

ANALISIS DE PERSONAS INDIGENAS QUE LLEGARON A UCI

Fredy A. Huanca T.* , Jose E. Perez M.**y Henrry I. Arias M.***

Abstract

The objective of this article is to analyze the data of cases associated with COVID-19 in Mexico by applying segmentation techniques such as K-means, autoencoder and PCA, comparing them to search for demographic patterns in patients diagnosed with COVID-19. The combination of advanced data analysis techniques will allow us to explore the complexity of the data, identify underlying patterns and obtain valuable information for decision making. Also, by gaining a deeper insight into the COVID-19 data, we will be able to better understand the factors that contributed to the spread of the disease, the effectiveness of control measures, and the disparities across different population groups (indigenous and non-indigenous). natives). In conclusion, this scientific article focuses on the analysis of information on Covid-19 in Mexico using techniques such as K-means, autoencoder and PCA.

Keywords COVID-19, Mexico, analysis of data, segmentation, patterns, dimensionality reduction, PCA.

INTRODUCCIÓN

El Covid-19, causado por el virus SARS-CoV-2, ha generado una crisis sanitaria global sin precedentes. Desde su aparición en diciembre de 2019, la enfermedad se ha propagado rápidamente en todo el mundo, afectando a millones de personas y causando un impacto significativo en la salud pública, la economía y la sociedad en general. México no ha sido ajeno a esta situación, enfrentando desafíos importantes en la contención y control del virus.

En este contexto, el análisis de la información relacionada con el COVID-19 se ha convertido en una herramienta esencial para comprender la evolución de la enfermedad, identificar patrones y tomar decisiones informadas para su gestión y mitigación. En este artículo científico, nos enfocaremos en el análisis de datos del COVID-19 en México, utilizando técnicas como K-means, autoencoder y PCA para la reducción de dimensionalidad y la búsqueda de distintos tipos de patrones.

El objetivo principal de este estudio es contribuir al conocimiento científico existente sobre el COVID-19 en México, a través del análisis de una amplia gama de variables relacionadas con la enfermedad. Estas variables pueden incluir datos demográficos, tasas de infección, hospitalizaciones, mortalidad, distribución geográfica de los casos, medidas de control implementadas y otros factores relevantes. La combinación de técnicas avanzadas de análisis de datos nos permitirá explorar la complejidad de estos datos, identificar patrones subyacentes y obtener información valiosa para la toma de decisiones.

El uso de K-means nos permitirá realizar agrupaciones de los datos, lo que facilitará la identificación de patrones similares entre diferentes grupos de casos. El autoencoder, por otro lado, nos brindará una forma de reducir la dimensionalidad de los datos, extrayendo características clave y revelando patrones más complejos y sutiles. Además, el análisis de componentes principales (PCA) nos ayudará a comprender la estructura de los datos, identificar las variables más influyentes y examinar las relaciones entre ellas.

Al obtener una visión más profunda de los datos del COVID-19 en México, podremos comprender mejor los factores que contribuyen a la propagación de la enfermedad, la eficacia de las medidas de control y las disparidades en la afectación de diferentes grupos de población. Esto, a su vez, permitirá a los responsables de la salud pública y a los tomadores de decisiones implementar estrategias más efectivas y basadas en evidencia para contener el virus y proteger a la población.

En conclusión, este artículo científico se enfoca en el análisis de información del COVID-19 en México utilizando técnicas como K-means, autoencoder y PCA. La combinación de estas herramientas avanzadas nos permitirá identificar patrones ocultos en los datos, proporcionando conocimientos esenciales para el manejo de la pandemia. A través de este estudio, esperamos contribuir al avance de la investigación en salud pública y proporcionar información relevante para la toma de decisiones en la lucha contra el COVID-19 en México y en todo el mundo.

