۴	۱.۶۳۴۴	۴.۳۶۴۴	۴.۰۱۱۵
	RNN	۰.۷۰۷۰	۰.۴۸۰۳	۳.۶۱۵۱	۳.۵۸۹۱
	Transformer	۱.۵۴۰۲	۱.۱۷۸۵	۲.۸۸۱۷	۲.۷۷۶۴

در این مقایسه، عملکرد سه مدل **MLP**، **RNN** و **Transformer** بر اساس معیارهای **MAE**، **MAPE** و **RMSE** بررسی شده است. نتایج نشان می‌دهد پیاده‌سازی ما در تمام معیارها عملکرد بهتری نسبت به مقاله داشته است.

**MAE**: در مدل MLP، مقدار خطا از ۳.۱۰۶۷ (مقاله) به ۱.۳۵۵۷ (Huber) کاهش یافته است. مدل‌های RNN و Transformer نیز به ترتیب کاهش‌های قابل توجهی حدود ۸۵٪ و ۵۰٪ در MAE نشان داده‌اند.

**MAPE**: در مقاله مقادیر به صورت خام (غیر درصدی) گزارش شده‌اند. پس از تبدیل به درصد، برای مدل MLP از ۴.۲۲٪ به ۲.۱۱۵۵٪ کاهش یافت. در RNN و Transformer نیز مقادیر به ترتیب از ۴.۲۲٪ و ۳.۲۲٪ به ۰.۶۶۰۰٪ و ۱.۴۸۱۷٪ کاهش یافته‌اند.

**RMSE**: برای مدل MLP، خطا از ۴.۰۱۱۵ (مقاله) به ۱.۶۳۴۴ کاهش یافته است. مدل‌های RNN و Transformer نیز بیش از ۸۰٪ و ۵۰٪ بهبود در این معیار نشان داده‌اند.

نتایج نشان‌دهنده بهبود قابل توجه پیاده‌سازی ما در تمامی مدل‌ها و معیارها بوده و برتری روش‌های بهینه‌سازی ما را تایید می‌کند. عملکرد بهتر مدل احتمالاً به دلیل پیش‌پردازش بهتر داده‌ها مثل استفاده

از autoencoder برای کاهش نویز، تنظیم دقیق تر معماری مدل ها با انتخاب hyperparameter های مناسب است.

### ۳-۳-۶. مقادیر بهینه ابرپارامتر با استفاده از slime mould (امتیازی)

محتوای بخش امتیازی

در این پروژه، هدف استفاده از مدل **Transformer** برای پیش بینی سری های زمانی و بهینه سازی ابرپارامترهای آن با استفاده از الگوریتم **Slime Mould Algorithm (SMA)** است. مدل **Transformer** یکی از پیشرفته ترین مدل های یادگیری عمیق برای پردازش داده های توالی محور است که در این پروژه برای پیش بینی چندگامی به کار گرفته شده است.

#### معماری مدل **Transformer**

مدل **Transformer** شامل بخش های کلیدی زیر است:

- **Positional Encoding**: برای اضافه کردن اطلاعات زمانی به داده های ورودی
  - **Multi-head Attention**: برای یادگیری وابستگی های پیچیده در داده ها
  - **FeedForward Network**: برای پردازش ویژگی های استخراج شده
- این مدل به طور خاص برای پیش بینی چند مرحله ای طراحی شده است و می تواند چندین مقدار خروجی را به صورت همزمان پیش بینی کند.

#### ابرپارامترهای بهینه شده

ابرپارامترهای کلیدی که در این پروژه بهینه سازی شده اند، عبارتند از:

- **Nhead**: تعداد سرهای توجه
- **d\_model**: ابعاد مدل
- **num\_layers**: تعداد لایه ها
- **dim\_feedforward**: اندازه لایه **FeedForward**
- **dropout**: نرخ ریزش
- **learning rate**: نرخ یادگیری

## الگوریتم Slime Mould Algorithm (SMA)

الگوریتم SMA یک روش فراابتکاری الهام گرفته از رفتار کپک مخاطی است که برای بهینه‌سازی مسائل پیچیده استفاده می‌شود. در این پروژه، از این الگوریتم برای تنظیم مقادیر بهینه ابرپارامترهای مدل Transformer استفاده شده است.

### مراحل اجرای SMA

۱. ایجاد جمعیت اولیه: مقادیر اولیه ابرپارامترها به صورت تصادفی تولید می‌شوند.
۲. آموزش مدل: مدل Transformer با هر ترکیب از ابرپارامترها آموزش داده می‌شود.
۳. ارزیابی عملکرد: خطای اعتبارسنجی (Validation Loss) برای هر ترکیب محاسبه می‌شود.
۴. به‌روزرسانی جمعیت:
  - ترکیب‌ها براساس خطای اعتبارسنجی مرتب می‌شوند.
  - اعضای جمعیت براساس منطق الگوریتم SMA به سمت بهترین ترکیب حرکت می‌کنند.
  - مقادیر خارج از محدوده به مقادیر مجاز بازگردانده می‌شوند.
۵. تکرار مراحل تا همگرایی: فرآیند تا زمانی که به بهترین مقادیر برسیم، تکرار می‌شود.

```

class HyperparameterOptimizer:
    def __init__(self, train_loader, val_loader, input_dim: int, device: str = 'cuda'):
        self.train_loader = train_loader
        self.val_loader = val_loader
        self.input_dim = input_dim
        self.device = device

        # Modified bounds to ensure even numbers and proper divisibility
        self.param_bounds = {
            'd_model': (32, 256), # Will be adjusted to be divisible by nhead
            'nhead': (2, 8), # Must be power of 2
            'num_layers': (1, 6),
            'dim_feedforward': (64, 512),
            'dropout': (0.0, 0.5),
            'learning_rate': (1e-4, 1e-2)
        }

    def decode_solution(self, position: np.ndarray) -> Dict:
        # Get number of heads first
        nhead = 2 ** int(np.interp(position[1], [0, 1], [1, 3])) # This gives 2, 4, or 8

        # Get d_model and ensure it's divisible by nhead and even
        d_model = int(np.interp(position[0], [0, 1], self.param_bounds['d_model']))
        d_model = max(32, ((d_model // nhead) * nhead) // 2 * 2)

        return {
            'd_model': d_model,
            'nhead': nhead,
            'num_layers': int(np.interp(position[2], [0, 1], self.param_bounds['num_layers'])),
            'dim_feedforward': int(np.interp(position[3], [0, 1], self.param_bounds['dim_feedforward'])),
            'dropout': float(np.interp(position[4], [0, 1], self.param_bounds['dropout'])),
            'learning_rate': float(np.interp(position[5], [0, 1], self.param_bounds['learning_rate']))
        }

```

شکل ۲۰ - کد تعریف کلاس بهینه‌سازی هایپرپارامترها

## ویژگی‌های کد

- تنظیم خودکار ابرپارامترها: الگوریتم SMA بهترین ترکیب ابرپارامترها را با کمترین خطای اعتبارسنجی ارائه می‌دهد.
- توقف زودهنگام (Early Stopping): اگر خطای اعتبارسنجی بهبود نیابد، آموزش زودتر متوقف می‌شود.
- سازگاری مدل: اطمینان حاصل می‌شود که  $d\_model \setminus n\_head$  همواره بر  $n\_head$  بخش پذیر باشد.

