

Aplicação de Regressão Linear a Série Temporal*

Tiago C A Amorim (RA: 100675)^a, Taylon L C Martins (RA: 177379)^b

^aDoutorando no Departamento de Engenharia de Petróleo da Faculdade de Engenharia Mecânica, UNICAMP, Campinas, SP, Brasil

^bAluno especial, UNICAMP, Campinas, SP, Brasil

Keywords: Regressão Linear, Séries Temporais, Validação Cruzada

1. Introdução

Este relatório apresenta as principais atividades realizadas no desenvolvimento das atividades propostas na Lista 01 da disciplina IA048: Aprendizado de Máquina, primeiro semestre de 2024. O foco deste exercício é de construir modelos lineares para realizar a previsão de uma série temporal.

2. Tarefa Proposta

Trabalhar com a base de dados U.S. Airline Traffic Data, a qual contém informações referentes ao tráfego aéreo mensal norte-americano no período de 2003 a 2023, disponibilizadas pelo *U.S. Department of Transportation's (DOT) Bureau of Transportation Statistics*. Em particular, vamos explorar a série temporal do número total de voos (domésticos e internacionais).

Explorar um modelo linear para a previsão considerando que o horizonte de predição é $L = 1$ (passos à frente da série temporal).

(a) Exiba o gráfico da série temporal completa. Numa inspeção visual simples, é possível reconhecer ao menos três faixas distintas de comportamento aproximadamente “regular” na série:

- Jan/2003 a Ago/2008.
- Set/2008 a Dez/2019.
- Jan/2020 a Set/2023.

Discuta possíveis razões históricas/econômicas para as transições de comportamento.

(b) Divida a série em dois conjuntos:

- Treinamento** e **validação**: com amostras de 2003 a 2019.
- Teste**: com amostras de 2020 a 2023.

Faça a análise de desempenho do preditor linear ótimo, no sentido de quadrados mínimos irrestrito considerando:

- (1) A progressão do valor da raiz quadrada do erro quadrático médio (RMSE, do inglês *root mean squared error*), junto aos dados de validação, em função do número de entradas (**K**) do preditor (desde $K = 1$ a $K = 24$). Apresente o gráfico obtido e busque tecer conjecturas sobre os motivos subjacentes a seu comportamento.
 - (2) O gráfico com as amostras de teste da série temporal e as respectivas estimativas geradas pela melhor versão do preditor (i.e., usando o valor de **K** que levou ao mínimo erro de validação). Obtenha, também, o RMSE e o erro percentual absoluto médio (MAPE, do inglês *mean absolute percentage error*) para o conjunto de teste.
 - (3) O gráfico com as amostras apenas dos dois últimos anos (2022 e 2023) e as estimativas geradas pelo melhor preditor, além dos respectivos valores de RMSE e MAPE.
- (c) Repita o procedimento detalhado nos itens b1 e b2, mas adotando a seguinte divisão dos dados:
- (i) **Treinamento**: amostras de 2003 a 2019.
 - (ii) **Validação**: amostras de 2020 e 2021.
 - (iii) **Teste**: amostras de 2022 e 2023.
- Discuta os resultados obtidos e faça uma comparação com o cenário anterior (especialmente com o que foi obtido no item b3).

3. Aplicação

Toda a avaliação foi feita em um único *notebook* Jupyter, em Python. Foi feito o uso da biblioteca *Scikit-learn* [1] para fazer as diferentes manipulações nos dados. O código pode ser encontrado em https://github.com/TiagoCAAmorim/machine_learning.

3.1. Avaliação do Conjunto de Dados

O conjunto de dados utilizado nesta avaliação é o de tráfego aéreo dos EUA disponibilizado no Kaggle [2]. Segundo o autor este conjunto de dados fornece o tráfego aéreo mensal dos EUA de 2003 a 2023, incluindo o número de passageiros, o número de voos, as milhas de passageiros pagantes, as milhas de assentos disponíveis e o fator de ocupação.

*Relatório número 01 como parte dos requisitos da disciplina IA048: Aprendizado de Máquina.

Nesta avaliação o foco é no número total de voos. A figura 1 mostra que o período de janeiro/2003 a agosto/2008 (azul) é de certa normalidade na economia americana. O início do período em verde (setembro/2008) é marcado pelo início da crise financeira do *subprime*. Esta crise foi consequência do estouro da bolha imobiliária nos EUA devido a empréstimos sem lastro e preços de imóveis inflacionados. O período em vermelho corresponde à crise sanitária causada pela pandemia de Covid-19 e os anos que se seguiram. A queda mais significativa, que inicia em janeiro/2020, corresponde ao período de maior restrição de voos imposta durante a crise sanitária.

