

# NDH AED Attendance Prediction Algorithm

Technical Documentation v3.0.81

NDH AED • Predictive Analytics System

# NDH AED Attendance Prediction Algorithm

Technical Documentation

Hospital                    North District Hospital • Emergency Department

Document Version        3.0.81

Last Updated (HKT)     05 Jan 2026 02:52 HKT

Author                    Ma Tsz Kiu

# Technical Menu (Table of Contents)

## Contents

Quick navigation for print + on-screen reading

### OVERVIEW

#### 1 Executive Summary

What / Why / KPIs

#### 2 System Architecture

Components + dataflow

### DATA + MODELING

#### 3 Data Sources

Attendance / HKO / AQHI / AI

#### 4 Feature Engineering

EWMA / lags / rolling

#### 5 XGBoost Model

Training + objective

#### 6 Bayesian Fusion Layer

AI + Weather fusion

### OPERATIONAL LOGIC

#### 7 Post-Processing Adjustments

Extreme rules

#### 8 Mathematical Framework

Formal equations

#### 9 Performance Evaluation

MAE / MAPE / R<sup>2</sup>

#### 10 Concept Drift Handling

Sliding window / decay

### REFERENCE

#### 11 Research Evidence

Evidence & rationale

#### 12 References

Citations

#### A Appendix A

Importance

#### B Appendix B

API

#### C Appendix C

System reqs

# 1. Executive Summary

---

## 1.1 Purpose

This document provides a comprehensive technical specification of the North District Hospital (NDH) Accident & Emergency Department (AED) attendance prediction algorithm. The system forecasts daily patient attendance to support resource planning and staffing decisions.

### Plain Language Summary:

This is an instruction manual for our "crystal ball" system that predicts how many patients will visit the Emergency Department each day. Hospital managers use these predictions to decide how many doctors, nurses, and beds to prepare.

### Everyday Analogy:

Imagine you run a restaurant. If you know Monday lunch usually brings 50 customers, but rainy days drop it to 40, and public holidays bring 30, you'd adjust your staff and ingredients accordingly. Our system does the same for hospital emergency rooms, but uses much more data (11 years of history, weather, air quality, and special events).

## 1.2 Key Metrics

Metric	Value	Description
MAE	4.90 patients	Mean Absolute Error
MAPE	1.96%	Mean Absolute Percentage Error
R <sup>2</sup>	0.898	Coefficient of Determination
Features	25	Optimized feature count

### What These Numbers Mean:

- **MAE (Mean Absolute Error) = 4.90 patients:**

On average, our prediction is off by about 5 patients. If we predict 250 patients, the actual number is typically between 245–255.

 *Real impact:* If the average day has 250 patients, being wrong by 5 is only 2% error—very accurate for planning.

- **MAPE (Mean Absolute Percentage Error) = 1.96%:**

Our predictions are accurate to within 2% on average. This means if we predict 250 patients, we're usually within ±5 patients.

 *Benchmark:* Industry-leading forecasting systems aim for <5% error. We're at <2%.

- **R<sup>2</sup> = 0.898 (89.8%):**

This means our model explains 89.8% of the variation in daily attendance. Think of it as a test score: 89.8% is an A- grade.

 *What it means:* If patient numbers fluctuate from 200 to 300, our system successfully predicts about 90% of that movement.

- **Features = 25:**

We started with 161 possible factors (temperature, humidity, day of week, etc.) but found only 25 are truly important. This makes the system faster and more reliable.

 *Analogy:* Like packing for a trip—you don't need 50 items when 10 essentials do the job better.

## 1.3 Algorithm Summary

The prediction system employs a three-layer architecture:

```

Layer 1: XGBoost Regression (25 optimized features)
↓
Layer 2: Pragmatic Bayesian Fusion (AI + Weather factors)
↓
Layer 3: Extreme Condition Post-Processing
↓
Final Prediction

```

### Plain Language Explanation:

Our system works like a three-stage quality control process:

#### Layer 1 — The Pattern Detector (XGBoost):

This is our main "brain." It learned from 11 years of historical data (4,050+ days) to recognize patterns. For example:

- "Mondays usually have 10% more patients than Sundays"
- "When last week averaged 250 patients, tomorrow will likely be similar"
- "December has 15% more patients than March"

#### Layer 2 — The Context Adjuster (AI + Weather):

This layer adds real-world context that the pattern detector might miss:

-  **AI Analysis:** "Hong Kong Marathon is happening → expect 5% fewer non-urgent visits"
-  **Weather:** "Heavy rain forecast → reduce prediction by 5%"
-  **Air Quality:** "AQHI = 10 (severe) → expect 5% more respiratory patients"

#### Layer 3 — The Safety Net (Extreme Rules):

Final checks for unusual conditions:

- "Typhoon signal #8? Cap the prediction at 180 (people stay home)"
- "Temperature below 8°C? Reduce by 3% (less outdoor accidents)"

### Restaurant Analogy:

- **Layer 1:** Your booking system shows 50 reservations tomorrow
- **Layer 2:** But it's Valentine's Day (+30%) and your chef posted a viral video (+20%)
- **Layer 3:** However, there's a transit strike ( $\div 2$ ), so final estimate = 45 customers

