

# Previsão de Irradiância Solar de Curto Prazo Utilizando Modelo de Envelopes para os Preditores

Felipe P. Marinho \* Ajalmar R. R. Neto \*\* Paulo A. C. Rocha \*\*\*

\* Programa de Pós-Graduação em Engenharia de Teleinformática,  
Universidade Federal do Ceará, CE, (e-mail: fpmarinho@alu.ufc.br).  
\*\* Programa de Pós-Graduação em Ciências da Computação, Instituto  
Federal do Ceará, CE (e-mail: ajalmar@gmail.br)  
\*\*\* Departamento de Engenharia Mecânica, Universidade Federal do  
Ceará, CE (e-mail: paulo.rocha@ufc.br)

Português

---

**Abstract:** In this work, we employ the recent envelope estimation method to fit a multiple linear regression model in order to make predictions of solar irradiance at horizons of 10, 20 and 30 min ahead. The advantage of using this method lies in the fact that there is a reduction in the variance of the estimated coefficients. Such a reduction is obtained by enveloping the material part, while excluding the variation of the immaterial part, which also contributes to a reduction in dimensionality, since there is a decrease in the estimation variance without an increase in the number of observations. The model was used in a dataset formed by LDR luminosity signals and statistical attributes extracted from images of the sky. The integration of the LDRs and the camera was done using a *Raspberry Pi 3*. Its performance was compared to that of linear (LASSO and Ridge) and non-linear (Multi-Layer Perceptron) models, being evaluated by the metrics Mean Absolute Error (MAE), Root Mean Square Error (RMSE) and Root Mean Square Error Relative (rRMSE).

**Resumo:** Neste trabalho, empregamos o recente método de estimação por envelopes para ajustar um modelo de regressão linear múltipla com o objetivo de realizar previsões de irradiância solar nos horizontes de 10, 20 e 30 min a posteriori. A vantagem do uso deste método está no fato de se ter uma redução na variância dos coeficientes estimados. Tal redução é obtida pelo envelopamento da parte material, enquanto se exclui a variação da parcela imaterial, o que também contribui para uma redução de dimensionalidade, uma vez que se tem uma queda na variância da estimação sem um aumento no número de observações. O modelo foi empregado em um conjunto de dados formados por sinais de luminosidade LDR e atributos estatísticos extraídos de imagens do céu. A integração dos LDR's e da câmera foi feita pelo uso de um *Raspberry Pi 3*. Seu desempenho foi comparado ao de modelos lineares (LASSO e Ridge) e não-linear (Perceptron de Múltiplas Camadas) sendo a avaliação feita pelas métricas Erro Médio Absoluto (MAE), Erro Médio por Viés (MBE), Raiz do Erro Quadrático Médio (RMSE) e Raiz do Erro Quadrático Médio Relativo (rRMSE).

**Keywords:** Envelope Estimation; Multiple Linear Regression; Machine Learning; Solar energy; Solar Irradiance Forecast

**Palavras-chaves:** Estimação por Envelopes; Regressão Linear Múltipla; Aprendizagem de Máquina; Energia Solar; Previsão de Irradiância Solar.

---

## 1. INTRODUÇÃO

A substituição de fontes tradicionais de energia por renováveis como a solar e eólica é imprescindível para o desenvolvimento sustentável das nações. A região Nordeste brasileira têm grande potencial para o aproveitamento de tais fontes alternativas, eólica no litoral e solar no interior.

\* Os autores agradecem à Fundação Cearense de Apoio ao Desenvolvimento Científico e Tecnológico (FUNCAP) e ao projeto *Apple Developer Academy - IFCE* pelo suporte financeiro dado para a realização da pesquisa.

No que diz respeito a fonte solar, seu caráter intermitente e estocástico caracteriza-se por ser um entrave para o seu melhor aproveitamento em diversas aplicações, portanto, um sistema de previsão de irradiância solar contribui positivamente para uma melhor compreensão da variabilidade de tal recurso e, com isso, para uma melhor inserção da energia solar na matriz energética global.

O desenvolvimento de um dispositivo que forneça previsões confiáveis de irradiância solar é de grande valia para a realização de cultivo inteligente, agricultura de precisão e monitoramento meteorológico. Além de sua importância

no contexto agrário, pode-se destacar também sua relevância no aproveitamento de energia solar para geração de eletricidade para a rede, fornecendo ao operador informações úteis para prever quedas de fornecimento e auxiliar na programação de manutenções preventivas.

Neste sentido, muitas metodologias têm sido utilizadas para que previsões de irradiância em diferentes horizontes temporais sejam efetuadas. Uso de modelos *Numerical Weather Prediction* (NWP) (Nonnenmacher et al., 2016; Mejia et al., 2018), análise de séries temporais (Dong et al., 2013; Trapero et al., 2015), aplicações envolvendo processamento de imagens (Pedro et al., 2015; Pedro et al., 2018; Pawar et al., 2019) e utilização de algoritmos de aprendizagem de máquina (Koo et al., 2019; Yagli et al., 2019; Benali et al., 2019) destacam-se como abordagens comumente aplicadas para a obtenção de previsões de irradiância solar.

