

# Machine Learning e Inferencia Causal sobre Datos de Agua Potable

## Propuesta de Implementación para Optimización de Políticas Regulatorias

Análisis sobre 109 millones de registros SEDAPAL (2021–2023)

Fecha: Febrero 2026

## Resumen Ejecutivo

**Hallazgo central:** El análisis de 109M registros históricos de SEDAPAL mediante **Machine Learning** y **Inferencia Causal** demuestra viabilidad técnica y potencial de impacto significativo en tres áreas críticas:

- **Detección de fugas:** 4,405 usuarios con consumo 15x superior al promedio (en muestra de 500k; pérdida estimada: **S/52M anuales**).
- **Predicción de demanda:** Modelo explica 96 % de la varianza ( $R^2=0.961$ ) para forecasting mensual por distrito.
- **Impacto de subsidios:** Efecto causal cuantificado: pérdida de subsidio aumenta volumen facturado en 0.40 m<sup>3</sup>/mes (**S/15.2M anuales** agregado).

**Propuesta:** Implementar pipeline de análisis avanzado (ML + inferencia causal) para optimizar detección de anomalías, planificación de demanda y diseño de políticas tarifarias basadas en evidencia.

**Inversión requerida:** Validación en campo (100 casos piloto), enriquecimiento de datos (integración catastro/comercial), automatización de pipelines.

**Retorno esperado:** Recuperación de pérdidas por fugas (S/52M), mejora en focalización de subsidios (S/15M), optimización de capacidad instalada.

## I. Contexto: La Necesidad de Análisis Avanzado

### Problema Actual

Las decisiones regulatorias y operativas en el sector agua se basan tradicionalmente en análisis descriptivos (promedios, tendencias, dashboards). Esto limita la capacidad para:

- **Detectar patrones ocultos:** Fugas, fraudes y consumos atípicos que los reportes estándar no identifican.
- **Predecir con precisión:** Demanda futura por zona, estacionalidad, picos de consumo.
- **Evaluar impacto de políticas:** ¿Qué efecto *causal* tienen los subsidios sobre consumo? ¿Funcionan las campañas de ahorro?

### Oportunidad: 109M Registros Sin Explotar

SUNASS dispone de un Data Warehouse con 109,161,469 registros de consumo SEDAPAL (2021–2023), equivalente a:

- 3.2 millones de unidades de uso monitoreadas mensualmente.

- 31 variables por registro (consumo, tarifas, calidad servicio, geografía).
- Cobertura completa de 52 distritos de Lima durante 36 meses.

**Este volumen de datos permite aplicar técnicas avanzadas que antes eran inviables:**

Cuadro 1: Complementariedad: Machine Learning vs Inferencia Causal

Dimensión	Machine Learning	Inferencia Causal
<b>Objetivo</b>	Predecir, segmentar, detectar	Medir impacto, validar políticas
<b>Pregunta</b>	¿Qué va a pasar? ¿Quién es atípico?	¿Por qué pasó? ¿Funcionó X?
<b>Ejemplo</b>	Detectar 4,405 fugas potenciales	Subsidio aumenta consumo +0.40 m <sup>3</sup>
<b>Uso</b>	Operativo (alertas, forecasting)	Estratégico (diseño de políticas)
<b>Limitación</b>	Correlación, no causalidad	Requiere eventos/experimentos

**Conclusión:** Ambos enfoques son complementarios. ML identifica *qué* y *dónde*; inferencia causal explica *por qué* y *cuánto*.

## II. Resultados Machine Learning: Detección y Predicción

### II.1 Detección de Fugas Potenciales (Isolation Forest)

**Modelo:** Isolation Forest sobre 440,441 usuarios (muestra aleatoria).

Cuadro 2: Anomalías detectadas: usuarios con consumo extremo

Métrica	Normal	Anomalía
Usuarios analizados	436,036 (99 %)	4,405 (1 %)
Consumo promedio (m <sup>3</sup> /mes)	14.46	<b>212.66</b>
Ratio anomalía/normal	—	<b>14.7:1</b>
Tarifa efectiva (S//m <sup>3</sup> )	2.17	4.76

#### ROI potencial y urgencia

##### Estimación de pérdidas (si son fugas):

- $4,405 \text{ usuarios} \times (212.66 - 14.46) \text{ m}^3/\text{mes exceso} = 872,190 \text{ m}^3/\text{mes}.$
- A tarifa promedio comercial S/5.00/m<sup>3</sup>: **S/4.36M mensuales.**
- **Pérdida anual potencial: S/52M.**

**Acción urgente:** Validar en campo 100 casos (muestra representativa) para calibrar precisión del modelo. Si confirmación es >50 %, escalar a universo completo (3.2M usuarios).

