



PROGRAMA
ESTADO DE LA NACIÓN

Programa Estado de la Nación

Investigación especial

Modelo de machine learning para
estimar caudales mensuales
históricos de cuencas hidrográficas
en Costa Rica

Investigador:

Darío Rodríguez-García

San José | 2025



333.91
R696m

Rodríguez-García, Darío

Modelo de machine learning para estimar caudales mensuales históricos de cuencas hidrográficas en Costa Rica / Darío Rodríguez-García. -- San José, C.R. : PEN, 2025.

1 recurso en línea (75 páginas): archivos de texto PDF, 4250 KB

ISBN 978-9930-654-81-1
Investigación especial 2025

1. APRENDIZAJE AUTOMÁTICO. 2. CUENCAS HIDROGRÁFICAS. 3. PRECIPITACIÓN ATMOSFÉRICA. 4. CULTIVOS DE COBERTURA. I. Título.



Información de la persona autora:

Darío Rodríguez-García. <https://orcid.org/0000-0002-7875-8423>

Esta obra se comparte bajo la licencia Reconocimiento
– No Comercial – Compartir Igual (CC-BY-NC-SA)

Permite usar una obra para crear otra obra o contenido, modificando o no la obra original, siempre que se cite al autor, la obra resultante se comparte bajo el mismo tipo de licencia y no tenga fines comerciales

Permite usar una obra para crear otra obra o contenido, modificando o no la obra original, siempre que se cite al autor, la obra resultante se comparte bajo el mismo tipo de licencia y no tenga fines comerciales



Índice

Descargo de responsabilidad.....	5
Introducción.....	5
Objetivo	8
Metodología	8
Fuentes de información.....	10
Caudales.....	10
Precipitación y temperatura	12
Suelo	13
Elevación y pendiente	15
Cobertura del suelo	16
Regiones climáticas.....	17
Período de estudio	20
Métricas	20
Tipos de errores	21
Línea base nacional	22
Línea base internacional.....	23
Métricas esperadas	24
Análisis exploratorio de los datos	25
Caudales.....	25
Características de las cuencas	27
<i>Ubicación de las estaciones.....</i>	29
<i>Delimitación de cuencas.....</i>	29
<i>Área y perímetro</i>	30
<i>Pendiente y elevación media</i>	32
<i>Región climática</i>	33
Ingeniería de variables	34
Rezagos	34

Suelos	34
Fracción que infiltra por efecto de la cobertura del suelo (Kv)	35
Fracción que infiltra por efecto de la pendiente (Kp)	35
Coefficiente de infiltración	36
Resumen de variables	37
Conjuntos de entrenamiento, validación y prueba	39
Algoritmo predictivo y entrenamiento	40
Resultados del modelo	41
Errores de validación fuera de cuenca	41
Errores de prueba	43
Errores línea base nacional	44
Conclusiones	44
Interpretación del modelo	45
Impacto del cambio climático	48
Caso de estudio cuenca Chirripó Pacífico	53
Descripción de la cuenca	53
Precipitación	56
Temperatura	58
Caudales	60
Cambio climático	65
Conclusiones y recomendaciones	66
Bibliografía	69
Anexos	73

Descargo de responsabilidad

Esta investigación se realizó para el *Informe Estado de la Nación 2025*. El contenido es responsabilidad exclusiva de su autor, y las cifras pueden no coincidir con las consignadas en el capítulo respectivo, debido a revisiones posteriores. En caso de encontrarse diferencia entre ambas fuentes, prevalecen las publicadas en el Informe.

Introducción

El agua desempeña un papel fundamental en diversos aspectos, como el mantenimiento de los ecosistemas y biodiversidad, el abastecimiento de agua potable, la producción de alimentos, turismo, recreación, la generación de energía y la industria.

La gestión integral del recurso hídrico (GIRH) se refiere a un enfoque holístico que considera todos los aspectos relacionados con el agua, desde su captación y distribución hasta su uso y conservación, teniendo en cuenta tanto las necesidades humanas como las demandas ambientales. Esta gestión es particularmente difícil puesto que debe tomar en cuenta una pluralidad de factores como las comunidades locales, sector productivo, el medio ambiente y los complejos procesos hidrológicos como los efectos del cambio climático.

Una gestión sin información de calidad, oportuna y con el detalle adecuado es peligrosa y conlleva a toma de decisiones subjetivas que pueden ir en perjuicio de la sostenibilidad del recurso. Por la transversalidad del agua en las actividades humanas y naturales la información relacionada es extensa y variopinta, sin embargo, una de las más relevantes es la de los caudales fluviales.

El conocimiento de los caudales de los ríos a través del tiempo es particularmente importante ya que permite responder preguntas clave de diferentes sectores, tales como: ¿Cuánta agua pasa por este punto en verano?, ¿El río tiene la capacidad de brindar 3 l/s a un proyecto durante todo el año?, ¿Se aprecia una tendencia bajista del caudal?, ¿Cuánto caudal se debería reservar para mantener un ecosistema saludable?, ¿Cómo podemos organizar la disponibilidad regional de agua entre sus comunidades?, ¿Cuál es el potencial hidroeléctrico de un cauce?, ¿Cuál sería la concentración de contaminantes según su descarga?, ¿Durante que períodos y

con qué regularidad se debería brindar apoyos a ecosistemas vulnerables?, ¿Dónde se deben priorizar acciones de mitigación al cambio climático?, entre otros.

Conocer la disponibilidad de agua de manera directa es complejo y se necesitan instrumentos que midan el caudal a lo largo de un período de tiempo que permita capturar estacionalidades anuales y también fenómenos relevantes como tormentas tropicales, el fenómeno del Niño, entre otros. Además, estos instrumentos deben estar ubicados exactamente en el sitio de interés. Debido a que pocas veces se cuenta con un registro temporal y espacialmente pertinente, se recurre a técnicas indirectas que estiman los caudales como el traslado de caudales o balances hídricos.

Es común que los actores interesados recurran a estas técnicas para estimar los caudales que, si bien en muchas ocasiones cumplen su función, tienen ciertas limitaciones técnicas:

- Tienen a ser aplicables en un área pequeña y se desconoce la incertidumbre de la generalización del modelo en otros sitios.
- Tienen a solo utilizar información en las cercanías del punto de estudio y esta información tiende a tener una alta cantidad de datos faltantes.
- La información de caudales disponible es costosa y escasa en términos espaciales y temporales.
- La información meteorológica producida a nivel nacional no es pública.
- Al no existir un proceso estándar ni un único revisor de los estudios, los procesos de validación y prueba de los resultados no son homogéneos y pueden ser subjetivos.
- Los análisis pierden vigencia rápidamente por los cambios climatológicos y las condiciones de la cuenca.

A nivel administrativo y de acceso a la información:

- Tienen un alto costo económico.
- Los procesos de contratación, ejecución y análisis de resultados son largos y pueden retrasar la ejecución de proyectos relevantes que dependan de estos análisis.
- Generalmente son contratados por la persona interesada en el recurso hídrico y esto genera un conflicto de intereses.
- Los estudios tienden a ser archivados en las instituciones o empresas una vez que cumplen su propósito e impide que los descubrimientos sean utilizados en estudios posteriores.
- Organizaciones de base comunitaria e individuos no cuentan con los recursos económicos necesarios para la ejecución de balances hídricos.
- No existen herramientas de fácil acceso que permita a las comunidades y otros actores a acceder e interactuar con los modelos hidrológicos.
- El personal de las instituciones que no está altamente calificado normalmente no aprovecha de mejor manera los resultados de los estudios por su complejidad, alta barrera de acceso y dificultad de obtener resultados personalizados.

En este marco, la investigación busca probar la factibilidad de una herramienta que permita hacer uso de una vasta variedad de fuentes de información y estaciones hidrológicas, que sea escalable en una parte del territorio nacional, que se mantenga pertinente en el tiempo, que utilice técnicas modernas de aprendizaje automático y que cuente con información transparente respecto a la validez técnica de los resultados. Específicamente este estudio se concentrará en realizar y probar un modelo que permita estimar los caudales superficiales mensuales desde 1981 al 2021 -inclusive- en ríos que no cuentan con estaciones hidrológicas y que comparten similitudes con ríos que si las tienen a nivel nacional. Adicionalmente el estudio mostrará eventuales usos de este modelo como la aplicación de escenarios de cambio climático, análisis de tendencia y curvas de duración.

Objetivo

El objetivo principal de la consultoría es diseñar y construir un modelo de machine learning que estime los caudales mensuales históricos en cuencas hidrográficas en el país.

El modelo será validado en la cuenca Río Grande Candelaria (estación el Rey) y el Río General (estación El Brujo). Adicionalmente, se analizará su sensibilidad respecto al impacto del cambio climático y se aplicarán análisis hidrológicos derivados en una cuenca hidrológica.

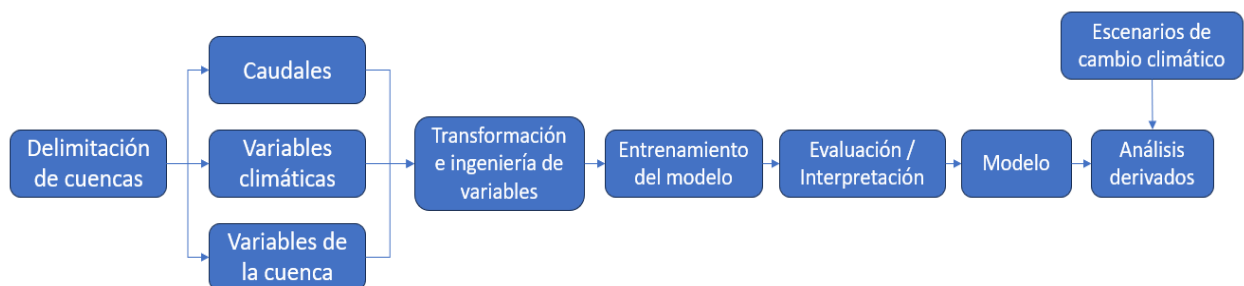
Tales análisis incluirán la estimación de los caudales históricos, los comportamientos del cauce por mes, las tendencias en el flujo, curvas de duración, entre otros.

Metodología

Esta investigación creará y probará modelos agregados por cuenca basados en datos que aprenden de comportamientos registrados en las cuencas y permiten predecir el caudal para un mes en específico. Modelos de esta categoría -aprendizaje automático- han demostrado ser en algunos casos una alternativa superior a los tradicionales en cuanto a la exactitud de sus predicciones (Kratzert, y otros, 2019). Específicamente se utilizarán modelos de tipo ensemble por su versatilidad y su aceptación en investigaciones hidrológicas en los últimos años (Zounemat-Kermani, Batelaan, Fadaee, & Hinkelmann, 2021).

A nivel general, el proceso se puede describir en la siguiente figura.

Figura 1
Metodología



Fuente: Elaboración propia.

En primer lugar, se deben delimitar las cuencas en estudio basadas en el sitio de la estación hidrológica. Esta delimitación permitirá en las siguientes etapas el cálculo de las variables climáticas y físicas de la cuenca.

Acto seguido, la investigación utilizará datos disponibles que se pueden agrupar en tres componentes, en primera instancia las mediciones de caudal de las estaciones de estudio, luego variables climáticas como la precipitación y la temperatura y por último variables que están relacionadas con las condiciones físicas de la cuenca como su topografía, geología y uso del suelo.

Esta información es procesada y transformada en variables que son útiles para el proceso de aprendizaje automático. En esta fase, se calcularon variables como el coeficiente de infiltración, retrasos de variables climáticas para diferentes períodos y se agregaron los resultados a nivel de cuenca.

Posteriormente el modelo es entrenado probando una serie de hiperparámetros que maximicen la función objetivo –Kling Gupta Efficiency-. Luego, el modelo es evaluado en cuanto a los errores reportados durante la validación y prueba y los aportes de las variables al proceso de aprendizaje. Según esta evaluación se puede iterar sobre las fases anteriores hasta llegar a una evaluación satisfactoria.

A continuación, el mejor modelo es seleccionado y comparado contra las métricas de la línea base.

Finalmente se utiliza el modelo para hacer predicciones, aplicar escenarios de cambio climático y realizar análisis hidrológicos derivados como curvas de duración y análisis de tendencias.

Fuentes de información

Los modelos de aprendizaje automático utilizan los datos como fuente principal para su aprendizaje y por ello la calidad y cantidad de información tendrá un impacto directo en la calidad de los resultados del modelo.

En esta investigación la información recopilada cuenta con varias características:

- Espacialmente adecuada: La información cuenta con un componente geográfico que le permite ser asociada a una cuenca. Debe abarcar Costa Rica y el sector de Panamá que está contenido dentro de la cuenca del Río Sixaola.
- Temporalidad 1981-2021: En la medida de lo posible, la información con componentes temporales debe tener un registro compatible con el período 1981 - 2021 a nivel mensual.
- Disponibilidad: La información de los predictores (toda a excepción de los caudales) deberá ser pública y de actualización continua con el fin de facilitar predicciones en el futuro, facilitar transparencia y reproducibilidad.
- Hidrológicamente pertinente: La información está asociada a diferentes procesos del ciclo hidrológico.

A continuación, se describen las fuentes de información utilizadas:

Caudales

Corresponde a un variable numérica con unidades de mm promedio mensual -transformada desde m^3/s - y es la información más relevante puesto que es la variable para predecir.

Esta información es la única que no está disponible públicamente. Para esta investigación, el Instituto Costarricense de Electricidad (ICE) compartió la información de 1981 al 2021 de 75 estaciones hidrológicas a lo largo del país. Lo anterior en el marco de un convenio entre el ICE y el Programa Estado de la Nación/Conare.

Luego de un proceso de consulta con el ICE, se decidió eliminar seis estaciones que no representan el comportamiento natural de un cauce -proceso a modelar- ya que estaban influenciadas por proyectos hidroeléctricos -Presa Peñas Blancas, Paso Hondo, El Congo, Bebedero, Hamburgo y Murcia-. Además, se eliminaron unos períodos de las estaciones Puente

la Virgen (>01/01/2019), Angostura (> 01/0/2021) y Pocosol (<01/01/2010) por afectaciones antropogénicas y datos anómalos.

Los caudales tienen la particularidad de que su medición es mediante métodos indirectos, discontinuos y en algunos casos difíciles de cuantificar y por lo tanto tienden a tener una incertidumbre importante.

La mayoría de los datos de caudales se derivan de una serie de medidas del nivel del río y una curva de descarga que relaciona el nivel con el flujo. Por lo tanto, los datos de flujo resultantes son estimaciones que contienen incertidumbres relacionadas con: las mediciones de los niveles, las mediciones de niveles y flujo utilizadas para derivar la curva de calibración, interpolación y extrapolación en el modelo de la curva de calibración, y cambios en la sección transversal del canal debido a erosión/relleno, vegetación y hielo, efectos de aguas abajo e histéresis, que provocan cambios en una curva de calibración determinística. (McMillan, y otros, 2017)

Esta condición de incertidumbre es particularmente relevante ya que transmite a la variable objetivo un componente de ruido que afectará las métricas del modelo en comparación con un modelo en donde la variable objetivo tiene una incertidumbre reducida. Ahora bien, puesto que se cuenta con más de 20 000 observaciones, se espera que el modelo sea capaz de reducir el componente de ruido al generalizar.

El ICE toma la información de caudal a nivel diario y lo promedia mensualmente. En algunas ocasiones, los meses no cuentan con toda la información diaria y en esos casos el ICE computa el promedio para meses con menos de 8 días faltantes o correlaciona la información diaria de otras estaciones.

Durante el análisis, toda la información disponible de caudales será utilizada independientemente si el promedio mensual fue generado con los días completos, con algunos días correlacionados o si tienen de uno a siete días sin datos.

Precipitación y temperatura

La precipitación y la temperatura normalmente se miden de manera directa en estaciones meteorológicas ubicadas en los sitios de interés. En Costa Rica esta información tiene varios limitantes como la distribución espacial de las estaciones, la cantidad de datos faltantes, la gobernanza dividida de la información y su acceso restringido.

A nivel internacional existen metodologías modernas que utilizan diversas fuentes de información para producir estimaciones de la precipitación y temperatura a nivel global, continuas y en horizontes temporales adecuados para esta investigación.

Para la precipitación, se utilizará la información de CHIRPS - Climate Hazards Group InfraRed Precipitation with Station data- por su amplio período de análisis, disponibilidad, completitud, escala y sus resultados satisfactorios en varios estudios nacionales (Venegas-Cordero, y otros, 2021), (Kaune, 2021), (Arciniega-Esparza, Birkel, Chavarría-Palma, Arheimer, & Agustín, 2022).

CHIRPS es un conjunto de datos de precipitaciones casi globales de más de 35 años. Cuenta con una extensión de 50°S-50°N -y todas las longitudes-, desde 1981 hasta casi el presente, CHIRPS incorpora algoritmos internos de los desarrolladores, imágenes satelitales de resolución de 0,05° -5 km aproximadamente- y datos de estaciones in situ para crear series cronológicas de lluvia cuadrículadas. (Funk, 2014). La versión utilizada es la 2.0 y fue descargada de las bases de datos de Google Earth Engine en agosto 2023.

En cuanto a la temperatura se utilizará el producto ERA5-LAND (Muñoz Sabater, 2019) producido por Copernicus. ERA5-Land es un conjunto de datos global de reanálisis que combina datos de modelos con observaciones de todo el mundo que proporciona variables terrestres desde 1950 a una resolución de 0,1° -9 km aproximadamente-. La variable utilizada corresponde a la temperatura de 2m por encima de la superficie. La información fue descargada del portal Climate Data Store en agosto 2023.

Suelo

En esta investigación se utilizará el producto SoilGrids250 (Hengl, y otros, 2017) que contiene información de propiedades del suelo a una resolución de 250 m. El producto fue construido aplicando algoritmos de aprendizaje automático en un *set* de datos globales de perfiles de suelos y de mediciones de sensores remotos.

El producto incluye:

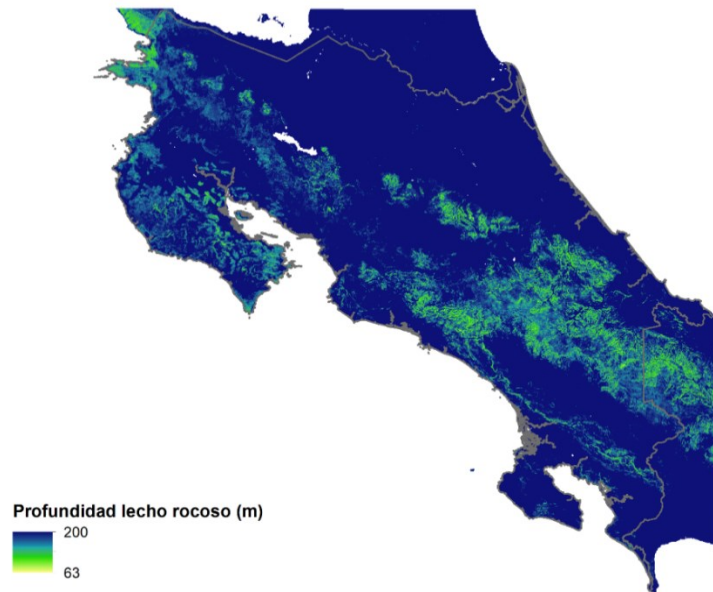
- Textura del suelo en 7 capas (0, 5, 15, 30, 60, 100, y 200 cm).
- Profundidad del suelo hasta el estrato rocoso.

Fracción volumétrica de fragmentos gruesos en 7 capas (0, 5, 15, 30, 60, 100, y 200 cm).

El producto fue comparado con varios conjuntos de datos de suelos existentes en el contexto de la aplicabilidad del modelo y se concluyó que SoilGrids250 es el estado del arte en datos de suelo (Dai, y otros, 2019).

Mapa 1

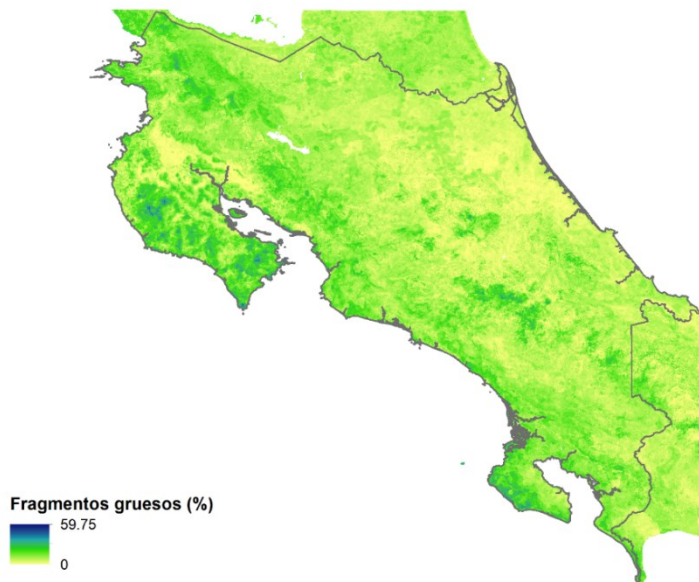
Profundidad del lecho rocoso



Fuente: SoilGrids250, 2017.