Estimação por Envelopes é um método recente empregado em regressão linear múltipla para se ter uma redução da variância na estimação dos coeficientes de regressão sem que seja necessário um aumento no número de observações de treinamento. Desde sua proposição muitos trabalhos foram desenvolvidos, principalmente com o intuito de melhorar a eficiência computacional da etapa de otimização do modelo (COOK et al., 2010).

A metodologia inicial realizava a otimização sobre uma variedade grassmaniana matricial, o que proporcionava um alto custo de processamento dependendo da dimensão dos envelopes empregados. Um novo algoritmo desenvolvido (Algoritmo 1D) (Cook et al., 2016) segmentava o problema original em vários problemas sobre variedades de dimensão 1, permitindo menor custo sem perda de desempenho.

A versão mais recente do método de estimação por envelopes utiliza um algoritmo que não requer otimização sobre uma grassmaniana. As simulações indicam maior velocidade e, tipicamente, melhor acurácia do que as demais abordagens (Cook et al., 2016). Por mais que muitos avanços tenham sido feitos no método, percebe-se uma carência de trabalhos que envolvam aplicações desta metodologia.

Assim, empregamos este modelo para realizar previsões de irradiância solar nos horizontes de 10, 20 e 30 minutos a posteriori por meio de sua aplicação em um conjunto de dados formados por sinais de luminosidade LDR e atributos estatísticos extraídos de imagens do céu.

A integração dos LDR's e da câmera foi feita pelo uso de um *Raspberry Pi 3*. Seu desempenho foi comparado ao de modelos lineares (LASSO e Ridge) e não-linear (Percéptron de Múltiplas Camadas) sendo a avaliação feita pelas métricas Erro Médio Absoluto (MAE), Erro Médio por Viés (MBE), Raiz do Erro Quadrático Médio (RMSE) e Raiz do Erro Quadrático Médio Relativo (rRMSE).

Dentre as contribuições do trabalho, destacam-se:

- Uso do recente modelo de estimação por envelopes para realizar previsões de irradiância solar;
- Comparação de seu desempenho com o de modelos lineares e não-lineares clássicos e bem postos;

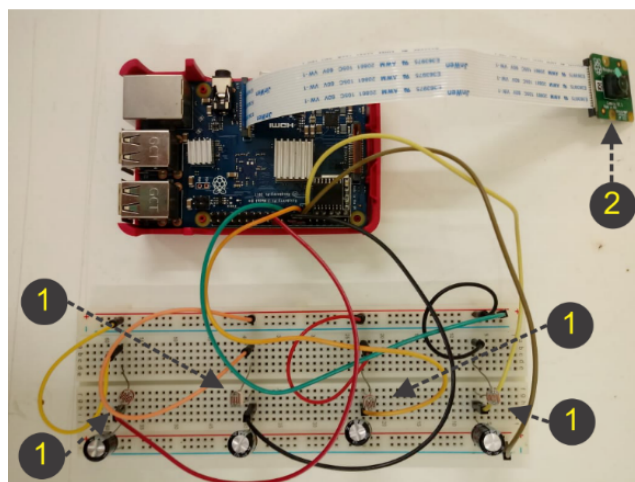


Figura 1. Sistema de aquisição montado.

- Uso de um dispositivo de baixo custo para a obtenção de previsões de irradiância solar em diversos horizontes de tempo.
- Desenvolvimento de um novo conjunto de dados para a comunidade de pesquisa em energia solar.

Vale destacar que o dispositivo de previsão desenvolvido nesta pesquisa utiliza sensores e uma plataforma *hardware* de baixo custo que fornece previsões de irradiância solar com acurácia competitiva quando comparado com as metodologias que fornecem os melhores desempenhos na literatura internacional, que utilizam para a construção de seus bancos de dados piranômetros, que são sensores de alto custo.

## 2. METODOLOGIA

### 2.1 Dados

O banco de dados é formado por variáveis de entrada relacionadas à sinais obtidos por sensores de luminosidade LDR's e atributos estatísticos (média aritmética, desvio padrão e entropia de Shannon), Equações (1), (2) e (3) calculados em cada canal de imagens do céu no formato RGB. A integração dos sensores foi feita pelo uso da plataforma *Raspberry Pi 3*. O Sistema de aquisição montado é indicado na Figura 1, onde o número 1 indica os sensores LDR's e o 2 a câmera.

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i \quad (1)$$

$$\sigma = \frac{1}{N-1} \sum_{i=1}^N (x_i - \mu)^2 \quad (2)$$

$$H(x) = - \sum_{i=0}^{255} p(x_i) \log(p(x_i)) \quad (3)$$

Onde,  $x_i$  é a intensidade de nível de cinza para o  $i$ -ésimo pixel e  $p(x_i)$  é a probabilidade da ocorrência da intensidade de nível de cinza do  $i$ -ésimo nível de cinza. Para cada imagem foram aplicados os filtros de suavização



Figura 2. Imagem utilizada para criação do conjunto de dados.

da mediana e o de aguçamento do laplaciano da gaussiana e novamente os atributos estatísticos foram calculados e adicionados ao conjunto de dados. A Figura 2 fornece um exemplo de imagem do céu utilizada para a obtenção deste banco dados.

