

INSTITUTO FEDERAL DO ESPÍRITO SANTO  
BACHARELADO EM ENGENHARIA DE CONTROLE E AUTOMAÇÃO

**VINICIUS DE MOURA SIQUEIRA**

**USO DE REGRESSÃO LINEAR (COM REDE NEURAL DE MULTICAMADAS)  
PARA GENERALIZAR A OBTENÇÃO DO COEFICIENTE DE DECAIMENTO DE  
TURBINAS EM PLANTAS DE PROPULSÃO NAVAL**

VINICIUS DE MOURA SIQUEIRA

**USO DE REGRESSÃO LINEAR (COM REDE NEURAL DE MULTICAMADAS)  
PARA GENERALIZAR A OBTENÇÃO DO COEFICIENTE DE DECAIMENTO DE  
TURBINAS EM PLANTAS DE PROPULSÃO NAVAL**

Trabalho apresentado na disciplina de Redes Neurais  
Artificiais no Ifes Campus Linhares

Orientador: Lucas de Assis Soares

Linhares  
2021

## **RESUMO**

O presente trabalho tem como objetivo analisar o desempenho de uma regressão linear definida por uma rede neural com MLP (Multi Layer Perceptron ou Perceptron Multicamada) para prever valores do coeficiente de decaimento de turbinas em plantas de propulsão naval. Como fonte de dados para essa implementação, usou-se uma base, disponibilizada no UCI Machine Learning Repository, que continha a leitura de 16 sensores (como entrada) e 2 resultados do coeficiente de decaimento obtido em cada situação (como saída).

Palavras-chave: Regressão Linear. MLP (Multi Layer Perceptron). Coeficiente de decaimento. Plantas de Propulsão Naval.

## LISTA DE FIGURAS

Figura 1 – Exibição do formato do dataframe com os dados importados . . . . .	8
Figura 2 – Gráfico da saída geral . . . . .	12
Figura 3 – Gráfico com os valores de C1 obtidos no conjunto de dados . . . . .	13
Figura 4 – Gráfico com os valores de C2 obtidos no conjunto de dados . . . . .	14

## SUMÁRIO

<b>1</b>	<b>INTRODUÇÃO . . . . .</b>	<b>4</b>
1.1	CONJUNTO DE DADOS UTILIZADO . . . . .	4
<b>2</b>	<b>DESENVOLVIMENTO . . . . .</b>	<b>6</b>
<b>3</b>	<b>RESULTADOS E DISCUSSÕES . . . . .</b>	<b>13</b>
3.1	ANÁLISE DA PREVISÃO PARA O COEFICIENTE DE DECAIMENTO DO COMPRESSOR . . . . .	13
3.2	ANÁLISE DA PREVISÃO PARA O COEFICIENTE DE DECAIMENTO DA TURBINA . . . . .	14
<b>4</b>	<b>CONCLUSÃO . . . . .</b>	<b>15</b>

## 1 INTRODUÇÃO

Considerando que a regressão linear pode ser usada para resumir a relação de dados de entrada com as saídas, faz-se útil o uso dessa técnica para generalizar a relação de informações obtidas por sensores com os resultados que esses valores, juntos, implicam em um objeto real.

### 1.1 CONJUNTO DE DADOS UTILIZADO

Uma base de dados disponibilizada pela UCI Machine Learning Repository é intitulada como "Condition Based Maintenance of Naval Propulsion Plants" ou, "Manutenção baseada na condição em plantas de propulsores navais". Esse conjunto de dados foi construído em um simulador sofisticado de turbinas a gás, montado em um navio destacado por ter uma propulsão do tipo "Combined Diesel e Letric And Gas (CODLAG)".

A base de dados conta com 11934 instâncias variáveis e possui 16 parâmetros de entrada para 2 variáveis de saída. O conjunto foi disponibilizado para uso público em 2014 e já possui mais de 75000 acessos.