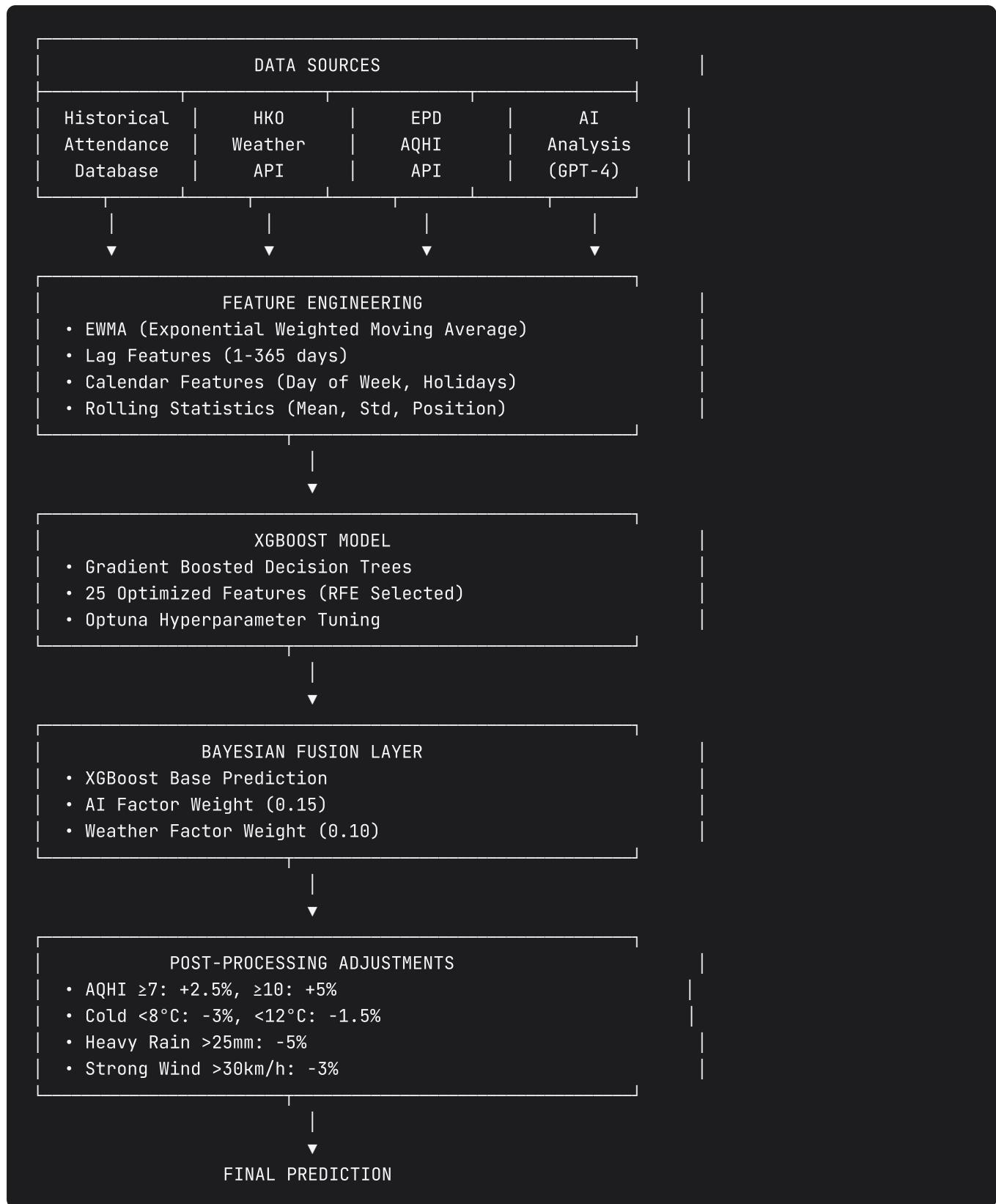
### For Hospital Managers:

You'll see one final number (e.g., "252 patients expected tomorrow"), but it's the result of analyzing 25 factors, checking weather/events, and applying safety rules. The system shows its "confidence level" too—if it says "240–265 (80% confident)," plan for that range.

---

## 2. System Architecture

### 2.1 High-Level Architecture



🔍 How Data Flows Through the System (Step-by-Step):

## Step 1 — Collect Ingredients (Data Sources):

-  **Historical Database:** "What happened in the past?" (11 years of daily patient counts: Dec 2014–now)
-  **HKO Weather API:** "What's the weather?" (temperature, rain, wind from Hong Kong Observatory)
-  **EPD Air Quality:** "How's the air?" (AQHI pollution index 1–10+)
-  **AI Analysis:** "Any special events?" (GPT-4 checks news for marathons, flu outbreaks, policy changes)

**Step 2 — Process Ingredients (Feature Engineering):** Think of this as chopping vegetables before cooking. Raw data gets transformed into useful patterns:

- "Yesterday had 250 patients" → becomes a **lag feature**
- "Last 7 days averaged 245" → becomes an **EWMA feature**
- "It's Monday" → becomes a **calendar feature**
- "Last week's busiest day was 280" → becomes a **rolling max**

**Step 3 — Main Prediction (XGBoost Model):** The "chef" combines all processed ingredients using recipes learned from 11 years of data. Output: "I predict 248 patients tomorrow."