Os dados foram coletados no Laboratório de Energia Solar e Gás Natural (LESGN) do Centro de Tecnologia da Universidade Federal do Ceará, na cidade de Fortaleza-CE, Brasil  $3^{\circ}43'6''S$  e  $38^{\circ}32'36''O$  entre os dias 03/06/2019 à 07/06/2019, caracterizados por serem dias de céu claro com baixa nebulosidade, englobando o intervalo de 8:00 às 17:30. Resultando em um banco de dados que consiste em uma planilha com 1466 observações e 21 variáveis.

## 2.2 Treinamento

Para o treinamento dos modelos, o conjunto de dados foi dividido aleatoriamente de tal forma que 70% ficou para o treinamento dos modelos e o restante foi utilizado para o teste. O ajuste dos métodos foi feito por meio de validação cruzada 10-Fold no conjunto de treino (JAMES et al. 2013).

## 2.3 Pré-processamento

O procedimento realizado consistiu em aplicar uma transformação logarítmica afim de tornar as distribuições das variáveis de entrada mais simétricas, já que para muitos modelos de aprendizagem a hipótese de distribuição normal é relevante.

Para fins de simplificação, todos os histogramas desenvolvidos serão limitados aos preditores de maior correlação positiva ou negativa com a irradiância que no caso foram a entropia do canal azul das imagens ( $ent_b$ ), entropia do canal vermelho ( $ent_r$ ) e a entropia da imagem filtrada com um filtro da mediana ( $ent_{mediana}$ ), com correlações de 0,706; 0,684; 0,693; respectivamente.

A matriz de correlação na forma de mapa de calor é representada na Figura 3, onde as variáveis ( $RAD_{GLOBAL}$ ) e ( $H_{30}$ ) indicam a irradiância global no instante atual e à 30 minutos a posteriori, respectivamente.

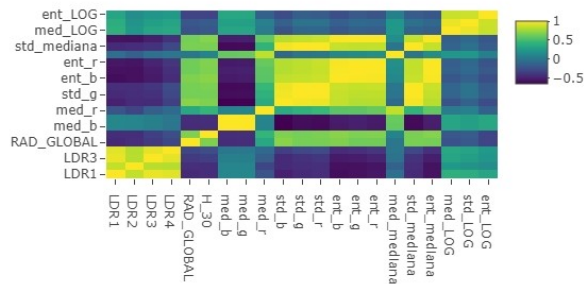


Figura 3. Matriz de correlação.

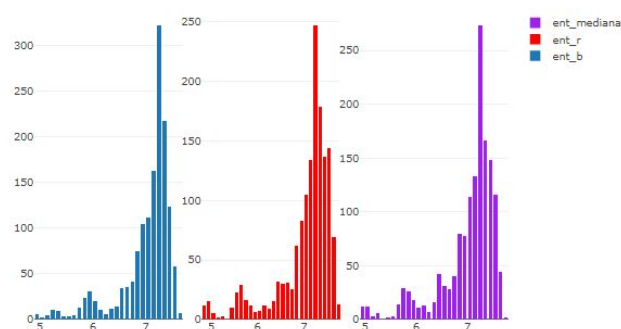


Figura 4. Histogramas das variáveis originais.

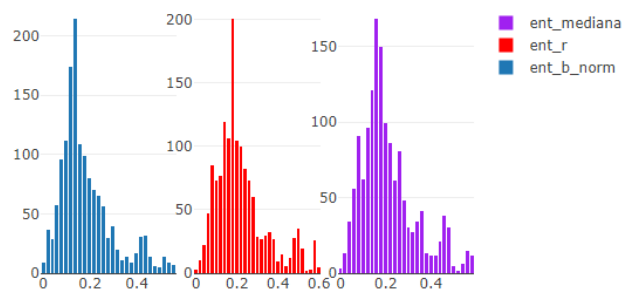


Figura 5. Histogramas das variáveis transformadas.

De fato, as distribuições dos preditores relacionados às entropias são bem assimétricas e distantes da distribuição normal como observado na Figura 4, os *skewness* para os preditores  $ent_b$ ,  $ent_r$  e  $ent_{mediana}$ , são -1,897; -1,791 e -1,720, respectivamente. Indicando, uma assimetria mais voltada para a direita. Visando a redução dessa assimetria, foi aplicada uma transformação logarítmica, como dada pela Equação (4).

$$\tilde{x} = \log((\max x + 1) - x) \quad (4)$$

Onde,  $\tilde{x}$  é o preditor transformado. Os histogramas para as variáveis de entrada transformadas são reportados na Figura 5. Percebe-se como a transformação logarítmica melhora consideravelmente a simetria dos histogramas.

## 2.4 Regressão Ridge

A regressão Ridge é similar à regressão linear, ela também parte da hipótese de que a saída e as variáveis de entrada estão relacionadas por meio de uma função linear e aditiva.

A diferença é observada na estimação do vetor de parâmetros, pois neste caso a minimização é realizada sobre a soma dos quadrados dos resíduos (RSS) adicionado com um termo de regularização, como se nota pela Equação (5).

$$\min_{\beta} \sum_{i=1}^N Y_i - (\beta_0 + \dots + \beta_p x_p) + \lambda \sum_{j=1}^p \beta_j^2 \quad (5)$$

O hiperparâmetro  $\lambda$  foi estimado utilizando validação cruzada 10-Fold, o mesmo atua como um parâmetro de penalidade que ajusta o nível de flexibilidade/variância na estimativa do vetor  $\beta$  quanto maior o valor de  $\lambda$  menor será a norma do vetor  $\beta$  estimado, vale destacar que por mais que esta abordagem limite os valores dos parâmetros livres estimados, a mesma não anula nenhum  $\beta_i$  ou seja, não há seleção de preditores.