## نتایج

۱. بهترین مقادیر برای ابرپارامترها، از جمله:

- $d\_model$
- $n\_head$

- num\_layers ○
- dim\_feedforward ○
- dropout ○
- learning rate ○

به همراه کمترین خطای اعتبارسنجی گزارش می‌شود.

۲. زمان کل بهینه‌سازی محاسبه شده و کارایی الگوریتم مشخص می‌شود.

۳. بهینه‌سازی نشان داده که استفاده از الگوریتم‌های فراابتکاری مانند SMA برای تنظیم ابرپارامترهای مدل‌های پیچیده‌ای مانند Transformer می‌تواند منجر به دستیابی به عملکرد بهینه شود.

استفاده از الگوریتم‌های فراابتکاری مانند SMA می‌تواند راه‌حلی کارآمد برای تنظیم ابرپارامترهای مدل‌های پیچیده باشد. مدل Transformer با ابرپارامترهای بهینه، توانست پیش‌بینی‌های چندمرحله‌ای با دقت بالا انجام دهد و الگوریتم SMA نقش مهمی در یافتن مقادیر بهینه داشت.

در نهایت با پیاده‌سازی نتیجه به فرم زیر در آمده و هایپرپارامترهای بهینه مدل اینگونه گزارش می‌شوند.

```
Iteration 1/5, Best Loss: 0.009315
Iteration 2/5, Best Loss: 0.009315
Iteration 3/5, Best Loss: 0.009315
Iteration 4/5, Best Loss: 0.009315
Iteration 5/5, Best Loss: 0.009315
```

```
Optimization completed!
Time taken: 2.11 minutes
```

```
Best Hyperparameters:
d_model: 164
nhead: 4
num_layers: 1
dim_feedforward: 151
dropout: 0.022613644455269033
learning_rate: 0.003320770274556317
Best Validation Loss: 0.009315
```

شکل ۲۱ - نتایج بهینه کردن هایپرپارامترهای مدل مبدل

### ۱-۳-۷. پیاده سازی مدل مبدل با استفاده از دو تابع loss برای حالت multi step

برای پیاده سازی این حالت معماری کلی مدل تغییری نمی کند (غیر از خروجی) بلکه چگونگی بکار گیری مدل و ارسال داده به مدل متفاوت می شود. تا پیش از این و برای حالت single step ما یک پنجره ۱۴۴ تایی از ورودی ها را به همراه  $y$  گام بعد به شبکه می دادیم و تا آموزش برای همان یک گام بعد از ۱۴۴ گام صورت گیرد و در بخش ارزیابی هم همین روند را پیش بردیم. اما در این بخش با تغییر پارامتر forecast\_step یا forecast\_horizon به جای عدد یک، یک لیست از اعداد شام ۴ و ۸ و ۱۶ را خواهیم داشت و این بدان معناست که در آموزش مدل ۱۴۴ گام داده می شود و برای  $y$  مقادیر ۴ گام بعد، ۸ گام بعد و ۱۶ گام بعد هم به مدل داده می شود تا الگوی این فرایند را آموزش ببیند. بدین صورت مدل توانایی پیشبینی در این گام ها را خواهد داشت.

در این رویکرد:

مدل پیش بینی ها را برای تمام افق های مشخص شده در یک مرحله محاسبه (forward pass) تولید می کند مثلاً یک بردار با اندازه ۳ برای  $t+4$ ،  $t+8$ ، و  $t+16$

برخلاف روش های پیش بینی تکراری، مدل پیش بینی های خود را برای پیش بینی مراحل بعدی به ورودی باز نمی گرداند، که این موضوع از انباشت خطا جلوگیری می کند.

هدف آموزشی تضمین می کند که مدل الگوها را در تمام افق های زمانی به صورت هم زمان یاد بگیرد و توانایی خود را برای پیش بینی چند مرحله ای بهینه کند.

این ساختار به مدل این امکان را می دهد که مستقیماً مقادیر مربوط به مراحل زمانی آینده را بر اساس پنجره ۱۴۴ مرحله ای ورودی پیش بینی کند.

- همچنین با توجه به اینکه هر گام در دیتاست معادل ۱۰ دقیقه است هر کدام از این گام های گفته شده  $t+4$ ،  $t+8$  و  $t+16$  به ترتیب معادل ۴۰ دقیقه، ۸۰ دقیقه و ۱۶۰ دقیقه بعد هستند. و در حالت تک گام تنها در حال پیش بینی ۱۰ دقیقه بعد بودیم.



```

class MultiStepTransformer(nn.Module):
    def __init__(self, input_dim, output_steps, d_model=64, nhead=4, num_layers=3, dim_feedforward=128, max_len=144):
        """
        Multi-step Transformer model for predicting multiple time steps ahead.

        Args:
            input_dim (int): Number of input features per time step
            output_steps (int): Number of time steps to predict
            d_model (int): Dimension of the model
            nhead (int): Number of attention heads
            num_layers (int): Number of transformer layers
            dim_feedforward (int): Dimension of feedforward network
            max_len (int): Maximum sequence length
        """
        super(MultiStepTransformer, self).__init__()
        self.d_model = d_model
        self.output_steps = output_steps
        # Input embedding
        self.embedding = nn.Linear(input_dim, d_model)
        self.pos_encoder = PositionalEncoding(d_model, max_len)
        # Transformer encoder
        encoder_layers = nn.TransformerEncoderLayer(d_model=d_model, nhead=nhead,
                                                    dim_feedforward=dim_feedforward)
        self.transformer_encoder = nn.TransformerEncoder(encoder_layers, num_layers=num_layers)

        # Output layer modified for multi-step prediction
        self.fc_out = nn.Linear(d_model, output_steps)
    def forward(self, x):
        """
        Forward pass of the model.

        Args:
            x (torch.Tensor): Input tensor of shape (batch_size, seq_len, input_dim)

        Returns:
            torch.Tensor: Predictions of shape (batch_size, output_steps)
        """
        # No need to unsqueeze if input_dim > 1
        x = self.embedding(x) # (batch_size, seq_len, d_model)
        x = x.permute(1, 0, 2) # (seq_len, batch_size, d_model)
        x = self.pos_encoder(x)
        x = self.transformer_encoder(x)
        x = x[-1, :, :] # Take the last sequence element
        return self.fc_out(x) # (batch_size, output_steps)

