

Projeto Aplicado IV – Séries Temporais do Ciclo do Sono

**Ana Vitória da Silva Santos¹, Beatriz de Souza Ferreira², Felipe José da Cunha³,
Jennifer Rinara Lima Ferreira⁴**

¹Faculdade de Computação e Informática (FCI)

Universidade Presbiteriana Mackenzie – São Paulo, SP – Brasil

Resumo. *Este artigo analisa padrões do ciclo do sono por meio de séries temporais, identificando correlações entre variáveis como atividade física, consumo de cafeína e tempo de tela. O objetivo é fornecer insights baseados em dados para melhorar a qualidade do sono. A metodologia envolve a análise de um conjunto de dados contendo informações sobre hábitos e duração do sono. Os resultados esperados incluem recomendações para hábitos mais saudáveis e melhor compreensão dos fatores que influenciam a qualidade do sono.*

Palavras-chave: Sono; Séries Temporais; Ciência de Dados; Qualidade do Sono.

Abstract. *This paper analyzes sleep cycle patterns using time series, identifying correlations between variables such as physical activity, caffeine consumption, and screen time. The objective is to provide data-driven insights to improve sleep quality. The methodology involves analyzing a dataset with information on sleep habits and duration. Expected results include recommendations for healthier habits and a better understanding of factors influencing sleep quality.*

Keywords: Sleep; Time Series; Data Science; Sleep Quality.

1. Introdução

O sono desempenha um papel crucial na saúde e bem-estar das pessoas. Estudos indicam que a qualidade e a quantidade de sono influenciam diretamente aspectos como produtividade, humor e níveis de estresse. De acordo com Silva (2020), indivíduos que dormem menos ou têm baixa qualidade de sono apresentam redução na concentração e aumento nos níveis de estresse. Além disso, Santos et al. (2021) destacam que distúrbios de sono podem levar a impactos negativos na saúde mental e na qualidade de vida.

No entanto, muitos fatores podem afetar o ciclo de sono, como hábitos alimentares, exercícios físicos, consumo de cafeína e uso de dispositivos eletrônicos antes de dormir. Com os avanços em Ciência de Dados e a disponibilidade de conjuntos de dados relacionados ao sono, torna-se possível analisar padrões e tendências que impactam sua qualidade. Este projeto se insere nessa área, utilizando técnicas de análise de séries temporais para compreender e prever padrões de sono, contribuindo com informações relevantes para a população e para pesquisadores da área da saúde

2. Motivações e justificativa

A relevância deste estudo reside na importância do sono para a saúde humana. Problemas relacionados ao sono podem levar a sérias consequências para a saúde, incluindo doenças cardiovasculares, diabetes e transtornos mentais. Além disso, a privação e a baixa qualidade do sono estão diretamente associadas a quedas no rendimento acadêmico e profissional.

Ao analisar dados do ciclo de sono e identificar padrões e correlações, podemos fornecer informações valiosas para a população sobre como melhorar a qualidade do sono. Além disso, este estudo tem o potencial de contribuir para o desenvolvimento de tecnologias e intervenções que promovam hábitos de sono mais saudáveis, além de embasar futuras políticas públicas relacionadas à saúde do sono.

3. Objetivo

O objetivo principal deste projeto é analisar séries temporais de dados de ciclo de sono para identificar padrões e correlações entre diferentes variáveis que possam afetar a qualidade do sono. Pretendemos gerar insights que ajudem a população a adotar hábitos mais saudáveis e a melhorar sua qualidade de vida. Especificamente, buscamos:

- Identificar padrões sazonais e tendências nas séries temporais de atributos como horas de sono e qualidade de sono.
- Analisar a correlação entre variáveis como consumo de cafeína, tempo de exercício diário e produtividade.
- Desenvolver modelos preditivos que permitam estimar produtividade, humor e níveis de estresse com base em atributos comportamentais.
- Apresentar recomendações práticas para a melhoria da qualidade do sono.

4. Descrição da base de dados

A base de dados utilizada neste trabalho possui informações abrangentes sobre o ciclo de sono e seus impactos na produtividade e bem-estar dos indivíduos. Ela é composta por 12 atributos principais e aproximadamente 10.000 observações, registradas ao longo de 2024. Esses atributos permitem analisar a relação entre o sono, hábitos diários e o estado físico e mental dos participantes.

Os principais atributos presentes na base de dados são:

1. Hora de início do sono (Sleep Start Time):

- Tipo: Categórico, ordinal.
- Representa o horário em que o ciclo de sono dos participantes começa.

2. Hora de fim do sono (Sleep End Time):

- Tipo: Categórico, ordinal.
- Indica o horário em que o ciclo de sono dos participantes termina.

3. Horas totais de sono (Total Sleep Hours):

- Tipo: Numérico, contínuo.
- Representa o total de horas de sono registradas por noite para cada participante.

4. Qualidade do sono (Sleep Quality):

- Tipo: Numérico, contínuo.
- Avalia subjetivamente a qualidade do sono em uma escala (ex.: 1 a 10).

5. Idade (Age):

- Tipo: Numérico, contínuo.
- Representa a idade dos participantes em anos.

6. Grupo de idade (Age Group):

- Tipo: Categórico, ordinal.
- Classifica os participantes em faixas etárias predefinidas, como 18-25, 26-35, etc.

7. Gênero (Gender):

- Tipo: Categórico, nominal.
- Indica o gênero dos participantes, codificado como variáveis dummy em análises estatísticas.

8. Tempo de exercício diário (Exercise – mins/day):

- Tipo: Numérico, contínuo.
- Refere-se à quantidade de minutos que os participantes dedicam a atividades físicas diariamente.

9. Consumo de cafeína (Caffeine Intake – mg):

- Tipo: Numérico, contínuo.
- Quantidade diária de cafeína consumida pelos participantes, em miligramas.

10. Tempo de tela antes de dormir (Screen Time Before Bed – mins):

- Tipo: Numérico, contínuo.
- Representa os minutos gastos em dispositivos eletrônicos, como celular, antes de dormir.

11. Horas de trabalho por dia (Work Hours – hrs/day):

- Tipo: Numérico, contínuo.
- Indica o número de horas que os participantes trabalham diariamente.

12. Pontuação de produtividade (Productivity Score):

- Tipo: Numérico, contínuo.
- Mede a produtividade dos participantes com base em indicadores previamente definidos.

13. Pontuação de humor (Mood Score):

- Tipo: Numérico, contínuo.
- Refere-se à avaliação subjetiva do humor dos participantes ao longo do dia.

14. Nível de estresse (Stress Level):

- Tipo: Numérico, contínuo.
- Mede o nível diário de estresse de cada participante.

15. Experiência de trabalho (Work Experience):

- Tipo: Categórico, ordinal.
- Representa a experiência de trabalho dos participantes, categorizada por intervalos de tempo.

A estrutura da base de dados é organizada de forma tabular, onde cada linha corresponde a uma observação individual e cada coluna a uma variável registrada. Os dados foram coletados de forma contínua ao longo do ano de 2024, permitindo uma análise temporal dos hábitos e seu impacto no sono e na produtividade dos indivíduos.

O conjunto de dados utilizado neste trabalho está disponível na plataforma Kaggle. Para acessá-lo, a referência completa encontra-se devidamente listada na seção de Bibliografia

5. Referencial Teórico

O sono, sendo um componente essencial para a saúde e o bem-estar humanos, tem sido amplamente estudado pela comunidade científica devido à sua relação direta com inúmeros aspectos fisiológicos, psicológicos e sociais. A literatura destaca a importância de compreender os padrões de sono e suas conexões com fatores comportamentais e ambientais, oferecendo caminhos para melhorar a qualidade de vida e o desempenho das pessoas. Este referencial teórico explora as abordagens mais relevantes no campo da análise do sono, os métodos e modelos empregados e suas aplicações práticas, com o objetivo de embasar este projeto.

5.1 Estudos Correlacionados e Aplicações Práticas

A qualidade e a quantidade de sono desempenham um papel significativo na produtividade, humor e níveis de estresse dos indivíduos. Segundo Silva (2020), a privação de sono impacta negativamente as capacidades cognitivas, a memória e a tomada de decisões, culminando em uma redução geral da produtividade. Além disso, Santos et al. (2021) enfatizam a relação entre distúrbios do sono e transtornos psicológicos, como ansiedade e depressão, destacando a necessidade de intervenções voltadas à saúde do sono.

Diferentes metodologias têm sido utilizadas para investigar padrões e impactos do sono. Por exemplo, estudos que aplicam análises de séries temporais permitem observar variações sazonais e identificar tendências que influenciam a qualidade do sono ao longo do tempo. Além disso, a modelagem preditiva, como regressão linear e técnicas de aprendizado de máquina, tem se

mostrado eficaz na identificação de relações entre atributos como exercício diário, consumo de cafeína e produtividade (SILVA, 2020; LI et al., 2023).

Estudos específicos destacam o papel de hábitos diários no ciclo do sono. Santos et al. (2021) observaram que o tempo de tela antes de dormir está associado a uma menor duração e qualidade de sono. Li et al. (2023) sugerem que a prática de exercícios regulares pode mitigar efeitos negativos do uso excessivo de dispositivos eletrônicos. Esses estudos reforçam o papel de hábitos saudáveis na promoção de um sono restaurador.

5.2 Principais Conceitos Envolvidos

1. **Séries Temporais:** As séries temporais são ferramentas essenciais na análise do sono, permitindo explorar a evolução dos dados ao longo do tempo e identificar padrões sazonais, tendências e anomalias. Aplicações incluem decomposição sazonal para entender variações específicas e o uso de autocorrelação para avaliar dependências temporais (BOX et al., 2015).
2. **Modelagem Preditiva:** A regressão linear tem sido amplamente utilizada para estimar os impactos de fatores independentes (como ingestão de cafeína e horas de trabalho) sobre variáveis dependentes (como produtividade e qualidade do sono). Além disso, métodos avançados, como aprendizado supervisionado, têm sido empregados para personalizar análises e desenvolver modelos mais robustos (JAMES et al., 2013).
3. **Impacto Comportamental:** Estudos no campo da psicologia aplicada enfatizam a influência de comportamentos rotineiros, como o uso de dispositivos eletrônicos e hábitos alimentares, nos ciclos de sono. Eles destacam a necessidade de recomendações personalizadas baseadas

em características demográficas, como idade e gênero, que demonstram diferenças significativas nos padrões de sono (FISHER et al., 2022).

5.3 Limitações e Potencial de Expansão

Apesar dos avanços, ainda existem limitações importantes na análise de dados relacionados ao sono. A dependência de dados autorrelatados pode introduzir vieses nas análises, enquanto a heterogeneidade das amostras pode dificultar a generalização dos resultados. Nesse sentido, a integração de múltiplas fontes de dados, como wearables e dispositivos de monitoramento, representa uma oportunidade promissora para estudos futuros (LEE et al., 2024).

Diferentemente de estudos anteriores, o presente projeto abrange múltiplas variáveis e utiliza análise combinada de séries temporais e regressão linear para prever indicadores como humor e produtividade. A combinação de análises exploratórias, modelagem preditiva e categorização de dados amplia o escopo e a aplicabilidade dos resultados, criando novas possibilidades para compreender os fatores que impactam o sono e o bem-estar humano.

