

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

MODELOS DE MACHINE LEARNING PARA ESTIMAR LA RADIACIÓN SOLAR HORIZONTAL EN LA PAMPA HÚMEDA CON INFORMACIÓN SATELITAL MULTIESCALA

Paula Iturbide¹, Ximena Orsi¹, María José Denegri^{1,2}, Santiago Fioretti¹, Pablo Ruiz¹, Sergio Luza¹, Valeria Stern¹, Rodrigo Alonso-Suárez³, Franco Ronchetti^{4,5}

¹ Grupo de Estudios de la Radiación Solar (GERSolar), Instituto de Ecología y Desarrollo Sustentable (INEDES). Univ. Nacional de Luján, CP 6700, Buenos Aires, Argentina.

² Departamento de Tecnología, Universidad Nacional de Luján, Buenos Aires, Argentina.

³ Laboratorio de Energía Solar, Dpto. de Física del CENUR Litoral Norte, Udelar, Uruguay.

⁴ Instituto de Investigación en Informática LIDI, Universidad Nacional de La Plata, Buenos Aires, Argentina.

⁵ Comisión de Investigaciones Científicas de la Pcia. de Buenos Aires (CICPBA), Buenos Aires, Argentina.

e-mail: paula.itur@gmail.com

RESUMEN: La falta de precisión en los datos de radiación solar puede tener un gran impacto en la rentabilidad de los proyectos de energía solar. Las redes de medición terrestre ofrecen información limitada por su distribución esparza en el territorio. Esto lleva a desarrollar modelos de estimación por imágenes satelitales, los cuales resuelven la espacialidad si son ajustados a mediciones terrestres de calidad. En este estudio, se desarrollan y validan modelos empíricos de aprendizaje automático para la estimación por satélite de radiación solar global horizontal, demostrando su utilidad y precisión en la región analizada. Estos modelos se alimentan con variables provenientes de imágenes satelitales GOES-16 y variables geométricas. Los resultados sugieren que, para ciertas combinaciones de variables satelitales de entrada, la información geométrica puede ser utilizada en forma implícita para realizar estimaciones precisas de la radiación solar. Debido al volumen de la información satelital disponible, desarrollamos un análisis de componentes principales para reducir la dimensionalidad. Para comparar el modelo propuesto adaptamos localmente las estimaciones del Heliosat-4 y del CIM-ESRA al sitio, y también implementamos el modelo CIM-McClearn. Los resultados muestran una superioridad de desempeño del modelo de aprendizaje automático propuesto, demostrando que es capaz de extraer información de la multiescala espacial satelital. Por otro lado, la mejora de desempeño obtenida es leve, lo que muestra la dificultad en seguir mejorando el desempeño de la estimación satelital de radiación solar.

Palabras clave: Radiación solar, aprendizaje automático, Imágenes satelitales. GOES16, GHI.

INTRODUCCIÓN

En el ámbito de la estimación de radiación solar por satélite coexisten tres enfoques: el físico, el estadístico y el híbrido. Los modelos físicos resuelven las ecuaciones de transferencia radiante en la atmósfera, utilizando información sobre los componentes atmosféricos que interactúan con la radiación solar (Perez R. et al., 2013). Los modelos estadísticos se basan en una serie de coeficientes ajustados empíricamente a datos terrestres tomando como entrada la información satelital. Por último, en los modelos híbridos o semi-empíricos la formulación del modelo tiene una base física, pero dependen de una serie de parámetros ajustables.

En la región de la Pampa Húmeda se han desarrollado modelos satelitales híbridos específicamente ajustados, como el CIM-ESRA y CIM-McClearn (Laguarda et al., 2020). Otro modelo relevante en la región es el modelo físico Heliosat-4 (Qu et al., 2017), que ha sido extensamente evaluado (Gonzalez et

50 al., 2019; Laguarda et al., 2020, 2021; Sarazola et al., 2023). Los modelos CIMs, al utilizar imágenes
51 GOES16 y estar ajustados localmente a la región, han sido evaluados con significativo mejor desempeño
52 en la Pampa Húmeda que el Heliosat-4, que utiliza imágenes del satélite europeo Meteosat, obteniendo
53 desvíos cuadráticos medios entre 16-17% (relativo a la media de las medidas) para estimaciones de
54 irradiancia solar global horizontal (GHI) a escala 10-minutal. Reducir la incertidumbre por debajo de
55 este límite ha demostrado ser un desafío.

56 Una posible forma de reducir esta incertidumbre es a través del uso de información satelital multi-escala
57 espacial y algoritmos de aprendizaje automático como redes neuronales artificiales (RN), k vecinos más
58 cercanos (kNN), vectores soporte de regresión (SVR), máquinas de aprendizaje extremas (ELM) y
59 ensamblajes de árboles como random forest (RF) y gradient boosting (GB). En la actualidad es común en
60 varias áreas el uso de estos algoritmos. En particular, para la estimación de la radiación solar se han
61 utilizado principalmente con base en mediciones terrestres de otras variables meteorológicas como
62 presión, temperatura de aire ambiente, heliofanía, humedad, precipitación, nubosidad vista desde tierra,
63 velocidad del viento y/o evaporación, etc., y variables auxiliares como el día del año, latitud, longitud,
64 altitud, modelos de cielo claro, y/o variables geométricas, entre otras (Raichijk, 2008; Sayago et al.,
65 2011; Jiménez et al., 2017; Olivera et al., 2020). Las investigaciones que emplean estas técnicas con
66 información satelital de nubosidad son escasas (Verbois et al., 2023).

67 Este trabajo es la continuación del artículo de Iturbide et al. (2023), donde se implementaron los
68 algoritmos RN, RF y regresión lineal simple para la estimación de la GHI por satélite. Se utilizaron 38
69 variables de entrada de las cuales 19 estaban relacionadas con el factor de reflectancia y 19 con la
70 reflectancia planetaria. Estas variables abarcaban diversas resoluciones espaciales, variando entre 0,01
71 y 0,9 grados de latitud y longitud. Además, se incluyeron el coseno del ángulo cenital y el modelo de
72 cielo claro McClear como parte de las variables de entrada. Los resultados de dicho artículo mostraron
73 un rendimiento superior por parte de la RN, seguida por el enfoque de RF, incluso después de la
74 eliminación de variables como el modelo de cielo claro y el coseno del ángulo cenital. En contraste, la
75 regresión lineal simple demostró un desempeño insuficiente al excluir estas variables, ya que carecía de
76 la capacidad para reconstruir la referencia temporal esencial para la estimación de la GHI. Cabe
77 mencionar que la comparación con modelos preexistentes en Iturbide et al. (2023) se llevó a cabo sin
78 realizar la adaptación al sitio.