```

شکل ۲۲- کد ساخت مدل transformer برای multi step

```

# Split data into train and test
train_size = int(0.8 * len(features_scaled))
features_train, features_test = features_scaled[:train_size], features_scaled[train_size:]
target_train, target_test = target_scaled[:train_size], target_scaled[train_size:]

# Define parameters
window_size = 144
forecast_steps = [4, 8, 16]

# Create windows
X_train, y_train_dict = create_multistep_windows(features_train, target_train, window_size, forecast_steps)
X_test, y_test_dict = create_multistep_windows(features_test, target_test, window_size, forecast_steps)

print(f"X_train shape: {X_train.shape}") # Should be (num_samples, seq_len, num_features)
print(f"y_train_dict keys: {y_train_dict.keys()}")
print(f"y_train_dict[4] shape: {y_train_dict[4].shape}")

X_train shape: (993, 144, 8)
y_train_dict keys: dict_keys([4, 8, 16])
y_train_dict[4] shape: (993,)

```

شکل ۲۳ - آماده‌سازی مدل برای آموزش mutli step

```

history = train_multi_step_transformer(
    model_multi,
    train_loader,
    test_loader,
    target_scaler, # Pass the scaler for inverse transformation
    device,
    num_epochs=150,
    learning_rate=0.001
)

```

Epoch [1/150]  
Training Loss: 0.1196

Before inverse transform for horizon 4:  
Predictions range: 0.5752 to 0.7024  
Targets range: 0.6556 to 0.8828

After inverse transform for horizon 4:  
Predictions range: 59.4664 to 72.7043  
Targets range: 67.8323 to 91.4916

Before inverse transform for horizon 8:  
Predictions range: 0.6535 to 0.7231  
Targets range: 0.6081 to 0.8669

After inverse transform for horizon 8:  
Predictions range: 67.6173 to 74.8677  
Targets range: 62.8918 to 89.8348

Before inverse transform for horizon 16:  
Predictions range: 0.5164 to 0.5880  
Targets range: 0.5989 to 0.8291

After inverse transform for horizon 16:  
Predictions range: 53.3415 to 60.7904  
Targets range: 61.9333 to 85.9018

Before inverse transform for horizon 4:  
Predictions range: 0.5063 to 0.6563  
Targets range: 0.5989 to 0.8339

After inverse transform for horizon 4:  
Predictions range: 52.2817 to 67.9119  
Targets range: 61.9333 to 86.4015

#### شکل ۲۴ - بخشی از فرایند آموزش مدل **mutli step**

در نهایت مدل را بعد از آموزش با داده تست ارزیابی کرده و جدولی مشابه جدول ۵ مقاله ایجاد کردیم. البته چون در این بخش سوال الزامی به ایجاد جدول دقیقا مشابه مقاله و در کد نشده بود صرفا جدولی که داده‌ها را به صورت منظم گردآوری کند نمایش می‌دهیم.

Horizon	RMSE	MAE	MAPE
t+4	10.4042	8.7134	11.4797
t+8	14.2471	11.7196	15.7167
t+16	22.7925	19.5234	27.6430

شکل ۲۵ - نتایج مدل **transformer** در حالت **mutli step**

مدل در پیش‌بینی گام‌های کوتاه‌مدت بسیار موفق عمل کرده و دقت بالایی ارائه داده است. با این حال، با افزایش بازه پیش‌بینی، خطاها به طور قابل توجهی افزایش می‌یابند. این موضوع نشان می‌دهد که مدل نیاز به بهبود در شناسایی وابستگی‌های بلندمدت دارد. به صورت کلی روند کاهش دقت با افزایش زمان پیش‌بینی قابل مشاهده است که منطقی هم بنظر می‌رسد.

#### نتایج مقاله:

- $t+4$ : MAE=58.65, MAPE=0.70%, RMSE=91.80
- $t+8$ : MAE=78.50, MAPE=1.12%, RMSE=110.22
- $t+16$ : MAE=135.60, MAPE=3.11%, RMSE=166.21

#### مقایسه:

مدل من در کاهش خطاهای MAE و RMSE عملکرد بهتری داشته است. برای مثال، RMSE برای  $t+4$  در مدل من ۱۰.۴۰ و در مقاله ۹۱.۸۰ بوده که بهبود قابل توجهی است. اما MAPE در مقاله برای تمام گام‌ها پایین‌تر گزارش شده است؛ به عنوان مثال، MAPE برای  $t+4$  در مدل من ۱۱.۴۸٪ و در مقاله ۰.۷۰٪ بوده است. این نشان می‌دهد که در مقاله ممکن است استفاده از ورودی‌های متنوع‌تر به کاهش خطای نسبی کمک کرده باشد. البته به احتمال زیاد مقادیر مقاله درصدی گزارش نشده است.

## پرسش ۲. استفاده از ViT برای طبقه‌بندی تصاویر گلوبولهای سفید

### ۱-۲. آماده‌سازی داده‌ها

ابتدا از هر کلاس یک تصویر به صورت نمونه نمایش داده‌ایم.

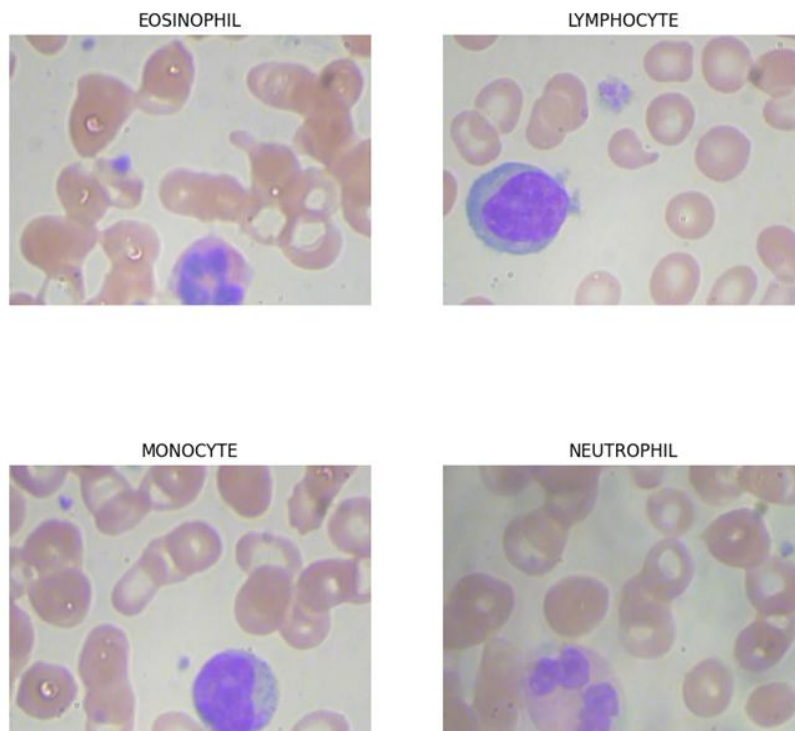
```
class_folders = ['EOSINOPHIL', 'LYMPHOCYTE', 'MONOCYTE', 'NEUTROPHIL']
class_counts = {}

plt.figure(figsize=(10, 10))
for i, class_folder in enumerate(class_folders):
    class_images = os.listdir(class_folder)
    class_counts[class_folder] = len(class_images)
    random_image_name = np.random.choice(class_images)
    random_image_path = os.path.join(class_folder, random_image_name)
    img = Image.open(random_image_path)

    plt.subplot(2, 2, i+1)
    plt.imshow(img)
    plt.title(class_folder)
    plt.axis('off')

plt.show()
```