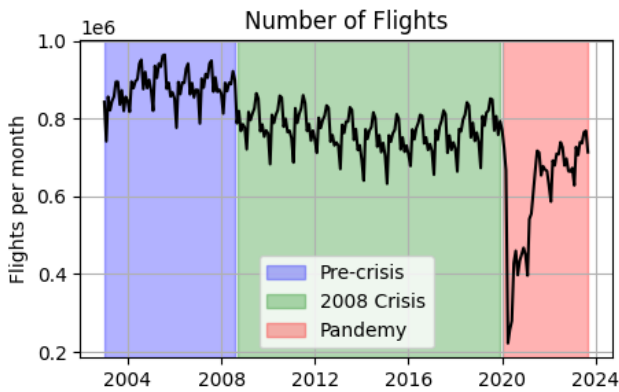


Figura 1: Série com total de voos nos EUA.

3.1.1. Pré-processamento

Alguns exemplos na figura 2 mostram que a série original apresenta um caráter cíclico (com frequência anual) e *serrilhado* (não-suave). Foram avaliadas duas alternativas para deixar a série temporal mais suave e facilitar o processo de regressão. Como os meses do ano não tem o mesmo número de dias, uma primeira tentativa foi de calcular o número médio de voos por dia para cada mês da série. A segunda tentativa foi de calcular o número de voos diários considerando apenas os dias úteis de cada mês.

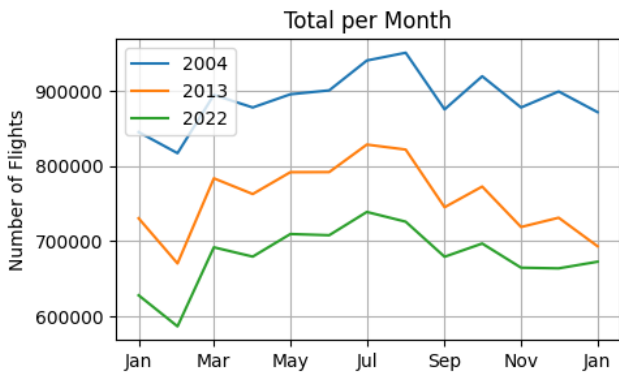


Figura 2: Exemplos de total de voos nos EUA ao longo do ano.

A figura 3 mostra que a série temporal de média diária de voos por mês se mostra mais suave que a série original.

O efeito é mais pronunciado no mês de fevereiro. A série temporal de total mensal de voos dividido pelo número de dias úteis, na figura 4, não teve uma resposta adequada.

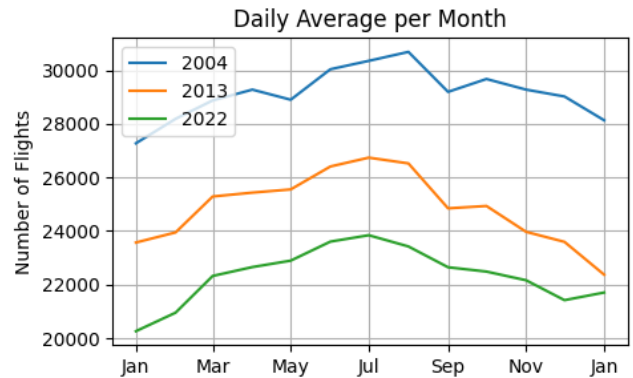


Figura 3: Exemplos de média diária do total de voos nos EUA.

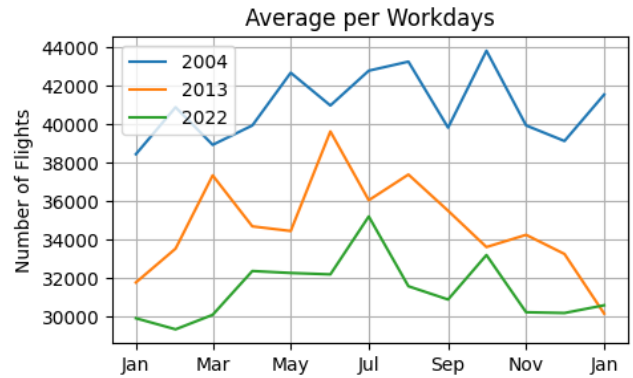


Figura 4: Exemplos do total de voos nos EUA pelo número de dias úteis no mês.

Optou-se por utilizar a média mensal de voos por dia nas análises posteriores (figura 5). Todos os valores de RMSE e MAPE foram calculados em função desta nova série.