A planta utilizada para realizar as simulações recebeu constantes melhorias, e teve a veracidade incrementada com diversas comparações a plantas de propulsão reais similares ao que estava sendo proposto no ambiente de simulação. Assim, os dados disponibilizados estão de acordo com o que poderia acontecer em uma implementação com um navio real.

As possibilidades do estado de degradação podem ser descritas por uma combinação de três variáveis: velocidade do navio, coeficiente de degradação do compressor e coeficiente de degradação da turbina. A variação da velocidade do navio pode variar entre 3 nós até 27 nós.

Para as variáveis de saída, os coeficientes de decaimento foram analisados em dois domínios diferentes, sendo  $[1; 0.95]$  para o coeficiente de decaimento do compressor e  $[1; 0.975]$  para o coeficiente de decaimento da turbina.

Uma lista completa dos parâmetros de entrada (da forma como foi fornecido no conjunto de dados) pode ser visualizada abaixo:

- Lever position (lp)
- Ship speed (v)
- Gas Turbine (GT) shaft torque

- GT rate of revolutions (GTn)
- Gas Generator rate of revolutions (GGn)
- Starboard Propeller Torque (Ts)
- Port Propeller Torque (Tp)
- High Pressure (HP) Turbine exit temperature (T48)
- GT Compressor inlet air temperature
- GT Compressor outlet air temperature
- HP Turbine exit pressure (P48)
- GT Compressor inlet air pressure (P1)
- GT Compressor outlet air pressure (P2)
- GT exhaust gas pressure (Pexh)
- Turbine Injection Control (TIC)
- Fuel flow (mf)

Para os dois dados de saída, tem-se:

- GT Compressor decay state coefficient
- GT Turbine decay state coefficient

Os dados foram fornecidos em um documento de texto nomeado como "data.txt" que contém os 18 dados espaçados igualmente.

## 2 DESENVOLVIMENTO

Para desenvolver o trabalho, usou-se o suporte da máquina virtual disponibilizada no Google Colaboratory, acessada diretamente pelo Google Drive.

```

1 # Importação das bibliotecas
2
3 import numpy as np #operacoes matemáticas
4 import matplotlib.pyplot as plt #plot de gráficos
5
6 import pandas as pd # manipulação de arquivos (csv, txt)
7
8 from sklearn.model_selection import train_test_split
9 from sklearn.metrics import mean_squared_error
10
11 from sklearn.preprocessing import MinMaxScaler
12
13 from sklearn.neural_network import MLPRegressor

```

O pacote *numpy* foi importado para possibilitar operações numéricas no programa. Também foi utilizado o *matplotlib.pyplot*, que possibilita a exibição de gráficos da forma como é feito originalmente no *MATLAB*, facilitando a manipulação. O pacote *pandas* foi utilizado para permitir a manipulação de arquivos com a extensão ".txt", que contém os dados de interesse para criar a regressão linear no contexto descrito.

Na linha 6 do trecho acima, há a importação do *train\_test\_split* que dividirá o conjunto de dados disponíveis em partes de treino e teste (para verificar que a rede não está apenas decorando os valores, mas sim, generalizando da maneira correta).

Em seguida, há a importação do *mean\_squared\_error* que possibilita o cálculo do erro médio quadrático, medida que indica a diferença entre o valor estimado e do parâmetro ao quadrado.

O modelo *MinMaxScaler* permite o processamento diferenciado dos intervalos de existência dos dados, para que a visualização fique coerente no gráfico, sem deixar os pontos muito distantes e desconexos.

Por fim, há a importação do modelo de regressão linear. Ele carrega consigo as características de uma regressão linear propriamente dita, além de todos os métodos de ajuste para que ela se enquadre nos parâmetros necessários. Nesse trabalho, utilizou-se um modelo de rede neural com multicamadas, o *MLPRegressor*, buscando melhorar o desempenho (com relação a uma simples regressão linear).