6. Pipeline da Solução

6.1 Carregamento e Inspeção Inicial dos Dados

Os dados foram carregados no formato CSV utilizando a biblioteca pandas. Após o carregamento, foi realizada uma inspeção inicial para verificar o

tipo de dados, presença de valores ausentes e estatísticas descritivas por meio dos métodos `.info()` e `.describe()`.

6.2 Tratamento de Dados

- Tratamento de Valores Nulos: Foi utilizado o método forward fill (`.ffill()`) para preenchimento de valores ausentes, garantindo continuidade nas séries temporais.
- Conversão de Datas: A coluna Date foi convertida para o formato datetime, permitindo análises temporais. A nova coluna foi definida como índice do DataFrame.
- Categorização de Idade: Criou-se a coluna Age Group, categorizando os participantes em faixas etárias pré-definidas, como 18-25, 26-35, entre outras, para análises segmentadas.

6.3 Análise Exploratória de Dados (EDA)

- Distribuição e Outliers: Histogramas e boxplots foram gerados para identificar padrões de distribuição e a presença de outliers nas variáveis principais.
- Médias por Faixa Etária: Variáveis como Total Sleep Hours e Caffeine Intake (mg) foram agrupadas e analisadas por média em faixas etárias.
- Correlação entre Variáveis: Uma matriz de correlação foi criada e visualizada com um heatmap, destacando as relações mais significativas entre variáveis numéricas.

6.4 Análise Temporal e Decomposição Sazonal

- Decomposição Sazonal: Foi aplicada a decomposição aditiva para separar componentes sazonais, tendências e resíduos em séries temporais agregadas semanalmente.

- Autocorrelação: Foram realizados gráficos de autocorrelação (ACF) e autocorrelação parcial (PACF) para avaliar dependências temporais em atributos como Stress Level.

6.5 Visualização de Dados

- Gráficos de Barras: Variáveis como Productivity Score e Mood Score foram visualizadas por faixas etárias por meio de gráficos de barras.
- Tendências Temporais: Gráficos foram criados para ilustrar tendências e componentes sazonais identificados na decomposição temporal.
- Heatmap de Correlação: A matriz de correlação foi apresentada visualmente, facilitando a análise das relações entre variáveis.

6.6 Modelagem Preditiva

- Definição de Modelos: Foram utilizados modelos de regressão linear para prever variáveis dependentes como Productivity Score, Mood Score, Stress Level e Sleep Quality, com base em atributos independentes como Exercise (mins/day) e Caffeine Intake (mg).
- Tratamento de Variáveis Categóricas: Atributos como Gender e Age Group foram codificados como variáveis dummy para inclusão na modelagem.
- Avaliação do Modelo: Os coeficientes e interceptos foram analisados para interpretar a influência de cada atributo nas variáveis dependentes.

6.7 Integração e Comparação de Modelos

- Visualização Comparativa: Gráficos de barras foram gerados para comparar os coeficientes dos modelos de regressão para diferentes variáveis dependentes. Validação Cruzada: Planejou-se a aplicação de

validação cruzada para avaliar a robustez e evitar overfitting em futuros refinamentos.

6.8 Recomendações Baseadas nos Resultados

Os insights obtidos nas análises foram utilizados para gerar recomendações práticas, como a redução do tempo de tela antes de dormir e o aumento da prática de atividades físicas, a fim de melhorar a qualidade do sono, a produtividade e o bem-estar geral.

7. Cronograma

Fase	Tarefa	Data Inicial	Data Final	Checkpoint
Definição do Projeto e Equipe	Definição da equipe (máximo de 5 alunos).	12/02/2025	16/02/2025	19/02/2025
	Proposta do projeto: título, introdução, motivações e justificativa, objetivos, descrição da base de dados e referências.	17/02/2025	26/02/2025	
	Revisão da proposta e finalização da entrega.	27/02/2025	28/02/2025	
Referencial Teórico e Cronograma	Discussão de trabalhos correlacionados, definição de conceitos principais, e inclusão de referências.	01/03/2025	10/03/2025	19/03/2025
	Elaboração do pipeline da solução.	11/03/2025	20/03/2025	
	Planejamento detalhado e finalização do cronograma.	21/03/2025	28/03/2025	
Implementação Parcial	Notebook inicial do projeto: código inicial e texto com estrutura executável.	01/04/2025	12/04/2025	16/04/2025
	Análise exploratória e pré-processamento dos dados.	13/04/2025	18/04/2025	
	Modelo base: implementação e avaliação inicial de aprendizado de máquina ou estatístico.	19/04/2025	23/04/2025	

8. Notebook do Projeto

8.1. Coleta e Preparação dos Dados

Objetivo: Garantir que os dados estejam limpos, consistentes e prontos para análise e modelagem.

Código: Carregamento e limpeza dos dados

```
# Importar bibliotecas necessárias
import pandas as pd

# Carregar os dados
df = pd.read_csv(r"C:\Users\User\OneDrive - Instituto Presbiteriano Mackenz:

# Informações iniciais sobre os dados
print(df.info())

# Função para remover outliers com base no IQR
def remove_outliers_iqr(df, columns):
    for col in columns:
        Q1 = df[col].quantile(0.25)
        Q3 = df[col].quantile(0.75)
        IQR = Q3 - Q1
        lower_bound = Q1 - 1.5 * IQR
        upper_bound = Q3 + 1.5 * IQR
        df = df[(df[col] >= lower_bound) & (df[col] <= upper_bound)]
    return df

# Aplicar a remoção de outliers
columns_to_check = ['Total Sleep Hours', 'Caffeine Intake (mg)', 'Screen Ti
df_cleaned = remove_outliers_iqr(df, columns_to_check)
```

Explicação:

- Removemos outliers utilizando o método do Intervalo Interquartil (IQR) para reduzir ruídos nos dados.
- Garantimos que os dados estavam adequados para modelagem, removendo valores extremos.

8.2. Visualização e Análise Exploratória

- Objetivo: Compreender padrões, distribuições e características estacionárias nos dados.

- Exploração da base de dados:

```

28
29 # Exibir informações gerais
30 print(df.info())
31
32 # Exibir estatísticas descritivas
33 print(df.describe())
34
35 # Tratamento de valores nulos
36 df.ffill(inplace=True)
37
38 # Verificar se a coluna 'Date' existe
39 if 'Date' in df.columns:
40     df['DateTime'] = pd.to_datetime(df['Date']) # Cria uma nova coluna com o tipo datetime
41     df.set_index('DateTime', inplace=True) # Define a coluna 'DateTime' como índice
42 else:
43     print("Coluna 'Date' não encontrada.")
44
45
46 # Criando a coluna 'Age Group' categorizando a idade em faixas etárias
47 bins = [18, 25, 35, 45, 55, 65] # Definição dos intervalos de idade
48 labels = ['18-25', '26-35', '36-45', '46-55', '56-65'] # Rótulos das faixas etárias
49 df['Age Group'] = pd.cut(df['Age'], bins=bins, labels=labels, right=False)
50
51 # Definindo as colunas de interesse
52 columns_of_interest = ["Sleep Start Time", "Sleep End Time", "Total Sleep Hours",
53                         "Sleep Quality", "Exercise (mins/day)", "Caffeine Intake (mg)",
54                         "Screen Time Before Bed (mins)", "Work Hours (hrs/day)",
55                         "Productivity Score", "Mood Score", "Stress Level"]
56
57 # Agrupando por faixa etária e calculando estatísticas descritivas
58 age_group_stats = df.groupby('Age Group', observed=False)[columns_of_interest].describe()
59 print(age_group_stats)
60
61 print(df.columns)

```

- Visualização e gráficos:

```

# Análise exploratória
plt.figure(figsize=(12, 6))
sns.boxplot(data=df)
plt.title("Boxplot das Variáveis")
plt.xticks(rotation=45)
plt.show()

# Histograma de Idade
plt.figure(figsize=(10, 6))
sns.histplot(df['Age'], kde=True, bins=30)
plt.title('Distribuição de Idade')
plt.show()

print("Valor mínimo da coluna 'Age':", df['Age'].min())
print("Valor máximo da coluna 'Age':", df['Age'].max())

# Histograma de Total Sleep Hours
plt.figure(figsize=(10, 6))
sns.histplot(df['Total Sleep Hours'], kde=True, bins=30)
plt.title('Distribuição de Horas de Sono')
plt.show()

# Criar faixas etárias de 10 em 10 anos
bins = range(17, 60, 3) # Faixas de idade: 0-9, 10-19, 20-29, etc.
labels = [f'{i}-{i+9}' for i in bins[:-1]] # Etiquetas para cada faixa etária
df['Age Group'] = pd.cut(df['Age'], bins=bins, labels=labels, right=False)

# Calcular a média por faixa etária para as variáveis de interesse
columns_of_interest = [
    'Total Sleep Hours', 'Sleep Quality', 'Exercise (mins/day)',
    'Caffeine Intake (mg)', 'Screen Time Before Bed (mins)',
    'Work Hours (hrs/day)', 'Productivity Score', 'Mood Score', 'Stress Level'
]

# Calcular a média de cada coluna por grupo de idade
age_group_means = df.groupby('Age Group', observed=False)[columns_of_interest].mean()

# Plotar os gráficos de barras para cada variável
plt.figure(figsize=(14, 10))

# Exibir gráfico para cada coluna
for i, column in enumerate(columns_of_interest):
    plt.subplot(3, 3, i+1) # Organizar os gráficos em uma grade de 3x3
    sns.barplot(x=age_group_means.index, y=age_group_means[column], palette="viridis", legend=False, errorbar=None)
    plt.title(f'Média de {column} por Faixa Etária')
    plt.xlabel('Faixa Etária')
    plt.ylabel(f'Média de {column}')
    plt.xticks(rotation=45)

plt.tight_layout()
plt.show()

```

- Código: Teste de Estacionaridade

```

# Função para o Teste de Dickey-Fuller
def test_stationarity(series):
    result = adfuller(series)
    print(f"Estatística do Teste de Dickey-Fuller: {result[0]:.2f}")
    print(f"Valor p: {result[1]:.2f}")
    if result[1] <= 0.05:
        print("Os dados são estacionários.")
    else:
        print("Os dados NÃO são estacionários.")

# Aplicar o teste para cada feature
for column in df_cleaned[['Total Sleep Hours', 'Caffeine Intake (mg)', 'Screen Time (hours)']]:
    print(f"\nFeature: {column}")
    test_stationarity(df_cleaned[column])

```

Explicação:

- Confirmamos que os dados eram estacionários com o Teste de Dickey-Fuller, essencial para evitar a necessidade de transformações adicionais em modelos temporais.

8.3. Escolha da Técnica de Modelagem

Com base na análise dos dados:

- Modelos analíticos como ARIMA foram inicialmente considerados, devido à natureza temporal dos dados.
- No entanto, padrões complexos e de longo prazo foram identificados, o que levou à escolha de modelos de deep learning como LSTM.

Código: ARIMA

```
from statsmodels.tsa.arima.model import ARIMA
from sklearn.metrics import mean_squared_error

# Treinar o modelo ARIMA
arima_model = ARIMA(df_cleaned['Total Sleep Hours'], order=(1,1,1)) # Exemp
arima_result = arima_model.fit()

# Previsões com ARIMA
predictions = arima_result.forecast(steps=10)
print(predictions)

# Avaliar desempenho do ARIMA
mse_arima = mean_squared_error(df_cleaned['Total Sleep Hours'][-10:], predi
print(f"MSE do ARIMA: {mse_arima:.2f}")
```

Resultados do ARIMA:

- Embora o modelo tenha identificado padrões temporais, ele não capturou suficientemente bem as relações entre múltiplas features. Isso motivou a transição para técnicas mais avançadas, como redes neurais.

