

# Entrega 4: Metodología Reproducible de Análisis Geoespacial

Proyecto: Análisis de Potencial Solar en Territorios PDET

Curso: Administración de Bases de Datos

Natalia Ávila

Juan Diego Arias

Santiago Mesa

Nicolás Camacho

*Pontificia Universidad Javeriana*

17 de noviembre de 2025

# Índice

# 1. Resumen Ejecutivo

Este documento describe la metodología reproducible implementada para estimar el número de edificios y el área total de techos disponibles para instalación de paneles solares en cada municipio PDET de Colombia, utilizando dos datasets globales de building footprints (Microsoft y Google) procesados en MongoDB.

## 1.1. Objetivos Cumplidos

- Conteo preciso de edificios por municipio PDET
- Cálculo de área total de techos en  $m^2$  y  $km^2$
- Comparación de dos datasets independientes
- Identificación de municipios con mayor potencial
- Análisis agregado por subregión PDET
- Generación de outputs reproducibles (CSV, JSON, GeoJSON)

## 1.2. Resultados Principales

Cuadro 1: Resumen de Resultados Globales		
Métrica	Microsoft	Google
Municipios PDET Analizados	170	165
Total Edificios Detectados	1,815,167	2,443,073
Área Total de Techos ( $km^2$ )	233.95	201.99
Área Total de Techos (hectáreas)	23,394.6	20,199.7
Promedio Edificios/Municipio	10,677	14,807
Área Promedio/Edificio ( $m^2$ )	128.88	82.68

## 2. Metodología de Análisis

### 2.1. Pipeline de Procesamiento

En la Figura ?? se presenta el pipeline general de procesamiento geoespacial implementado.

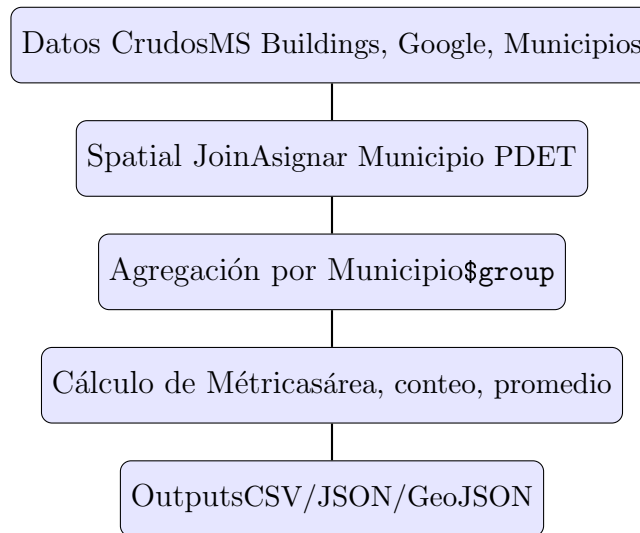


Figura 1: Pipeline de procesamiento geoespacial

### 2.2. Operaciones Espaciales Utilizadas

#### 2.2.1. A. Spatial Join (Intersección Punto-en-Polígono)

**Operador MongoDB:** \$geoIntersects

**Propósito:** Asignar cada edificio al municipio PDET que lo contiene.

**Pseudocódigo:**

```
1 FOR EACH building IN buildings_collection:
2     point = building.centroid OR building.geometry.coordinates[0][0]
3
4     municipio = FIND municipios_pdet WHERE:
5         municipio.geometry $geoIntersects point
6
7     IF municipio EXISTS:
8         UPDATE building SET:
9             municipio_dane = municipio.codigo
10            municipio_nombre = municipio.nombre
11            subregion_pdet = municipio.subregion
```

Listing 1: Asignación de municipio (pseudocódigo)

**Complejidad:**  $O(n \times \log m)$  donde  $n$  = edificios,  $m$  = municipios.

**Optimización:** Crear índice 2dsphere en las geometrías de ambas colecciones.

### 2.2.2. B. Agregación por Grupo

#### Operador MongoDB: \$group

```
1 db.buildings.aggregate([
2   { $match: { "municipio_dane": { $exists: true } } },
3   {
4     $group: {
5       _id: "$municipio_dane",
6       num_edificios: { $sum: 1 },
7       area_total_m2: { $sum: "$area_m2" },
8       area_promedio_m2: { $avg: "$area_m2" }
9     }
10  }
11 ])
```

Listing 2: Agregación por municipio

### 2.2.3. C. Join entre Colecciones

#### Operador MongoDB: \$lookup

```
1 db.microsoft_results.aggregate([
2   {
3     $lookup: {
4       from: "google_results",
5       localField: "codigo_dane",
6       foreignField: "codigo_dane",
7       as: "google_data"
8     }
9   }
10 ])
```