79 El objetivo de este artículo es mejorar el desempeño de los modelos de aprendizaje automático mediante
80 la incorporación de variables de entrada previamente no consideradas, tales como el índice de nubosidad
81 y un cálculo mejorado de la reflectancia planetaria. Además, se introduce un análisis de componentes
82 principales que contribuye a reducir la dimensionalidad del conjunto de entrada, mejorando la eficacia
83 de los modelos y condensando la información satelital en un conjunto reducido de variables. También
84 se agrega el modelo GB para comparar su desempeño. Se procede a comparar el modelo resultante con
85 los disponibles para la región, adaptados al sitio, y se incluye la implementación con ajuste local a
86 medidas del modelo CIM-McCclear. Para ello, se usa el mismo esquema de ajuste y testeo que el utilizado
87 para los modelos de aprendizaje automático, de modo de realizar una comparativa justa.

88 **METODOLOGÍA**

89 *Medidas en tierra y pre procesamiento de datos*

90 En este estudio se utilizaron datos de la estación Luján (de la red GERSolar) correspondientes al periodo
91 2019-2021, adquiridos con piranómetro de la firma Kipp & Zonen modelo CMP21 - equipo de Clase A
92 según la norma de clasificación de equipamiento para la medida de radiación solar (ISO 9060:2018)-, y
93 un adquisidor de datos Campbell Scientific modelo CR1000. Se adoptó la escala temporal de 10
94 minutos, que es la cadencia temporal de las imágenes capturadas por el satélite GOES-16. Las integrales
95 10 minutas (en W/m^2) fueron sometidas a un algoritmo de control de calidad que consta de los cuatro
96 filtros secuenciales mostrados en la Tabla 1, seguido de una revisión visual de las series para descartar
97 períodos afectados por sombras u otros fallos. En Iturbide et al. (2023) se detalla lo que impone cada
98 filtro. En la Tabla 2 se muestran los resultados del filtrado para la estación Luján, donde se indica además
99 su ubicación precisa.

100

Tabla 1: Filtros aplicados a las medidas en tierra

Filtro	Criterio	Descripción
1	$\alpha_s > 7^\circ$	Mínima altura solar
2	$-2W/m^2 < I_h < I_0 \cdot 1,2 \cdot \cos\theta_z^{1,2} + 50W/m^2$	Cotas de la BSRN (Long y Shi, 2008)
3	$0W/m^2 < I_h < I_h^{ESRA} (TL = 1,8)$	Cotas de un modelo de cielo claro
4	$k_{tp} < 0,89$	Cota del índice de claridad de Pérez

101

102

103

104

Tabla 2: Ubicación de la estación de medida y sus equipos de medición. Los N totales (luego de los filtros) corresponden a medidas integradas 10-minutales y el período corresponde a 2019-2021

Estación	Latitud (grados)	Longitud (grados)	Equipo	N Total
Luján, ARG	-34,558	-59,062	CMP21	62.592

105

106

107

Información satelital

108

109

110

111

112

113

114

115

116

117

118

Se utilizan las imágenes del canal visible (C02, centrado en $0,64 \mu m$) del satélite meteorológico geostacionario GOES-16. Este satélite forma parte de la red de satélites geostacionarios para la observación de la Tierra que cubre todo el globo terráqueo y es administrado por la National Oceanic and Atmospheric Administration (NOAA) de los Estados Unidos. Desde el año 2018 este satélite genera imágenes para todo el continente americano con una cadencia temporal regular de entre 10 y 15 minutos. Se encuentra ubicado sobre el ecuador terrestre en la longitud $-75^\circ W$. Su resolución espacial es variable a lo largo de la imagen, siendo de 500 m en su nadir. Sobre la región de la Pampa Húmeda el tamaño del píxel varía entre 1 y 3 km. El canal visible es el adecuado para la estimación de radiación solar debido a que la nubosidad diurna es claramente reconocible y cuantificable. Las nubes son típicamente más reflectivas que el fondo (la superficie terrestre), y por tanto, distinguibles.

119

120

121

122

123

124

125

126

127

Las dos variables típicas que se calculan a partir del canal visible de una imagen satelital son el factor de reflectancia (FR) y la reflectancia planetaria (RP). Esta última cantidad es también conocida como Albedo terrestre. El FR es una normalización de la radiancia medida por el satélite proveniente de cada píxel respecto al máximo que es capaz de medir (es decir, la radiación solar que incide sobre el tope de la atmósfera normalizada por la respuesta espectral del radiómetro en órbita). Se encuentra por tanto en el intervalo $[0, 1]$ y contiene, además de información sobre nubosidad, información espacial sobre la iluminación variable del Sol sobre la Tierra. La cantidad RP contiene la normalización necesaria para eliminar esta dependencia espacial y es efectivamente la reflectividad de la Tierra, en su sentido físico estricto. Esta normalización se obtiene dividiendo a FR por el coseno del ángulo cenital solar.

128

129

La reflectancia planetaria es también normalizada para obtener el índice de nubosidad N de la siguiente manera:

130

$$N = (R - R_0) / (R_{max} - R_0) \quad (1)$$

131

132

133

134

donde R_0 es la reflectancia planetaria de fondo asociado a condiciones de cielo claro para cada celda y el parámetro R_{max} se asocia a condiciones de nubosidad total. Se utiliza aquí un modelo de fondo para el cálculo de R_0 específicamente ajustado al píxel objetivo (Alonso-Suárez et al., 2011) y un valor fijo de 0,8 para R_{max} , que fue optimizado para la estimación de GHI en la región (Laguarda et al., 2018).

135

136

Estas variables permiten caracterizar la nubosidad a partir de las imágenes satelitales y son las que se usan para estimar la radiación solar en toda condición de cielo.