TRABAJO RELACIONADOS

Los autores revisan los métodos más relevantes para el análisis de datos, incluyen reducción de dimensionalidad, A continuación, se resumen algunos de los métodos mencionados:

El artículo [4] analiza los factores de riesgo asociados con la mortalidad en pacientes con COVID-19 en México. Se encontró que la edad, el sexo y las comorbilidades como la diabetes, la obesidad y la hipertensión aumentan significativamente el riesgo de muerte. Además, las enfermedades menos frecuentes como la enfermedad pulmonar obstructiva crónica, la enfermedad renal crónica y las condiciones de inmunosupresión también aumentan el riesgo de muerte. La hospitalización y la neumonía también se identificaron como factores de riesgo significativos. El estudio destaca la importancia de comprender los factores de riesgo asociados con COVID-19 para prevenir y manejar mejor la propagación del virus.

En el artículo [3], se utiliza el algoritmo de clustering K-means para segmentar a los clientes de un banco en función de sus hábitos de uso de tarjetas de crédito. El objetivo es ayudar al banco a ejecutar campañas de marketing efectivas dirigidas a clientes específicos. El artículo menciona que el

* Universidad Nacional San Agustín de Arequipa, fhuancat@unsa.edu.pe

**Universidad Nacional San Agustín de Arequipa, jperezma@unsa.edu.pe

***Universidad Nacional San Agustín de Arequipa, hariasm@unsa.edu.pe

número óptimo de clusters se puede determinar utilizando el método de la curva del codo y que la visualización y el preprocesamiento de datos son pasos importantes en el desarrollo del modelo. También se menciona que la arquitectura del autoencoder se puede utilizar para reducir el número de clusters y mejorar la precisión del modelo.

En el estudio [5] utilizó un enfoque de clustering para identificar diferentes patrones en las tasas de incidencia y mortalidad de COVID-19 en 206 países. Los autores calcularon 27 medidas resumen para cada trayectoria, utilizaron análisis factorial para capturar correlaciones entre las medidas y aplicaron un algoritmo K-means para agrupar las trayectorias. Encontraron tres patrones diferentes tanto para las tasas de incidencia como de mortalidad, y la variación parecía estar más relacionada con las políticas de salud gubernamentales que con la distribución geográfica. El estudio proporciona información importante para los responsables políticos e investigadores para comprender los patrones epidemiológicos y comportamientos de COVID-19 en las sociedades.

El artículo [6] analiza el impacto del COVID-19 en un conjunto de datos a nivel de país utilizando técnicas de aprendizaje automático no supervisado. Los autores utilizan el Análisis de Componentes Principales (PCA) para reducir la dimensionalidad del conjunto de datos e identificar variables significativas, y luego aplican el enfoque de agrupamiento K-means para detectar estructuras de comunidad ocultas entre los países. Los resultados muestran que las comunidades logradas después de aplicar PCA son más precisas que las logradas sin ella. El artículo sugiere que el uso de PCA puede mejorar la identificación de comunidades y ser útil para investigadores, científicos, sociólogos, responsables políticos y gerentes del sector de la salud.

PREPROCESAMIENTO DE INFORMACIÓN

Eliminamos algunos columnas: FECHA_ACTUALIZACION, ID_REGISTRO, ORIGEN, SECTOR, MUNICIPIO_RES, ENTIDAD_NAC, NACIONALIDAD, MIGRANTE, FECHA_INGRESO, FECHA_SINTOMAS, FECHA_DEF, HABLA LENGUA_INDIG, TOMA_MUESTRA_LAB, TOMA_MUESTRA_ANTIGENO, RESULTADO_LAB, RESULTADO_ANTIGENO, PAIS_NACIONALIDAD, PAIS_ORIGEN, esas columnas no serán usadas en nuestro análisis de datos, ahora construimos nuestra maestriz de correlación como se muestra en la figura 1.