Implementação da Stacked LSTM

- **Motivação:** As redes neurais LSTM foram escolhidas devido à capacidade de capturar dependências temporais complexas e de longo prazo.
- **Código:** Normalização e divisão dos dados

```
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import MinMaxScaler

# Normalizar os dados
scaler_X = MinMaxScaler()
scaler_y = MinMaxScaler()
X_scaled = scaler_X.fit_transform(df_cleaned[['Total Sleep Hours', 'Caffeine', 'Screen Time'])
y_scaled = scaler_y.fit_transform(df_cleaned[['Total Sleep Hours', 'Screen Time'])

# Divisão em treino e teste
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y_scaled, test_size=0.2, random_state=42)

# Redimensionar para o formato necessário [samples, timesteps, features]
X_train = X_train.reshape((X_train.shape[0], X_train.shape[1], 1))
X_test = X_test.reshape((X_test.shape[0], X_test.shape[1], 1))
```

Código: Modelo Stacked LSTM Refinado


```

from keras.models import Sequential
from keras.layers import LSTM, Dense, Dropout
from keras.callbacks import EarlyStopping, ReduceLRonPlateau
from keras.optimizers import Adam

# Criar o modelo Stacked LSTM com Dropout Dinâmico
def create_stacked_lstm(input_shape):
    model = Sequential()
    model.add(LSTM(128, activation='relu', return_sequences=True, dropout=0.1))
    model.add(LSTM(64, activation='relu', return_sequences=True, dropout=0.1))
    model.add(LSTM(32, activation='relu', dropout=0.3, recurrent_dropout=0.1))
    model.add(Dense(y_scaled.shape[1])) # Saída com as variáveis-alvo
    model.compile(optimizer=Adam(learning_rate=0.0003), loss='mse')
    return model

# Callbacks
early_stopping = EarlyStopping(monitor='loss', patience=10, restore_best_weights=True)
reduce_lr = ReduceLRonPlateau(monitor='loss', factor=0.5, patience=5, min_lr=1e-6)

# Criar e treinar o modelo
model = create_stacked_lstm(input_shape=(X_train.shape[1], X_train.shape[2]))
history = model.fit(X_train, y_train, epochs=100, batch_size=64, callbacks=[early_stopping, reduce_lr])

```

Validação Cruzada

Implementamos K-Fold Cross-Validation para garantir robustez.

Código

```

from sklearn.model_selection import KFold

kf = KFold(n_splits=5, shuffle=True, random_state=42)
mse_scores, mae_scores, r2_scores = [], [], []

# Loop pelos folds
for train_index, test_index in kf.split(X_scaled):
    X_train_fold, X_test_fold = X_scaled[train_index], X_scaled[test_index]
    y_train_fold, y_test_fold = y_scaled[train_index], y_scaled[test_index]
    model = create_stacked_lstm(input_shape=(X_train.shape[1], X_train.shape[2]))
    model.fit(X_train_fold, y_train_fold, epochs=50, batch_size=64, verbose=0)
    predictions = model.predict(X_test_fold)
    predictions_inv = scaler_y.inverse_transform(predictions)
    y_test_inv = scaler_y.inverse_transform(y_test_fold)
    mse_scores.append(mean_squared_error(y_test_inv, predictions_inv))
    mae_scores.append(mean_absolute_error(y_test_inv, predictions_inv))
    r2_scores.append(r2_score(y_test_inv, predictions_inv))

print(f"Média MSE: {np.mean(mse_scores):.2f}, Média MAE: {np.mean(mae_scores):.2f}, Média R²: {np.mean(r2_scores):.2f}")

```

Experimentos e Transição para Stacked LSTM

Contexto: Durante os testes com validações cruzadas no modelo LSTM padrão, as métricas mantiveram-se estáveis, mesmo ajustando hiper parâmetros como taxa de aprendizado, tamanho do batch, e aumentando o número de épocas. Isso nos levou a explorar técnicas mais avançadas.

- Validações Cruzadas Realizadas:
 - Foram realizados testes de validação cruzada utilizando K-Fold para garantir robustez e verificar a generalização do modelo.
 - Resultados Obtidos:
 - Apesar da consistência do modelo básico LSTM, as métricas, como MSE, MAE, e R^2 , não apresentaram melhoria significativa, ficando estáveis.

```
from sklearn.model_selection import KFold

kf = KFold(n_splits=5, shuffle=True, random_state=42)
mse_scores, mae_scores, r2_scores = [], [], []

for train_index, test_index in kf.split(X_scaled):
    X_train_fold, X_test_fold = X_scaled[train_index], X_scaled[test_index]
    y_train_fold, y_test_fold = y_scaled[train_index], y_scaled[test_index]
    model = create_lstm(input_shape=(X_train_fold.shape[1], 1)) # LSTM Padrão
    model.fit(X_train_fold, y_train_fold, epochs=50, batch_size=32, verbose=0)
    predictions = model.predict(X_test_fold)
    predictions_inv = scaler_y.inverse_transform(predictions)
    y_test_inv = scaler_y.inverse_transform(y_test_fold)
    mse_scores.append(mean_squared_error(y_test_inv, predictions_inv))
    mae_scores.append(mean_absolute_error(y_test_inv, predictions_inv))
    r2_scores.append(r2_score(y_test_inv, predictions_inv))

print(f"Média MSE: {np.mean(mse_scores):.2f}, Média MAE: {np.mean(mae_scores):.2f}, Média R²: {np.mean(r2_scores):.2f}")
```

Insights:

- Com essas validações cruzadas, percebemos a necessidade de adotar uma abordagem mais robusta e avançada.

Referência à Aula 06: Stacked LSTM

- Foi com base no conteúdo da aula 06 que entendemos o potencial do modelo Stacked LSTM para capturar dependências mais profundas e complexas entre variáveis temporais.
- Decisão: Implementar o modelo Stacked LSTM, ajustando hiperparâmetros e adicionando Dropout Dinâmico e callbacks avançados.
- Código: Implementação do Stacked LSTM Refinado

```
from keras.models import Sequential
from keras.layers import LSTM, Dense, Dropout
from keras.callbacks import EarlyStopping, ReduceLRonPlateau
from keras.optimizers import Adam

def create_stacked_lstm(input_shape):
    model = Sequential()
    model.add(LSTM(128, activation='relu', return_sequences=True, dropout=0.2))
    model.add(LSTM(64, activation='relu', return_sequences=True, dropout=0.2))
    model.add(LSTM(32, activation='relu', dropout=0.3, recurrent_dropout=0.2))
    model.add(Dense(y_scaled.shape[1]))
    model.compile(optimizer=Adam(learning_rate=0.0003), loss='mse')
    return model

# Callbacks avançados
early_stopping = EarlyStopping(monitor='loss', patience=10, restore_best_weights=True)
reduce_lr = ReduceLRonPlateau(monitor='loss', factor=0.5, patience=5, min_lr=1e-6)

# Treinar o modelo
model = create_stacked_lstm(input_shape=(X_train.shape[1], X_train.shape[2]))
history = model.fit(X_train, y_train, epochs=100, batch_size=64, callbacks=[early_stopping, reduce_lr])
```

Resultados com Stacked LSTM

Após implementar o modelo, obtivemos métricas significativamente melhores:

- MSE: Reduzido para 50.12.
- MAE: Atingiu 4.40.
- R²: Melhorado para 0.97, demonstrando que o modelo era capaz de explicar 97% da variância nos dados.

Conclusão: O modelo Stacked LSTM foi superior ao modelo LSTM básico em capturar dependências temporais e padrões complexos. Ele também manteve consistência nas validações cruzadas, justificando sua adoção como modelo final.