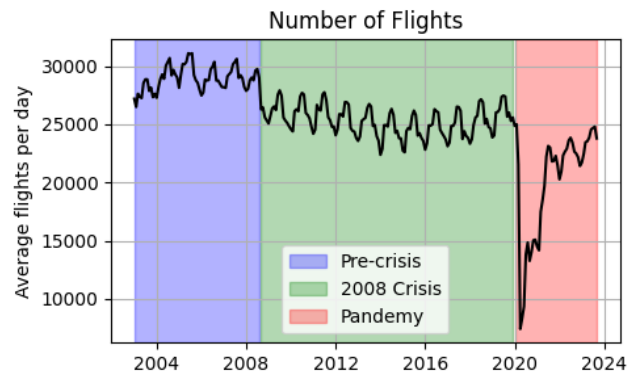


Figura 5: Média diária do total de voos nos EUA.

3.2. Primeiro Modelo de Regressão

Antes de separar os dados em treino+validação e teste, foram construídas séries temporais com os valores passados associado a cada entrada na série de voos¹. Definindo a série de média diária de voos como \mathbf{y} , as séries dos valores passados foram definidas como:

$$\mathbf{x}_k(t_i) = \mathbf{y}(t_{i-k}) \quad \text{com } k = 1, \dots, 24 \quad (3.1)$$

Onde t_i são os índices dos dados da série temporal.

De modo a ter uma melhor comparação entre os diferentes modelos que serão construídos, foram retiradas todas as linhas com valores não definidos (as 24 primeiras entradas da tabela de dados). Desta forma todos os processos de treino e validação são feitos com os mesmos conjuntos de dados. Com este filtro o total de dados passa de 249 para 225.

Aplicando os limites propostos no item b1, o conjunto de dados de treino+validação fica com 180 elementos e o de teste com 45. Os intervalos de dados de treino+validação e teste apresentam um comportamento cíclico, com suas variações ao longo dos anos (figura 6). A exceção no conjunto de treino+validação é no início de 2008, quando eclodiu a crise financeira nos EUA. O início do intervalo de teste corresponde ao período crítico de restrição de vôos da pandemia de Covid-19. É esperado que o modelo apresente dificuldades em prever a série temporal no período mais crítico da pandemia.

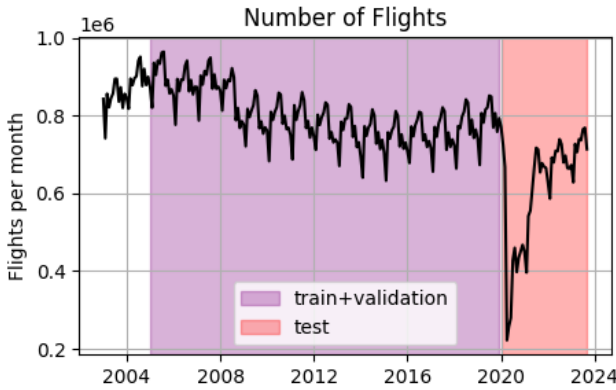


Figura 6: Períodos de treino+validação e teste.

3.2.1. Normalização dos Dados

Como todas as variáveis de entrada dos modelos lineares são de mesma natureza, uma modificação na escala dos dados não tem qualquer impacto nos resultados dos ajustes dos modelos lineares. Foi realizado um rápido teste assumindo um modelo linear sem regularização. Como conjunto de teste foram utilizados os primeiros 140 dados do conjunto de treino+validação, e os demais como conjunto

de validação. Foram utilizados \mathbf{x}_1 e \mathbf{x}_2 como variáveis do problema.

A tabela 1 mostra que, mesmo gerando modelos distintos, os resultados de RMSE são os mesmos quando as variáveis de entrada estão no intervalo $[0; 1]^2$ e quando estão com os valores originais.

RMSE	Valores Originais	Com Mudança de Escala
Treino	788.91	788.91
Validação	854.33	854.33

Tabela 1: Efeito da normalização dos dados de entrada.

Durante as iterações da rotina de busca do melhor modelo linear o algoritmo teve dificuldades de convergência com valores altos de \mathbf{K} (>15). Ao modificar a escala dos dados de entrada o algoritmo não apresentou problemas. Deste modo, as análises seguintes foram feitas com os dados de entrada escalados para o intervalo $[0; 1]$.

3.2.2. Validação Cruzada

Como existe dependência entre os dados observados, o processo de validação cruzada *clássico*, que define os conjuntos de treino e validação sem levar em conta a sua ordem, não pode ser utilizado. Os dados de validação ($I(v)$) são definidos em índices posteriores aos dos dados de treinamento ($I(t)$) [3]:

$$\min_{i \in I(t), j \in I(v)} |i - j| > h > 0 \quad (3.2)$$

Foi utilizada a rotina **TimeSeriesSplit**, implementada no pacote *Scikit-learn*, e um total de 4 pastas (*folds*). A rotina do *Scikit-learn* equivale a utilizar $h = 1$. A figura 7 mostra a divisão dos dados entre treino e validação para cada pasta.