Mapa 2

Porcentaje de fragmentos gruesos

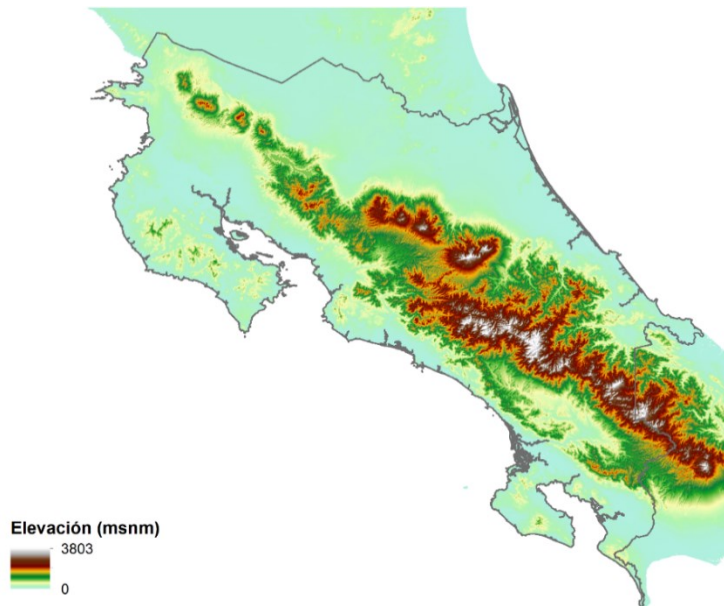


Fuente: SoilGrids250, 2017.

Elevación y pendiente

Corresponde a los modelos de elevación digital del *Shuttle Radar Topography Mission* (Tom G. Farr, 2007) el cual es un esfuerzo internacional en una escala casi global, tiene una resolución de 30m. La información fue descargada de *Google Earth Engine* en agosto de 2023.

Mapa 3
Elevación

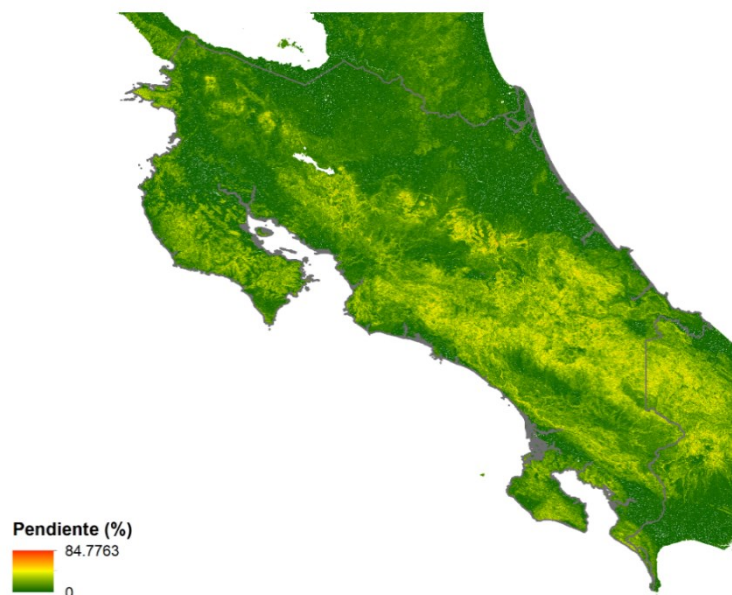


Fuente: Shuttle Radar Topography Mission, Farr y otros. 2007.

Con base en esa información se calcularon las pendientes a nivel nacional:

Mapa 4

Pendiente (%). Calculada del modelo de elevación digital de Shuttle Radar Topography Mission, Far y otros. 2007



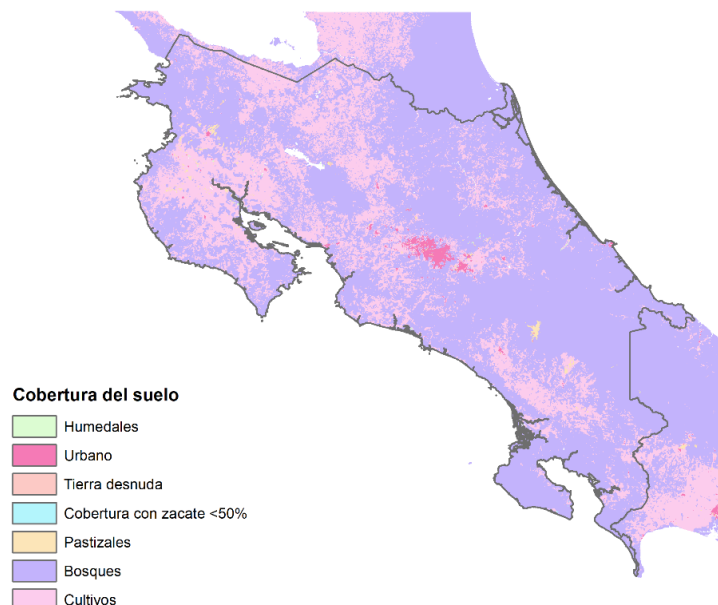
Fuente: Shuttle Radar Topography Mission, Far y otros. 2007.

Cobertura del suelo

Corresponde a los mapas de cobertura del suelo realizados por *Copernicus Climate Change Service* (Copernicus Climate Change Service, 2019). La base de datos describe la tierra en 22 clases definidas por la Organización de las Naciones Unidas para la Alimentación y Agricultura. La información tiene una escala global con una resolución de 300m e incluye los datos desde 1992 hasta el presente. La información fue descargada en agosto de 2023. Se utilizó la versión 2.0.7 para los años 1992 a 2015 y del 2015 en adelante la versión 2.1.1.

Por la limitación temporal, el uso del suelo de 1992 se utilizó para los años anteriores -1981 a 1992-. Las 22 clases fueron resumidas en siete -humedales, urbano, tierra desnuda, cobertura con zacate <50%, pastizales, boques y cultivos-, las equivalencias se basaron en las propuestas por *Land Cover Climate Change Partnership* (Land Cover Climate Change Initiative partnership, 2017) para el Grupo Intergubernamental de Expertos sobre el cambio climático para detección de cambio del uso del suelo.

Mapa 5
Cobertura del suelo. 2022



Fuente: Copernicus Climate Change Service, 2019.

Regiones climáticas

Las regiones climáticas de Costa Rica presentadas por Solano y Villalobos dividen el territorio nacional basado en características pluviométricas y térmicas. A continuación, se presenta una descripción de cada una de las regiones.

Región Pacífico Norte: Comprende la provincia de Guanacaste y los cantones de Esparza y Montes de Oro de la provincia de Puntarenas; y los cantones de Orotina y San Mateo de la provincia de Alajuela. Esta región pertenece al régimen de precipitación del Pacífico, conocido por la presencia de un período seco y otro lluvioso bien definidos (Solano & Villalobos).

Región Pacífico Central: Situada en la parte central de la Vertiente del Pacífico de Costa Rica, se extiende desde Playa Herradura o Jacó, hasta Dominical, Cerros de Herradura, Cerro Turrubares, Cerro Cangreja, partes bajas (pie de monte) de la Fila Costeña; comprende el poblado de Tinamaste, todo el Valle de Parrita, Quepos y Manuel Antonio. Esta región, al igual que el resto de las regiones del Pacífico nacional, se caracterizan por presentar el régimen de precipitación del Pacífico. Su posición geográfica al sureste, con la protección al norte por la cadena montañosa de la Fila Costera o Brunqueña, impiden la incursión de vientos alisios del noreste, estableciéndose una modificación de este régimen del Pacífico, presentándose una

caracterización propia de la región, como es, un clima tropical con estación seca con un período lluvioso muy severo y largo, y un período seco corto y moderado (Solano & Villalobos).

Región Pacífico Sur: Es una región extensa que se ubica al sureste del Pacífico Central, se extiende desde Punta Uvita, San Isidro del General, estribaciones de la Cordillera de Talamanca del lado del Pacífico, Cerro Darí, hasta Cerro Echandí, límite fronterizo con Panamá, hasta Punta Burica. Comprende todo el Valle del General, la Península de Osa, Valle de Coto Colorado, Valle de Coto Brus, Golfito. Esta región se ubica en la parte más sureste del Pacífico de Costa Rica, los contraste geográficos, entre ellos, los extensos valles, la barrera montañosa de la Cordillera de Talamanca al norte, como su influencia oceánica, generan en esta región un régimen de lluvias sumamente contrastado con relación al resto de las regiones de la vertiente, esta heterogeneidad se percibe en un clima en donde el período seco es muy favorable y corto y el lluvioso intenso, además, aparecen áreas pequeñas con clima tropical húmedo y lluvias todo el año (Solano & Villalobos).

Región Central -Montañosa del Sur-: Situada al sur del Valle Intermontano Central, o al sur de los Cerros del Tablazo, Candelaria y Puriscal. Al Norte del Pacífico Central. En esta pequeña región montañosa sobresalen los valles de altura, es una región intermedia entre el Valle Intermontano Central y la región del Pacífico Central, sus características, como las de cada región, son muy propias, donde su relieve montañoso y alturas medias de 800 a 1000 msnm la hacen ser un punto climático intermedio entre las lluvias moderadas del Valle Intermontano y las lluvias torrenciales del Pacífico Central. El clima es modificado por factores citados anteriormente para otras regiones, así se registran temperaturas cálidas en las partes bajas y frías en las partes altas (Solano & Villalobos).

Región Norte: Limita al Norte con la Cordillera Volcánica Central, al Oeste el límite de esta región lo forman la Cordillera de Guanacaste y la Cordillera de Tilarán. El Río Chirripó forma el límite convencional entre la Región Norte y la Región Norte del Este. Esta región pertenece al régimen de precipitación del Caribe, al cual se le identifica como lluvioso todo el año, con una disminución de las lluvias en los meses de febrero, marzo y octubre. Esta región presenta un clima tropical húmedo (típico ecuatorial desplazado), el cual presenta dos rasgos esenciales

que son: 1) ningún mes del año tiene temperaturas inferiores a los 22C; 2) no presenta promedios pluviométricos mensuales superiores a los 75 mm. Es una región de contrastes en la lluvia, ya que en ella interactúan tanto elementos climáticos como factores geográficos propios de la región, debido a su relieve montañoso, llanuras extensas y la influencia al nornoroeste, del lago de Nicaragua, estableciéndose una serie de pequeñas subregiones climáticas (Solano & Villalobos).

Región Atlántica: Comprende toda la provincia de Limón, y la parte oriental de la provincia de Cartago (de Turrialba hacia el Este). En esta región de clima tropical húmedo, la lluvia es abundante, siendo más acentuada en las partes montañosas donde llueve todo el año. La región presenta una serie de subregiones pequeñas como producto de la misma diversidad de factores de la Región Norte, así se encuentran áreas con clima lluvioso, principalmente en las llanuras y en alturas inferiores a los 60 metros sobre el nivel del mar (msnm), áreas con clima de las faldas de la Cordillera Volcánica del Norte del lado Caribe, en alturas de 600 a 1600 msnm (Solano & Villalobos).

Mapa 6
Regiones climáticas de Costa Rica



Fuente: SIG – Departamento de Desarrollo, IMN, 2022.

Período de estudio

El período de estudio será del año 1981 al 2021. Este intervalo fue elegido principalmente por ser el período en donde la mayoría de los datos cuentan con información disponible. Años anteriores a 1981 no cuentan con información pública de alta definición de precipitación.

Métricas

El objetivo de esta investigación es predecir los caudales promedio mensuales desde 1980 al 2021. En problemas de regresión supervisado varias métricas son comúnmente utilizadas tales como el error cuadrático medio, error absoluto medio, porcentaje del error absoluto medio, entre otras. En el ámbito de la hidrología estas métricas son también utilizadas, sin embargo, una de las más métricas con mayor aceptación es la Eficiencia Kling Gupta -KGE- (Gupta, Kling, Yilmaz, & Martinez, 2009) y su versión modificada -KGE'- (Kling, Fuchs, & Paulin, 2012).

Esta métrica se computa según la siguiente ecuación:

$$KGE = 1 - ED$$

$$ED = \sqrt{(r - 1)^2 + (\gamma - 1)^2 + (\beta - 1)^2}$$

$$\gamma = \frac{\sigma_s / \mu_s}{\sigma_o / \mu_o}$$

Donde:

r coeficiente de correlación de Pearson entre el flujo observado y simulado.

β la división entre la media del flujo simulado y el observado. Razón de sesgo.

γ razón de variabilidad.

σ_s desviación estándar del flujo simulado u observado (σ_s)

μ_s media del flujo simulado u observado (μ_o)

La métrica KGE' y sus tres componentes -correlación, la razón de sesgo y la razón de variabilidad son adimensionales y tienen un valor óptimo de 1. Por su definición el KGE' siempre será como máximo el menor valor de sus componentes. Esto garantiza que valores altos de KGE' reflejen una buena correspondencia entre los flujos simulados y observados.

El KGE' puede tener valores entre $-\infty$ y 1. Siendo 1 un ajuste perfecto. Valores de KGE superiores a -0,41 denota que el modelo muestra una capacidad predictiva superior que la

media de la serie temporal. En otras palabras, si $-0,41 < KGE \leq 1$ el modelo tiene un rendimiento “razonable” en el entendido que supera la estimación del flujo medio (Knoben, Freer, & Woods, 2019).

Adicional a la métrica KGE se reportará el Error Nash-Sutcliffe -NSE- ya que es una de las métricas más comunes en hidrología (Kling, Fuchs, & Paulin, 2012) puesto que sirven como punto de comparación con otros estudios realizados a nivel nacional e internacional.

Tipos de errores

A la hora de reportar métricas se debe tener en cuenta el proceso de modelación en el que se encuentra. A nivel, general se espera que los errores aumenten cuando las observaciones evaluadas se alejan de la distribución de los datos con los que el modelo fue entrenado.

Es común que los errores se comporten de esta manera:

$$\text{error calibración} < \text{error validación} < \text{error prueba}$$

Donde:

Error de calibración

Llamado también error de entrenamiento. Es el error que se obtiene durante el proceso de calibración o entrenamiento del modelo. Tiende a ser el error más pequeño y no se debe utilizar para reportar errores para el usuario puesto que probablemente están sobre ajustados al conjunto de entrenamiento y no representan necesariamente la capacidad de generalización del modelo.

Error de validación

Una vez que el modelo ha sido calibrado, se aplica sobre datos que no han sido utilizados durante la calibración. Ese error puede ser calculado en una cuenca que se usó para la calibración, pero en un período que el modelo no ha visto o en una cuenca que del todo no ha sido utilizada durante el entrenamiento. Lo usual es que el error sea menor en cuencas en donde el modelo ha sido calibrado que en cuencas no utilizadas durante ese proceso.

Error de prueba

Una vez que el modelo fue entrenado y validado, se calcula el error de prueba en un conjunto de datos no vistos por el modelo en ninguna fase previa. Este error representa el error que un usuario esperará al utilizar el modelo en una cuenca o espacio temporal totalmente nuevo para el modelo.

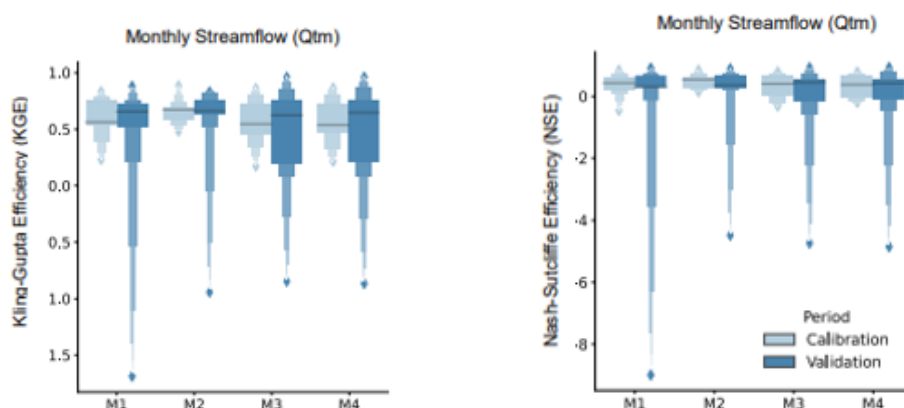
Línea base nacional

En Costa Rica no existe un reporte generalizado de las métricas alcanzadas de varios modelos predictivos de caudales mensuales en las estaciones hidrológicas del ICE y aunque existieran la información base y períodos de aplicación son distintos y por ende dificulta la comparación de los modelos.

Arciniega et al (2022) reportaron valores de calibración medios de KGE de 0,67 y una mediana NSE de 0,43 para 13 cuencas en Costa Rica -modelo M2-, las cuencas fueron calibradas en el período 1991-1999 y validadas en 2000-2023 y utilizaron CHIRPS como base para la precipitación. A pesar de que los períodos de análisis son diferentes, este estudio es un útil marco de referencia puesto que comparte 10 cuencas con la investigación y usa la misma fuente de precipitación.

Gráfico 1

Métricas obtenidas para diferentes configuraciones de modelos en calibración, 1991-1999 y validación. 2000-2003



Nota: Los valores corresponden a la media \pm desviación estándar utilizando todas las cuencas del estudio.
Fuente: Elaboración propia con datos de (Arciniega-Esparza, Birkel, Chavarría-Palma, Arheimer, & Agustín, 2022).

Recientemente, el Instituto Costarricense de Acueductos y Alcantarillados (AyA) realizó estudios de balance hídrico (Kaune, 2021) en las cuencas definidas por la estación El Brujo,

Electronia¹ y El Rey. Estos trabajos tienen la peculiaridad de que tienen un horizonte de análisis prácticamente igual al de esta investigación y que comparten la información base de precipitación, temperatura y caudales. Por este motivo son un buen candidato para comparar los resultados obtenidos por este modelo. Kaune (2022) utilizó el modelo hidrológico distribuido *Spatial Process in Hydrology* -SPHY- específicamente calibrado para esas estaciones.

Los resultados reportados en la calibración son:

Cuadro 1

Resultados del cálculo de caudales mensuales reportado por Kaune, 2021 durante calibración

Cuenca	Período	NSE
El Rey	1990-1999	0,74
El Brujo	1990-1994	0,68

Fuente: Elaboración propia.

Se debe mencionar que ambos estudios calculan sus métricas en la misma cuenca donde se calibró el modelo.

Línea base internacional

Conocer una línea base adecuada para el proyecto es una tarea difícil especialmente por la variabilidad en la información hidrológica, las diferencias en las cuencas y la incertidumbre asociada a la información base. En general, el modelo debería ser evaluado en relación con lo que es posible con la información disponible. A pesar estas limitaciones, se utilizarán como referencia los criterios emitidos por Moriasi y otros (2007) para la evaluación de modelos de descargas.

¹ Esta estación no fue utilizada en esta investigación por no cumplir con una cantidad suficiente de datos.

Cuadro 2
Categorización de métricas de modelos de caudales mensuales

Desempeño	NSE
Muy bueno	$0,75 < NSE \leq 1.0$
Bueno	$0.65 \leq NSE \leq 0,75$
Satisfactorio	$0,50 < NSE \leq 0,65$
Insatisfactorio	$NSE \leq 0,5$

Fuente: Elaboración propia con datos de (Moriasi, y otros, 2007).

Métricas esperadas

El modelo se considerará exitoso si cumple cuatro condiciones simultáneamente:

- En todos los casos la mediana de la KGE de validación fuera de cuenca es superior a -0,41.
- La mediana del error de prueba KGE es superior a -0,41 y los valores individuales se encuentran dentro de 1,5 desviaciones estándar de la distribución de los errores de validación.
- La mediana de los errores de validación fuera de cuenca supera los criterios de satisfacción establecidos por Moriasi y otros (2007).
- El modelo tiene un valor de NSE no menor a los valores obtenidos por Kaune (2021) y los valores de NSE y KGE son similares o superiores a los obtenidos por Arciniega y otros (2022).

Cumplir con los criterios implicaría que el modelo en todos los casos cuenta con una capacidad predictiva, segundo que el modelo generaliza bien y es congruente con la validación en cuencas nunca antes vistas por el modelo y por último que el modelo tiene métricas iguales o superiores a estudios específicos realizados con modelos físicos y datos base similares.

Análisis exploratorio de los datos

Caudales

En esta sección se analizará la información hidrológica desde el punto de disponibilidad de datos y de valores atípicos.

