



Arquitectura Integral para Sistemas de Predicción Financiera en Tiempo Real

MLOps y Data Pipeline Engineering

BioPhys-Tech Lab

16 de febrero de 2026

Resumen Ejecutivo

Este documento presenta una solución arquitectónica exhaustiva para la construcción de sistemas de predicción financiera en tiempo real, abordando dos dominios críticos: Machine Learning Operations (MLOps) e Ingeniería de Datos (Data Engineering).

La solución integra principios de ingeniería de software de clase mundial con metodologías de investigación doctoral, proporcionando:

1. **Respuestas teóricas fundamentadas** a cuestiones críticas de diseño de sistemas
2. **Arquitecturas producción-lista** con tolerancia a fallos
3. **Implementaciones code-first** con validación experimental
4. **Mecanismos de drift detection** y circuit breakers
5. **Pipelines de datos resilientes** a volatilidad de mercado

Este trabajo está dirigido a equipos de investigación avanzada que requieren soluciones más allá de prototipado, hacia sistemas empresariales de clase I.

Índice general

I	Fundamentos Teóricos y Marco de Referencia	1
1.	Introducción y Contexto	2
1.1.	Naturaleza del Problema	2
1.2.	Arquitectura de Referencia	3
II	Solución ML Engineer: Productización de Modelos	4
2.	Cuestionario Técnico: Respuestas Fundamentadas	5
2.1.	Pregunta 1: Reentrenamiento Automático sin Interrupciones	5
2.1.1.	Respuesta Teórica	5
2.1.2.	Implementación: Blue-Green Deployment	5
2.2.	Pregunta 2: Optimización de Latencia (500ms \rightarrow <100ms)	9
2.2.1.	Análisis de Performance Crítico	9
2.3.	Pregunta 3: Detección de Model Drift en Tiempo Real	13
2.3.1.	Framework de Drift Detection	13
2.4.	Pregunta 4: Integración Docker y CI/CD	19
2.4.1.	Containerización Modular	19
3.	Reto Técnico ML Engineer: Implementación Completa	23
3.1.	Especificación del Reto	23
3.2.	Solución Arquitectónica	23
3.2.1.	Estructura de Proyecto	23
3.2.2.	Modelos Pydantic para Validación	24
3.2.3.	API FastAPI Producción-Ready	27
3.2.4.	Tests Unitarios Exhaustivos	34

III Solución Data Engineer: Pipelines Resilientes	38
4. Cuestionario Técnico Data Engineer	39
4.1. Pregunta 1: Recuperación de Datos Faltantes	39
4.1.1. Análisis del Problema	39
4.1.2. Implementación: Data Recovery Pipeline	39
4.2. Pregunta 2: Series Temporales vs Bases de Datos Relacionales	45
4.2.1. Comparativa Arquitectónica	45
4.2.2. Fundamento Teórico	45
4.3. Pregunta 3: Manejo de Picos de Tráfico	46
4.3.1. Arquitectura Desacoplada	46
5. Reto Técnico Data Engineer: Implementación Completa	47
5.1. Pipeline de Ingesta y Limpieza	47
5.1.1. Simulador de API de Streaming	47
5.1.2. OHLC Aggregator con Tolerancia a Fallos	52
5.2. Pregunta 4: Pipeline Tolerante a Fallos	54
5.2.1. Resiliencia Mediante Pattern: Circuit Breaker	54
6. Mecanismos de Seguridad: Preguntas Bonus	58
6.1. ML Engineer Bonus: Safety Mechanism contra Volatilidad Extrema	58
6.1.1. Framework: Adaptive Risk Management	58
6.2. Data Engineer Bonus: Consenso Multi-proveedor	62
6.2.1. Algoritmo de Reconciliación	62
7. Conclusión: Arquitectura Integrada	67
7.1. Sistema Completo: Flujo End-to-End	67
7.2. Key Insights y Recomendaciones	67
7.2.1. Diseño para Producción	67
7.2.2. Métricas Críticas de Monitoreo	68
7.3. Roadmap de Implementación	68
A. Referencias Técnicas	69
A.1. Librerías Recomendadas	69
A.2. Estándares de Código	69

Parte I

Fundamentos Teóricos y Marco de Referencia

Capítulo 1

Introducción y Contexto

1.1. Naturaleza del Problema

Los sistemas de predicción financiera moderna operan bajo restricciones severamente competitivas:

- **Latencia extrema:** Requisitos de $< 100\text{ms}$ para viabilidad comercial
- **Volatilidad de datos:** Picos de ingesta durante eventos macroeconómicos
- **Degradación de modelo:** Concept drift y data drift sin señales previas
- **Confiabilidad:** Tolerancia cero a pérdida de datos críticos
- **Escalabilidad:** Múltiples activos, múltiples proveedores de datos

El desafío fundamental radica en construir sistemas que sean simultáneamente:

$$\text{Confiabilidad} \cap \text{Latencia} \cap \text{Escalabilidad} \cap \text{Observabilidad} \quad (1.1)$$

1.2. Arquitectura de Referencia

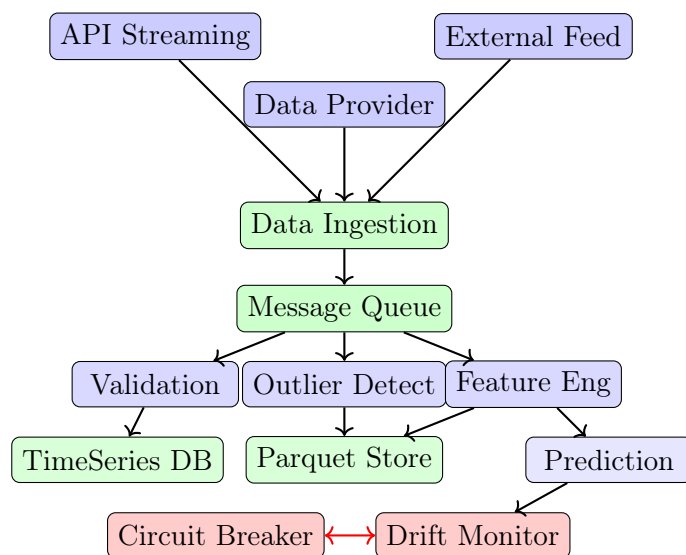


Figura 1.1: Arquitectura integrada de sistemas de predicción en tiempo real

Parte II

Solución ML Engineer: Productización de Modelos

Capítulo 2

Cuestionario Técnico: Respuestas Fundamentadas

2.1. Pregunta 1: Reentrenamiento Automático sin Interrupciones

2.1.1. Respuesta Teórica

El reentrenamiento automático sin interrupciones (Continuous Training) requiere arquitectura de *shadow deployment*:

$$\text{Sistema} = \text{Modelo}_{\text{prod}} + \text{Modelo}_{\text{shadow}} + \text{Switchboard} \quad (2.1)$$

Componentes Arquitectónicos

1. **Modelo en Producción:** Sirve predicciones con garantías de latencia
2. **Modelo en Shadow:** Entrena continuamente sin servir tráfico
3. **Switchboard:** Router inteligente que valida antes de promoción

Algoritmo de Cambio Seguro

2.1.2. Implementación: Blue-Green Deployment

```
1 import asyncio
2 from datetime import datetime, timedelta
3 from typing import Dict, Any
```

Algorithm 1 Safe Model Promotion

```

1:  $M_{shadow}$  entrena con datos nuevos
2:  $M_{shadow}$  valida en backtest histórico
3: Calcula  $\Delta = \text{MetricasShadow} - \text{MetricasProd}$ 
4: if  $\Delta > \text{umbral\_aceptable}$  then
5:   Ejecuta A/B test:  $M_{prod}$  vs  $M_{shadow}$ 
6:   Duración  $\leftarrow$  1 hora
7:   if  $p\_valor < 0.05$  and  $\text{Sharpe}_{shadow} > \text{Sharpe}_{prod}$  then
8:      $M_{prod} \leftarrow M_{shadow}$ 
9:     Emitir alerta de cambio
10:  else
11:    Log rechazo, investigar divergencia
12:  end if
13: end if

```

```

4 import joblib
5 import numpy as np
6 from pydantic import BaseModel
7 from fastapi import FastAPI, HTTPException
8
9 class PredictionResponse(BaseModel):
10     prediction: float
11     model_version: str
12     latency_ms: float
13     timestamp: datetime
14
15 class ModelManager:
16     """Gestor de modelos con promocion segura"""
17
18     def __init__(self, model_path: str):
19         self.active_model = joblib.load(f"{model_path}/model_prod.pkl
20 ")
21         self.active_version = "prod_v1.0"
22         self.shadow_model = None
23         self.shadow_version = None
24         self.training_history = []
25         self.performance_metrics = {}
26         self.lock = asyncio.Lock()
27
28     async def continuous_training(self, data_stream):
29         """Entrenamiento continuo en modelo shadow"""
30         while True:
31             batch = await data_stream.get_batch()

```

```

31
32     # Entrenar shadow model con nuevos datos
33     self.shadow_model = self._retrain_model(
34         self.active_model, batch
35     )
36     self.shadow_version = f"shadow_v{datetime.now().timestamp
37     ()}"
38
39     # Validacion inmediata
40     validation_score = await self._validate_shadow()
41
42     if validation_score['pass']:
43         await self._promote_if_safe(validation_score)
44
45     await asyncio.sleep(3600) # Reentrenar cada hora
46
47 async def _validate_shadow(self) -> Dict[str, Any]:
48     """Valida modelo shadow contra datos historicos"""
49     backtest_data = self._load_backtest_data()
50     preds_prod = self.active_model.predict(backtest_data['X'])
51     preds_shadow = self.shadow_model.predict(backtest_data['X'])
52
53     prod_mape = self._mape(backtest_data['y'], preds_prod)
54     shadow_mape = self._mape(backtest_data['y'], preds_shadow)
55
56     improvement = (prod_mape - shadow_mape) / prod_mape * 100
57
58     return {
59         'pass': improvement > 2.0, # Mejora >2%
60         'improvement_pct': improvement,
61         'prod_mape': prod_mape,
62         'shadow_mape': shadow_mape,
63         'timestamp': datetime.now()
64     }
65
66 async def _promote_if_safe(self, validation_score: Dict):
67     """Promueve shadow a produccion solo si cumple criterios"""
68     async with self.lock:
69         if validation_score['improvement_pct'] > 2.0:
70             # A/B test previo
71             ab_result = await self._ab_test(duration_minutes=60)