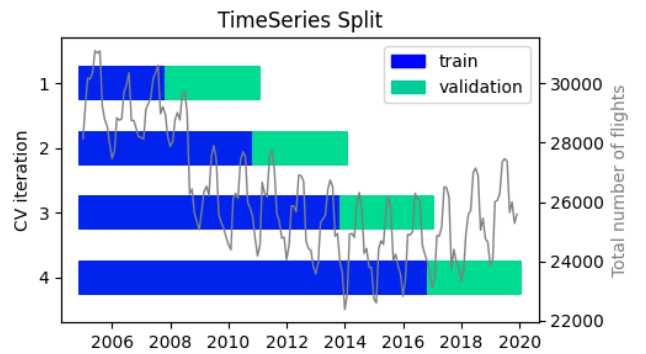


Figura 7: Divisão entre treino e validação para cada pasta.

3.2.3. Modelos Simples

Antes de seguir com o teste de diversos modelos lineares, foram propostos três modelos simples (*naïve*):

¹Ao longo destes tópicos a série que se busca ajustar é a da média diária de voos.

²Foi utilizada a classe **MinMaxScaler** do *Scikit-learn*

1. Igual ao passo de tempo anterior:
 $\hat{y}(t) = y(t - 1)$
2. Igual ao passo de tempo 12 meses atrás:
 $\hat{y}(t) = y(t - 12)$
3. Igual ao passo de tempo 12 meses atrás mais variação entre 24 e 12 meses atrás:
 $\hat{y}(t) = y(t - 12) + [y(t - 12) - y(t - 24)]$

Os modelos *naïve* foram propostos com base na natureza cíclica da série temporal em estudo. Os valores de RMSE e MAPE destes modelos servirão como referência para os demais modelos a serem construídos.

Foi aplicada validação cruzada a cada um dos modelos e computada a média dos RMSE. Os resultados de RMSE são relativamente próximos: 861.25, 765.83 e 950.38, para o primeiro, segundo e terceiro modelos *naïve* propostos, respectivamente. Em uma inspeção gráfica (figura 8) a primeira proposta de modelo parece ser melhor que as outras duas. Uma avaliação do gráfico do erro cometido por estes modelos (figura 9) demonstra que na verdade os modelos tem níveis de erro parecidos, mas com diferentes frequências.

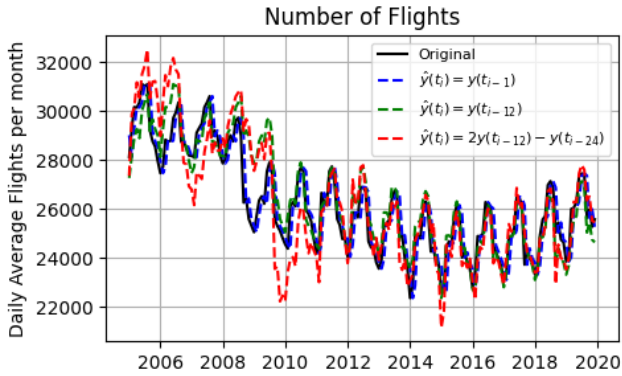


Figura 8: Modelos simples (*naïve*).

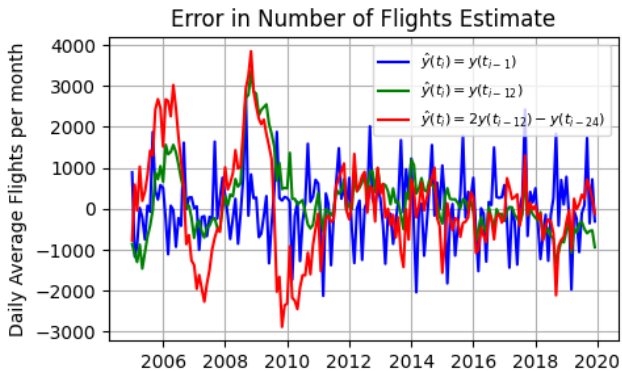


Figura 9: Erros dos modelos simples.

O segundo modelo *naïve* tem o menor valor de RMSE médio, e apresenta um gráfico de erro melhor *comportado* que os demais. Este segundo modelo foi o escolhido para ser utilizado como referência na busca pelo modelo *ótimo*.

3.2.4. Busca pelo Melhor Modelo

Em um primeiro momento foi definido utilizar um modelo linear regularizado do tipo *Elastic-Net*:

$$\frac{1}{2n} \|y - Xw\|_2^2 + \alpha l_{ratio} \|w\|_1 + \frac{\alpha}{2} (1 - l_{ratio}) \|w\|_2^2 \quad (3.3)$$

Onde n é o número de amostras e w são os pesos do modelo linear.

