

UNIVERSIDAD DE SANTIAGO DE CHILE

Facultad de Ingeniería

Departamento de Ingeniería Informática

Manual del Profesor

Clase 04: Pipeline de Desarrollo Geoespacial

Geoinformática - Segundo Semestre 2025

Prof. Francisco Parra O.

2 de septiembre de 2025

Índice

1. Resumen Ejecutivo	3
1.1. Objetivo de la Clase	3
1.2. Resultados de Aprendizaje Esperados	3
1.3. Duración y Estructura	3
2. Preparación Pre-Clase	3
2.1. Requerimientos Técnicos	3
2.2. Ambiente de Desarrollo	4
3. Marco Teórico y Conceptual	4
3.1. Principios SOLID en Pipelines Geoespaciales	4
3.1.1. Single Responsibility Principle (SRP)	4
3.1.2. Open/Closed Principle (OCP)	4
3.1.3. Liskov Substitution Principle (LSP)	5
3.1.4. Interface Segregation Principle (ISP)	5
3.1.5. Dependency Inversion Principle (DIP)	5
3.2. Teoría de Grafos Aplicada a Redes Viales	5
3.2.1. Conceptos Fundamentales	5
3.2.2. Métricas de Red	5
3.2.3. Algoritmos Clave	5
3.3. Paradigmas de Procesamiento Geoespacial	6
3.3.1. Batch Processing	6
3.3.2. Stream Processing	6
3.3.3. Micro-batch Processing	6
3.4. Testing y Calidad en Pipelines Geoespaciales	7
3.4.1. Pirámide de Testing	7
3.4.2. Tests Específicos Geoespaciales	7
3.5. Monitoreo y Observabilidad	7
3.5.1. Los Tres Pilares de la Observabilidad	7
3.5.2. Métricas Clave para Pipelines Geoespaciales	8
3.6. Desafíos Comunes y Soluciones	8
3.6.1. Desafío 1: Volumen de Datos	8
3.6.2. Desafío 2: Heterogeneidad de Datos	8
3.6.3. Desafío 3: Tiempo Real	9
3.6.4. Desafío 4: Calidad de Datos	9
3.7. Patrones Arquitectónicos para Escalabilidad	9
3.7.1. Lambda Architecture	9
3.7.2. Kappa Architecture	10
3.7.3. Event Sourcing	10
3.8. Estrategias de Escalamiento	10
3.8.1. Escalamiento Vertical	10
3.8.2. Escalamiento Horizontal	11
4. Guión Detallado de la Clase	11
4.1. Introducción (5 minutos)	11
4.2. Sección 1: Arquitectura y Principios (20 minutos)	11
4.3. Sección 2: Flujo de Trabajo Geoespacial (15 minutos)	12
4.4. Sección 2: Conexión a Fuentes de Datos Reales (20 minutos)	13
4.5. Sección 3: Análisis Espacial Aplicado y Testing (25 minutos)	15
4.6. Sección 4: Optimización, Escalabilidad y Monitoreo (20 minutos)	19

4.7. Sección 5: Deployment con Arquitectura de Microservicios (20 minutos)	23
4.8. Sección 6: Desafíos Comunes y Soluciones (10 minutos)	30
4.9. Cierre y Síntesis (10 minutos)	31
5. Material de Apoyo	33
5.1. Troubleshooting Común	33
5.2. Configuración de Ambiente	33
5.2.1. Docker Compose para desarrollo	33
5.2.2. Script de inicialización de base de datos	34
5.3. Ejercicios para Estudiantes	36
5.3.1. Ejercicio 1: Pipeline ETL Completo	36
5.3.2. Ejercicio 2: API de Geocodificación con Caché	36
5.3.3. Ejercicio 3: Dashboard de Análisis Inmobiliario	36
6. Evaluación y Rúbrica	36
6.1. Participación en Clase	36
6.2. Proyecto Pipeline (Tarea)	37
7. Notas Post-Clase	37
7.1. Seguimiento	37
7.2. Material Adicional	37
7.3. Preparación Próxima Clase	37

1. Resumen Ejecutivo

1.1. Objetivo de la Clase

Esta clase tiene como objetivo proporcionar a los estudiantes una comprensión integral de los pipelines de desarrollo geoespacial, combinando principios arquitectónicos sólidos con implementaciones prácticas. Se enfoca en construir soluciones escalables y mantenibles desde la conceptualización hasta el despliegue.

1.2. Resultados de Aprendizaje Esperados

Al finalizar la clase, los estudiantes serán capaces de:

- Aplicar principios SOLID al diseño de pipelines geoespaciales
- Comprender y aplicar teoría de grafos en redes viales
- Seleccionar paradigmas de procesamiento apropiados (batch, stream, micro-batch)
- Diseñar arquitecturas de microservicios geoespaciales
- Implementar estrategias de testing y monitoreo
- Identificar y resolver desafíos comunes en pipelines geoespaciales
- Estructurar profesionalmente un proyecto geoespacial
- Conectarse y obtener datos de APIs reales (OSM, Google Maps, IDE Chile)
- Implementar análisis espaciales complejos (clustering, interpolación)
- Optimizar consultas y operaciones para grandes volúmenes de datos
- Desplegar una aplicación geoespacial escalable

1.3. Duración y Estructura

- **Duración total:** 90 minutos
- **Formato:** Presencial con demos en vivo
- **Metodología:** Learning by doing - cada concepto se acompaña de código ejecutable

2. Preparación Pre-Clase

2.1. Requerimientos Técnicos

Importante

Es CRUCIAL verificar estos elementos antes de la clase para evitar retrasos técnicos.

- **Software instalado:**
 - Python 3.9+ con ambiente virtual preparado
 - PostgreSQL con PostGIS habilitado
 - Docker Desktop funcionando
 - VS Code o Jupyter Lab

■ Cuentas y API Keys:

- Cuenta Google Cloud con \$300 de crédito gratuito
- API Key de OpenWeatherMap (gratuita)
- Cuenta en GitHub

■ Datasets de prueba:

- Shapefile de comunas de Chile (IDE Chile)
- CSV con direcciones de Santiago para geocodificar
- GeoJSON de red de metro de Santiago

2.2. Ambiente de Desarrollo

Crear un ambiente Python con todas las dependencias:

```
1 # Crear ambiente virtual
2 python -m venv geo_env
3 source geo_env/bin/activate # Linux/Mac
4 # o
5 geo_env\Scripts\activate # Windows
6
7 # Instalar dependencias
8 pip install geopandas pandas numpy
9 pip install folium streamlit fastapi uvicorn
10 pip install osmnx networkx scikit-learn
11 pip install sqlalchemy psycopg2-binary
12 pip install python-dotenv requests
13 pip install dask[complete] dask-geopandas
14 pip install pykrige rtree
```

3. Marco Teórico y Conceptual

3.1. Principios SOLID en Pipelines Geoespaciales

3.1.1. Single Responsibility Principle (SRP)

Cada componente del pipeline debe tener una única responsabilidad bien definida:

- Módulo de geocodificación: solo convierte direcciones a coordenadas
- Servicio de routing: solo calcula rutas óptimas
- Componente de visualización: solo renderiza mapas

Ejemplo práctico: En lugar de tener una clase `GeoProcessor` que hace todo, separar en `GeoValidator`, `GeoTransformer`, y `GeoAnalyzer`.

3.1.2. Open/Closed Principle (OCP)

El sistema debe ser abierto para extensión pero cerrado para modificación:

- Usar interfaces para definir contratos
- Permitir agregar nuevas fuentes de datos sin modificar código existente
- Implementar plugins para nuevos tipos de análisis

3.1.3. Liskov Substitution Principle (LSP)

Las implementaciones deben ser intercambiables:

- Diferentes geocoders (Google, Nominatim, Mapbox) con misma interfaz
- Múltiples backends de almacenamiento (PostGIS, MongoDB, BigQuery)
- Proveedores de tiles intercambiables

3.1.4. Interface Segregation Principle (ISP)

No forzar dependencias innecesarias:

- APIs específicas para diferentes tipos de clientes
- Separar interfaces de lectura y escritura
- Métodos granulares en lugar de interfaces monolíticas

3.1.5. Dependency Inversion Principle (DIP)

Depender de abstracciones, no de implementaciones concretas:

- Inyección de dependencias para servicios
- Configuración externa (archivos YAML, variables de entorno)
- Factories para crear objetos complejos

3.2. Teoría de Grafos Aplicada a Redes Viales

3.2.1. Conceptos Fundamentales

- **Nodos (Vértices):** Representan intersecciones, puntos de interés
- **Aristas (Edges):** Representan segmentos de calle, conexiones
- **Pesos:** Distancia, tiempo de viaje, costo, tráfico
- **Dirección:** Calles de uno o doble sentido

3.2.2. Métricas de Red

- **Centralidad de grado:** Número de conexiones de un nodo
- **Centralidad de cercanía:** Distancia promedio a otros nodos
- **Centralidad de intermediación:** Frecuencia en caminos más cortos
- **Coefficiente de clustering:** Conectividad local