شکل ۲۶- کد نمایش تصادفی یک تصویر از هر کلاس



شکل ۲۷- نمایش تصادفی یک تصویر از هر کلاس

سپس از تصاویر هر کلاس را شمرده‌ایم که در کلاس EOSINOPHIL تعداد ۸۸، در کلاس LYMPHOCYTE تعداد ۳۳، در کلاس MONOCYTE تعداد ۲۰ و در کلاس NEUTROPHIL تعداد ۲۰۶ تصویر قرار دارد.

```
for class_folder, count in class_counts.items():
    print(f"{class_folder}: {count} images")
```

```
EOSINOPHIL: 88 images
LYMPHOCYTE: 33 images
MONOCYTE: 20 images
NEUTROPHIL: 206 images
```

شکل ۲۸- تعداد تصاویر در هر کلاس

وقتی تعداد نمونه‌های موجود در هر کلاس تقریباً برابر باشد، مجموعه داده متوازن است. و وقتی تعداد نمونه‌های یک یا چند کلاس به طور قابل توجهی کمتر یا بیشتر از کلاس‌های دیگر باشد، مجموعه داده نامتوازن است. حال با توجه به تعریف و تعداد داده‌های هر کلاس به این نتیجه می‌رسیم که تعداد داده‌های موجود در دیتاست نامتوازن است و باید متوازن شود. تعداد تصاویر هر کلاس را به با تقویت داده به تعداد حداکثر یعنی همان ۲۰۶ تصویر می‌رسانیم.

```

max_count = max(class_counts.values())

augment_transforms = transforms.Compose([
    transforms.RandomHorizontalFlip(),
    transforms.RandomRotation(20),
    transforms.ColorJitter(brightness=0.2, contrast=0.2, saturation=0.2, hue=0.1),
    transforms.Resize((224, 224)),
    transforms.ToTensor(),
    transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]),
])

resized_normal_transforms = transforms.Compose([
    transforms.Resize((224, 224)),
    transforms.ToTensor(),
    transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]),
])

all_images = []
all_labels = []

for idx, cls in enumerate(class_folders):
    class_path = os.path.join(cls)
    images = [f for f in os.listdir(class_path) if os.path.isfile(os.path.join(class_path, f))]

    if len(images) < max_count:
        num_to_generate = max_count - len(images)
        print(f"Generating {num_to_generate} images for class {cls} ...")

        for image in images:
            img_path = os.path.join(class_path, image)
            img = Image.open(img_path).convert("RGB")
            augmented_img = resized_normal_transforms(img)
            all_images.append(augmented_img)
            all_labels.append(idx)

        for i in range(num_to_generate):
            img_name = images[i % len(images)]
            img_path = os.path.join(class_path, img_name)
            img = Image.open(img_path).convert("RGB")
            augmented_img = augment_transforms(img)
            all_images.append(augmented_img)
            all_labels.append(idx)
    else:
        print(f"No Generating images for class {cls}. It has {max_count} Images.")
        for image in images:
            img_path = os.path.join(class_path, image)
            img = Image.open(img_path).convert("RGB")
            augmented_img = resized_normal_transforms(img)
            all_images.append(augmented_img)
            all_labels.append(idx)

Generating 118 images for class EOSINOPHIL ...
Generating 173 images for class LYMPHOCYTE ...
Generating 186 images for class MONOCYTE ...
No Generating images for class NEUTROPHIL. It has 206 Images.

```

شکل ۲۹- متوازن کردن کلاس ها با تقویت داده

مدل‌های ViT نیاز دارند که تصاویر ورودی به ابعاد خاصی تنظیم شوند، زیرا این مدل‌ها تصاویر را به پچ‌های کوچک تقسیم می‌کنند. ابعاد تصویر معمولاً به صورت  $H \times W$  مشخص می‌شود، که  $H$  و  $W$  عرض و ارتفاع تصویر هستند. ابعاد استاندارد  $224 \times 224$  هستند. برای همین تمام تصاویر را به ابعاد گفته شده تبدیل کرده‌ایم.

```
print(f"Total images: {len(all_images)}")

train_data, val_data, train_labels, val_labels = train_test_split(all_images, all_labels, test_size=0.1, random_state=42)

print(f"Train images: {len(train_data)}")
print(f"Validation images: {len(val_data)}")
```

Total images: 824  
Train images: 741  
Validation images: 83

شکل ۳۰- تقسیم داده به دو قسمت آموزش و ارزیابی

## ۲-۲. آموزش مدل‌ها

مدل Vision Transformer (ViT) از معماری ترانسفورمر برای تحلیل داده‌های تصویری استفاده می‌کند. برخلاف مدل‌های کانولوشن عصبی (CNN)، این مدل تصاویر را به صورت توالی‌هایی از پچ‌ها تقسیم کرده و از مکانیزم توجه (Attention) برای پردازش آن‌ها استفاده می‌کند. ما می‌توانیم از PyTorch و کتابخانه transformers استفاده کنیم تا مدل ViT گوگل را بارگذاری کرده و آن را برای دسته‌بندی ۴ کلاس تنظیم کنیم



```

model = ViTForImageClassification.from_pretrained(
    "google/vit-base-patch16-224-in21k",
    num_labels=4, # number of classes
)
print(model)