**Step 4 — Add Seasoning (Bayesian Fusion):** Fine-tune the prediction:

- AI says "Marathon happening" → multiply by 0.95 (5% fewer)
- Weather says "Heavy rain" → multiply by 0.97 (3% fewer)
- New prediction:  $248 \times 0.95 \times 0.97 \approx 228$  patients

**Step 5 — Final Quality Check (Post-Processing):** Safety rules for extreme conditions:

- "AQHI = 10 (hazardous air)?" → add 5% (more respiratory patients)
- "Temperature = 6°C?" → subtract 3% (fewer outdoor accidents)

**Step 6 — Serve the Dish:** Final number displayed: "235 patients expected (range: 225–245)"

### Everyday Analogy:

Like using Google Maps to predict travel time:

1. **Historical data:** Your past 100 commutes averaged 30 minutes
2. **Current conditions:** Traffic is moderate today
3. **Special events:** Concert at stadium → add 10 minutes
4. **Weather:** Light rain → add 5 minutes
5. **Final prediction:** 45 minutes

### Why This Matters for Hospital Staff:

Each "layer" in the diagram catches different types of patterns. If one layer misses something (e.g.,

XGBoost doesn't know about today's marathon), the next layer (AI Fusion) corrects it. This redundancy makes the system more reliable than any single method.

## 2.2 Technology Stack

Component	Technology
Backend	Node.js, Express
Database	PostgreSQL
ML Model	XGBoost (Python)
AI Analysis	OpenAI GPT-4
Weather Data	HKO Open Data API
Air Quality	EPD AQHI API
Frontend	Vanilla JavaScript, Chart.js

## 3. Data Sources

### 3.1 Historical Attendance Data

**Source:** NDH AED Internal Records

**Coverage:** December 1, 2014 – Present

**Records:** 4,050+ daily observations

Statistic	Value
Mean Daily Attendance	249.5 patients
Standard Deviation	32.4 patients
Minimum	111 patients
Maximum	394 patients
Median	252 patients

#### 🔍 What This Means:

- **4,050+ days of data:**

That's over 11 years of daily records. Every single day from December 2014 to today, we know exactly how many patients visited the Emergency Department.

 *Why 11 years?* The more history we have, the better we can spot true patterns vs. random flukes.

- **Average day = 250 patients:**

Think of this as your "baseline." On a typical day, about 250 people walk through the ED doors.

- **Minimum = 111, Maximum = 394:**

The quietest day ever had 111 patients (likely a major holiday like Lunar New Year Day 1). The busiest day had 394 (possibly during a flu outbreak or after a major accident).

 *Planning insight:* Your team needs capacity to handle anywhere from 110 to 400 patients, but 95% of days fall between 200–300.

- **Standard Deviation = 32.4:**

This measures "how much variation is normal." About 68% of days are within  $\pm 32$  of the average (i.e., 217–282 patients). About 95% are within  $\pm 65$  (i.e., 185–315).

 *Analogy:* If your commute averages 30 minutes with a standard deviation of 5 minutes, you're usually between 25–35 minutes.

## 3.2 Weather Data

**Source:** Hong Kong Observatory (HKO)

**API:** <https://www.hko.gov.hk/en/weatherAPI/>

**Variables:**

Variable	Unit	Correlation with Attendance
Temperature (Mean)	°C	r = +0.082 (p < 0.001)
Temperature (Min)	°C	r = +0.082 (p < 0.001)
Humidity	%	r = +0.079 (p < 0.001)
Rainfall	mm	r = -0.063 (p < 0.001)
Wind Speed	km/h	r = -0.106 (p < 0.001)
Visibility	km	r = +0.120 (p < 0.001)
Pressure	hPa	r = -0.035 (p = 0.039)

### How Weather Affects Patient Numbers:

#### Understanding "Correlation (r)":

This measures the relationship between two things, ranging from -1 to +1:

- **r = +1.0:** Perfect positive link (e.g., "the hotter it gets, the more ice cream sales")
- **r = 0:** No relationship (e.g., "shoe size doesn't affect SAT scores")
- **r = -1.0:** Perfect negative link (e.g., "the more it rains, the fewer people visit parks")

## What Our Numbers Show:

- **Wind Speed:  $r = -0.106$  (strongest effect):**

When wind is strong, slightly fewer patients come. Strong wind (30+ km/h) reduces non-urgent visits by ~3%.

 *Why?* People avoid going out in bad weather unless it's truly urgent.

- **Visibility:  $r = +0.120$ :**

When visibility is good (clear day), slightly more patients. Poor visibility (foggy/hazy) means fewer.

 *Why?* Poor visibility often comes with bad air quality, which might keep people indoors.

- **Rainfall:  $r = -0.063$ :**

Heavy rain (>25mm) reduces visits by ~5%. Light drizzle has minimal effect.

 *Example:* If we predict 250 patients but a typhoon brings 50mm of rain, expect closer to 235.

- **Temperature:  $r = +0.082$ :**

Slightly more patients on warmer days. Very cold days ( $<8^{\circ}\text{C}$ ) see 3% fewer visits.