Em seguida, deve-se incluir os dados disponíveis no conjunto para uma variável que será utilizada internamente.

```

1  titulos_tabela = np.array([
2      "lp", #Lever Position
3      "v", #Ship Speed
4      "GTT", #Gas Turbine Shaft Torque
5      "GTn", #Gas Turbine Rate of Revolutions
6      "GGn", #Gas Generator Rate of Revolutions
7      "Ts", #Starboard Propeller Torque
8      "Tp", #Port Propeller Torque
9      "T48", #HP Turbine Exit Temperature
10     "T1", #GT Compressor Inlet air Pressure
11     "T2", #GT Compressor Outlet air Temperature
12     "P48", #HP Turbine Exhaust Gas Pressure
13     "P1", #GT Compressor Inlet Air Pressure
14     "P2", #GT Compressor Outlet Air Pressure
15     "Pexh", #Gas Turbine Exhaust Gas Pressure
16     "TIC", #Turbine Injecton Control
17     "mf", #Fuel Flow
18     "C1", #GT Compressor Decay State Coefficient
19     "C2"]) #GT Turbine Decay State Coefficient
20
21 dataframe = pd.read_csv(
22     "data.txt",
23     sep="   ",
24     header=None,
25     names=titulos_tabela)

```

O primeiro passo foi nomear as 18 colunas que serão importadas do arquivo que armazena as informações e, em seguida, utilizar uma função da biblioteca *pandas* para realizar a leitura desses dados em um dataframe adequado que permitirá a manipulação desses dados. O parâmetro *sep* se refere ao critério de separação das informações dentro do documento, que nesse caso, são 3 espaços entre as colunas.

Na Figura 1 é possível visualizar o formato dos dados identificados pelas abreviações introduzidas no início deste trabalho. As duas saídas, coeficiente de decaimento do compressor e da turbina foram identificadas como *C1* e *C2*, respectivamente.

Em seguida, os dados são extraídos do dataframe para que possam ser tratados individualmente (entradas e saídas). No trabalho em questão, as entradas foram tratadas como *dados* e as saídas como *resultados*.

	lp	v	GTT	GTn	GGn	Ts	Tp	T48	T1	T2	P48	P1	P2	Pexh	TIC	mf	C1	C2
0	1.138	3.0	289.964	1349.489	6677.380	7.584	7.584	464.006	288.0	550.563	1.096	0.998	5.947	1.019	7.137	0.082	0.95	0.975
1	2.088	6.0	6960.180	1376.166	6828.469	28.204	28.204	635.401	288.0	581.658	1.331	0.998	7.282	1.019	10.655	0.287	0.95	0.975
2	3.144	9.0	8379.229	1386.757	7111.811	60.358	60.358	606.002	288.0	587.587	1.389	0.998	7.574	1.020	13.086	0.259	0.95	0.975
3	4.161	12.0	14724.395	1547.465	7792.630	113.774	113.774	661.471	288.0	613.851	1.658	0.998	9.007	1.022	18.109	0.358	0.95	0.975
4	5.140	15.0	21636.432	1924.313	8494.777	175.306	175.306	731.494	288.0	645.642	2.078	0.998	11.197	1.026	26.373	0.522	0.95	0.975
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
11929	5.140	15.0	21624.934	1924.342	8470.013	175.239	175.239	681.658	288.0	628.950	2.087	0.998	10.990	1.027	23.803	0.471	1.00	1.000
11930	6.175	18.0	29763.213	2306.745	8800.352	245.954	245.954	747.405	288.0	658.853	2.512	0.998	13.109	1.031	32.671	0.647	1.00	1.000
11931	7.148	21.0	39003.867	2678.052	9120.889	332.389	332.389	796.457	288.0	680.393	2.982	0.998	15.420	1.036	42.104	0.834	1.00	1.000
11932	8.206	24.0	50992.579	3087.434	9300.274	438.024	438.024	892.945	288.0	722.029	3.594	0.998	18.293	1.043	58.064	1.149	1.00	1.000
11933	9.300	27.0	72775.130	3560.400	9742.950	644.880	644.880	1038.411	288.0	767.595	4.531	0.998	22.464	1.052	86.067	1.704	1.00	1.000