3.2.3. Algoritmos Clave

- **Dijkstra:** Camino más corto desde un origen
- **A*:** Dijkstra con heurística para mayor eficiencia
- **Bellman-Ford:** Maneja pesos negativos
- **Floyd-Warshall:** Todos los caminos más cortos
- **Kruskal/Prim:** Árbol de expansión mínima

3.3. Paradigmas de Procesamiento Geoespacial

3.3.1. Batch Processing

Características:

- Procesamiento de grandes volúmenes en bloques
- Alta latencia aceptable
- Optimizado para throughput

Casos de uso:

- Análisis histórico de patrones de movilidad
- Generación de reportes mensuales
- Procesamiento de imágenes satelitales

Tecnologías: Apache Spark, PostGIS batch operations, GDAL

3.3.2. Stream Processing

Características:

- Procesamiento continuo de datos en tiempo real
- Baja latencia crítica
- Optimizado para respuesta rápida

Casos de uso:

- Tracking GPS en tiempo real
- Alertas de geofencing
- Detección de anomalías espaciales

Tecnologías: Apache Kafka, Apache Flink, Redis Streams

3.3.3. Micro-batch Processing

Características:

- Híbrido entre batch y stream
- Procesamiento en ventanas de tiempo pequeñas
- Balance entre latencia y throughput

Casos de uso:

- Actualización de dashboards cada minuto
- Agregaciones temporales de sensores
- Cálculo de métricas de tráfico

Tecnologías: Spark Streaming, Storm Trident

3.4. Testing y Calidad en Pipelines Geoespaciales

3.4.1. Pirámide de Testing

- **Unit Tests (70 %):** Pruebas de funciones individuales
 - Validación de geometrías
 - Transformaciones de coordenadas
 - Cálculos de distancia
- **Integration Tests (20 %):** Pruebas de componentes conectados
 - Conexión a PostGIS
 - Llamadas a APIs externas
 - Pipeline ETL completo
- **System Tests (8 %):** Pruebas del sistema completo
 - Flujo completo de usuario
 - Rendimiento bajo carga
 - Recuperación ante fallos
- **E2E Tests (2 %):** Pruebas de extremo a extremo
 - Escenarios reales de usuario
 - Pruebas en ambiente de staging

3.4.2. Tests Específicos Geoespaciales

- **Validación de geometrías:** Polígonos cerrados, sin auto-intersecciones
- **Proyecciones correctas:** Verificar CRS antes y después
- **Topología consistente:** Nodos conectados, sin gaps
- **Precisión espacial:** Tolerancias y redondeo apropiados
- **Cobertura de área:** Verificar límites geográficos

3.5. Monitoreo y Observabilidad

3.5.1. Los Tres Pilares de la Observabilidad

- **Logs:** Eventos discretos con timestamp
 - Errores de geocodificación
 - Tiempos de respuesta de APIs
 - Validaciones fallidas
- **Métricas:** Valores numéricos agregados
 - Features procesadas por segundo
 - Latencia P50, P95, P99
 - Tasa de error por servicio
- **Traces:** Flujo de requests a través del sistema

- Path completo de una geocodificación
- Cuellos de botella en el pipeline
- Dependencias entre servicios

3.5.2. Métricas Clave para Pipelines Geoespaciales

- **Latencia de procesamiento:** Tiempo por geometría
- **Throughput:** Features/segundo, MB/segundo
- **Tasa de error:** Geometrías inválidas, timeouts
- **Uso de recursos:** CPU, RAM, almacenamiento, ancho de banda
- **Cache hit rate:** Eficiencia del cacheo espacial
- **Precisión de geocoding:** Confianza en resultados

3.6. Desafíos Comunes y Soluciones

3.6.1. Desafío 1: Volumen de Datos

Problema:

- Terabytes de imágenes satelitales
- Millones de puntos GPS por día
- Datasets que no caben en memoria

Soluciones:

- Particionamiento espacial (QuadTree, GeoHash)
- Procesamiento paralelo con Dask/Spark
- Streaming en lugar de batch
- Compresión y formatos eficientes (Parquet, COG)

3.6.2. Desafío 2: Heterogeneidad de Datos

Problema:

- Múltiples formatos (Shapefile, GeoJSON, KML, etc.)
- Diferentes sistemas de coordenadas
- Calidad variable de datos

Soluciones:

- ETL robusto con validación
- Estandarización a formato común
- Transformación automática de CRS
- Data quality scoring

3.6.3. Desafío 3: Tiempo Real

Problema:

- Actualizaciones constantes de sensores
- Baja latencia requerida ($\leq 100\text{ms}$)
- Consistencia eventual vs inmediata

Soluciones:

- Stream processing (Kafka, Flink)
- Caché distribuido (Redis, Hazelcast)
- Arquitectura CQRS
- WebSockets para push notifications

3.6.4. Desafío 4: Calidad de Datos

Problema:

- Geometrías inválidas o corruptas
- Datos faltantes o incompletos
- Duplicados y outliers

Soluciones:

- Validación automática con tolerancias
- Reparación de geometrías (`buffer(0)`)
- Imputación espacial (Kriging, IDW)
- Detección de anomalías con ML

3.7. Patrones Arquitectónicos para Escalabilidad

3.7.1. Lambda Architecture

Componentes:

- **Batch Layer:** Procesamiento completo y preciso
- **Speed Layer:** Procesamiento rápido pero aproximado
- **Serving Layer:** Combina resultados de ambas capas

Aplicación Geoespacial:

- Batch: Análisis histórico de movilidad
- Speed: Tracking en tiempo real
- Serving: Dashboard unificado

3.7.2. Kappa Architecture

Concepto:

- Todo es un stream
- Reprocesar desde el log si necesario
- Simplifica la arquitectura (una sola capa)

Aplicación Geoespacial:

- Stream de eventos de sensores IoT
- Actualizaciones de posición GPS
- Cambios en OpenStreetMap

3.7.3. Event Sourcing

Principio:

- Guardar eventos, no estados
- Estado actual = suma de todos los eventos
- Permite "time travel" auditoría completa

Aplicación Geoespacial:

- Historial de movimientos de vehículos
- Cambios en el uso del suelo
- Evolución de la red vial

3.8. Estrategias de Escalamiento

3.8.1. Escalamiento Vertical

Cuándo aplicar:

- Base de datos PostGIS principal
- Cálculos complejos que requieren todos los datos
- Operaciones que no se pueden paralelizar

Límites:

- Costo exponencial
- Límite físico del hardware
- Single point of failure

3.8.2. Escalamiento Horizontal

Cuándo aplicar:

- APIs stateless
- Procesamiento paralelo de tiles
- Cache distribuido
- Microservicios independientes

Consideraciones:

- Particionamiento de datos (sharding)
- Consistencia eventual
- Overhead de coordinación

4. Guión Detallado de la Clase

4.1. Introducción (5 minutos)

Demo en vivo

Mostrar un dashboard geoespacial en producción como ejemplo de lo que construirán.

[00:00 - 00:02] Apertura:

"Buenos días. En la clase anterior vimos los fundamentos de datos geoespaciales. Hoy vamos directo a la práctica: cómo construir una solución geoespacial real de principio a fin."

[00:02 - 00:05] Motivación:

"Varios de ustedes están trabajando en proyectos como valoración inmobiliaria o optimización de rutas. Hoy aprenderán exactamente cómo implementar estas soluciones profesionalmente. Vean este dashboard [MOSTRAR DEMO] - al final de la clase sabrán cómo construir algo así."

4.2. Sección 1: Arquitectura y Principios (20 minutos)

[00:05 - 00:10] Aplicación de SOLID:

Demo en vivo

Mostrar diagrama UML de arquitectura siguiendo principios SOLID.

."Antes de escribir código, entendamos la arquitectura. Un pipeline mal diseñado es deuda técnica desde el día uno."

Ejemplos concretos para cada principio:

- **SRP:** ."Este servicio SOLO geocodifica, nada más."
- **OCP:** ."Podemos agregar Mapbox sin tocar el código existente."
- **LSP:** ."Cualquier geocoder funciona igual: dirección entra, coordenadas salen."

- **ISP:** "El frontend no necesita saber de PostGIS."
- **DIP:** "Cambiamos de PostGIS a MongoDB solo cambiando configuración."