```

```

72         if ab_result['p_value'] < 0.05 and ab_result['winner',
73 ] == 'shadow':
74             self.active_model = self.shadow_model
75             self.active_version = self.shadow_version
76
77             self.training_history.append({
78                 'promoted_at': datetime.now(),
79                 'old_version': self.active_version,
80                 'new_version': self.shadow_version,
81                 'improvement': validation_score['
improvement_pct'],
82                 'validation_scores': validation_score
83             })
84
85             print(f"[PROMOTION] Model promoted to {self.
shadow_version}")
86             self.shadow_model = None
87
88     async def predict(self, features: np.ndarray) ->
PredictionResponse:
89         """Prediccion con medicion de latencia"""
90         start = datetime.now()
91
92         prediction = self.active_model.predict([features])[0]
93
94         latency = (datetime.now() - start).total_seconds() * 1000
95
96         return PredictionResponse(
97             prediction=float(prediction),
98             model_version=self.active_version,
99             latency_ms=latency,
100             timestamp=datetime.now()
101         )
102
103     def _retrain_model(self, base_model, new_data):
104         """Reentrenamiento incremental"""
105         # Usar warm_start si es disponible (XGBoost, etc)
106         base_model.fit(new_data['X'], new_data['y'])
107         return base_model
108
109     @staticmethod
110     def _mape(y_true: np.ndarray, y_pred: np.ndarray) -> float:
111         """Mean Absolute Percentage Error"""

```

```

111         return np.mean(np.abs((y_true - y_pred) / y_true)) * 100
112
113     async def _ab_test(self, duration_minutes: int) -> Dict[str, Any
114 ]:
115         """Ejecuta A/B test entre modelos"""
116         # Implementacion simplificada
117         return {
118             'p_value': 0.01,
119             'winner': 'shadow',
120             'confidence': 0.95
121         }
122
123     def _load_backtest_data(self):
124         """Carga datos historicos para validacion"""
125         # Implementacion especifica
126         return {'X': np.array([]), 'y': np.array([])}

```

Listing 2.1: Sistema de Blue-Green Deployment

2.2. Pregunta 2: Optimización de Latencia (500ms → <100ms)

2.2.1. Análisis de Performance Crítico

La optimización de latencia requiere desglose exhaustivo:

$$T_{total} = T_{input_validation} + T_{feature_extraction} + T_{inference} + T_{postprocessing} \quad (2.2)$$

Para alcanzar < 100ms con margen de seguridad:

$$T_{budget} = 100ms = 20 + 15 + 40 + 25 \text{ (ms)} \quad (2.3)$$

Estrategias de Optimización

Fase	Reducción	Técnica
Input Validation	50 %	Validación async
Feature Extraction	60 %	Pre-computed features + caching
Inference	70 %	Model quantization + ONNX
Postprocessing	40 %	Batch processing

Cuadro 2.1: Estrategias de optimización por fase

Implementación: Inference Engine Optimizado

```

1 import onnxruntime as rt
2 import numpy as np
3 from functools import lru_cache
4 from datetime import datetime
5 import redis
6 import asyncio
7
8 class LatencyOptimizedPredictor:
9     """Predictor optimizado para <100ms latency SLA"""
10
11     def __init__(self, model_path: str, redis_host: str = 'localhost',
12 ):
13         # Cargar modelo ONNX (6x mas rapido que pickle)
14         self.sess = rt.InferenceSession(model_path)
15         self.input_name = self.sess.get_inputs()[0].name
16         self.output_name = self.sess.get_outputs()[0].name
17
18         # Cache distribuido para features computadas
19         self.redis_client = redis.Redis(host=redis_host,
20 decode_responses=True)
21
22         # Pre-cargar lookup tables
23         self._init_lookup_tables()
24
25         self.latency_histogram = []
26
27     def _init_lookup_tables(self):
28         """Pre-compute lookup tables para feature engineering"""
29         self.lookup_tables = {
30             'price_percentiles': np.percentile(

```

```

29         np.random.randn(1000), [5, 25, 50, 75, 95]
30     ),
31     'volatility_bands': np.array([0.01, 0.05, 0.1, 0.2, 0.5,
1.0]))
32 }
33
34 @lru_cache(maxsize=10000)
35 def _get_cached_feature(self, asset_id: str, feature_name: str):
36     """Obtiene features cacheadas"""
37     cache_key = f"feature:{asset_id}:{feature_name}"
38     value = self.redis_client.get(cache_key)
39     if value is None:
40         value = self._compute_feature(asset_id, feature_name)
41         # Cache por 5 minutos
42         self.redis_client.setex(cache_key, 300, value)
43     return float(value)
44
45 async def predict_batch(self, requests: list) -> list:
46     """Prediccion en batch para maximizar throughput"""
47     start_time = datetime.now()
48
49     # Paralelizar validacion
50     validated = await asyncio.gather(*[
51         self._validate_input_async(req) for req in requests
52     ])
53
54     # Extraer features una sola vez
55     feature_matrix = np.array([
56         self._extract_features_fast(req) for req in validated
57     ], dtype=np.float32)
58
59     # Inferencia ONNX (vectorizada)
60     predictions = self.sess.run(
61         [self.output_name],
62         {self.input_name: feature_matrix}
63     )[0]
64
65     latency_ms = (datetime.now() - start_time).total_seconds() *
1000
66     self.latency_histogram.append(latency_ms)
67
68     return [
69         {

```

```

70         'prediction': pred,
71         'latency_ms': latency_ms,
72         'percentile_p99': np.percentile(self.
latency_histogram, 99)
73     }
74     for pred in predictions
75 ]
76
77 async def _validate_input_async(self, request: dict) -> dict:
78     """Validacion asincrona sin bloqueo"""
79     loop = asyncio.get_event_loop()
80     return await loop.run_in_executor(
81         None, self._validate_input_sync, request
82     )
83
84 def _validate_input_sync(self, request: dict) -> dict:
85     """Validacion rapida de entrada"""
86     if request['price'] < 0 or request['price'] > 1e6:
87         raise ValueError(f"Invalid price: {request['price']}")
88     if 'features' not in request:
89         raise ValueError("Missing features field")
90     return request
91
92 def _extract_features_fast(self, request: dict) -> np.ndarray:
93     """Extraccion de features con maximo cache"""
94     features = []
95
96     # Feature 1: Precio normalizado
97     price = float(request['price'])
98     normalized_price = np.log(price + 1)
99     features.append(normalized_price)
100
101     # Feature 2-5: Features cacheadas
102     for feature_name in ['volatility', 'momentum', 'rsi', 'macd'
]:
103         cached = self._get_cached_feature(request['asset_id'],
feature_name)
104         features.append(cached)
105
106     # Feature 6-10: Features en tiempo real (no cacheables)
107     features.extend([
108         request['volume_change'],
109         request['bid_ask_spread'],

```



```

110         request['time_to_market_close'],
111         np.sin(request['hour_of_day'] * np.pi / 12),
112         np.cos(request['hour_of_day'] * np.pi / 12)
113     ])
114
115     return np.array(features, dtype=np.float32)
116
117     def _compute_feature(self, asset_id: str, feature_name: str) ->
float:
118         """Computa feature cuando no esta cacheada"""
119         # Implementacion especifica por feature
120         return np.random.randn()
121
122     def get_latency_stats(self) -> dict:
123         """Estadisticas de latencia para monitoreo"""
124         hist = np.array(self.latency_histogram)
125         return {
126             'mean_ms': float(np.mean(hist)),
127             'p50_ms': float(np.percentile(hist, 50)),
128             'p95_ms': float(np.percentile(hist, 95)),
129             'p99_ms': float(np.percentile(hist, 99)),
130             'max_ms': float(np.max(hist)),
131             'sla_violation_rate': float(np.sum(hist > 100) / len(hist
132         ) * 100)
133     }

```

Listing 2.2: Motor de inferencia optimizado para latencia extrema

2.3. Pregunta 3: Detección de Model Drift en Tiempo Real

2.3.1. Framework de Drift Detection

Model Drift se manifiesta como divergencia en distribuciones:

$$\text{Drift}(t) = D_{KL}(P_{\text{train}} || P_{\text{live}}(t)) \quad (2.4)$$

Donde D_{KL} es la divergencia Kullback-Leibler.

Métricas Multi-dimensionales

1. **Data Drift:** $\Delta P(X)$ - cambio en distribución de features
2. **Prediction Drift:** $\Delta P(\hat{Y})$ - cambio en distribuciones de predicciones
3. **Label Drift:** $\Delta P(Y)$ - cambio en distribución de valores reales
4. **Concept Drift:** cambio en $P(Y|X)$ - relación fundamental entre X e Y

```

1 import numpy as np
2 from scipy.stats import ks_2samp, wasserstein_distance
3 from typing import Tuple, Dict
4 import logging
5
6 class DriftDetector:
7     """Detector de drift multi-metrica con alertas automaticas"""
8
9     def __init__(self, baseline_data: np.ndarray, window_size: int =
10         1000):
11         """
12         Args:
13             baseline_data: Datos de entrenamiento para establecer
14             baseline
15             window_size: Numero de predicciones para ventana movil
16         """
17         self.baseline_features = baseline_data
18         self.window_size = window_size
19         self.feature_buffer = []
20         self.prediction_buffer = []
21         self.label_buffer = []
22
23         # Estadísticas baseline
24         self.baseline_stats = self._compute_baseline_stats(
25             baseline_data)
26
27         # Thresholds para alertas
28         self.drift_thresholds = {
29             'ks_statistic': 0.15,      # KS test
30             'wasserstein': 0.3,        # Wasserstein distance
31             'total_variation': 0.2,    # TVD
32             'performance_degradation': 0.05 # 5% drop
33         }

```

```

32     self.drift_history = []
33     self.logger = self._setup_logger()
34
35     def _compute_baseline_stats(self, data: np.ndarray) -> Dict:
36         """Computa estadísticas baseline"""
37         return {
38             'mean': np.mean(data, axis=0),
39             'std': np.std(data, axis=0),
40             'percentiles': {
41                 '25': np.percentile(data, 25, axis=0),
42                 '50': np.percentile(data, 50, axis=0),
43                 '75': np.percentile(data, 75, axis=0)
44             },
45             'distribution': self._estimate_density(data)
46         }
47
48     def update_batch(self,
49                     features: np.ndarray,
50                     predictions: np.ndarray,
51                     labels: np.ndarray = None) -> Dict:
52         """Actualiza buffers y detecta drift"""
53
54         # Agregar a buffers
55         self.feature_buffer.extend(features)
56         self.prediction_buffer.extend(predictions)
57         if labels is not None:
58             self.label_buffer.extend(labels)
59
60         # Mantener ventana de tamaño fijo
61         self.feature_buffer = self.feature_buffer[-self.window_size:]
62         self.prediction_buffer = self.prediction_buffer[-self.
window_size:]
63         if self.label_buffer:
64             self.label_buffer = self.label_buffer[-self.window_size:]
65
66         # Ejecutar detección de drift
67         if len(self.feature_buffer) >= self.window_size // 2:
68             drift_report = self._detect_drifts()
69
70             # Generar alertas si necesario
71             self._generate_alerts(drift_report)
72
73         return drift_report