137

139 *Definición de los conjuntos de entrenamiento y testeo.* El propósito principal es que los algoritmos
140 adquieran la capacidad de estimar la GHI mediante el ajuste a mediciones terrestres. Dado el carácter
141 estacional anual de la GHI, se empleó un conjunto de datos abarcando dos de los tres años disponibles
142 para el entrenamiento, mientras que el tercer año se reservó para el testeo. Se aplicaron todas las posibles
143 combinaciones de años para evitar tres posibles sesgos. En primer lugar, se evitó un sesgo relacionado
144 con la distribución aleatoria de los datos en los conjuntos de entrenamiento y testeo, ya que los datos de
145 momentos consecutivos podrían presentar similitudes. En segundo lugar, se mitigó un sesgo vinculado
146 a las particularidades de cada año, pues un año podría diferir significativamente de otro en términos de
147 radiación solar. En tercer lugar, se abordó un sesgo relacionado con los datos faltantes: el año 2020
148 registró un hueco de dos meses, y para los años 2019 y 2021 se carece de información correspondiente
149 al mes de diciembre.

150 *Variables de entrada.* Las variables de entrada utilizadas en los modelos de aprendizaje automático son
151 el coseno del ángulo cenital ($\cos z$), el factor de reflectancia (FR), la reflectancia planetaria y el índice
152 de nubosidad. Se trabaja con dos índices de nubosidad N1 y N2 y dos factores de reflectancia R y RC.
153 Las variables $\cos z$, FR y R son las mismas que se utilizaron en el artículo de Iturbide et al. (2023). La
154 distinción entre las reflectancias planetarias (R y RC), radica en la manera en que son normalizadas.
155 Para la normalización de RC se utiliza la expresión de la masa de aire de Young (1994), lo que produce
156 mejores resultados al inicio y final del día que la normalización simple por coseno del ángulo cenital,
157 utilizada para R. La diferencia entre N1 y N2 está en el modelo de brillo de fondo utilizado. El valor de
158 N1 fue calculado con el modelo de fondo original de Alonso-Suárez et al. (2011) utilizado desde la
159 generación vieja de satélites GOES. En cambio, el N2 se calculó con un modelo de fondo actualizado y
160 ajustado al satélite GOES16. Algunas pruebas empíricas sugieren que este último presenta mejoras en
161 comparación con el primero.

162 Las variables satelitales FR, R, RC, N1 y N2 se consideran a distintas escalas espaciales con
163 nomenclatura del 01 al 20: FR01-FR20, R01-R20, RC01-RC20, N1_01-N1_20, N2_01-N2_20, siendo
164 promedios espaciales en celdas cuadradas de dimensiones crecientes. El espaciado entre tamaños no
165 sigue una relación lineal; en los tamaños más pequeños, el espaciado es más detallado, y viceversa
166 para los tamaños mayores.

167 *Modelos de aprendizaje supervisado.* En continuación al estudio previo, se amplió el conjunto de
168 modelos de aprendizaje automático supervisado. Además de la regresión lineal, RN (100 neuronas
169 ocultas y función de activación ReLu) y el RF (30 estimadores) que fueron abordados en el trabajo
170 anterior, se incorpora ahora el modelo de GB. Este nuevo modelo busca enriquecer el análisis al ofrecer
171 una perspectiva adicional en la estimación de la radiación solar.

172 **RESULTADOS**

173 *Ajuste local de los modelos CIM-ESRA y CAMS*

174 Se descargaron los estimativos de GHI de los modelos satelitales CIM-ESRA y Heliosat-4 para el
175 periodo 2019-2021 para la estación Luján de cada sitio web <http://les.edu.uy/online/stack-loc/> (CIM-
176 ESRA, portal LES) y <https://www.soda-pro.com/web-services> (Heliosat-4, portal CAMS). Los
177 estimativos de CIM-ESRA están disponibles en una escala 10-minutal para diferentes estaciones
178 latinoamericanas una de las cuales corresponde a Luján, Argentina. Los estimativos de CAMS no están
179 disponibles en la resolución temporal 10-minutal trabajada en este artículo, por lo que se descargaron
180 datos minutales y luego se integraron a la escala 10-minutal. Estos modelos se utilizaron para comparar
181 sus indicadores de desempeño respecto a las medidas en tierra con los de los modelos estadísticos
182 desarrollados en este trabajo.

183
184 La adaptación local de los modelos CIM-ESRA y CAMS fue ejecutada mediante cuatro enfoques
185 distintos (Salazar et al., 2021). Las tres primeras estrategias involucraron diversas combinaciones de
186 ajustes mediante regresión lineal simple, mientras que la cuarta opción se basó en un enfoque de mapeo

187 cuantílico. La disparidad entre las metodologías lineales y el mapeo cuantílico radica en su enfoque para
 188 abordar el sesgo promedio de las estimaciones. Las estrategias lineales buscan mitigar el sesgo promedio
 189 ajustando una regresión de primer orden que relaciona las estimaciones con las mediciones reales. En
 190 contraste, el enfoque del mapeo cuantílico modifica las estimaciones satelitales para lograr una mayor
 191 aproximación de la función de probabilidad acumulada (CDF) a la de los datos medidos. La elección de
 192 la adaptación específica se basó en la optimización de las métricas de rendimiento resultantes de cada
 193 ajuste para los respectivos modelos. Cabe mencionar que este artículo no tiene como objetivo abordar
 194 en detalle este punto.

195
 196 El desempeño de estos modelos (respecto de las mediciones en tierra) se analizó con las siguientes métricas:
 197 *MBE*, *RMSE*, *MAE* y R^2 (*MBE* es el desvío promedio (o sesgo), *RMSE* es el error cuadrático medio,
 198 *MAE* es el error absoluto medio y R^2 el coeficiente de determinación). Para los primeros tres se informan
 199 también sus correspondientes valores relativos como porcentaje de la media de las medidas terrestres,
 200 los cuales nombramos respectivamente: *MBEn*, *RMSEn* y *MAEn*. El valor de normalización en este
 201 trabajo es de 420,3 W/m². Los resultados de la evaluación de estos modelos con y sin su ajuste local se
 202 pueden ver en la Tabla 3.