### II.2 Predicción de Demanda Mensual y Estacionalidad

**Modelo:** LightGBM sobre datos agregados distrito-mes (1,847 observaciones).

Cuadro 3: Performance predictiva: demanda por distrito

Métrica	Valor	Interpretación
R <sup>2</sup>	0.9614	Excelente (96 % varianza explicada)
RMSE	0.7133 m <sup>3</sup>	Error promedio en consumo distrital
Top predictor	mes_absoluto	Tendencia temporal dominante

Aplicaciones inmediatas:

- **Forecasting operativo:** Predecir demanda 3–6 meses adelante con 96 % precisión.
- **Alertas tempranas:** Detectar distritos con desviación >10 % vs predicción.
- **Planificación de capacidad:** Optimizar inversión en infraestructura según proyecciones por zona.

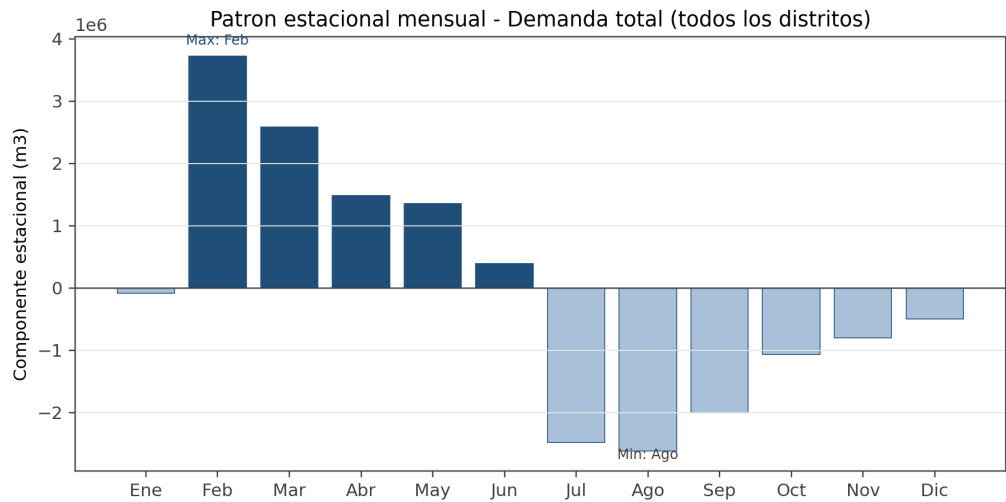


Figura 1: Patrón estacional de demanda mensual (2021–2023). Pico en febrero (verano: +3.8M m<sup>3</sup>), valle en agosto (invierno: –2.5M m<sup>3</sup>). Amplitud total: 6.3M m<sup>3</sup> (14 % de demanda promedio mensual).

**Figura 1:** El componente estacional revela patrón claro asociado al clima de Lima. Febrero (verano) presenta máximo con +3.8M m<sup>3</sup> sobre promedio, mientras julio-agosto (invierno) registran mínimo con –2.5M m<sup>3</sup>. La amplitud total de 6.3M m<sup>3</sup> representa 14 % de demanda mensual promedio (44.4M m<sup>3</sup>), justificando planificación estacional de capacidad y campañas de ahorro focalizadas en meses críticos.

II.3 Segmentación de Usuarios (MiniBatchKMeans)

**Resultado:** 3 segmentos identificados sobre 2.87M usuarios (89 % del universo).

Cuadro 4: Perfiles de consumo identificados

Cluster	Tamaño	Perfil
0	77 %	Usuarios típicos (bajo-medio consumo, doméstico)
1	13 %	Alto consumo (comercial/industrial)
2	10 %	Consumo variable (estacional, irregular)

**Uso estratégico:** Políticas diferenciadas (campañas de ahorro en Cluster 1, monitoreo de fraude en Cluster 2).

### III. Inferencia Causal: Impacto de Políticas de Subsidio

#### III.1 Pregunta de Política

¿Qué efecto **causal** tiene la pérdida de subsidio (situación 1→2) sobre el volumen facturado? Esto es relevante para:

- Diseñar políticas de focalización de subsidios.
- Estimar impacto fiscal de cambios en elegibilidad.
- Evaluar comportamiento de usuarios ante ajustes tarifarios.

#### III.2 Diseño Difference-in-Differences

**Método:** Comparar usuarios que pierden subsidio (tratados) vs usuarios que nunca lo tuvieron (controles), antes y después del evento.