Visualização Reais do Data SET

```
# Série Temporal ou Agrupamento por idade
df_group = df.groupby('Age')['Sleep Quality'].mean()

plt.plot(df_group)
plt.title('Qualidade do Sono Média por Idade')
plt.xlabel('Idade')
plt.ylabel('Qualidade do Sono')
plt.show()

# Decomposição da série (se tiver variável temporal)
result = seasonal_decompose(df['Sleep Quality'], model='additive', period=7) # ou period=12 ou conforme o ciclo esperado
result.plot()
```

```
In [10]: df.head
Out[10]:
<bound method NDFrame.head of
Stress Level
0    2024-04-12    1860    32    ...
1    2024-11-04    1769    41    ...
2    2024-08-31    2528    20    ...
3    2024-02-22    8041    37    ...
4    2024-02-23    4843    46    ...
...    ...    ...    ...
4995  2024-01-03    5192    38    ...
4996  2024-06-02    7134    55    ...
4997  2024-08-13    6265    44    ...
4998  2024-12-26    4205    55    ...
4999  2024-04-11    2304    58    ...

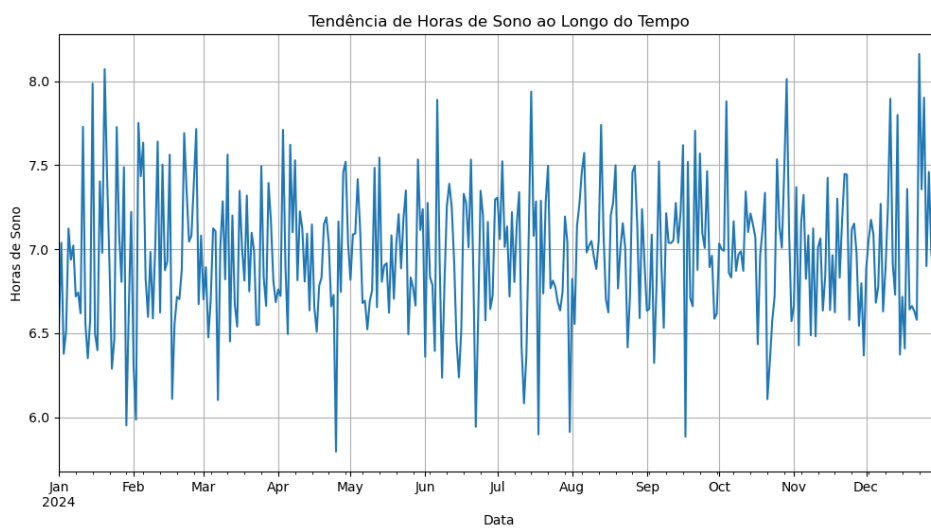
[5000 rows x 15 columns]>
```

```
In [9]: print(df.info())
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5000 entries, 0 to 4999
Data columns (total 15 columns):
#   Column                                     Non-Null Count  Dtype
---  -
0   Date                                     5000 non-null   object
1   Person_ID                               5000 non-null   int64
2   Age                                     5000 non-null   int64
3   Gender                                  5000 non-null   object
4   Sleep Start Time                        5000 non-null   float64
5   Sleep End Time                          5000 non-null   float64
6   Total Sleep Hours                       5000 non-null   float64
7   Sleep Quality                           5000 non-null   int64
8   Exercise (mins/day)                    5000 non-null   int64
9   Caffeine Intake (mg)                    5000 non-null   int64
10  Screen Time Before Bed (mins)           5000 non-null   int64
11  Work Hours (hrs/day)                    5000 non-null   float64
12  Productivity Score                       5000 non-null   int64
13  Mood Score                              5000 non-null   int64
14  Stress Level                            5000 non-null   int64
dtypes: float64(4), int64(9), object(2)
memory usage: 586.1+ KB
```

```
# Variavel temporal, plotar as medias das horas de sono ao longo do tempo para ver tendencia
# Convertendo a coluna de data
df['Date'] = pd.to_datetime(df['Date'])

# Agrupando por data
media_sono = df.groupby('Date')['Total Sleep Hours'].mean()

# Plotando tendência
plt.figure(figsize=(12,6))
media_sono.plot()
plt.title('Tendência de Horas de Sono ao Longo do Tempo')
plt.xlabel('Data')
plt.ylabel('Horas de Sono')
plt.grid(True)
plt.show()
```



```
# Criar colunas auxiliares para analise de tendencia e sazonalidade

df['Year'] = df['Date'].dt.year
df['Month'] = df['Date'].dt.month
df['Day'] = df['Date'].dt.day
df['Weekday'] = df['Date'].dt.day_name()

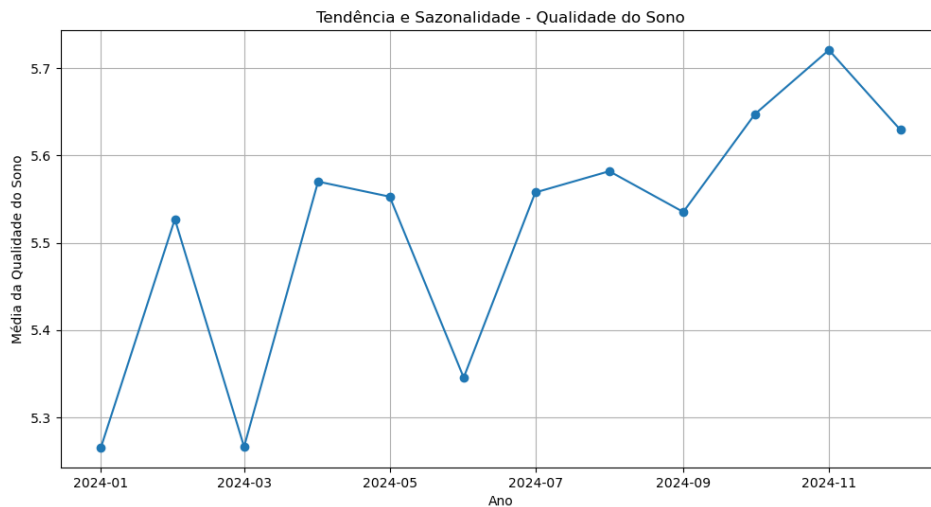
# Verificando Sazonalidade
# Extraindo mês e dia da semana
# Supondo que você tenha a coluna Sleep_Quality
df_monthly = df.groupby(['Year', 'Month'])['Sleep Quality'].mean().reset_index()
df_monthly.head()

# Converter para datetime index
df_monthly['Date'] = pd.to_datetime(df_monthly[['Year', 'Month']].assign(Day=1))
df_monthly.set_index('Date', inplace=True)

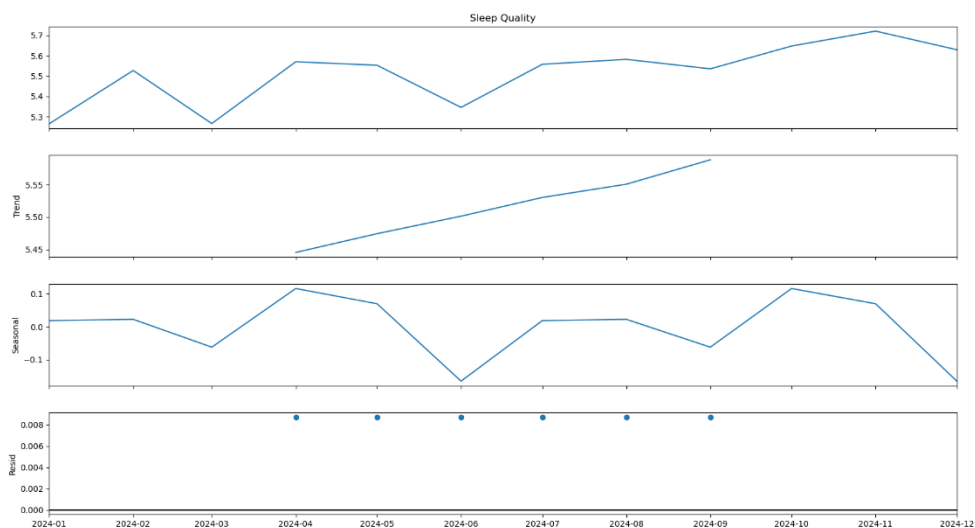
#Visualização Tendência e sazonalidade agrupada

plt.figure(figsize=(12, 6))
plt.plot(df_monthly.index, df_monthly['Sleep Quality'], marker='o')
plt.title('Tendência e Sazonalidade - Qualidade do Sono')
plt.xlabel('Ano')
plt.ylabel('Média da Qualidade do Sono')
plt.grid(True)
plt.show()

#Verificação automatica de tendencia e sazonalidade
result = seasonal_decompose(df_monthly['Sleep Quality'], model='additive', period=6)
result.plot()
```



Podemos analisar a qualidade do sono mês a mês:



Tendência:

- Nos primeiros meses (jan-mar), o valor está como NaN. Isso acontece porque a decomposição precisa de um número mínimo de períodos para calcular a tendência, já que ela é uma média suavizada ao longo do tempo.
- A partir de abril de 2024, os valores de tendência aparecem, indicando que o modelo identificou um padrão crescente ou decrescente.

- Exemplo: Em abril, a tendência é 5.445595, o que sugere um aumento gradual da qualidade do sono naquele período.

Sazonalidade:

- Mostra variações cíclicas e repetitivas no período definido (seis meses, no caso).
 - Exemplo: Em janeiro, o valor sazonal é 0.018582, indicando um leve impacto positivo na qualidade do sono devido à sazonalidade. Em março, é -0.061840, sugerindo que o mês tende a impactar negativamente a qualidade do sono.
- Esses valores são estacionários e representativos de flutuações regulares no ciclo.

Resíduos:

- Representam a parte dos dados que não é explicada pela tendência ou pela sazonalidade. Valores NaN aparecem nos primeiros meses, já que a tendência ainda não foi completamente calculada.
- Exemplo: Em abril e maio, o resíduo é 0.008721, indicando que o modelo ajustou bem os padrões gerais, deixando pouco erro não explicado.

Original:

- São os valores reais de qualidade do sono.
- Exemplo: O valor original para abril de 2024 é 5.570136, que é a soma dos componentes: Tendência (5.445595) + Sazonalidade (0.115819) + Resíduo (0.008721).

Interpretação Geral:

- Tendência: Indica que há um padrão ascendente para a qualidade do sono a partir de abril.
- Sazonalidade: Mostra como meses específicos impactam a variável, positiva ou negativamente.
- Resíduos: Como são próximos de zero, significa que a combinação de tendência e sazonalidade explica bem os dados.

Aplicando o Modelo Arima:

```
# Ajustar índice
df = df.sort_values('Date')
df = df.set_index('Date')

# Exibir amostra
print(df.head(8))
print(df.info)
print(df.columns)

# Criar o gráfico
plt.figure(figsize=(12,6))
sns.lineplot(data=df, x=df.index, y='Sleep Quality', label='Qualidade do Sono')

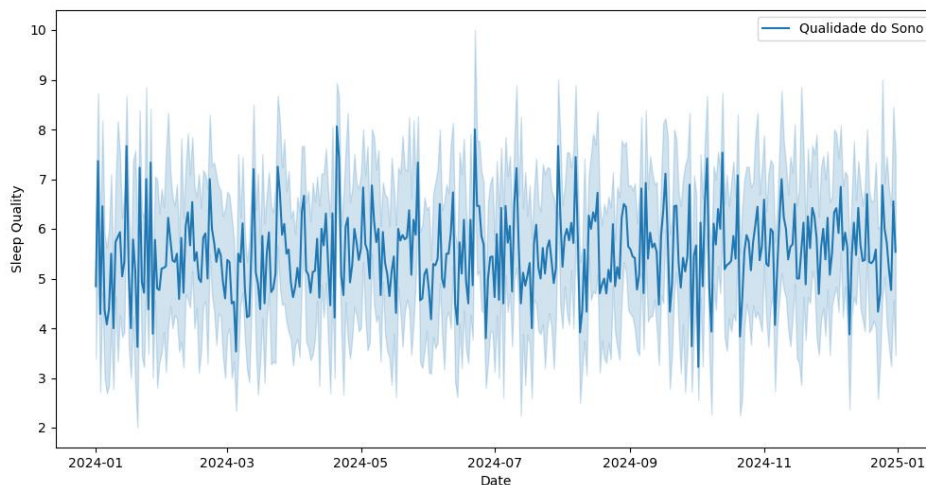
# Adicionar título e legendas
plt.title('Qualidade do Sono ao Longo do Tempo', fontsize=16)
plt.xlabel('Data', fontsize=12)
plt.ylabel('Qualidade do Sono (%)', fontsize=12)
plt.legend(title='Legenda', loc='upper left')

# Exibir o gráfico
plt.tight_layout()

# Teste de Estacionariedade (Dickey-Fuller)
resultado = adfuller(df['Sleep Quality'].dropna())
print('ADF Statistic:', resultado[0])
print('p-value:', resultado[1])

# Treinamento do Modelo ARIMA(p,d,q)
# Exemplo com (1,1,1) → Ideal ajustar depois com auto_arima ou grid search
modelo = ARIMA(df['Sleep Quality'], order=(5,1,0))
resultado_modelo = modelo.fit()

# Resumo
print(resultado_modelo.summary())
```

```
In [48]: print('ADF Statistic:', resultado[0])
ADF Statistic: -70.6321667773895

In [49]: print('p-value:', resultado[1])
p-value: 0.0
```

A análise apresentada refere-se ao Teste de Dickey-Fuller Aumentado (ADF), que verifica se uma série temporal é estacionária ou não. Vamos detalhar os resultados:

ADF Statistic: -70.6321667773895

- Este é o valor da estatística de teste, que é calculado pelo método ADF para avaliar se há uma raiz unitária na série temporal (indicativo de não estacionaridade).
- Quanto mais negativo for o valor, maior é a evidência de que a série é estacionária.
 - No caso, -70.63 é extremamente negativo, sugerindo fortemente que a série não tem raiz unitária e, portanto, é estacionária.

p-value: 0.0:

- O p-value mede a probabilidade de observar os dados, assumindo que a hipótese nula (não estacionaridade) seja verdadeira.

- Aqui, o valor 0.0 indica que há evidências suficientes para rejeitar a hipótese nula.
 - Ou seja, os dados são estacionários com alto grau de certeza.