A busca pelo melhor modelo para cada valor de K foi feita com uma busca em grade (*grid search*) por parâmetros ótimos de α (entre 0.01 e 1.0) e l_{ratio} (entre 0.1 e 0.9). O parâmetro de busca do melhor modelo foi a média do RMSE.

Todos os 24 modelos encontrados utilizaram o menor valor de α que fez parte da busca em grade (0.01). Em face destes resultados a avaliação foi refeita utilizando um modelo linear simples (sem regularização). Como o novo modelo não tem parâmetros, não foi preciso realizar uma busca em grade, e foi testado apenas o valor de K (figura 10).

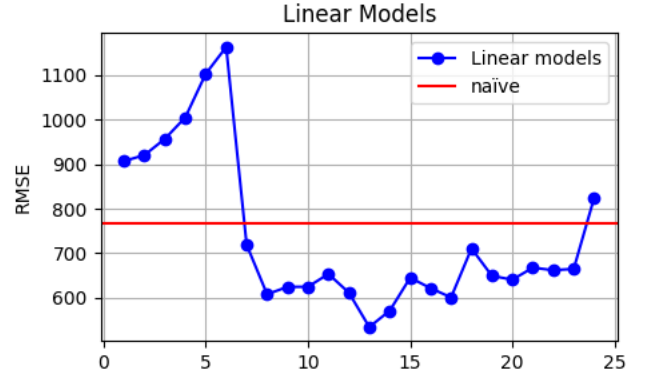


Figura 10: RMSE de validação médio dos modelos lineares em função de K .

O melhor resultado foi com K igual a 13 (validação cruzada com RMSE = 534.25 e MAPE = 0.02). Com o melhor valor de K definido, foi utilizado todo o conjunto de dados de treino+validação para ajustar os parâmetros do modelo escolhido.

Foi ajustado um segundo modelo com K igual a 13, mas sem normalizar as variáveis. Com este modelo fica mais fácil analisar os valores dos coeficientes do modelo. O modelo linear pode ser interpretado como uma média ponderada dos valores nos meses anteriores, já que as variáveis e o valor de saída tem todos a mesma natureza. A soma dos coeficientes com o valor relativo da constante³ é próxima da unidade: $0.9687 + 0.0305 = 0.9992$.

A figura 11 mostra os valores dos coeficientes utilizados. As variáveis de maior peso foram o passo de tempo anterior (x_1), um ano antes (x_{12}) e o anterior (x_{13}). Este

³Foi definido valor relativo da constante como a divisão do termo constante do modelo linear pela média dos valores de saída.

resultado guarda certa correspondência com dois dos modelos *naïve* propostos, dando maior peso ao valor anterior (tendência de curto prazo) e ao valor um ano atrás (efeito da sazonalidade).

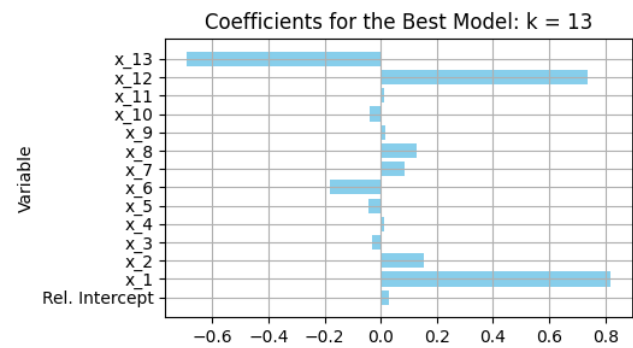


Figura 11: Coeficientes do melhor modelo linear ($K = 13$).

3.2.5. Erros com os Dados de Teste

O modelo construído teve um resultado ruim nos dados de teste, com $RMSE = 3\,187$. e $MAPE = 0.1126$. Na figura 12 observa-se que o conjunto de teste (área em vermelho na figura) coincide com o período da pandemia e os meses posteriores. Por se tratar de um evento singular e que não apareceu nos dados de treino e validação, já era esperado que o modelo tivesse dificuldade em prever o comportamento da série temporal. O alto valor de RMSE associado ao conjunto de teste demonstra a dificuldade do modelo.

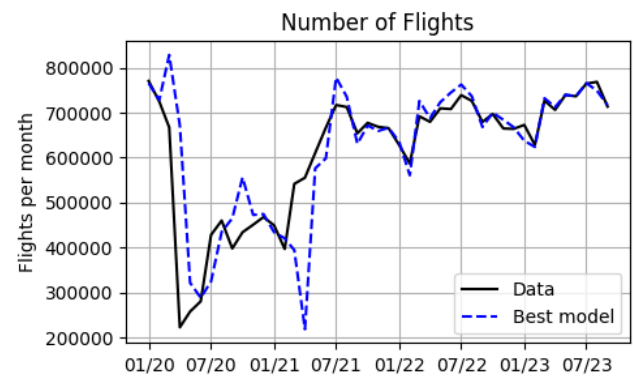


Figura 12: Ajuste do melhor modelo linear ($K = 13$).