La disponibilidad de información es relevante para el estudio desde el punto de vista de entrenamiento y validación. Cuanto mayor sea el período de información disponible para entrenamiento mejor será la generalización del modelo, a su vez si el modelo es validado en estaciones con pocas observaciones el resultado no será representativo del período de estudio.

Para evaluar la disponibilidad de información se calculó el porcentaje de datos faltantes como

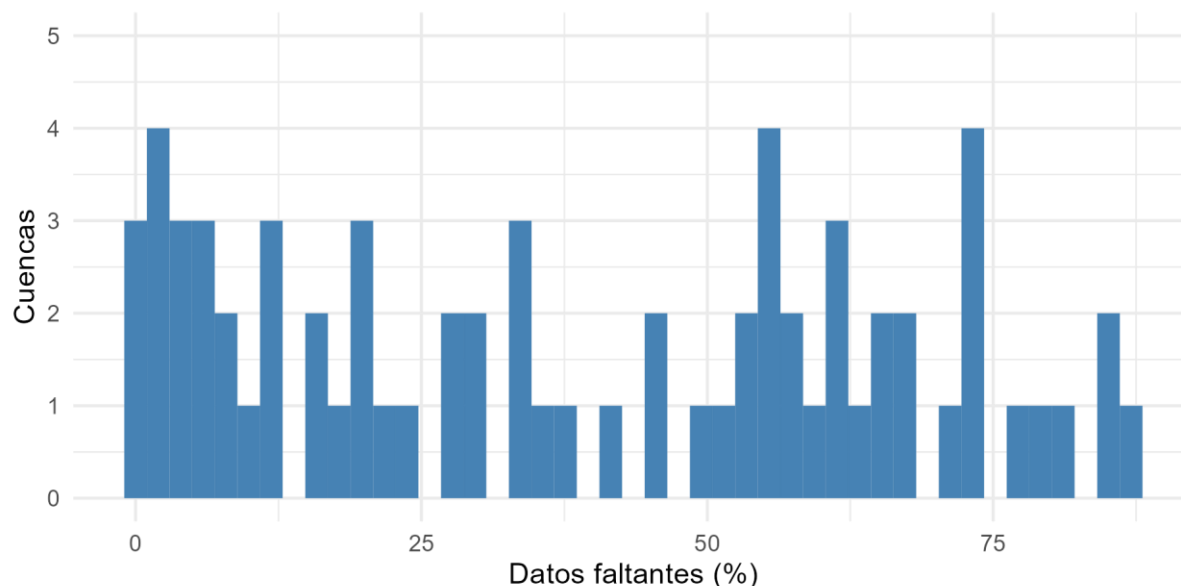
$$\text{Datos faltantes} = 1 - \frac{\text{Cantidad de datos de estación}}{\text{Cantidad máxima de datos posible}}$$

Donde:

Cantidad máxima de datos posible = 495.

Al aplicar la fórmula a cada estación se obtuvo el gráfico 2.

Gráfico 2
Porcentaje de datos faltantes por estación



Fuente: Elaboración propia.

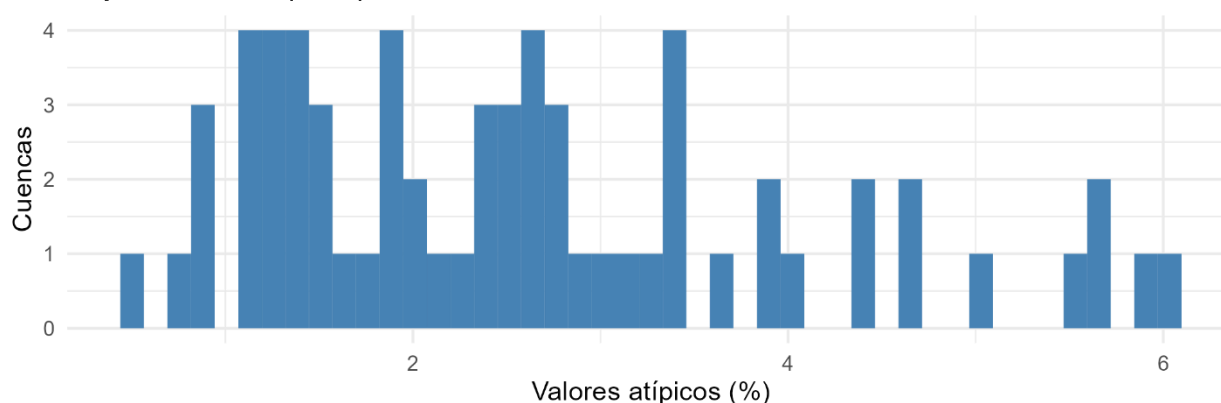
No se tomarán en cuenta cuencas que tengan más de un 80% de datos faltantes, es decir solo se utilizarán estaciones con más de 100 mediciones a lo largo del período de estudio. Las cuencas que se retirarán son las de Chirripó, Cola Embalse Reventazón, Agua Gata y Puente la Virgen.

Valores atípicos son aquellas mediciones que se alejan de manera significativa de la distribución de la variable en estudio. Los valores atípicos pueden deberse a diversos factores como una mala lectura del fenómeno, mala transcripción, eventos extremos, liberación de caudal de una represa, represamiento, entre otros.

Los valores fueron identificados utilizando la metodología de rango inter-cuartil la cual es un procedimiento univariado robusto a valores atípicos en donde clasifica una medición como atípica cuando esta se encuentra más lejos que la mediana de la muestra más 2,5 veces el rango inter-cuartil.

Este procedimiento se aplicó a cada mes de cada cuenca y finalmente se contabilizaron los valores atípicos como porcentaje del total de observaciones.

Gráfico 3
Porcentaje de valores atípicos por cuenca



Fuente: Elaboración propia.

En total se contabilizaron y removieron 509 valores atípicos que corresponde a un 2,45% de los datos.

Características de las cuencas

En esta sección se analizarán las 65 cuencas asociadas con las estaciones seleccionadas. Las estaciones y su correspondiente cuenca se muestran en la siguiente lista:

- | | | |
|-----------------------|----------------------|------------------|
| 1. Angostura | 8. Cabagra | 13. Desagüe |
| 2. Bajos de Chilamate | 9. Caracucho | 14. División |
| 3. Balsa | 10. Casa Máquinas El | 15. Dos Montañas |
| 4. Barbilla | Brujo 1 | 16. Dota |
| 5. Belén | 11. Cola Embalse | 17. El Brujo |
| 6. Bijagual | Angostura | 18. El Humo |
| 7. Bratsi | 12. Coyolar | 19. El Rey |

- | | |
|-----------------------|---------------------|
| 20. El Salado | 50. Remolino |
| 21. Guapinol | 51. Río Segundo |
| 22. Guardia | 52. Rivas |
| 23. Javillos | 53. Santa Lucía |
| 24. La Cuesta | 54. Santo Domingo |
| 25. La Isla | 55. Savegre |
| 26. Las Juntas | 56. Sitio de Presa |
| 27. Limonal Viejo | Guayabo |
| 28. Londres | 57. Sitio de Presa |
| 29. Los Llanos | Savegre |
| 30. Nagatac | 58. Sixaola |
| 31. Nuestro Amo | 59. Tabacales |
| 32. Oriente | 60. Tapantí Arriba |
| 33. Pacuare | 61. Terrón Colorado |
| 34. Palmar | 62. Toro |
| 35. Palomo | 63. Turrialba |
| 36. Pascua | 64. Veracruz |
| 37. Pejibaye Caribe | 65. Verge |
| 38. Pejibaye Pacífico | |
| 39. Peñas Blancas | |
| 40. Peralta | |
| 41. PH Ceibo | |
| 42. Piedras Blancas | |
| 43. Pocosol | |
| 44. Providencia | |
| 45. Puente de Hamaca | |
| 46. Puente Negro | |
| 47. Puerto Viejo | |
| 48. Rancho Horizontes | |
| 49. Rancho Rey | |

Ubicación de las estaciones

Las estaciones se distribuyen en el territorio nacional de la siguiente manera:

Mapa 7

Ubicación de estaciones hidrológicas



Fuente: Elaboración propia.

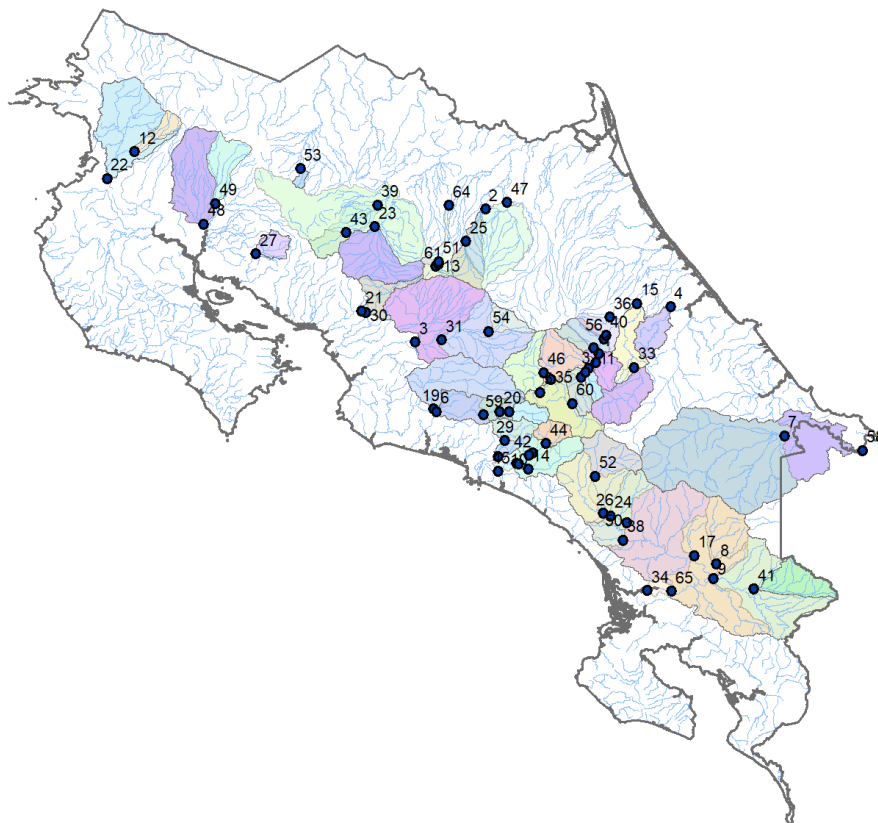
En el mapa 7 se puede ver como las estaciones no se encuentran distribuidas homogéneamente en el país. Se aprecian vacíos de información en la Península de Nicoya, Península de Osa, Valle de Coto Colorado, Llanuras del Norte y en la Cordillera de Talamanca. También hay concentración de estaciones en sitios de importancia hidroeléctrica como la cuenca del Río Reventazón, Río Savegre, Pirrís, Río Grande de Térraba, Río Toro, Sarapiquí y Naranjo.

Delimitación de cuencas

Se delimitaron las cuencas hidrográficas asociadas con las estaciones hidrológicas. Estas cuencas se pueden ver en el siguiente mapa:

Mapa 8

Cuencas de estaciones hidrológicas Polígonos de colores indican las cuencas delimitadas por las estaciones (puntos azules)



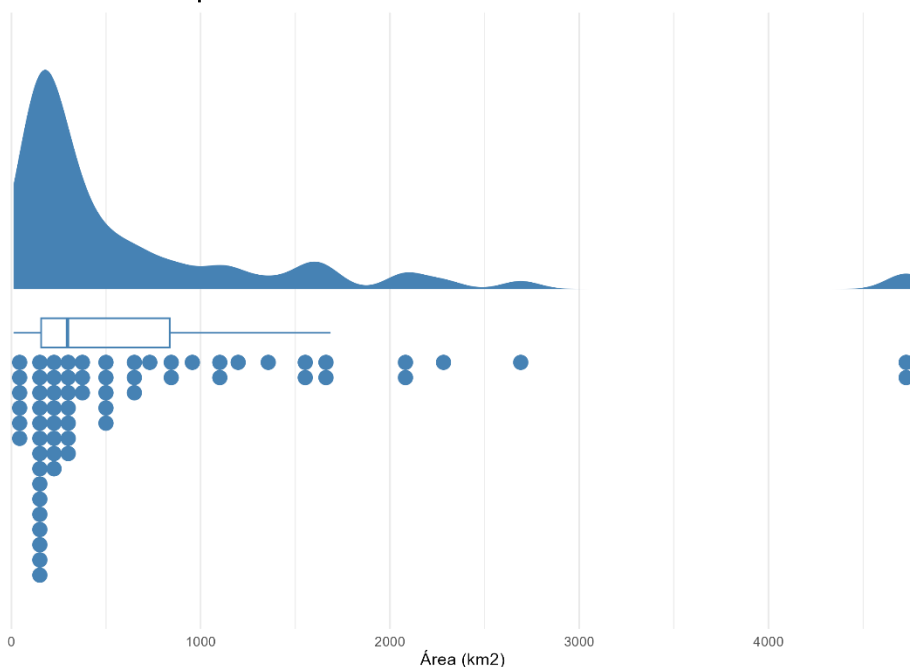
Fuente: Elaboración propia.

Como se puede apreciar en el mapa 8, las cuencas disponibles tienen características muy variadas que serán analizadas a continuación:

Área y perímetro

La distribución del área de las cuencas tiene un sesgo a la izquierda. La gran mayoría de las cuencas se encuentra por debajo de los 1000 km². Resalta la cuenca de Vergel y Palmar con áreas superiores a los 4.500 km² -gráfico 4-.

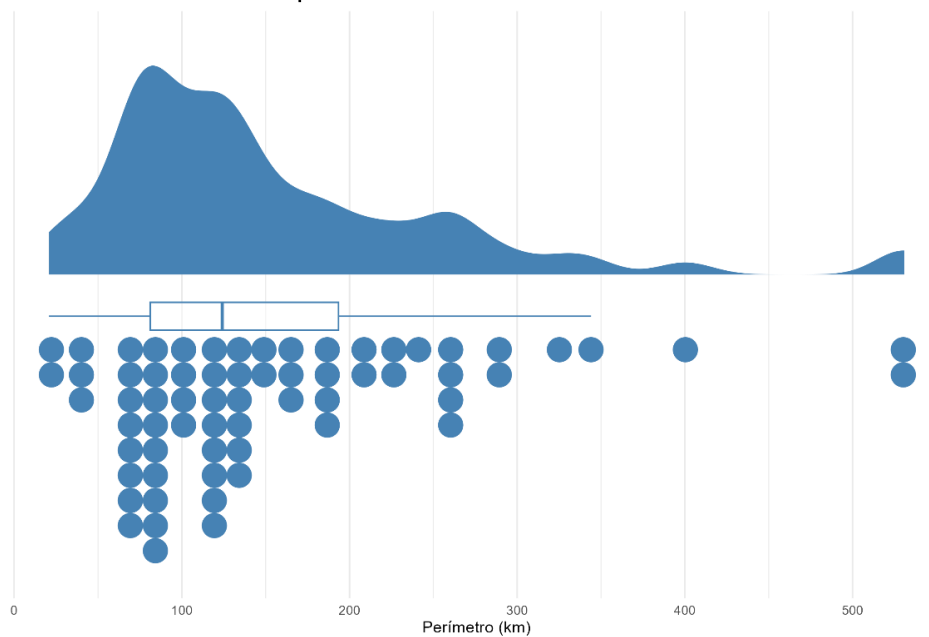
Gráfico 4
Área de las cuencas por analizar



Fuente: Elaboración propia.

Como es de esperar el perímetro tiene una distribución similar, pero con un sesgo no tan marcado -gráfico 5 -.

Gráfico 5
Perímetro de las cuencas por analizar

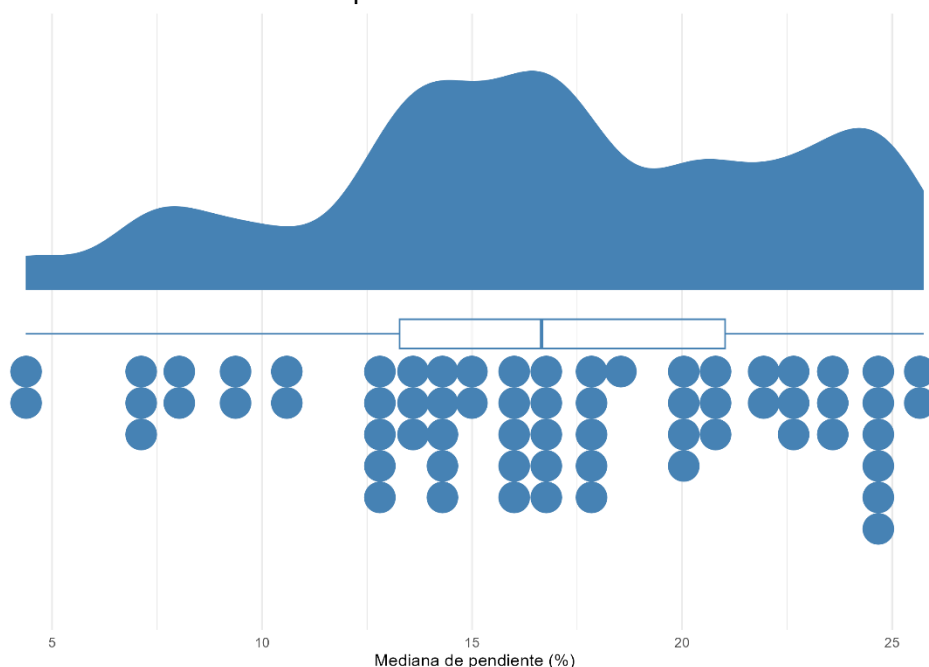


Fuente: Elaboración propia.

Pendiente y elevación media

La pendiente y elevación de una cuenca están fuertemente correlacionados puesto que cuencas elevadas tienen una pendiente mayor por encontrarse más cerca de las áreas montañosas del país. En las cuencas por analizar la mediana de la pendiente tiene una distribución con un pico principal entre 12,5% y 17,5%, las cuencas por encima de la pendiente media de 25% son Rivas y Piedras Blancas. Por el contrario, las cuencas con menor pendiente media son Guardia y Rancho Horizontes con pendientes cercanas al 4%.

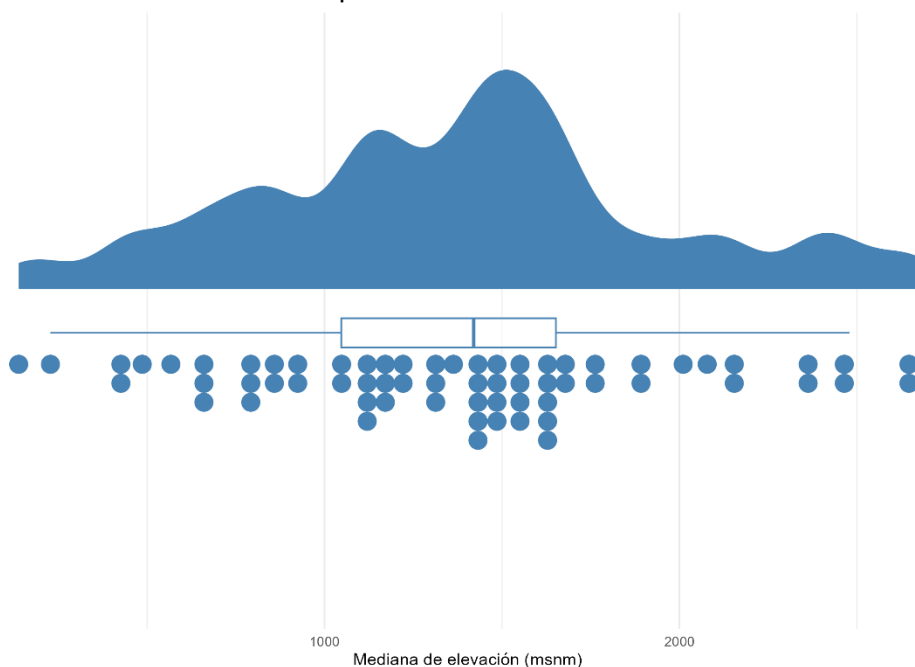
Gráfico 6
Pendiente media de las cuencas por analizar



Fuente: Elaboración propia.

La elevación media va desde los 138 msnm -Rancho Horizontes- a 2.667 msnm -Tapantí Arriba-. La distribución tiene un pico entre 1.250 msnm y 1.750 msnm.