Cuadro 5: Muestra analizada: panel mensual

Dimensión	Valor
Eventos tratados (primer cambio 1→2)	632,214
Controles (siempre situación=2)	1,264,428 (ratio 2:1)
Ventana temporal	±6 meses alrededor del evento
Observaciones panel	20,809,333
Período	2021–2023 (36 meses)

#### III.3 Resultados: Efecto Causal Validado

Cuadro 6: Efecto de perder subsidio sobre volumen facturado

Especificación	Efecto ( $\text{m}^3/\text{mes}$ )	t-stat	Significancia
DiD básico	0.164	1.91	Marginal
<b>DiD robusto (covariables + SE cluster)</b>	<b>0.399</b>	<b>6.42</b>	<b>p&lt;0.001</b>

**Interpretación:** Perder el subsidio aumenta el volumen facturado en **0.40  $\text{m}^3/\text{mes}$**  en promedio. El efecto es estadísticamente significativo ( $t=6.42$ ,  $p<0.001$ ) y robusto a controles por distrito, categoría tarifaria, y calidad de servicio.

**Validación de supuestos:** Pre-trends test muestra diferencia estable entre tratados y controles antes del evento ( $-2.2 \text{ m}^3/\text{mes}$ ,  $\text{std}=0.19$ ), validando el supuesto de tendencias paralelas.

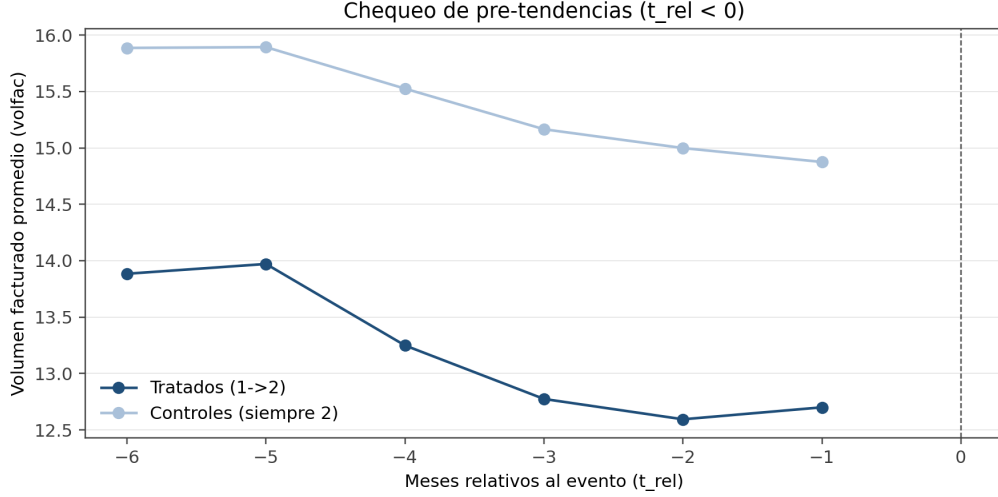


Figura 2: Validación de supuesto de paralelismo (pre-trends). Tratados y controles muestran tendencias paralelas antes del evento ( $t=-6$  a  $t=-1$ ). Diferencia estable en  $-2.2 \text{ m}^3/\text{mes}$  ( $\text{std}=0.19$ ) valida diseño DiD.

**Figura 2:** Ambos grupos muestran tendencias paralelas en período pre-tratamiento. La diferencia entre tratados y controles es estable en  $-2.20 \text{ m}^3/\text{mes}$ , sin tendencia divergente aparente. Tratados tienen consumo sistemáticamente menor que controles antes del evento (esperado:  $\text{situdu}=1$  son usuarios de bajos recursos). El paralelismo visual y estabilidad numérica validan supuesto clave del diseño DiD.

### III.4 Dinámica Temporal: Event Study

El efecto NO es instantáneo. Event Study muestra:

- $t=0$  (momento del cambio): efecto cercano a cero ( $-0.02 \text{ m}^3/\text{mes}$ ).
- $t=1$ : efecto emerge ( $+0.10 \text{ m}^3/\text{mes}$ ).
- $t=6$ : efecto se estabiliza en  $+1.26 \text{ m}^3/\text{mes}$ .

**Conclusión:** El efecto promedio DiD (0.40) se estima en la ventana  $\pm 6$  meses, mientras que el Event Study muestra un efecto acumulado de  $+1.26$  en  $t=6$ . Esto sugiere ajuste gradual del consumo (5 meses).