[00:10 - 00:15] Teoría de Grafos para Redes:

```

1 # Conceptos de grafos aplicados
2 import networkx as nx
3 import osmnx as ox
4
5 # El grafo representa la red vial
6 G = ox.graph_from_place("Las Condes, Santiago", network_type='drive')
7
8 # Nodos = intersecciones
9 print(f"Intersecciones: {len(G.nodes)}")
10
11 # Aristas = calles
12 print(f"Segmentos de calle: {len(G.edges)}")
13
14 # Métricas de red
15 centrality = nx.betweenness_centrality(G)
16 top_nodes = sorted(centrality.items(), key=lambda x: x[1], reverse=True)[:5]
17 print("Intersecciones más importantes (mayor flujo de tráfico):")
18 for node, score in top_nodes:
19     print(f"    Nodo {node}: {score:.3f}")

```

[00:15 - 00:20] Paradigmas de Procesamiento:

Comparación práctica:

- **Batch:** "Procesamos todos los datos de ayer a las 2 AM"
- **Stream:** "Procesamos cada dato cuando llega, sin esperar"
- **Micro-batch:** "Procesamos cada minuto lo que llegó"

4.3. Sección 2: Flujo de Trabajo Geoespacial (15 minutos)

[00:05 - 00:10] Pipeline completo:

Demo en vivo

Crear estructura de carpetas en tiempo real usando terminal.

```

1 # Demostración en vivo
2 mkdir proyecto_geo
3 cd proyecto_geo
4 mkdir -p data/{raw,processed,cache}
5 mkdir -p src/{etl,analysis,api,visualization}
6 mkdir -p tests config docker notebooks docs
7 touch config/config.yaml
8 touch src/__init__.py
9 touch README.md .gitignore .env

```

Explicar cada carpeta mientras se crea:

- **data/raw:** "Aquí van los datos originales, NUNCA los modificamos"
- **data/processed:** "Datos limpios y listos para análisis"
- **src/etl:** "Scripts de extracción y transformación"
- **config:** "Toda la configuración centralizada"

[00:10 - 00:15] Control de versiones con DVC:**Tip**

Si los estudiantes no conocen DVC, hacer analogía con Git LFS pero más poderoso.

```

1 # Demo DVC
2 pip install dvc
3 dvc init
4 dvc add data/comunas_santiago.gpkg
5 git add data/comunas_santiago.gpkg.dvc .gitignore
6 git commit -m "Add comunas dataset"
7
8 # Explicar el archivo .dvc generado
9 cat data/comunas_santiago.gpkg.dvc

```

Preguntas frecuentes en esta sección:

- "¿Por qué no usar Git directamente?" → Git no maneja bien archivos ¿100MB
- "¿Qué pasa con las credenciales?" → Usar .env y NUNCA commitear

4.4. Sección 2: Conexión a Fuentes de Datos Reales (20 minutos)**[00:15 - 00:20] OpenStreetMap con OSMnx:****Importante**

Esta demo requiere conexión a internet estable. Tener datos cached como backup.

```

1 # notebook: 01_osm_data.ipynb
2 import osmnx as ox
3 import geopandas as gpd
4 import matplotlib.pyplot as plt
5
6 # Configurar cache para evitar re-descargas
7 ox.config(use_cache=True, log_console=True)
8
9 # Demo 1: Obtener red vial
10 print("Descargando red vial de Las Condes...")
11 G = ox.graph_from_place("Las Condes, Santiago, Chile",
12                        network_type='drive')
13 print(f"Nodos: {len(G.nodes)}, Aristas: {len(G.edges)}")
14
15 # Convertir a GeoDataFrame
16 nodes, edges = ox.graph_to_gdfs(G)
17 edges[['name', 'highway', 'maxspeed', 'length']].head()
18
19 # Demo 2: Obtener POIs
20 print("Buscando hospitales...")
21 tags = {'amenity': 'hospital'}
22 hospitales = ox.geometries_from_place("Santiago, Chile", tags)
23 print(f"Encontrados: {len(hospitales)} hospitales")
24
25 # Visualizar
26 fig, ax = plt.subplots(figsize=(12, 8))
27 edges.plot(ax=ax, color='gray', linewidth=0.5)
28 hospitales.plot(ax=ax, color='red', markersize=100)
29 plt.title("Red vial y hospitales en Santiago")
30 plt.show()

```

[00:20 - 00:25] PostGIS - Base de datos espacial:**Demo en vivo**

Mostrar pgAdmin o DBeaver con la base de datos ya configurada.

```

1 # notebook: 02_postgis_connection.ipynb
2 from sqlalchemy import create_engine
3 import geopandas as gpd
4 import pandas as pd
5 from dotenv import load_dotenv
6 import os
7
8 load_dotenv()
9
10 # Conexión segura usando variables de entorno
11 engine = create_engine(
12     f"postgresql://{os.getenv('DB_USER')}:{os.getenv('DB_PASS')}@"
13     f"@{os.getenv('DB_HOST')}/{os.getenv('DB_NAME')}"
14 )
15
16 # Demo: Consulta espacial compleja
17 sql = """
18 WITH metro_buffer AS (
19     SELECT ST_Buffer(geom::geography, 500)::geometry as buffer
20     FROM estaciones_metro
21 )
22 SELECT
23     p.id,
24     p.direccion,
25     p.precio,
26     p.m2,
27     p.precio / p.m2 as precio_m2,
28     CASE
29         WHEN EXISTS(SELECT 1 FROM metro_buffer mb
30                     WHERE ST_Intersects(p.geom, mb.buffer))
31         THEN 'Cerca de metro'
32         ELSE 'Lejos de metro'
33     END as acceso_metro
34 FROM propiedades p
35 WHERE p.comuna = 'Las Condes'
36 ORDER BY precio_m2 DESC
37 LIMIT 10
38 """
39
40 df = gpd.read_postgis(sql, engine, geom_col='geom')
41 print(df[['direccion', 'precio_m2', 'acceso_metro']])

```

[00:25 - 00:30] APIs comerciales - Google Maps:**Tip**

Enfatizar el costo de las APIs comerciales y la importancia del rate limiting.

```

1 # notebook: 03_google_maps_api.ipynb
2 import googlemaps
3 import pandas as pd
4 from datetime import datetime
5 import time
6
7 # Cliente con API key
8 gmaps = googlemaps.Client(key=os.getenv('GOOGLE_API_KEY'))

```

```

9
10 # Demo: Geocodificación batch con manejo de errores
11 direcciones = [
12     "Av. Apoquindo 3000, Las Condes",
13     "Providencia 1234, Providencia",
14     "Estado 10, Santiago Centro"
15 ]
16
17 resultados = []
18 for direccion in direcciones:
19     try:
20         print(f"Geocodificando: {direccion}")
21         result = gmaps.geocode(f"{direccion}, Santiago, Chile")
22         if result:
23             location = result[0]['geometry']['location']
24             resultados.append({
25                 'direccion': direccion,
26                 'lat': location['lat'],
27                 'lon': location['lng'],
28                 'formatted': result[0]['formatted_address']
29             })
30             time.sleep(0.1) # Rate limiting
31     except Exception as e:
32         print(f"Error en {direccion}: {e}")
33         resultados.append({'direccion': direccion, 'error': str(e)})
34
35 df_geocoded = pd.DataFrame(resultados)
36 print(df_geocoded)

```

4.5. Sección 3: Análisis Espacial Aplicado y Testing (25 minutos)

[00:30 - 00:35] Geocodificación masiva robusta:

Importante

Esta sección es crítica para proyectos con direcciones. Dedicar tiempo a explicar el manejo de errores.

```

1 # notebook: 04_geocoding_pipeline.ipynb
2 from geopy.geocoders import Nominatim
3 from geopy.extra.rate_limiter import RateLimiter
4 import pandas as pd
5 import geopandas as gpd
6 from shapely.geometry import Point
7
8 # Cargar dataset de prueba
9 df = pd.read_csv('data/raw/direcciones_propiedades.csv')
10 print(f"Total direcciones a geocodificar: {len(df)}")
11
12 # Configurar geocoder con rate limiting automatico
13 geolocator = Nominatim(user_agent="clase_geoinformatica")
14 geocode = RateLimiter(geolocator.geocode, min_delay_seconds=1)
15
16 # Función robusta con fallback
17 def geocodificar_chile(direccion, comuna=None, retry=3):
18     """Geocodifica con mltiples intentos y variaciones"""
19     attempts = [
20         f"{direccion}, {comuna}, Santiago, Chile" if comuna else None,
21         f"{direccion}, Santiago, Chile",
22         f"{direccion}, Chile",
23         direccion

```



```

24 ]
25
26 for attempt in filter(None, attempts):
27     for i in range(retry):
28         try:
29             location = geocode(attempt)
30             if location:
31                 # Validar que est en Chile
32                 if -36 < location.latitude < -17 and -76 < location.
33 longitude < -66:
34                 return {
35                     'lat': location.latitude,
36                     'lon': location.longitude,
37                     'confianza': 'alta' if comuna in attempt else 'media
38 ',
39                     'direccion_usada': attempt
40                 }
41             except Exception as e:
42                 print(f"Intento {i+1} fall : {e}")
43                 time.sleep(2 ** i) # Exponential backoff
44
45         return {'lat': None, 'lon': None, 'confianza': 'nula', 'error': 'No
46 geocodificado'}
47
48 # Aplicar a subset para demo (en producci n usar Dask)
49 sample = df.head(20)
50 results = sample.apply(
51     lambda row: geocodificar_chile(row['direccion'], row.get('comuna')),
52     axis=1
53 )
54
55 # Crear GeoDataFrame con resultados
56 geocoded = pd.concat([sample, pd.DataFrame(list(results))], axis=1)
57 geocoded = geocoded[geocoded['lat'].notna()]
58
59 geometry = [Point(xy) for xy in zip(geocoded['lon'], geocoded['lat'])]
60 geo_df = gpd.GeoDataFrame(geocoded, geometry=geometry, crs='EPSG:4326')
61
62 # Estad sticas
63 print(f"Exitosas: {len(geocoded)}/{len(sample)} ({len(geocoded)/len(sample)
64 *100:.1f}%)")
65 print(f"Confianza alta: {(geocoded['confianza']=='alta').sum()}")
66 print(f"Confianza media: {(geocoded['confianza']=='media').sum()}")