Gráfico 7
Elevación media de las cuencas por analizar

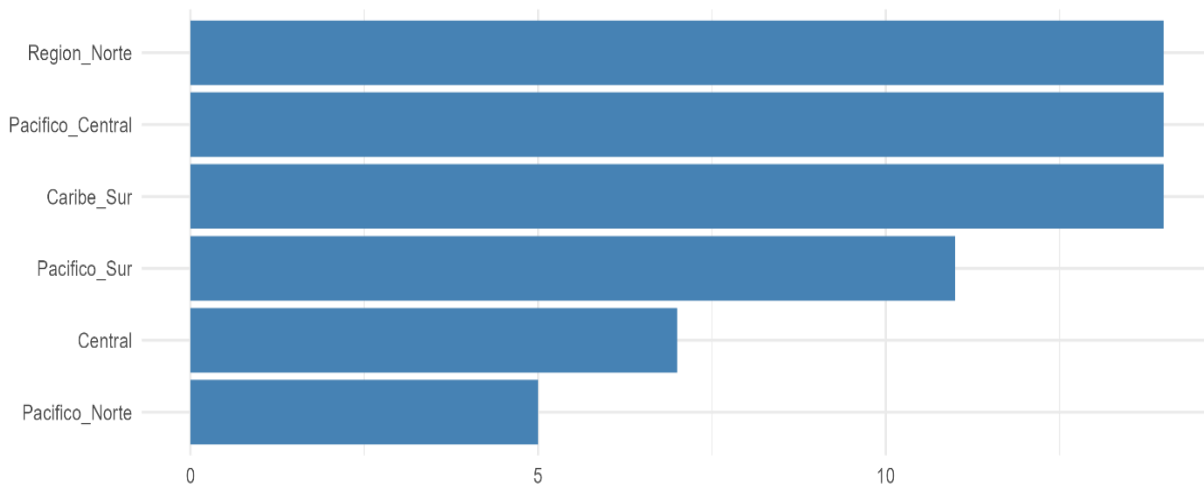


Fuente: Elaboración propia.

Región climática

A cada una de las cuencas, se les asignó una región climática. La asignación fue realizada con base en el centroide de la cuenca. Luego de aplicar este procedimiento, la distribución de las cuencas dentro de las regiones se puede ver en gráfico 8.

Gráfico 8
Cantidad de cuencas por región



Fuente: Elaboración propia.

El gráfico anterior muestra como la distribución de cuencas no es homogénea y que la región Caribe Norte no cuenta con cuencas asociadas. Las regiones con mayor cantidad de cuencas son la Región Norte, Pacífico Central y Caribe Sur con 14 seguidas por Pacífico Sur con 11, Central con 7 y Pacífico Norte con 5.

Ingeniería de variables

Rezagos

El caudal mensual de los ríos depende tanto de condiciones hidrológicas pasadas como a las correspondientes al mes presente. Para que el modelo tenga información relacionada a eventos anteriores se aplicarán diversos rezagos a la precipitación mensual.

- Precipitación en T-1, T-2, T-3 y T-4, donde T es el mes actual.
- Precipitación media de T-5 a T-8 y de T-8 a T-12, donde T es el mes actual.
- Temperatura en T-1 donde T es el mes actual.

Suelos

La textura de los suelos está asociada a diferentes características físicas involucradas en el ciclo hidrológico. A la textura del suelo se le asignaron las siguientes características recomendadas por Clapp y Hornberger (1978):

Cuadro 3

Propiedades del suelo según su textura

Textura	Contenido volumétrico de agua (cm ³ /cm ³)	Conductividad hidráulica saturada (cm/min)
Arcilloso	0,482	0,0077
Arcilloso limoso	0,492	0,0062
Arcilloso arenoso	0,426	0,013
Franco arcilloso	0,476	0,0147
Franco arcilloso limoso	0,477	0,0102
Franco arcilloso arenoso	0,42	0,0378

	Contenido	Conductividad
Textura	volumétrico de agua (cm ³ /cm ³)	hidráulica saturada (cm/min)
Franco	0,451	0,0417
Franco limoso	0,485	0,0432
Franco arenoso	0,435	0,208
Limoso	0,485	0,0432
Arenoso	0,395	1,056

Fuente: Elaboración propia con datos de (Clapp & Hornberger, 1978).

Fracción que infiltra por efecto de la cobertura del suelo (Kv)

La cobertura del suelo incide directamente en la respuesta de la cuenca frente a un evento de precipitación. Superficies impermeables hacen que el efecto de una precipitación se transfiera más rápidamente en escorrentía y viceversa.

Tomando como base a Schosinsky y Losilla (2000), la fracción que infiltra por efecto de la cobertura del suelo (Kv) se define en el siguiente cuadro:

Cuadro 4

Fracción que infiltra por efecto de la cobertura del suelo

Cobertura	Kv
Cobertura con zacate menor a 50%	0,09
Cultivos	0,10
Bosques	0,20
Pastizales	0,18
Tierra desnuda ^{a/}	0,07
Urbano ^{a/}	0,05

a/ Añadidas por el autor.

Fuente: Elaboración propia con datos de Schosinsky & Losilla, 2000.

Fracción que infiltra por efecto de la pendiente (Kp)

Tomando como base a Schosinsky y Losilla (2000), la fracción que infiltra por efecto de la pendiente (Kp) se define en el siguiente cuadro:

Cuadro 5

Fracción que infiltra por efecto de la pendiente

Pendiente	Rango	Rango	Kp
	inferior	superior	
Extremo baja	0	0,06	0,3
Muy baja	0,06	0,3	0,25
Baja	0,3	0,4	0,2
Media	0,4	1	0,175
Alta	1	2	0,15
Muy alta	2	7	0,1
Extremo alto	7	100	0,06

Fuente: Schosinsky & Losilla, 2000.

Coeficiente de infiltración

El coeficiente de infiltración es un valor que representa la fracción de precipitación mensual que infiltra hacia el suelo. En esta investigación se calculará este valor de acuerdo con las recomendaciones de Schosinsky (Schosinsky, 2006) en donde el coeficiente de infiltración (C_i) depende de tres componentes:

$$C_i = k_p + k_v + k_{fc}$$

$$\text{si } C_i > 1 : C_i = 1$$

Donde:

C_i = Coeficiente de infiltración [adimensional].

K_p = Fracción que infiltra por efecto de pendiente [adimensional]

K_v = Fracción que infiltra por efecto de cobertura vegetal [adimensional]

K_{fc} = Fracción que infiltra por textura del suelo.

Donde K_{fc} se calcula de la siguiente manera:

Fracción que infiltra por textura del suelo (Kfc)

(Schosinsky & Losilla, 2000)

$$\begin{aligned} \text{si } 16 < fc < 1568: Kfc &= 0,267 \ln(fc) - 0,000154fc - 0,723 \\ \text{si } fc < 16: Kfc &= \frac{0,0148fc}{16} \\ \text{si } fc > 1568: Kfc &= 1 \end{aligned}$$

Donde:

Kfc [adimensional] = Coeficiente de infiltración (fracción que infiltra por textura del suelo).

fc [mm/día] = Permeabilidad del suelo saturado, en los primeros 30 centímetros de profundidad.

El valor de Kfc se calculó a nivel nacional basados en la textura del suelo y los valores recomendados por Hornberger y Clapp (Clapp & Hornberger, 1978).

Resumen de variables

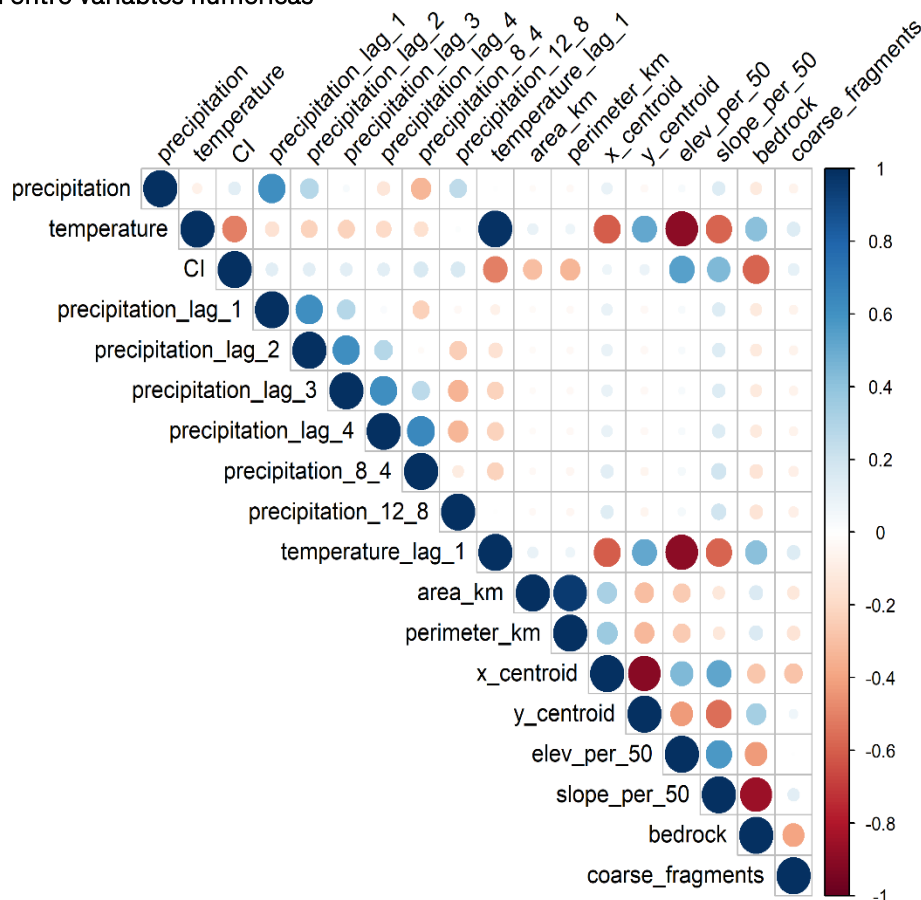
Luego de las diversas transformaciones se utilizarán un total de 20 variables predictoras:

1. Precipitación_t (mm)
2. Precipitación_{t-1} (mm)
3. Precipitación_{t-2} (mm)
4. Precipitación_{t-3} (mm)
5. Precipitación_{t-4} (mm)
6. Precipitación_{t-4 a t-8} (mm)
7. Precipitación_{t-8 a t-12} (mm)
8. Temperatura_t (C°)
9. Temperatura_{t-1} (C°)
10. Coeficiente de infiltración (CI)
11. Área (km)
12. Perímetro (km)
13. Centroide en y (m)
14. Centroide en x (m)
15. Elevación percentil 50 (msnm)
16. Pendiente percentil 50 (%)
17. Profundidad a lecho rocoso (m)
18. Fragmentos rocosos (%)
19. Mes
20. Región climática

En general todas las variables son numéricas a excepción de la región climática y mes. Estas variables fueron desagregadas en variables binarias para cada una de sus categorías.

Las variables no tienen valores faltantes y todas fueron agregadas a nivel de cuenca. Por su naturaleza, las variables están correlacionadas entre sí, en especial aquellas que tienen una relación en el tiempo como la temperatura y la precipitación. Hay otras variables que están relacionadas como la altura y la temperatura, condiciones geológicas y parámetros físicos de la cuenca. Una visualización de la correlación entre variables se puede ver en el siguiente gráfico.

Gráfico 9
Correlación entre variables numéricas



Fuente: Elaboración propia.

Conjuntos de entrenamiento, validación y prueba

Antes de iniciar el proceso de entrenamiento se separó aleatoriamente una cuenca por cada región y las mediciones de las cuencas de la estación El Rey y El Brujo que fueron utilizadas para conocer las métricas en los estudios realizados por Kaune, 2022.

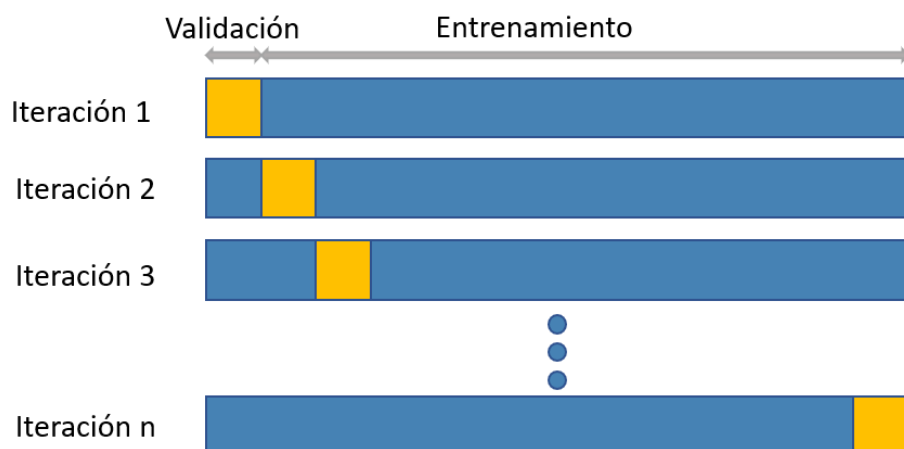
Cuadro 6
Cuenkas seleccionadas para el conjunto de prueba

Cuenca	Período	Región	Observaciones de prueba	Total de observaciones	Porcentaje de observaciones de prueba
Javillos	1981-2021	Región Norte	486	486	100
Dos Montañas	1981-2000	Caribe Sur	228	228	100
El Salado	2002-2021	Pacífico Central	219	219	100
Rivas	1981-2021	Pacífico Sur	386	386	100
Limal Viejo	1981-1996	Pacífico Norte	180	180	100
Puente Negro	1998-2021	Central	267	267	100
El Brujo	1990-2000	Pacífico Sur	60	528	11,4
El Rey	1990-1995	Pacífico Central	117	415	28,2

Fuente: Elaboración propia.

El restante de las observaciones será utilizado para entrenamiento. La metodología de validación a utilizar es conocida como *leave-one-out cross validation* en donde el algoritmo se ajusta a todas las cuencas menos una y se valida con la cuenca que se dejó de lado.

Figura 2
Validación cruzada leave-one-out



Fuente: Elaboración propia.

Este procedimiento permite obtener los errores de validación para cada cuenca de manera independiente, lo cual será útil para analizar las métricas en todas las cuencas.

Algoritmo predictivo y entrenamiento

En esta investigación se utilizó el algoritmo *Extreme Gradient Boosting* (XGBoost) el cual es un algoritmo de aprendizaje automático supervisado que utiliza árboles de decisión y *boosting* - combinación de modelos simples en uno más complejo y de mejor rendimiento-para mejorar gradualmente la precisión del modelo. Es ampliamente utilizado en la comunidad de aprendizaje automático debido a su funcionalidad para manejar problemas de regresión y clasificación, y su capacidad para regularizar y optimizar el rendimiento del modelo.

En el proceso de entrenamiento se optimizaron los hiperparámetros del modelo mediante optimización bayesiana cuya función objetivo era el NSE. Los parámetros que se optimizaron fueron:

- `mtry`: El número de predictores que serán aleatoriamente seleccionados en cada partición cuando se crea un modelo de árbol.
- `trees`: El número de árboles contenidos en el ensamble.

- **min_n**: La cantidad mínima de observaciones en un nodo que son requeridos para dividir el nodo una vez más.
- **tree_depth**: Profundidad máxima del árbol.
- **learn_rate**: La tasa a la que el algoritmo de *boosting* se adapta tras cada iteración.
- **loss_reduction**: La reducción en la función de pérdida para dividir el nodo una vez más.
- **sample_size**: Cantidad de datos expuestos a la rutina de entrenamiento.

Resultados del modelo

Una vez que se probaron más de un centenar de combinaciones de hiperparámetros en el conjunto de validación, se eligió un modelo con la mejor métrica KGE en las cuencas seleccionadas para la validación. El mejor modelo obtuvo una mediana de KGE de 0,7 en el conjunto de validación y un 0,67 en el conjunto de prueba. En cuanto a los valores NSE se obtuvo una mediana de 0,68 en validación y un 0,6 en los conjuntos de prueba.

Respecto a los errores de validación intra-cuenca se obtuvo un NSE de 0,85 y 0,79 para la estación El Brujo y el Rey, respectivamente. En las siguientes secciones se detallan los resultados del modelo.

Errores de validación fuera de cuenca

Todas las métricas KGE son mayores a 0,25 indicando que el modelo tiene poder predictivo muy por encima de un modelo cuyas predicciones son los valores medios del flujo.

Cuadro 7
Resultados de validación fuera de cuenca KGE

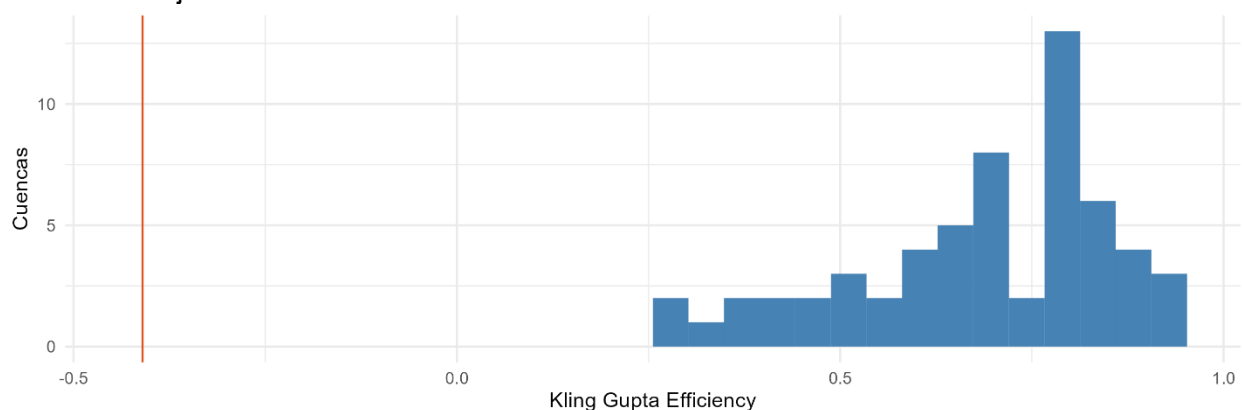
Rango	Cuencas	Porcentaje
KGE < 0,26	0	0
0,26 <KGE< 0,5	9	15,2
0,5 <KGE< 0,75	23	39,0
0,75 <KGE< 0,94	27	45,8

Fuente: Elaboración propia.

La distribución de los errores se puede apreciar en gráfico 10.

Gráfico 10

KGE en los conjuntos de validación^{a/}



a/ Línea naranja indica. KGE= -0.41. Valores superiores a la línea naranja se consideran satisfactorios en relación con el flujo medio.

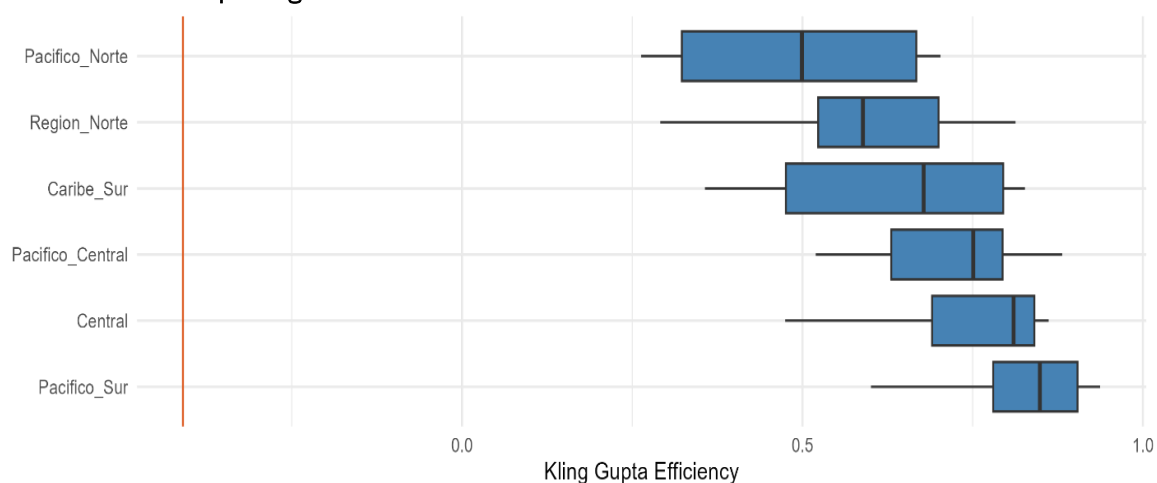
Fuente: Elaboración propia.

Según el Cuadro 7 y gráfico 10 los valores de KGE tienen una distribución por encima de 0,263 y un máximo de 0,94. Dentro del rango de -0,41 a 1, la distribución está sesgada a la derecha en donde se aprecia que la gran mayoría de las cuencas (85%) están por encima de 0,50.

A nivel regional los errores se distribuyen de la siguiente manera:

Gráfico 11

Distribución de KGE por región^{a/}



a/ Línea naranja indica KGE= -0.41. Valores superiores a la línea naranja se consideran satisfactorios en relación con el flujo medio

Fuente: Elaboración propia.