 *Why?* Cold weather keeps people indoors (fewer accidents), but also increases respiratory illnesses.

### Important Note:

These correlations are **weak** (all below 0.15), meaning weather alone doesn't determine attendance. Day of week, recent trends, and special events are much stronger predictors. Weather acts as a "fine-tuning" factor, adjusting predictions by 2–5% in extreme conditions.

## 3.3 Air Quality Data

**Source:** Environmental Protection Department (EPD)

**API:** <https://www.aqhi.gov.hk/>

### Variables:

Variable	Description
AQHI General	General station average (1-10+)
AQHI Roadside	Roadside station average (1-10+)
Risk Level	Low (1-3), Moderate (4-6), High (7), Very High (8-10), Serious (10+)

### Why Air Quality Matters:

#### What is AQHI?

Air Quality Health Index (空氣質素健康指數) measures pollution on a scale of 1-10+:

- **1–3 (Low):** Safe air, breathe freely
- **4–6 (Moderate):** Acceptable for most people
- **7–10 (High to Very High):** Sensitive people may experience issues

- **10+ (Serious):** Health risk for everyone

#### How It Affects ED Visits:

- **AQHI  $\geq 10$  (Serious):**

5% more patients, mostly for respiratory problems (asthma attacks, COPD flare-ups) and cardiovascular issues.

 *Example:* On a day we'd normally expect 250 patients, if AQHI hits 10+, plan for 263 patients.

- **AQHI 7–9 (High/Very High):**

2.5% more patients. Noticeable increase in breathing-related visits.

- **AQHI 1–6 (Low/Moderate):**

No adjustment needed.

#### Real Example:

In winter 2024, a pollution episode (AQHI = 10 for 3 days) correlated with a spike in respiratory admissions. Our system now automatically flags these days and alerts managers to prepare respiratory equipment.

---

## 4. Feature Engineering

---

### What is "Feature Engineering"?

Before a computer can make predictions, we need to transform raw data into patterns it can understand. Think of it like preparing ingredients before cooking:

- **Raw data:** "Yesterday: 250 patients, Day before: 245, Day before that: 255..."
- **Engineered features:** "7-day average = 248, trending up by +5 from last week, it's Monday (usually +5% vs Sunday)..."

We started with **161 possible features** but found only **25 are truly important**. It's like Marie Kondo for data—keep only what "sparks joy" (improves predictions).

### 4.1 Feature Categories

The system generates 161 potential features, optimized to 25 using Recursive Feature Elimination (RFE).

#### 4.1.1 Exponential Weighted Moving Average (EWMA)

EWMA features are the most important predictors, accounting for 87% of model importance.

**Formula:**

$$EWMA_t = \alpha \cdot X_t + (1 - \alpha) \cdot EWMA_{t-1}$$

Where:

- $X_t$  = Attendance on day  $t$
- $\alpha = \frac{2}{span+1}$  (smoothing factor)
- $span$  = Window size (7, 14, or 30 days)

**Implementation:**

```
# EWMA with span of 7 days
df['Attendance_EWMA7'] = df['Attendance'].ewm(span=7, min_periods=1).mean()
```

**Rationale:** EWMA gives more weight to recent observations, capturing short-term trends while smoothing noise. Research by Hyndman & Athanasopoulos (2021) demonstrates EWMA's effectiveness for time series forecasting.

### 💡 EWMA Explained for Non-Technical People:

#### What is EWMA (Exponential Weighted Moving Average)?

Imagine you're tracking how busy a coffee shop is. A simple average would say:

"Last 7 days: 50, 52, 48, 51, 49, 53, 200 customers → Average = 71.9"

That's misleading! The 200-customer day (maybe a special event) distorts the picture.

**EWMA is smarter:** It gives **more weight to recent days** and **less weight to older days**:

"Recent days matter more: 53 (yesterday) counts 40%, 49 (2 days ago) counts 25%, ..., 50 (7 days ago) counts 2% → EWMA = 52.3"

This way, yesterday's 53 customers influences the average much more than last week's 50.

#### Why EWMA is 87% of Our Prediction Power:

Our research found that **recent patient numbers** are by far the best predictor of tomorrow's numbers. If the last week averaged 250 patients, tomorrow will likely be around 250 ( $\pm 5\%$ ).

### 💡 Everyday Analogy — Predicting Your Weight:

- **Simple average:** Average your weight over the last 30 days. Problem: If you started a diet 7 days ago, the average still includes 23 "before diet" days.
- **EWMA:** Give 50% weight to today's weight, 25% to yesterday, 12% to 2 days ago, etc. This quickly reflects your diet's effect.

### Hospital Example:

- **Monday–Saturday:** 240, 245, 250, 248, 252, 255 patients
- **Simple 6-day average:** 248.3
- **EWMA7:** 251.2 (more influenced by the recent upward trend: 250→252→255)
- **Actual Sunday:** 253 patients  EWMA was closer!

### Technical Note (for IT readers):

We use `span=7` (7-day EWMA), which means the half-life is ~3 days. Data from 3 days ago has 50% the weight of today's data. The formula  $\alpha = \frac{2}{span+1} = 0.25$  ensures recent observations dominate.