11934 rows x 18 columns

Figura 1 – Exibição do formato do dataframe com os dados importados

```

1 dados = dataframe.iloc[:, :16].values
2 #Obtencao dos dados do compressor na saída C1
3 resultados_C1 = dataframe.iloc[:, 16].values
4
5 #Obtencao dos dados da turbina na saída C2
6 resultados_C2 = dataframe.iloc[:, 17].values

```

A variável de entrada foi povoada com todos os 11934 registros, com os dados das 16 primeiras colunas, identificadas na Figura 1 como o intervalo de variáveis entre *lp* e *mf*. Já a variável de saída recebeu todos os 11934 registros, assim como a variável de entrada, porem contém apenas os dois últimos dados da tabela, *C1* e *C2*.

Em seguida é necessário dividir os dados de treino e teste para possibilitar a validação da capacidade de generalização da regressão linear criada.

```

1 #Divisao dos dados de treino e teste para C1
2 dados_train_C1, dados_test_C1, resultados_train_C1, resultados_test_C1 =
   train_test_split(dados, resultados_C1, test_size=0.33)
3
4 #Divisao dos dados de treino e teste para C2
5 dados_train_C2, dados_test_C2, resultados_train_C2, resultados_test_C2 =
   train_test_split(dados, resultados_C2, test_size=0.33)

```

Para definir a nova divisão dos dados (treino e teste), deve-se especificar o parâmetro *test\_size* com a proporção para os dados de teste. Nesse caso, 33% dos dados serão separados para teste.

Como os dados estão em unidades diferentes, é esperado que as medidas estejam muito distantes uma das outras (quando a unidade é desconsiderada). Para que facilite a visualização e o entendimento do resultado criado, é necessário normalizar os dados em um

intervalo que abrangerá desde os valores mínimos até os valores máximos de cada variável.

```

1 #Tratamento da escala para C1
2 scalerDados_C1 = MinMaxScaler()
3 scalerDados_C1.fit(dados_train_C1)
4 dadosTrain_norm_C1 = scalerDados_C1.transform(dados_train_C1)
5 dadosTest_norm_C1 = scalerDados_C1.transform(dados_test_C1)
6
7 scalerResultados_C1 = MinMaxScaler()
8 scalerResultados_C1.fit(resultados_train_C1.reshape(-1, 1))
9 resultadosTrain_norm_C1 = scalerResultados_C1.transform(resultados_train_C1
   .reshape(-1, 1))
10 resultadosTest_norm_C1 = scalerResultados_C1.transform(resultados_test_C1.
   reshape(-1, 1))
11
12 #Tratamento da escala para C1
13 scalerDados_C2 = MinMaxScaler()
14 scalerDados_C2.fit(dados_train_C2)
15 dadosTrain_norm_C2 = scalerDados_C2.transform(dados_train_C2)
16 dadosTest_norm_C2 = scalerDados_C2.transform(dados_test_C2)
17
18 scalerResultados_C2 = MinMaxScaler()
19 scalerResultados_C2.fit(resultados_train_C2.reshape(-1, 1))
20 resultadosTrain_norm_C2 = scalerResultados_C2.transform(resultados_train_C2
   .reshape(-1, 1))
21 resultadosTest_norm_C2 = scalerResultados_C2.transform(resultados_test_C2.
   reshape(-1, 1))

```

Para isso, uma nova escala é criada com o tipo *MinMaxScaler()*, que armazena os limites superior e inferior do conjunto de dados a que se refere (uma para os dados de entrada, e outra para os dados de saída).

Em seguida, as novas variáveis são criadas para armazenar todos esses valores normalizados (dentro de um mesmo intervalo) para a entrada e saída. Elas serão utilizadas para gerar um modelo de regressão linear e, em seguida, será feito um processo para retornar os valores para a escala normal.

Para criar o modelo de regressão linear que será adaptado aos dados do conjunto utilizado, utilizamos o modelo importado anteriormente, o *MLPRegressor*, que criará uma regressão linear utilizando um modelo multicamadas.

