

Distributed Geospatial Analytics of Agricultural Land Use from ENA Data using Spark + GIS (GeoMesa, GeoSpark, RasterFrames) in RStudio

Research Statement

Katia Daishy Ticona Casa

Resumen

Este Research Statement describe una línea de trabajo centrada en el análisis distribuido de patrones de uso de tierra agrícola en el Perú, tomando como fuente la Encuesta Nacional Agropecuaria (ENA) y la variable principal P104_SUP_ha (superficie total de parcelas en hectáreas). El enfoque combina Apache Spark con herramientas GIS distribuidas (GeoSpark/Sedona, GeoMesa, RasterFrames) integrado desde RStudio para permitir análisis espacial a escala nacional y producir evidencia útil para políticas públicas.

1. Introducción y Motivación

El uso del suelo agrícola es crítico para la seguridad alimentaria y el desarrollo rural. La ENA proporciona información detallada a nivel de unidades productivas; la variable P104_SUP_ha mide la superficie total agrícola por productor y permite analizar la estructura productiva territorial. El volumen y la complejidad espacial de estos datos exigen soluciones escalables: Apache Spark con extensiones GIS (GeoSpark/Sedona, GeoMesa, RasterFrames) ofrece un marco para análisis distribuido reproducible desde RStudio, habilitando estadísticas espaciales a gran escala y visualización interactiva [1–3].

2. Estado del Arte

La literatura reciente ha avanzado en dos ejes: (i) herramientas de procesamiento geoespacial distribuidas sobre Spark (GeoSpark/Sedona, GeoMesa, RasterFrames) [4–7], y (ii) aplicaciones de análisis espacial y aprendizaje a gran escala en dominios ambientales y agrícolas [8–13]. Revisiones modernas discuten retos y oportunidades del Big Spatial Data y su rol en agricultura de precisión [14–17].

3. Preguntas de Investigación

1. ¿Qué patrones espaciales globales y locales existen en `P104_SUP_ha` a nivel distrital/provincial en el Perú?
2. ¿Cómo se distribuye la autocorrelación espacial (Moran's I, LISA) entre ecosistemas (costa, sierra, selva)?
3. ¿Qué combinación de GeoSpark/Sedona, GeoMesa y RasterFrames optimiza el trade-off entre eficiencia y precisión?
4. ¿Cuál es la ganancia práctica (tiempo, memoria, escalabilidad) frente a análisis convencionales con `sf/spdep` en R?
5. ¿Cómo pueden estos resultados orientar políticas de ordenamiento territorial y apoyo a pequeños productores?

4. Objetivos

Objetivo general: Diseñar e implementar un pipeline distribuido para análisis espacial del uso agrícola con ENA (`P104_SUP_ha`), usando Spark + GIS desde RStudio.

Objetivos específicos:

- Ingestar y preprocesar ENA en Spark (sparklyr / interfaces).
- Asociar registros con límites administrativos y/o coordenadas de parcela.
- Implementar autocorrelación espacial distribuida (Moran's I global, LISA local).
- Integrar análisis ráster (NDVI, cobertura) con RasterFrames.
- Validar y comparar rendimiento con métodos convencionales en R.

5. Metodología propuesta

5.1. Datos

- Fuente principal: Datos ENA (2022). Variable clave: `P104_SUP_ha`.
- Capas auxiliares: límites administrativos (INEI), imágenes ráster (Sentinel/Landsat), variables climáticas.

5.2. Plataforma y herramientas

- Apache Spark cluster.
- GeoSpark / Sedona [4].
- GeoMesa (indexación espacial/temporal distribuida) [5].
- RasterFrames (manejo de ráster en Spark) [6].
- RStudio (sparklyr, reticulate).

5.3. Flujo de trabajo

1. Ingesta y limpieza de ENA.
2. Unión geográfica con shapefiles.
3. Indexación espacial con GeoMesa.
4. Cálculo distribuido de Moran's I y LISA.
5. Validación con métodos locales en R.
6. Visualización con mapas y dashboards.