### 2.5 LASSO

O LASSO segue a mesma metodologia da regressão Ridge, diferindo apenas no termo de regularização da etapa de estimação do vetor de parâmetros, aqui utiliza-se a norma  $L_1$  ao invés da norma  $L_2$  como acontece no Ridge, isto é ilustrado na Equação (6).

$$\min_{\beta} \sum_{i=1}^N Y_i - (\beta_0 + \dots + \beta_p x_p) + \lambda \sum_{j=1}^p |\beta_j| \quad (6)$$

De forma equivalente a regressão Ridge, o hiperparâmetro  $\lambda$  foi estimado por meio de uma validação cruzada 10-Fold, este atua como um fator de penalidade que ajusta o nível de flexibilidade/variância na estimativa do vetor  $\beta$  quanto maior o valor de  $\lambda$  menor será a norma do vetor  $\beta$  estimado. Diferentemente do que ocorre no Ridge, no LASSO tem-se a restrição sobre os valores dos parâmetros livres com a possibilidade de que pelo menos um deles seja zero, caracterizando assim uma seleção de preditores.

### 2.6 Percéptron de Múltiplas Camadas

Os percéptrons de múltiplas camadas (MLP) são o tipo mais comum de redes neurais artificiais (RNA) e são considerados modelos de *benchmarking* devido ao seu bom desempenho em muitas aplicações.

Tal modelo tem inspiração nos neurônios biológicos. Tem-se na terminologia de RNA que o vetor  $\mathbf{x} = [x_1, x_2, x_3]$  é o vetor de entradas,  $w = [w_1, w_2, w_3]$  é o vetor de pesos sinápticos,  $f(\cdot)$  é a função de ativação e  $y$  é muitas vezes chamado de campo local induzido (Haykin, 2001), que é definido pela Equação (7).

$$y = \beta + \sum_{i=1}^3 w_i x_i \quad (7)$$

Onde  $\beta$  representa um viés externo. Então, a saída  $z$  é determinada quando se calcula o valor da função de ativação no campo local induzido, assim  $z = f(y)$ . existem diversos tipos de funções de ativação que podem selecionadas (Haykin, 2001), mas a função sigmóide ou logística

é geralmente a mais aplicada devido às suas atrativas propriedades matemáticas, tais como, monotonicidade, continuidade e diferenciabilidade. Assim, a saída  $z$  é calculada pela Equação (8), tal função foi empregada em todas as camadas da MLP.

$$z = \frac{1}{1 + \exp(\beta + \sum_{i=1}^3 w_i x_i)} \quad (8)$$

Com o objetivo de reduzir o custo computacional associado ao uso da MLP, uma vez que o problema envolve 20 variáveis de entrada e uma saída, empregou-se uma arquitetura com uma camada oculta e as camadas de entrada e saída para MLP, onde o ajuste é realizado sobre o número de neurônios da camada oculta. Para a implementação da rede foi utilizado o pacote *neuralnet* do software R.

### 2.7 Modelo de Envelope para preditores

O Modelo de Envelope para preditores permite ganhos de eficiência na estimação dos coeficientes de regressão por meio da redução da variância. Para descrever o Modelo de Envelope a regressão linear múltipla é reformulada para ser consistente com a convenção usada em (Cook et al., 2013) como dado pela Equação (9).

$$\mathbf{Y} = \mu + \beta^T (\mathbf{X} - \mu_x) + \epsilon \quad (9)$$

Onde  $\mathbf{Y} \in \mathbf{R}^r$  é o vetor das saídas e  $\mathbf{X} \in \mathbf{R}^p$  é o vetor de preditores com média  $\mu_x$  e covariância  $\Sigma_x$ . O vetor de erros  $\epsilon$  tem média 0 e matriz de covariância  $\Sigma_{X|Y}$ . Os coeficientes de regressão estão contidos na matriz  $\beta \in \mathbf{R}^{r \times p}$ .

Sendo  $\mathcal{S}$  um subespaço de  $\mathbf{R}^p$ . É realizado uma decomposição em  $\mathbf{X}$  obtendo assim uma parte material  $\mathbf{P}_{\mathcal{S}}\mathbf{X}$  e uma parte imaterial  $\mathbf{Q}_{\mathcal{S}}\mathbf{X}$ . tal que as duas condições são satisfeitas.

- (1)  $cov(\mathbf{Y}, \mathbf{Q}_{\mathcal{S}}\mathbf{X} | \mathbf{P}_{\mathcal{S}}\mathbf{X}) = 0$
- (2)  $cov(\mathbf{P}_{\mathcal{S}}\mathbf{X}, \mathbf{Q}_{\mathcal{S}}\mathbf{X}) = 0$

Onde o operador  $\mathbf{P}_{\mathcal{S}}$  representa a projeção sobre o subespaço  $\mathcal{S}$  e  $\mathbf{Q}_{\mathcal{S}} = \mathbf{I} - \mathbf{P}_{\mathcal{S}}$  seu complementar.