Para criar essa regressão linear, deve-se definir o tamanho de cada camada para que a rede consiga performar adequadamente. Testando diversas combinações de funções de ativação e métodos de solução, chegou-se a conclusão de que a combinação da função logística (*logistic*) de ativação e o método *lbfgs* para a solução chega a um resultado satisfatório.

```

1 #Regressao linear para C1 utilizando o MLPRegressor
2
3 modelo_C1 = MLPRegressor(hidden_layer_sizes=20, activation='logistic',
4                           max_iter = 1000, solver="lbfgs")
5
6 #Regressao linear para C2 utilizando o MLPRegressor
7
8 modelo_C2 = MLPRegressor(hidden_layer_sizes=20, activation='logistic',
9                           max_iter = 1000, solver="lbfgs")
10 modelo_C2.fit(dadosTrain_norm_C2, resultadosTrain_norm_C2)

```

Em seguida, utilizam-se os dados normalizados para criar a regressão linear.

Uma nova variável deve ser criada para armazenar os dados previstos pela regressão linear (ainda normalizados) e, em seguida, utilizar uma função de transformação inversa que retornará os valores das variáveis normalizadas para o intervalo inicial (armazenado com o modelo *MinMaxScaler*()).

```

1 resultadosPred_norm_C1 = modelo_C1.predict(dadosTest_norm_C1).reshape(-1,1)
2
3 resultados_pred_C1 = scalerResultados_C1.inverse_transform(
4     resultadosPred_norm_C1)
5
6 resultadosPred_norm_C2 = modelo_C2.predict(dadosTest_norm_C2).reshape(-1,1)
7
8 resultados_pred_C2 = scalerResultados_C2.inverse_transform(
9     resultadosPred_norm_C2)

```

Para mensurar a diferença média com relação aos dados já existentes, usa-se a função *mean\_squared\_error*, que calculará o erro médio quadrático entre os resultados (saídas) dos dados de teste e os que foram previstos pela regressão linear.

```

1 #Calculo do Erro Médio Quadrático para C1
2 erro_C1 = mean_squared_error(resultados_test_C1, resultados_pred_C1)
3 print("O erro médio quadrático de C1 é: %.7f" % np.sqrt(erro_C1))
4 # erro médio quadrático gerado para C1 nessa instância de testes:
5     0.0008116
6
7 #Calculo do Erro Médio Quadrático para C2
8 erro_C2 = mean_squared_error(resultados_test_C2, resultados_pred_C2)
9 print("O erro médio quadrático de C2 é: %.7f" % np.sqrt(erro_C2))
10 # erro médio quadrático gerado para C2 nessa instância de testes: 0.0005617

```

```

10 |
11 | #houve uma melhora nos dois indicadores com relacao a regressao linear do
    | trabalho anterior.
12 | #C1: 0.0058707 -> 0.0008116
13 | #C2: 0.0022257 -> 0.0005617

```

Para visualizar a previsão das duas variáveis de saída ( $C1$  e  $C2$ ), plotam-se as variáveis correspondentes aos resultados previstos com os resultados do conjunto de testes (usa-se o conjunto de testes pois a regressão linear foi criada com base nos dados disponíveis no conjunto dos dados de treino).

```

1 | # Plot do gráfico comparando as informações do resultado 1 (#GT Compressor
    | Decay State Coefficient )
2 |
3 | plt.figure(figsize=(10, 10))
4 |
5 | plt.scatter(resultados_pred_C1, resultados_test_C1, color = "#2E6AB8")
6 | plt.scatter(resultados_pred_C2, resultados_test_C2, color = "#2EFAF8")
7 |
8 | plt.plot(resultados_pred_C1, resultados_pred_C1, color="#1B3E6B")
9 | plt.grid()

```

A reta na diagonal da Figura 2 representa uma função em que  $y = x$ , e na relação em que o eixo vertical corresponde aos dados previstos e o eixo horizontal aos dados do conjunto de teste, dados espalhados seguindo exatamente uma função  $y = x$  representaria 100% de acerto. Os dados em azul escuro correspondem à saída  $C1$ , e os em azul claro representam a saída  $C2$ .