El gráfico 11 muestra como la Pacífico Norte es la región con la mediana más baja de todas las regiones 0,5 y un rango inter-cuartil amplio, la Región Norte tiene una mediana de 0,59 pero con un rango inter-cuartil más reducido. Caribe Sur tiene una distribución amplia y sesgada a la derecha con una mediana de 0,68. Las regiones Pacífico Central, Central y Pacífico Sur son las regiones con las medianas más elevadas y con distribuciones prácticamente por encima de 0,5 en todas las cuencas.

Errores de prueba

Finalmente, el modelo se aplicó al conjunto de prueba para conocer la KGE en datos nunca vistos por el modelo. Los resultados se muestran en el siguiente cuadro.

Cuadro 8
Valores de NSE en el conjunto de prueba

Cuenca	Región	KGE	Z-score absoluto	Acorde con la distribución ²
Javillos	Región Norte	0,74	0,34	Sí
Dos Montañas	Caribe Sur	0,63	0,33	Sí
El Salado	Pacífico Central	0,75	0,376	Sí
Limal Viejo	Pacífico Norte	0,51	1,04	Sí
Rivas	Pacífico Sur	0,48	1,22	Sí
Puente Negro	Central	0,71	0,136	Sí

Fuente: Elaboración propia.

Los resultados en el conjunto de prueba son acordes con los del conjunto de validación, todos tienen un valor de z-score menor a 1,5. Además, en la mitad de los casos los resultados son mejores que la mediana del conjunto de validación lo que indica que no hay un sobreajuste evidente a los conjuntos de entrenamiento y validación.

² Se encuentra a menos de 1,5 veces la desviación estándar de la media.

Errores línea base nacional

Adicionalmente a las cuencas de prueba, se dejaron de lado un subconjunto de datos de las cuencas El Brujo y El Rey para obtener una métrica de prueba en cuencas que fueron utilizadas por el modelo. Los resultados de esta evaluación se muestran en el Cuadro 9.

Cuadro 9

Comparación de línea base nacional con el modelo

Cuenca	Región	KGE	NSE modelo (prueba intra- cuenca)	NSE línea base (calibración)	Supera la línea base
El Brujo	Pacífico Sur	0,81	0,85	0,68	Sí
El Rey	Pacífico Central	0,83	0,79	0,74	Sí

Fuente: Elaboración propia.

Con base en el cuadro anterior, los resultados del modelo utilizado en esta investigación son superiores a los del modelo de la línea base. Es interesante que esta diferencia se hace presente inclusive cuando se comparan errores de calibración -estudio de línea base- contra los errores de prueba intra-cuenca.

Estas métricas positivas tanto en NSE y KGE demuestra que el modelo de esta investigación podría ser utilizado para completar información faltante en cuencas que fueron utilizadas por el modelo de manera satisfactoria.

A pesar de que la comparación no debe tomarse rigurosamente por la diferencia en las cuencas utilizadas, información base y los períodos analizados, la investigación tiene valores de KGE similares y NSE superiores a los reportados por Archiega y otros (2022) y Kaune (2021).

Conclusiones

En relación con la métrica KGE y NSE el modelaje de los caudales promedio mensuales de 65 cuencas a nivel nacional se considera exitoso puesto que cumple con las métricas esperadas de la investigación:

- En todos los casos la KGE es superior a -0,41, de hecho, la mediana de la KGE en validación es muy superior con un valor 0,7.

- La KGE de prueba es de 0,67, superior a -0,41.
- Respecto al NSE, el modelo presenta resultados satisfactorios, buenos y muy buenos en el 71% de las cuencas analizadas con base las recomendaciones de Moriasi y otros (2007).
- Los valores de KGE de prueba muestran resultados acordes a las métricas de validación y por tanto se considera que, si el modelo es puesto en producción, los valores esperados coincidirán con los presentados en este estudio.
- Las métricas KGE y NSE del modelo son comparables o superiores a otros modelos realizados a nivel nacional en dos estudios recientes.

Interpretación del modelo

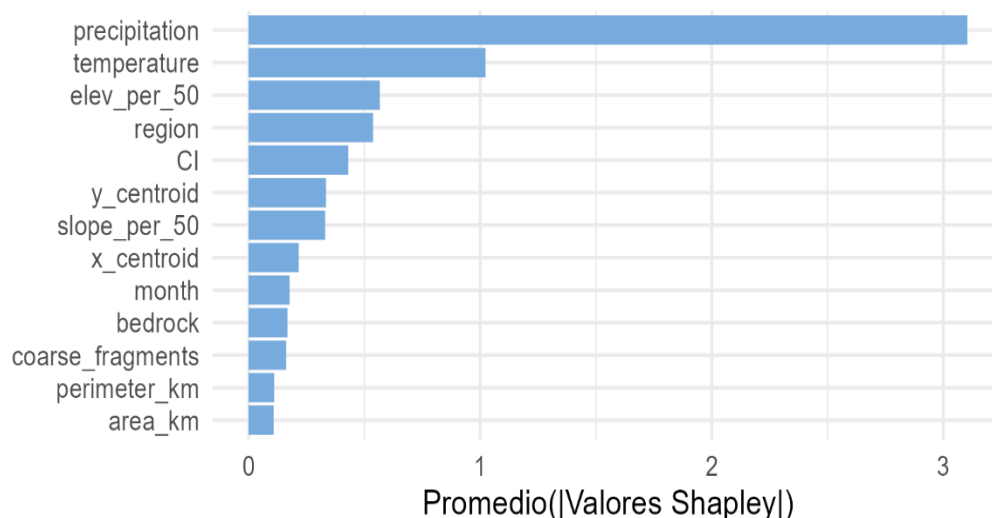
La interpretación de un modelo de *machine learning* es esencial para comprender su funcionamiento, generar confianza, corregir errores, optimizar su rendimiento y facilitar la comunicación efectiva entre los diferentes actores involucrados en su desarrollo y uso.

Para este caso se aplicará la metodología *Shapley Values* la cual es ampliamente utilizada en el campo de *Machine Learning* por proporcionar una interpretación de la contribución de las variables locales y globales, ser agnóstica al modelo de aprendizaje automático y que cuenta con una sólida base matemática.

Los valores Shapley corresponden a la contribución marginal promedio de una variable a través de todas las posibles combinaciones de las variables. En este caso, nos interesa conocer los aportes de cada una de las variables en el proceso de predicción del caudal mensual.

El promedio de los valores absolutos de los valores Shapley permite brindar una estimación global de las variables de su contribución normal. Esto se puede apreciar en la siguiente figura³:

Gráfico 12
Promedio del valor absoluto de valores Shapley



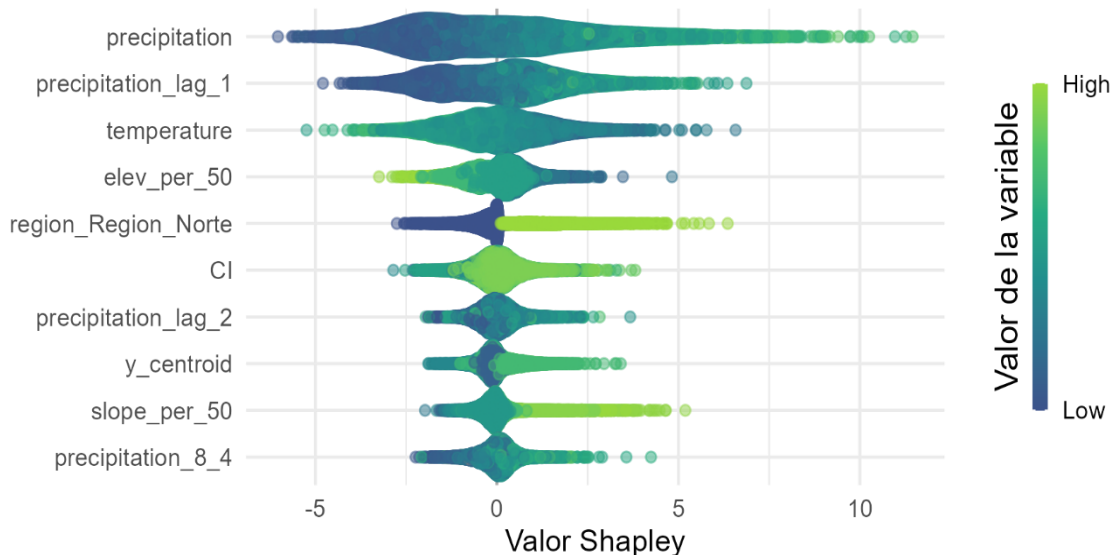
Fuente: Elaboración propia.

Como se puede observar en el gráfico anterior, precipitación, la temperatura y la elevación media son las variables que contribuyen en mayor cuantía a la predicción de caudales.

A nivel de variable y los valores de sus observaciones se puede apreciar la siguiente figura en donde se visualiza en el eje “x” la contribución marginal de las diez variables con mayor impacto promedio en el modelo. Las contribuciones positivas indican un incremento en el caudal esperado y viceversa. Por otro lado, los valores altos de las variables están coloreados en verde claro y los valores bajos en azul. La cantidad de observaciones de cada variable se ve reflejada en la amplitud del eje “y”.

³ Las variables de precipitación, temperatura, región y mes fueron agrupadas en una sola variable por grupo para facilitar la interpretación.

Gráfico 13
Valores Shapley por observación y variable

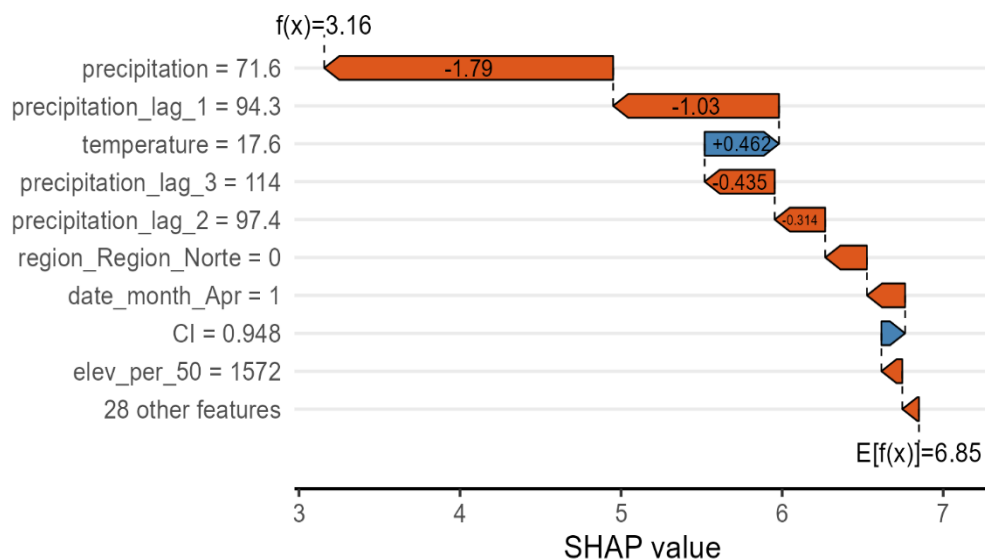


Fuente: Elaboración propia.

El gráfico 13 muestra las relaciones esperadas de las variables predictoras con la variable a predecir. En general se ve como valores altos de precipitación contribuyen a un aumento de caudales, lo contrario ocurre con la temperatura. En el caso del coeficiente de infiltración -CI-, a mayor coeficiente mayor aporte en los caudales, esto se puede explicar principalmente en los meses secos en donde los aportes de flujo base son más relevantes. En resumen, las variables más importantes tienen un comportamiento acorde con la representación física del proceso hidrológico, lo cual valida los resultados obtenidos.

Los valores Shapley también permiten analizar una predicción individual. Esto permite facilitar la interpretabilidad del modelo. En la siguiente figura se muestra una predicción $f(x)$ de la cuenca Angostura en abril de 1989 y los aportes de cada variable respecto al valor esperado $E[f(x)]$.

Gráfico 14
Contribuciones de las variables a una observación



Fuente: Elaboración propia.

En este caso, el valor de la predicción del caudal es de 3,16 mm, este valor es inferior al esperado de los datos en general (6,85 mm) principalmente por características de la precipitación y la temperatura. Por ejemplo, la precipitación en el mes en estudio fue de 71,6 mm que es muy inferior al promedio de precipitación de los datos-269 mm- lo que hace que en conjunto con las otras variables su aporte sea negativo, lo mismo se puede decir de la temperatura a la inversa en donde la temperatura en la cuenca fue inferior al promedio de los datos -18,76 °C-.

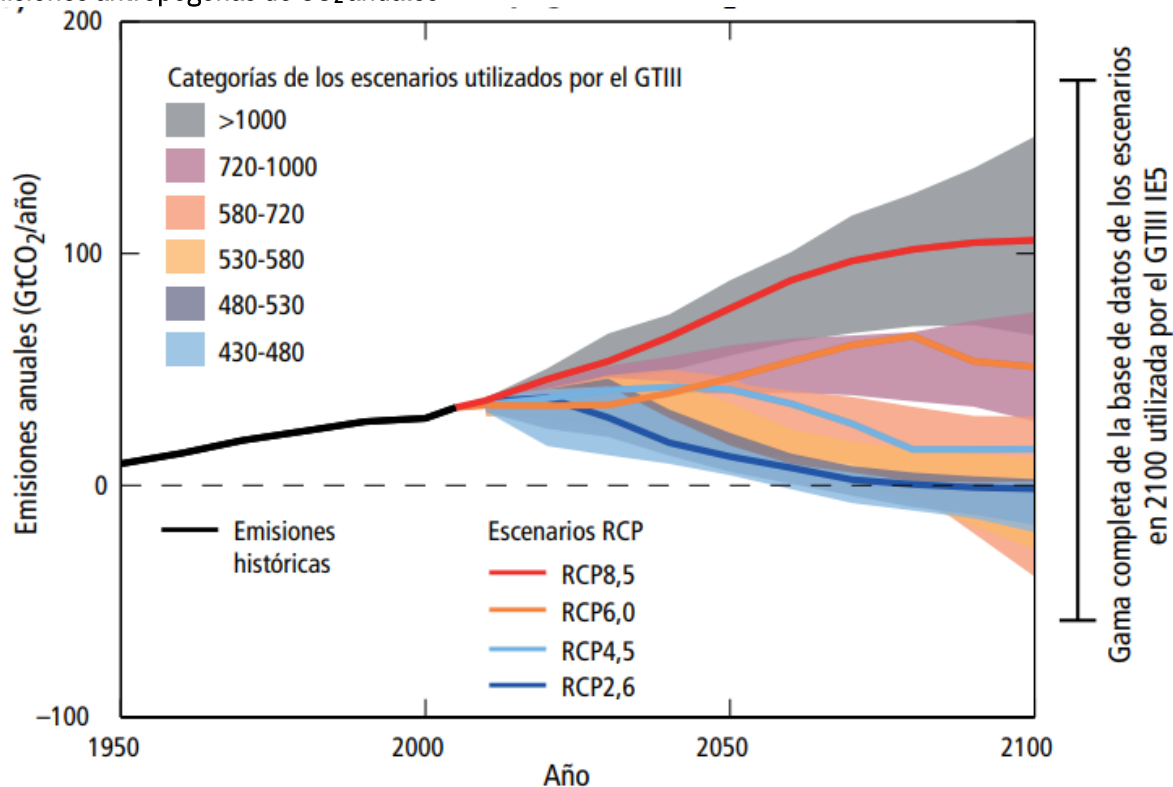
Impacto del cambio climático

El cambio climático tiene una serie de efectos significativos en el clima a nivel mundial. Puesto que estas variaciones en los patrones climáticos afectaran directamente variables relevantes para la predicción de la escorrentía tales como la temperatura y precipitación, es relevante analizar el impacto producto del cambio climático en las predicciones.

El documento Proyecciones de Cambio Climático regionalizadas para Costa Rica (Alvarado, 2021) realizado por el Instituto Meteorológico Nacional (IMN) brinda estimaciones de los cambios en precipitación, lluvia, humedad relativa, radiación solar y velocidad del viento para

los escenarios futuros RCP 2,6 y RCP 8,5 en tres períodos climáticos (2010-2039, 2040-2069 y 2070-2099). Es importante indicar que los escenarios utilizados buscan simular bajas emisiones de gases de efecto invernadero -RCP 2,6- y altas emisiones RCP 8,5- tal y como se muestra en el gráfico 15.

Gráfico 15
Emisiones antropógenas de CO₂ anuales^{a/}



a/ Emisiones antropógenas de CO₂ históricas y proyecciones de las emisiones totales de los escenarios de trayectorias de concentración representativas (RCP, por sus siglas en inglés) de gases de efecto invernadero (IPCC, 2014).

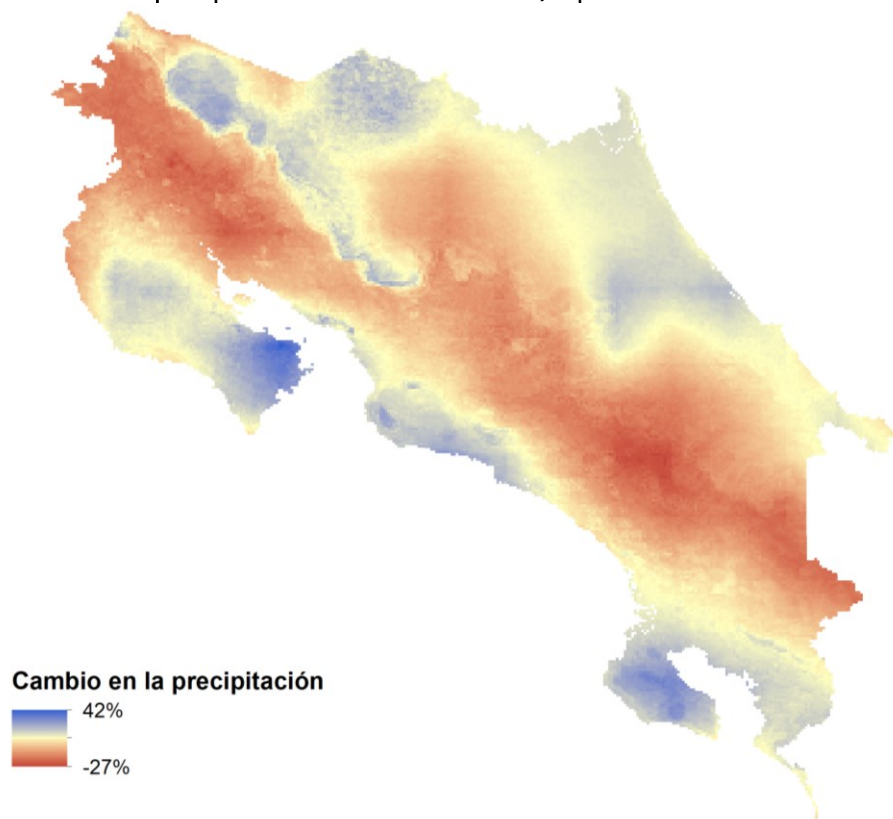
Fuente: Elaboración propia con datos de IPCC, 2014.

Para aplicar los resultados de las proyecciones climáticas, en primera instancia se calculó el diferencial porcentual entre la proyección y el dato base para cada escenario y período de tiempo.

Seguidamente este porcentaje se aplicó a las variables de precipitación y temperatura del modelo para recalcular los caudales bajo estos escenarios. Las diferencias entre los valores de los caudales permitirán estimar el impacto del cambio climático en el modelo.