Estas duas condições indicam que  $\mathbf{Q}_{\mathcal{S}}\mathbf{X}$  não afeta a distribuição de  $\mathbf{Y}$  ou tendo associação com  $\mathbf{P}_{\mathcal{S}}\mathbf{X}$ . A interseção de todos os  $\mathcal{S}$  que satisfazem as condições (1) e (2) é o  $\Sigma_x$ -envelope de  $\beta$ , denotado por  $\mathcal{E}_{\Sigma_x}(\beta)$  ou  $\mathcal{E}$ . Sua dimensão  $u$  é o único hiperparâmetro do modelo.

Seja  $\Gamma \in \mathbf{R}^{p \times u}$  uma base ortonormal de  $\mathcal{E}$  e  $\Gamma_0 \in \mathcal{R}^{p \times (p-u)}$  seu complemento ortogonal, sob as condições (1) e (2) o problema (9) pode ser formulado como dado pelas Equações (10), (11).

$$\mathbf{Y} = \mu + \eta^T \Omega^{-1} \Gamma^T (\mathbf{X} - \mu_x) + \epsilon \quad (10)$$

$$\Sigma_x = \Gamma \Omega \Gamma^T + \Gamma_0 \Omega_0 \Gamma_0^T \quad (11)$$

Onde  $\beta = \Gamma \Omega^{-1} \eta \in \mathbf{R}^{u \times r}$ ,  $\Omega^{-1} \eta$  representa as coordenadas de  $\beta$  com relação a  $\Gamma$ , as matrizes  $\Omega \in \mathbf{R}^{u \times u}$  e  $\Omega_0 \in \mathbf{R}^{(p-u) \times (p-u)}$  são definidas positivas. O modelo dado pelas Equações (10) e (11) é o modelo de envelope para

os preditores, onde para  $u = p$  o modelo se resume na regressão linear padrão.

Uma vez que se conhece o envelope, fica fácil de obter o estimador, basta determinar os coeficientes utilizando a metodologia padrão e projetá-los no subespaço  $\mathcal{E}$  por meio da base ortonormal  $\Gamma$ . Para estimar o envelope  $\mathcal{E}$  é necessário resolver um problema de otimização restrita sobre uma variedade grassmaniana o que pode ser consideravelmente lento se o problema é de tamanho razoável.

Agora considerando a base ortonormal  $\Omega \in \mathbf{R}^{a \times b}$  de  $\mathcal{E}$ , pode-se sem perda de generalidade considerar que a matriz  $\Gamma_1$  formada pelas  $b$  linhas de  $\Gamma$  é não singular. Seja  $\Gamma_2$  a matriz  $(a-b) \times b$  formada pelas linhas restantes de  $\Gamma$ . Então  $\Gamma$  pode ser particionada como representada pela Equação (12)

$$\Gamma = \begin{pmatrix} \Gamma_1 \\ \Gamma_2 \end{pmatrix} = \begin{pmatrix} I_b \\ \Gamma_2 \Gamma_1^{-1} \end{pmatrix} \Gamma_1 = \begin{pmatrix} I_b \\ A \end{pmatrix} \Gamma_1 = G_A \Gamma_1 \quad (12)$$

Esta reparametrização constrói uma correspondência um a um entre  $A$  e  $\mathcal{E}$  e daí com essa manipulação é possível estimar  $A$  através da minimização da função objetivo indicada na Equação (13) sobre a matriz  $A \in \mathbf{R}^{(a-b) \times b}$ .

$$\log |G_A^T M G_A| + \log |G_A^T (M + U) G_A| - 2 \log |G_A^T G_A| \quad (13)$$

Assim a otimização é irrestrita e realizada sobre uma matriz. Para problemas onde  $(r-u)u$  é grande o algoritmo de coordenadas descendente por blocos pode ser uma boa opção como detalhado em (COOK et al., 2016).

### 2.8 Métricas de Erro

As métricas de erro utilizadas para realizar a avaliação do desempenho dos modelos de aprendizagem de máquina utilizados são apresentadas nesta seção.

**Erro médio absoluto:** O MAE calcula a média das diferenças absolutas entre o valor previsto,  $\hat{y}_i$  e o valor real,  $y_i$  isso é, não leva em consideração se o erro é para mais ou para menos e às diferenças absolutas não é atribuído peso, como indica a Equação (14).

$$MAE = \frac{1}{N} \sum_{i=1}^N |\hat{y}_i - y_i| \quad (14)$$

**Erro Médio po Viés:** MBE se assemelha ao cálculo do MAE, mas se diferencia deste por considerar o sinal do erro, isto é, não calcula o valor absoluto das diferenças. Desse modo, requer prudência em sua análise, uma vez que permite a compensação de erros (erros com sinais distintos). A mesma é dada pela Equação (15).

$$MBE = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i) \quad (15)$$

**Raiz do erro quadrático médio:** O RMSE calcula a magnitude da média do erro pela raiz quadrada da média dos quadrados dos erros. Desse modo atribui um peso maior

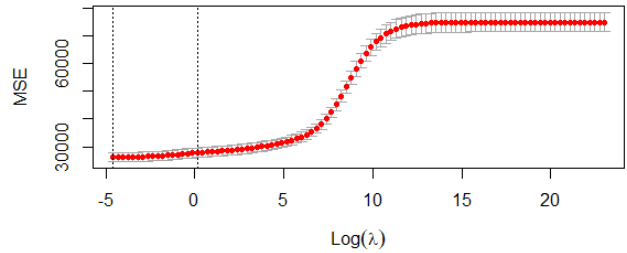


Figura 6. MSE em termos de  $\log(\lambda)$  no Ridge.