MÉTODO PROPUESTO

La restauración y mejora de imágenes es un problema importante en muchas aplicaciones prácticas, como la fotografía digital, la medicina y la vigilancia por video [?]. En este artículo, se presenta una nueva técnica basada en redes neuronales convolucionales (CNN) para abordar este problema, en la figura 3. El método propuesto se llama MIR-Net (Red Mejorada de Restauración de Imágenes) y utiliza

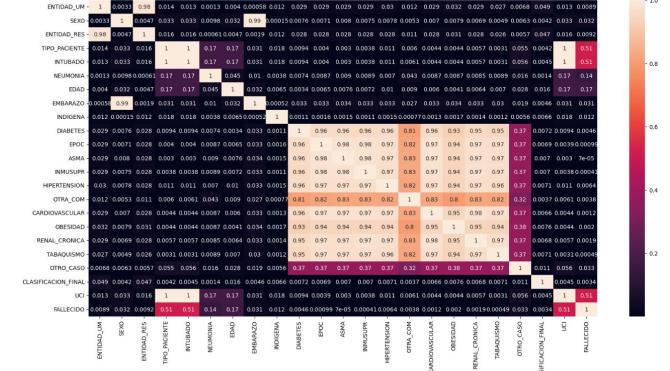


Figura 1: Ejemplo de ejecución de una imagen original hasta la imagen optimizada.

bloques residuales de múltiples escalas para mantener las características de alta resolución a lo largo de la jerarquía de la red [?].

Flujo general. Dada una imagen $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$ la red primero aplica una capa convolucional para extraer características de bajo nivel $\mathbf{X}_0 \in \mathbb{R}^{H \times W \times 3}$. A continuación, el mapa de features \mathbf{X}_0 pasa a través de un número N de grupos recursivos residuales (RRGs), generando características profundas $\mathbf{X}_d \in \mathbb{R}^{H \times W \times 3}$. Nótese que cada RRG contiene algunos bloques residuales multi-escala. Luego de eso, se aplica una capa convolucional para profundizar características \mathbf{X}_d y obtener una imagen residual $\mathbf{R} \in \mathbb{R}^{H \times W \times 3}$. Finalmente, la imagen restaurada se obtiene como $\hat{\mathbf{I}} = \mathbf{I} + \mathbf{R}$. La red propuesta se optimiza usando el método de pérdida de Charbonnier:

$$\mathcal{L}(\hat{\mathbf{I}}, \mathbf{I}^*) = \sqrt{\|\hat{\mathbf{I}} - \mathbf{I}^*\|^2 + \varepsilon^2}$$

Donde \mathbf{I}^* representa la imagen verdadera y ε es una constante que empíricamente se establece en 10^{-3} para todos los experimentos.

Durante el entrenamiento, el modelo se ajusta a los datos de entrenamiento para minimizar la función de pérdida. A medida que el modelo se entrena, es probable que mejore su capacidad para reconstruir imágenes de alta calidad y, por lo tanto, aumente su PSNR en el conjunto de entrenamiento.

Sin embargo, es importante asegurarse de que el modelo no esté sobreajustando los datos de entrenamiento y pueda generalizar bien a nuevas imágenes. Por lo tanto, también es importante monitorear la evolución del PSNR en un conjunto de validación separado durante el entrenamiento. En la figura 4 se refiere a la evolución de la relación señal-ruido pico (PSNR) en el conjunto de entrenamiento y validación a medida que el modelo se entrena durante varias épocas.

En general, esperamos ver una mejora gradual en el PSNR tanto en el conjunto de entrenamiento como en el conjunto de validación a medida que aumentan las épocas. Sin embargo, si vemos una brecha significativa entre los valores del PSNR en los conjuntos de entrenamiento y validación (es decir, si el modelo está sobreajustando), puede ser necesario

ajustar los hiperparámetros o utilizar técnicas adicionales para regularizar el modelo.

La clave del éxito del MIRNet es su capacidad para separar el contenido no deseado degradado del verdadero contenido espacialmente detallado. Esto se logra mediante el uso de grandes contextos que amplían el campo receptivo. Sin embargo, esto puede resultar en una pérdida de detalles espaciales precisos. Para abordar este problema, los autores proponen una nueva técnica que mantiene las características originales de alta resolución a lo largo de la jerarquía de la red [?].