```

Some weights of ViTForImageClassification were not initialized from the model checkpoint at google/vit-base-patch16-224-in21k and are newly initialized: You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

```

ViTForImageClassification(
  (vit): ViTModel(
    (embeddings): ViTEmbeddings(
      (patch_embeddings): ViTPatchEmbeddings(
        (projection): Conv2d(3, 768, kernel_size=(16, 16), stride=(16, 16))
      )
      (dropout): Dropout(p=0.0, inplace=False)
    )
    (encoder): ViTEncoder(
      (layer): ModuleList(
        (0-11): 12 x ViTLayer(
          (attention): ViTSdpaAttention(
            (attention): ViTSdpaSelfAttention(
              (query): Linear(in_features=768, out_features=768, bias=True)
              (key): Linear(in_features=768, out_features=768, bias=True)
              (value): Linear(in_features=768, out_features=768, bias=True)
              (dropout): Dropout(p=0.0, inplace=False)
            )
            (output): ViTSelfOutput(
              (dense): Linear(in_features=768, out_features=768, bias=True)
              (dropout): Dropout(p=0.0, inplace=False)
            )
          )
          (intermediate): ViTIntermediate(
            (dense): Linear(in_features=768, out_features=3072, bias=True)
            (intermediate_act_fn): GELUActivation()
          )
          (output): ViTOutput(
            (dense): Linear(in_features=3072, out_features=768, bias=True)
            (dropout): Dropout(p=0.0, inplace=False)
          )
          (layernorm_before): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
          (layernorm_after): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
        )
      )
    )
    (layernorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
  )
  (classifier): Linear(in_features=768, out_features=4, bias=True)
)

```

شکل ۳۱- بارگذاری مدل ViT و نمایش معماری آن

خروجی بالا معماری مدل Vision Transformer (ViT) گوگل برای دسته‌بندی تصویر را نشان می‌دهد. این مدل برای ۴ کلاس داده تنظیم شده و شامل اجزای زیر است:

- **Embedding:** تصاویر ورودی به صورت قطعات یا پچ‌های  $16 \times 16 \times 16$  برش داده می‌شوند. هر پچ به کمک یک کانولوشن (Conv2d) به بردارهایی با ۷۶۸ ویژگی تبدیل می‌شود.
- **Encoder:** شامل ۱۲ لایه Transformer است که هر کدام دارای موارد زیر هستند:
  - **Self-Attention Mechanism:** شامل لایه‌های query, key, و value برای محاسبه توجه
  - **Intermediate Layer:** یک لایه Fully Connected که ورودی را از ۷۶۸ به ۳۰۷۲ ویژگی گسترش داده و با تابع فعال‌سازی GELU اعمال می‌شود.
  - **Output Layer:** ورودی ۳۰۷۲ ویژگی را به ۷۶۸ کاهش می‌دهد.

○ Normalization Layers: هر لایه دارای Normalization قبل و بعد از عملیات پردازشی است.

- LayerNorm: یک لایه نرمال سازی نهایی با ۷۶۸ ویژگی اعمال می شود.
- Classifier: یک لایه Fully Connected نهایی که ویژگی های خروجی از Encoder را به ۴ کلاس دسته بندی می کند.

این معماری ViT برای یادگیری ویژگی های تصویری در حوزه بینایی ماشین بسیار قدرتمند است و با استفاده از مکانیزم توجه (Attention)، روی ویژگی های مهم تصویر تمرکز می کند.

در حالت اول فقط دسته بند Classifier قابل آموزش است. پس تمام لایه های مدل را فریز کردیم و فقط Classifier را آزاد گذاشتیم. تعداد کل پارامترها ۸۵۸۰۱۷۳۲ و تعداد پارامترهای قابل آموزش ۳۰۷۶ است.

```

model = ViTForImageClassification.from_pretrained("google/vit-base-patch16-224-in21k", num_labels=4)
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
model.to(device)
optimizer = Adam(model.parameters(), lr=0.0001)
loss_function = CrossEntropyLoss()

for param in model.parameters():
    param.requires_grad = False
for param in model.classifier.parameters():
    param.requires_grad = True

total_params = sum(p.numel() for p in model.parameters())
trainable_params = sum(p.numel() for p in model.parameters() if p.requires_grad)

print(f"Total parameters: {total_params}")
print(f"Trainable parameters: {trainable_params}")

train_losses_1 = []
val_losses_1 = []
train_accuracies_1 = []
val_accuracies_1 = []

for epoch in range(num_epochs):
    model.train()
    running_loss = 0.0
    correct = 0
    total = 0

    for inputs, labels in train_loader:
        inputs, labels = inputs.to(device), labels.to(device)
        optimizer.zero_grad()
        outputs = model(inputs).logits
        loss = loss_function(outputs, labels)
        loss.backward()
        optimizer.step()
        running_loss += loss.item()
        _, predicted = torch.max(outputs, 1)
        correct += (predicted == labels).sum().item()
        total += labels.size(0)

    epoch_train_loss = running_loss / len(train_loader)
    epoch_train_accuracy = 100 * correct / total
    train_losses_1.append(epoch_train_loss)
    train_accuracies_1.append(epoch_train_accuracy)

    model.eval()
    val_loss = 0.0
    correct = 0
    total = 0
    with torch.no_grad():
        for inputs, labels in val_loader:
            inputs, labels = inputs.to(device), labels.to(device)
            outputs = model(inputs).logits
            loss = loss_function(outputs, labels)
            val_loss += loss.item()
            _, predicted = torch.max(outputs, 1)
            correct += (predicted == labels).sum().item()
            total += labels.size(0)

    epoch_val_loss = val_loss / len(val_loader)
    epoch_val_accuracy = 100 * correct / total
    val_losses_1.append(epoch_val_loss)
    val_accuracies_1.append(epoch_val_accuracy)

    print(f"Epoch {epoch+1}/{num_epochs} - Train Loss: {epoch_train_loss:.4f} - Train Accuracy: {epoch_train_accuracy:.2f}% - Val Loss: {epoch_val_loss:.4f} - Val Accuracy: {epoch_val_accuracy:.2f}%")

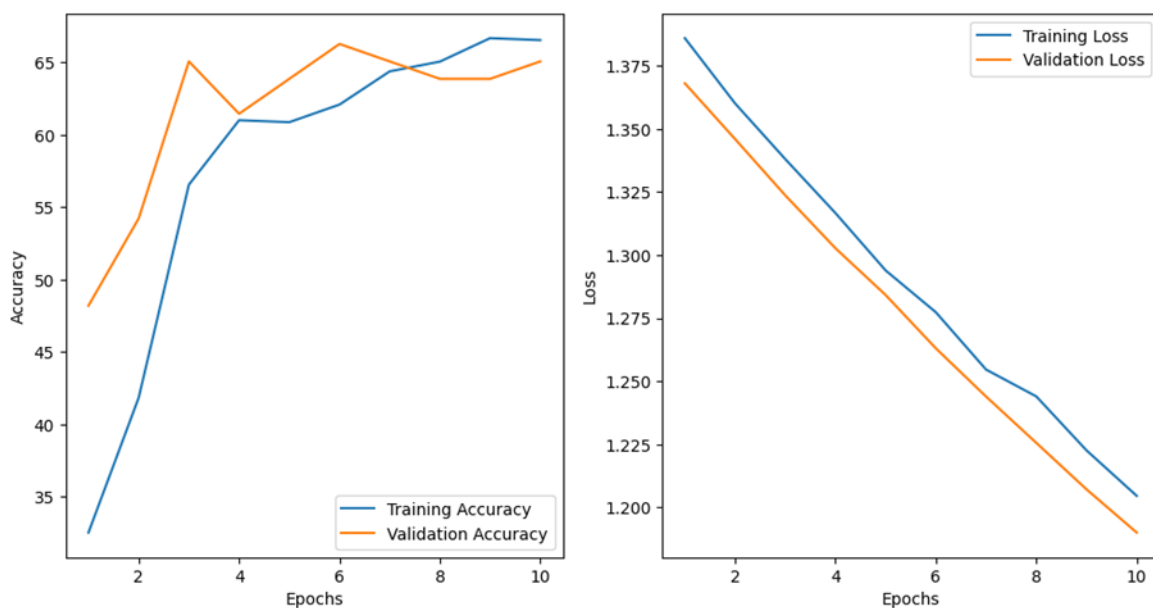
plt.figure(figsize=(12, 6))
plt.subplot(1, 2, 1)
plt.plot(range(1, num_epochs+1), train_accuracies_1, label='Training Accuracy')
plt.plot(range(1, num_epochs+1), val_accuracies_1, label='Validation Accuracy')
plt.xlabel('Epochs')
plt.ylabel('Accuracy')
plt.legend()

plt.subplot(1, 2, 2)
plt.plot(range(1, num_epochs+1), train_losses_1, label='Training Loss')
plt.plot(range(1, num_epochs+1), val_losses_1, label='Validation Loss')
plt.xlabel('Epochs')
plt.ylabel('Loss')
plt.legend()

plt.show()

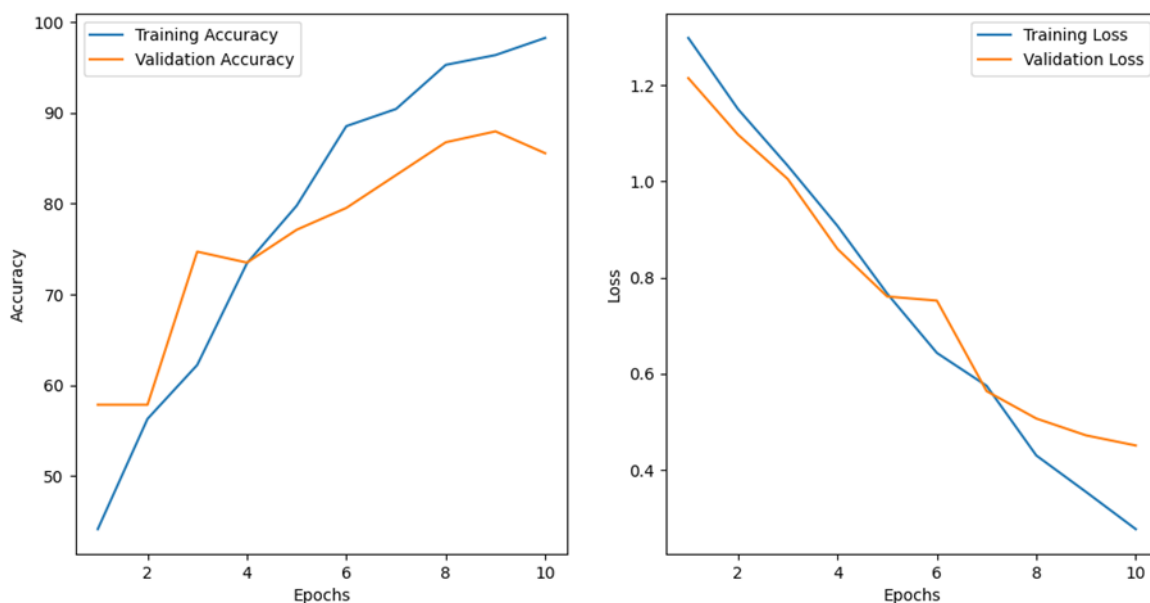
```