Os dados de  $C2$  estão corretamente dispostos em um intervalo menor, pois os registros para o coeficiente de decaimento da turbina foram medidos em um domínio menor.

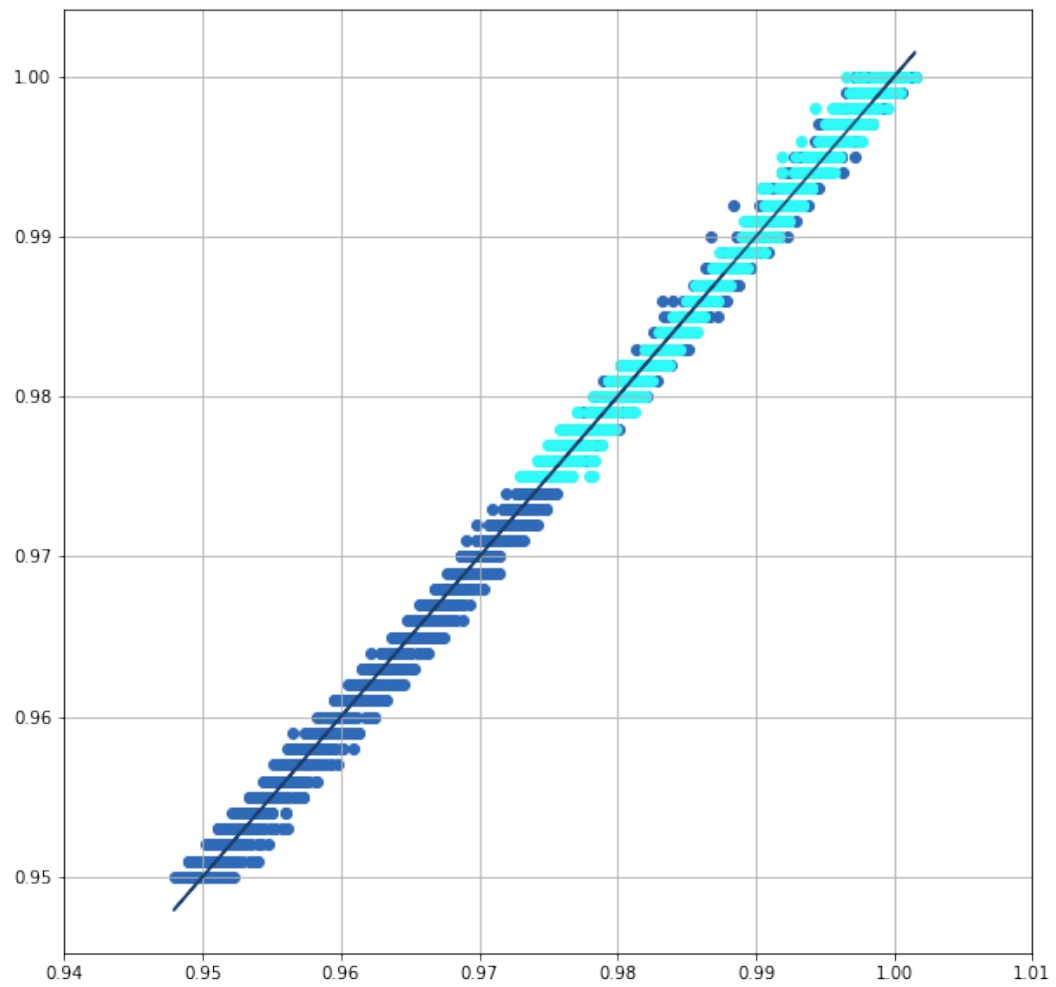


Figura 2 – Gráfico da saída geral

### 3 RESULTADOS E DISCUSSÕES

Para analisar o desempenho do modelo criado (com regressão linear), deve-se obter os dados sobre a raiz do erro médio quadrático.