Quando o período de teste é limitado aos anos de 2022 e 2023 os resultados são bem melhores ($RMSE = 575.85$ e $MAPE = 0.0190$). O período de 2022 a 2023 corresponde ao final da pandemia de Covid-19, quando muitas das barreiras sanitárias já estavam sendo levantadas. O comportamento da série temporal neste intervalo já se assemelha mais ao comportamento observado nos dados de treino e validação (figura 13). Um valor de RMSE mais próximo do observado no processo de validação cruzada

mostra que o modelo teve um resultado aceitável para este período (figura 14).

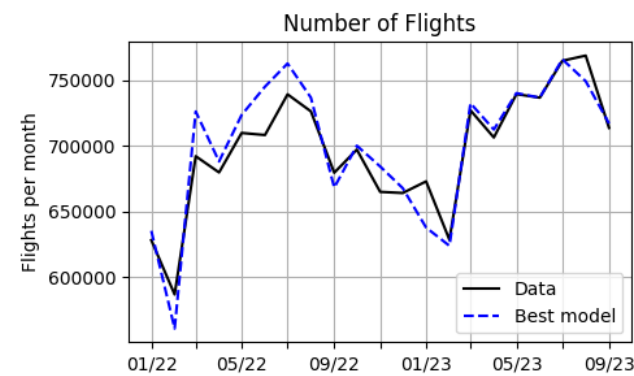


Figura 13: Ajuste do melhor modelo linear ($K = 13$) ao anos de 2022 e 2023.

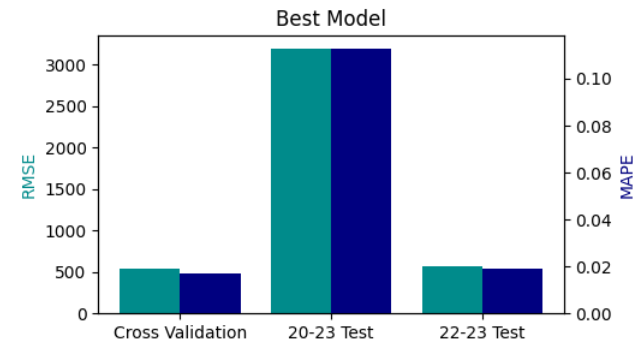


Figura 14: Resumo dos erros do melhor modelo linear ($K = 13$).

3.3. Segundo Modelo de Regressão

Para o segundo modelo de regressão a proposta é utilizar outros intervalos para treino, validação e teste (figura 15). O intervalo de validação coincide com o período de maior restrição de deslocamento da pandemia de Covid-19. Dada a discrepância entre o comportamento da série temporal no período de treino e no de validação, é esperado que as variáveis que correspondem a dados mais recentes tenham maior peso. Adicionalmente, como o intervalo de validação está fixo, não será possível fazer validação cruzada.

Foram testados os mesmos modelos simples (*naïve*) propostos anteriormente (em 3.2.3). O modelo de melhor resultado foi o que utiliza o valor do mês anterior para prever o próximo intervalo de tempo. Os resultados de RMSE foram 3304., 11023. e 17523., para o primeiro, segundo e terceiro modelos *naïve* propostos, respectivamente. Como o intervalo de validação corresponde a um período de significativas variações na série temporal, os modelos que se baseiam em horizontes de tempo mais longo tem mais dificuldade de prever este comportamento.

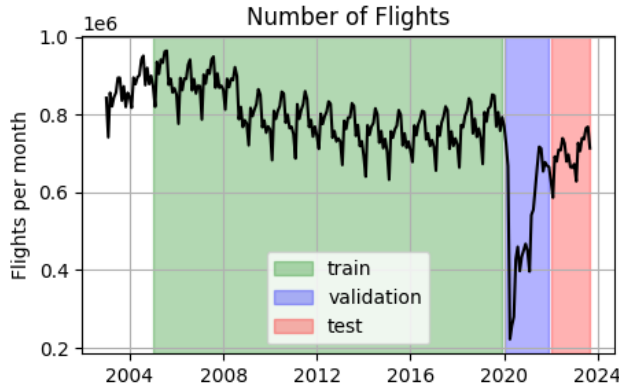


Figura 15: Série temporal com novos limites dos dados de treino, validação e teste.