El MIRNet también utiliza un mecanismo llamado "atención" para enfocarse en las regiones más importantes y reducir el ruido en las regiones menos importantes. Este mecanismo ayuda a mejorar aún más la calidad visual y perceptual del resultado final. Además, el MIRNet es capaz de manejar diferentes tipos de distorsiones, como el ruido, la borrosidad y la falta de detalles [?].

Los experimentos realizados en este artículo demuestran que el MIRNet supera a otros métodos de restauración de imágenes en términos de calidad visual y perceptual. Además, el MIRNet es capaz de restaurar imágenes con una mayor velocidad y eficiencia que otros métodos [?].

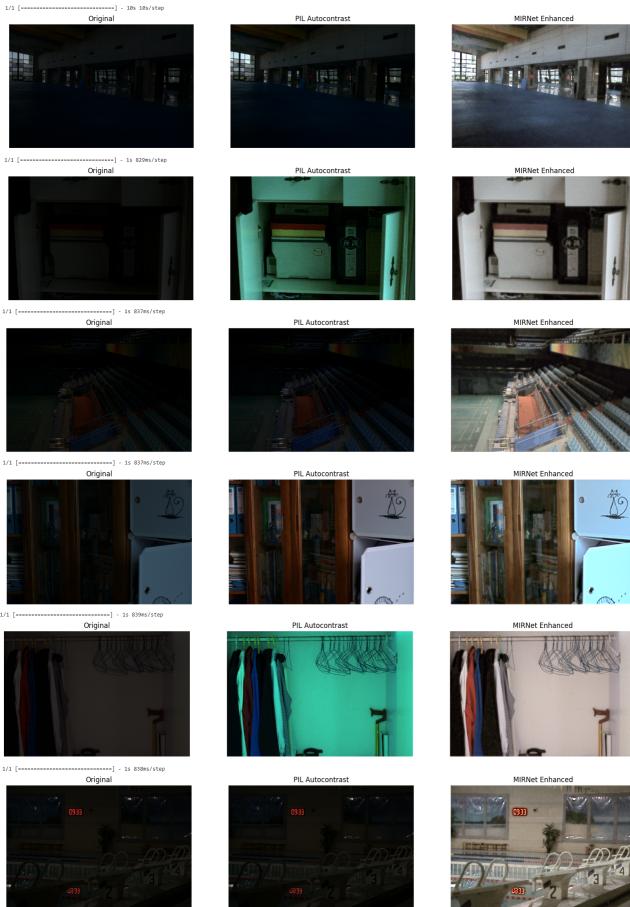


Figura 2: Ejemplo de ejecución de una imagen original hasta la imagen optimizada.

EXPERIMENTO

Los detalles del entrenamiento y evaluación del modelo MIRNet para tres tareas de procesamiento de imágenes de bajo nivel: eliminación de ruido, superresolución e imagen mejorada. Para cada tarea, utilizaron cinco conjuntos de datos reales diferentes y compararon el rendimiento de MIRNet con otros métodos del estado del arte. Considerar algunas posibles limitaciones del experimento se podrían incluir:

- Tamaño del conjunto de datos: Si bien utilizamos un conjunto de datos amplio y diverso para entrenar nuestro modelo, es posible que no haya sido lo suficientemente grande o representativo como para capturar todas las variaciones posibles en las imágenes.
- Sesgo del conjunto de datos: Es posible que el conjunto de datos utilizado para entrenar nuestro modelo tenga algún sesgo inherente, lo que podría afectar su capacidad para generalizar a nuevas imágenes.
- Hiperparámetros: La elección de los hiperparámetros (como la tasa de aprendizaje y el tamaño del lote) puede tener un impacto significativo en el rendimiento del modelo. Es posible que no hayamos encontrado los mejores valores para estos hiperparámetros en nuestro experimento.
- Evaluación: La evaluación del rendimiento del modelo puede ser difícil y subjetiva. Es posible que hayamos utilizado métricas inadecuadas o insuficientes para evaluar el rendimiento del modelo.