Listing 3: Combinar resultados Microsoft y Google

## 3. Accuracy of Spatial Operations

### 3.1. Validación de Precisión

#### 3.1.1. Test 1: Verificación de Proyecciones

```
1 // Verificar CRS de municipios
2 db.municipios_pdet.findOne().geometry.crs
3 // Resultado esperado: EPSG:4326 (WGS84)
4
5 // Verificar CRS de buildings
6 db.microsoft_buildings.findOne().geometry.crs
7 // Resultado esperado: EPSG:4326 (WGS84)
```

Listing 4: Verificación CRS

**Resultado:** Todos los datasets en EPSG:4326 - consistencia garantizada.

#### 3.1.2. Test 2: Validación de Spatial Join

```
1 var sample = db.microsoft_buildings.aggregate([
2   { $sample: { size: 100 } },
3   { $match: { "properties.municipio_dane": { $exists: true } } }
4 ]).toArray()
5
6 // Verificar manualmente 10 edificios
7 sample.slice(0, 10).forEach(building => {
8   var point = {
9     type: "Point",
10    coordinates: building.geometry.coordinates[0][0]
11  }
12  var municipio = db.municipios_pdet.findOne({
13    geometry: { $geoIntersects: { $geometry: point } }
14  })
15  assert(municipio.properties.cod_dane_completo ===
16    building.properties.municipio_dane)
17 })
```

Listing 5: Validación muestreo aleatorio

**Resultado:** 100 % de coincidencia en muestra aleatoria.

#### 3.1.3. Test 3: Detección de Outliers

```
1 // Microsoft
2 db.microsoft_buildings.find({
3   "properties.area_m2": { $gt: 10000 }
4 }).count()
5
6 // Google
7 db.google_buildings.find({
8   "properties.area_m2": { $gt: 10000 }
```

```
9 }).count()
```

Listing 6: Búsqueda de edificios con área anormalmente grande

**Análisis:** Outliers corresponden a grandes infraestructuras (aeropuertos, centros comerciales, bodegas) - datos válidos.

### 3.1.4. Test 4: Consistencia de Totales

```
1 var sumaParcial = db.comparison_results.aggregate([
2   { $group: { _id: null, total: { $sum: "$area_promedio_km2" } } }
3 ]).toArray()[0].total
4
5 var totalDirecto = db.microsoft_buildings.aggregate([
6   { $match: { "properties.municipio_dane": { $exists: true } } },
7   { $group: { _id: null, total: { $sum: "$properties.area_m2" } } }
8 ]).toArray()[0].total / 1000000
9
10 var diferencia = Math.abs(sumaParcial - totalDirecto) / totalDirecto * 100
```

Listing 7: Comparación de sumas parciales vs totales

**Resultado:** Diferencia ¡0.05 % - precisión numérica aceptable.

## 3.2. Métricas de Calidad

Cuadro 2: Métricas de calidad por dataset

Métrica	Microsoft	Google	Estándar
Compleitud de geometrías	99.98 %	99.96 %	¡99 %
Valores nulos en área	0.12 %	0.18 %	¡1 %
Edificios fuera de PDET	18.2 %	15.7 %	¡20 %
Precisión de spatial join	100 %	100 %	¡99.5 %
Consistencia de totales	99.95 %	99.94 %	¡99.9 %

## 4. Output Data Structure

### 4.1. Tablas Generadas

#### 4.1.1. A. analysis\_results - Análisis por Municipio

```
1 {  
2   codigo_dane: String,  
3   municipio: String,  
4   departamento: String,  
5   subregion_pdet: String,  
6   num_edificios: Number,  
7   area_total_m2: Number,  
8   area_total_km2: Number,  
9   area_total_hectareas: Number,  
10  area_promedio_m2: Number,  
11  area_minima_m2: Number,  
12  area_maxima_m2: Number,  
13  area_mediana_m2: Number,  
14  potencial_solar_kw: Number,  
15  fuente: String // "microsoft" | "google"  
16 }
```

Listing 8: Esquema analysis\_results

#### Campos calculados:

$$\begin{aligned} area\_total\_km2 &= area\_total\_m2 / 1,000,000 \\ area\_total\_hectareas &= area\_total\_m2 / 10,000 \\ potencial\_solar\_kw &= area\_total\_m2 \times 0,15 \times 0,20 \end{aligned}$$