O mesmo comportamento observado entre os modelos *naïve* apareceu na seleção de modelos em função da quantidade de passos de tempo passados (figura 16). Os modelos que se restringem a dados de passos de tempo mais recentes apresentaram melhores resultados. Este resultado é devido ao conjunto de dados de validação utilizado. O melhor modelo é o que se baseia apenas nos dois passos de tempo anteriores, e com maior impacto do passo de tempo imediatamente anterior.

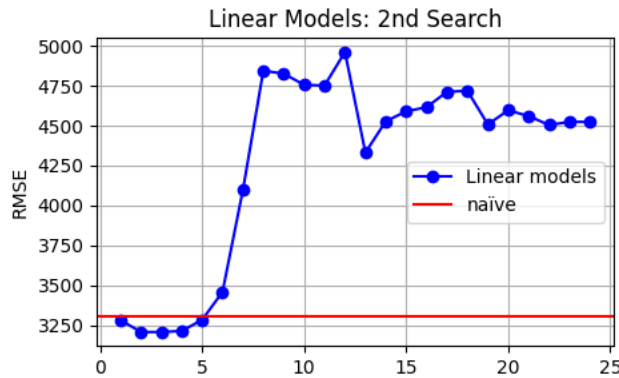


Figura 16: RMSE de validação dos diferentes modelos testados com o novo conjunto de dados.

Foi aplicado o mesmo processo de reconstruir o modelo *ótimo* (com $K = 2$) em todo o conjunto de treino+validação. Também foi construída uma versão deste modelo sem normalizar as variáveis do problema, de forma a melhor analisar os seus coeficientes (figura 17). Fica claro que o modelo é próximo do modelo *naïve*, com um peso grande do valor imediatamente anterior, acrescido de um uma *correção* do valor dois meses atrás. Novamente o modelo pode ser interpretado como um média ponderada. A soma dos coeficientes com o valor relativo da constante é aproximadamente a unidade: $0.9211 + 0.0778 = 0.9989$. Como esperado, este novo modelo se ajusta melhor aos dados no período da pandemia (figura 18).

Os dados de teste deste novo modelo correspondem aos

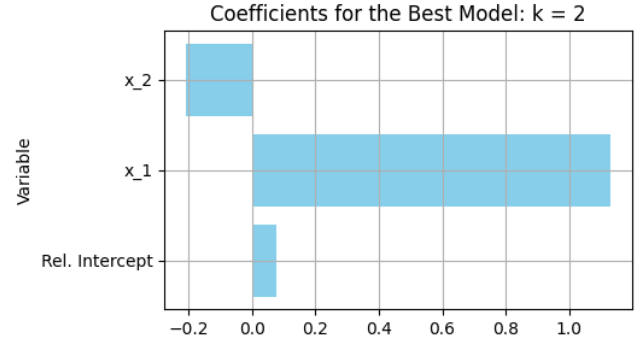


Figura 17: Coeficientes do novo melhor modelo linear ($K = 2$).

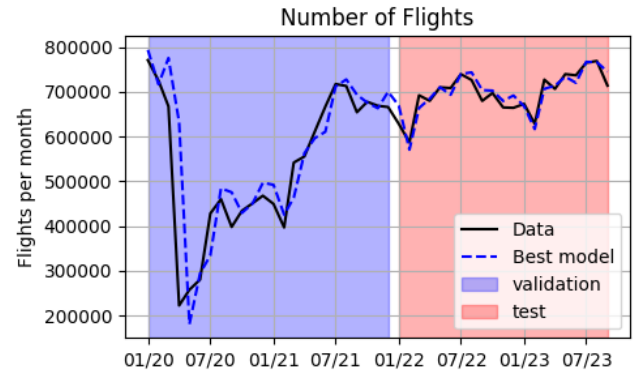


Figura 18: Ajuste do novo melhor modelo linear ($K = 2$).

anos de 2022 e 2023, com valores de erro RMSE = 590.57 e MAPE = 0.0208. Para os anos de 2022 e 2023 os dois modelos construídos mostraram resultados próximos (figuras 19 e 20), com o primeiro modelo *ligeiramente* melhor..

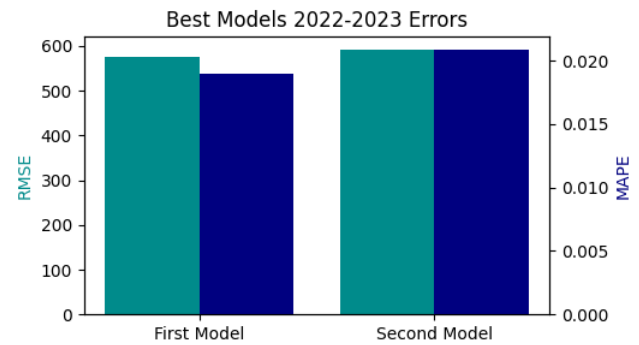


Figura 19: Erros para os dois modelos construídos.