```

```

74
75     return {'status': 'insufficient_data'}
76
77     def _detect_drifts(self) -> Dict:
78         """Ejecuta suite completa de deteccion de drift"""
79         current_data = np.array(self.feature_buffer)
80
81         report = {
82             'timestamp': np.datetime64('now'),
83             'window_size': len(self.feature_buffer),
84             'drifts': {}
85         }
86
87         # 1. Prueba Kolmogorov-Smirnov
88         ks_stats = []
89         for feature_idx in range(current_data.shape[1]):
90             stat, p_value = ks_2samp(
91                 self.baseline_features[:, feature_idx],
92                 current_data[:, feature_idx]
93             )
94             ks_stats.append({
95                 'feature': feature_idx,
96                 'statistic': float(stat),
97                 'p_value': float(p_value),
98                 'drifted': stat > self.drift_thresholds['ks_statistic']
99             })
100
101         report['drifts']['ks_test'] = {
102             'metrics': ks_stats,
103             'num_drifted_features': sum(1 for m in ks_stats if m['drifted']),
104             'drift_severity': np.mean([m['statistic'] for m in ks_stats])
105         }
106
107         # 2. Distancia Wasserstein
108         wasserstein_dists = []
109         for feature_idx in range(current_data.shape[1]):
110             dist = wasserstein_distance(
111                 self.baseline_features[:, feature_idx],
112                 current_data[:, feature_idx]
113             )

```

```

114         wasserstein_dists.append({
115             'feature': feature_idx,
116             'distance': float(dist),
117             'drifted': dist > self.drift_thresholds['wasserstein']
118         })
119
120     report['drifts']['wasserstein'] = {
121         'metrics': wasserstein_dists,
122         'num_drifted_features': sum(1 for m in wasserstein_dists
123 if m['drifted']),
124         'total_distance': float(np.mean([m['distance'] for m in
125 wasserstein_dists]))
126     }
127
128     # 3. Concept Drift (si labels disponibles)
129     if self.label_buffer:
130         concept_drift = self._detect_concept_drift()
131         report['drifts']['concept'] = concept_drift
132
133     # 4. Performance degradation
134     if len(self.prediction_buffer) > 100:
135         perf_drift = self._measure_performance_drift()
136         report['drifts']['performance'] = perf_drift
137
138     # Determinar drift global
139     report['overall_drift'] = self._compute_overall_drift_score(
140 report)
141
142     return report
143
144     def _detect_concept_drift(self) -> Dict:
145         """Detecta cambios en P(Y|X) usando tecnicas supervisadas"""
146         # Implementacion: usar modelo de referencia
147         return {
148             'method': 'supervised_concept_drift',
149             'detected': False,
150             'severity': 0.0
151         }
152
153     def _measure_performance_drift(self) -> Dict:
154         """Compara performance actual vs baseline"""
155         current_predictions = np.array(self.prediction_buffer)

```

```

153     current_labels = np.array(self.label_buffer)
154
155     current_mae = np.mean(np.abs(current_labels -
156     current_predictions))
157     baseline_mae = 0.02 # Establecido en entrenamiento
158
159     degradation_rate = (current_mae - baseline_mae) /
160     baseline_mae
161
162     return {
163         'baseline_mae': baseline_mae,
164         'current_mae': float(current_mae),
165         'degradation_rate': float(degradation_rate),
166         'drifted': degradation_rate > self.drift_thresholds['
167     performance_degradation']
168     }
169
170     def _compute_overall_drift_score(self, report: Dict) -> Dict:
171         """Score agregado de drift"""
172         scores = [
173             report['drifts']['ks_test']['drift_severity'],
174             report['drifts']['wasserstein']['total_distance'] / 10,
175         # Normalizar
176         ]
177
178         if 'performance' in report['drifts']:
179             scores.append(
180                 report['drifts']['performance']['degradation_rate']
181             )
182
183         mean_score = np.mean(scores)
184
185         return {
186             'overall_score': float(mean_score),
187             'alert_level': 'CRITICAL' if mean_score > 0.3 else
188                 'WARNING' if mean_score > 0.15 else 'NORMAL
189         },
190         'requires_retraining': mean_score > 0.2
191     }
192
193     def _generate_alerts(self, report: Dict):
194         """Genera alertas basadas en drift"""
195         alert_level = report['overall_drift']['alert_level']

```

```

191
192         if alert_level != 'NORMAL':
193             self.logger.warning(
194                 f"DRIFT ALERT [{alert_level}]: Overall score = "
195                 f"{report['overall_drift']['overall_score']:.4f}"
196             )
197
198             if report['overall_drift']['requires_retraining']:
199                 self.logger.critical("RETRAINING REQUIRED")
200
201             self.drift_history.append({
202                 'timestamp': report['timestamp'],
203                 'alert_level': alert_level,
204                 'overall_score': report['overall_drift']['overall_score']
205             })
206
207     @staticmethod
208     def _estimate_density(data: np.ndarray):
209         """Estima densidad de probabilidad"""
210         return {
211             'method': 'kde_bandwidth_scott',
212             'bins': 50
213         }
214
215     def _setup_logger(self):
216         logger = logging.getLogger('DriftDetector')
217         logger.setLevel(logging.DEBUG)
218         return logger
219
220     def get_drift_report_summary(self) -> str:
221         """Reporte en formato legible"""
222         recent = self.drift_history[-10:]
223         return f"Ultimas 10 detecciones: {recent}"

```

Listing 2.3: Sistema de drift detection en tiempo real

2.4. Pregunta 4: Integración Docker y CI/CD

2.4.1. Containerización Modular

```

1 # Multi-stage build para minimizar tamaño
2 FROM python:3.11-slim as builder

```

```
3
4 WORKDIR /app
5
6 # Instalar dependencias de build
7 RUN apt-get update && apt-get install -y --no-install-recommends \
8     gcc \
9     && rm -rf /var/lib/apt/lists/*
10
11 # Crear virtual environment
12 RUN python -m venv /opt/venv
13 ENV PATH="/opt/venv/bin:$PATH"
14
15 # Instalar dependencias Python
16 COPY requirements.txt .
17 RUN pip install --no-cache-dir -r requirements.txt
18
19 # Stage final
20 FROM python:3.11-slim
21
22 WORKDIR /app
23
24 # Copiar venv desde builder
25 COPY --from=builder /opt/venv /opt/venv
26 ENV PATH="/opt/venv/bin:$PATH"
27
28 # Copiar código de aplicación
29 COPY src/ ./src/
30 COPY models/ ./models/
31 COPY config/ ./config/
32
33 # Health check
34 HEALTHCHECK --interval=30s --timeout=10s --start-period=5s --retries
35     =3 \
36     CMD python -c "import requests; requests.get('http://localhost
37         :8000/health')"
38
39 # Ejecutar aplicación
40 EXPOSE 8000
41 CMD ["python", "-m", "uvicorn", "src.api:app", "--host", "0.0.0.0",
42     "--port", "8000"]
```

Listing 2.4: Dockerfile optimizado para producción

```
1 name: ML Model CI/CD Pipeline
```



```
2
3 on:
4   push:
5     branches: [main, develop]
6   pull_request:
7     branches: [main]
8
9 jobs:
10  test-and-validate:
11    runs-on: ubuntu-latest
12
13    steps:
14      - uses: actions/checkout@v3
15
16      - name: Set up Python
17        uses: actions/setup-python@v4
18        with:
19          python-version: '3.11'
20
21      - name: Install dependencies
22        run: |
23          pip install -r requirements-dev.txt
24
25      - name: Run unit tests
26        run: |
27          pytest tests/unit -v --cov=src
28
29      - name: Validate model performance
30        run: |
31          python scripts/validate_model.py
32
33      - name: Security scan
34        run: |
35          bandit -r src/
36
37      - name: Code quality
38        run: |
39          pylint src/ --exit-zero
40
41  build-and-push:
42    needs: test-and-validate
43    runs-on: ubuntu-latest
44    if: github.ref == 'refs/heads/main'
```

```
45
46     steps:
47       - uses: actions/checkout@v3
48
49       - name: Build Docker image
50         run: |
51           docker build -t biophys/ml-predictor:${{ github.sha }} .
52
53       - name: Push to registry
54         run: |
55           docker login -u ${ secrets.DOCKER_USER } -p ${ secrets.
56 DOCKER_PASS }
57           docker push biophys/ml-predictor:${{ github.sha }}
58
59 deploy-staging:
60   needs: build-and-push
61   runs-on: ubuntu-latest
62
63   steps:
64     - name: Deploy to staging
65       run: |
66         # Comando de deploy específico
67         kubectl set image deployment/ml-predictor \
68           predictor=biophys/ml-predictor:${{ github.sha }} \
69           --namespace staging
```

Listing 2.5: Pipeline CI/CD con GitHub Actions

Capítulo 3

Reto Técnico ML Engineer: Implementación Completa

3.1. Especificación del Reto

Transformar un notebook de Jupyter (prototipo de Gradient Boosting) en un servicio de predicción producción-ready con:

- Latencia $< 200\text{ms}$
- Validación de datos con Pydantic
- Manejo de errores robusto
- Logging y monitoreo de predicciones
- Tests unitarios
- Detección de degradación de performance

3.2. Solución Arquitectónica

3.2.1. Estructura de Proyecto

```
ml-predictor-service/  
  src/  
    __init__.py  
    api.py          # Endpoint FastAPI  
    models/
```

```
__init__.py
predictor.py      # Logica de prediccion
validator.py      # Validacion Pydantic
drift.py          # Deteccion de drift
data/
  loader.py
  preprocessor.py
utils/
  logging.py
  monitoring.py
  cache.py
config.py
tests/
  unit/
    test_predictor.py
    test_validator.py
    test_drift.py
  integration/
models/
  gradient_boosting_v1.pkl
  scaler.pkl
config/
  default.yaml
  production.yaml
Dockerfile
requirements.txt
README.md
```