203 *Tabla 3: Resultados de las métricas de desempeño de los modelos CAMS y CIM-ESRA con y sin*
 204 *adaptación al sitio respecto de las mediciones en tierra. Las medidas MBE, RMSE y MAE están medidas*
 205 *en W/m²*

	MBEn	RMSE	RMSEn	MAE	MAEn	R²
CAMS sin adaptar	-0,99	93,38	22,22	54,94	13,07	0,894
CAMS adaptado	0	91,77	21,83	55,39	13,68	0,886
CIM-ESRA sin adaptar	1,74	76,12	17,32	50,06	11,39	0,926
CIM-ESRA adaptado	0	75,29	17,13	49,31	11,22	0,931

206
 207 Se observa que los estimados del modelo CIM-ESRA demuestran una adaptación más precisa a la región
 208 en comparación con los resultados obtenidos mediante el modelo Heliosat-4. La superioridad del
 209 desempeño de CIM-ESRA se origina por dos razones: en primer lugar, hace uso de información del
 210 satélite GOES-16 en lugar de MSG, lo cual ofrece ángulos de visión más favorables para la región de la
 211 Pampa Húmeda. En segundo lugar, CIM-ESRA se caracteriza por ser un modelo semi-empírico, cuyos
 212 parámetros ajustables fueron determinados específicamente para la región en base a datos de 10
 213 ubicaciones durante el periodo 2010-2017. En consecuencia, la referencia de rendimiento para este
 214 trabajo la establece en el modelo CIM-ESRA, el cual, además, emplea información del mismo satélite
 215 que los datos de entrada utilizados en los algoritmos de aprendizaje automático presentados en este
 216 estudio. Por otro lado, se evidencia que entre las métricas examinadas, aquella que muestra una mejora
 217 es el sesgo, el cual tiende a aproximarse a 0. Mientras tanto, las demás métricas tienden a mejorar
 218 levemente, salvo en el caso del MAE del modelo CAMS donde se observa una pequeña desmejoría. En
 219 cualquier caso, la ganancia observada de los métodos de post-proceso para adaptación al sitio es
 220 limitada, estando en general por debajo del 1%.

221 *Implementación del modelo CIM-McClear*

222 Se implementó el modelo CIM-McClear con el propósito de contar con una referencia adicional, más
 223 exigente, para comparar el rendimiento frente a de los modelos de aprendizaje automático. Además, su
 224 forma de ajuste y testeo es la misma que la utilizada para los algoritmos de aprendizaje automático. Los
 225 modelos pertenecientes a la familia CIM (Cloud Index Model) se caracterizan por una estructura que
 226 combina un modelo de cielo despejado con un factor de atenuación que considera el efecto de las nubes.
 227 Este factor de atenuación, denotado como F, se rige por una función lineal que se vincula al índice de
 228 nubosidad derivado de datos satelitales como se ve en la ecuación (2) y (3).

$$229 \quad GHI = GHI_{CS} * (a + b(1 - N)) \quad (2)$$

230

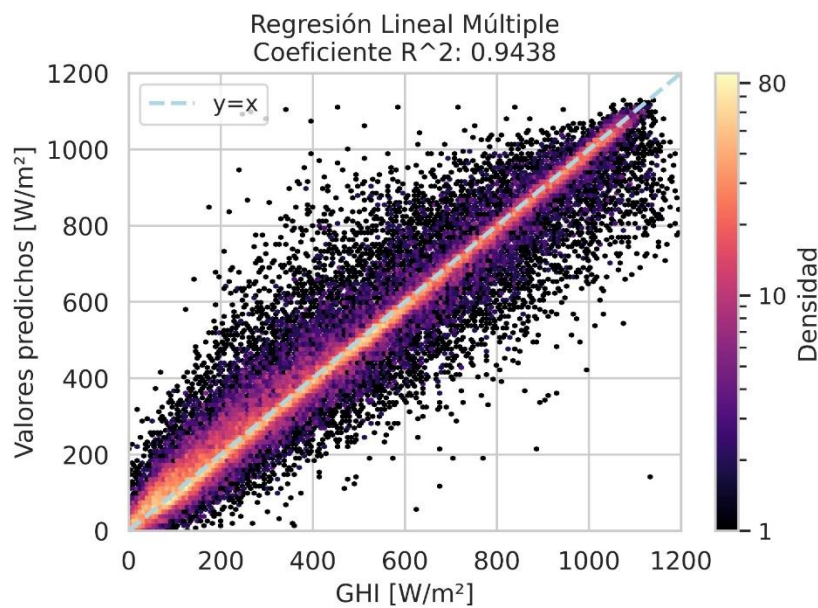
$$F(N) = a + b(1 - N) \tag{3}$$

231 donde a y b son parámetros que se ajustan localmente, GHI_{CS} es la irradiancia en condiciones de cielo
 232 despejado (en este artículo se utilizó el modelo McClear) y N es el índice de nubosidad. Se utilizó en
 233 este caso el índice N1, según la implementación usual. Para realizar la implementación, se tomaron los
 234 años 2019, 2020 y 2021, asignando dos de ellos para el entrenamiento y reservando el restante para la
 235 validación. El índice de nubosidad N corresponde a una resolución espacial de 100 x 100 (variable
 236 N1_08) la cual resultó ser la óptima. La figura 1 muestra la distribución y densidad de los puntos,
 237 comparando los valores reales con sus respectivas predicciones. Los promedios de las métricas obtenidas
 238 durante estos tres años se presentan en la tabla 4.

239 *Tabla 4: Resultados de los promedios de las métricas de desempeño del modelo implementado CIM-
 240 McClear, las medidas MBE, RMSE y MAE están medidas en W/m2*

	MBEn	RMSE	RMSEn	MAE	MAEn	R ²
CIM-McClear implementado	0	70,74	16,15	42,37	13,07	0,943

241



242

243 *Figura 1: Distribución y densidad de los puntos para los valores reales vs. valores predicho.*

244

245

246 **Implementación de los modelos de aprendizaje automático**

247 Se utilizaron como entrada las variables mencionadas en la sección *Modelos estadísticos con*
 248 *aprendizaje automático*. Para los modelos implementados se probaron diferentes combinaciones de las
 249 variables de entrada, encontrándose que para los algoritmos de aprendizaje automático fue suficiente
 250 contar con la información de las variables provenientes de imágenes satelitales. Se presentan las métricas
 251 para cada año de validación, siendo el ajuste con los otros dos años que completan el periodo 2019-
 252 2021. En la última columna se presenta el promedio de desempeño entre los 3 años de validación, lo
 253 que se usa como valor de comparación con la referencia del CIM-ESRA de la Tabla 3.

254 El primer objetivo fue saber si las correcciones en el índice de nubosidad N1 que da lugar al índice N2
 255 y la corrección a la reflectancia planetaria R que da lugar a la reflectancia RC llevan a mejoras en los
 256 resultados. Con las variables disponibles se armaron 4 distintos conjuntos de datos para encontrar la
 257 combinación de variables más adecuada. Cabe aclarar que cada conjunto contiene todo el set de variables
 258 de cada tipo en sus distintas resoluciones.