aos erros de maior magnitude, e peso menor aos erros de menor magnitude. É obtido na mesma unidade da variável em análise, sendo definida pela Equação (16)

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2} \quad (16)$$

**Raiz do erro quadrático médio relativo:** O rRMSE é a razão entre o RMSE e a média dos valores da variável em análise, como dado pela Equação (17).

$$rRMSE = \frac{RMSE}{\sum_{i=1}^N y_i} \quad (17)$$

Tal métrica foi utilizada porque fornece faixas de classificação do desempenho das previsões:  $rRMSE < 10\%$  excelente,  $10\% < rRMSE < 20\%$  bom,  $20\% < rRMSE < 30\%$  razoável,  $rRMSE > 30\%$  ruim (LI et al., 2013).

## 3. RESULTADOS

### 3.1 Treinamento

Nesta seção são reportados alguns resultados relacionados a etapa de ajuste dos métodos. Ao realizar o treinamento para a regressão Ridge, foi selecionado um *grid* com 100 valores de  $\lambda$  iniciando em  $10^{10}$  e finalizando em  $10^{-2}$ , o valor ótimo encontrado foi de  $\lambda^* = 0,01$ , os valores do Erro Quadrático Médio (MSE) em termos de  $\log \lambda$  são reportados nas Figura 6. O intuito de utilizar o logaritmo se deve ao fato de estar considerando valores muito elevados de  $\lambda$  assim para facilitar a visualização gráfica aplicou-se tal transformação.

De fato, percebe-se que o valor de  $\lambda^* = 0,01$  corresponde ao menor valor do *grid*, indicando que para o banco de dados em estudo, a penalidade sobre a estimativa do vetor de parâmetros  $\beta$  é a menor possível, ou seja, uma maior flexibilidade/variância nos valores estimados dos parâmetros.

Para o LASSO, o valor de  $\lambda^* = 0,01$  foi obtido utilizando o mesmo *grid* e validação cruzada 10-Fold no conjunto de treinamento, os valores do Erro Quadrático Médio (MSE) em termos de  $\log \lambda$  são reportados nas Figura 7.

Os números na parte superior da Figura 7, indicam a quantidade de preditores considerados na regressão, percebe-se que a partir de  $\log(5)$  apenas o intercepto é considerado, o que é um indicativo de baixa influência das variáveis de entrada sobre a variável de saída o que é condizente

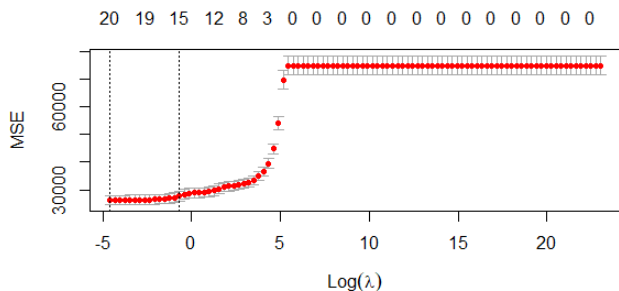


Figura 7. MSE em termos de  $\log(\lambda)$  no LASSO.

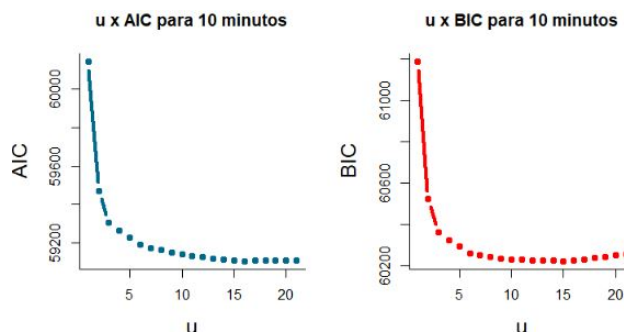


Figura 8. AIC e BIC em termos de  $u$  para todos os horizontes de previsão.

com a matriz de correlação obtida. Novamente, nota-se que uma maior flexibilidade/variação nos valores estimados dos parâmetros é mais recomendada para o conjunto de dados em estudo.

Para a obtenção da dimensão ótima do envelope  $u^*$  foi utilizado o Critério de Informação de Akaike (AIC) e o Critério de Informação Bayesiano (BIC), que são coeficientes baseados em teoria da informação e geralmente utilizados para seleção ótima de parâmetros em modelos paramétricos. variando a dimensão em um *grid* de 1 à 21 como indicado na Figura 8 para o horizonte de 10 minutos, os demais horizontes apresentam o mesmo comportamento. Neste caso, os menores valores para o AIC e BIC considerando o horizonte de 10 minutos ocorreram nas dimensões  $u = 14$  e  $u = 15$ , para 20 minutos  $u = 11$  e  $u = 17$  e para 30 minutos  $u = 14$  e  $u = 17$  respectivamente. Assim foi escolhido sempre a menor dimensão para fins de redução de custo computacional.

para a rede MLP foi escolhida uma arquitetura padrão com três camadas: uma entrada, uma oculta e uma de saída. O número de neurônios na camada oculta foi determinado por um procedimento de *Hold-Out* no conjunto de treino o *grid* selecionado foi de 1 à 10 neurônios na camada oculta, onde uma maior faixa não foi utilizada para se ter uma redução no custo computacional associado. Os resultados são ilustrados na Figura 9. Onde nota-se que três neurônios na camada oculta fornece os menores valores de RMSE e MAE, assim a arquitetura final consiste de uma camada oculta com três neurônios.