شکل ۳۲- فقط دسته-بند Classifier قابل آموزش باشد

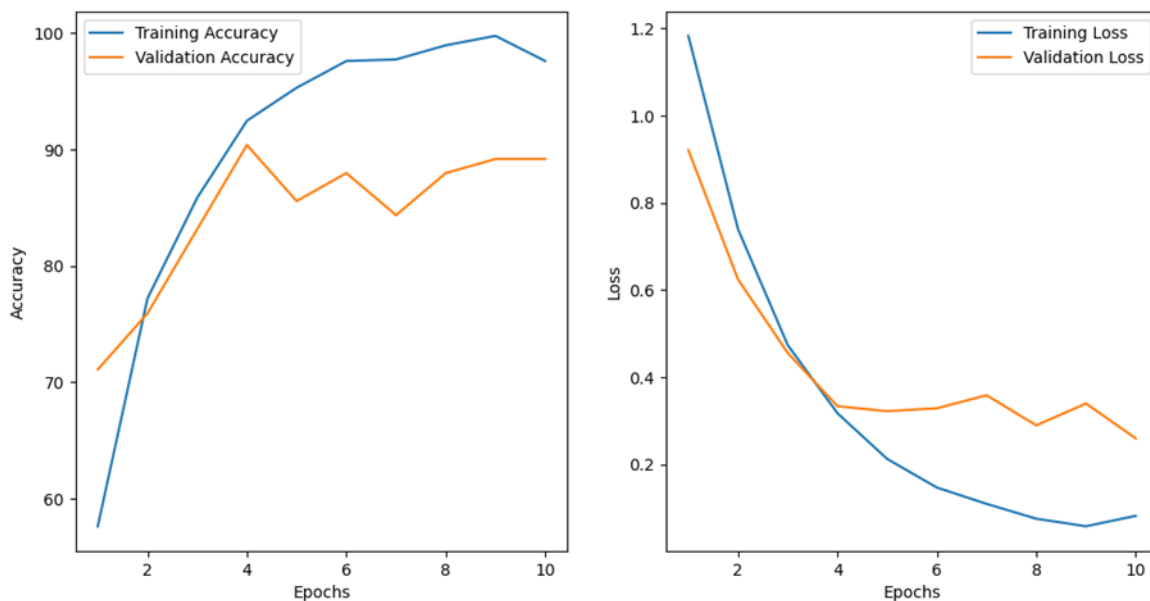


شکل ۳۳- نمودار **Accuracy** و **Loss** آموزش و ارزیابی برای حالت یک

در حالت دوم فقط دو لایه اول Encoder قابل آموزش است. لذا همه لایه‌ها را ابتدا فریز کردیم و سپس دو لایه اول Encoder و دسته‌بند را آزاد گذاشتیم. تعداد کل پارامترها ۸۵۸۰۱۷۳۲ و تعداد پارامترهای قابل آموزش ۱۴۹۲۱۴۷۶ است.