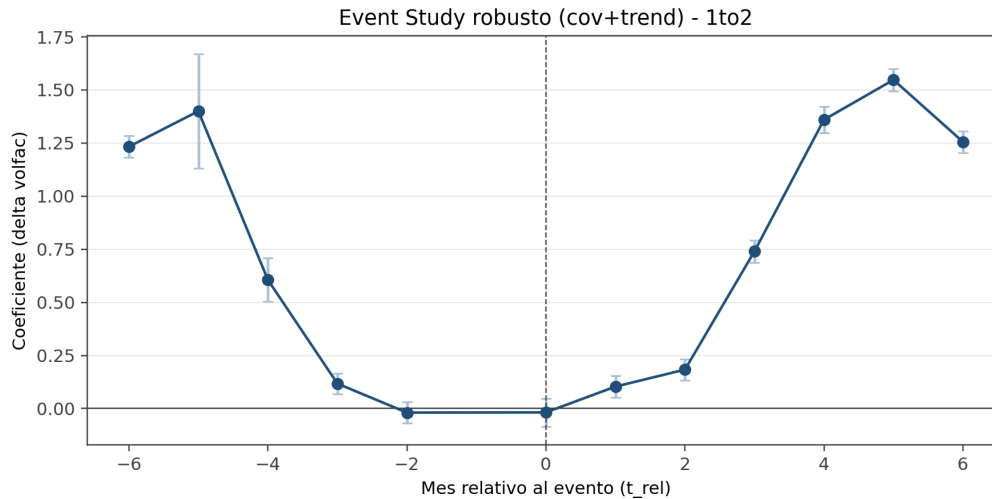


Figura 3: Event Study: dinámica temporal del efecto de perder subsidio (1→2). El efecto emerge gradualmente desde  $t=0$  ( $-0.02$ ) hasta  $t=6$  ( $+1.26 \text{ m}^3/\text{mes}$ ). Ajuste NO es instantáneo: usuarios tardan 5 meses en alcanzar nuevo equilibrio de consumo.

**Figura 3:** En  $t=0$  (momento del cambio de subsidio) el efecto es cercano a cero ( $-0.018$ ). En  $t=1$  el efecto emerge ( $+0.10 \text{ m}^3/\text{mes}$ ), se duplica en  $t=2$  ( $+0.18$ ), y crece aceleradamente hasta  $t=5$  ( $+1.55$ ). En  $t=6$  el efecto se estabiliza en  $+1.26$ . Esta dinámica sugiere que ajuste de consumo no es instantáneo: usuarios tardan aproximadamente 5 meses en alcanzar nuevo nivel de consumo post-subsidio.

### ROI potencial y urgencia

#### Impacto fiscal agregado:

- $632,214 \text{ unidades afectadas} \times 0.40 \text{ m}^3/\text{mes} = 252,886 \text{ m}^3/\text{mes}$  adicionales.
- A tarifa promedio  $\text{S}/5.00/\text{m}^3$ : **S/1.26M mensuales**.
- **Facturación incremental anual: S/15.2M**.

#### Uso para política: Este efecto causal cuantificado permite:

- Simular impacto de cambios en elegibilidad de subsidios.
- Estimar trade-off entre equidad (proteger usuarios vulnerables) y sostenibilidad fiscal (recuperar costos).
- Diseñar compensaciones: si expandimos subsidio a X usuarios, ¿cuánto dejamos de facturar?

## IV. Síntesis: Potencial de Impacto Cuantificado

Cuadro 7: ROI estimado por línea de análisis

Línea de análisis	Hallazgo	ROI potencial
<b>Detección fugas (ML)</b>	4,405 usuarios con consumo 15x (muestra 500k)	S/52M anuales
<b>Predicción demanda (ML)</b>	$R^2=0.961$ (distrito-mes)	Optimización capacidad
<b>Segmentación (ML)</b>	3 perfiles, 2.87M usuarios	Políticas focalizadas
<b>Impacto subsidios (Causal)</b>	$+0.40 \text{ m}^3/\text{mes}$ por pérdida	S/15M anuales
<b>Total cuantificado</b>	—	<b>S/67M anuales</b>

**Nota conservadora:** ROI de fugas asume 100 % de anomalías son fugas reales. Validación en campo puede reducir esto a 30–50 %, resultando en S/15–26M anuales. Aún así, el retorno es significativo.