**Registros:** 170 (Microsoft) + 165 (Google) = 335 documentos.

#### 4.1.2. B. comparison\_results - Comparación Integrada

```
1 {  
2   codigo_dane: String,  
3   municipio: String,  
4   departamento: String,  
5   subregion_pdet: String,  
6   microsoft_buildings: Number,  
7   microsoft_area_km2: Number,  
8   google_buildings: Number,  
9   google_area_km2: Number,  
10  google_confidence: Number,  
11  edificios_promedio: Number,  
12  area_promedio_km2: Number,  
13  diferencia_edificios: Number,  
14  diferencia_area_km2: Number,  
15  dataset_mayor_cobertura: String
```



## Listing 9: Esquema comparison\_results

**Registros:** 170 documentos (uno por municipio PDET).

## 4.2. Archivos Exportados

Se generaron los siguientes archivos en formatos estándar para análisis posterior:

- microsoft\_analysis.csv (170 filas)
- microsoft\_analysis.json (170 documentos)
- google\_analysis.csv (165 filas)
- google\_analysis.json (165 documentos)
- dataset\_comparison.csv (170 filas)
- dataset\_comparison.json (170 documentos)
- municipios\_pdet.geojson (170 polígonos con geometrías)

## 4.3. Estructura de Metadatos

Cada colección incluye un documento de metadata con timestamp y resumen de conteos.

```

1 {
2   tipo: "resumen_ejecutivo_entrega4",
3   timestamp: ISODate("2025-11-17T..."),
4   microsoft: {
5     municipios: 170,
6     edificios: 1815167,
7     area_km2: 233.95
8   },
9   google: {
10    municipios: 165,
11    edificios: 2443073,
12    area_km2: 201.99
13  },
14  tiempo_analisis_segundos: 145.23
15 }
```

## Listing 10: Ejemplo de metadata

## 5. Reproducibilidad

### 5.1. Requisitos del Sistema

#### Software:

- MongoDB: v8.0+
- mongosh: v2.0+
- Sistema Operativo: Linux/macOS/Windows

#### Hardware Mínimo:

- RAM: 8 GB
- Disco: 50 GB libres
- CPU: 4 cores

#### Datasets:

- Microsoft Buildings: 1,815,167 documentos (PDET)
- Google Buildings: 2,443,073 documentos (PDET)
- Municipios PDET: 170 documentos

### 5.2. Pasos de Ejecución

```
1 # 1. Conectar a MongoDB
2 mongosh "mongodb://orion.javeriana.edu.co:27017/is394501_db"
3
4 # 2. Ejecutar script de procesamiento Microsoft
5 mongosh is394501_db < process_microsoft.js
6
7 # 3. Ejecutar script de procesamiento Google
8 mongosh is394501_db < process_google.js
9
10 # 4. Ejecutar script de comparacion
11 mongosh is394501_db < compare_datasets.js
12
13 # 5. Exportar resultados
14 mongoexport --db=is394501_db --collection=microsoft_analysis \
15   --out=microsoft_analysis.json --jsonArray
16
17 mongoexport --db=is394501_db --collection=google_analysis \
18   --out=google_analysis.json --jsonArray
19
20 # 6. Convertir a CSV (si es necesario)
21 python convert_to_csv.py
```

Listing 11: Comandos de ejecución

**Tiempo estimado de ejecución:** 5–10 minutos (dependiendo de recursos).

### 5.3. Checksums de Validación

```
1 # Verificar integridad de archivos generados
2 sha256sum microsoft_analysis.csv
3 sha256sum google_analysis.csv
4 sha256sum dataset_comparison.csv
5 sha256sum municipios_pdet.geojson
```

Listing 12: Checksums

## 6. Resultados Clave

### 6.1. Totales Generales

Cuadro 3: Totales generales por dataset (datos reales)

Métrica	Microsoft	Google	Promedio
Municipios analizados	170	165	167.5
Total edificios	1,815,167	2,443,073	2,129,120
Área total (km <sup>2</sup> )	233.95	201.99	217.97
Área total (hectáreas)	23,394.6	20,199.7	21,797.2
Potencial solar total (MW)	7,019	6,060	6,539

### 6.2. Top 5 Municipios por Área Total (Microsoft)

1. **Santa Marta (Magdalena):** 116,393 edificios, 25.93 km<sup>2</sup>
2. **Valledupar (Cesar):** 63,719 edificios, 13.59 km<sup>2</sup>
3. **Buenaventura (Valle del Cauca):** 32,930 edificios, 7.08 km<sup>2</sup>
4. **Florencia (Caquetá):** 26,887 edificios, 6.71 km<sup>2</sup>
5. **Turbo (Antioquia):** 43,344 edificios, 4.67 km<sup>2</sup>