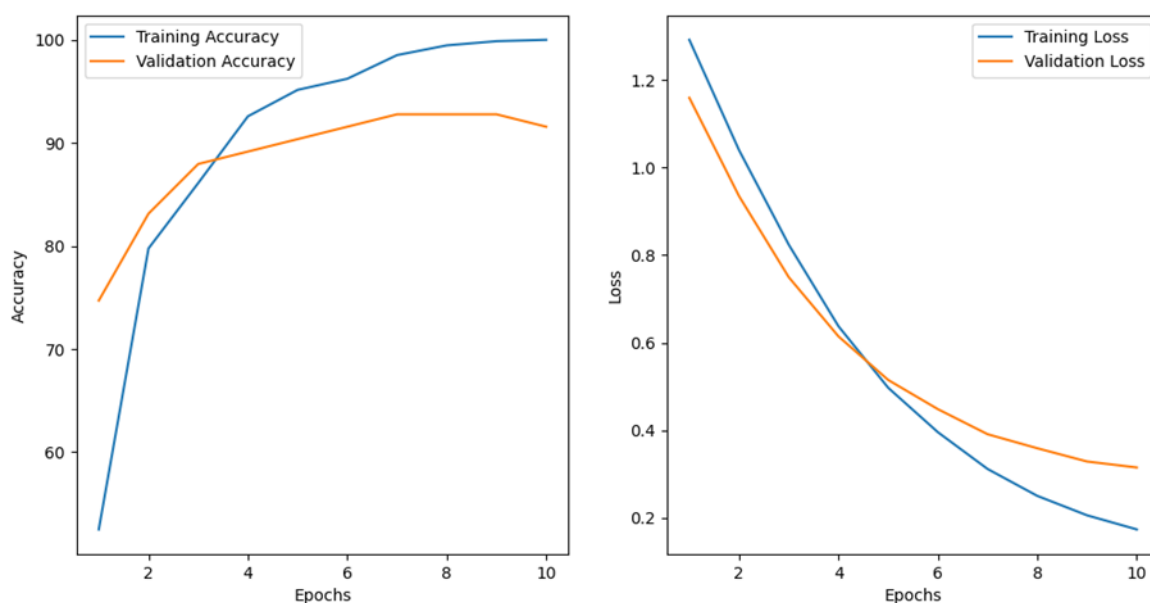


شکل ۳۴- نمودار **Accuracy** و **Loss** آموزش و ارزیابی برای حالت دو



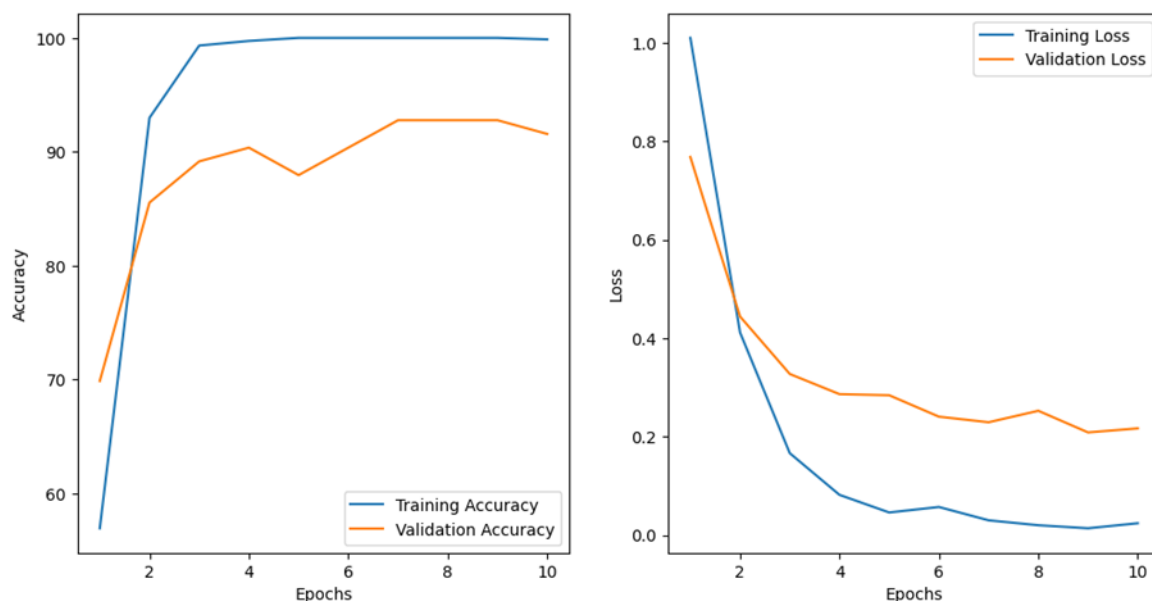
شکل ۳۵- نمودار **Accuracy** و **Loss** آموزش و ارزیابی برای حالت سه

در حالت چهارم تمام لایه‌ها قابل آموزش است. لذا همه لایه‌ها آزاد می‌کنیم و مدل را به طور کامل آموزش می‌دهیم. تعداد کل پارامترها ۸۵۸۰۱۷۳۲ و تعداد پارامترهای قابل آموزش ۸۵۸۰۱۷۳۲ است.



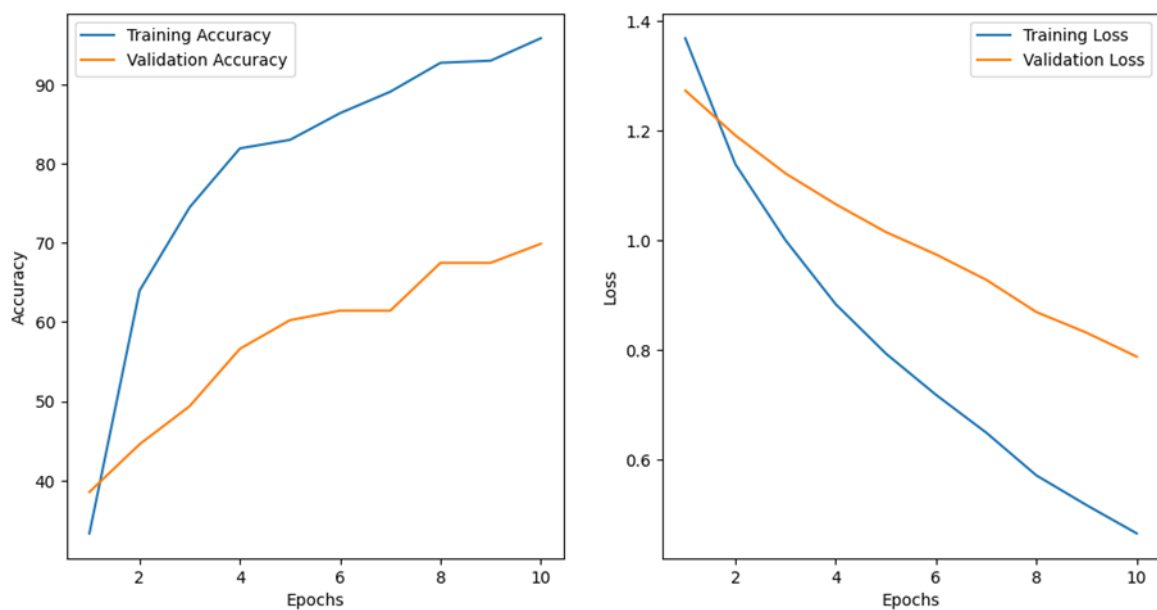
شکل ۳۶- نمودار **Accuracy** و **Loss** آموزش و ارزیابی برای حالت چهار