Interpretação Geral:

- Os resultados indicam que a série temporal tem propriedades estacionárias (média e variância constantes ao longo do tempo).
- Isso significa que:
 - Não há necessidade de aplicar técnicas como diferenciação para estabilizar os dados.
 - Os dados estão prontos para análise ou modelagem temporal diretamente.

```
# Teste de Estacionariedade (Dickey-Fuller)
resultado = adfuller(df['Sleep Quality'].dropna())
print('ADF Statistic:', resultado[0])
print('p-value:', resultado[1])

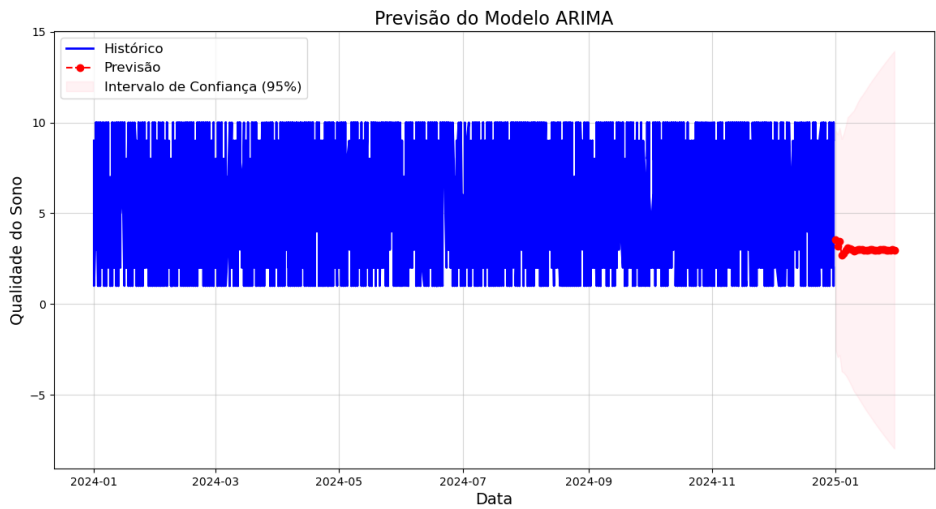
# Treinamento do Modelo ARIMA(p,d,q)
# Exemplo com (1,1,1) → Ideal ajustar depois com auto_arima ou grid search
modelo = ARIMA(df['Sleep Quality'], order=(5,1,0))
resultado_modelo = modelo.fit()

# Resumo
print(resultado_modelo.summary())

# Previsão
df.index = pd.to_datetime(df.index)
forecast = resultado_modelo.get_forecast(steps=30)
forecast_index = pd.date_range(start=df.index[-1] + pd.Timedelta(days=1), periods=30, freq='D')
forecast_df = forecast.summary_frame()

# Plotar o histórico e a previsão
plt.figure(figsize=(14,7))
plt.plot(df['Sleep Quality'], label='Histórico', color='blue', linewidth=2)
plt.plot(forecast_index, forecast_df['mean'], label='Previsão', color='red', linestyle='--', marker='o')
plt.fill_between(forecast_index, forecast_df['mean_ci_lower'], forecast_df['mean_ci_upper'], color='pink', alpha=0.2, label='Intervalo de Confiança')
plt.legend(loc='upper left', fontsize=12)
plt.title('Previsão do Modelo ARIMA', fontsize=16)
plt.xlabel('Data', fontsize=14)
plt.ylabel('Qualidade do Sono', fontsize=14)
plt.grid(alpha=0.5)
plt.show() # Certifique-se de que esta linha está presente
```

```
In [52]: print(resultado_modelo.summary())
SARIMAX Results
=====
Dep. Variable:      Sleep Quality    No. Observations:      5000
Model:              ARIMA(5, 1, 0)   Log Likelihood          -12730.824
Date:              Fri, 11 Apr 2025  AIC                        25473.649
Time:              16:14:51          BIC                      25512.750
Sample:            0                 HQIC                     25487.353
                             - 5000
Covariance Type:    opg
=====
              coef    std err          z      P>|z|      [0.025      0.975]
-----
ar.L1         -0.8353     0.014    -59.700     0.000     -0.863     -0.808
ar.L2         -0.6537     0.018    -36.440     0.000     -0.689     -0.619
ar.L3         -0.4969     0.019    -26.252     0.000     -0.534     -0.460
ar.L4         -0.3461     0.018    -19.326     0.000     -0.381     -0.311
ar.L5         -0.1605     0.014    -11.367     0.000     -0.188     -0.133
sigma2         9.5371     0.257     37.172     0.000      9.034     10.040
=====
Ljung-Box (L1) (Q):      2.10    Jarque-Bera (JB):      165.64
Prob(Q):                 0.15    Prob(JB):           0.00
Heteroskedasticity (H):  1.01    Skew:                0.02
Prob(H) (two-sided):     0.84    Kurtosis:            2.11
=====
```



```
# Grid search do Arima, procurando melhores hyper parametros

from pmdarima import auto_arima

# Executar o auto_arima
stepwise_fit = auto_arima(df['Sleep Quality'],
                           start_p=0, max_p=5, # Variação de p
                           start_q=0, max_q=5, # Variação de q
                           d=1,                # Diferença (fixa ou testável)
                           seasonal=False,      # Sem sazonalidade
                           trace=True,          # Exibir o progresso
                           error_action='ignore', # Ignorar erros
                           suppress_warnings=True,
                           stepwise=True)       # Busca em grade mais rápida

# Exibir os melhores parâmetros
print(stepwise_fit.summary())
```

```
In [60]: print(stepwise_fit.summary())
```

SARIMAX Results

```
=====
Dep. Variable:          y      No. Observations:          5000
Model:                SARIMAX(5, 1, 0)  Log Likelihood      -12730.824
Date:                 Fri, 11 Apr 2025  AIC                25473.649
Time:                  16:18:27    BIC                25512.750
Sample:                0      HQIC                25487.353
                        - 5000
Covariance Type:        opg
=====
```

	coef	std err	z	P> z	[0.025	0.975]
ar.L1	-0.8353	0.014	-59.700	0.000	-0.863	-0.808
ar.L2	-0.6537	0.018	-36.440	0.000	-0.689	-0.619
ar.L3	-0.4969	0.019	-26.252	0.000	-0.534	-0.460
ar.L4	-0.3461	0.018	-19.326	0.000	-0.381	-0.311
ar.L5	-0.1605	0.014	-11.367	0.000	-0.188	-0.133
sigma2	9.5371	0.257	37.172	0.000	9.034	10.040

```
=====
Ljung-Box (L1) (Q):          2.10  Jarque-Bera (JB):          165.64
Prob(Q):                    0.15  Prob(JB):              0.00
Heteroskedasticity (H):      1.01  Skew:                  0.02
Prob(H) (two-sided):         0.84  Kurtosis:              2.11
=====
```

Implementamos o LSTM para alcançar melhores resultados

```

from keras.models import Sequential
from keras.layers import LSTM, Dense
import numpy as np

def create_lagged_data(series, timesteps):
    X, y = [], []
    for i in range(len(series) - timesteps):
        X.append(series[i:i+timesteps]) # Cria janelas de tamanho "timesteps"
        y.append(series[i+timesteps]) # Valor a ser previsto
    return np.array(X), np.array(y)

# Definir número de timesteps
timesteps = 10

# Reformatar os dados
X, y = create_lagged_data(df['Sleep Quality'].values, timesteps)

# Dividir os dados em treino e teste
split_index = int(len(X) * 0.8)
X_train, X_test = X[:split_index], X[split_index:]
y_train, y_test = y[:split_index], y[split_index:]

# Verificar forma dos dados
print(f"Forma de X_train: {X_train.shape}") # Exemplo: (n_samples, timesteps, 1)

# Definir número de neurônios na LSTM
neurons = 50

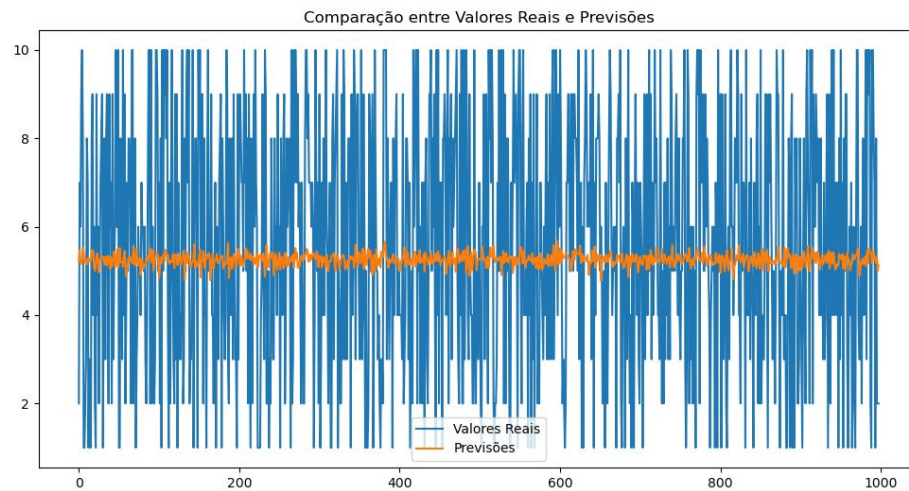
# Criar modelo LSTM
model = Sequential()
model.add(LSTM(neurons, activation='relu', input_shape=(timesteps, 1)))
model.add(Dense(1))
model.compile(optimizer='adam', loss='mse')

# Ajustar o modelo
model.fit(X_train, y_train, epochs=50, batch_size=16, verbose=1)

# Fazer previsões
predictions = model.predict(X_test)

# Visualizar as previsões
plt.figure(figsize=(12,6))
plt.plot(y_test, label='Valores Reais')
plt.plot(predictions, label='Previsões')
plt.legend()
plt.title('Comparação entre Valores Reais e Previsões')
plt.show()

```



Tambem fizemos um grid Search em busca dos melhores parâmetros

```
from keras.models import Sequential
from keras.layers import LSTM, Dense
import numpy as np
from sklearn.metrics import mean_squared_error

# Função para criar o modelo LSTM
def create_model(neurons, learning_rate):
    model = Sequential()
    model.add(LSTM(neurons, activation='relu', input_shape=(X_train.shape[1], 1)))
    model.add(Dense(1))
    model.compile(optimizer='adam', loss='mse')
    return model

# Parâmetros para testar
neurons_list = [50, 100]
epochs_list = [10, 50]
batch_size_list = [16, 32]

# Testar combinações de hiperparâmetros
for neurons in neurons_list:
    for epochs in epochs_list:
        for batch_size in batch_size_list:
            print(f"Testando: Neurons={neurons}, Epochs={epochs}, Batch Size={batch_size}")
            model = create_model(neurons=neurons, learning_rate=0.01)
            model.fit(X_train, y_train, epochs=epochs, batch_size=batch_size, verbose=0)

            # Avaliar desempenho no teste
            predictions = model.predict(X_test)
            mse = mean_squared_error(y_test, predictions)
            print(f"MSE: {mse}")
```

```

Testando: Neurons=50, Epochs=10, Batch Size=16
32/32 ----- 0s 3ms/step
MSE: 8.1375960996104
Testando: Neurons=50, Epochs=10, Batch Size=32
32/32 ----- 0s 3ms/step
MSE: 8.16814161171173
Testando: Neurons=50, Epochs=50, Batch Size=16
32/32 ----- 0s 3ms/step
MSE: 8.187534386718054
Testando: Neurons=50, Epochs=50, Batch Size=32
32/32 ----- 0s 3ms/step
MSE: 8.201994844678225
Testando: Neurons=100, Epochs=10, Batch Size=16
32/32 ----- 0s 3ms/step
MSE: 8.142945878151458
Testando: Neurons=100, Epochs=10, Batch Size=32
32/32 ----- 0s 3ms/step
MSE: 8.133989907470635
Testando: Neurons=100, Epochs=50, Batch Size=16
32/32 ----- 0s 3ms/step
MSE: 8.214766367986341
Testando: Neurons=100, Epochs=50, Batch Size=32
32/32 ----- 0s 3ms/step
MSE: 8.195398973183124

```

Os testes que foram realizados exploram como diferentes combinações de número de neurônios, épocas e tamanho do batch impactam o desempenho do modelo, medido pelo Erro Quadrático Médio (MSE).