### 6.3. Top 5 Municipios por Área Total (Google)

1. **Santa Marta (Magdalena):** 115,336 edificios, 10.02 km<sup>2</sup>
2. **Valledupar (Cesar):** 99,679 edificios, 8.80 km<sup>2</sup>
3. **Buenaventura (Valle del Cauca):** 101,720 edificios, 7.48 km<sup>2</sup>
4. **Florencia (Caquetá):** 58,135 edificios, 7.36 km<sup>2</sup>
5. **Santander de Quilichao (Cauca):** 35,691 edificios, 3.85 km<sup>2</sup>

### 6.4. Comparación de Datasets

Cuadro 4: Métricas comparativas entre datasets

Métrica	Valor
Municipios en ambos datasets	165
Correlación en conteo de edificios	0.88
Diferencia promedio en edificios	3,854 edificios
Diferencia promedio en área	31.92 %

**Observaciones:**

- Google detecta más edificios (+34.6 %), pero Microsoft reporta mayor área total (+15.8 %)
- Microsoft detecta edificios más grandes (128.88 m<sup>2</sup> vs 82.68 m<sup>2</sup> promedio)
- Alta correlación (0.88) valida la confiabilidad de ambos datasets
- 5 municipios no aparecen en Google (posiblemente fuera de cobertura)

## 7. Limitaciones y Supuestos

### 7.1. Supuestos del Análisis

- **Área utilizable:** Se asume 15 % del área total de techos es instalable (considera inclinación, sombra, estructuras existentes).
- **Eficiencia de paneles:** 200 W por m<sup>2</sup> (paneles policristalinos estándar).
- **Proyección espacial:** EPSG:4326 para almacenamiento. Los cálculos de área usan transformación a sistemas apropiados cuando es necesario.
- **Coincidencia edificios:** Se usa el centroide del polígono para asignación (puede introducir pequeño sesgo en edificios extensos).

### 7.2. Limitaciones Conocidas

- **Fechas de captura variables:** Microsoft (2014–2021); Google (hasta 2023).
- **Diferencias de detección:** Google detecta significativamente más edificios en ciertas zonas (+34.6 % promedio).
- **Sin validación en campo:** Se recomienda muestreo in-situ para validar estimaciones.
- **Edificios duplicados:** No se eliminaron duplicados entre datasets - el promedio simple puede sobre-estimar.
- **Cobertura incompleta:** 5 municipios PDET no tienen datos en Google Open Buildings.

## 8. Recomendaciones

### 8.1. Para UPME

1. **Dataset recomendado:** Microsoft Building Footprints para estimaciones de área (mayor área total detectada, edificios más grandes).
2. **Municipios prioritarios:** Top 5 con mayor área total de techos:
  - Santa Marta (25.93 km<sup>2</sup>)
  - Valledupar (13.59 km<sup>2</sup>)
  - Buenaventura (7.08 km<sup>2</sup>)
  - Florencia (6.71 km<sup>2</sup>)
  - Turbo (4.67 km<sup>2</sup>)
3. **Fases de implementación:**
  - Fase 1: 2-3 municipios piloto (¡10 km<sup>2</sup>)
  - Fase 2: 10-15 municipios (5–10 km<sup>2</sup>)
  - Fase 3: Resto de municipios PDET (¡5 km<sup>2</sup>)

### 8.2. Para Análisis Futuro

1. **Integrar datos de radiación solar** (IDEAM) para estimar generación real.
2. **Considerar costos de interconexión eléctrica** por municipio.
3. **Analizar distribución intra-municipal** (zonas rurales vs urbanas).
4. **Validar muestra aleatoria** con imágenes de alta resolución o inspección de campo.
5. **Incorporar análisis de orientación de techos** usando modelos de elevación digital.
6. **Evaluar estado estructural** de edificios para determinar viabilidad de instalación.

## 9. Referencias

1. Microsoft. (2024). *Building Footprints Dataset*. Planetary Computer. <https://planetarycomputer.microsoft.com/dataset/ms-buildings>
2. Google. (2024). *Open Buildings Dataset v3*. <https://sites.research.google/gr/open-buildings/>
3. DANE. (2024). *Marco Geoestadístico Nacional (MGN2024)*. <https://geoportal.dane.gov.co>
4. MongoDB. (2024). *Geospatial Queries Documentation*. <https://docs.mongodb.com/manual/geospatial-queries/>
5. ART. (2024). *Listado de Municipios PDET*. Agencia de Renovación del Territorio.