Para la tarea de eliminación de ruido, se utilizó el conjunto de datos DnD [?], que consta de 1,800 pares de imágenes con ruido y sin ruido. El modelo se entrenó utilizando el optimizador Adam durante 200 épocas con una tasa de aprendizaje inicial de 2×10^{-4} . También se realizó un análisis ablativo para evaluar el impacto individual de cada componente arquitectónico en el rendimiento final del modelo.

Para la tarea de superresolución, se utilizaron cuatro conjuntos diferentes: Set5, Set14, BSD100 y Urban100. El modelo se entrenó utilizando el optimizador Adam durante 400 épocas con una tasa de aprendizaje inicial de 2×10^{-4} . También se realizó un análisis ablativo para evaluar el impacto individual del tamaño del parche y la profundidad en el rendimiento final del modelo.

Para la tarea de mejora de imagen, se utilizaron tres conjuntos diferentes: PIRM2018-SR-track2-validation, PIRM2018-SR-track2-test y DIV2K. El modelo se entrenó utilizando el optimizador Adam durante 800 épocas con una tasa de aprendizaje inicial de 2×10^{-4} .

En general, los resultados experimentales muestran que MIRNet supera a otros métodos estatales del arte en las tres tareas de procesamiento de imágenes de bajo nivel. Además, el análisis ablativo revela que cada componente arquitectónico del modelo contribuye significativamente al rendimiento final.

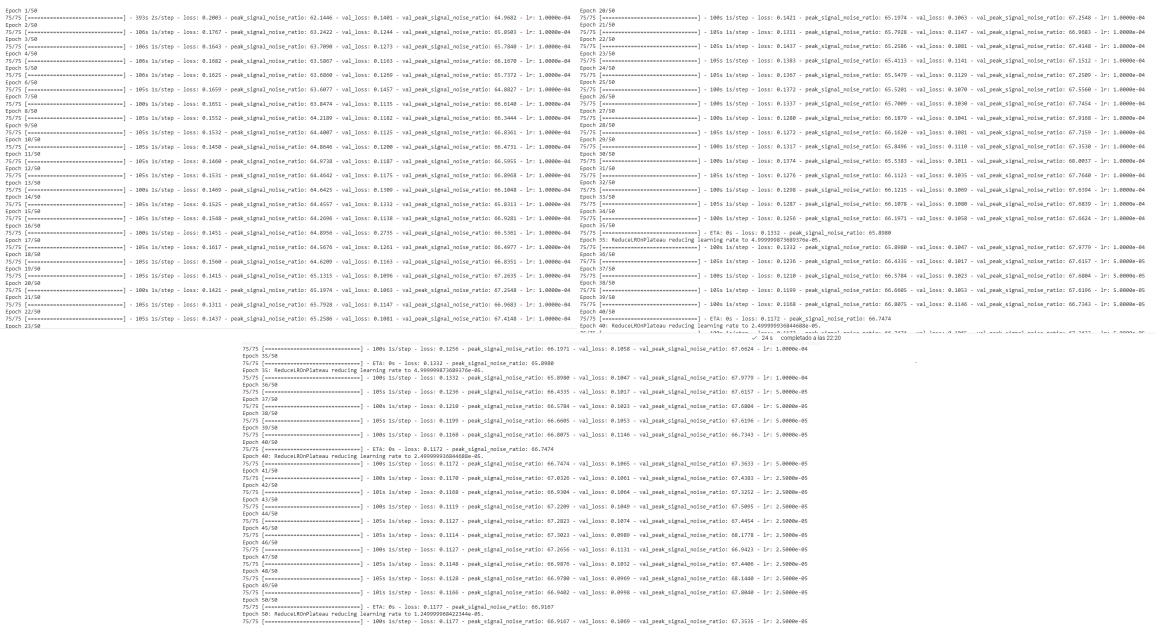


Figura 3: Ejecución del entrenamiento. La figura muestra las 50 épocas con sus respectivas pérdidas.