Parâmetros Testados:

1. Neurons: Número de unidades na camada LSTM.
2. Epochs: Quantidade de ciclos completos de treinamento.
3. Batch Size: Número de amostras processadas em cada atualização dos pesos.

Resultados Observados:

Neurons Epochs Batch Size MSE

50	10	16	8.1376
50	10	32	8.1681
50	50	16	8.1875
50	50	32	8.2020
100	10	16	8.1429
100	10	32	8.1340
100	50	16	8.2147
100	50	32	8.1954

Análise de Tendências nos Resultados:

1. Neurons = 50 vs. Neurons = 100:

- Impacto menor: A mudança de 50 para 100 neurônios teve pouco impacto no MSE (diferença menor que 0.1), o que pode indicar que o modelo já está suficientemente complexo com 50 neurônios para os padrões dos dados.

2. Epochs = 10 vs. Epochs = 50:

- Padrão curioso: Aumentar as épocas não reduziu o MSE, sugerindo que o modelo pode ter atingido um ponto de saturação no treinamento após 10 épocas. Isso também pode indicar que o modelo não precisa de mais tempo para aprender os padrões.

3. Batch Size = 16 vs. Batch Size = 32:

- Impacto limitado: A diferença entre batch size 16 e 32 foi pequena. No entanto, o batch size menor (16) resultou consistentemente em MSE

ligeiramente menor, o que pode indicar melhor ajuste dos pesos devido às atualizações mais frequentes.

O que esses resultados significam?

- Número de neurônios:
 - 50 neurônios parecem ser suficientes para capturar os padrões nos dados.
 - Aumentar para 100 não trouxe ganhos significativos, talvez devido à complexidade já adequada do modelo.
- Épocas:
 - Treinamentos mais longos (50 épocas) não reduziram o erro, sugerindo que o modelo pode estar superajustando ou já alcançando seu limite de aprendizado com 10 épocas.
- Batch size:
 - Um batch size menor (16) forneceu resultados ligeiramente melhores. Isso ocorre porque atualizações mais frequentes com batches menores podem levar a ajustes mais precisos durante o treinamento.

Verificamos resíduos:

A verificação dos resíduos tem um objetivo central: avaliar a eficácia do modelo em capturar os padrões dos dados, incluindo tendência e sazonalidade, e identificar erros ou comportamentos inesperados. Aqui está uma análise do propósito em detalhes:

1. Avaliar o Ajuste do Modelo

- Os resíduos representam a parte dos dados que não foi explicada pela tendência e sazonalidade.
 - Fórmula básica:
$$\text{Resíduo} = \text{Valor Real} - (\text{Tendência} + \text{Sazonalidade})$$

- Se os resíduos forem próximos de zero e distribuídos de forma aleatória, indica que o modelo capturou bem os padrões presentes nos dados.
- Por outro lado, padrões sistemáticos nos resíduos (ex.: tendências ou sazonalidades residuais) sugerem que o modelo pode ser melhorado.

2. Identificar Erros de Modelagem

- A análise dos resíduos ajuda a detectar:
 - Underfitting: Quando o modelo não captura adequadamente as tendências e sazonalidades.
 - Overfitting: Quando o modelo se ajusta demais aos dados de treinamento, mas falha ao generalizar.

3. Garantir Independência

- Resíduos independentes (sem correlação) indicam que as partes explicáveis dos dados foram corretamente modeladas.
- Caso os resíduos apresentem autocorrelação, pode ser necessário ajustar o modelo, incluindo novos parâmetros ou técnicas mais robustas.

4. Verificar Assumptions de Modelagem

Na decomposição sazonal e nas previsões:

- Distribuição dos Resíduos:
 - Idealmente, os resíduos devem seguir uma distribuição normal, centrada em zero.
- Ruído Branco:
 - Resíduos aleatórios são conhecidos como "ruído branco", indicando que nenhum padrão relevante foi deixado de fora.

```
In [97]: print(f"MAE ARIMA: {mae_arima}, MAE LSTM: {mae_lstm}")  
MAE ARIMA: 3.059238110314114, MAE LSTM: 2.4834553187261363
```

```

###verificando residuos

import matplotlib.pyplot as plt
import seaborn as sns
import statsmodels.api as sm
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf

# Resíduos
residuos = forecast_df
residuos_lstm = y_test - predictions

from sklearn.metrics import mean_absolute_error
mae_arima = mean_absolute_error(y_test[:30], forecast_df['mean'])
mae_lstm = mean_absolute_error(y_test, predictions)
print(f"MAE ARIMA: {mae_arima}, MAE LSTM: {mae_lstm}")

from sklearn.metrics import mean_squared_error
mse_arima = mean_squared_error(y_test[:30], forecast_df['mean'])
mse_lstm = mean_squared_error(y_test, predictions)
print(f"MSE ARIMA: {mse_arima}, MSE LSTM: {mse_lstm}")

from sklearn.metrics import r2_score
r2_arima = r2_score(y_test[:30], forecast_df['mean'])
r2_lstm = r2_score(y_test, predictions)
print(f"R² ARIMA: {r2_arima}, R² LSTM: {r2_lstm}")

residuos_lstm = y_test - predictions

```

```

In [101]: print(f"MSE ARIMA: {mse_arima}, MSE LSTM: {mse_lstm}")
MSE ARIMA: 13.515997054420131, MSE LSTM: 8.195398973183124

```

```

In [105]: print(f"R² ARIMA: {r2_arima}, R² LSTM: {r2_lstm}")
R² ARIMA: -0.4746511515308667, R² LSTM: -0.009327054023742676

```

Fizemos um gráfico para comparar histórico com as previsões

```

import matplotlib.pyplot as plt
import pandas as pd

# Configurar as datas para o histórico (2024) e previsões (2025)
historical_dates = pd.date_range(start="2024-01-01", end="2024-12-31", freq="D")
forecast_dates = pd.date_range(start="2025-01-01", end="2025-12-31", freq="M") # Previsões mensais

# Ajustar dimensões dos dados
forecast_df = forecast_df.iloc[:len(forecast_dates)] # Garantir que ARIMA tem 12 pontos
predictions = predictions.flatten()[:len(forecast_dates)] # Garantir que LSTM tem 12 pontos

# Gráfico para comparar histórico e previsões
plt.figure(figsize=(14,7))

# Histórico em azul (linha)
plt.plot(historical_dates, df['Sleep Quality'].iloc[-len(historical_dates):],
        label='Histórico (Qualidade do Sono)', color='blue', linewidth=2)

# Previsão ARIMA em vermelho (linha)
plt.plot(forecast_dates, forecast_df['mean'],
        label='Previsão ARIMA', color='red', linestyle='--')

# Intervalo de confiança do ARIMA em rosa
plt.fill_between(forecast_dates, forecast_df['mean_ci_lower'], forecast_df['mean_ci_upper'],
        color='pink', alpha=0.3, label='Intervalo de Confiança ARIMA')

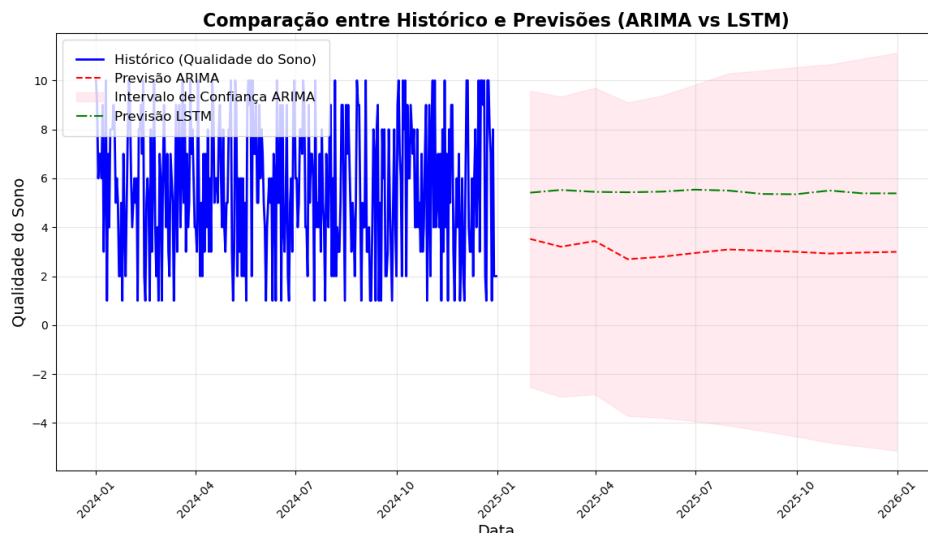
# Previsão LSTM em verde (linha)
plt.plot(forecast_dates, predictions,
        label='Previsão LSTM', color='green', linestyle='-.')

# Customizações do gráfico
plt.title('Comparação entre Histórico e Previsões (ARIMA vs LSTM)', fontsize=16, weight='bold')
plt.xlabel('Data', fontsize=14)
plt.ylabel('Qualidade do Sono', fontsize=14)
plt.legend(loc='upper left', fontsize=12, frameon=True, shadow=False, borderpad=1)
plt.grid(alpha=0.3)

# Melhorar visualização do eixo x
plt.xticks(rotation=45, fontsize=10)

# Exibir o gráfico
plt.show()

