

# Especificación de API Python - Universal Predictor

Ingeniería de Software

February 18, 2026

## Contents

<b>1</b>	<b>Introducción</b>	<b>2</b>
<b>2</b>	<b>Estructuras de Datos (Tipado)</b>	<b>2</b>
2.1	Configuración ( $\Lambda$ ) . . . . .	2
2.2	Entrada Operativa ( $y_t, y_{target}, \tau$ ) . . . . .	2
2.3	Salida del Sistema . . . . .	2
<b>3</b>	<b>Arquitectura Multitenencia (Stateless/Functional Pattern)</b>	<b>3</b>
3.1	Maximización de Throughput (Batching Vectorizado) . . . . .	3
<b>4</b>	<b>Clase Principal: UniversalPredictor (Stateful Wrapper)</b>	<b>4</b>
4.1	Inicialización . . . . .	4
4.2	Método de Ejecución (Paso $t \rightarrow t + 1$ ) . . . . .	4
<b>5</b>	<b>Persistencia (Atomic Snapshotting)</b>	<b>5</b>
<b>6</b>	<b>Manejo de Errores y Excepciones</b>	<b>6</b>

# 1 Introducción

Este documento detalla la implementación en Python de la interfaz abstracta I/O definida en *Predictor\_Estocastico\_IO*. La API expone la clase `UniversalPredictor`, diseñada para entornos de alto rendimiento utilizando JAX para la aceleración numérica.

## 2 Estructuras de Datos (Tipado)

Se utilizan `dataclasses` y `jaxtyping` para garantizar la inmutabilidad y el tipado dimensional estricto de los tensores.

### 2.1 Configuración ( $\Lambda$ )

```
1 from dataclasses import dataclass
2 from typing import Optional
3 from jaxtyping import Float, Array, Bool
4
5 @dataclass(frozen=True)
6 class PredictorConfig:
7     """Vector de Hiperparámetros Lambda."""
8     epsilon: float = 1e-3          # Regularización Entrópica (Sinkhorn)
9     learning_rate: float = 0.01    # Tasa de Aprendizaje JKO
10    log_sig_depth: int = 3        # Profundidad de Firma (Kernel D)
11    wtmm_buffer_size: int = 128   # Memoria WTMM (N_buf)
12    besov_cone_c: float = 1.5     # Cono de Influencia de Besov
13    holder_threshold: float = 0.4 # Umbral Circuit Breaker (H_min)
14    cusum_h: float = 5.0          # Umbral Drift (h)
15    cusum_k: float = 0.5          # Slack (k)
16    volatility_alpha: float = 0.1 # Decaimiento EWMA de Varianza
17
18    # Política de Abandono y Anti-Aliasing
19    staleness_ttl_ns: int = 500_000_000      # TTL Latencia (500ms)
20    besov_nyquist_interval_ns: int = 100_000_000 # Límite Nyquist (100ms) para
estabilidad WTMM
```

### 2.2 Entrada Operativa ( $y_t, y_{target}, \tau$ )

```
1 @dataclass(frozen=True)
2 class MarketObservation:
3     price: Float[Array, "1"]       # y_t (Normalizado o Absoluto)
4     target: Float[Array, "1"]      # y_target (Generalmente price actual)
5     timestamp_ns: int            # Unix Epoch (Nanosegundos)
6
7     def validate_domain(self, sigma_bound: float = 20.0, sigma_val: float = 1.0) -> bool
8     :
9         """Detección de Outliers Catastróficos (> 20 sigma)."""
10        return abs(self.price) <= (sigma_bound * sigma_val)
```

### 2.3 Salida del Sistema

```
1 @dataclass(frozen=True)
2 class PredictionResult:
3     predicted_next: Float[Array, "1"]    # y_{t+1} (Espacio Z-Score)
4
5     # Telemetría de Estado (S_risk)
6     holder_exponent: Float[Array, "1"]   # H_t
7     cusum_drift: Float[Array, "1"]        # G^+
8     distance_toCollapse: Float[Array, "1"] # h - G^+
9     free_energy: Float[Array, "1"]        # F (Energía JKO)
10
11    # Estado del Orquestador
12    weights: Float[Array, "4"]           # [rho_A, rho_B, rho_C, rho_D] (Simplex)
13
14    # Flags de Salud y Control
15    sinkhorn_converged: Bool[Array, "1"] # Convergencia JKO
```

```

16     is_stable: bool                      # is_degraded (False si viola TTL)
17     mode: str                           # "Standard" | "Robust"

```

### 3 Arquitectura Multitenencia (Stateless/Functional Pattern)

Para soportar cientos de activos (Multi-Asset) en un solo servidor, la API soporta un modo puramente funcional. Esto permite gestionar el estado en bases de datos externas de baja latencia (Redis) y compartir el grafo de computación JAX compilado (el Predictor) entre todos los activos.