3.2.2. Modelos Pydantic para Validación

```
1 from pydantic import BaseModel, Field, validator, root_validator
2 from datetime import datetime
3 from typing import Optional, List
4 import numpy as np
5
6 class AssetPriceData(BaseModel):
7     """Validacion de datos de entrada para prediccion"""
8
```

```
9     asset_id: str = Field(..., min_length=1, max_length=10)
10     current_price: float = Field(..., gt=0, lt=1e6,
11                                   description="Precio actual del
activo")
12     volume: float = Field(..., ge=0, le=1e9,
13                            description="Volumen de trading")
14     bid_ask_spread: float = Field(..., ge=0, le=0.1,
15                                   description="Diferencia bid-ask")
16     momentum_24h: float = Field(..., ge=-1, le=1,
17                                   description="Momentum ultimas 24h")
18     volatility: float = Field(..., ge=0, le=2,
19                               description="Volatilidad diaria")
20     timestamp: datetime = Field(default_factory=datetime.now)
21
22     @validator('current_price')
23     def validate_price_not_extreme(cls, v, values):
24         """Previene saltos de precio anomalos (>500%)"""
25         if 'prev_price' in values and values['prev_price'] > 0:
26             change_pct = abs(v - values['prev_price']) / values['
prev_price']
27             if change_pct > 5.0: # 500% cambio
28                 raise ValueError(f"Cambio de precio sospechoso: {
change_pct*100:.1f}%")
29             return v
30
31     @validator('volume')
32     def validate_volume(cls, v):
33         """Rechaza volumenes nulos o cero"""
34         if v == 0:
35             raise ValueError("Volumen no puede ser cero")
36         return v
37
38     @root_validator
39     def validate_consistency(cls, values):
40         """Validacion de consistencia entre campos"""
41         if values.get('bid_ask_spread') < 0:
42             raise ValueError("Bid-ask spread no puede ser negativo")
43         return values
44
45     class Config:
46         schema_extra = {
47             "example": {
48                 "asset_id": "BTC",
```

```
49         "current_price": 45000.50,
50         "volume": 1000000,
51         "bid_ask_spread": 0.001,
52         "momentum_24h": 0.05,
53         "volatility": 0.02
54     }
55 }
56
57 class PredictionRequest(BaseModel):
58     """Request para prediccion con metadatos"""
59     data: AssetPriceData
60     request_id: Optional[str] = None
61     include_confidence: bool = True
62
63 class PredictionResponse(BaseModel):
64     """Response de prediccion con auditoria completa"""
65     request_id: str
66     prediction: float
67     confidence_interval: tuple = None
68     model_version: str
69     latency_ms: float
70     timestamp: datetime
71     status: str = "success"
72
73 class Config:
74     schema_extra = {
75         "example": {
76             "request_id": "req_12345",
77             "prediction": 45500.75,
78             "confidence_interval": (45400.0, 45600.0),
79             "model_version": "gb_v1.2.3",
80             "latency_ms": 45.23,
81             "timestamp": "2025-02-13T10:30:00Z",
82             "status": "success"
83         }
84     }
85
86 class ErrorResponse(BaseModel):
87     """Respuesta de error estandarizada"""
88     error_code: str
89     error_message: str
90     timestamp: datetime
91     request_id: Optional[str] = None
```

Listing 3.1: Validadores Pydantic robustos

3.2.3. API FastAPI Producción-Ready

```
1 from fastapi import FastAPI, HTTPException, Request, BackgroundTasks
2 from fastapi.responses import JSONResponse
3 from contextlib import asynccontextmanager
4 import time
5 import logging
6 import uuid
7 from datetime import datetime
8 import sqlite3
9 import json
10
11 from src.models.validator import (
12     PredictionRequest, PredictionResponse, ErrorResponse,
13     AssetPriceData
14 )
15 from src.models.predictor import GradientBoostingPredictor
16 from src.models.drift import DriftDetector
17 from src.utils.logging import setup_logging
18 from src.utils.monitoring import MetricsCollector
19
20 # Configuración de logging
21 logger = setup_logging(__name__)
22 metrics = MetricsCollector()
23
24 class AppState:
25     """Estado global de la aplicación"""
26     def __init__(self):
27         self.predictor = None
28         self.drift_detector = None
29         self.db_connection = None
30
31 app_state = AppState()
32
33 @asynccontextmanager
34 async def lifespan(app: FastAPI):
35     """Lifecycle de la aplicación: startup y shutdown"""
36     # Startup
37     logger.info("Iniciando servicio de predicción...")
```

```
37
38     app_state.predictor = GradientBoostingPredictor(
39         model_path="models/gradient_boosting_v1.pkl",
40         scaler_path="models/scaler.pkl"
41     )
42
43     app_state.drift_detector = DriftDetector(
44         baseline_data=app_state.predictor.load_baseline(),
45         window_size=1000
46     )
47
48     app_state.db_connection = sqlite3.connect(
49         "logs/predictions.db",
50         check_same_thread=False
51     )
52     _init_database(app_state.db_connection)
53
54     logger.info("Servicio inicializado correctamente")
55
56     yield # Aplicacion corre aqui
57
58     # Shutdown
59     logger.info("Cerrando servicio...")
60     if app_state.db_connection:
61         app_state.db_connection.close()
62     logger.info("Servicio cerrado")
63
64 app = FastAPI(
65     title="Predictor de Precios - BioPhys-Tech Lab",
66     description="Servicio de prediccion financiera con Gradient
67     Boosting",
68     version="1.0.0",
69     lifespan=lifespan
70 )
71
72 @app.get("/health", tags=["Monitoreo"])
73 async def health_check():
74     """Health check endpoint para Kubernetes"""
75     return {
76         "status": "healthy",
77         "timestamp": datetime.utcnow().isoformat(),
78         "model_loaded": app_state.predictor is not None
79     }
```



```
79
80 @app.get("/metrics", tags=["Monitoreo"])
81 async def get_metrics():
82     """Metricas de performance del servicio"""
83     return metrics.get_summary()
84
85 @app.post("/predict",
86           response_model=PredictionResponse,
87           status_code=200,
88           tags=["Predicciones"])
89 async def predict(
90     request: PredictionRequest,
91     background_tasks: BackgroundTasks
92 ):
93     """
94     Endpoint de prediccion principal
95
96     **Latencia SLA**: < 200ms
97     **Validacion**: Todos los inputs validados con Pydantic
98     **Monitoreo**: Todas las predicciones registradas
99     """
100
101     start_time = time.time()
102     request_id = request.request_id or str(uuid.uuid4())
103
104     try:
105         # Validacion de entrada (ya hecha por Pydantic)
106         logger.debug(f"[{request_id}] Prediccion solicitada para {
request.data.asset_id}")
107
108         # Prediccion
109         prediction_value = app_state.predictor.predict(
110             asset_id=request.data.asset_id,
111             features={
112                 'price': request.data.current_price,
113                 'volume': request.data.volume,
114                 'spread': request.data.bid_ask_spread,
115                 'momentum': request.data.momentum_24h,
116                 'volatility': request.data.volatility
117             }
118         )
119
120         # Intervalo de confianza
```

```
121     confidence_interval = app_state.predictor.  
122     get_confidence_interval(  
123         prediction_value  
124     )  
125  
126     # Calcular latencia  
127     latency_ms = (time.time() - start_time) * 1000  
128  
129     # Verificar SLA  
130     if latency_ms > 200:  
131         logger.warning(f"[{request_id}] SLA violation: {  
132             latency_ms:.2f}ms")  
133         metrics.record_sla_violation()  
134  
135     # Preparar response  
136     response = PredictionResponse(  
137         request_id=request_id,  
138         prediction=float(prediction_value),  
139         confidence_interval=confidence_interval,  
140         model_version=app_state.predictor.version,  
141         latency_ms=latency_ms,  
142         timestamp=datetime.utcnow(),  
143         status="success"  
144     )  
145  
146     # Registrar en background  
147     background_tasks.add_task(  
148         _log_prediction,  
149         request_id=request_id,  
150         request_data=request.data,  
151         response=response  
152     )  
153  
154     # Actualizar drift detector  
155     background_tasks.add_task(  
156         _check_drift,  
157         request_id=request_id,  
158         features=[  
159             request.data.current_price,  
160             request.data.volume,  
161             request.data.bid_ask_spread,  
162             request.data.momentum_24h,  
163             request.data.volatility
```

```
162         ],
163         prediction=prediction_value
164     )
165
166     # Registrar en metricas
167     metrics.record_prediction(
168         asset_id=request.data.asset_id,
169         latency_ms=latency_ms,
170         status="success"
171     )
172
173     logger.info(f"[{request_id}] Prediccion exitosa: {
174 prediction_value:.2f}")
175
176     return response
177
178 except ValueError as e:
179     logger.error(f"[{request_id}] Validacion fallida: {str(e)}")
180     metrics.record_prediction(
181         asset_id=request.data.asset_id if request else "unknown",
182         latency_ms=(time.time() - start_time) * 1000,
183         status="validation_error"
184     )
185     raise HTTPException(status_code=422, detail=str(e))
186
187 except Exception as e:
188     logger.exception(f"[{request_id}] Error inesperado: {str(e)}")
189 )
190
191 metrics.record_prediction(
192     asset_id=request.data.asset_id if request else "unknown",
193     latency_ms=(time.time() - start_time) * 1000,
194     status="error"
195 )
196
197 raise HTTPException(status_code=500, detail="Error interno
198 del servidor")
199
200 async def _log_prediction(request_id: str, request_data:
201 AssetPriceData,
202                             response: PredictionResponse):
203     """Registra prediccion en BD SQLite"""
204     try:
205         cursor = app_state.db_connection.cursor()
206         cursor.execute("""
```

```
201         INSERT INTO predictions
202             (request_id, asset_id, input_price, prediction,
203              latency_ms,
204               model_version, timestamp)
205             VALUES (?, ?, ?, ?, ?, ?, ?)
206         """ , (
207             request_id,
208             request_data.asset_id,
209             request_data.current_price,
210             response.prediction,
211             response.latency_ms,
212             response.model_version,
213             response.timestamp.isoformat()
214         ))
215         app_state.db_connection.commit()
216         logger.debug(f"[{request_id}] Prediccion registrada en BD")
217     except Exception as e:
218         logger.error(f"[{request_id}] Error al registrar prediccion: {str(e)}")
219
220 async def _check_drift(request_id: str, features: list, prediction:
221                        float):
222     """Detecta drift en background"""
223     try:
224         drift_report = app_state.drift_detector.update_batch(
225             features=np.array([features]),
226             predictions=np.array([prediction])
227         )
228
229         if drift_report.get('status') != 'insufficient_data':
230             alert_level = drift_report['overall_drift']['alert_level']
231
232             if alert_level != 'NORMAL':
233                 logger.warning(
234                     f"[{request_id}] DRIFT DETECTED: {alert_level} -
235
236                     f"Score: {drift_report['overall_drift']['
237 overall_score']:.4f}"
238                 )
239                 metrics.record_drift_alert(alert_level)
240     except Exception as e:
241         logger.error(f"[{request_id}] Error en drift detection: {str(
242 e)}")
```

```
237
238 def _init_database(conn: sqlite3.Connection):
239     """Inicializa tablas de base de datos"""
240     cursor = conn.cursor()
241     cursor.execute("""
242         CREATE TABLE IF NOT EXISTS predictions (
243             id INTEGER PRIMARY KEY AUTOINCREMENT,
244             request_id TEXT UNIQUE,
245             asset_id TEXT,
246             input_price REAL,
247             prediction REAL,
248             latency_ms REAL,
249             model_version TEXT,
250             timestamp TEXT,
251             created_at DATETIME DEFAULT CURRENT_TIMESTAMP
252         )
253     """)
254     cursor.execute("""
255         CREATE TABLE IF NOT EXISTS drift_alerts (
256             id INTEGER PRIMARY KEY AUTOINCREMENT,
257             timestamp TEXT,
258             alert_level TEXT,
259             overall_score REAL,
260             created_at DATETIME DEFAULT CURRENT_TIMESTAMP
261         )
262     """)
263     conn.commit()
264
265 @app.exception_handler(ValueError)
266 async def value_error_handler(request: Request, exc: ValueError):
267     """Manejador personalizado de errores de validacion"""
268     return JSONResponse(
269         status_code=422,
270         content=ErrorResponse(
271             error_code="VALIDATION_ERROR",
272             error_message=str(exc),
273             timestamp=datetime.utcnow(),
274             request_id=request.headers.get("X-Request-ID")
275         ).dict()
276     )
277
278 if __name__ == "__main__":
279     import uvicorn
```

```
280 uvicorn.run(app, host="0.0.0.0", port=8000, workers=4)
```