#### 3.1 ANÁLISE DA PREVISÃO PARA O COEFICIENTE DE DECAIMENTO DO COMPRESSOR

Para a primeira saída ( $C1$ ), a raiz do erro médio quadrático foi de 0.0008116. O resultado foi ainda melhor do que a regressão linear criada utilizando um modelo simples de regressão, o *LinearRegression*, que apresentou a raiz do erro médio quadrático de 0.0058707.

Considerando que o erro médio quadrático pode ser interpretado como a diferença entre o valor previsto e o parâmetro já existente, deve-se considerar que a escala (domínio) de existência para as saídas é considerado preciso até a faixa de  $10^{-3}$ . Assim, o erro médio quadrático obtido com a regressão linear pode ser considerado bom, pois indica que as diferenças entre os valores são, em média, pouco impactantes para o resultado (numericamente) por indicarem variações na faixa de  $10^{-4}$ .

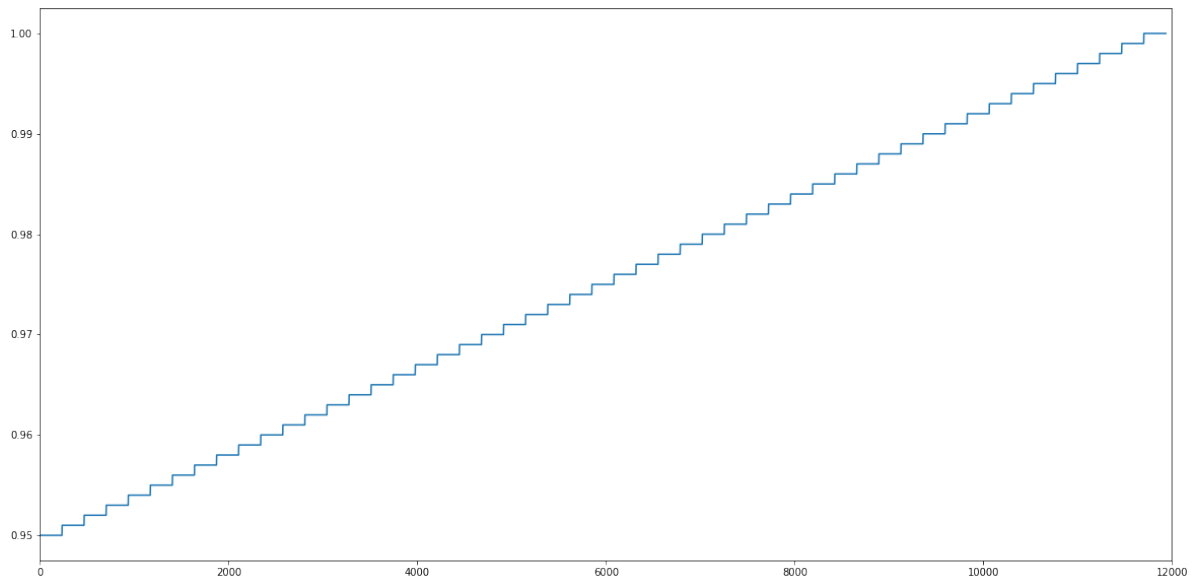


Figura 3 – Gráfico com os valores de  $C1$  obtidos no conjunto de dados

Analisando a Figura 3 é possível perceber que ela representa uma excelente aproximação de reta. Se o ruído for desconsiderado, a taxa de crescimento seria muito bem interpretada como constante. Também graças à isso, o valor da raiz do erro médio quadrático é menos satisfatório que o que pode ser verificado a partir de uma análise do segundo caso (para  $C2$ ). Isso acontece pois o valor da média já pode ser considerado como válido por se tratar de uma reta (e a regressão linear também formularia uma reta).

### 3.2 ANÁLISE DA PREVISÃO PARA O COEFICIENTE DE DECAIMENTO DA TURBINA

Para a segunda saída ( $C2$ ), o erro médio quadrático foi de 0.0005617. Os resultados para a segunda saída também sofreram uma melhora perceptível com relação ao resultado obtido com o modelo baseado no modelo do tipo *LinearRegression*, quando o erro médio quadrático havia sido de 0.0058707.