Mapa 9

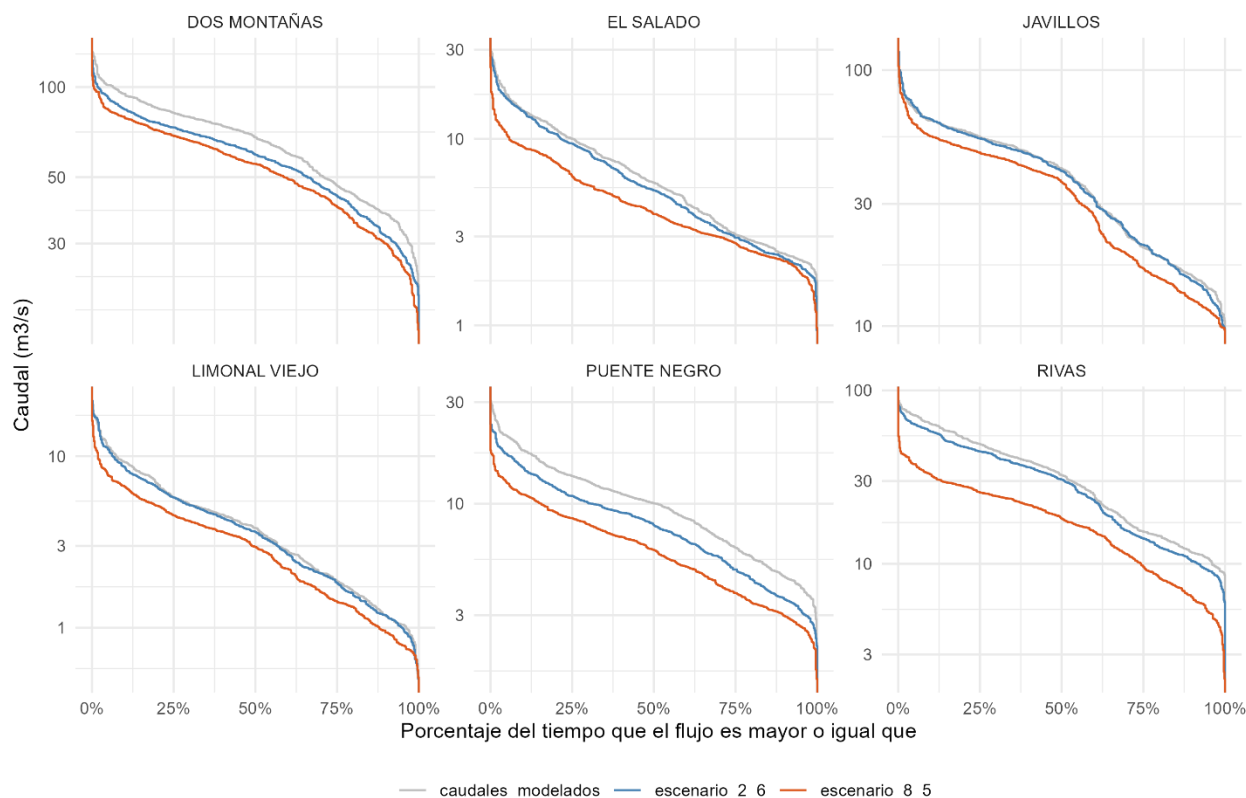
Cambio en la precipitación en escenario RCP 8,5. período 2077-2099



Fuente: Elaboración propia con datos del IMN (Alvarado, 2021).

Para ejemplificar el impacto del cambio climático se presentan las curvas de duración de las cuencas de prueba en el período de 2070-2099 para los dos escenarios y la línea base del modelo.

Gráfico 16
Curvas de duración para conjunto de cuencas de prueba^{a/}



a/ Gris: caudales modelados escenario actual. Azul: caudales estimados para escenario RCP 2,6 período 2077-2099. Naranja: caudales estimados para escenario RCP 8,5 período 2077-2099. Fuente: Elaboración propia.

A nivel general, los escenarios de cambio climático RCP 2,6 y 8,5 impactan reduciendo los caudales esperados a la baja. En Limonal Viejo, Javillos, El Salado y Rivas el escenario 2.6 afecta levemente los caudales y por ende la línea gris (caudales actuales) y la línea azul (caudales escenario 2,6) son similares.

En este caso, todos los caudales del escenario RCP 8,5 van en detrimento de la disponibilidad hídrica principalmente porque las cuencas se ubican en zonas del territorio en donde los cambios en la temperatura y precipitación combinados junto con las otras variables del modelo reducen la capacidad hídrica, sin embargo, cuando las condiciones del cambio climático son inversas su efecto resultará en un aumento del caudal disponible.

Las estadísticas de la diferencia en los caudales medios de estas cuencas se describen a continuación:

Cuadro 10

Cambios en los caudales medios estimados en RCP 2,6 y 8,5. período 2077-2099

Cuenca	Región	Q modelado medio	Q medio RCP 2,6	Q medio RCP 8,5	Porcentaje cambio RCP 2,6	Porcentaje cambio RCP 8,5
Dos Montañas	Caribe Sur	8,79	7,81	7,24	-11,12	-17,59
El Salado	Pacifico Central	3,62	3,38	2,44	-6,77	-32,52
Javillos	Región Norte	6,43	6,39	5,56	-0,64	-13,58
Limal Viejo	Pacifico Norte	3,01	2,89	2,25	-4,19	-25,31
Puente Negro	Central	3,28	2,64	2,02	-19,37	-38,31
Rivas	Pacifico Sur	9,45	8,57	5,19	-9,24	-45,11

Fuente: Elaboración propia.

Como es de esperar los resultados observados en las curvas de duración se corresponden con los caudales medios. En el Cuadro 10 se puede apreciar que bajo el escenario RCP 2,6 las variaciones en los caudales varían entre -0,6% y -19,37% mientras que en el escenario RCP 8,5 entre -13,58% y -38,31%.

Es importante mencionar que los resultados de los impactos del cambio climático deben ser utilizados como estimaciones con alta incertidumbre tanto por las proyecciones de los modelos climáticos como por las estimaciones del modelo presentado en esta investigación. Las estimaciones suponen que el cambio porcentual en la precipitación y temperatura se aplican igual en todos los meses y años, asumen que el clima de referencia del IMN es comparable con el utilizado por esta investigación, buscan reflejar el impacto del cambio climático en el modelo el cual no es una representación de los procesos físicos, a pesar de que las variables más importantes del modelo son la precipitación y la temperatura, existen otras variables

relacionadas que podrían influir en el resultado y que son aplicables siempre y cuando existan observaciones similares en los datos de entrenamiento.

Caso de estudio cuenca Chirripó Pacífico

La transición de una investigación a la aplicación es un paso no trivial y fundamental para impactar en el proceso de desarrollo de la sociedad. Este trabajo no solo busca desarrollar y probar modelos a nivel teórico si no que tiene como fin último democratizar la información, facilitar la barrera de ingreso al modelo automatizando una serie de procesos como la creación y caracterización de la cuenca, obtención de información base, aplicación del modelo y escenarios de cambio climático y análisis hidrológicos derivados en una sola aplicación.

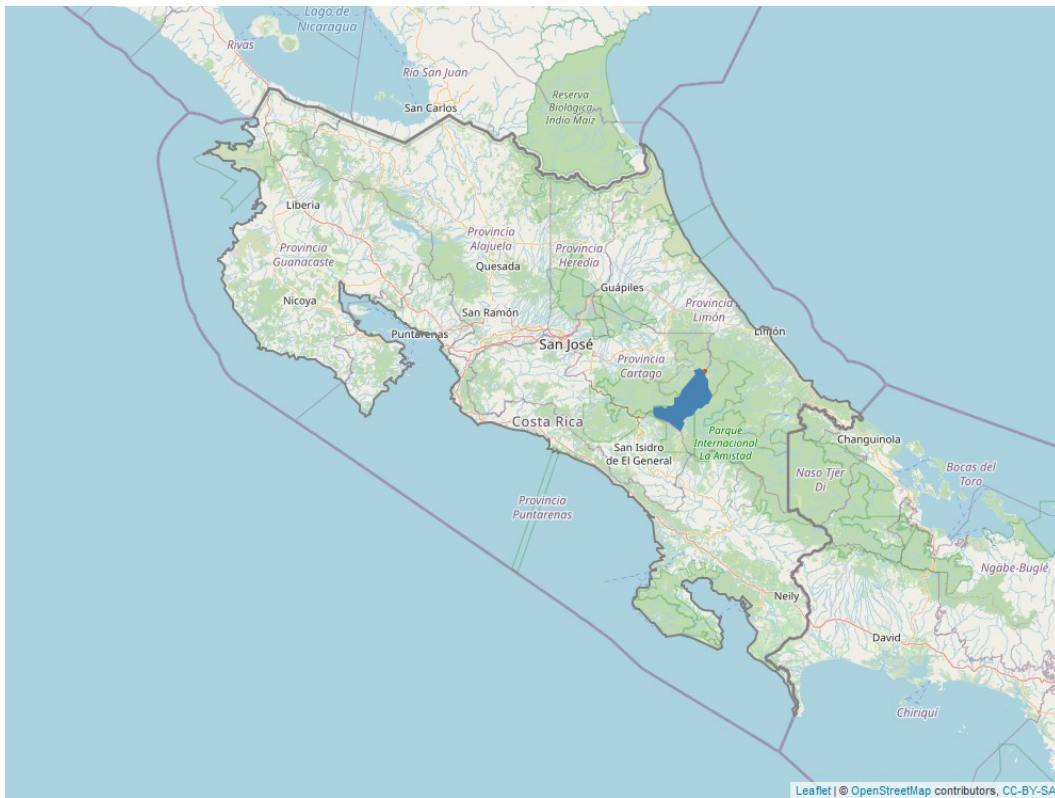
Como caso de estudio, se aplicarán los resultados de la investigación a una cuenca sin medición de caudales con el fin de ilustrar la información y posibilidades que este trabajo puede ofrecer a la sociedad en temas de la Gestión Integrada del Recurso Hídrico. Este apartado es una muestra fidedigna de lo que se le podría ofrecer al usuario con tan solo contar con unas coordenadas geográficas del río que quieren analizar.

Descripción de la cuenca

La cuenca fue delimitada a partir del punto de control fijado en las coordenadas 1 082 811 (y) y 570 022 (x) CRTM05.

Mapa 10

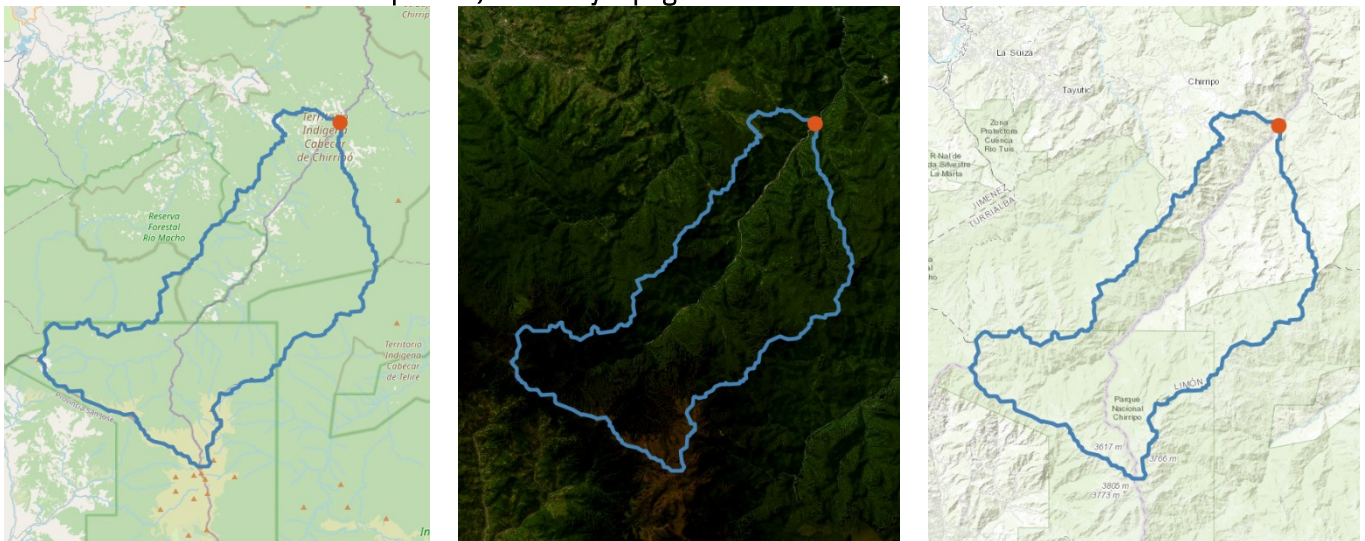
Ubicación de la cuenca a nivel nacional



Fuente: OpenStreetMap, 2023.

Mapa 11

Visualizaciones de la cuenca política, satelital y topográfica



Fuente: OpenStreetMap y Esri, 2023.

La cuenca tiene un área de 487 km² y un perímetro de 163,8 km y está ubicada en la región climática Caribe Sur dentro de la Gran Cuenca Chirripó Caribe.

Un resumen de las características más relevantes de la cuenca se detalla en el Cuadro 11.

Cuadro 11

Características de la cuenca

Variable	Valor
Coordenada punto de control en X	570022
Coordenada punto de control en Y	1082811
Área (km ²)	487
Perímetro (km)	163,8
Coord. centroide en X	559508,9
Coord. centroide en Y	1065582
Región climática	Caribe Sur
Gran Cuenca	Chirripó Caribe
Elevación mínima (m.s.n.m.)	517,44
Elevación percentil 5 (m.s.n.m.)	828,67
Elevación percentil 25 (m.s.n.m.)	1313,11
Elevación percentil 50 (m.s.n.m.)	1783,44
Elevación percentil 75 (m.s.n.m.)	2498,44
Elevación percentil 95 (m.s.n.m.)	3204,47
Elevación máxima (m.s.n.m.)	3769,67
Pendiente percentil 5 (%)	9,1
Pendiente percentil 25 (%)	17,8
Pendiente percentil 50 (%)	24,24
Pendiente percentil 75 (%)	30,61
Pendiente percentil 95 (%)	38,86
Profundidad media al lecho recoso (m)	158,67

Variable	Valor
Porcentaje de fragmentos gruesos en primeros 30cm de suelo (%)	10,36
Área bajo la curva hipsométrica	0,42

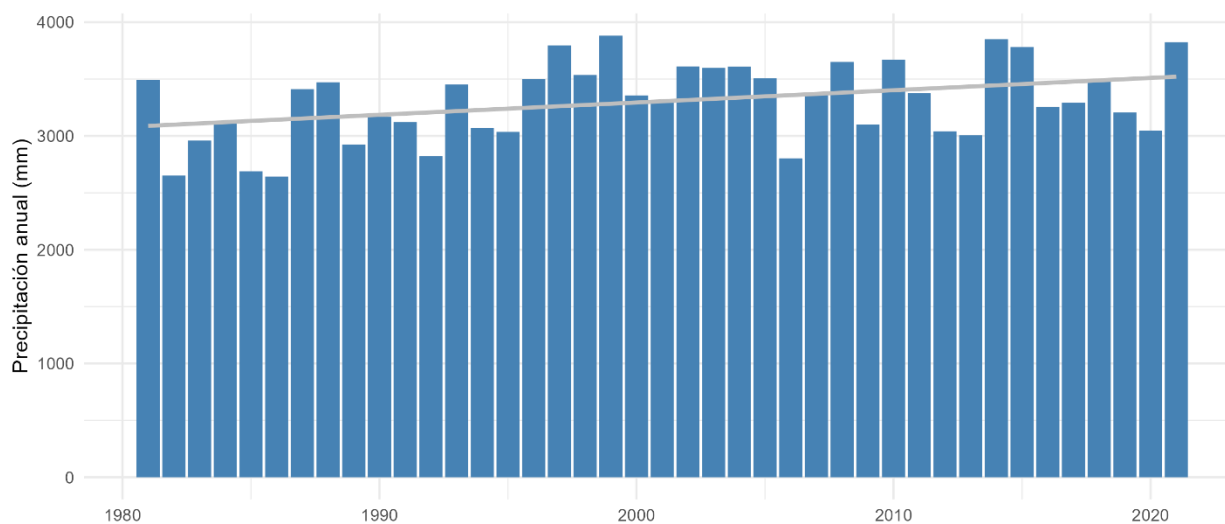
Fuente: Elaboración propia.

Precipitación

La cuenca tiene una precipitación media anual de 3305 mm. Si bien la precipitación tiene una tendencia al alza a una razón de 10,80 mm/año este aumento no es significativo -p valor 0,014.

Los años más secos⁴ registrados fueron 1982, 1985 y 1986 con un promedio de 2661 mm/año y los años húmedos⁵ fueron 1999, 2014 y 2021 con un promedio de 3852 mm/año.

Gráfico 17
Precipitación anual
(mm)



Fuente: Elaboración propia.

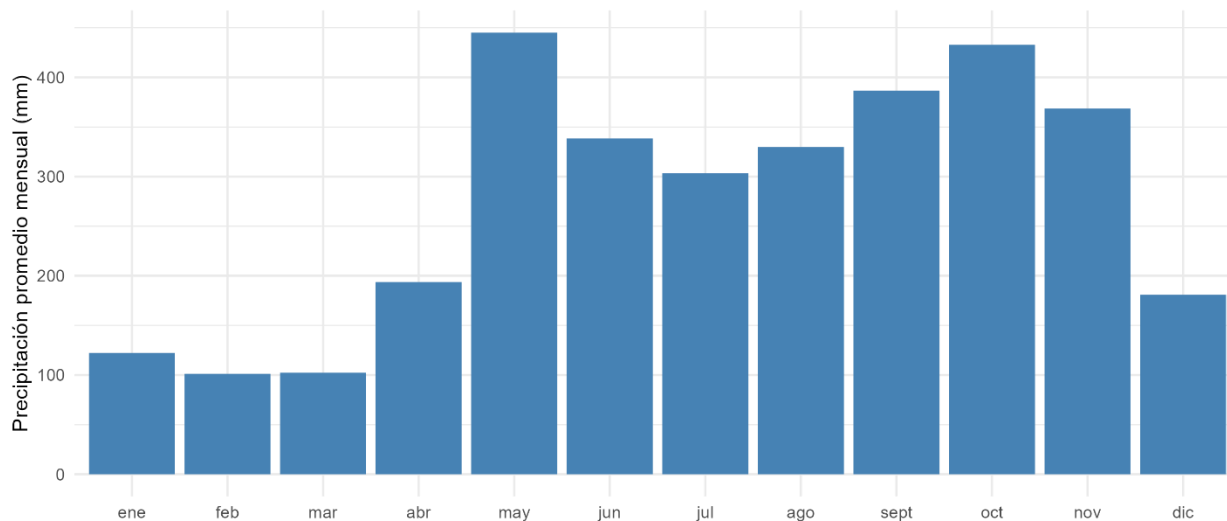
⁴ Con valores de z-score por debajo de 1,5.

⁵ Con valores de z-score por encima de 1,5.

A nivel mensual los meses con menor precipitación promedio son marzo y febrero y los meses más lluviosos son octubre y mayo.

Gráfico 18

Precipitación media mensual

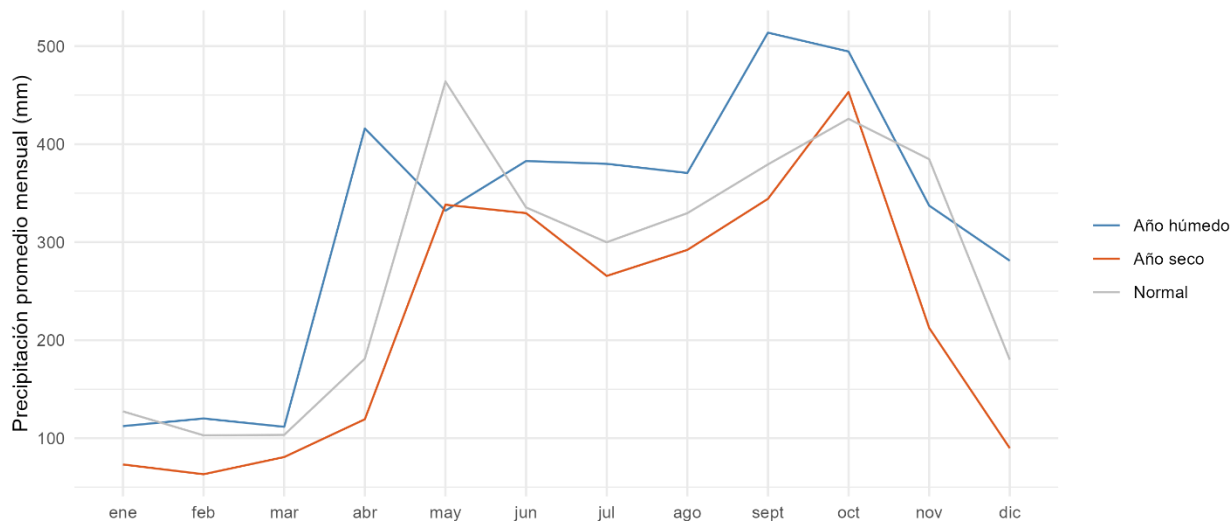


Fuente: Elaboración propia.