```

[00:35 - 00:42] Áreas de influencia reales con isócronas:

Demo en vivo

Esta demo es visualmente impactante. Mostrar la diferencia entre buffer circular e isócrona real.

```

1 # notebook: 05_isochrones.ipynb
2 import osmnx as ox
3 import networkx as nx
4 import geopandas as gpd
5 from shapely.geometry import Point, Polygon
6 import matplotlib.pyplot as plt
7
8 # Punto de inter s: Hospital Salvador
9 hospital_coords = (-70.6356, -33.4569)
10 hospital = Point(hospital_coords)
11
12 # Descargar red vial alrededor del hospital

```

```

13 G = ox.graph_from_point(hospital_coords, dist=2000, network_type='walk')
14 G = ox.project_graph(G)
15
16 # Encontrar nodo m s cercano
17 hospital_node = ox.distance.nearest_nodes(G, hospital_coords[0], hospital_coords
    [1])
18
19 # Calcular is cronas para 5, 10 y 15 minutos caminando
20 walk_speed = 4.5 # km/h
21 times = [5, 10, 15] # minutos
22 isochrones = {}
23
24 for trip_time in times:
25     # Distancia m xima en metros
26     meters = walk_speed * 1000 / 60 * trip_time
27
28     # Subgrafo alcanzable
29     subgraph = nx.ego_graph(G, hospital_node, radius=meters, distance='length')
30
31     # Extraer nodos del subgrafo
32     node_points = [Point((data['x'], data['y']))
33                     for node, data in subgraph.nodes(data=True)]
34
35     # Crear pol gono convexo (simplificado)
36     if len(node_points) >= 3:
37         from shapely.ops import unary_union
38         from shapely.geometry import MultiPoint
39         isochrones[trip_time] = MultiPoint(node_points).convex_hull
40
41 # Visualizaci n comparativa
42 fig, axes = plt.subplots(1, 2, figsize=(15, 7))
43
44 # Buffer circular (m todo naive)
45 ax1 = axes[0]
46 for minutes in times:
47     radius = walk_speed * 1000 / 60 * minutes
48     circle = hospital.buffer(radius)
49     gpd.GeoSeries([circle]).plot(ax=ax1, alpha=0.3,
50                                 label=f'{minutes} min')
51 ax1.plot(*hospital.xy, 'r*', markersize=15, label='Hospital')
52 ax1.set_title('M todo Buffer Circular (Incorrecto)')
53 ax1.legend()
54
55 # Is crona real considerando calles
56 ax2 = axes[1]
57 colors = ['green', 'yellow', 'red']
58 for i, (minutes, geom) in enumerate(isochrones.items()):
59     gpd.GeoSeries([geom]).plot(ax=ax2, alpha=0.3, color=colors[i],
60                                 label=f'{minutes} min')
61 ax2.plot(*hospital.xy, 'r*', markersize=15, label='Hospital')
62 ax2.set_title('Is crona Real (Correcto)')
63 ax2.legend()
64
65 plt.suptitle('Comparaci n: Buffer vs Is crona')
66 plt.tight_layout()
67 plt.show()
68
69 print("Observen c mo la is crona real sigue la red vial")
70 print("y no asume que se puede caminar a trav s de edificios")

```

[00:42 - 00:50] Clustering espacial para detectar patrones:

Tip

Relacionar con el proyecto de valoración inmobiliaria - detectar zonas de precios similares.

```

1 # notebook: 06_spatial_clustering.ipynb
2 from sklearn.cluster import DBSCAN
3 import numpy as np
4 import geopandas as gpd
5 import matplotlib.pyplot as plt
6 from matplotlib.patches import Circle
7
8 # Cargar datos de propiedades
9 props = gpd.read_file('data/processed/propiedades_santiago.geojson')
10 props = props[props['precio'].notna()]
11
12 # Preparar features para clustering
13 # Incluir ubicaci n Y caracter sticas
14 coords = np.array([[p.x, p.y] for p in props.geometry])
15 precios_norm = props['precio'].values / props['precio'].std()
16 m2_norm = props['m2'].values / props['m2'].std()
17
18 # Feature matrix ponderada
19 X = np.column_stack([
20     coords[:, 0] * 100, # Peso alto a ubicaci n X
21     coords[:, 1] * 100, # Peso alto a ubicaci n Y
22     precios_norm,       # Precio normalizado
23     m2_norm             # Tama o normalizado
24 ])
25
26 # DBSCAN - detectar clusters de propiedades similares
27 db = DBSCAN(eps=50, min_samples=5).fit(X)
28 props['cluster'] = db.labels_
29
30 # Estad sticas por cluster
31 print(f"Clusters encontrados: {len(set(db.labels_)) - (1 if -1 in db.labels_
32     else 0)}")
33 print(f"Puntos ruido: {list(db.labels_).count(-1)}")
34
35 # An lisis por cluster
36 for cluster_id in set(db.labels_):
37     if cluster_id != -1: # Ignorar ruido
38         cluster_data = props[props['cluster'] == cluster_id]
39         print(f"\nCluster {cluster_id}:")
40         print(f"  Propiedades: {len(cluster_data)}")
41         print(f"  Precio promedio: ${cluster_data['precio'].mean():.0f}")
42         print(f"  M2 promedio: {cluster_data['m2'].mean():.1f}")
43         print(f"  Precio/M2: ${(cluster_data['precio']/cluster_data['m2']).mean
44             ():.0f}")
45
46 # Visualizaci n
47 fig, ax = plt.subplots(figsize=(12, 10))
48
49 # Plot por clusters
50 for cluster_id in set(db.labels_):
51     if cluster_id == -1:
52         color = 'gray'
53         label = 'Ruido'
54     else:
55         color = plt.cm.Set1(cluster_id)
56         label = f'Cluster {cluster_id}'
57
58     cluster_data = props[props['cluster'] == cluster_id]

```

```

57     cluster_data.plot(ax=ax, color=color, label=label,
58                       markersize=20, alpha=0.6)
59
60 ax.set_title('Clusters de Propiedades Similares')
61 ax.legend()
62 plt.show()
63
64 # Crear polígonos de mercado para cada cluster
65 from shapely.ops import unary_union
66 market_zones = []
67 for cluster_id in set(db.labels_):
68     if cluster_id != -1:
69         cluster_data = props[props['cluster'] == cluster_id]
70         zone = unary_union(cluster_data.geometry).convex_hull
71         market_zones.append({
72             'cluster': cluster_id,
73             'geometry': zone,
74             'precio_m2_promedio': (cluster_data['precio']/cluster_data['m2']).
75                                     mean()
76         })
77 zones_gdf = gpd.GeoDataFrame(market_zones)
78 zones_gdf.to_file('data/processed/zonas_mercado.geojson', driver='GeoJSON')
79 print("\nZonas de mercado exportadas a GeoJSON")

```

4.6. Sección 4: Optimización, Escalabilidad y Monitoreo (20 minutos)

[00:50 - 00:55] Testing de Componentes Geoespaciales:

Importante

Esta sección es crítica para la calidad del software geoespacial.

```

1 # tests/test_spatial_operations.py
2 import pytest
3 import geopandas as gpd
4 from shapely.geometry import Point, Polygon
5 import numpy as np
6
7 class TestSpatialValidation:
8     """Tests para validación de geometrías"""
9
10    def test_polygon_validity(self):
11        """Verificar que detectamos polígonos inválidos"""
12        # Polígono con auto-intersección (inválido)
13        coords = [(0,0), (2,2), (2,0), (0,2), (0,0)]
14        invalid_poly = Polygon(coords)
15        assert not invalid_poly.is_valid
16
17        # Reparar con buffer(0)
18        fixed_poly = invalid_poly.buffer(0)
19        assert fixed_poly.is_valid
20
21    def test_crs_transformation(self):
22        """Verificar transformación de coordenadas"""
23        # Punto en Santiago (WGS84)
24        point = Point(-70.65, -33.45)
25        gdf = gpd.GeoDataFrame([1], geometry=[point], crs='EPSG:4326')
26
27        # Transformar a UTM Zone 19S
28        gdf_utm = gdf.to_crs('EPSG:32719')