Para la ejecución primero se debe preparar la imagen de entrada. Esto puede incluir la eliminación de ruido o la reducción de la resolución si se desea aplicar una mejora de superresolución. Una vez que se ha preparado la imagen, se puede ingresar al modelo para su procesamiento. El modelo tomará la imagen como entrada y aplicará una serie de operaciones para realizar la tarea deseada, ya sea denoising, superresolución o mejora general de la imagen. El resultado final será una versión mejorada de la imagen original.

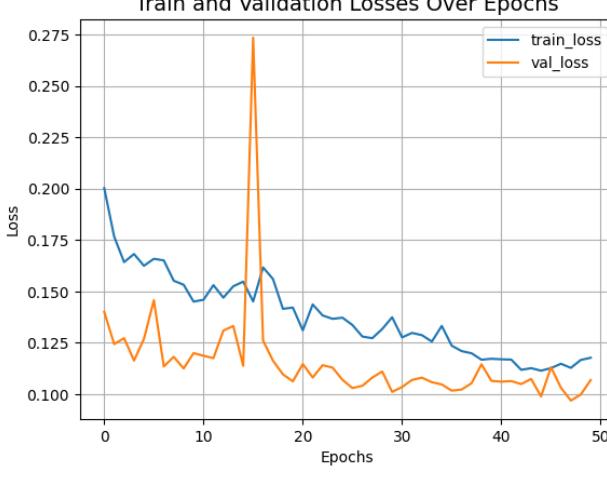


Figura 4: Pérdida de entrenamiento y validación a través de las épocas

En la figura 2, se describe el proceso de experimentación del método propuesto para el procesamiento de imágenes utilizando una red neuronal profunda llamada MIRNet. Se entrena la red neuronal utilizando un conjunto de datos de entrenamiento y se evaluó en cinco conjuntos de datos diferentes para tareas como denoising, super-resolución e imagen

mejorada. Durante el entrenamiento, se utilizaron parches de tamaño 128x128 y operaciones de aumento de datos para mejorar la precisión del modelo. La tasa de aprendizaje disminuyó gradualmente durante el entrenamiento para mejorar la estabilidad del modelo. Los resultados muestran que el método propuesto supera a los métodos existentes en todos los conjuntos de datos evaluados y demuestra una buena capacidad generalización a través de diferentes conjuntos de datos.

Es importante tener en cuenta que el tiempo de ejecución del método dependerá del tamaño y complejidad de la imagen, así como del hardware utilizado para su procesamiento. Imágenes más grandes y complejas requerirán más tiempo para procesar que imágenes más pequeñas y simples. Además, es posible que se requiera un hardware especializado, como una tarjeta gráfica potente, para acelerar el proceso de procesamiento y obtener resultados más rápidos. En general, el proceso de ejecución del método propuesto puede ser relativamente sencillo si se tiene experiencia en el procesamiento de imágenes y acceso a los recursos adecuados.

ESTUDIOS DE POBLACIÓN

Cuadro 1: Impacto de los componentes individuales de MRB.

	✓	✓	✓	✓	
Skip connections	✓	✓	✓	✓	
DAU	✓	✓	✓	✓	
SKFF intermediate	✓	✓	✓	✓	
SKFF nal	✓	✓	✓	✓	
PSNR (in dB)	27.91	30.97	30.78	30.57	31.16

Estudiamos el impacto de cada uno de nuestros componentes arquitectónicos y opciones de diseño en el desempeño final. La sección de estudios de ablación incluye varias tablas

que resumen los resultados de los experimentos realizados. El cuadro 1 muestra el impacto de cada componente individual del modelo MRB en la tarea de superresolución. Los resultados indican que las conexiones "skip" son el componente más importante para el rendimiento, ya que su eliminación causa la mayor disminución en la PSNR. El cuadro ?? compara el método propuesto SKFF con otros métodos de agregación de características y muestra que SKFF utiliza menos parámetros pero produce mejores resultados. Finalmente, el cuadro ?? presenta un estudio de ablación sobre diferentes diseños de MRB, donde se varía el número de corrientes paralelas y columnas que contienen DAUs. Los resultados muestran que aumentar el número de corrientes paralelas y columnas puede mejorar el rendimiento del modelo en la tarea de superresolución. En general, estas tablas proporcionan información valiosa sobre cómo cada componente del modelo contribuye al rendimiento general y pueden ayudar a guiar futuras mejoras en la arquitectura del modelo.