Al agrupar la información por tipo de año se obtiene la siguiente distribución de precipitación:

Gráfico 19

Precipitación media mensual por tipo de año



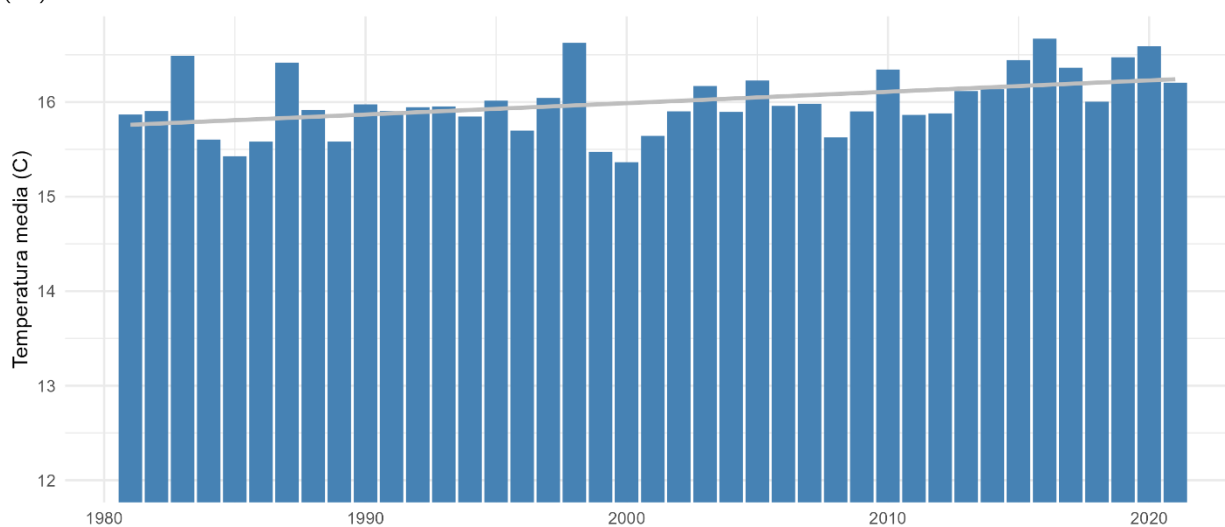
Fuente: Elaboración propia.

Temperatura

La cuenca tiene una temperatura media anual de 16,00 °C. Si bien la temperatura tiene una tendencia al alza a una razón de 0,12 °C /año este aumento no es significativo -p valor 0,0047-.

Los años más fríos⁶ registrados fueron 1985, 1999 y 2000 con un promedio de 15,42 °C y los calientes⁷ fueron 1998, 2016 y 2020 con un promedio de 16,63 °C.

Gráfico 20
Temperatura anual
(°C)



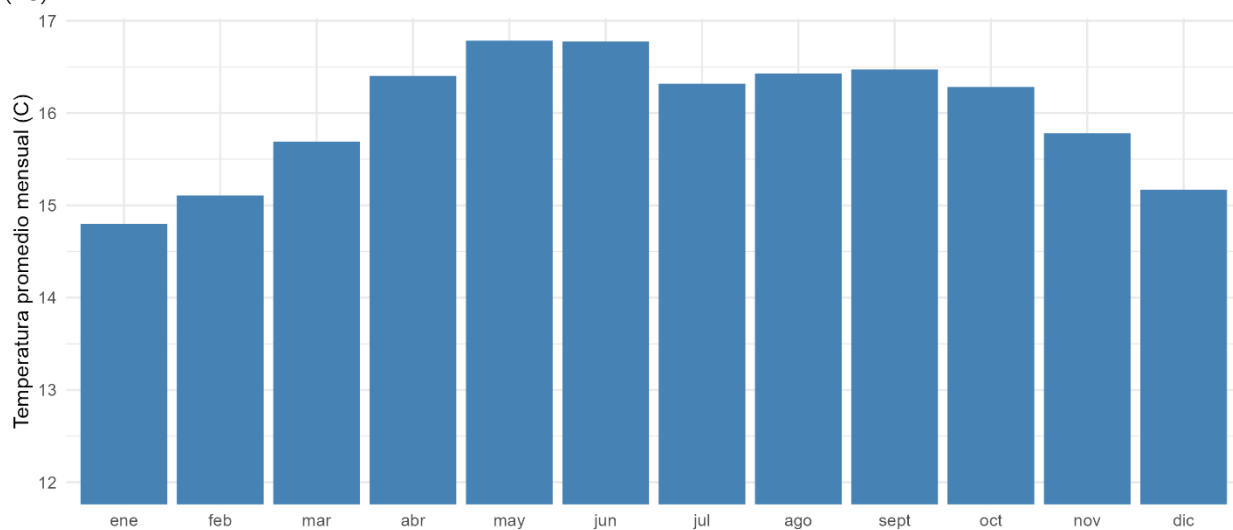
Fuente: Elaboración propia.

A nivel mensual los meses con menor temperatura promedio son enero y febrero y los meses más calientes son junio y mayo.

⁶ Con valores de z-score por debajo de 1,5.

⁷ Con valores de z-score por encima de 1,5.

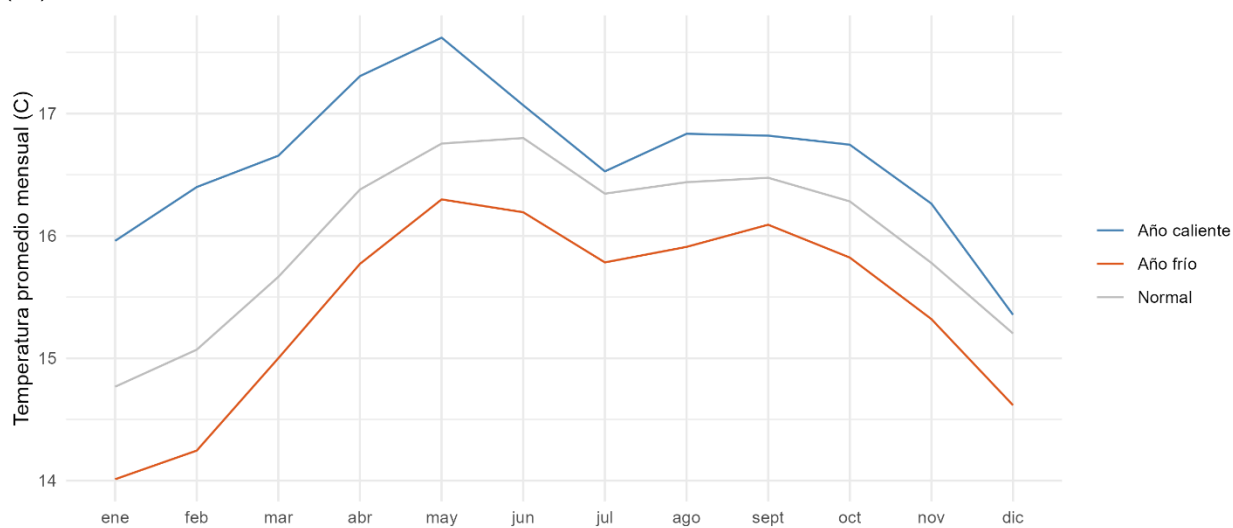
Gráfico 21
Temperatura media mensual
(°C)



Fuente: Elaboración propia.

Al agrupar la información por tipo de año se obtiene la siguiente distribución de temperatura:

Gráfico 22
Temperatura media mensual por tipo de año
(°C)



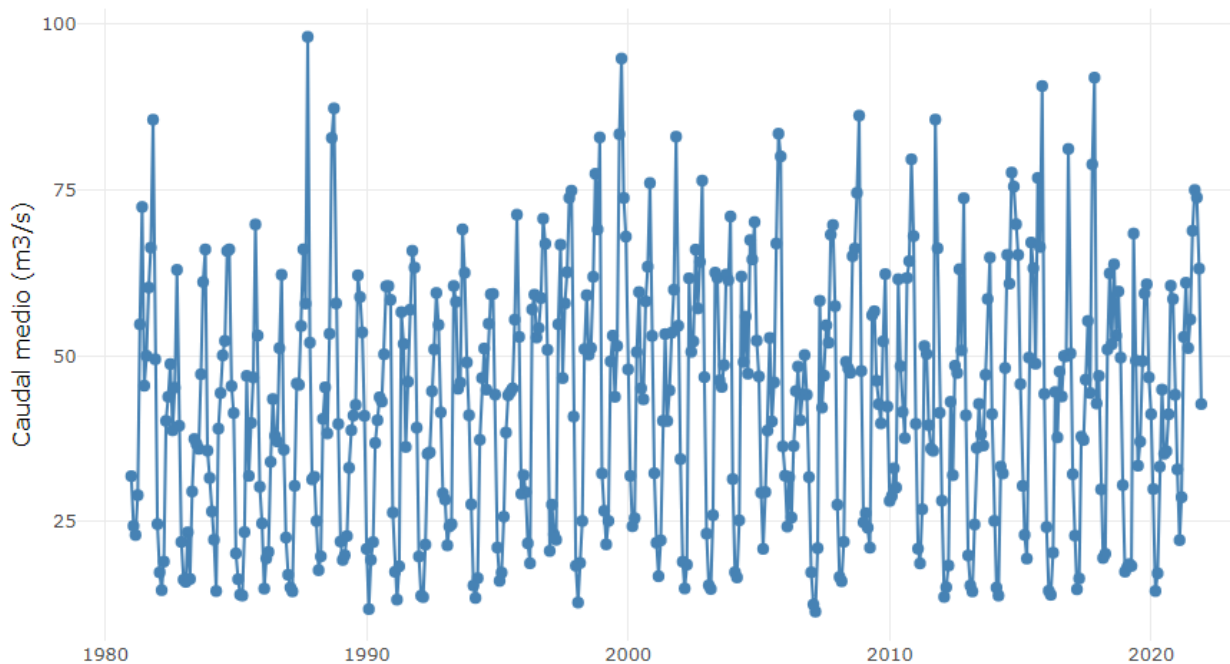
Fuente: Elaboración propia.

Caudales

Los caudales mensuales aquí presentados son estimaciones del modelo de aprendizaje automático desarrollado en la investigación y está sujeto a una incertidumbre tal y como se detalla en la sección 0 Resultados del modelo.

Al aplicar el modelo sobre la cuenca se obtiene la siguiente serie de tiempo:

Gráfico 23
Caudales promedio mensuales estimados
(m³/s)

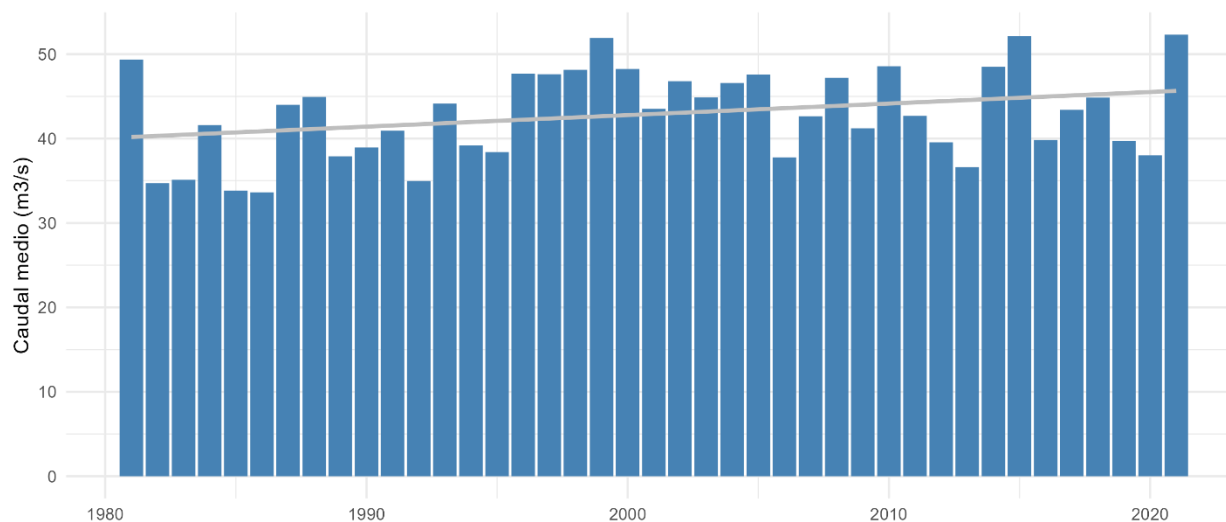


Fuente: Elaboración propia.

La cuenca tiene un caudal medio anual de 42,62 m³/s. Si bien el caudal tiene una tendencia al alza a una razón de 0,14 m³/s año este aumento no es significativo -p valor 0,047-.

Los años con menores caudales⁸ estimados fueron 1982, 1985 y 1986 con un promedio de 34,06 m³/s y los más caudalosos⁹ fueron 1999, 2014 y 2021 con un promedio de 50,96 m³/s.

Gráfico 24
Caudal promedio anual
(m³/s)



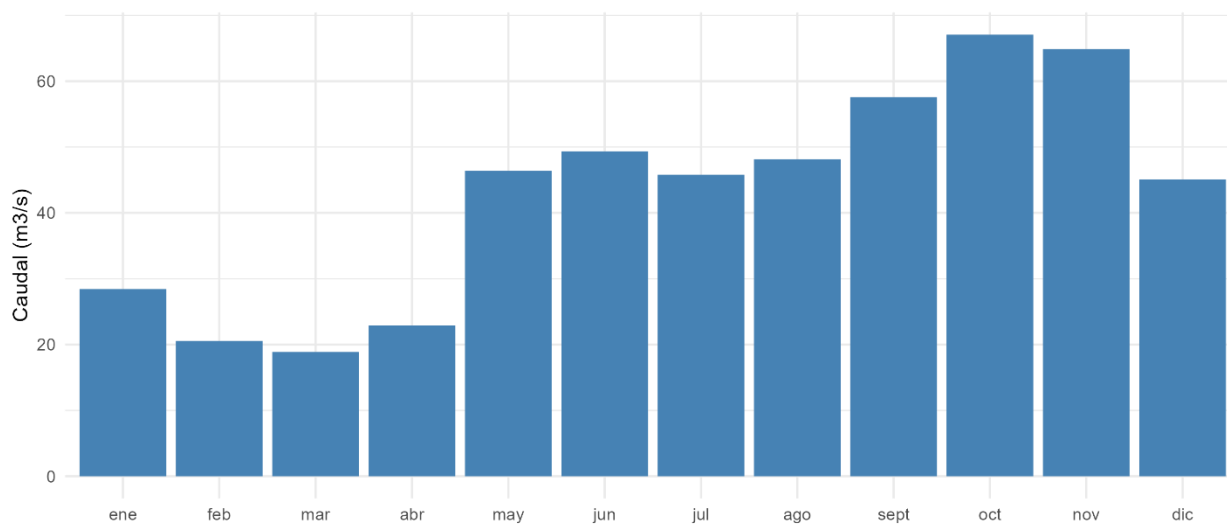
Fuente: Elaboración propia.

A nivel mensual los meses con menor caudal promedio son marzo y febrero y los meses más caudalosos son octubre y noviembre.

⁸ Con valores de z-score por debajo de 1,5.

⁹ Con valores de z-score por encima de 1,5.

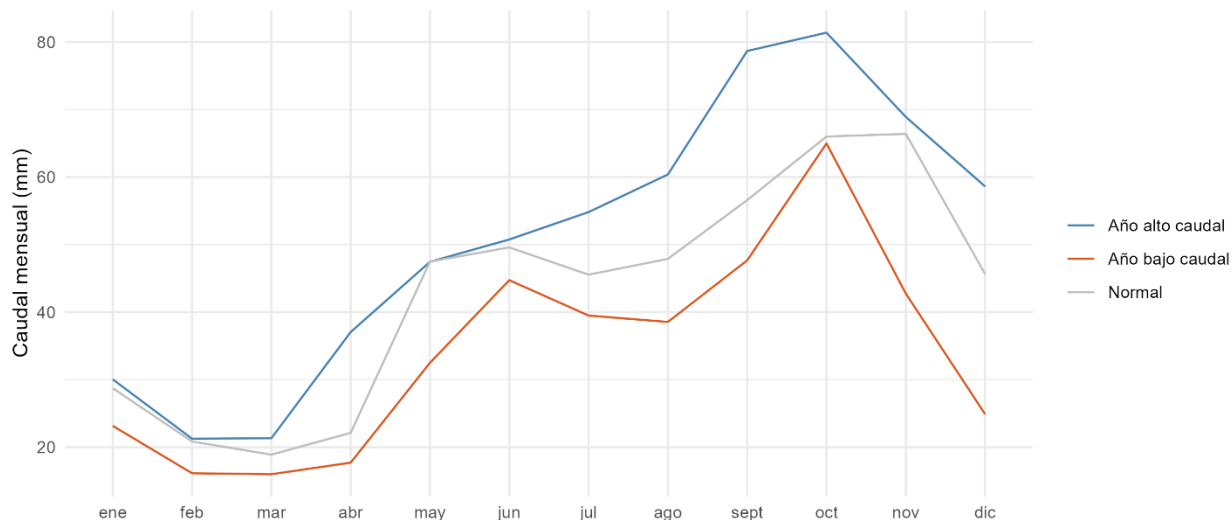
Gráfico 25
Caudal medio mensual
(m³/s)



Fuente: Elaboración propia.

En diversos escenarios es relevante conocer el comportamiento mensual de los caudales en casos de años secos, húmedos o normales. Al agrupar la información por tipo de año se obtiene la siguiente distribución de caudales:

Gráfico 26
Caudal medio mensual por tipo de año
(m³/s)



Fuente: Elaboración propia.

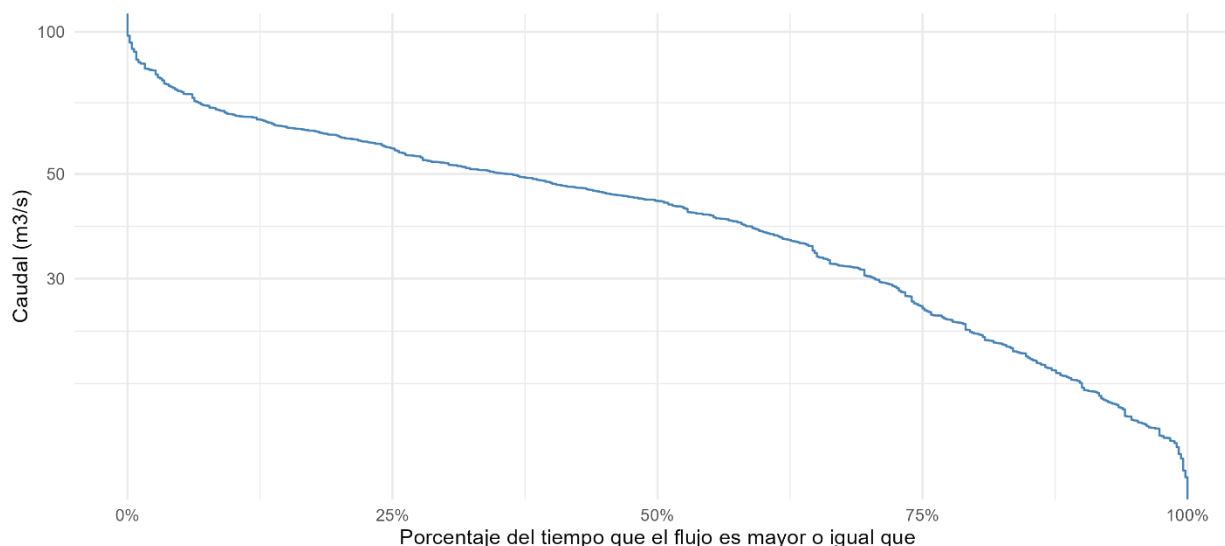
Las curvas de duración facilitan la toma de decisiones estratégicas como la gestión eficiente de recursos hídricos y la planificación en la capacidad de producción de un cauce. Las curvas de duración muestran el porcentaje del tiempo que los caudales mensuales son igualados o excedidos en el período de estudio.

En esta curva los valores de los caudales se trazan desde los más altos hasta los más bajos en el eje “Y” y en el eje “X” se grafica el porcentaje del tiempo (meses) en el que flujo medio mensual es menor o igual que un caudal específico. En este caso, la curva de duración corresponde al promedio del período 1981-2021. Por tener 40 años de extensión se considera que engloba diferentes fenómenos climáticos relevantes. Si se conoce que en la cuenca ha habido o habrá cambios significativos en el comportamiento natural de la cuenca como derivaciones para riego, represamiento u otros, las estimaciones aquí presentadas no deberían ser usadas como referencia.

El gráfico 27 muestra la curva de duración mensual para la cuenca seleccionada.

Gráfico 27

Curva de duración mensual. 1981-2021



Fuente: Elaboración propia.

A manera de ejemplo, el 75% del tiempo de la cuenca tendrá caudales medios mensuales menores o iguales a 56,7 m³/s. Para analizar en detalle la curva de duración referirse al

Anexos

Anexo 1 caso de estudio

Cuadro 13

Cambio climático

Los resultados presentados hasta el momento consideran los caudales históricos desde 1980 hasta el 2021. Puesto que el cambio climático tendrá afectaciones en el ciclo hidrológico, se aplicarán dos escenarios analizados en el documento Proyecciones de Cambio Climático regionalizadas para Costa Rica (Alvarado, 2021) realizado por el Instituto Meteorológico Nacional (IMN). Los escenarios brindan estimaciones de los cambios en precipitación, lluvia, humedad relativa, radiación solar y velocidad del viento para los escenarios futuros RC P2,6 y RCP 8,5 en tres períodos climáticos (2010-2039, 2040-2069 y 2070-2099). Es importante indicar que los escenarios utilizados buscan simular bajas emisiones de gases de efecto invernadero -RCP 2.6- y altas emisiones RCP 8.5- tal y como se muestra en gráfico 15.