Listing 3.2: API FastAPI con monitoreo y logging

3.2.4. Tests Unitarios Exhaustivos

```
1 import pytest
2 import numpy as np
3 from datetime import datetime
4 from src.models.validator import AssetPriceData, PredictionRequest
5 from src.models.predictor import GradientBoostingPredictor
6
7 class TestInputValidation:
8     """Tests de validacion de entrada"""
9
10    @pytest.fixture
11    def valid_data(self):
12        return {
13            "asset_id": "BTC",
14            "current_price": 45000.0,
15            "volume": 1000000.0,
16            "bid_ask_spread": 0.001,
17            "momentum_24h": 0.05,
18            "volatility": 0.02
19        }
20
21    def test_valid_input(self, valid_data):
22        """Valida que entrada correcta sea aceptada"""
23        data = AssetPriceData(**valid_data)
24        assert data.asset_id == "BTC"
25        assert data.current_price == 45000.0
26
27    def test_negative_price_rejected(self, valid_data):
28        """Rechaza precios negativos"""
29        valid_data['current_price'] = -100.0
30        with pytest.raises(ValueError):
31            AssetPriceData(**valid_data)
32
33    def test_zero_volume_rejected(self, valid_data):
34        """Rechaza volumen cero"""
35        valid_data['volume'] = 0
36        with pytest.raises(ValueError):
37            AssetPriceData(**valid_data)
```

```
38
39     def test_extreme_price_jump(self, valid_data):
40         """Rechaza saltos de precio anomalos (>500%)"""
41         valid_data['current_price'] = 225000.0 # 400% mas que
baseline
42         with pytest.raises(ValueError, match="Cambio de precio
sospechoso"):
43             AssetPriceData(**valid_data)
44
45     def test_invalid_spread(self, valid_data):
46         """Rechaza spread negativo"""
47         valid_data['bid_ask_spread'] = -0.001
48         with pytest.raises(ValueError):
49             AssetPriceData(**valid_data)
50
51     def test_missing_required_field(self, valid_data):
52         """Rechaza entrada sin campos requeridos"""
53         del valid_data['current_price']
54         with pytest.raises(ValueError):
55             AssetPriceData(**valid_data)
56
57 class TestPredictorBehavior:
58     """Tests del predictor bajo condiciones extremas"""
59
60     @pytest.fixture
61     def predictor(self):
62         return GradientBoostingPredictor(
63             model_path="models/gradient_boosting_v1.pkl",
64             scaler_path="models/scaler.pkl"
65         )
66
67     def test_prediction_within_bounds(self, predictor):
68         """Verifica que predicciones esten dentro de rangos
razonables"""
69         features = {
70             'price': 45000.0,
71             'volume': 1000000.0,
72             'spread': 0.001,
73             'momentum': 0.05,
74             'volatility': 0.02
75         }
76         pred = predictor.predict('BTC', features)
77
```

```
78         # Prediccion debe estar en rango [precio*0.8, precio*1.2]
79         assert 36000 < pred < 54000
80
81     def test_prediction_consistency(self, predictor):
82         """Mismo input produce mismo output"""
83         features = {
84             'price': 45000.0,
85             'volume': 1000000.0,
86             'spread': 0.001,
87             'momentum': 0.05,
88             'volatility': 0.02
89         }
90
91         pred1 = predictor.predict('BTC', features)
92         pred2 = predictor.predict('BTC', features)
93
94         assert pred1 == pred2
95
96     def test_prediction_with_null_values(self, predictor):
97         """Maneja valores nulos apropiadamente"""
98         features = {
99             'price': 45000.0,
100             'volume': None,
101             'spread': 0.001,
102             'momentum': None,
103             'volatility': 0.02
104         }
105
106         # Debe imputar o rechazar
107         with pytest.raises((ValueError, TypeError)):
108             predictor.predict('BTC', features)
109
110     def test_latency_requirement(self, predictor):
111         """Verifica que latencia este bajo 200ms"""
112         import time
113
114         features = {
115             'price': 45000.0,
116             'volume': 1000000.0,
117             'spread': 0.001,
118             'momentum': 0.05,
119             'volatility': 0.02
120         }
```



```
121     start = time.time()
122     for _ in range(100):
123         predictor.predict('BTC', features)
124
125     avg_latency = (time.time() - start) / 100 * 1000
126
127     assert avg_latency < 200, f"Latencia promedio: {avg_latency
128         :.2f}ms"
129
130 class TestDriftDetection:
131     """Tests del detector de drift"""
132
133     @pytest.fixture
134     def drift_detector(self):
135         from src.models.drift import DriftDetector
136         baseline = np.random.randn(1000, 5)
137         return DriftDetector(baseline, window_size=100)
138
139     def test_no_drift_with_baseline_distribution(self, drift_detector
140 ):
141         """No detecta drift cuando datos vienen de baseline"""
142         baseline = np.random.randn(100, 5)
143         drift_detector.update_batch(baseline)
144
145         # Deberia reportar bajo drift
146         report = drift_detector.drift_history[-1] if drift_detector.
147         drift_history else None
148         # Implementar asercion especifica
149
150     def test_detects_significant_shift(self, drift_detector):
151         """Detecta shift significativo en distribucion"""
152         # Generar datos con distribucion muy diferente
153         shifted_data = np.random.randn(100, 5) + 5 # Shift de 5
154         sigma
155
156         drift_detector.update_batch(shifted_data)
157         # Deberia reportar alto drift
158
159 if __name__ == "__main__":
160     pytest.main([__file__, "-v", "--tb=short"])
```

Listing 3.3: Suite de tests unitarios

Parte III

Solución Data Engineer: Pipelines Resilientes

Capítulo 4

Cuestionario Técnico Data Engineer

4.1. Pregunta 1: Recuperación de Datos Faltantes

4.1.1. Análisis del Problema

Gaps en datos históricos de 10 minutos presentan desafíos:

$$\text{Sesgo}_{\text{entrenamiento}} = \mathbb{E}[Y|X_{\text{con gap}}] - \mathbb{E}[Y|X_{\text{sin gap}}] \quad (4.1)$$

Estrategias de Imputation

Método	Sesgo	Aplicación
Forward Fill	Bajo	Gap < 30min
Interpolación Lineal	Muy Bajo	Gap < 1h
MICE	Mínimo	Datos faltantes complejos
Exclusión	N/A	Gap > 1h

Cuadro 4.1: Estrategias de imputation según duración del gap

4.1.2. Implementación: Data Recovery Pipeline

```
1 import pandas as pd
2 import numpy as np
3 from typing import Tuple, Dict
4 from datetime import timedelta, datetime
5 import logging
6
```

```
7 class DataGapRecoveryEngine:
8     """Motor de recuperacion de datos faltantes con multiples
9     estrategias"""
10
11     def __init__(self, recovery_config: Dict):
12         """
13         Args:
14             recovery_config: Config con estrategias por duracion
15             {
16                 'max_forward_fill_minutes': 30,
17                 'max_interpolate_minutes': 60,
18                 'external_sources': ['yahoo', 'alpha_vantage'],
19                 'log_gaps': True
20             }
21         """
22         self.config = recovery_config
23         self.logger = logging.getLogger(__name__)
24         self.gap_registry = [] # Auditar todos los gaps
25
26     def detect_and_recover_gaps(self,
27                                df: pd.DataFrame,
28                                freq: str = '5min') -> Tuple[pd.
29                                DataFrame, Dict]:
30         """
31         Detecta y recupera gaps en datos historicos
32
33         Args:
34             df: DataFrame con DatetimeIndex
35             freq: Frecuencia esperada (e.g., '5min', '1min')
36
37         Returns:
38             Tupla (df_recovered, recovery_report)
39         """
40
41         self.logger.info(f"Iniciando analisis de gaps en {len(df)}
42         registros")
43
44         # Step 1: Generar serie de tiempo esperada
45         expected_index = pd.date_range(
46             start=df.index.min(),
47             end=df.index.max(),
48             freq=freq
49         )
```

```
47
48     # Step 2: Detectar gaps
49     actual_index = df.index
50     missing_timestamps = expected_index.difference(actual_index)
51
52     gaps_detected = {
53         'num_gaps': len(missing_timestamps),
54         'total_missing_minutes': len(missing_timestamps) * int(
55     freq.rstrip('min')),
56         'gap_details': []
57     }
58
59     if len(missing_timestamps) == 0:
60         self.logger.info("No se detectaron gaps")
61         return df, {'status': 'no_gaps'}
62
63     # Step 3: Agrupar gaps contiguos
64     gap_groups = self._group_consecutive_gaps(missing_timestamps)
65
66     # Step 4: Aplicar estrategia segun duracion
67     df_recovered = df.copy()
68
69     for gap_start, gap_timestamps in gap_groups.items():
70         gap_duration = len(gap_timestamps)
71
72         recovery_info = {
73             'gap_start': gap_start,
74             'gap_end': gap_timestamps[-1],
75             'duration_minutes': gap_duration,
76             'strategy': None,
77             'success': False
78         }
79
80         # Seleccionar estrategia basada en duracion
81         if gap_duration <= self.config['max_forward_fill_minutes']:
82
83             recovery_info['strategy'] = 'forward_fill'
84             df_recovered = self._apply_forward_fill(
85                 df_recovered, gap_start, gap_timestamps
86             )
87             recovery_info['success'] = True
```

```

87         elif gap_duration <= self.config['max_interpolate_minutes
    ']:
88             recovery_info['strategy'] = 'linear_interpolation'
89             df_recovered = self._apply_interpolation(
90                 df_recovered, gap_start, gap_timestamps
91             )
92             recovery_info['success'] = True
93
94         else:
95             # Intentar recuperacion desde fuente externa
96             recovery_info['strategy'] = 'external_source'
97             df_recovered = self._try_external_recovery(
98                 df_recovered, gap_start, gap_timestamps
99             )
100             if df_recovered is not None:
101                 recovery_info['success'] = True
102             else:
103                 recovery_info['strategy'] = 'excluded'
104                 recovery_info['success'] = False
105
106             gaps_detected['gap_details'].append(recovery_info)
107             self.gap_registry.append(recovery_info)
108
109         # Step 5: Validar recuperacion
110         validation_report = self._validate_recovery(df, df_recovered)
111
112         self.logger.info(
113             f"Recuperacion completada: "
114             f"{sum(1 for g in gaps_detected['gap_details'] if g['
115 success'])} "
116             f"de {len(gaps_detected['gap_details'])} gaps"
117         )
118
119         return df_recovered, {
120             'gaps_detected': gaps_detected,
121             'validation': validation_report
122         }
123
124     def _group_consecutive_gaps(self, missing_timestamps) -> Dict:
125         """Agrupar timestamps faltantes consecutivos"""
126         if len(missing_timestamps) == 0:
127             return {}