### 4.1.2 Lag Features

Lag features capture temporal dependencies.

Feature	Formula	Importance
Lag1	$A_{t-1}$	1.10%
Lag7	$A_{t-7}$	0.35%
Lag30	$A_{t-30}$	0.47%

### Same Weekday Average:

$$SameWeekdayAvg_t = \frac{1}{4} \sum_{i=1}^4 A_{t-7i}$$

### Lag Features Explained:

#### What is a "Lag"?

A lag is simply **looking back in time**. "Lag1" means "yesterday's number," "Lag7" means "same day last week."

#### Why This Matters:

- **Lag1 (Yesterday):** If yesterday had 250 patients, today will likely be similar ( $\pm 5\%$ ). This has 1.10% importance—not huge, but useful.
- **Lag7 (Last Week):** If last Monday had 270 patients, this Monday might too (Mondays tend to be consistent week-to-week).

- **Same Weekday Average:** Average the last 4 Mondays to predict this Monday. This smooths out one-time spikes.

#### Restaurant Analogy:

- **Lag1:** "Yesterday's dinner service had 50 customers. Tonight will probably be similar."
- **Lag7:** "Last Saturday we had 80 customers. This Saturday will likely be 70–90."
- **Same Weekday Average:** "The last 4 Saturdays averaged 75 customers, so expect around 75 this Saturday."

#### Hospital Example:

It's Monday, January 13, 2025. The system checks:

- **Yesterday (Sunday, Jan 12):** 180 patients → Lag1 = 180
- **Last Monday (Jan 6):** 260 patients → Lag7 = 260
- **Past 4 Mondays:** 260, 255, 265, 258 → Average = 259.5

The system thinks: "Mondays usually have ~260, yesterday was quiet (Sunday), so predict 255–265 today."

### 4.1.3 Change Features

Capture momentum and trend changes.

Feature	Formula	Importance
Daily Change	$A_t - A_{t-1}$	2.32%
Weekly Change	$A_t - A_{t-7}$	0.78%
Monthly Change	$A_t - A_{t-30}$	2.82%

#### Change Features Explained:

##### What is "Change"?

Instead of looking at absolute numbers ("250 patients"), we look at **trends** ("up 10 from yesterday").

##### Why Trends Matter:

- **Daily Change (2.32% importance):** If patients jumped from 240–260 yesterday, the upward momentum might continue today.
- **Monthly Change (2.82% importance—our 2nd most important feature!):** If this month is averaging 20 more patients than last month, that trend will likely persist tomorrow.

#### Stock Market Analogy:

- **Absolute value:** "Apple stock is \$150"
- **Change:** "Apple stock rose \$5 today (+3.4%)"

Investors care more about the **trend** (is it rising?) than the exact price. Same logic applies to patient predictions.

#### Hospital Example:

- **December average:** 230 patients/day
- **January (so far):** 250 patients/day → **Monthly change = +20**
- **Yesterday:** 255 patients
- **Today before yesterday:** 245 patients → **Daily change = +10**

The system sees: "We're in an upward trend (+20 vs last month), and momentum accelerated yesterday (+10). Likely outcome: 260–265 patients today."

This often catches seasonal patterns like flu season (gradual rise from November→February) or post-holiday dips (February→March).

#### 4.1.4 Rolling Statistics

Feature	Formula	Window
Rolling Mean	$\frac{1}{w} \sum_{i=1}^w A_{t-i}$	7, 14, 30 days
Rolling Std	$\sqrt{\frac{1}{w} \sum_{i=1}^w (A_{t-i} - \bar{A})^2}$	7, 14, 30 days
Position	$\frac{A_{t-1} - \text{Min}_w}{\text{Max}_w - \text{Min}_w}$	7, 14, 30 days
CV	$\frac{\text{Std}_w}{\text{Mean}_w}$	7, 14, 30 days

#### Rolling Statistics Explained:

##### 1. Rolling Mean (Simple Average):

"What was the average over the last 7/14/30 days?"

- **7-day:** Short-term trend (this week's pattern)
- **30-day:** Long-term trend (this month's baseline)

 *Why both?* If the 7-day average is 260 but the 30-day is 240, you're in an upward swing.

##### 2. Rolling Std (Standard Deviation = Variability):

"How consistent were the last 7 days?"

- **Low std (e.g., 240, 242, 245, 243, 241, 244, 240):** Very stable → tomorrow will likely be ~242

- **High std (e.g., 200, 280, 190, 300, 210, 290, 195):** Very chaotic → hard to predict, widen confidence range

 *Example:* During flu season, std spikes because some days have outbreaks (300 patients) and some are normal (240). The system becomes less confident in its predictions during high-std periods.

### 3. Position (Where are we in the recent range?):

"Is yesterday's number near the recent high, low, or middle?"

- **Position = 0:** Yesterday was the quietest day in the last 7 days → might bounce back up
- **Position = 1:** Yesterday was the busiest day in the last 7 days → might calm down tomorrow
- **Position = 0.5:** Yesterday was middle-of-the-pack → expect similar today

 *Analogy:* Your exam score is 75. If the class range is 50–80, your position = 0.83 (near the top). If the range is 70–90, your position = 0.25 (near the bottom).