Seguindo a mesma justificativa dada para o para a dimensão do coeficiente em  $C1$ , o erro médio quadrático obtido com a regressão linear pode ser considerado bom.

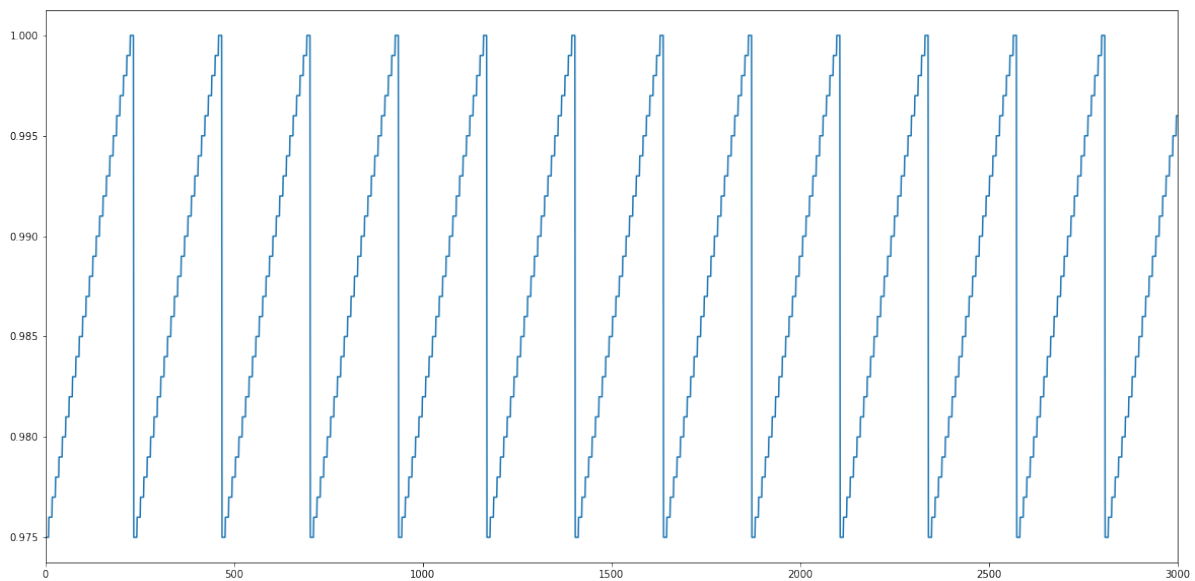


Figura 4 – Gráfico com os valores de  $C2$  obtidos no conjunto de dados

Analisando a Figura 4 é possível perceber que ela representa uma função periódica. Mesmo que contra-intuitivamente, a raiz do erro médio quadrático foi ainda melhor para esse segundo caso. Isso deve acontecer pois, apesar de se tratar de uma função periódica (que se repete ciclicamente), ainda apresenta uma taxa de variação constante com relação à subida, e praticamente uma descontinuidade com relação ao retorno ao ponto inicial (0.975), por retornar quase que instantaneamente. Sendo assim, os valores oscilam em um domínio muito menor (50% do domínio de  $C1$ , que era de 0.950 até 1), explicando o motivo da variação ser menor para o segundo coeficiente do que para o primeiro.

Como uma regressão linear formula uma reta e a função do coeficiente de decaimento do motor também pode ser aproximada de uma reta (quando desconsidera-se o ruído), é compreensível que a diferença entre os valores previstos sejam menores para o caso em que o domínio de variação para a previsão do coeficiente de rendimento é menor (que é o caso da variação do  $C2$ ).



## 4 CONCLUSÃO

Baseando-se nos resultados e valores evidenciados pelos gráficos plotados, é possível afirmar que a regressão linear criada pode ser considerada satisfatória, e ainda melhor nessa aplicação utilizando o modelo *MLPRegressor*. Mesmo com uma pequena variação entre as regressões para as saídas *C1* e *C2*, as duas apresentaram valores com pouco desvio e, realizando a avaliação com o cálculo da raiz do erro médio quadrático, podem ser consideradas satisfatórias.