در حالت پنجم یک مدل CNN مثل مدل DenseNet-121 را بارگذاری کردیم و آن را به طور کامل آموزش دادیم. تعداد کل پارامترها ۶۹۵۷۹۵۶ و تعداد پارامترهای قابل آموزش ۶۹۵۷۹۵۶ است.



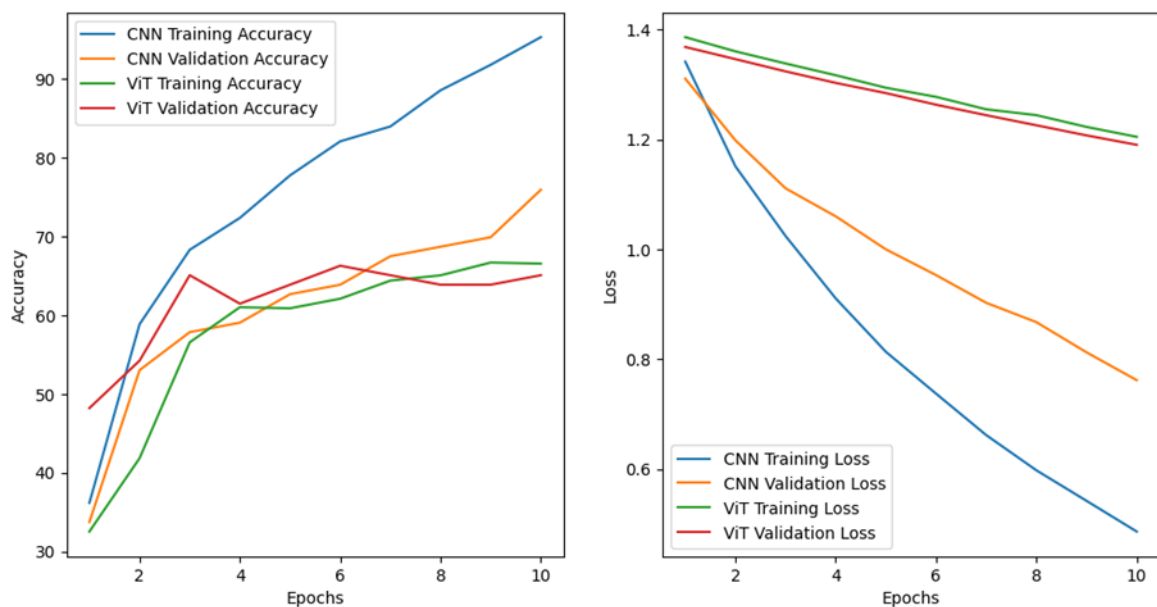
شکل ۳۷- نمودار Accuracy و Loss آموزش و ارزیابی برای حالت پنجم

مرحله امتیازی: در حالت ششم یک مدل CNN مثل DenseNet-121 را بارگذاری کردیم و فقط لایه Classifier مدل CNN را آموزش دادیم. تعداد کل پارامترها ۶۹۵۷۹۵۶ و تعداد پارامترهای قابل آموزش ۴۱۰۰ است.



شکل ۳۸- نمودار **Accuracy** و **Loss** برای مدل **CNN** آموزش فقط در لایه طبقه بند

در آخر نموداری برای مقایسه مدل **CNN** و **ViT** که فقط آموزش در لایه **Classification** انجام شده است ترسیم کردیم. با توجه به نمودار بدست آمده واضح است که مدل **CNN** بهتر آموزش دیده است.



شکل ۳۹- مقایسه مدل **CNN** و **ViT** با آموزش در لایه طبقه بند

CNN ها برای استخراج ویژگی‌های محلی (Local Features) طراحی شده‌اند. آنها به خوبی قادر به شناسایی الگوهای محلی در داده‌های تصویری هستند. در داده‌هایی که الگوهای مکانی خاصی دارند (مانند لبه‌ها، خطوط یا الگوهای تکرارشونده)، CNN ها معمولاً عملکرد خوبی نشان می‌دهند. در عوض مدل مبتنی بر مکانیسم توجه (Attention Mechanism) است و تمرکز بیشتری بر ویژگی‌های سراسری (Global Features) در تصویر دارد. ViT برای تحلیل تصاویر نیاز به داده‌های بیشتر و پیش‌پردازش مناسب دارد، چرا که ممکن است با داده‌های کوچک و تعداد محدود تصاویر، مدل دچار Overfitting شود و اینکه به خوبی ویژگی‌ها را استخراج نکند.

## ۲-۳. تحلیل و نتیجه گیری

در این تمرین ما دو مدل CNN و ViT را با حالت‌های متفاوت آموزش داده‌ایم. با توجه به این که تعداد تصاویر ما در هر کلاس ۲۰۶ مورد بود عملکرد در مدل شبکه عصب کانولوشنی DenseNet121 بهتر از سایر مدل‌ها بود. زیرا آموزش با این مدل در شرایطی که تعداد داده‌ها محدود است، عملکرد بهتری از خود نشان می‌دهد. و در این مدل به خوبی ویژگی‌های محلی را استخراج می‌شوند و برای تصاویر با ابعاد کوچک (۲۲۴\*۲۲۴) بسیار موثر است.

در مدل ViT به جای استفاده از لایه‌های کانولوشنی، از تکنیک‌های توجه (Attention) برای پردازش تصاویر استفاده می‌شود. ViT برای عملکرد مناسب معمولاً به داده‌های آموزشی بیشتری نیاز دارد. ViT به طور کلی برای عملکرد بهینه به حجم داده‌های زیادی نیاز دارد، زیرا مدل نیاز به یادگیری ویژگی‌های پیچیده‌تر و روابط غیرخطی بین بخش‌های مختلف تصویر دارد. در صورتی که تعداد تصاویر محدود است، مدل ViT نتوانسته است به خوبی ویژگی‌های تصویری را استخراج کند، چون به اندازه کافی داده برای آموزش تمام پارامترهایش را نداشته است.