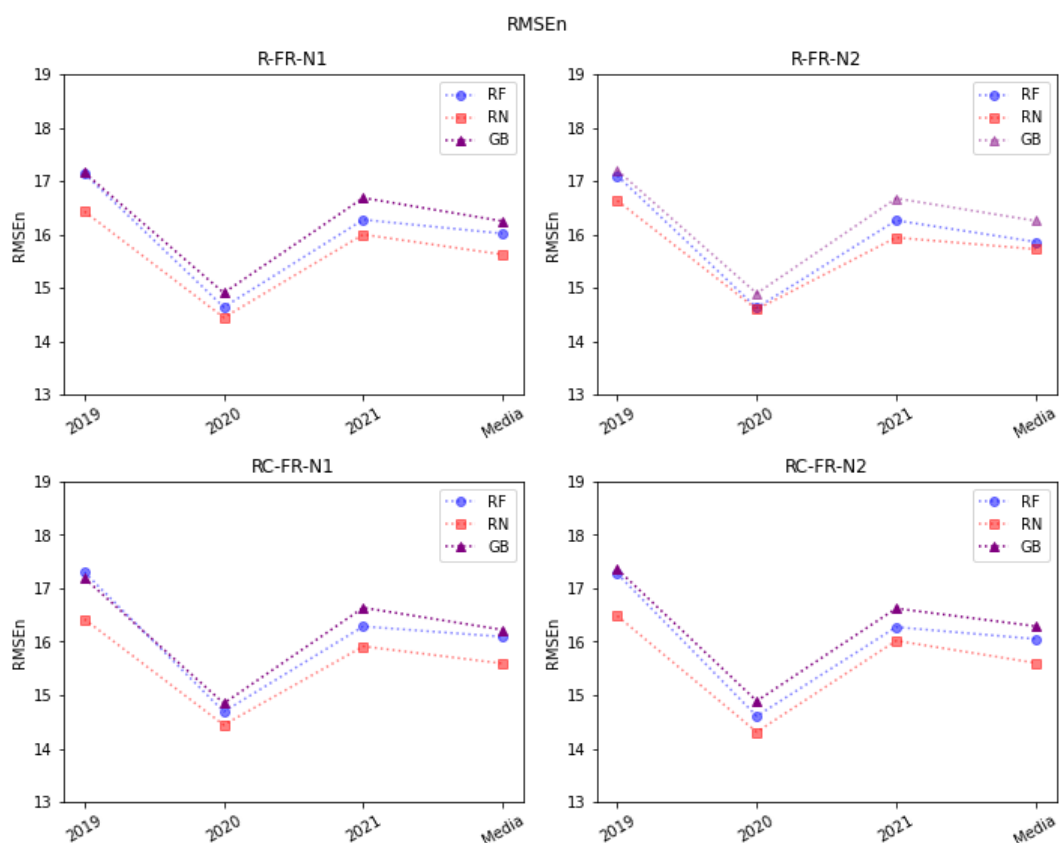
259 Conjunto 1: R - FR - N1

260 Conjunto 2: R - FR - N2

261 Conjunto 3: RC- FR - N1

262 Conjunto 4: RC - FR - N2

263 Se aplicaron los modelos de aprendizaje automático RF, GB y RN. El desempeño fue muy similar y no
264 permitió concluir qué conjunto es más adecuado. El mejor desempeño se obtuvo con RN, seguido de
265 RF y por último GB, como se muestra en la Figura 2. Con RN y RF se logró en todos los casos una
266 pequeña mejoría respecto a los resultados (Iturbide, et al. 2023) donde las variables utilizadas fueron FR
267 y R.



268

269 *Figura 2: Error cuadrático medio porcentual para los distintos modelos empleados*

270

271 **Reducción de dimensionalidad del conjunto de datos**

272 El conjunto de variables N1, N2, R y RC está altamente correlacionado. Se realizó un análisis de
273 componentes principales para disminuir la dimensión de los conjuntos, reteniendo hasta la componente
274 principal 5 inclusive. En todos los conjuntos el porcentaje de varianza acumulada con estas 5
275 componentes es mayor al 99.5% y se recuperan los valores RMSEn obtenidos al trabajar con todas las
276 variables (diferencias menores al 0,05%). Tomar más componentes principales no mejora los
277 desempeños (ver tabla 4).

278 Se buscaron qué variables eran las más relevantes para el modelo. Para eso, se analizó primero qué
279 variable presenta mayor correlación con la variable a predecir. En todos los casos (FR, R, RC, N1 y N2)
280 estas variables se encuentran alrededor de la resolución media (variables 9 o 10). Se buscaron 2 variables

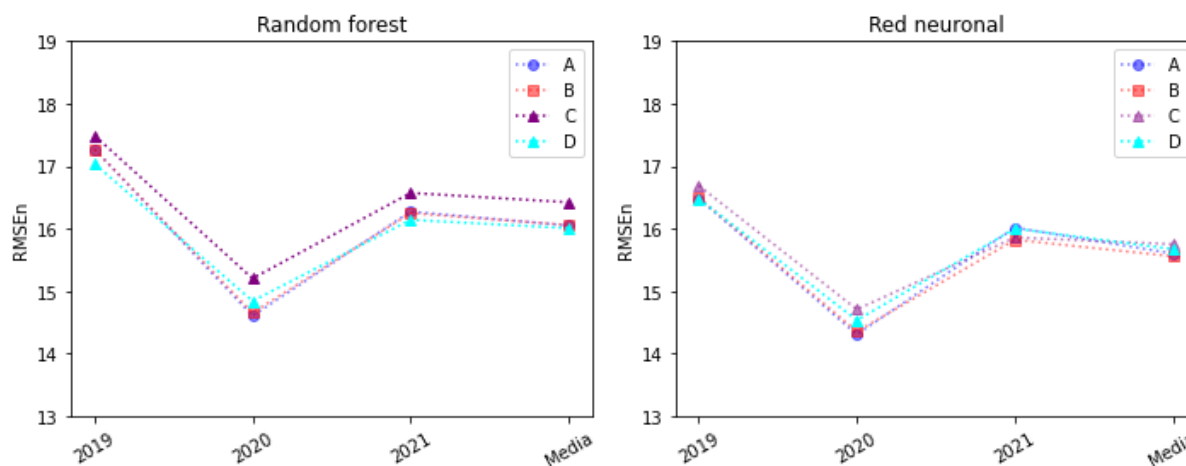
281 más en los extremos de las resoluciones. La búsqueda no fue exhaustiva y se encontraron mejores
 282 resultados con RN que con RF. Cuando se reduce la dimensión del conjunto de esta manera, RF aumenta
 283 0,4% aproximadamente respecto al desempeño del algoritmo sobre el conjunto de datos que contiene a
 284 todas las variables, mientras que la RN da valores más cercanos y aumenta 0,2% como máximo. Para
 285 mejorar el desempeño se agregó la variable geométrica cosz. Con la referencia geométrica los dos
 286 algoritmos se aproximan a los valores de desempeño cuando se utilizan todas las variables.