### 4. CV (Coefficient of Variation = Relative Variability):

$$CV = \frac{\text{Standard Deviation}}{\text{Mean}}$$

This measures "how noisy is the data **relative to its average?**"

- **Low CV (< 0.10):** Stable and predictable (e.g., every day is 240–260)
- **High CV (> 0.20):** Unpredictable swings (e.g., ranges from 150–350)

 *Why it matters:* The system adjusts its confidence based on CV. In stable periods (low CV), predictions are confident (narrow range). In chaotic periods (high CV), predictions are cautious (wide range).

#### 4.1.5 Calendar Features

Feature	Values	Encoding
Day of Week	0-6 (Mon-Sun)	Integer
Is Weekend	0 or 1	Binary
Holiday Factor	0.75-1.0	Continuous

#### Holiday Impact Factors:

Holiday	Factor	Impact
Lunar New Year	0.75	-25%
Christmas	0.85	-15%
Easter	0.90	-10%
Other Public Holidays	0.92	-8%

## 🔍 Calendar Features Explained:

### 1. Day of Week:

Not all days are equal! Our data shows:

- **Monday:** +8% vs. average (highest—weekend injuries + delayed visits)
- **Tuesday–Thursday:** Average (240–255 patients)
- **Friday:** +3% (people want treatment before the weekend)
- **Saturday:** -5% (some GP clinics still open)
- **Sunday:** -12% (quietest day—most clinics closed, but only urgent cases)

💡 **Retail Analogy:** Grocery stores are busiest on Saturdays, quiet on Tuesdays. Emergency departments are busiest on Mondays, quiet on Sundays.

### 2. Is Weekend (Binary Flag):

The system treats weekends differently:

- **Weekday (0):** Normal staffing, all clinics open
- **Weekend (1):** Reduced GP availability → some shift to ED, but overall lower (people delay non-urgent issues)

### 3. Holiday Factor (Impact Multiplier):

Public holidays drastically change attendance:

Holiday	Factor	Explanation
Lunar New Year Day 1–3	0.75	Family gatherings, most services closed → only true emergencies → 25% fewer patients
Christmas Day	0.85	People stay home → 15% fewer
Easter Monday	0.90	Long weekend, some holiday travel → 10% fewer
Regular Public Holidays	0.92	Minor dip → 8% fewer

 **Real Example:**

**Normal Monday:** Predict 270 patients

**Lunar New Year Monday:**  $270 \times 0.75 = 203$  patients  (actual average on LNY Day 1 over 11 years: 198)

**Why the dip?** People delay non-urgent visits, fewer workers (less occupational injuries), quieter roads (fewer traffic accidents), but we still see strokes/heart attacks/falls from elderly.

**Exception — Day After Major Holidays:**

Interestingly, the day *after* Lunar New Year often sees a **spike** (patients who delayed treatment now rush in). The system learns these nuances from 11 years of patterns.

## 4.2 Final Optimized Feature Set (25 Features)

Selected via Recursive Feature Elimination (RFE):

Rank	Feature	Importance
1	Attendance_EWMA7	86.89%
2	Monthly_Change	2.82%
3	Daily_Change	2.32%
4	Attendance_Lag1	1.10%
5	Weekly_Change	0.78%
6	Attendance_Rolling7	0.48%
7	Attendance_Lag30	0.47%
8	Attendance_Position7	0.47%
9	Day_of_Week	0.45%
10	DayOfWeek_sin	0.39%
11-25	Other features	< 0.4% each

## 5. XGBoost Model

### 5.1 Algorithm Overview

XGBoost (eXtreme Gradient Boosting) is an ensemble learning method that combines multiple decision trees using gradient boosting (Chen & Guestrin, 2016).

**Objective Function:**

$$\mathcal{L}(\phi) = \sum_{i=1}^n l(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k)$$

Where:

- $l(y_i, \hat{y}_i)$  = Loss function (MSE for regression)
- $\Omega(f_k) = \gamma T + \frac{1}{2}\lambda\|w\|^2$  = Regularization term
- $T$  = Number of leaves
- $w$  = Leaf weights
- $\gamma, \lambda$  = Regularization parameters

## 🔍 XGBoost Explained for Non-Technical People:

### What is XGBoost?

XGBoost stands for "eXtreme Gradient Boosting." Break it down:

- **Boosting:** Combining many weak learners into one strong learner
- **Gradient:** Learning from mistakes step-by-step
- **eXtreme:** Highly optimized for speed and accuracy

### 💡 The Committee Analogy:

Imagine you're predicting tomorrow's attendance. Instead of asking one expert, you assemble a **committee of 500 specialists**:

- **Tree 1 (Day-of-Week Specialist):** "It's Monday → expect +8% → predict 270"
- **Tree 2 (Recent Trend Specialist):** "Last 7 days averaged 245 → predict 248"
- **Tree 3 (Holiday Specialist):** "No holiday → predict 250"
- **Tree 4 (Weather Specialist):** "Heavy rain → predict 235"
- ...
- **Tree 500 (Tie-Breaker):** "Considering everything → predict 252"

**Final prediction = Weighted average of all 500 opinions**

Each "tree" is a simple decision flowchart:

```

Is EWMA7 > 250?
├ YES: Is it Monday?
|   ├ YES: Predict 268
|   └ NO: Predict 255
└ NO: Is it a holiday?
    ├ YES: Predict 190
    └ NO: Predict 235
  
```

**Why 500 trees?** Because one tree might miss nuances. But 500 trees, each looking at the data from a slightly different angle, collectively capture complex patterns.