En este caso, se aplicarán al modelo de aprendizaje automático los dos escenarios para el período 2070-2099 al modificar las variables de precipitación y temperatura según la variación esperada en la cuenca producto del cambio climático.

Cuadro 12

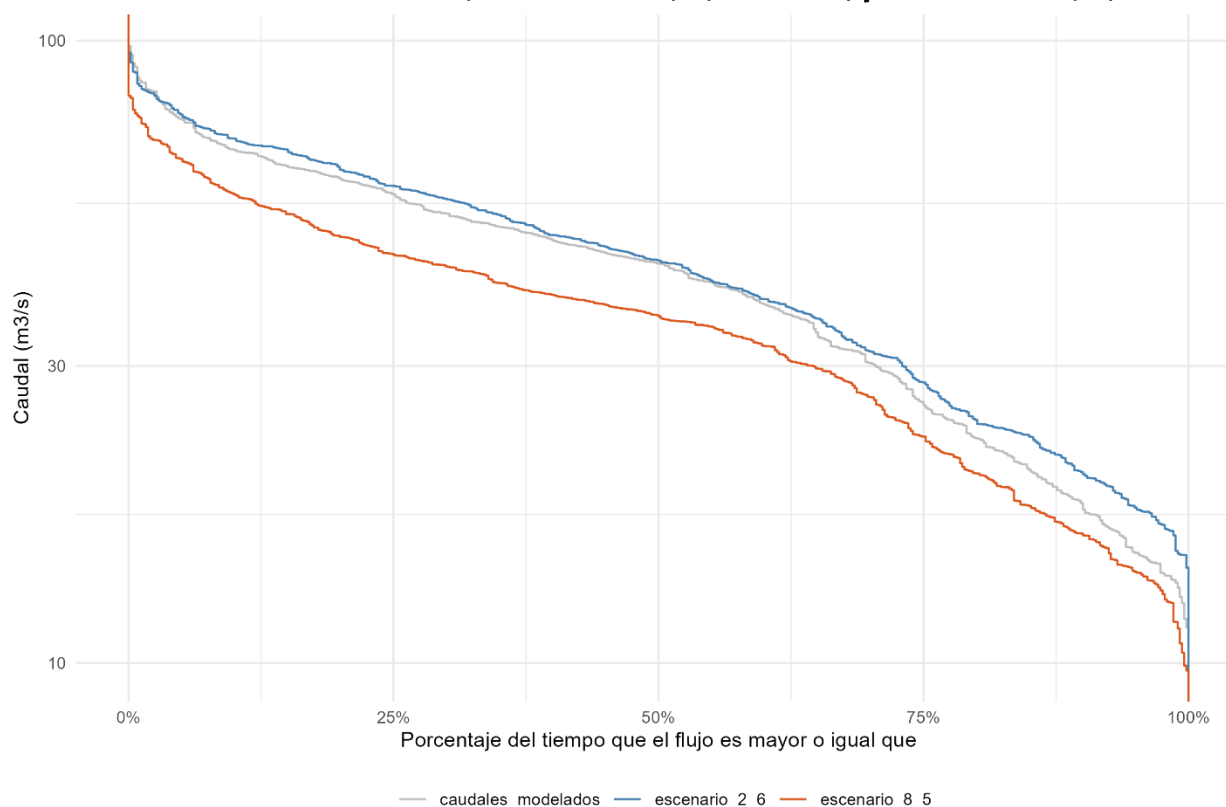
Cambios en los caudales medios estimados en RCP 2,6 y 8,5. período 2077-2099

Q modelado medio sin CC (m ³ /s)	Escenario CC	Cambio temp. (%)	Cambio prec. (%)	Q medio con CC (m ³ /s)	Diferencia en caudales (%)
42,916	2,6	-20,14	-7,15	44,65	4,05
42,916	8,5	-6,59	-16,09	36,13	-15,81

Fuente: Elaboración propia.

Gráfico 28

Curva de duración mensual 1981-2021, escenario RCP 2,6 (2070-2099) y escenario RCP 8,5 (2070-2099)



Fuente: Elaboración propia.

Según estos hallazgos, durante el periodo 2070-2099, se anticipa un aumento en la capacidad hídrica de la cuenca analizada, con un incremento del 4,5% en el caudal promedio, en caso de que se materialice el escenario de bajas emisiones de gases de efecto invernadero (RCP 2,6). En cambio, si se llega a dar el escenario RCP 8,5, se proyecta una disminución de aproximadamente un 15,81% en la cantidad de recurso hídrico disponible.

Conclusiones y recomendaciones

Este estudio ha logrado desarrollar con éxito un modelo de aprendizaje automático para predecir los caudales mensuales de 65 cuencas en todo el país, con métricas de desempeño notables, como una mediana de KGE de 0,7 en el conjunto de validación y de 0,67 en el conjunto de prueba. Estas cifras son comparables o superiores a modelos hidrológicos tradicionales realizados anteriormente en el país. Por sus características, la aplicabilidad del

modelo se puede extender a cuencas similares y será especialmente útil en unidades sin mediciones hidrológicas.

La utilización de información abierta como predictores ha resultado efectiva, destacando la viabilidad de este enfoque en el modelaje hidrológico en Costa Rica. La escalabilidad, diversidad de cuencas y la disponibilidad de información abierta permiten que los resultados de esta investigación sean generalizables a un amplio conjunto de regiones.

Por estas características, se propone la creación de una herramienta en línea basada en estos resultados, abierta al público, que democratice la información, acelere la obtención de datos, reduzca costos económicos y sirva como base para iniciar proyectos hídricos.

En un esfuerzo conjunto entre instituciones, la herramienta propuesta puede integrarse con otros instrumentos complementarios, formando un conjunto de materiales para proporcionar información más completa y valiosa a los usuarios. Los datos base de caudal son vitales para el desarrollo de estos artefactos y, por ende, sería de gran utilidad para el desarrollo científico del país contar con ellos de manera abierta.

La posibilidad de actualizar la herramienta periódicamente con nuevos datos meteorológicos garantiza la vigencia de las estimaciones en el corto y mediano plazo. Sin embargo, para asegurar su relevancia a largo plazo, se recomienda reajustar el modelo para capturar cambios en el comportamiento de las cuencas y el clima, junto con la actualización de los datos de caudales.

Se aconseja a los usuarios siempre tener en cuenta la incertidumbre inherente al modelado. Aunque este modelo ha demostrado ofrecer resultados satisfactorios, es esencial reconocer que existe un grado de incertidumbre en sus predicciones y en sus datos base. Se sugiere a los usuarios que evalúen esta incertidumbre y la consideren según sus necesidades y contextos específicos.

En el proceso de aplicar el modelo a otras regiones, se recomienda verificar la similitud física y climatológica de las cuencas seleccionadas con aquellas utilizadas en el estudio original. Esta consideración desempeña un papel fundamental en el fortalecimiento de la generalización del modelo.

La eficacia del modelo para la imputación de datos de caudal faltantes se evidencia mediante los sólidos resultados obtenidos en el conjunto de prueba intra-cuenca (El Brujo y El Rey), con valores de KGE de 0,81 y 0,83. La implementación de esta herramienta en cuencas con información incompleta podría ser altamente beneficiosa, facilitando la recuperación de registros faltantes y permitiendo la continuación de investigaciones adicionales.

Se insta a futuros estudios a explorar nuevas variables predictoras y otras técnicas de aprendizaje automático, como redes neuronales *Long Short-Term Memory* o *Transformers*, que también han demostrado utilidad en estos escenarios.

En conclusión, la implementación de los hallazgos de este estudio ofrecerá al país una herramienta valiosa para la Gestión Integral del Recurso Hídrico. Este instrumento agilizará y simplificará la adquisición de datos hidrológicos a un costo mínimo, brindando a los distintos actores una mejora significativa en la comprensión y control de los caudales de los ríos en una amplia región del territorio nacional.

Bibliografía

- Alvarado, L. F. (2021). *Proyecciones de Cambio Climático regionalizadas para Costa Rica*. San José: IMN-PNUD.
- Arciniega-Esparza, S., Birkel, C., Chavarría-Palma, A., Arheimer, B., & Agustín. (2022). *Remote sensing-aided rainfall-runoff modelling in the tropics of Costa*. Hydrology Earth Systems Science.
- Clapp, R., & Hornberger, G. (1978). *Empirical equations for some soil hydraulic properties*. Virginia: Department of Environmental Sciences, University of Virginia.
- Copernicus Climate Change Service. (2019). *Land cover classification gridded maps from 1992 to present derived from satellite observation*. Copernicus Climate Change Service (C3S) Climate Data Store (CDS).
- Dai, Y., Shangguan, W., Wei, N., Xin, Q., Yuan, H., Zhang, S., . . . Yan, F. (2019). *A review of the global soil property maps for Earth system models*. SOIL.
- Farr, T., Rosen, P., Caro, E., Crippen, R., Duren, R., Hensley, S., . . . Alsdorf, D. (2007). *The shuttle radar topography mission: Reviews of Geophysics*.
- Funk, C. P. (2014). *A quasi-global precipitation time series for drought monitoring*. U.S. Geological Survey Data Series 832.
- Gupta, H. V., Kling, H., Yilmaz, K. K., & Martinez, G. F. (2009). *Decomposition of the mean squared error and NSE performance criteria: Implications for improving hydrological modelling*. Journal of Hydrology, Elsevier.
- Hengl, T., Jesús, J. M., Heuvelink, G., González, M., Kilibarda, M., Blagotić, A., . . . Leenars, J. (2017). *SoilGrids250m: Global gridded soil information based on machine learning*. PLOS ONE.
- IPCC. (2014). *Cambio climático 2014: Informe de síntesis. Contribución de los Grupos de trabajo I, II y III al Quinto Informe de Evaluación del Grupo*. Ginebra.

- Jackson, E. K., Roberts, W., Nelsen, B., Williams, G. P., Nelson, E. J., & Ames, D. P. (2019). *Introductory overview: Error metrics for hydrologic modelling – A review of common practices and an open source library to facilitate use and adoption*. Environmental Modelling and Software, Elsevier.
- Kaune, A. (2021). *Cuenca El Brujo: Estudios de ingeniería para evaluar el impacto de la variabilidad climática en varias cuencas donde el AyA aprovecha recurso hídrico superficial*. Instituto Costarricense de Acueductos y Alcantarillados (AyA).
- Kaune, A. (2021). *Cuenca El Rey: Estudios de ingeniería para evaluar el impacto de la variabilidad climática en varias cuencas donde el AyA aprovecha recurso hídrico superficial*. Instituto Costarricense de Acueductos y Alcantarillados (AyA).
- Kling, H., Fuchs, M., & Paulin, M. (2012). *Runoff conditions in the upper Danube basin under an ensemble of climate change scenarios*. Journal of Hydrology, Elsevier.
- Kling, H., Fuchs, M., & Paulin, M. (2012). *Runoff conditions in the upper Danube basin under an ensemble of climate change scenarios*. Journal of Hydrology, Elsevier.
- Knoben, W. J., Freer, J. E., & Woods, R. A. (2019). *Technical note: Inherent benchmark or not? Comparing Nash-Sutcliffe and Kling-Gupta efficiency scores*. Hydrology and Earth System Sciences.
- Kratzert, F., Klotz, D., Herrnegger, M., Sampson, A. K., Hochreiter, S., & Nearing, G. S. (2019). *Toward Improved Predictions in Ungauged Basins: Exploiting the Power of Machine Learning*. Water Resources Research, Wiley.
- L. Crochemore, K. I. (2020). *Lessons learnt from checking the quality of openly accessible river flow data worldwide*. Hydrological Sciences Journal.
- Land Cover Climate Change Initiative partnership. (2017). *Land Cover CCI Product User Guide Version 2.0*. Bélgica: UCL-Geomatics.

- Lundberg, S. M., & Lee, S.-I. (2017). *A Unified Approach to Interpreting Model Predictions*. Curran Associates, Inc.
- McMillan, H., Seibert, J., Petersen-Overleir, A., Lang, M., White, P., Snelder, T., . . . Kiang, J. (2017). *How uncertainty analysis of streamflow data can reduce costs and promote robust decisions in water management applications*. Water Resources Research.
- Mendez, M., Calvo-Valverde, L.-A., Imbach, P., Maathuis, B., Hein-Grigg, D., Hidalgo-Madriz, J.-A., & Alvarado-Gamboa, L.-F. (2022). *Hydrological Response of Tropical Catchments to Climate Change as Modeled by the GR2M Model: A Case Study in Costa Rica*. Sustainability.
- Moriasi, D. N., Arnold, J. G., Liew, M. W., Bingner, R. L., Harmel, R. D., & Veith, T. L. (2007). *Model evaluation guidelines for systematic quantification of accuracy in watershed simulations*. American Society of Agricultural and Biological Engineers.
- Muñoz Sabater, J. (2019). *ERA5-Land monthly averaged data from 1950 to present*. Copernicus Climate Change Service (C3S) Climate Data Store (CDS).
- Schosinsky, G. (2006). *Cálculo de la recarga potencial de acuíferos mediante un balance hídrico de suelos*. San José: Escuela Centroamericana de Geología, Universidad de Costa Rica, .
- Schosinsky, G., & Losilla, M. (2000). *Modelo analítico para determinar la infiltración con base en la lluvia mensual*. San José: Escuela Centroamericana de Geología, Universidad de Costa Rica.
- Solano, J., & Villalobos, R. (s.f.). *Regiones y subregiones climáticas de Costa Rica*. Instituto Meteorológico Nacional.
- Tom G. Farr, P. A. (2007). *The Shuttle Radar Topography Mission*. Reviews of Geophysics.
- Venegas-Cordero, N., Birkel, C., Giraldo-Osorio, J., Correa-Barahona, A., Durán-Quesada, A. M., Arce-Mesen, R., & Nauditt, A. (2021). *Can hydrological drought be efficiently predicted*

by conceptual rainfall-runoff models with global data products? Journal of Natural Sciences and Development.

Zounemat-Kermani, M., Batelaan, O., Fadaee, M., & Hinkelmann, R. (2021). *Ensemble machine learning paradigms in hydrology: A review*. Elsevier.

Anexos

Anexo 1 caso de estudio

Cuadro 13

Valores de curva de duración mensual

Porcentaje del tiempo que el flujo es mayor o igual que	Caudal (m ³ /s)	Caudal escenario cambio climático 8.5 (m ³ /s)	Caudal escenario cambio climático 2.6 (m ³ /s)
1	13.4	11.6	15.1
5	15	14	17.9
10	18.1	16.2	20.3
15	20.5	17.9	23.2
20	23	20.2	24.6
25	26.3	23.1	28.3
30	30.4	26.8	31.7
35	34.1	29.8	35.8
40	37.8	32.3	38.4
45	41	34.7	41.2
50	43.8	36.2	44.4
55	45.7	37.7	46.7
60	48	39.1	48.7
65	50.3	40.8	52.5
70	52.8	43.5	55.7
75	56.7	45.4	58.4
80	60.2	48.5	62.4
85	63	52.7	66.9
90	66.9	56.8	69.6
95	74.9	64.7	76.2
99	87.5	75.7	85.5

Fuente: Elaboración propia.

*Modelo de machine learning para estimar caudales
mensuales históricos de cuencas hidrográficas en Costa Rica*

Cuadro 14

Caudales medios modelados

Año	Enero	Feb.	Marzo	Abril	Mayo	Junio	Julio	Ago.	Sept.	Oct.	Nov.	Dic.
1981	31.8	24.3	23.0	29.0	54.7	72.4	45.5	50.0	60.3	66.3	85.6	49.5
1982	24.6	17.3	14.6	18.9	40.2	43.8	48.7	38.8	45.2	63.0	39.4	21.9
1983	16.2	15.9	23.4	16.3	29.5	37.4	36.7	35.9	47.2	61.1	66.0	35.7
1984	31.6	26.5	22.2	14.5	39.0	44.4	50.1	52.2	65.8	66.0	45.4	41.4
1985	20.2	16.3	14.0	13.8	23.4	47.0	31.9	39.9	46.7	69.8	53.0	30.2
1986	24.7	14.9	19.4	20.4	34.0	43.5	37.9	37.1	51.1	62.2	35.8	22.6
1987	17.0	15.0	14.5	30.4	45.8	45.6	54.5	66.1	57.9	98.1	52.0	31.3
1988	31.7	25.1	17.6	19.7	40.5	45.2	38.3	53.3	82.8	87.3	57.9	39.7
1989	22.0	19.2	19.9	22.8	33.1	38.8	41.0	42.6	62.1	58.9	53.5	40.9
1990	20.9	11.8	19.2	21.9	36.8	40.3	43.7	43.1	50.2	60.5	60.5	58.4
1991	26.3	17.4	13.2	18.2	56.6	51.8	36.2	46.1	56.9	65.8	63.3	39.1
1992	19.7	13.8	13.6	21.5	35.2	35.4	44.7	51.0	59.5	54.6	41.5	29.2
1993	28.3	21.4	24.2	24.6	60.5	58.1	45.0	45.9	69.1	62.5	49.0	41.0
1994	27.6	15.3	13.5	16.5	37.3	46.6	51.1	44.9	54.8	59.3	59.3	44.1
1995	21.1	16.0	17.3	25.7	38.4	44.0	44.5	45.1	55.4	71.3	52.8	29.2
1996	32.0	29.4	21.7	18.7	57.0	59.2	52.7	54.2	58.7	70.7	66.9	50.9
1997	20.6	27.5	23.0	22.2	54.7	66.7	46.6	57.9	62.6	73.8	74.9	40.8
1998	18.3	12.8	18.7	25.1	51.0	59.1	50.1	51.2	61.9	77.4	69.1	82.9
1999	32.2	26.6	21.5	25.1	49.2	53.0	43.8	51.5	83.4	94.8	73.8	68.0
2000	47.9	31.9	24.3	25.5	50.6	59.6	45.1	43.4	58.2	63.4	76.0	53.0
2001	32.3	21.8	16.7	22.2	40.2	53.3	40.2	44.8	53.5	60.0	83.0	54.5
2002	34.4	18.9	14.9	18.5	61.7	50.6	52.2	66.0	57.1	64.2	76.4	46.8
2003	23.2	15.4	14.8	25.9	62.5	61.8	46.3	45.3	48.6	62.3	61.4	71.0
2004	31.4	17.3	16.6	25.2	61.9	49.1	55.9	47.3	67.5	64.5	70.2	52.3
2005	46.9	29.3	20.9	29.4	38.7	52.7	40.1	46.0	66.9	83.5	80.1	36.3
2006	31.9	24.2	31.7	25.6	36.4	44.7	48.4	40.3	44.2	50.1	44.1	31.7

*Modelo de machine learning para estimar caudales
mensuales históricos de cuencas hidrográficas en Costa Rica*

Año	Enero	Feb.	Marzo	Abril	Mayo	Junio	Julio	Ago.	Sept.	Oct.	Nov.	Dic.
2007	17.3	12.5	11.4	21.0	58.3	42.2	47.0	54.6	52.0	68.3	69.7	57.5
2008	27.5	16.6	16.0	21.9	49.1	48.1	47.4	65.0	66.2	74.6	86.2	47.7
2009	24.9	26.3	24.1	21.0	56.1	56.7	46.2	42.7	39.8	52.2	62.3	42.3
2010	28.1	28.8	33.0	30.1	61.5	48.4	41.5	37.6	61.7	64.3	79.6	68.0
2011	39.7	20.9	18.7	26.8	51.5	50.3	39.5	36.0	35.7	85.6	66.2	41.4
2012	28.1	13.6	15.1	18.3	43.1	32.0	48.5	47.4	63.1	50.8	73.8	41.0
2013	19.9	15.3	14.4	24.5	36.1	42.7	38.2	36.5	47.1	58.6	64.8	41.2
2014	25.1	15.0	13.8	33.3	32.2	48.1	65.2	60.9	77.6	75.5	69.8	65.2
2015	45.7	30.4	23.0	19.4	49.7	67.1	63.2	48.8	76.8	66.4	90.7	44.2
2016	24.2	14.5	13.9	20.3	44.5	37.7	47.6	43.8	50.0	49.9	81.2	50.3
2017	32.1	22.9	14.8	16.4	37.8	37.3	46.4	55.3	44.4	78.9	91.9	42.8
2018	47.0	29.8	19.5	20.1	51.0	62.4	51.7	63.8	53.0	59.7	49.7	30.5
2019	17.4	18.1	18.6	18.3	68.4	49.3	33.4	37.0	49.2	59.4	60.8	46.7
2020	41.2	29.9	14.5	17.2	33.3	44.9	35.2	35.6	41.2	60.6	58.6	44.1
2021	32.9	22.2	28.7	52.8	61.0	51.1	55.5	68.9	75.0	73.9	63.1	42.7

Fuente: Elaboración propia.