287 Los resultados de los desempeños promediados en los 3 años se muestran en la Figura 3 para los
 288 algoritmos RF y RN en los conjuntos de datos formados con: todas las variables (A), las 5 primeras
 289 componentes principales (B), una selección de variables (C) y la selección anterior más cosz (D). En
 290 particular se muestran los resultados del conjunto 4: RC - FR - N2. La selección C corresponde a FR2,
 291 FR10, FR14, N2_1, N2_9, N2_14, RC1, RC10 y RC14.

292 *Tabla 4: Resultados de los modelos de ML utilizando las variables FR-RC-N2.*

	MBEn	RMSE	RMSEn	MAE	MAEn	R ²
RF A	-0,08	70,25	16,05	42,86	9,80	0,939
RN A	0,03	68,32	15,60	41,88	9,57	0,943
RF B	-0,07	70,29	16,06	42,92	9,81	0,939
RN B	0,05	68,14	15,56	41,44	9,47	0,944
RF C	-0,04	71,92	16,42	43,84	10,02	0,937
RN C	-0,34	69,01	15,75	41,29	9,43	0,943
RF D	0,05	70,11	16,01	41,95	9,59	0,941
RN D	0,04	68,65	15,67	41,78	9,55	0,943

293
294

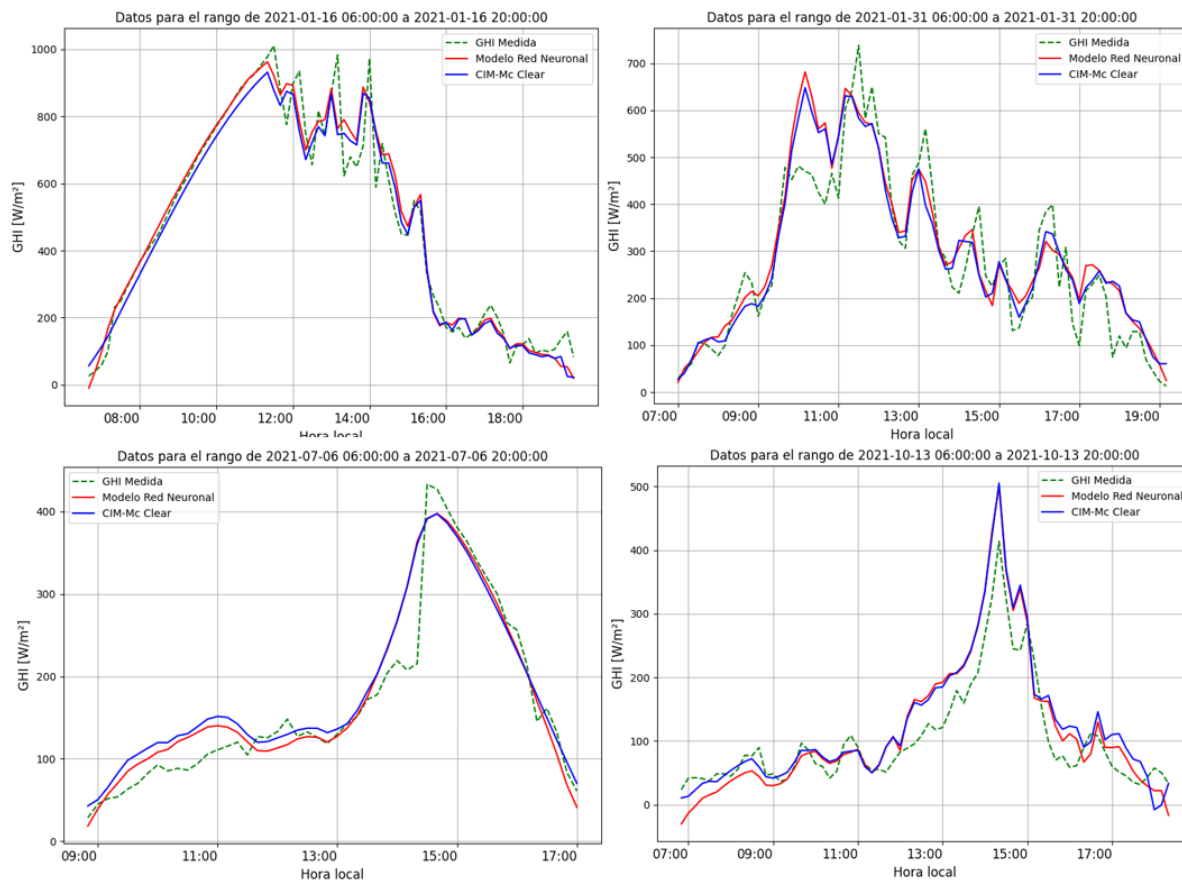


295 *Figura 3: Error cuadrático medio porcentual para los algoritmos RF y RN de las 4 selecciones de*
 296 *datos utilizadas para las variables RC - FR - N2*
 297