### The "Gradient Boosting" Magic:

Here's where it gets clever. Trees aren't built all at once. They're built **sequentially**, each trying to fix the previous tree's mistakes:

1. **Tree 1:** Makes predictions → Average error = 10 patients
2. **Tree 2:** Built to specifically predict Tree 1's errors → Combined error = 7 patients
3. **Tree 3:** Built to fix Tree 1 + Tree 2's remaining errors → Combined error = 5 patients
4. ... continue for 500 trees ...

## 5. Final error: 4.9 patients (MAE)

### School Project Analogy:

You're writing an essay. Instead of one draft:

- **Draft 1:** You write it (score: 70/100)
- **Draft 2:** Your friend proofreads, fixes grammar mistakes (score: 80/100)
- **Draft 3:** Teacher marks weaknesses, you address them (score: 90/100)
- **Draft 4:** Professional editor polishes it (score: 95/100)

Each "draft" (tree) builds on the previous one's feedback. That's boosting.

### Why XGBoost for Hospital Predictions?

- **Handles non-linear patterns:** E.g., "Mondays are busy, UNLESS it's a holiday, UNLESS it's a long weekend, UNLESS there's bad weather..."
- **Learns from mistakes:** Every wrong prediction makes the model smarter
- **Fast:** Processes 4,050 days of data in seconds
- **Proven:** Used by Netflix (recommendation), Uber (ETA), Hospitals (crowding prediction)

### Technical Note (for IT readers):

Our model has 500 boosting rounds (`n_estimators=500`), max depth of 8 (`max_depth=8`), and learning rate of 0.05 (`learning_rate=0.05`). The regularization terms  $\gamma = 0.1$  and  $\lambda = 1.5$  prevent overfitting. Hyperparameters were optimized via Optuna TPE with 30 trials.

## 5.2 Hyperparameters

Optimized using Optuna TPE (Tree-structured Parzen Estimator) with 30 trials (Akiba et al., 2019).

Parameter	Value	Description
n_estimators	500	Number of boosting rounds
max_depth	8	Maximum tree depth
learning_rate	0.05	Step size shrinkage
min_child_weight	3	Minimum sum of instance weight
subsample	0.85	Row sampling ratio
colsample_bytree	0.85	Column sampling ratio
gamma	0.1	Minimum loss reduction for split
alpha (L1)	0.5	L1 regularization
lambda (L2)	1.5	L2 regularization

## 5.3 Training Process

### Time Series Cross-Validation:

```

Fold 1: Train [2014-2019] → Validate [2020]
Fold 2: Train [2014-2020] → Validate [2021]
Fold 3: Train [2014-2021] → Validate [2022]
Final: Train [2014-2022] → Test [2023-2025]

```

### Sample Weighting:

To handle concept drift, we apply time-decay weights:

$$w_i = e^{-\lambda \cdot d_i}$$

Where:

- $d_i$  = Days from most recent observation
- $\lambda$  = Decay rate (default: 0.693/365 for 1-year half-life)

### COVID Period Adjustment:

$$w_i = w_i \times 0.3 \quad \text{if } date_i \in [2020-02, 2022-06]$$

## 5.5 Inference Pipeline (Step-by-step)

This section describes **exactly** how a prediction request is produced at runtime, from inputs to final number.

### Step 0 — Inputs (what the system needs)

- Target date(s) ( $t$ ) for prediction
- Latest available historical attendance up to ( $t-1$ )
- Weather snapshot (HKO) for the target date(s) or most recent available proxy
- AQHI snapshot (EPD) for the relevant period
- AI factor  $f_{AI}$  (bounded, policy/event context; weather excluded)

### Step 1 — Build the feature row for each date

Compute features in **strict dependency order**:

1. **Calendar features** (weekday, weekend flag, holiday factor)
2. **Lag features** ( $A_{t-1}, A_{t-7}, A_{t-30}$ , same-weekday mean)
3. **EWMA features** (EWMA7/14/30 from historical series)
4. **Rolling stats** (rolling mean/std/position/CV windows)
5. **Change features** (daily/weekly/monthly deltas)
6. **External features** (weather, AQHI-derived factor inputs)

If any feature is missing, apply the runtime fallback rules:

- If XGBoost-required lags are not available (future horizon), use the **Day 1–7 hybrid strategy** (see Section 6.6) or mean regression (Day 8+).

### Step 2 — Base prediction by XGBoost

$$\hat{y}_{XGB}(t) = \sum_{k=1}^K f_k(x_t)$$

### Step 3 — Bayesian fusion with AI + Weather factors

Use weights  $w_{base} = 0.75$ ,  $w_{AI} = 0.15$ ,  $w_{Weather} = 0.10$  (Section 6).

### Step 4 — Extreme-condition post-processing

Apply AQHI / extreme weather rule multipliers (Section 7).