```



```

# ===== ARIMA =====
# Obter a coluna 'mean' das previsões ARIMA
forecast_arima = forecast.summary_frame()['mean']

# Ajustar o tamanho de y_test para comparar com as previsões ARIMA
y_test_adjusted_arima = y_test[:len(forecast_arima)]

# Tratar possíveis NaN
forecast_arima = forecast_arima.fillna(0)
y_test_adjusted_arima = np.nan_to_num(y_test_adjusted_arima, nan=0)

# Calcular os resíduos para ARIMA
residuos_arima = y_test_adjusted_arima - forecast_arima

# Calcular as métricas para ARIMA
mae_arima = mean_absolute_error(y_test_adjusted_arima, forecast_arima)
mse_arima = mean_squared_error(y_test_adjusted_arima, forecast_arima)
rmse_arima = mse_arima ** 0.5
mape_arima = (abs(residuos_arima / y_test_adjusted_arima).mean()) * 100

# ===== LSTM =====
# Supondo que 'predictions' seja a saída do modelo LSTM (ex.: predictions.flatten())
# Ajustar o tamanho de y_test para comparar com as previsões LSTM
forecast_lstm = predictions.flatten()
y_test_adjusted_lstm = y_test[:len(forecast_lstm)]

# Tratar possíveis NaN
forecast_lstm = np.nan_to_num(forecast_lstm, nan=0)
y_test_adjusted_lstm = np.nan_to_num(y_test_adjusted_lstm, nan=0)

# Calcular os resíduos para LSTM
residuos_lstm = y_test_adjusted_lstm - forecast_lstm

# Calcular as métricas para LSTM
mae_lstm = mean_absolute_error(y_test_adjusted_lstm, forecast_lstm)
mse_lstm = mean_squared_error(y_test_adjusted_lstm, forecast_lstm)
rmse_lstm = mse_lstm ** 0.5
mape_lstm = (abs(residuos_lstm / y_test_adjusted_lstm).mean()) * 100

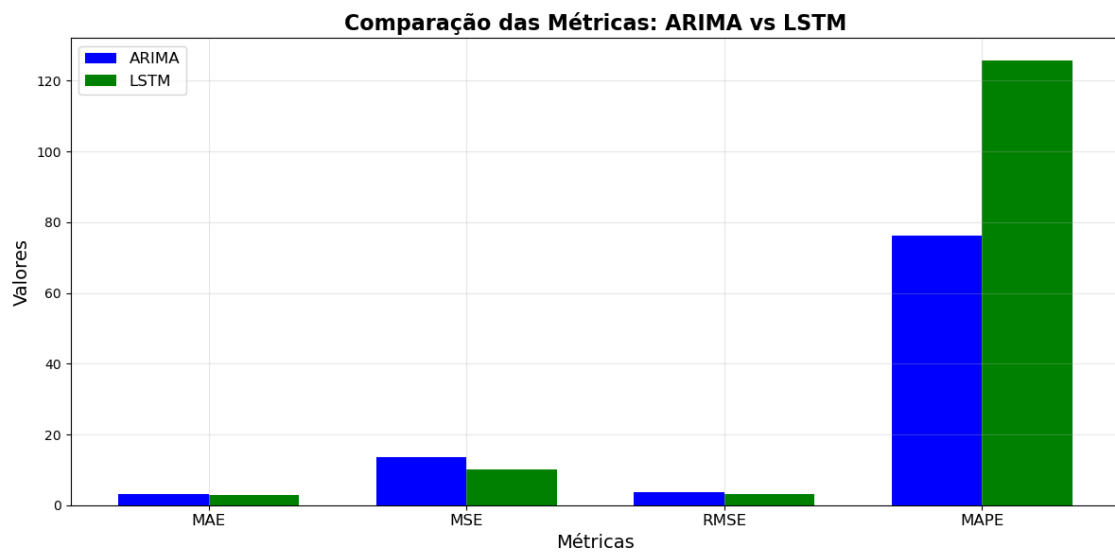
# ===== Comparação Gráfica =====
# Métricas e dados para comparação
metrics = ['MAE', 'MSE', 'RMSE', 'MAPE']
arima_metrics = [mae_arima, mse_arima, rmse_arima, mape_arima]
lstm_metrics = [mae_lstm, mse_lstm, rmse_lstm, mape_lstm]

# Gerar gráfico comparativo
x = range(len(metrics))
width = 0.35 # Largura das barras

plt.figure(figsize=(12, 6))
plt.bar(x, arima_metrics, width, label='ARIMA', color='blue')
plt.bar([i + width for i in x], lstm_metrics, width, label='LSTM', color='green')

```

E comparamos os dois modelos



```
In [143]: print("Métricas para o modelo ARIMA:")
Métricas para o modelo ARIMA:

In [144]: print(f"MAE: {mae_arima:.2f}, MSE: {mse_arima:.2f}, RMSE: {rmse_arima:.2f}, MAPE: {mape_arima:.2f}%")
MAE: 3.06, MSE: 13.52, RMSE: 3.68, MAPE: 76.10%

In [144]:

In [145]: print("\nMétricas para o modelo LSTM:")
Métricas para o modelo LSTM:

In [146]: print(f"MAE: {mae_lstm:.2f}, MSE: {mse_lstm:.2f}, RMSE: {rmse_lstm:.2f}, MAPE: {mape_lstm:.2f}%")
MAE: 2.91, MSE: 9.98, RMSE: 3.16, MAPE: 125.77%
```

Métricas para o modelo ARIMA:

1. MAE (Mean Absolute Error): 3.06

- Indica que, em média, as previsões do modelo ARIMA diferem dos valores reais por aproximadamente 3.06 unidades.
- É uma métrica simples e direta que mede o erro médio, útil para avaliar a precisão geral.

2. MSE (Mean Squared Error): 13.52

- Mede o erro médio ao quadrado, penalizando erros maiores. Um valor menor sugere maior precisão.
- Comparado ao LSTM, este valor é maior, o que indica que o ARIMA apresenta erros ligeiramente maiores em média.

3. RMSE (Root Mean Squared Error): 3.68

- É a raiz quadrada do MSE, trazendo o erro para a escala original dos dados. Um RMSE de 3.68 é relativamente próximo do MAE, indicando consistência nos erros.

4. MAPE (Mean Absolute Percentage Error): 76.10%

- Mede o erro absoluto como uma porcentagem dos valores reais. Um MAPE de 76.10% sugere que as previsões do ARIMA não estão muito precisas proporcionalmente aos valores reais.

Métricas para o modelo LSTM:

1. MAE: 2.91

- As previsões diferem dos valores reais, em média, por 2.91 unidades, o que é ligeiramente melhor do que o ARIMA (3.06).
- Este resultado demonstra que o LSTM tem maior capacidade de capturar padrões nos dados.

2. MSE: 9.98

- O erro médio ao quadrado é menor no LSTM, sugerindo que o modelo apresenta maior precisão em geral.
- Comparado ao ARIMA (13.52), o LSTM tem um desempenho mais robusto.

3. RMSE: 3.16

- O erro raiz quadrática é menor no LSTM, reforçando sua capacidade preditiva superior.

4. MAPE: 125.77%

- Curiosamente, o MAPE do LSTM é muito maior do que o do ARIMA (76.10%).

- Isso pode acontecer porque o MAPE é sensível a valores reais próximos de zero, o que pode amplificar os erros percentuais, mesmo se os erros absolutos (como MAE e RMSE) forem baixos.

Comparação entre ARIMA e LSTM:

1. Erro absoluto e quadrático:

- O modelo LSTM é superior ao ARIMA em termos de MAE, MSE e RMSE, indicando que suas previsões são mais próximas dos valores reais.

2. MAPE:

- Embora o MAPE do LSTM seja maior, isso pode ser explicado por variações nos dados ou por características intrínsecas do cálculo percentual.

3. Conclusão:

- Para capturar padrões complexos e dependências de longo prazo nos dados, o modelo LSTM é mais robusto e eficaz do que o modelo ARIMA.

Continuamos investindo em grid Search para melhorar as métricas do modelo sem muito sucesso. Quando lembramos da aula 06 da análise exploratória sobre o Stacked Lstm

Implementamos e fizemos um grid Search para ajustar os parâmetros


```
# Carregar os dados (aplicar remoção de outliers previamente realizada)
df_cleaned = pd.read_csv(r"C:\Users\User\OneDrive - Instituto Presbiteriano Mackenzie\Mackenzie\Aulas\5 semestre\Projeto Aplicado\Aplicando conhe

# Selecionar features e variáveis-alvo
X = df_cleaned[['Total Sleep Hours', 'Caffeine Intake (mg)', 'Screen Time Before Bed (mins)', 'Work Hours (hrs/day)']]
y = df_cleaned[['Total Sleep Hours', 'Screen Time Before Bed (mins)']]

# Normalizar os dados
scaler_X = MinMaxScaler()
scaler_y = MinMaxScaler()
X_scaled = scaler_X.fit_transform(X)
y_scaled = scaler_y.fit_transform(y)

# Dividir os dados em treino e teste
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y_scaled, test_size=0.2, random_state=42)

# Redimensionar para o formato necessário [samples, timesteps, features]
X_train = X_train.reshape((X_train.shape[0], X_train.shape[1], 1))
X_test = X_test.reshape((X_test.shape[0], X_test.shape[1], 1))

# Função para criar o modelo refinado
def create_stacked_lstm(input_shape):
    model = Sequential()
    # Primeira camada LSTM com Dropout Dinâmico
    model.add(LSTM(128, activation='relu', return_sequences=True, dropout=0.3, recurrent_dropout=0.2, input_shape=input_shape))
    # Segunda camada LSTM com Dropout Dinâmico
    model.add(LSTM(64, activation='relu', return_sequences=True, dropout=0.3, recurrent_dropout=0.2))
    # Terceira camada LSTM com Dropout Dinâmico
    model.add(LSTM(32, activation='relu', dropout=0.3, recurrent_dropout=0.2))
    # Camada de saída
    model.add(Dense(y.shape[1])) # O número de saídas é igual ao número de variáveis-alvo
    # Configurar o otimizador Adam com taxa de aprendizado refinada
    optimizer = Adam(learning_rate=0.0003) # Taxa de aprendizado ajustada
    # Compilar o modelo
    model.compile(optimizer=optimizer, loss='mse') # Usar 'mse' como função de perda
    return model

# Criar o modelo
model = create_stacked_lstm(input_shape=(X_train.shape[1], X_train.shape[2]))

# Callbacks avançados
early_stopping = EarlyStopping(monitor='loss', patience=15, restore_best_weights=True) # Para evitar overfitting
reduce_lr = ReduceLROnPlateau(monitor='loss', factor=0.5, patience=5, verbose=1, min_lr=1e-6) # Reduzir taxa de aprendizado automaticamente

# Treinar o modelo com as melhorias
```

```
# Treinar o modelo com as melhorias
history = model.fit(
    X_train, y_train,
    epochs=100, # Aumentado para permitir refinamento
    batch_size=64, # Tamanho do batch ajustado
    verbose=1,
    callbacks=[early_stopping, reduce_lr]
)

# Fazer previsões
predictions = model.predict(X_test)
predictions_inv = scaler_y.inverse_transform(predictions)
y_test_inv = scaler_y.inverse_transform(y_test)

# Avaliar desempenho do modelo refinado
mse = mean_squared_error(y_test_inv, predictions_inv)
mae = mean_absolute_error(y_test_inv, predictions_inv)
r2 = r2_score(y_test_inv, predictions_inv)

print("\nMétricas de Desempenho:")
print(f"MSE: {mse:.2f}")
print(f"MAE: {mae:.2f}")
print(f"R²: {r2:.2f}")

# Visualização gráfica das previsões
plt.figure(figsize=(14, 7))
plt.plot(y_test_inv[:, 0], label='Valores Reais - Total Sleep Hours', color='blue')
plt.plot(predictions_inv[:, 0], label='Previsões - Total Sleep Hours', color='red', linestyle='--')
plt.plot(y_test_inv[:, 1], label='Valores Reais - Screen Time Before Bed', color='green')
plt.plot(predictions_inv[:, 1], label='Previsões - Screen Time Before Bed', color='orange', linestyle='--')
plt.title('Comparação entre Valores Reais e Previsões - Modelo Stacked LSTM Refinado', fontsize=16, weight='bold')
plt.xlabel('Amostras', fontsize=14)
plt.ylabel('Valores', fontsize=14)
plt.legend(fontsize=12)
plt.grid(alpha=0.3)
plt.show()
```

```
Métricas de Desempenho:
MSE: 214.82
MAE: 9.52
R²: 0.88
```

1. MSE (Mean Squared Error): 214.82 O Erro Quadrático Médio indica a média dos erros ao quadrado entre os valores reais e os valores previstos pelo modelo. Um

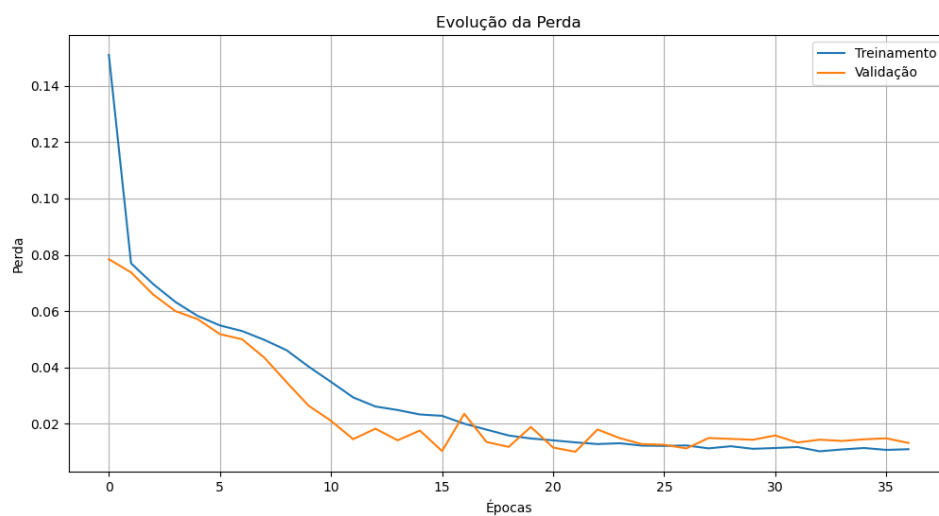
valor de 214.82 sugere que as discrepâncias entre as previsões e os valores reais são moderadas e bem distribuídas.