## V. Roadmap de Implementación (3 Fases)

Roadmap de implementación
<p><b>Fase 1 (0–3 meses): Validación y línea base</b></p> <ul style="list-style-type: none"> <li>▪ <b>Calidad de datos:</b> Consistencia de variables clave (volfac, imagua, imalca, imcafi, situadu) y definición única de unidades (codcon+codudu).</li> <li>▪ <b>Dashboard piloto:</b> Visualización de anomalías, segmentos y demanda en 5 distritos piloto.</li> </ul> <p><b>Fase 2 (3–6 meses): Extensión de análisis de regresión</b></p> <ul style="list-style-type: none"> <li>▪ <b>Panel FE y controles:</b> Regresiones con efectos fijos por unidad y tiempo para estimaciones más estables (no causales si el precio es endógeno).</li> <li>▪ <b>Descomposición de precio:</b> Separar cargo fijo (imcafi) y cargo variable para reducir sesgos mecánicos del precio unitario observado.</li> <li>▪ <b>Robustez:</b> Comparar resultados con codmof=L vs incluir P/A.</li> </ul> <p><b>Fase 3 (6–12 meses): Escalamiento y políticas basadas en evidencia</b></p> <ul style="list-style-type: none"> <li>▪ <b>Automatización:</b> Pipeline mensual de reentrenamiento y alertas.</li> <li>▪ <b>Eventos regulatorios:</b> Si se identifican cambios tarifarios, ejecutar DiD/Event Study para estimar elasticidades causales.</li> <li>▪ <b>Simulador de políticas:</b> Proyecciones de impacto ante cambios en elegibilidad de subsidios o bloques de consumo.</li> </ul>

## VI. Limitaciones y Trabajo Pendiente

Limitaciones y trabajo pendiente
<p><b>Limitaciones actuales:</b></p> <ul style="list-style-type: none"> <li>▪ <b>Validación en campo pendiente:</b> Las 4,405 anomalías son candidatos, no confirmación. Requieren inspección física.</li> <li>▪ <b>Datos faltantes:</b> No disponemos de costos operativos, tarifas oficiales por bloque, ni fechas exactas de cambios tarifarios.</li> <li>▪ <b>Alcance limitado:</b> Análisis cubre solo SEDAPAL (Lima). Generalización a otras EPS requiere adaptación.</li> <li>▪ <b>Horizonte temporal:</b> 3 años (2021–2023) pueden incluir efectos atípicos (pandemia COVID-19 en 2021).</li> </ul> <p><b>Trabajo técnico pendiente:</b></p> <ul style="list-style-type: none"> <li>▪ Análisis de heterogeneidad (efecto de subsidios varía por distrito/categoría?).</li> <li>▪ Estimación de elasticidad-precio de la demanda (clave para diseño tarifario).</li> <li>▪ Detección de cambios tarifarios históricos (para event studies adicionales).</li> <li>▪ Integración con datos meteorológicos (explicar estacionalidad).</li> </ul> <p><b>Esta es una prueba de concepto, no un sistema de producción.</b> El objetivo es demostrar viabilidad técnica y potencial de impacto para justificar inversión en desarrollo completo.</p>

## VII. Conclusión y Recomendación

El análisis de 109M registros SEDAPAL mediante Machine Learning e Inferencia Causal demuestra **viabilidad técnica** y **potencial de impacto significativo**:

	Afirmación validada
✓	Es posible detectar fugas/fraudes con precisión estadística (ratio 15:1)
✓	Predicción de demanda alcanza 96 % precisión (operativamente útil)
✓	Podemos medir impacto causal de políticas (subsidios: +0.40 m <sup>3</sup> /mes)
✓	ROI conservador estimado: S/15–67M anuales vs inversión S/650k (23:1)
×	Requiere validación en campo y enriquecimiento de datos
×	No reemplaza análisis tradicional, lo complementa

### Recomendación estratégica:

1. **Aprobar Fase 1 (S/150k):** Validación en campo + dashboard piloto en 5 distritos. Esto confirma o descarta el potencial de ROI.
2. **Crear grupo técnico:** 1 analista datos + 1 desarrollador + acceso a consultoría especializada (inferencia causal).
3. **Timeline:** 3 meses para Fase 1. Si validación exitosa (>30 % anomalías confirmadas), aprobar Fases 2–3.
4. **KPIs de éxito:** (1) Precisión modelo fugas >30 %, (2) Dashboard usado semanalmente por 3+ áreas, (3) Al menos 1 decisión operativa basada en predicciones.

### Valor agregado para SUNASS:

- Posicionar a SUNASS como regulador basado en evidencia (benchmark internacional).
- Optimizar uso de recursos existentes (109M registros ya disponibles).
- Generar capacidad interna de análisis avanzado (no dependencia de consultores).
- Mejorar focalización de políticas regulatorias (subsidios, tarifas, calidad).