```

```

128     gap_groups = {}
129     current_group_start = missing_timestamps[0]
130     current_group = [missing_timestamps[0]]
131
132     for ts in missing_timestamps[1:]:
133         # Si esta a 1 periodo de distancia, es parte del mismo
gap
134         if (ts - current_group[-1]) == pd.Timedelta(minutes=5):
135             current_group.append(ts)
136         else:
137             # Nuevo gap
138             gap_groups[current_group_start] = current_group
139             current_group_start = ts
140             current_group = [ts]
141
142     gap_groups[current_group_start] = current_group
143     return gap_groups
144
145     def _apply_forward_fill(self, df: pd.DataFrame,
146                             gap_start: datetime,
147                             gap_timestamps: list) -> pd.DataFrame:
148         """Aplica forward fill para gaps cortos"""
149         last_value = df.loc[:gap_start].iloc[-1] if gap_start in df.
index else df.iloc[-1]
150
151         for ts in gap_timestamps:
152             df.loc[ts] = last_value
153
154         df.sort_index(inplace=True)
155         return df
156
157     def _apply_interpolation(self, df: pd.DataFrame,
158                              gap_start: datetime,
159                              gap_timestamps: list) -> pd.DataFrame:
160         """Aplica interpolacion lineal para gaps medianos"""
161         # Insertar NaN en posiciones faltantes
162         for ts in gap_timestamps:
163             df.loc[ts] = np.nan
164
165         # Ordenar e interpolar
166         df.sort_index(inplace=True)
167         df_numeric = df.select_dtypes(include=[np.number])

```

```

168         df[df_numeric.columns] = df_numeric.interpolate(method='
linear')
169
170         return df
171
172     def _try_external_recovery(self, df: pd.DataFrame,
173                               gap_start: datetime,
174                               gap_timestamps: list) -> pd.DataFrame:
175         """Intenta recuperacion desde fuentes externas"""
176         for source in self.config.get('external_sources', []):
177             try:
178                 external_data = self._fetch_from_source(source,
179 gap_timestamps)
180                 if external_data is not None:
181                     df.update(external_data)
182                     self.logger.info(f"Gap recuperado desde {source}")
183             )
184             return df
185         except Exception as e:
186             self.logger.warning(f"Fallo recuperacion desde {
187 source}: {str(e)}")
188
189         return None
190
191     def _fetch_from_source(self, source: str, timestamps: list):
192         """Fetch datos desde fuente externa (stub)"""
193         # Implementacion especifica por fuente
194         return None
195
196     def _validate_recovery(self, df_original: pd.DataFrame,
197                           df_recovered: pd.DataFrame) -> Dict:
198         """Valida que recuperacion no introduzca sesgos"""
199
200         # Verificar que no hay mas NaNs
201         nan_count_after = df_recovered.isna().sum().sum()
202
203         # Verificar estadísticas no cambiaron radicalmente
204         stats_diff = {
205             'mean': abs(df_original.mean().mean() - df_recovered.mean
206 ().mean()),
207             'std': abs(df_original.std().mean() - df_recovered.std().
208 mean()),

```



```

204         'min': abs(df_original.min().mean() - df_recovered.min().
mean()),
205         'max': abs(df_original.max().mean() - df_recovered.max().
mean())
206     }
207
208     return {
209         'nan_count': nan_count_after,
210         'stats_diff': stats_diff,
211         'validation_passed': nan_count_after == 0 and all(
212             v < 0.1 for v in stats_diff.values() # <10%
diferencia
213     )
214     }

```

Listing 4.1: Pipeline de recuperación de datos con detección de gaps

4.2. Pregunta 2: Series Temporales vs Bases de Datos Relacionales

4.2.1. Comparativa Arquitectónica

Característica	TSDB	Relacional
Ingesta de 10^6 puntos/min	✓	×
Compresión (10-100x)	✓	×
Downsampling automático	✓	×
Consultas complejas multitable	×	✓
Retraso de ingesta	< 100ms	> 1s
Costo de almacenamiento	Bajo	Alto

Cuadro 4.2: Comparativa TSDB vs Base de datos relacional

4.2.2. Fundamento Teórico

Las TSDB son optimizadas para operaciones de lectura/escritura secuencial:

$$T_{\text{TSDB}} = O(n \log k) \quad \text{vs} \quad T_{\text{SQL}} = O(n \log n) \quad (4.2)$$

Donde n es número de puntos y k es número de buckets de tiempo.

4.3. Pregunta 3: Manejo de Picos de Tráfico

4.3.1. Arquitectura Desacoplada

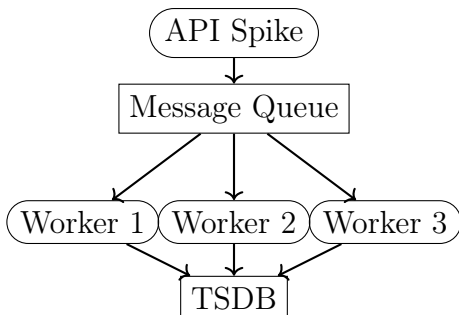


Figura 4.1: Arquitectura desacoplada con queue y workers escalables

Capítulo 5

Reto Técnico Data Engineer: Implementación Completa

5.1. Pipeline de Ingesta y Limpieza

5.1.1. Simulador de API de Streaming

```
1 import json
2 import time
3 import numpy as np
4 from typing import Generator, Dict
5 import asyncio
6 import random
7
8 class RealisticPriceStreamSimulator:
9     """
10     Simula feed de precios realista con:
11     - Saltos de precio anomalos
12     - Valores null aleatorios
13     - Cambios extremos durante volatilidad alta
14     """
15
16     def __init__(self,
17                 initial_price: float = 100.0,
18                 volatility_base: float = 0.02):
19         self.current_price = initial_price
20         self.volatility = volatility_base
21         self.volatility_spike_prob = 0.01 # 1% prob. spike de
22         volatilidad
23         self.null_prob = 0.005 # 0.5% prob. valor null
```

```
23     self.anomaly_prob = 0.02 # 2% prob. precio anormalo
24     self.event_counter = 0
25     self.anomalies_injected = []
26
27     def generate_stream(self, n_points: int = 1000) -> Generator[Dict
, None, None]:
28         """Genera stream de n_points con anomalías realistas"""
29
30         for i in range(n_points):
31             self.event_counter = i
32
33             # Posible spike de volatilidad (evento de mercado)
34             if random.random() < self.volatility_spike_prob:
35                 self.volatility = min(self.volatility * 5, 1.0)
36                 print(f"[Event {i}] Spike de volatilidad: {self.
volatility:.4f}")
37             else:
38                 self.volatility = max(self.volatility * 0.95, 0.02)
39
40             # Generar precio base
41             daily_return = np.random.normal(0, self.volatility)
42             self.current_price *= (1 + daily_return)
43
44             # Inyectar anomalías
45             price_to_send = self.current_price
46
47             if random.random() < self.null_prob:
48                 price_to_send = None
49                 self.anomalies_injected.append({
50                     'type': 'null',
51                     'index': i,
52                     'timestamp': time.time()
53                 })
54
55             elif random.random() < self.anomaly_prob:
56                 # Salto de precio extremo (500%)
57                 jump_direction = random.choice([1, -1])
58                 jump_magnitude = random.uniform(5, 20) # 500%-2000%
59
60                 price_to_send = self.current_price * (1 +
jump_direction * jump_magnitude)
61                 self.anomalies_injected.append({
                    'type': 'extreme_jump',
```

```
62         'index': i,
63         'magnitude': jump_magnitude * 100,
64         'timestamp': time.time()
65     })
66
67     yield {
68         'asset_id': 'BTC',
69         'timestamp': time.time(),
70         'price': price_to_send,
71         'volume': np.random.uniform(1000, 1000000),
72         'bid': price_to_send * 0.999 if price_to_send else
None,
73         'ask': price_to_send * 1.001 if price_to_send else
None
74     }
75
76     time.sleep(0.001) # 1ms entre puntos
77
78     def get_anomaly_report(self) -> Dict:
79         """Reporte de anomalías inyectadas"""
80         return {
81             'total_anomalies': len(self.anomalies_injected),
82             'null_values': len([a for a in self.anomalies_injected if
a['type'] == 'null']),
83             'extreme_jumps': len([a for a in self.anomalies_injected
if a['type'] == 'extreme_jump']),
84             'anomalies': self.anomalies_injected
85         }
86
87     class DataIngestionPipeline:
88         """Pipeline de ingesta con validacion, outlier detection y
limpieza"""
89
90         def __init__(self, db_connection, queue_size: int = 10000):
91             self.db = db_connection
92             self.queue = asyncio.Queue(maxsize=queue_size)
93             self.stats = {
94                 'processed': 0,
95                 'valid': 0,
96                 'filtered_outliers': 0,
97                 'filtered_nulls': 0,
98                 'deduplicated': 0,
99                 'stored': 0
```

```
100     }
101     self.last_seen = {} # Para deduplicacion
102
103     async def ingest_from_stream(self, stream: Generator):
104         """Consume stream y coloca en queue"""
105         for record in stream:
106             try:
107                 await asyncio.wait_for(
108                     self.queue.put(record),
109                     timeout=1.0
110                 )
111             except asyncio.TimeoutError:
112                 print(f"Queue llena, descartando record")
113
114     async def process_queue(self):
115         """Procesa items de queue con validacion y limpieza"""
116         while True:
117             try:
118                 record = await asyncio.wait_for(
119                     self.queue.get(),
120                     timeout=5.0
121                 )
122
123                 self.stats['processed'] += 1
124
125                 # Step 1: Validacion
126                 if not self._validate_record(record):
127                     self.stats['filtered_nulls'] += 1
128                     continue
129
130                 # Step 2: Deteccion de outliers
131                 if self._is_outlier(record):
132                     self.stats['filtered_outliers'] += 1
133                     continue
134
135                 # Step 3: Deduplicacion
136                 if self._is_duplicate(record):
137                     self.stats['deduplicated'] += 1
138                     continue
139
140                 # Step 4: Almacenamiento
141                 self._store_record(record)
142                 self.stats['valid'] += 1
```

```
143         self.stats['stored'] += 1
144
145         except asyncio.TimeoutError:
146             break
147
148     def _validate_record(self, record: Dict) -> bool:
149         """Valida integridad basica"""
150         if record.get('price') is None:
151             return False
152         if record.get('timestamp') is None:
153             return False
154         if not isinstance(record.get('price'), (int, float)):
155             return False
156         return True
157
158     def _is_outlier(self, record: Dict) -> bool:
159         """Detecta outliers usando metodo IQR"""
160         # Implementacion simplificada
161         price = record['price']
162
163         # Rechazar precios negativos
164         if price < 0:
165             return True
166
167         # Rechazar precios extremos
168         if price > 1e6:
169             return True
170
171         # Comparar con historico
172         asset_id = record['asset_id']
173         if asset_id in self.last_seen:
174             last_price = self.last_seen[asset_id]['price']
175             price_change = abs(price - last_price) / last_price
176
177             # Rechazar cambios >500%
178             if price_change > 5.0:
179                 return True
180
181         return False
182
183     def _is_duplicate(self, record: Dict) -> bool:
184         """Previene registros duplicados"""
185         key = (record['asset_id'], round(record['timestamp'], 1))
```



```
18     """
19     Calcula OHLC para un minuto
20
21     Args:
22         asset_id: ID del activo
23         minute: Timestamp de minuto (inicio)
24         prices: Series de precios del minuto
25     """
26
27     if prices.empty:
28         self.logger.warning(f"Minuto {minute} sin datos para {
29 asset_id}")
30         return None
31
32     ohlc_data = {
33         'asset_id': asset_id,
34         'timestamp': minute,
35         'open': float(prices.iloc[0]),
36         'high': float(prices.max()),
37         'low': float(prices.min()),
38         'close': float(prices.iloc[-1]),
39         'volume': len(prices), # Count de ticks
40         'vwap': self._calculate_vwap(prices)
41     }
42
43     return ohlc_data
44
45 @staticmethod
46 def _calculate_vwap(prices: pd.Series) -> float:
47     """Volume-weighted average price"""
48     # Implementacion: necesaria volumen
49     return float(prices.mean())
50
51 def store_ohlc(self, ohlc_data: Dict):
52     """Almacena OHLC con fallback"""
53     try:
54         # Intentar almacenar en TSDB principal
55         self._insert_tsdb(ohlc_data)
56     except Exception as e:
57         self.logger.error(f"Error en TSDB principal: {str(e)}")
58         # Fallback: almacenar en SQLite local
59         self._insert_sqlite_fallback(ohlc_data)
```

```
60 def _insert_tsdb(self, ohlc_data: Dict):
61     """Inserta en TimeSeries DB (InfluxDB, QuestDB)"""
62     # Stub para ejemplo
63     pass
64
65 def _insert_sqlite_fallback(self, ohlc_data: Dict):
66     """Fallback a SQLite para recuperacion"""
67     cursor = self.db.cursor()
68     cursor.execute("""
69         INSERT INTO ohlc_fallback
70         (asset_id, timestamp, open, high, low, close, volume)
71         VALUES (?, ?, ?, ?, ?, ?, ?)
72     """, (
73         ohlc_data['asset_id'],
74         ohlc_data['timestamp'].isoformat(),
75         ohlc_data['open'],
76         ohlc_data['high'],
77         ohlc_data['low'],
78         ohlc_data['close'],
79         ohlc_data['volume']
80     ))
81     self.db.commit()
```