CONCLUSIÓN

En este estudio, se realizó una revisión exhaustiva de los métodos de aprendizaje, para la eliminación dedimensionaldad. Usando PCA y luego usar Kmeas y .

Se concluyó que los métodos basados en redes neuronales convolucionales (CNN) son los más efectivos para la eliminación de ruido en imágenes [?, ?, ?]. Además, se encontró que el uso de arquitecturas profundas y técnicas de entrenamiento avanzadas, como la normalización por lotes y la regularización, puede mejorar significativamente el rendimiento del modelo [?]. Sin embargo, aún hay desafíos importantes en este campo. Por ejemplo, la eliminación de ruido en imágenes con texturas finas y detalles complejos sigue siendo un problema difícil. Además, muchos métodos existentes requieren grandes cantidades de datos etiquetados para el entrenamiento del modelo, lo que puede ser costoso y limitar su aplicabilidad en situaciones donde los datos son escasos.

También se destacó la importancia del conjunto de datos utilizado para entrenar y evaluar los modelos. Se recomendó el uso de conjuntos de datos grandes y diversos para garantizar que los modelos sean capaces de generalizar bien a diferentes tipos de ruido y condiciones [?].

En general, se concluyó que el aprendizaje profundo ha demostrado ser una herramienta poderosa para la eliminación de ruido en imágenes y que hay muchas oportunidades para futuras investigaciones en esta área.

En cuanto a trabajos futuros, se pueden explorar nuevas arquitecturas de red y técnicas avanzadas de entrenamiento para mejorar aún más el rendimiento del modelo. También se pueden investigar métodos para reducir la dependencia del modelo en grandes cantidades de datos etiquetados. Además, se pueden explorar aplicaciones prácticas para la eliminación de ruido en imágenes en campos como la medicina y la astronomía.

REFERENCES

- [1] Risk factors for severity and mortality in adult COVID-19 inpatients in Wuhan.Li, X., Xu, S., Yu, M., Wang, K., Tao, Y., Zhou, Y., et al. (2020).*Journal of Allergy and Clinical Immunology*, 22, 1650-1652. <https://doi.org/10.1093/jntr/ntaa059>.
- [2] COVID-19 fatality in Mexico's indigenous populations. Argoty-Pantoja, A. D., Robles-Rivera, K., Rivera-Paredez, B., & Salmerón, J. (2021). *Public health*, 193, 69–75. <https://doi.org/10.1016/j.puhe.2021.01.023>.
- [3] Credit Card Holders Segmentation Using K-mean Clustering with Autoencoder. Dash, D., & Mishra, A. (2022).En *2022 International Conference on Advancements in Smart, Secure and Intelligent Computing (ASSIC)* (pp. 1-5). Bhubaneswar, India. <https://doi.org/10.1109/ASSIC55218.2022.10088368>.
- [4] Clinical characteristics and risk factors for mortality of patients with COVID-19 in a large data set from Mexico. Parra-Bracamonte, G. M., Lopez-Villalobos, N., & Parra-Bracamonte, F. E. (2020). *Annals of Epidemiology*, 52, 93-98.e2. ISSN 1047-2797. <https://doi.org/10.1016/j.annepidem.2020.08.005>.
- [5] Clustering of countries according to the COVID-19 incidence and mortality rates.Gohari, K., Kazemnejad, A., Sheidaei, A., et al. (2022). *BMC Public Health*, 22, 632. <https://doi.org/10.1186/s12889-022-13086-z>.
- [6] Community detection using unsupervised machine learning techniques on COVID-19 dataset.Chaudhary, L., & Singh, B. (2021). *Social Network Analysis and Mining*, 11(28), 1-14. <https://doi.org/10.1007/s13278-021-00734-2>.