#### 3.1 Maximización de Throughput (Batching Vectorizado)

Esta arquitectura habilita el uso de `jax.vmap` para procesar lotes de estados de múltiples activos en una sola llamada al hardware, minimizando el impacto del GIL de Python y maximizando la ocupación de la GPU.

```

1 class FunctionalPredictor:
2     """
3         Implementación Stateless para JAX Core.
4         Permite escalar a miles de predictores compartiendo la misma estructura
5             computacional.
6     """
7
8     def __init__(self, config: PredictorConfig):
9         # Compilación JIT única para todos los activos
10        # Habilita vectorización automática (vmap) sobre la dimensión del batch (activos
11    )
12        self.config = config
13        self._core_step = self._core_update_step
14        self._jit_update = jax.jit(self._core_step)
15        self._vmap_update = jax.jit(jax.vmap(self._core_step, in_axes=(0, 0, 0, 0)))
16
17    def init_state(self):
18        """Genera un estado cero inicial (cold state structure)."""
19        return self._initialize_state_structure()
20
21    def step(self, state, obs: MarketObservation) -> tuple[object, PredictionResult]:
22        """
23            Transición de Estado Pura: (S_t, Obs_t) -> (S_{t+1}, Pred_{t+1})
24        """
25        # 1. Validaciones (Outlier, Staleness, Nyquist) logic idéntica a
26        UniversalPredictor
27        # ... logic for freeze_weights flag calculation ...
28
29        # 2. Ejecución Kernel JAX
30        # Zero-Copy: La actualización de búferes ocurre dentro de XLA (
31        dynamic_update_slice)
32        new_state, raw_result = self._jit_update(
33            state, # Estado injectado explícitamente desde Redis/Memoria
34            obs.price,
35            obs.target,
36            freeze_weights=should_freeze
37        )
38
39        # 3. Mapeo de Resultados
40        result = PredictionResult(
41            predicted_next=raw_result.y_next,
42            # ... resto de campos ...
43        )
44
45        return new_state, result
46
47    def step_batch(self, states, obs_batch: MarketObservation):
48        """
49            Procesamiento vectorizado para N activos simultáneos.
50            Utiliza vmap para paralelizar la inferencia y actualización.
51        """
52        # ... logic for batch flags ...
53        new_states, results = self._vmap_update(states, obs_batch.price, obs_batch.
54            target, freeze_flags)
55        return new_states, results

```

## 4 Clase Principal: UniversalPredictor (Stateful Wrapper)

Esta clase envuelve el patrón funcional para casos de uso de un solo activo (Single-Tenant), manteniendo el estado en memoria local (`self._state`).

### 4.1 Inicialización

```
1 class UniversalPredictor:
2     def __init__(self, config: PredictorConfig):
3         """
4             Inicializa el grafo de cómputo JAX (XLA JIT compilation).
5             Asigna memoria estática para los búferes en el dispositivo (VRAM).
6             El estado interno (self._state) contiene los 'jnp.array' persistentes (rolling
7             buffers)
8                 que se actualizarán mediante operaciones funcionales (jnp.roll, lax.
9             dynamic_update)
10                para eliminar la latencia de transferencia de memoria (Zero-Copy).
11                """
12                self.config = config
13                self._state = self._initialize_state() # Estado interno JAX (resident un GPU)
14                self._jit_update = jax.jit(self._core_update_step)
15                self._last_timestamp_ns = 0 # Para cálculo de frecuencia
16
17    def fit_history(self, history: list[float]) -> bool:
18        """
19            Bootstrapping inicial (Protocolo de Cold Start).
20            Procesa el lote histórico para estabilizar los pesos JK0 y llenar los búferes.
21            Requiere un mínimo de N_buf muestras.
22
23            Returns:
24                bool: True si el sistema alcanzó convergencia estable (Sinkhorn + CUSUM).
25            Raises:
26                ValueError: Si el historial es insuficiente (< wtmm_buffer_size).
27                RuntimeError: Si el sistema diverge tras el calentamiento.
28                """
29                if len(history) < self.config.wtmm_buffer_size:
30                    raise ValueError(f"Historial insuficiente. Requerido: {self.config.
31                        wtmm_buffer_size}")
32
33                # Ejecución batch acelerada (jax.lax.scan) para calentar el estado
34                # Simula el paso del tiempo para llenar colas y estabilizar gradientes
35                self._state, final_metrics = self._jit_scan_history(self._state, jnp.array(
36                    history))
37
38                # Validación de Convergencia
39                is_converged = final_metrics.sinkhorn_converged
40                is_stable = final_metrics.cusum_drift < self.config.cusum_h
41
42                if not (is_converged and is_stable):
43                    logger.warning("Cold Start finalizado sin convergencia estable.")
44                    return False
45
46                return True
```