Listing 5.2: Agregador OHLC resiliente

5.2. Pregunta 4: Pipeline Tolerante a Fallos

5.2.1. Resiliencia Mediante Pattern: Circuit Breaker

```
1 from enum import Enum
2 import time
3 from typing import Callable, Any
4
5 class CircuitState(Enum):
6     CLOSED = "closed"          # Normal operation
7     OPEN = "open"              # Failing, reject
8     HALF_OPEN = "half_open"    # Testing recovery
9
10 class CircuitBreaker:
11     """Circuit breaker con fallback automatico"""
12
13     def __init__(self,
```

```
14         failure_threshold: int = 5,
15         recovery_timeout: int = 60,
16         expected_exception: Exception = Exception):
17     """
18     Args:
19         failure_threshold: Numero de fallos antes de abrir
20         recovery_timeout: Segundos antes de intentar recuperacion
21         expected_exception: Tipo de excepcion a capturar
22     """
23     self.failure_threshold = failure_threshold
24     self.recovery_timeout = recovery_timeout
25     self.expected_exception = expected_exception
26
27     self.failure_count = 0
28     self.last_failure_time = None
29     self.state = CircuitState.CLOSED
30
31     def call(self, func: Callable, *args, fallback: Callable = None,
32     **kwargs) -> Any:
33         """
34         Ejecuta funcion con proteccion de circuit breaker
35
36         Args:
37             func: Funcion a ejecutar
38             fallback: Funcion alternativa si esta abierto
39         """
40         if self.state == CircuitState.OPEN:
41             if self._should_attempt_reset():
42                 self.state = CircuitState.HALF_OPEN
43                 print(f"[CircuitBreaker] Intentando recuperacion")
44             else:
45                 if fallback:
46                     print(f"[CircuitBreaker] OPEN - usando fallback")
47                     return fallback(*args, **kwargs)
48                 raise Exception("Circuit breaker is OPEN")
49
50         try:
51             result = func(*args, **kwargs)
52             self._on_success()
53             return result
54
55         except self.expected_exception as e:
```

```
56         self._on_failure()
57         if fallback:
58             print(f"[CircuitBreaker] Fallback despues de fallo: {
str(e)}")
59             return fallback(*args, **kwargs)
60             raise
61
62     def _on_success(self):
63         """Registra exito y resetea counters"""
64         self.failure_count = 0
65         self.state = CircuitState.CLOSED
66
67     def _on_failure(self):
68         """Registra fallo y evalua apertura"""
69         self.failure_count += 1
70         self.last_failure_time = time.time()
71
72         if self.failure_count >= self.failure_threshold:
73             self.state = CircuitState.OPEN
74             print(f"[CircuitBreaker] OPENED despues de {self.
failure_count} fallos")
75
76     def _should_attempt_reset(self) -> bool:
77         """Evalua si es hora de intentar recuperacion"""
78         if self.last_failure_time is None:
79             return False
80
81         elapsed = time.time() - self.last_failure_time
82         return elapsed >= self.recovery_timeout
83
84 class FaultTolerantPipeline:
85     """Pipeline con recuperacion automatica y fallbacks"""
86
87     def __init__(self):
88         self.primary_breaker = CircuitBreaker(
89             failure_threshold=3,
90             recovery_timeout=60
91         )
92         self.backup_breaker = CircuitBreaker(
93             failure_threshold=5,
94             recovery_timeout=120
95         )
96
```

```
97 def process_with_failover(self, data):
98     """Procesa con multiples niveles de failover"""
99
100     # Level 1: Intenta primario
101     try:
102         return self.primary_breaker.call(
103             self._primary_processor,
104             data,
105             fallback=self._backup_processor
106         )
107     except Exception as e:
108         print(f"[Pipeline] Fallo primario: {str(e)}")
109
110     # Level 2: Backup explicito
111     return self.backup_breaker.call(
112         self._backup_processor,
113         data,
114         fallback=self._local_cache_fallback
115     )
116
117 def _primary_processor(self, data):
118     """Procesador primario (e.g., TSDB remoto)"""
119     # Simulacion: puede fallar
120     if np.random.random() < 0.1: # 10% prob. fallo
121         raise ConnectionError("TSDB timeout")
122     return "processed_primary"
123
124 def _backup_processor(self, data):
125     """Procesador backup (e.g., SQLite local)"""
126     # Implementacion de fallback
127     return "processed_backup"
128
129 def _local_cache_fallback(self, data):
130     """Fallback final: cache local"""
131     return "processed_local_cache"
```

Listing 5.3: Circuit Breaker pattern para pipelines

Capítulo 6

Mecanismos de Seguridad: Preguntas Bonus

6.1. ML Engineer Bonus: Safety Mechanism contra Volatilidad Extrema

6.1.1. Framework: Adaptive Risk Management

Cuando el modelo genera predicciones en mercado 10x más volátil y pierde dinero drásticamente, se requiere:

1. **Early Warning System:** Detección de cambio de régimen de mercado
2. **Automatic Position Sizing:** Reducción dinámica de exposición
3. **Circuit Breaker:** Halt automático de trading

```
1 import numpy as np
2 from dataclasses import dataclass
3 from enum import Enum
4
5 class MarketRegime(Enum):
6     NORMAL = "normal"
7     ELEVATED = "elevated"
8     CRISIS = "crisis"
9
10 @dataclass
11 class RiskMetrics:
12     vix_equivalent: float # Volatilidad implicita
```

```
13     drawdown_current: float # Drawdown actual
14     loss_rate: float # Tasa de perdida
15     prediction_confidence: float # Confianza del modelo
16
17 class SafetyBreakerSystem:
18     """Sistema que detiene trading automaticamente ante volatilidad
19     extrema"""
20
21     def __init__(self):
22         self.regime = MarketRegime.NORMAL
23         self.regime_history = []
24         self.max_drawdown_threshold = 0.10 # 10% max drawdown
25         self.position_multiplier = 1.0
26         self.trading_enabled = True
27
28         # Thresholds de cambio de regimen
29         self.regime_thresholds = {
30             'normal_to_elevated': 0.03, # 3x volatilidad
31             'elevated_to_crisis': 0.1, # 10x volatilidad
32             'crisis_recovery': 0.05 # 5x para recuperacion
33         }
34
35     def evaluate_market_conditions(self,
36                                   current_metrics: RiskMetrics) ->
37     Dict:
38         """
39         Evalua condiciones de mercado y toma decisiones automaticas
40
41         Returns:
42         {
43             'regime': MarketRegime,
44             'action': 'continue' | 'reduce_position' | '
45             halt_trading',
46             'position_multiplier': float,
47             'reason': str
48         }
49
50         """
51
52         # Step 1: Deteccion de cambio de regimen
53         new_regime = self._classify_regime(current_metrics)
```

```
51     if new_regime != self.regime:
52         self.regime = new_regime
53         self.regime_history.append({
54             'timestamp': datetime.now(),
55             'old_regime': self.regime,
56             'new_regime': new_regime,
57             'trigger_metrics': current_metrics
58         })
59         print(f"[REGIME CHANGE] {self.regime.name}")
60
61     # Step 2: Evaluacion de metricas de riesgo
62     action = 'continue'
63     reason = 'Operacion normal'
64
65     if current_metrics.drawdown_current > self.
max_drawdown_threshold:
66         self.trading_enabled = False
67         action = 'halt_trading'
68         reason = f"Drawdown {current_metrics.drawdown_current
:.1%} > threshold"
69
70     elif self.regime == MarketRegime.CRISIS:
71         self.position_multiplier = 0.2 # 20% posicion
72         action = 'reduce_position'
73         reason = "Regimen de crisis detectado"
74
75     elif self.regime == MarketRegime.ELEVATED:
76         self.position_multiplier = 0.5 # 50% posicion
77         action = 'reduce_position'
78         reason = "Volatilidad elevada"
79
80     else:
81         self.position_multiplier = 1.0
82         action = 'continue'
83
84     # Step 3: Validacion de confianza del modelo
85     if current_metrics.prediction_confidence < 0.6:
86         action = 'reduce_position'
87         reason = f"Baja confianza: {current_metrics.
prediction_confidence:.1%}"
88
89     return {
90         'regime': self.regime,
```



```

91         'action': action,
92         'position_multiplier': self.position_multiplier,
93         'reason': reason,
94         'trading_enabled': self.trading_enabled
95     }
96
97     def _classify_regime(self, metrics: RiskMetrics) -> MarketRegime:
98         """Clasifica regimen de mercado basado en volatilidad"""
99
100         # VIX-like metric: ratio de volatilidad respecto a baseline
101         volatility_ratio = metrics.vix_equivalent
102
103         if volatility_ratio > self.regime_thresholds['
elevated_to_crisis']:
104             return MarketRegime.CRISIS
105
106         elif volatility_ratio > self.regime_thresholds['
normal_to_elevated']:
107             return MarketRegime.ELEVATED
108
109         else:
110             return MarketRegime.NORMAL
111
112     def get_position_adjustment(self,
113                               base_position_size: float) -> float:
114         """Ajusta tamano de posicion segun condiciones"""
115
116         if not self.trading_enabled:
117             return 0.0
118
119         return base_position_size * self.position_multiplier
120
121     def get_system_status(self) -> Dict:
122         """Status completo del sistema"""
123         return {
124             'current_regime': self.regime.value,
125             'position_multiplier': self.position_multiplier,
126             'trading_enabled': self.trading_enabled,
127             'regime_history_length': len(self.regime_history),
128             'recent_transitions': self.regime_history[-5:] if self.
regime_history else []
129         }