6. Contribuciones esperadas

- Pipeline reproducible en RStudio que integra Spark + GIS distribuidos.
- Estudio empírico de patrones espaciales de P104_SUP_ha.
- Evaluación comparativa de eficiencia entre soluciones distribuidas vs. tradicionales.
- Evidencia para políticas agrícolas basadas en datos.

7. Líneas futuras

Series temporales multi-censo, modelos predictivos espaciales (ML distribuido), optimización de índices espaciales y arquitecturas de streaming para monitoreo en tiempo casi real [10, 12, 16].

8. Conclusión

Combinar ENA y técnicas distribuidas (Spark + GIS) desde RStudio permite abordar la brecha entre datos nacionales y capacidad analítica, ofreciendo resultados de rigor estadístico con aplicabilidad práctica para el manejo agrícola.

Referencias

- [1] Matei Zaharia et al. Apache spark: a unified engine for big data processing. *Communications of the ACM*, 59(11):56–65, 2016.
- [2] Shashi Shekhar et al. Spatial computing and big data: Challenges and opportunities. *IEEE Computer*, 48(11):1–8, 2015.
- [3] M.O. Mete et al. Geospatial big data analytics for sustainable smart governance. *ISPRS Archives / Proceedings*, 2023.
- [4] Jia Yu, Jianzhong Wu, and Mohamed Sarwat. Geospark: A cluster computing framework for processing large-scale spatial data. In *ACM SIGSPATIAL*, pages 1–4, 2015.
- [5] LocationTech GeoMesa Project. Geomesa: Distributed spatio-temporal indexing and data store (docs). <https://www.geomesa.org>, 2020. Project / software documentation.
- [6] Luke Russell et al. Rasterframes: Dataframe-centric spatiotemporal queries and map algebra over apache spark. <https://rasterframes.io>, 2019. Project / software.
- [7] Ming Tang and et al. A distributed spatial autocorrelation algorithm for big spatial data. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 14:1–13, 2020.
- [8] Zhenhua Huang and et al. A scalable framework for spatial autocorrelation analysis. *IEEE Transactions on Big Data*, 7(4):765–779, 2019.
- [9] Zhenlong Li and et al. A spark-based parallel approach for spatial autocorrelation. *ISPRS International Journal of Geo-Information*, 8(5):225, 2019.
- [10] Pankaj Kumar and et al. Scalable machine learning for geospatial data analysis. *Computers and Electronics in Agriculture*, 182:105992, 2021.
- [11] Myeongcheol Hwang and et al. Spatial data management for smart agriculture. *IEEE Big Data and Smart Computing Proceedings*, 2019.
- [12] FirstName1 LastName1, FirstName2 LastName2, and ... Big data and precision agriculture: a novel spatio-temporal semantic iot data management framework for improved interoperability. *Journal of Big Data*, 10(52):1–30, 2023. doi: 10.1186/s40537-023-00729-0.
- [13] Grupo de Investigación en Ciencia de la Información Geoespacial et al. Spatiotemporal modeling of rural agricultural land use change and area forecasts in historical time series after covid-19 pandemic, using google earth engine in peru. *Sustainability*, 16(17):7755, 2024. doi: 10.3390/su16177755.
- [14] Leilani Battle et al. Dynamic workload balancing for spatial data analytics. *VLDB Endowment*, 9(13):1509–1520, 2016.

- [15] Noel Gorelick et al. Google earth engine: Planetary-scale geospatial analysis. *Remote Sensing of Environment*, 202:18–27, 2017.
- [16] Wenwen Li and et al. Distributed deep learning for geographic data analysis. *IEEE Geoscience and Remote Sensing Magazine*, 8(3):27–42, 2020.
- [17] Anonymous. Geospatial big data: Survey and challenges. *arXiv preprint*, 2024. survey of geospatial big data technologies incl. Spark.