```

```

29
30     # Verificar que las coordenadas cambiaron
31     assert gdf_utm.geometry[0].x != point.x
32     assert abs(gdf_utm.geometry[0].x - 347000) < 1000 # Aprox
33
34     def test_spatial_join_performance(self):
35         """Test de rendimiento para spatial join"""
36         import time
37
38         # Crear datasets de prueba
39         points = [Point(np.random.uniform(-71, -70),
40                         np.random.uniform(-34, -33))
41                  for _ in range(1000)]
42         polygons = [Point(x, y).buffer(0.01)
43                    for x in np.linspace(-71, -70, 10)
44                    for y in np.linspace(-34, -33, 10)]
45
46         gdf_points = gpd.GeoDataFrame(geometry=points, crs='EPSG:4326')
47         gdf_polygons = gpd.GeoDataFrame(geometry=polygons, crs='EPSG:4326')
48
49         # Medir tiempo con ndice espacial
50         start = time.time()
51         result = gpd.sjoin(gdf_points, gdf_polygons, predicate='within')
52         time_with_index = time.time() - start
53
54         assert time_with_index < 1.0 # Debe ser rápido
55         assert len(result) > 0 # Debe haber matches
56
57 # Ejecutar con: pytest tests/ -v --cov=src

```

[00:55 - 01:00] Monitoreo y Observabilidad:

```

1 # src/monitoring/metrics.py
2 from prometheus_client import Counter, Histogram, Gauge
3 import time
4 import functools
5
6 # Métricas Prometheus
7 geo_requests = Counter('geo_requests_total',
8                        'Total geocoding requests',
9                        ['service', 'status'])
10 geo_latency = Histogram('geo_request_duration_seconds',
11                        'Geocoding request latency')
12 active_connections = Gauge('postgis_connections_active',
13                           'Active PostGIS connections')
14
15 def monitor_performance(func):
16     """Decorator para monitorear funciones geoespaciales"""
17     @functools.wraps(func)
18     def wrapper(*args, **kwargs):
19         start = time.time()
20         try:
21             result = func(*args, **kwargs)
22             geo_requests.labels(service=func.__name__,
23                               status='success').inc()
24             return result
25         except Exception as e:
26             geo_requests.labels(service=func.__name__,
27                               status='error').inc()
28             raise e
29         finally:
30             geo_latency.observe(time.time() - start)
31     return wrapper
32

```

```

33 @monitor_performance
34 def geocode_address(address):
35     """Geocodificar con monitoreo automatico"""
36     # Tu c digo de geocodificaci n aqu
37     pass

```

[01:00 - 01:05] Procesamiento eficiente con Dask:

Importante

Si los estudiantes no conocen Dask, explicar que es "Pandas paralelo para big data".

```

1 # notebook: 07_optimization.ipynb
2 import dask.dataframe as dd
3 import dask_geopandas as dgpd
4 import geopandas as gpd
5 import time
6
7 # Comparaci n: Pandas vs Dask
8 print("Dataset grande: 1 mill n de puntos")
9
10 # M todo tradicional (lento)
11 start = time.time()
12 gdf = gpd.read_file('data/raw/million_points.geojson')
13 gdf['buffer_100'] = gdf.buffer(100)
14 tradicional_time = time.time() - start
15 print(f"Tiempo con GeoPandas: {tradicional_time:.2f} segundos")
16
17 # M todo con Dask (r pido)
18 start = time.time()
19 ddf = dgpd.read_file('data/raw/million_points.geojson', npartitions=8)
20 ddf['buffer_100'] = ddf.geometry.buffer(100)
21 result = ddf.compute() # Ejecutar en paralelo
22 dask_time = time.time() - start
23 print(f"Tiempo con Dask: {dask_time:.2f} segundos")
24 print(f"Speedup: {tradicional_time/dask_time:.2f}x m s r pido")
25
26 # Procesamiento por chunks para memoria limitada
27 def procesar_chunk(chunk_df):
28     """Procesa un chunk de datos"""
29     chunk_df['area'] = chunk_df.geometry.area
30     chunk_df['perimetro'] = chunk_df.geometry.length
31     chunk_df['compacidad'] = 4 * np.pi * chunk_df['area'] / chunk_df['perimetro']
32     return chunk_df[chunk_df['compacidad'] > 0.7] # Filtrar formas compactas
33
34 # Procesar archivo enorme sin cargar todo en memoria
35 chunks = []
36 for chunk in pd.read_csv('huge_file.csv', chunksize=10000):
37     chunk['geometry'] = chunk.apply(lambda x: Point(x['lon'], x['lat']), axis=1)
38     chunk_gdf = gpd.GeoDataFrame(chunk, crs='EPSG:4326')
39     processed = procesar_chunk(chunk_gdf)
40     chunks.append(processed)
41     print(f"Procesado chunk con {len(processed)} registros")
42
43 result = pd.concat(chunks, ignore_index=True)
44 print(f"Total procesado: {len(result)} registros")

```

[01:05 - 01:10] Índices espaciales y caché:

```

1 # notebook: 08_spatial_index.ipynb
2 from rtree import index
3 import pickle
4 import time

```

```

5
6 class SpatialCache:
7     """Cache espacial con R-tree para b squeda ultra-r pida"""
8
9     def __init__(self):
10         # Propiedades del ndice R-tree
11         p = index.Property()
12         p.dimension = 2
13         p.variant = index.RT_Star # Mejor algoritmo
14         p.fill_factor = 0.7
15
16         self.idx = index.Index(properties=p)
17         self.cache = {}
18         self.stats = {'hits': 0, 'misses': 0}
19
20     def load_data(self, gdf):
21         """Cargar GeoDataFrame al ndice """
22         print(f"Indexando {len(gdf)} features...")
23         start = time.time()
24
25         for idx, row in gdf.iterrows():
26             bounds = row.geometry.bounds
27             self.idx.insert(idx, bounds)
28             self.cache[idx] = row.to_dict()
29
30         print(f"Indexado en {time.time()-start:.2f} segundos")
31
32     def query_bbox(self, bbox, use_cache=True):
33         """B squeda por bounding box"""
34         if use_cache:
35             self.stats['hits'] += 1
36             candidates = list(self.idx.intersection(bbox))
37             return [self.cache[i] for i in candidates]
38         else:
39             self.stats['misses'] += 1
40             # Sin ndice - b squeda lineal (lenta)
41             results = []
42             for key, item in self.cache.items():
43                 # Verificaci n manual (muy lenta)
44                 geom_bounds = item['geometry'].bounds
45                 if (bbox[0] <= geom_bounds[2] and bbox[2] >= geom_bounds[0] and
46                     bbox[1] <= geom_bounds[3] and bbox[3] >= geom_bounds[1]):
47                     results.append(item)
48             return results
49
50     def print_stats(self):
51         print(f"Cache hits: {self.stats['hits']}")
52         print(f"Cache misses: {self.stats['misses']}")
53         hit_rate = self.stats['hits'] / (self.stats['hits'] + self.stats['misses']
54         ']) * 100
55         print(f"Hit rate: {hit_rate:.1f}%")
56
57 # Demo de uso
58 cache = SpatialCache()
59 comunas = gpd.read_file('data/processed/comunas.geojson')
60 cache.load_data(comunas)
61
62 # B squeda con ndice (r pida)
63 bbox_santiago = [-70.7, -33.5, -70.5, -33.4]
64 start = time.time()
65 results_indexed = cache.query_bbox(bbox_santiago, use_cache=True)
66 time_indexed = time.time() - start

```

```

67 # B squeda sin ndice (lenta)
68 start = time.time()
69 results_linear = cache.query_bbox(bbox_santiago, use_cache=False)
70 time_linear = time.time() - start
71
72 print(f"\nResultados encontrados: {len(results_indexed)}")
73 print(f"Tiempo CON ndice : {time_indexed*1000:.2f} ms")
74 print(f"Tiempo SIN ndice : {time_linear*1000:.2f} ms")
75 print(f"Speedup: {time_linear/time_indexed:.1f}x m s r pido")
76
77 cache.print_stats()

```

4.7. Sección 5: Deployment con Arquitectura de Microservicios (20 minutos)

[01:00 - 01:07] API REST con FastAPI:

Demo en vivo

Ejecutar la API y probarla con Postman o curl en tiempo real.

```

1 # archivo: src/api/main.py
2 from fastapi import FastAPI, HTTPException, Query
3 from fastapi.middleware.cors import CORSMiddleware
4 from pydantic import BaseModel, validator
5 import geopandas as gpd
6 from shapely.geometry import Point
7 import json
8 from typing import Optional, List
9 from datetime import datetime
10 import redis
11 import hashlib
12
13 # Inicializar FastAPI
14 app = FastAPI(
15     title="GeoAPI Inmobiliaria",
16     description="API para an lisis geoespacial de propiedades",
17     version="1.0.0"
18 )
19
20 # CORS para permitir frontend
21 app.add_middleware(
22     CORSMiddleware,
23     allow_origins=["*"],
24     allow_methods=["*"],
25     allow_headers=["*"],
26 )
27
28 # Cache Redis
29 cache = redis.Redis(host='localhost', port=6379, decode_responses=True)
30
31 # Cargar datos al iniciar
32 print("Cargando datos geoespaciales...")
33 comunas = gpd.read_file("data/processed/comunas.geojson")
34 propiedades = gpd.read_file("data/processed/propiedades.geojson")
35 print(f"Datos cargados: {len(comunas)} comunas, {len(propiedades)} propiedades")
36
37 # Modelos Pydantic para validaci n
38 class LocationRequest(BaseModel):
39     lat: float
40     lon: float
41
42     @validator('lat')