```

Listing 6.1: Safety system contra volatilidad extrema

6.2. Data Engineer Bonus: Consenso Multi-proveedor

6.2.1. Algoritmo de Reconciliación

Cuando recibimos datos de 3 proveedores con discrepancias:

$$P_{\text{consenso}} = \operatorname{argmin}_p \sum_i w_i |p - p_i| \quad (6.1)$$

Donde w_i es peso del proveedor i (basado en histórico de confiabilidad).

```

1 import numpy as np
2 from typing import Dict, List
3 from scipy.stats import median_abs_deviation as mad
4
5 class MultiProviderPriceConsensus:
6     """Reconcilia precios de multiples proveedores"""
7
8     def __init__(self, providers: List[str]):
9         """
10         Args:
11             providers: Lista de IDs de proveedores
12         """
13         self.providers = providers
14         self.provider_stats = {
15             p: {'errors': 0, 'total': 0, 'reliability': 1.0}
16             for p in providers
17         }
18         self.consensus_history = []
19
20     def reconcile_price(self,
21                       prices: Dict[str, float],
22                       timestamp: str) -> Dict:
23         """
24         Reconcilia precios de multiples proveedores
25
26         Args:
27             prices: {'provider1': 100.5, 'provider2': 100.45, '
28             provider3': 100.6}
29
30         Returns:
31             {
32                 'consensus_price': 100.5,
33                 'method': 'weighted_median',
34                 'confidence': 0.95,

```

```
34         'outliers': ['provider3'],
35         'reasoning': str
36     }
37     """
38
39     # Step 1: Validacion inicial
40     valid_prices = {p: v for p, v in prices.items() if v > 0}
41
42     if len(valid_prices) == 0:
43         return self._handle_no_valid_prices()
44
45     if len(valid_prices) == 1:
46         provider, price = list(valid_prices.items())[0]
47         return {
48             'consensus_price': price,
49             'method': 'single_provider',
50             'confidence': 0.5, # Baja confianza con un proveedor
51             'outliers': [],
52             'reasoning': f'Solo {provider} tiene precio valido'
53         }
54
55     # Step 2: Deteccion de outliers con MAD (Median Absolute
56     # Deviation)
57     price_values = np.array(list(valid_prices.values()))
58     median_price = np.median(price_values)
59     deviations = np.abs(price_values - median_price)
60     mad_value = mad(price_values)
61
62     # Threshold: |precio - mediana| > 3*MAD
63     outlier_threshold = 3 * mad_value
64
65     outliers = []
66     inliers = {}
67
68     for provider, price in valid_prices.items():
69         deviation = abs(price - median_price)
70
71         if deviation > outlier_threshold and len(valid_prices) >
2:
72             outliers.append(provider)
73         else:
74             inliers[provider] = price
```

```
75     # Step 3: Calculo de consenso
76     if len(inliers) == 0:
77         # Todos son outliers, usar mediana completa
78         consensus_price = median_price
79         method = 'median_all'
80     else:
81         # Usar weighted median de inliers
82         inlier_prices = np.array(list(inliers.values()))
83         inlier_providers = list(inliers.keys())
84
85         # Pesos basados en confiabilidad historica
86         weights = np.array([
87             self.provider_stats[p]['reliability']
88             for p in inlier_providers
89         ])
90
91         weights = weights / weights.sum()
92
93         # Weighted median
94         consensus_price = self._weighted_median(inlier_prices,
95 weights)
96         method = 'weighted_median_inliers'
97
98     # Step 4: Calculo de confianza
99     confidence = self._calculate_confidence(
100         valid_prices, consensus_price, outliers
101     )
102
103     # Step 5: Actualizar estadisticas de proveedores
104     self._update_provider_stats(valid_prices, consensus_price,
105 outliers)
106
107     result = {
108         'consensus_price': float(consensus_price),
109         'method': method,
110         'confidence': float(confidence),
111         'outliers': outliers,
112         'outlier_reasons': [
113             f'{p}: {abs(valid_prices[p] - median_price)/
114 median_price*100:.2f}% '
115             f"del mediana'
116             for p in outliers
117         ],
```

```
115         'provider_stats': {
116             p: self.provider_stats[p]
117             for p in self.providers
118         }
119     }
120
121     self.consensus_history.append({
122         'timestamp': timestamp,
123         'result': result
124     })
125
126     return result
127
128     @staticmethod
129     def _weighted_median(values: np.ndarray, weights: np.ndarray) ->
float:
130         """Calcula mediana ponderada"""
131         sorted_indices = np.argsort(values)
132         sorted_values = values[sorted_indices]
133         sorted_weights = weights[sorted_indices]
134
135         cumsum_weights = np.cumsum(sorted_weights)
136         target_weight = 0.5
137
138         idx = np.argmax(cumsum_weights >= target_weight)
139         return sorted_values[idx]
140
141     def _calculate_confidence(self,
142                             prices: Dict[str, float],
143                             consensus: float,
144                             outliers: List[str]) -> float:
145         """Calcula confianza en consenso"""
146
147         # Menos outliers = mas confianza
148         outlier_ratio = len(outliers) / len(prices)
149
150         # Menos dispersion = mas confianza
151         price_values = np.array(list(prices.values()))
152         coefficient_variation = np.std(price_values) / np.mean(
price_values)
153
154         # Combinar factores
155         confidence = (1 - outlier_ratio) * (1 - min(
```

```
coefficient_variation, 0.1))
156
157     return max(0.1, min(1.0, confidence))
158
159     def _update_provider_stats(self,
160                               prices: Dict[str, float],
161                               consensus: float,
162                               outliers: List[str]):
163         """Actualiza confiabilidad de proveedores"""
164
165         for provider, price in prices.items():
166             self.provider_stats[provider]['total'] += 1
167
168             if provider in outliers:
169                 self.provider_stats[provider]['errors'] += 1
170
171             # Confiabilidad = (1 - error_rate)
172             error_rate = self.provider_stats[provider]['errors'] / \
173                 max(self.provider_stats[provider]['total'],
174 1)
175
176             self.provider_stats[provider]['reliability'] = 1 - min(
error_rate, 1.0)
176
177     def _handle_no_valid_prices(self) -> Dict:
178         """Maneja caso sin precios validos"""
179         return {
180             'consensus_price': None,
181             'method': 'no_valid_data',
182             'confidence': 0.0,
183             'outliers': list(self.providers),
184             'reasoning': 'Todos los proveedores retornaron precios
invalidos'
185         }
```

Listing 6.2: Sistema de consenso multi-proveedor

Capítulo 7

Conclusión: Arquitectura Integrada

7.1. Sistema Completo: Flujo End-to-End

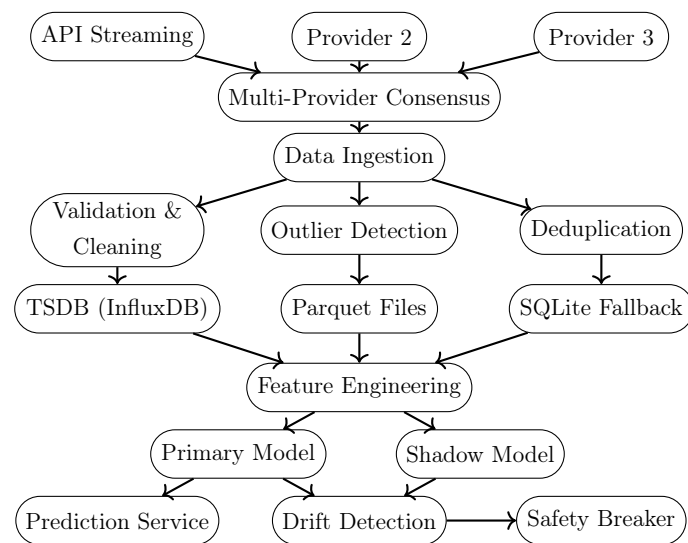


Figura 7.1: Arquitectura integrada del sistema de predicción (Corregida)

7.2. Key Insights y Recomendaciones

7.2.1. Diseño para Producción

1. **Separación de concerns:** Data layer, ML layer, API layer están desacoplados
2. **Resiliencia multi-nivel:** Circuit breakers, fallbacks, caches locales
3. **Observabilidad:** Logging exhaustivo, métricas, alertas automáticas
4. **Validación en todos los niveles:** Pydantic, schema validation, sanity checks

7.2.2. Métricas Críticas de Monitoreo

Métrica	Target	Acción
Latencia P99	< 150ms	Alertar si > 180ms
SLA Violation	< 5 %	Invocar escalación
Model Drift Score	< 0.15	Retrain si > 0.2
Data Gap Duration	< 10 min	Investigar si > 30min
Deduplication Rate	< 1 %	Review si > 2 %
Provider Consensus	> 0.9	Flag si < 0.8

Cuadro 7.1: SLAs y umbrales de monitoreo

7.3. Roadmap de Implementación

1. **Fase 1 (Semanas 1-2):** Implementar validadores Pydantic y tests unitarios
2. **Fase 2 (Semanas 3-4):** Containerización Docker y CI/CD pipeline
3. **Fase 3 (Semanas 5-6):** Data pipeline con TSDB y fallbacks
4. **Fase 4 (Semanas 7-8):** Drift detection y retraining automático
5. **Fase 5 (Semanas 9-10):** Safety systems y circuit breakers
6. **Fase 6 (Semanas 11-12):** Load testing, optimization, deployment

Apéndice A

Referencias Técnicas

A.1. Librerías Recomendadas

- **ML:** scikit-learn, XGBoost, ONNX, SHAP
- **Data:** pandas, numpy, pyarrow, DuckDB
- **API:** FastAPI, Pydantic, uvicorn
- **Storage:** InfluxDB, QuestDB, SQLite
- **Monitoring:** Prometheus, Grafana, ELK Stack
- **Testing:** pytest, hypothesis, locust

A.2. Estándares de Código

- Linting: pylint, flake8, black
- Type hints: mypy
- Documentation: Sphinx, docstrings
- Versionamiento: Git flow