### 4.2 Método de Ejecución (Paso $t \rightarrow t + 1$ )

```
1     def step(self, obs: MarketObservation) -> PredictionResult:
2         """
3             Ejecuta un ciclo completo de predicción.
4             Maneja internamente la validación de dominio y TTL.
5             """
6
7             # 1. Validación de Dominio (Outlier Check)
8             if not obs.validate_domain():
9                 logger.error("Outlier Catastrófico detectado. Ignorando tick.")
10                return self._last_valid_result # Mantiene inercia
11
12                # 2. Check de Abandono (Staleness) y Frecuencia (Anti-Aliasing)
13                current_time = time.time_ns()
14                latency = current_time - obs.timestamp_ns
15                is_stale = latency > self.config.staleness_ttl_ns
```

```

15     # Validación de Frecuencia Nyquist (WTMM Stability)
16     dt_arrival = obs.timestamp_ns - self._last_timestamp_ns
17     is_sparse = (self._last_timestamp_ns > 0) and (dt_arrival > self.config.
18     besov_nyquist_interval_ns)
19
20     if is_sparse:
21         logger.warning(f"FrequencyWarning: Event interval {dt_arrival}ns > Nyquist
22 limit. WTMM spectrum might alias.")
23
24     self._last_timestamp_ns = obs.timestamp_ns
25
26     # 3. Actualización Core (JAX) - Zero-Copy State Management
27     # IMPORTANTE: El buffer de señal reside en GPU/TPU (self._state.signal_buffer).
28     # La actualización se realiza "in-place" funcionalmente usando jax.lax.
29     dynamic_update_slice
30     # o jnp.roll dentro del kernel compilado para evitar transferencias CPU <-> VRAM
31
32     # Si hay staleness o sparsity excesiva, se congelan pesos para no degradar la
33     # geometría.
34     should_freeze = is_stale or is_sparse
35
36     new_state, result_data = self._jit_update(
37         self._state,
38         obs.price,
39         obs.target,
40         freeze_weights=should_freeze,
41         # No se pasa history_buffer explícitamente, ya vive en _state
42     )
43
44     self._state = new_state
45
46     # 4. Empaquetado de Resultados
47     return PredictionResult(
48         predicted_next=result_data.y_next,
49         holder_exponent=result_data.H_t,
50         sinkhorn_converged=result_data.converged,
51         is_stable=not (is_stale or is_sparse),
52         # ... mapeo resto de campos
53     )

```

## 5 Persistencia (Atomic Snapshotting)

El sistema implementa persistencia binaria protegida por checksum.

```

1 import hashlib
2 import msgpack
3
4 def save_snapshot(self, filepath: str):
5     """
6         Exporta el estado interno Sigma_t a formato binario (MessagePack).
7         Incluye Checksum SHA-256 al final del archivo.
8     """
9     # Serialización de tensores JAX a bytes
10    state_dict = self._serialize_jax_state(self._state)
11
12    # Segmentación Modular (K-Blocks)
13    payload = {
14        "timestamp": time.time_ns(),
15        "config": asdict(self.config),
16        "global": state_dict["global"], # rho, G+, ema
17        "kernels": {
18            "A": state_dict["kernel_a"],
19            "B": state_dict["kernel_b"],
20            "C": state_dict["kernel_c"],
21            "D": state_dict["kernel_d"]
22        }
23    }
24
25    data_bytes = msgpack.packb(payload)
26    checksum = hashlib.sha256(data_bytes).hexdigest()

```

```

27
28     with open(filepath, "wb") as f:
29         f.write(data_bytes)
30         f.write(checksum.encode('utf-8')) # Append hash
31
32     def load_snapshot(self, filepath: str):
33         """
34             Carga estado. Valida SHA-256 antes de deserializar.
35             Lanza IntegrityError si falla la validación.
36         """
37         with open(filepath, "rb") as f:
38             content = f.read()
39
40             data_bytes = content[:-64] # Todo menos los últimos 64 bytes (SHA256 hex)
41             stored_checksum = content[-64:].decode('utf-8')
42
43             computed = hashlib.sha256(data_bytes).hexdigest()
44             if computed != stored_checksum:
45                 raise ValueError("Snapshot corrupto: Checksum mismatch.")
46
47             payload = msgpack.unpackb(data_bytes)
48             self._state = self._deserialize_jax_state(payload)

```

## 6 Manejo de Errores y Excepciones

- **DomainError:** Se lanza (o se loguea crítico) si  $y_t$  excede los límites (Outlier Catastrófico).
- **StalenessWarning:** Emitido mediante el sistema de logging estándar de Python cuando se activa la protección TTL.
- **FrequencyWarning:** Alerta si la tasa de arribo de eventos cae por debajo del límite de Nyquist para el análisis de Besov.
- **IntegrityError:** Fallo crítico en la carga de snapshot. El sistema debe abortar y solicitar reinicio en frío.