```



```

43     def validate_lat(cls, v):
44         if not -90 <= v <= 90:
45             raise ValueError('Latitud debe estar entre -90 y 90')
46         return v
47
48     @validator('lon')
49     def validate_lon(cls, v):
50         if not -180 <= v <= 180:
51             raise ValueError('Longitud debe estar entre -180 y 180')
52         return v
53
54 class PropertyFilter(BaseModel):
55     precio_min: Optional[float] = 0
56     precio_max: Optional[float] = 1e9
57     m2_min: Optional[float] = 0
58     m2_max: Optional[float] = 1000
59     comuna: Optional[str] = None
60     tipo: Optional[str] = None
61
62 # Endpoints
63 @app.get("/")
64 async def root():
65     return {
66         "message": "API Geoespacial Funcionando",
67         "endpoints": [
68             "/docs",
69             "/api/comuna",
70             "/api/propiedades/cercanas",
71             "/api/analisis/precio-m2",
72             "/api/isocronas"
73         ]
74     }
75
76 @app.post("/api/comuna")
77 async def get_comuna(location: LocationRequest):
78     """Obtener informaci n de la comuna para una ubicaci n"""
79
80     # Cache key
81     cache_key = f"comuna:{location.lat}:{location.lon}"
82     cached = cache.get(cache_key)
83     if cached:
84         return json.loads(cached)
85
86     point = Point(location.lon, location.lat)
87
88     for idx, comuna in comunas.iterrows():
89         if comuna.geometry.contains(point):
90             result = {
91                 "comuna": comuna['nombre'],
92                 "region": comuna['region'],
93                 "poblacion": int(comuna['poblacion']),
94                 "area_km2": comuna.geometry.area / 1e6,
95                 "timestamp": datetime.now().isoformat()
96             }
97
98     # Guardar en cache por 1 hora
99     cache.setex(cache_key, 3600, json.dumps(result))
100     return result
101
102     raise HTTPException(status_code=404, detail="Ubicaci n fuera de Chile")
103
104 @app.get("/api/propiedades/cercanas")
105 async def propiedades_cercanas(

```

```

106 lat: float = Query(..., description="Latitud"),
107 lon: float = Query(..., description="Longitud"),
108 radio: float = Query(500, description="Radio en metros"),
109 limit: int = Query(10, description="M ximo de resultados")
110 ):
111     """Encontrar propiedades dentro de un radio"""
112
113     point = Point(lon, lat)
114     point_utm = gpd.GeoSeries([point], crs='EPSG:4326').to_crs('EPSG:32719')[0]
115
116     # Buffer en metros
117     buffer = point_utm.buffer(radio)
118
119     # Filtrar propiedades
120     props_utm = propiedades.to_crs('EPSG:32719')
121     mask = props_utm.geometry.within(buffer)
122     cercanas = propiedades[mask].copy()
123
124     if len(cercanas) == 0:
125         return {"message": "No hay propiedades en el radio especificado", "
126         results": []}
127
128     # Calcular distancias
129     cercanas['distancia'] = props_utm[mask].geometry.distance(point_utm)
130     cercanas = cercanas.nsmallest(limit, 'distancia')
131
132     # Preparar respuesta
133     results = []
134     for idx, prop in cercanas.iterrows():
135         results.append({
136             "id": int(idx),
137             "direccion": prop.get('direccion', 'Sin direcci n'),
138             "precio": float(prop['precio']),
139             "m2": float(prop['m2']),
140             "precio_m2": float(prop['precio'] / prop['m2']),
141             "distancia": float(prop['distancia']),
142             "lat": prop.geometry.y,
143             "lon": prop.geometry.x
144         })
145
146     return {
147         "centro": {"lat": lat, "lon": lon},
148         "radio": radio,
149         "total": len(results),
150         "results": results
151     }
152
153 @app.post("/api/analisis/precio-m2")
154 async def analizar_precio_m2(filters: PropertyFilter):
155     """An lisis estad stico de precios por m2"""
156
157     # Aplicar filtros
158     filtered = propiedades.copy()
159     filtered = filtered[
160         (filtered['precio'] >= filters.precio_min) &
161         (filtered['precio'] <= filters.precio_max) &
162         (filtered['m2'] >= filters.m2_min) &
163         (filtered['m2'] <= filters.m2_max)
164     ]
165
166     if filters.comuna:
167         filtered = filtered[filtered['comuna'] == filters.comuna]

```

```

168     if filters.tipo:
169         filtered = filtered[filtered['tipo'] == filters.tipo]
170
171     if len(filtered) == 0:
172         raise HTTPException(status_code=404, detail="No hay propiedades con esos
173         filtros")
174
175     # Calcular estadísticas
176     filtered['precio_m2'] = filtered['precio'] / filtered['m2']
177
178     return {
179         "total_propiedades": len(filtered),
180         "estadísticas": {
181             "precio_m2_promedio": float(filtered['precio_m2'].mean()),
182             "precio_m2_mediana": float(filtered['precio_m2'].median()),
183             "precio_m2_min": float(filtered['precio_m2'].min()),
184             "precio_m2_max": float(filtered['precio_m2'].max()),
185             "precio_m2_std": float(filtered['precio_m2'].std())
186         },
187         "distribucion": {
188             "q25": float(filtered['precio_m2'].quantile(0.25)),
189             "q50": float(filtered['precio_m2'].quantile(0.50)),
190             "q75": float(filtered['precio_m2'].quantile(0.75)),
191             "q90": float(filtered['precio_m2'].quantile(0.90))
192         },
193         "filtros_aplicados": filters.dict()
194     }
195
196 # Ejecutar con:
197 # uvicorn main:app --reload --host 0.0.0.0 --port 8000

```

[01:07 - 01:15] Dashboard interactivo con Streamlit:

Tip

Streamlit es perfecto para prototipos rápidos. Mostrar cómo en 50 líneas tienen un dashboard.

```

1 # archivo: src/visualization/dashboard.py
2 import streamlit as st
3 import geopandas as gpd
4 import pandas as pd
5 import folium
6 from streamlit_folium import st_folium
7 import plotly.express as px
8 import requests
9
10 st.set_page_config(
11     page_title="Dashboard Inmobiliario",
12     page_icon="🏠",
13     layout="wide"
14 )
15
16 # CSS personalizado
17 st.markdown("""
18 <style>
19 .big-font {
20     font-size:20px !important;
21     font-weight: bold;
22 }
23 </style>
24 """, unsafe_allow_html=True)
25

```

```

26 # Título
27 st.title("      Dashboard de An lisis Inmobiliario")
28 st.markdown("### An lisis geoespacial de propiedades en Santiago")
29
30 # Sidebar
31 with st.sidebar:
32     st.header("      Filtros")
33
34     # Filtros
35     comuna = st.selectbox(
36         "Comuna",
37         ["Todas", "Las Condes", "Providencia", "Vitacura", "Santiago Centro"]
38     )
39
40     precio_range = st.slider(
41         "Rango de precio (UF)",
42         min_value=1000,
43         max_value=50000,
44         value=(5000, 15000),
45         step=500,
46         format="%d UF"
47     )
48
49     m2_range = st.slider(
50         "Superficie (m )",
51         min_value=30,
52         max_value=500,
53         value=(50, 200),
54         step=10
55     )
56
57     tipo_propiedad = st.multiselect(
58         "Tipo de propiedad",
59         ["Departamento", "Casa", "Oficina"],
60         default=["Departamento", "Casa"]
61     )
62
63     st.markdown("---")
64
65     # An lisis
66     if st.button("      Actualizar An lisis"):
67         st.experimental_rerun()
68
69 # Layout principal
70 col1, col2, col3, col4 = st.columns(4)
71
72 # KPIs
73 @st.cache_data
74 def load_data():
75     """Cargar y cachear datos"""
76     props = gpd.read_file('data/processed/propiedades.geojson')
77     return props
78
79 props = load_data()
80
81 # Aplicar filtros
82 if comuna != "Todas":
83     props = props[props['comuna'] == comuna]
84
85 props = props[
86     (props['precio_uf'] >= precio_range[0]) &
87     (props['precio_uf'] <= precio_range[1]) &
88     (props['m2'] >= m2_range[0]) &