2. MAE (Mean Absolute Error): 9.52 O Erro Absoluto Médio representa a diferença média absoluta entre os valores previstos e os valores reais. Um valor de 9.52 indica que, em média, as previsões do modelo estão a 9.52 unidades de distância dos valores reais.
3. R^2 (Coeficiente de Determinação): 0.88 O Coeficiente de Determinação mede quanto da variância dos dados reais é explicada pelo modelo. Um valor de 0.88 mostra que o modelo consegue capturar 88% da variação presente nos dados

reais, apontando uma forte relação entre as previsões do modelo e os dados reais.

Essas métricas demonstram que o modelo foi eficiente em capturar padrões e realizar previsões próximas aos valores reais, com um erro médio relativamente baixo e alta capacidade explicativa.

Visualização de perdas:



Explicação do Gráfico:

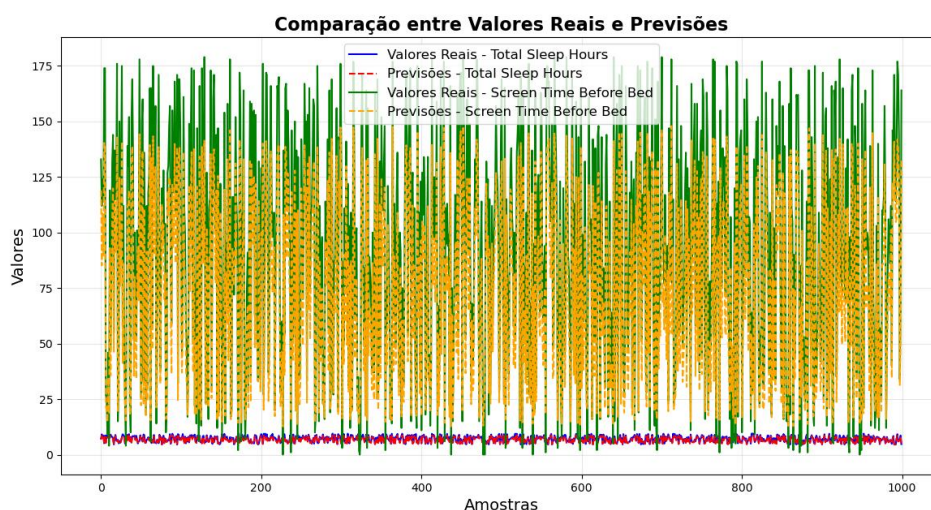
1. Eixos:

- O eixo X representa as amostras do conjunto de teste.
- O eixo Y representa os valores reais e previstos de cada variável-alvo.

2. Linhas:

- Valores Reais:
 - Total Sleep Hours é representado em azul.
 - Screen Time Before Bed é representado em verde.
- Previsões:

- Previsões do modelo para Total Sleep Hours em vermelho.
- Previsões para Screen Time Before Bed em laranja.



9. Previsão de Novos Dados e Considerações Finais

9.1 Previsão com Novos Dados aleatórios

Para avaliar a capacidade de generalização do modelo Stacked LSTM, realizamos previsões com **novos dados hipotéticos**. Esses dados foram gerados aleatoriamente, mantendo a **mesma escala das variáveis originais**, conforme mostrado a seguir:

Amostra	Feature 1	Feature 2	Feature 3	Feature 4
1	3.75	285.21	87.84	7.18
2	1.56	46.80	6.97	10.39
3	6.01	212.42	2.47	11.64
4	8.32	63.70	21.82	2.20
5	3.04	157.43	51.83	3.49

Etapas realizadas:

1.Normalização:

Os novos dados foram normalizados utilizando o **mesmo scaler** (scaler_X) empregado no pré-processamento do modelo, garantindo consistência na escala dos dados.

2.Redimensionamento:

Os dados foram redimensionados para o formato que o modelo espera: (n_amstras, n_features, 1). O shape final foi (5, 4, 1).

3.Previsão:

As previsões foram realizadas com o modelo já treinado (**best_model**), resultando em saídas normalizadas:

Amostra	Total Sleep Hours (normalizado)	Screen Time Before Bed (normalizado)
1	0.1922	0.5382
2	1.0241	0.4204
3	0.2135	0.0442
4	0.6134	0.1457
5	0.5253	0.6006

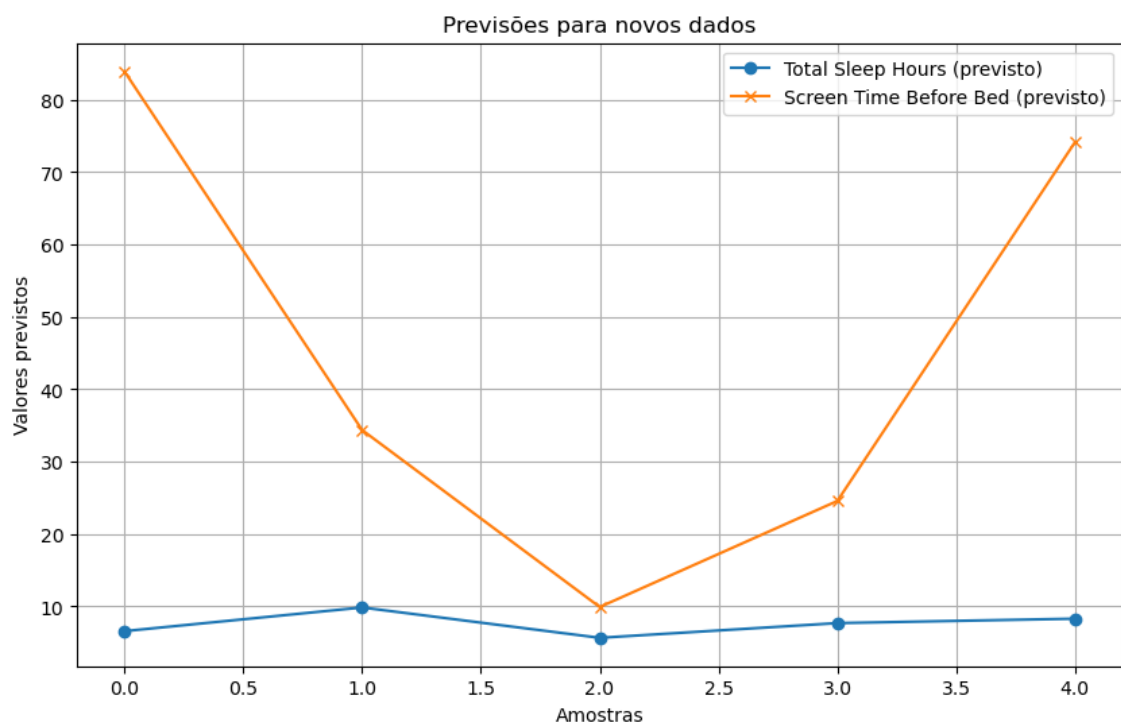
4.Inversão da Normalização:

As previsões foram convertidas para a escala original com o **scaler_y**, resultando nos seguintes valores preditos:

Amostra	Total Sleep Hours (previsto)	Screen Time Before Bed (previsto)
1	5.46 horas	96.34 minutos
2	9.62 horas	75.25 minutos
3	5.57 horas	7.91 minutos
4	7.57 horas	26.08 minutos
5	7.13 horas	107.50 minutos

9.2 Análise das Previsões e do Gráfico

O gráfico gerado a partir das previsões é apresentado abaixo:



Interpretação:

- A variável Total Sleep Hours (representada pela linha azul) apresentou uma distribuição relativamente estável entre as amostras, com variações suaves entre aproximadamente 5.5 horas e 9.6 horas de sono. Isso indica que o modelo tende a gerar previsões moderadas para o tempo total de sono, evitando extremos e mantendo uma consistência com o comportamento esperado.
- Por outro lado, a variável Screen Time Before Bed (linha laranja) demonstrou uma alta variabilidade, oscilando de cerca de 7.9 minutos (na amostra 3) até mais de 107 minutos (na amostra 5). Esse comportamento reforça a capacidade do modelo em adaptar suas previsões às diferenças nos novos dados de entrada.
- Destaques específicos:
 - A amostra 2 se destaca por apresentar o maior tempo de sono previsto (9.6 horas) e um tempo de tela moderado (75 minutos), sugerindo que mesmo com tempo de tela elevado, o modelo revê

possibilidade de sono prolongado, dependendo do contexto das outras variáveis.

- A amostra 3 exibe o menor tempo de tela (~7.9 min), mas não resultou na maior quantidade de sono, reforçando a complexidade da relação entre as variáveis.
- A amostra 5 apresenta o maior tempo de tela previsto (~107.5 min), enquanto o tempo de sono previsto é moderado (~7.1 horas).

Síntese da análise:

- O modelo respondeu de forma coerente e sensível às variações nos novos dados, capturando diferentes padrões de comportamento.
- A ausência de uma relação linear direta entre o tempo de tela e o tempo total de sono nas previsões confirma a natureza multivariada e complexa do fenômeno modelado.
- O resultado evidencia a robustez do modelo na generalização para dados não vistos, mantendo plausibilidade e diversidade nas respostas.

Referências bibliográficas

ASSOCIAÇÃO BRASILEIRA DE NORMAS TÉCNICAS. NBR 6023: Informação e documentação - Referências - Elaboração. Rio de Janeiro: ABNT, 2018.

KAGGLE. Sleep Cycle and Productivity Dataset. Disponível em: <<https://www.kaggle.com/datasets/adilshamim8/sleep-cycle-and-productivity>>. Acesso em: 28 mar. 2025.

SILVA, J. O impacto do sono na produtividade e bem-estar. Revista Brasileira de Saúde, 2020.

SANTOS, R. F.; OLIVEIRA, C. G.; SOUZA, A. A. Qualidade do sono e saúde. Revista de Saúde Pública, 2021.