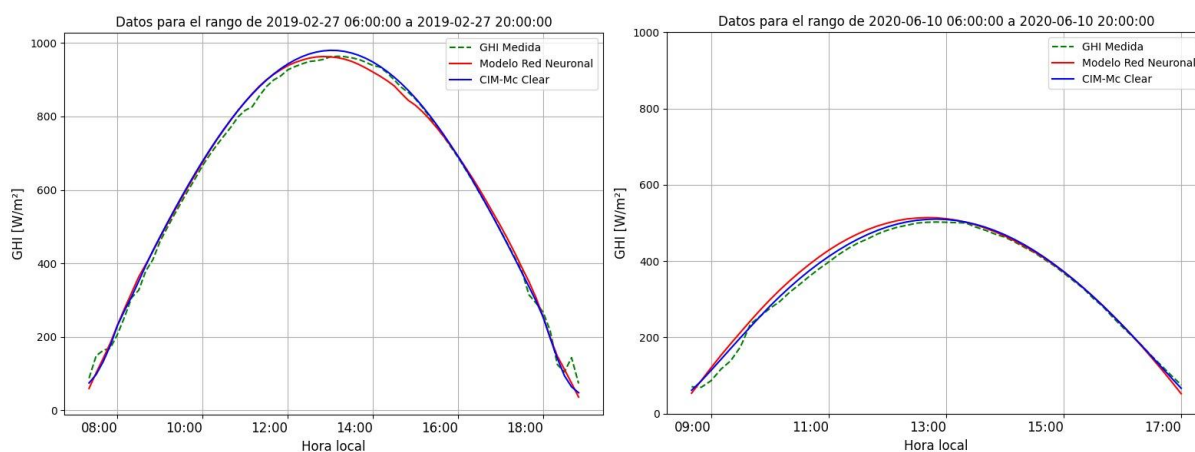
298
299
300 **Visualización comparativa de estimaciones**

301 Las Figuras 4 y 5 exhiben las proyecciones generadas por la RN (representadas por la línea roja) en
 302 contraste con los valores registrados en tierra (indicados por la línea punteada), junto con la estimación
 303 obtenida mediante la implementación del modelo CIM-McClear, que se distinguió como la más precisa

304 en términos de comparación, para días con nubosidad y claros, respectivamente. Existe notoria
305 concordancia entre las aproximaciones, aunque se aprecian sutiles diferencias en las mismas. Tanto la
306 RN como el modelo CIM-McClear logran de manera efectiva capturar las condiciones de nubosidad y
307 la GHI en términos generales.
308



309 *Figura 4: Gráficos comparativos entre las medidas de tierra, el modelo CIM-McClear implementado*
310 *ya la red neuronal para cuatro días de 2021 en condiciones de nubosidad.*
311
312



313 *Figura 5: Gráficos comparativos entre las medidas de tierra, el modelo CIM-McClear implementado*
314 *ya la red neuronal para dos días despejados de 2021.*
315
316

317 CONCLUSIONES

318 Entre los diversos algoritmos de Aprendizaje Automático utilizados en este estudio, se observó un mejor
319 rendimiento por parte de la RN, seguido del método de RF, mientras que el método de GB quedó

320 rezagado. La combinación de variables FR, RC y N2 mostró los resultados más favorables, obteniendo
321 un error cuadrático medio promedio porcentual del 15,56% después de la reducción de la
322 dimensionalidad mediante el análisis de componentes principales. Al comparar con el modelo de
323 referencia CIM-ESRA, los modelos empíricos de ML demostraron un rendimiento superior. Las
324 métricas comparativas entre CIM-ESRA y la RN favorecieron a esta última: MAEn de 9,7% vs. 11,3%,
325 RMSEn de 15,6% vs. 17,1%.

326 Una comparación interesante surge con el modelo CIM-McClearn, que logró un RMSEn promedio de
327 16,15% y un sesgo nulo. Cabe destacar que este modelo requiere solo la entrada del modelo de cielo
328 claro y una variable de nubosidad parametrizada mediante una función lineal. Su simplicidad ofrece un
329 rendimiento menor pero similar al aprendizaje automático, que además utiliza información espacial
330 multiescala.

331 El análisis de componentes principales no redujo de manera significativa los errores. En resumen, el
332 modelo de ML propuesto mejora las estimaciones para la región, superando a modelos ajustados al sitio
333 y al CIM-McClearn ajustado localmente, con métricas de desempeño levemente mejores.

334 En futuros trabajos se debe analizar el comportamiento del modelo empírico en otras áreas de la Pampa
335 Húmeda, extrapolando espacialmente mediante pruebas en una ubicación distinta. También sería
336 recomendable considerar otras variables satelitales relevantes, como los canales infrarrojos del satélite,
337 que brindan información adicional sobre el sistema Tierra-Atmósfera y podrían mejorar la precisión de
338 las estimaciones de radiación solar. Evaluar la inclusión de este canal en el modelo y su impacto en el
339 rendimiento sería valioso.

340 **AGRADECIMIENTOS**

341 R. Alonso-Suárez agradece a la Comisión Sectorial de Investigación Científica (CSIC), Udelar, por el
342 apoyo financiero al Laboratorio de Energía Solar a través de su programa de Grupos de I+D.

343 Los investigadores del GERSolar agradecen a la Secretaría de Ciencia y Tecnología (SCyT) de la UNLu
344 por el apoyo financiero a través de los Proyectos de Investigación para Investigadores en Formación
345 2022 (RESREC-LUJ: 266-22) y al Ing. Lucas Burgos por la recopilación y el control de calidad
346 realizado a los datos de la estación Luján utilizados en este trabajo.

347 **REFERENCIAS**

- 348 Abal, G., Aicardi, D., Alonso-Suárez, R., y Laguarda, A. (2017). Performance of empirical models for
349 diffuse fraction in Uruguay. *Solar Energy*, 141:166–181.
- 350 Alonso-Suárez, R., Abal, G., Siri, R., y Musé, P. (2012). Brightness-dependent Tarpley model for global
351 solar radiation estimation using GOES satellite images: application to Uruguay. *Solar Energy*, 86 ,
352 3205–3215. doi:10.1016/j.solener.2012.08.012.
- 353 Gonzalez, J., Teixeira-Branco, V., y Alonso-Suárez, R. (2019). Evaluation of the Heliosat-4 and
354 FLASH-Flux models for solar global daily irradiation estimate in Uruguay. En *ISES Conf.*
355 *Proceedings, Solar World Congress*.
- 356 Iturbide, P., Alonso-Suarez, R., Ronchetti, F. (2023). An Analysis of Satellite-Based Machine Learning
357 Models to Estimate Global Solar Irradiance at a Horizontal Plane. In: Naiouf, M., Rucci, E.,
358 Chichizola, F., De Giusti, L. (eds) *Cloud Computing, Big Data & Emerging Topics. JCC-BD&ET*
359 *2023. Communications in Computer and Information Science*, vol 1828. Springer, Cham.
360 https://doi.org/10.1007/978-3-031-40942-4_9
- 361 Jiménez, V. A., Will, A., & Rodríguez, S. (2017). Estimación de radiación solar horaria utilizando
362 modelos empíricos y redes neuronales artificiales. *Ciencia y tecnología*, (17), 29-45.
- 363 Laguarda, A., Iturbide, P., Orsi, X., Denegri, M. J., Luza, S., Burgos, B. L., Stern, V., y Alonso-Suárez,
364 R. (2021). Validación de modelos satelitales Heliosat-4 y CIM-ESRA para la estimación de
365 irradiancia solar en la Pampa Húmeda. *Energías Renovables y Medio Ambiente*, 48, 1-9.