```

```

89     (props['m2'] <= m2_range[1]) &
90     (props['tipo'].isin(tipo_propiedad))
91 ]
92
93 # Métricas
94 with col1:
95     st.metric(
96         "Total Propiedades",
97         f"{len(props):,}",
98         f"{len(props) - 100:+,} vs mes anterior"
99     )
100
101 with col2:
102     precio_promedio = props['precio_uf'].mean()
103     st.metric(
104         "Precio Promedio",
105         f"{precio_promedio:,.0f} UF",
106         "+2.3% vs mes anterior"
107     )
108
109 with col3:
110     m2_promedio = props['m2'].mean()
111     st.metric(
112         "M Promedio",
113         f"{m2_promedio:.0f} m ",
114         "-1.2% vs mes anterior"
115     )
116
117 with col4:
118     precio_m2 = (props['precio_uf'] / props['m2']).mean()
119     st.metric(
120         "UF/M Promedio",
121         f"{precio_m2:.1f}",
122         "+3.1% vs mes anterior"
123     )
124
125 st.markdown("---")
126
127 # Mapa y gráficos
128 col_mapa, col_graficos = st.columns([2, 1])
129
130 with col_mapa:
131     st.markdown('<p class="big-font"> Mapa de Propiedades</p>',
132                 unsafe_allow_html=True)
133
134     # Crear mapa
135     m = folium.Map(location=[-33.45, -70.65], zoom_start=11)
136
137     # Agregar marcadores con clustering
138     from folium.plugins import MarkerCluster
139     marker_cluster = MarkerCluster().add_to(m)
140
141     for idx, row in props.iterrows():
142         folium.Marker(
143             [row.geometry.y, row.geometry.x],
144             popup=f"""
145             <b>{row['tipo']}</b><br>
146             Precio: {row['precio_uf']:,.0f} UF<br>
147             Superficie: {row['m2']:,.0f} m <br>
148             Comuna: {row['comuna']}
149             """,
150             icon=folium.Icon(
151                 color='green' if row['precio_uf'] < precio_promedio else 'red',

```

```

152         icon='home'
153     )
154     ).add_to(marker_cluster)
155
156     # Mostrar mapa
157     st_folium(m, height=400, width=None)
158
159 with col_graficos:
160     st.markdown('<p class="big-font">    An lisis </p>',
161                 unsafe_allow_html=True)
162
163     # Gr fico 1: Distribuci n de precios
164     fig1 = px.histogram(
165         props,
166         x='precio_uf',
167         nbins=30,
168         title="Distribuci n de Precios",
169         labels={'precio_uf': 'Precio (UF)', 'count': 'Cantidad'})
170
171     fig1.update_layout(height=200)
172     st.plotly_chart(fig1, use_container_width=True)
173
174     # Gr fico 2: Precio por comuna
175     precio_comuna = props.groupby('comuna')['precio_uf'].mean().sort_values()
176     fig2 = px.bar(
177         x=precio_comuna.values,
178         y=precio_comuna.index,
179         orientation='h',
180         title="Precio Promedio por Comuna",
181         labels={'x': 'Precio (UF)', 'y': 'Comuna'})
182
183     fig2.update_layout(height=200)
184     st.plotly_chart(fig2, use_container_width=True)
185
186 st.markdown("---")
187
188 # Tabla detallada
189 st.markdown('<p class="big-font">    Detalle de Propiedades </p>',
190             unsafe_allow_html=True)
191
192 # Preparar datos para tabla
193 tabla = props[['tipo', 'comuna', 'precio_uf', 'm2', 'dormitorios', 'banos']].
194     copy()
195 tabla['precio_m2'] = tabla['precio_uf'] / tabla['m2']
196
197 # Mostrar tabla con formato
198 st.dataframe(
199     tabla,
200     use_container_width=True,
201     hide_index=True,
202     column_config={
203         "precio_uf": st.column_config.NumberColumn(
204             "Precio (UF)",
205             format="%,.0f UF"
206         ),
207         "m2": st.column_config.NumberColumn(
208             "Superficie",
209             format="%d m "
210         ),
211         "precio_m2": st.column_config.NumberColumn(
212             "UF/m ",
213             format="%.1f"

```

```

214     )
215 }
216 )
217
218 # Footer
219 st.markdown("---")
220 st.caption("Dashboard actualizado en tiempo real | Datos: Portal Inmobiliario")
221
222 # Ejecutar con:
223 # streamlit run dashboard.py

```

4.8. Sección 6: Desafíos Comunes y Soluciones (10 minutos)

[01:15 - 01:20] Resolución de Problemas Comunes:

Demo en vivo

Mostrar casos reales de errores y cómo solucionarlos.

```

1 # Problema 1: Mezcla de CRS
2 def safe_spatial_join(gdf1, gdf2, **kwargs):
3     """Spatial join con validaci n de CRS"""
4     if gdf1.crs != gdf2.crs:
5         print(f"CRS mismatch: {gdf1.crs} vs {gdf2.crs}")
6         print("Transformando al CRS del primer GeoDataFrame...")
7         gdf2 = gdf2.to_crs(gdf1.crs)
8
9     return gpd.sjoin(gdf1, gdf2, **kwargs)
10
11 # Problema 2: Geometr as inv lidas
12 def clean_geometries(gdf):
13     """Limpiar y reparar geometr as problem ticas"""
14     # Detectar geometr as inv lidas
15     invalid_mask = ~gdf.geometry.is_valid
16     print(f"Geometr as inv lidas: {invalid_mask.sum()}")
17
18     if invalid_mask.any():
19         # Intentar reparar con buffer(0)
20         gdf.loc[invalid_mask, 'geometry'] = \
21             gdf.loc[invalid_mask, 'geometry'].buffer(0)
22
23         # Verificar de nuevo
24         still_invalid = ~gdf.geometry.is_valid
25         if still_invalid.any():
26             print(f"No se pudieron reparar {still_invalid.sum()} geometr as")
27             # Opci n: eliminar o usar convex hull
28             gdf.loc[still_invalid, 'geometry'] = \
29                 gdf.loc[still_invalid, 'geometry'].convex_hull
30
31     return gdf
32
33 # Problema 3: Datos fuera de memoria
34 def process_large_file(filepath, chunksize=10000):
35     """Procesar archivo grande por chunks"""
36     results = []
37
38     for chunk_df in pd.read_csv(filepath, chunksize=chunksize):
39         # Convertir a GeoDataFrame
40         geometry = [Point(xy) for xy in zip(chunk_df.lon, chunk_df.lat)]
41         chunk_gdf = gpd.GeoDataFrame(chunk_df, geometry=geometry)
42
43         # Procesar chunk

```

```

44     processed = your_processing_function(chunk_gdf)
45     results.append(processed)
46
47     print(f"Procesado chunk con {len(chunk_gdf)} registros")
48
49     return pd.concat(results, ignore_index=True)
50
51 # Problema 4: Rate limiting en APIs
52 from tenacity import retry, stop_after_attempt, wait_exponential
53
54 @retry(stop=stop_after_attempt(3),
55       wait=wait_exponential(multiplier=1, min=4, max=10))
56 def geocode_with_retry(address):
57     """Geocodificar con reintentos y backoff exponencial"""
58     try:
59         result = geocoder.geocode(address)
60         return result
61     except RateLimitError:
62         print("Rate limit alcanzado, esperando...")
63         raise
64     except Exception as e:
65         print(f"Error geocodificando {address}: {e}")
66         raise

```

Matriz de Decisión Tecnológica:

Criterio	PostGIS	MongoDB	Elasticsearch	BigQuery
Consultas espaciales				
Escalabilidad				
ACID				
Costo				
Tiempo real				
Análisis complejo				

Cuadro 1: Comparación de tecnologías para almacenamiento geoespacial

4.9. Cierre y Síntesis (10 minutos)

[01:20 - 01:25] Demo integrada del pipeline completo:

Demo en vivo

Mostrar el flujo completo: desde dato crudo hasta dashboard en producción.

```

1 # Demo del pipeline completo
2
3 # 1. Obtener datos
4 python src/etl/download_osm.py --place "Santiago" --tipo "hospitals"
5
6 # 2. Procesar y limpiar
7 python src/etl/process_data.py --input raw/hospitals.json --output processed/
8
9 # 3. Análisis espacial
10 python src/analysis/clustering.py --data processed/hospitals.geojson
11
12 # 4. Levantar API
13 uvicorn src.api.main:app --reload &
14
15 # 5. Levantar dashboard
16 streamlit run src/visualization/dashboard.py &
17

```



```
18 # 6. Probar API
19 curl http://localhost:8000/api/comuna -X POST \
20   -H "Content-Type: application/json" \
21   -d '{"lat": -33.45, "lon": -70.65}'
22
23 echo "Pipeline completo funcionando!"
```

[01:25 - 01:28] Mejores prácticas y errores comunes:

Importante

Enfatizar estos puntos - son los errores más comunes en proyectos reales.

- **Error #1:** No validar CRS

”SIEMPRE verifiquen el CRS. Mezclar WGS84 con UTM es el error más común.”

- **Error #2:** No usar índices espaciales

”Sin índices, una consulta de 1 segundo puede tomar 1 minuto.”

- **Error #3:** Cargar todo en memoria

”Usen chunks o Dask. No intenten cargar 1GB de GeoJSON en pandas.”

- **Error #4:** No cachear resultados costosos

”Geocodificar la misma dirección 100 veces es desperdiciar dinero y tiempo.”

[01:28 - 01:30] Q&A y recursos:

Tip

Dejar tiempo para preguntas específicas de sus proyectos.