### 3.2 Teste

Os modelos ajustados foram aplicados no conjunto de teste e seus desempenhos foram dados em termos dos valores das

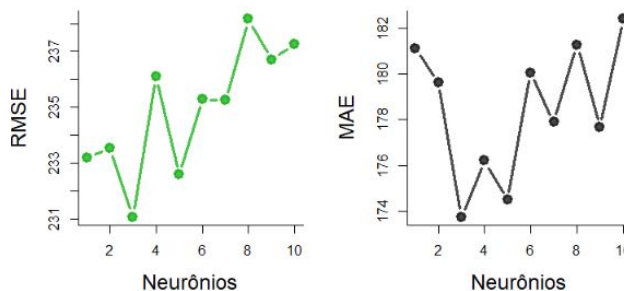


Figura 9. RMSE em termos de neurônios na camada oculta.

métricas de erro, os resultados são reportados na Tabela 1, 2 e 3, para os horizontes de 10, 20 e 30 minutos a posteriori, respectivamente.

Tabela 1. Resultados para o horizonte de 10 minutos.

Modelos	MAE	MBE	RMSE	rRMSE
Ridge	113,65	-8,27	164,29	27,07
LASSO	113,94	-8,50	164,63	27,13
<b>Envelope</b>	<b>112,35</b>	<b>-8,17</b>	<b>163,51</b>	<b>26,94</b>
MLP	177,11	-91,02	234,19	38,59

Tabela 2. Resultados para o horizonte de 20 minutos.

Modelos	MAE	MBE	RMSE	rRMSE
Ridge	112,41	5,34	155,02	25,48
<b>LASSO</b>	<b>112,08</b>	<b>5,14</b>	<b>154,67</b>	<b>25,42</b>
Envelope	113,87	4,72	156,56	25,73
MLP	170,43	-86,92	227,44	37,38

Tabela 3. Resultados para o horizonte de 30 minutos.

Modelos	MAE	MBE	RMSE	rRMSE
<b>Ridge</b>	<b>120,43</b>	<b>5,21</b>	<b>164,70</b>	<b>27,60</b>
LASSO	120,99	6,00	165,44	27,73
Envelope	122,23	6,58	167,39	28,06
MLP	172,85	-81,30	226,58	37,24

Como observado pelas Tabelas 1, 2 e 3 para o conjunto de dados em questão os modelos lineares forneceram melhores acurácias do que a rede MLP, de fato, a rede neural sofre com um problema de *overfitting*, obtendo convergência no treinamento, mas com problemas no teste. O modelo de envelopes fornece melhores resultados para o horizonte de previsão de 10 minutos posteriori.

Os modelos de regressão linear regularizados (Ridge e LASSO) apresentaram desempenho superior para os horizontes de 20 e 30 minutos a posteriori, respectivamente. Os valores das métricas RMSE e MAE são bem similares para as metodologias lineares, indicando que na média as mesmas obtiveram a mesma acurácia para o problema considerado.

Muitos trabalhos (Gutierrez et al., 2016; Rocha e Modolo, 2019; Madhidasan et al., 2021) demonstraram o alto poder preditivo das redes MLP na realização de previsões de irradiância solar. É interessante notar o baixo desempenho da rede MLP, o que pode ser um indicativo de que a arquitetura escolhida não é a mais adequada e que um

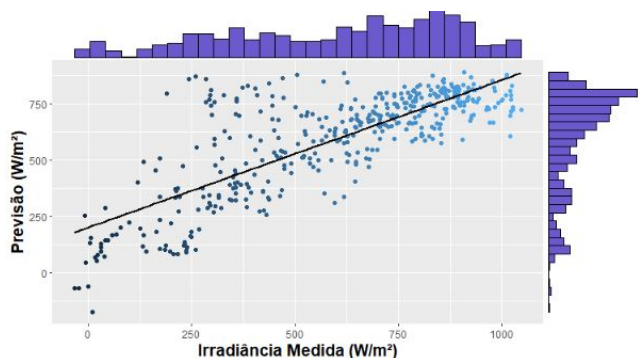


Figura 10. Previsão por valores medidos de irradiância para o horizonte de 30.

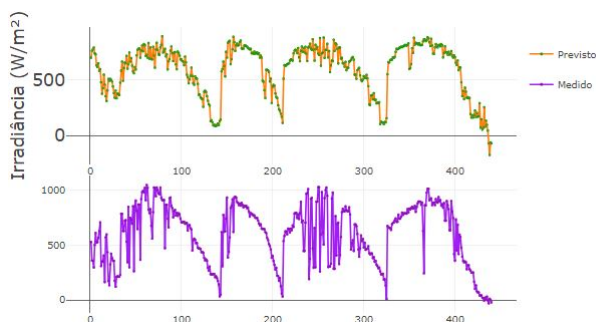


Figura 11. Previsão e valores medidos de irradiância para o horizonte de 30.

melhor procedimento de *tunning* deve ser realizado para se obter uma rede melhor ajustada.