- 366 Laguarda, A., Giacosa, G., Alonso-Suárez, R., y Abal, G. (2020). Performance of the site-adapted
367 CAMS database and locally adjusted cloud index models for estimating global solar horizontal
368 irradiation over the Pampa Húmeda region. *Solar Energy*, 199:295–307.
- 369 Laguarda, A., Alonso-Suárez, R., y Abal, G. (2018). Modelo semi-empírico de irradiación solar global
370 a partir de imágenes satelitales GOES. *Anales del VII Congreso Brasileiro de Energia Solar*.
- 371 Lefèvre, M., Oumbe, A., Blanc, P., Espinar, B., Qu, Z., Wald, L., Homscheidt, M. S., y Arola, A. (2013).
372 McClear: a new model estimating downwelling solar radiation at ground level in clear-sky
373 conditions. *Atmospheric Measurement Techniques*, European Geosciences Union, 6 , 2403–2418.
374 doi:10.5194/amt-6-2403-2013.
- 375 Long, C. N., & Shi, Y. (2008). An automated quality assessment and control algorithm for surface
376 radiation measurements. *The Open Atmospheric Science Journal*, 2(1).
- 377 McArthur, L. (2005). Baseline Surface Radiation Network (BSRN) Operations Manual. Td-no. 1274,
378 wrcp/wmo, World Meteorological Organization (WMO, www.wmo.org).
- 379 Olivera, L., Atia, J., Amet, L., Osio, J., Morales, M., & Cappelletti, M. (2020). Uso de redes neuronales
380 artificiales para la estimación de la radiación solar horaria bajo diferentes condiciones de cielo.
381 *Avances en Energías Renovables y Medio Ambiente-AVERMA*, 24, 232-243.
- 382 Perez, R., Ineichen, P., Seals, R., & Zelenka, A. (1990). Making full use of the clearness index for
383 parameterizing hourly insolation conditions. *Solar Energy*, 45(2), 111-114.
- 384 Perez, R., Cebecauer, T., & Šúri, M. (2013). Semi-empirical satellite models. *Solar energy forecasting
385 and resource assessment*, 21-48.
- 386 Qu, Z., Oumbe, A., Blanc, P., Espinar, B., Gesell, G., Gschwind, B., Klüser, L., Lefèvre, M., Saboret,
387 L., Schroedter-Homscheidt, M., y Wald, L. (2017). Fast radiative transfer parameterisation for
388 assessing the surface solar irradiance: The Heliosat-4 method. *Meteorologische Zeitschrift*,
389 26(1):33–57.
- 390 Raichijk, C. (2008). Estimación de la irradiación solar global en Argentina mediante el uso de redes
391 neuronales. *Energías Renovables y Medio Ambiente (ISSN 0328-932X)*. Vol. 22, pp. 1 - 6.
- 392 Salazar, G. A., Alonso-Suárez, R., Cirigliano, A. L., y Ledesma, R. D. (2021). Evaluación del proceso
393 de adaptación al sitio aplicado a la irradiancia solar global medida en la ciudad de Salta, Argentina.
394 *Avances en Energías Renovables y Medio Ambiente-AVERMA*, 25, 353-362.
- 395 Sarazola, I., Laguarda, A., Ceballos, J. C., y Alonso-Suárez, R. (2023). Benchmarking of modeled solar
396 irradiation data in Uruguay at a daily time scale. *IEEE Latin American Transactions*.
- 397 Sayago, S., Bocco, M., Ovando, G., & Willington, E. A. (2011). Radiación solar horaria: modelos de
398 estimación a partir de variables meteorológicas básicas. *Avances en Energías Renovables y Medio
399 Ambiente*, 15.
- 400 Verbois, H., Saint-Drenan, Y.-M., Becquet, V., Gschwind, B., Blanc, P. (2023). Retrieval of surface
401 solar irradiance from satellite using machine learning: pitfalls and perspectives, *EGUsphere*
402 [preprint], <https://doi.org/10.5194/egusphere-2023-243>.
- 403 Yang, D. (2020). Choice of clear-sky model in solar forecasting. *Journal of Renewable and Sustainable
404 Energy* 12, 026101, <https://doi.org/10.1063/5.0003495>.
- 405 Young, A.T. (1994). Air mass and refraction. *Applied optics*, 33 6, 1108-10 .
- 406

407 MACHINE LEARNING MODELS FOR ESTIMATING SOLAR HORIZONTAL RADIATION 408 IN THE PAMPA HÚMEDA WITH MULTISCALE SATELLITE INFORMATION 409

410 **ABSTRACT:** The lack of precision in solar radiation data impacts the solar energy projects risk. Ground
411 measurement networks provide limited information due to their sparse spatial distribution. This leads to
412 estimation models based on satellite imagery, solving the spatial issue if carefully adjusted to quality
413 ground measurements. In this article, we develop and validate an empirical Machine Learning (ML)
414 model for satellite-based solar radiation estimation, demonstrating its usefulness and accuracy in the
415 studied region. The models are fed with variables from GOES-16 satellite imagery, McClear model
416 estimates, and geometric data. Our results suggest that for certain proposed models, satellite information
417 is sufficient for accurately estimating solar radiation, by obtaining the temporal reference from implicit
418 relationships between the considered satellite variables. Given the size of the data set, we propose a
419 principal component analysis to reduce dimensionality. In order to compare the proposed model, we
420 adapt Heliosat-4 and CIM-ESRA estimates to the site and implement the CIM-McCclear model. The

421 results indicate that the proposed model outperforms others, although slightly, showing how difficult it
422 is to further improve solar radiation satellite-based estimation.

423 **Keywords:** Solar radiation, Machine Learning, Satellite images, GOES16, GHI.