Recursos esenciales para sus proyectos:

- **Datos Chile:** IDE.cl, datos.gob.cl, geoportal.cl
- **Documentación:** geopandas.org, postgis.net
- **Comunidad:** GIS StackExchange, r/gis
- **Mi email:** fparra@usach.cl para dudas específicas

5. Material de Apoyo

5.1. Troubleshooting Común

Problema	Solución
ImportError: No module named 'gdal'	<pre>conda install -c conda-forge gdal</pre> <pre># o</pre> <pre>sudo apt-get install gdal-bin python3-gdal</pre>
	CRS mismatch warning
# Siempre transformar al mismo CRS chunks o Dask	
gdf1 = gdf1.to_crs(gdf2.crs)	<pre>for chunk in pd.read_csv('big.csv', chunksize=10000):</pre> <pre> process(chunk)</pre>
Memory error con archivo grande	Geocoding rate limit
	Verificar que el PostgreSQL está corriendo
from geopy.extra.rate_limiter import RateLimiter	<pre>geocode = RateLimiter(geolocator.geocode)</pre> <pre>status = geolocator.status</pre> <pre>min_delay = 1</pre>
PostGIS connection refused	Verificar puerto y host en connection string

5.2. Configuración de Ambiente

5.2.1. Docker Compose para desarrollo

```

1 # docker-compose.yml
2 version: '3.8'
3
4 services:
5   postgis:
6     image: postgis/postgis:14-3.2
7     container_name: postgis_dev
8     environment:
9       POSTGRES_DB: geodata
10      POSTGRES_USER: geouser
11      POSTGRES_PASSWORD: geopass123
12      PGDATA: /var/lib/postgresql/data/pgdata
13     ports:
14       - "5432:5432"
15     volumes:
16       - pgdata:/var/lib/postgresql/data
17       - ./init.sql:/docker-entrypoint-initdb.d/init.sql
18     healthcheck:
19       test: ["CMD-SHELL", "pg_isready -U geouser"]
20       interval: 10s
21       timeout: 5s
22       retries: 5
23
24   redis:
25     image: redis:7-alpine
26     container_name: redis_cache
27     ports:
28       - "6379:6379"
29     command: redis-server --appendonly yes

```

```

30     volumes:
31         - redis_data:/data
32
33     jupyter:
34         build:
35             context: .
36             dockerfile: Dockerfile.jupyter
37         container_name: jupyter_geo
38         ports:
39             - "8888:8888"
40         volumes:
41             - ./notebooks:/home/jovyan/work
42             - ./data:/home/jovyan/data
43         environment:
44             JUPYTER_ENABLE_LAB: "yes"
45             GRANT_SUDO: "yes"
46         depends_on:
47             - postgis
48
49 volumes:
50     pgdata:
51     redis_data:

```

5.2.2. Script de inicialización de base de datos

```

1  -- init.sql
2  CREATE EXTENSION IF NOT EXISTS postgis;
3  CREATE EXTENSION IF NOT EXISTS postgis_topology;
4  CREATE EXTENSION IF NOT EXISTS fuzzystrmatch;
5  CREATE EXTENSION IF NOT EXISTS postgis_tiger_geocoder;
6
7  -- Tabla de comunas
8  CREATE TABLE IF NOT EXISTS comunas (
9      id SERIAL PRIMARY KEY,
10     nombre VARCHAR(100) NOT NULL,
11     region VARCHAR(100),
12     provincia VARCHAR(100),
13     poblacion INTEGER,
14     superficie DECIMAL(10,2),
15     geom GEOMETRY(Polygon, 4326)
16 );
17
18 -- ndices
19 CREATE INDEX idx_comunas_geom ON comunas USING GIST(geom);
20 CREATE INDEX idx_comunas_nombre ON comunas(nombre);
21
22 -- Tabla de propiedades
23 CREATE TABLE IF NOT EXISTS propiedades (
24     id SERIAL PRIMARY KEY,
25     direccion VARCHAR(255),
26     comuna_id INTEGER REFERENCES comunas(id),
27     tipo VARCHAR(50),
28     precio DECIMAL(12,2),
29     precio_uf DECIMAL(10,2),
30     m2 DECIMAL(8,2),
31     dormitorios INTEGER,
32     banos INTEGER,
33     estacionamientos INTEGER,
34     fecha_publicacion DATE,
35     geom GEOMETRY(Point, 4326)
36 );
37

```

```

38 -- ndices para propiedades
39 CREATE INDEX idx_propiedades_geom ON propiedades USING GIST(geom);
40 CREATE INDEX idx_propiedades_precio ON propiedades(precio);
41 CREATE INDEX idx_propiedades_comuna ON propiedades(comuna_id);
42 CREATE INDEX idx_propiedades_tipo ON propiedades(tipo);
43
44 -- Vista materializada para estadísticas por comuna
45 CREATE MATERIALIZED VIEW mv_stats_comuna AS
46 SELECT
47     c.id,
48     c.nombre,
49     COUNT(p.id) as total_propiedades,
50     AVG(p.precio_uf) as precio_promedio_uf,
51     AVG(p.m2) as m2_promedio,
52     AVG(p.precio_uf / NULLIF(p.m2, 0)) as precio_m2_promedio,
53     MIN(p.precio_uf) as precio_minimo,
54     MAX(p.precio_uf) as precio_maximo,
55     STDDEV(p.precio_uf) as precio_stddev
56 FROM comunas c
57 LEFT JOIN propiedades p ON c.id = p.comuna_id
58 GROUP BY c.id, c.nombre;
59
60 CREATE INDEX idx_mv_stats_comuna ON mv_stats_comuna(nombre);
61
62 -- Funci n para encontrar propiedades cercanas
63 CREATE OR REPLACE FUNCTION propiedades_cercanas(
64     lat FLOAT,
65     lon FLOAT,
66     radio_metros FLOAT,
67     limite INTEGER DEFAULT 10
68 )
69 RETURNS TABLE (
70     id INTEGER,
71     direccion VARCHAR,
72     precio DECIMAL,
73     distancia FLOAT
74 ) AS $$
75 BEGIN
76     RETURN QUERY
77     SELECT
78         p.id,
79         p.direccion,
80         p.precio,
81         ST_Distance(
82             p.geom::geography,
83             ST_SetSRID(ST_MakePoint(lon, lat), 4326)::geography
84         ) as distancia
85     FROM propiedades p
86     WHERE ST_DWithin(
87         p.geom::geography,
88         ST_SetSRID(ST_MakePoint(lon, lat), 4326)::geography,
89         radio_metros
90     )
91     ORDER BY distancia
92     LIMIT limite;
93 END;
94 $$ LANGUAGE plpgsql;
95
96 GRANT ALL PRIVILEGES ON ALL TABLES IN SCHEMA public TO geouser;
97 GRANT ALL PRIVILEGES ON ALL SEQUENCES IN SCHEMA public TO geouser;

```

5.3. Ejercicios para Estudiantes

5.3.1. Ejercicio 1: Pipeline ETL Completo

Implementar un pipeline que:

1. Descargue datos de hospitales desde OSM
2. Calcule áreas de servicio (isócronas de 10 min)
3. Identifique zonas sin cobertura
4. Exporte resultados a PostGIS
5. Cree visualización en mapa web

5.3.2. Ejercicio 2: API de Geocodificación con Caché

Crear una API que:

1. Reciba direcciones en formato chileno
2. Geocodifique usando múltiples servicios (fallback)
3. Implemente caché Redis con TTL
4. Valide resultados dentro de Chile
5. Retorne GeoJSON válido

5.3.3. Ejercicio 3: Dashboard de Análisis Inmobiliario

Desarrollar dashboard que muestre:

1. Mapa de calor de precios
2. Clustering de propiedades similares
3. Predicción de precios con features espaciales
4. Comparación entre comunas
5. Exportación de reportes PDF

6. Evaluación y Rúbrica

6.1. Participación en Clase

- **Excelente (7.0):** Implementa demos, hace preguntas relevantes
- **Bueno (5.5):** Sigue las demos, participa ocasionalmente
- **Regular (4.0):** Presente pero pasivo
- **Insuficiente (¡4.0):** No participa o ausente

6.2. Proyecto Pipeline (Tarea)

Criterio	Descripción	Ponderación
Estructura	Proyecto bien organizado con carpetas apropiadas	15 %
ETL	Pipeline funcional de extracción y transformación	20 %
Análisis	Implementa al menos 2 análisis espaciales complejos	25 %
Optimización	Usa índices, caché o procesamiento paralelo	15 %
API/Dashboard	Endpoint REST o visualización interactiva funcionando	20 %
Documentación	README claro con instrucciones de instalación	5 %

7. Notas Post-Clase

7.1. Seguimiento

- Crear canal Slack/Discord para dudas de implementación
- Compartir repositorio con código de los ejemplos
- Office hours para ayuda con proyectos específicos
- Grabar demos para estudiantes que falten

7.2. Material Adicional

- Video: "Building Production GeoAPIs PyCon 2024"
- Tutorial: "PostGIS Performance Tips Paul Ramsey"
- Curso: "Scalable Geospatial Analytics Coursera"
- Libro: "Geospatial Development By Example with Python"

7.3. Preparación Próxima Clase

La Clase 5 cubrirá "Análisis espacial y geoestadística". Los estudiantes deberían:

- Tener su pipeline básico funcionando
- Haber identificado datos para su proyecto
- Revisar conceptos de estadística espacial
- Instalar librerías: PySAL, scikit-gstat