Vale destacar o bom desempenho do modelo de envelopes, mesmo com a remoção da parte imaterial na regressão linear múltipla. Em todos os horizontes a acurácia foi similar aos demais métodos lineares com o benefício de redução da variância nas estimativas dos coeficientes sem adição de amostras e, portanto, redução do custo computacional. As Figuras 10 e 11 fornecem uma representação gráfica para o desempenho do modelo de envelopes na previsão de irradiância no horizonte de 30 minutos.

#### AGRADECIMENTOS

Coloque aqui seus agradecimentos.

#### 4. CONCLUSÃO

Neste trabalho, empregamos o recente método de estimação por envelopes para realizar previsões de irradiância solar nos horizontes de 10, 20 e 30 minutos a posteriori utilizando um conjunto de dados formados por sinais de luminosidade coletados por LDR's e atributos estatísticos extraídos de imagens do céu capturadas por uma câmera direcionada ao zênite, onde a integração dos sensores foi realizada por um *Raspberry Pi 3*. A acurácia desta abordagem foi comparada com a dos modelos lineares Ridge e LASSO e com o modelo não linear MLP. Para 10 minutos a frente o modelo de envelope obteve melhor desempenho, o que pode ser um indicativo de maior nível de variabilidade na saída quando se considera tal horizonte. Os métodos LASSO e Ridge obtiveram melhor acurácia para os horizontes de 20 e 30 minutos, respectivamente. Destaca-se o

desempenho ruim da rede MLP em todos os horizontes o que pode ser um indicativo de que a arquitetura escolhida não é a mais adequada e que um melhor procedimento de *tunning* deve ser realizado para se obter uma rede melhor ajustada.

#### 5. REFERÊNCIAS

- Benali, L., Notton, G., Fouilloy, A., Voyant, C., and Dizene, R. (2019). Solar radiation forecasting using artificial neural network and random forest methods: Application to normal beam, horizontal diffuse and global components. *Renewable energy*, 132, 871–884.
- Cook, R.D., Forzani, L., and Su, Z. (2016). A note on fast envelope estimation. *Journal of Multivariate Analysis*, 150, 42–54.
- Cook, R.D., Helland, I., and Su, Z. (2013). Envelopes and partial least squares regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 75(5), 851–877.
- Cook, R.D. and Zhang, X. (2016). Algorithms for envelope estimation. *Journal of Computational and Graphical Statistics*, 25(1), 284–300.
- Dong, Z., Yang, D., Reindl, T., and Walsh, W.M. (2013). Short-term solar irradiance forecasting using exponential smoothing state space model. *Energy*, 55, 1104–1113.
- Gutierrez-Corea, F.V., Manso-Callejo, M.A., Moreno-Rigidor, M.P., and Manrique-Sancho, M.T. (2016). Forecasting short-term solar irradiance based on artificial neural networks and data from neighboring meteorological stations. *Solar Energy*, 134, 119–131.
- Haykin, S. (2001). *Redes neurais: princípios e prática*. Bookman Editora.
- Koo, C., Li, W., Cha, S.H., and Zhang, S. (2019). A novel estimation approach for the solar radiation potential with its complex spatial pattern via machine-learning techniques. *Renewable energy*, 133, 575–592.
- Madhiarasan, M., Louzazni, M., and Roy, P.P. (2021). Novel cooperative multi-input multilayer perceptron neural network performance analysis with application of solar irradiance forecasting. *International Journal of Photoenergy*, 2021.
- Mejia, J.F., Giordano, M., and Wilcox, E. (2018). Conditional summertime day-ahead solar irradiance forecast. *Solar Energy*, 163, 610–622.
- Nonnenmacher, L., Kaur, A., and Coimbra, C.F. (2016). Day-ahead resource forecasting for concentrated solar power integration. *Renewable energy*, 86, 866–876.
- Pawar, P., Cortés, C., Murray, K., and Kleissl, J. (2019). Detecting clear sky images. *Solar Energy*, 183, 50–56.
- Pedro, H.T. and Coimbra, C.F. (2015). Nearest-neighbor methodology for prediction of intra-hour global horizontal and direct normal irradiances. *Renewable Energy*, 80, 770–782.
- Pedro, H.T., Coimbra, C.F., David, M., and Lauret, P. (2018). Assessment of machine learning techniques for deterministic and probabilistic intra-hour solar forecasts. *Renewable Energy*, 123, 191–203.

Rocha, P., Fernandes, J., Modolo, A., Lima, R., da Silva, M., and Bezerra, C. (2019). Estimation of daily, weekly and monthly global solar radiation using anns and a long data set: a case study of fortaleza, in brazilian northeast region. *International Journal of Energy and Environmental Engineering*, 10(3), 319–334.

Trapero, J.R., Kourentzes, N., and Martin, A. (2015). Short-term solar irradiation forecasting based on dynamic harmonic regression. *Energy*, 84, 289–295.

Yagli, G.M., Yang, D., and Srinivasan, D. (2019). Automatic hourly solar forecasting using machine learning models. *Renewable and Sustainable Energy Reviews*, 105, 487–498.