3.4. Extrapolação Recursiva

Para testar a capacidade preditiva destes modelos para um prazo maior, foi feita uma extrapolação recursiva nos períodos de teste, a partir do último valor do período de validação.

Como esperado, o primeiro modelo não consegue prever o comportamento da série durante a pandemia (figura 21).

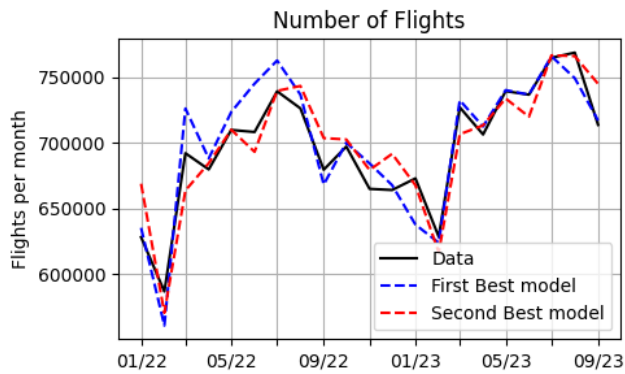


Figura 20: Ajuste em 2022 e 2023 dos dois modelos construídos.

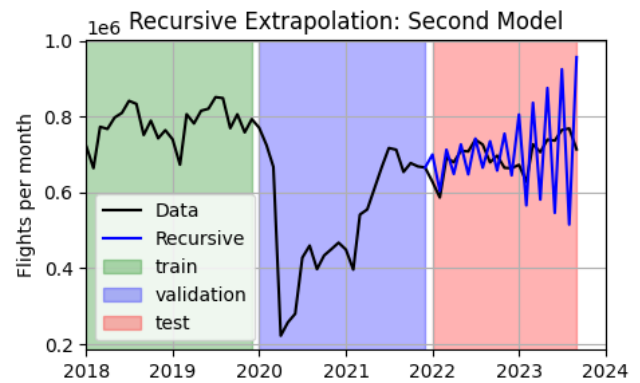


Figura 22: Extrapolação recursiva com o segundo modelo.

O modelo apresenta um comportamento *relativamente* estável nesta previsão recursiva. Depois de cerca de um ano o modelo começa a apresentar um comportamento errático, o que coincide com o passo de tempo em que realiza previsões baseadas apenas em dados estimados pelo próprio modelo.

O segundo modelo (figura 22), que se baseia unicamente nos dois valores mais recentes, rapidamente perde a capacidade de realizar qualquer previsão confiável.

Nenhum dos modelos pode ser utilizado para realizar previsões além do mês seguinte. É necessário realizar uma ajuste específico para valores de L maiores que um.

Referências

- [1] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, E. Duchesnay, Scikit-learn: Machine learning in Python, Journal of Machine Learning Research 12 (2011) 2825–2830.
- [2] YYXian, U.s. airline traffic data (2003-2023), acessado: 22/03/2024 (Jan 2024).
URL <https://www.kaggle.com/datasets/yyxian/u-s-airline-traffic-data/data>
- [3] S. Arlot, A. Celisse, A survey of cross-validation procedures for model selection, Statistics Surveys 4 (none) (2010) 40 – 79. doi: 10.1214/09-SS054.
URL <https://doi.org/10.1214/09-SS054>

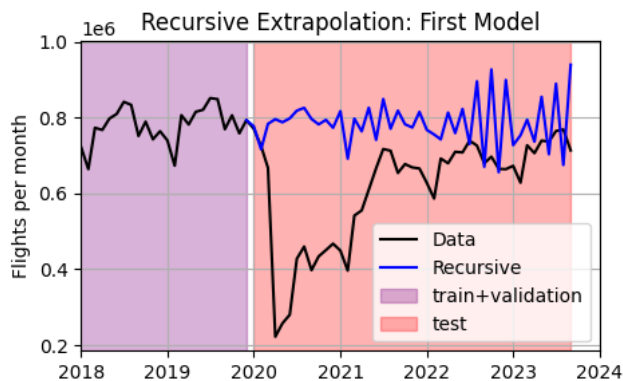


Figura 21: Extrapolação recursiva com o primeiro modelo.

4. Conclusão

As características dos períodos de treino e validação tiveram impacto no número de parâmetros de cada modelo. A definição dos períodos de teste e validação deve estar alinhada com o que se espera do modelo a ser construído.

Ambos modelos dependem muito do valor do mês anterior para estimar o próximo, e são incapazes de responder adequadamente a eventos com variações abruptas. Variáveis adicionais são necessárias para permitir que os modelos lineares consigam aproximar melhor eventos como as restrições impostas pela Pandemia de Covid-19.