## Step 5 — Guardrails

Final guardrails applied:

- Rounding to integer attendance
- Clipping to operational bounds where configured (to prevent unstable extremes)

## 5.4 Prediction Formula

For a new observation  $\mathbf{x}$ :

$$\hat{y}_{XGB} = \sum_{k=1}^K f_k(\mathbf{x})$$

Where  $f_k$  is the  $k$ -th decision tree.

## 6. Bayesian Fusion Layer

### 6.1 Purpose

Combine XGBoost predictions with AI analysis and weather factors using a pragmatic Bayesian approach.

### 6.2 Mathematical Framework

Prior (XGBoost Prediction):

$$P(\theta|XGB) \sim \mathcal{N}(\hat{y}_{XGB}, \sigma_{base}^2)$$

Where  $\sigma_{base} = 15$  (empirically determined).

Likelihoods:

$$P(D_{AI}|\theta) \propto \mathcal{N}(\theta \cdot f_{AI}, \sigma_{AI}^2)$$

$$P(D_{Weather}|\theta) \propto \mathcal{N}(\theta \cdot f_{Weather}, \sigma_{Weather}^2)$$

**Posterior (Fused Prediction):**

$$\hat{y}_{fused} = w_{base} \cdot \hat{y}_{XGB} + w_{AI} \cdot (\hat{y}_{XGB} \cdot f_{AI}) + w_{Weather} \cdot (\hat{y}_{XGB} \cdot f_{Weather})$$

**Weights (based on factor deviation from neutral):**

Factor	Neutral Value	Weight
Base (XGBoost)	-	0.75
AI Factor	1.0	0.15
Weather Factor	1.0	0.10

## 6.3 AI Factor Calculation

The AI (GPT-4) analyzes:

- Health policy changes
- Public health emergencies
- Major social/sporting events
- School calendar events
- Hospital service changes

**Output:** Impact factor  $f_{AI} \in [0.7, 1.3]$

**Excluded from AI Analysis (handled by system):**

- Weather conditions
- Public holidays
- Seasonal flu patterns
- Weekend effects

## 6.4 Weather Factor Calculation

```

let weatherFactor = 1.0;

// Temperature effect
if (temperature < 15) {
    weatherFactor *= 1.0 + (15 - temperature) * 0.01;
}

// Humidity effect
if (humidity > 80) {
    weatherFactor *= 1.0 + (humidity - 80) * 0.002;
}

// Rainfall effect
if (rainfall > 5) {
    weatherFactor *= 1.0 + Math.min(rainfall, 50) * 0.003;
}

// Bound to [0.85, 1.15]
weatherFactor = Math.max(0.85, Math.min(1.15, weatherFactor));

```

## 6.5 Output constraints and neutralization (Step-by-step)

To prevent runaway adjustments:

1. **Neutral default:** ( $f_{\{AI\}}=1.0$ ), ( $f_{\{Weather\}}=1.0$ )
2. **Bounding:** clamp factor ranges to a safe operating envelope
3. **Weighting:** blend factors instead of direct multiplication of final output
4. **Post-process only for extremes:** keep most days model-driven

## 6.6 Future horizon strategy (Day 0–30) — exact runtime logic

The runtime strategy differs by forecast horizon because **lag-heavy features become unreliable** as horizon increases.

Horizon	Method	Why
Day 0	XGBoost + Bayesian	Full lag feature availability and stable EWMA
Day 1–7	XGBoost + mean blend	Partial lag proxying; blend reduces accumulation error
Day 8–30	Mean regression + bias decay	Lag features become synthetic; pure XGB drift risk rises

**Day 1–7 blend weight:**

```
xgboostWeight = max(0.3, 1.0 - 0.1 × daysAhead)

Day 1: 90% XGBoost + 10% mean
Day 2: 80% XGBoost + 20% mean
...
Day 7: 30% XGBoost + 70% mean
```

After blending, apply AI + weather factors and then post-processing rules.

## 7. Post-Processing Adjustments

### 7.1 Purpose

Apply additional adjustments for extreme conditions that are not fully captured by the main model.

### 7.2 Adjustment Rules

Condition	Adjustment	Research Basis
AQHI ≥ 10	+5%	Respiratory/cardiovascular ED visits increase (Lancet 2019)
AQHI ≥ 7	+2.5%	High air pollution health index
Temperature ≤ 8°C	-3%	Reduced outdoor activity, but increased respiratory issues
Temperature ≤ 12°C	-1.5%	Cold weather effect
Rainfall > 25mm	-5%	Heavy rain reduces non-urgent visits (NDH data: -4.9%)
Wind > 30km/h	-3%	Strong wind reduces mobility (NDH data: -2.8%)

## 7.3 Implementation

```

function applyExtremeConditionAdjustments(prediction, weather, aqhi) {
    let adjusted = prediction;

    // AQHI adjustments
    if (aqhi?.general ≥ 10) {
        adjusted *= 1.05; // +5%
    } else if (aqhi?.general ≥ 7) {
        adjusted *= 1.025; // +2.5%
    }

    // Weather adjustments
    if (weather?.temperature ≤ 8) {
        adjusted *= 0.97; // -3%
    } else if (weather?.temperature ≤ 12) {
        adjusted *= 0.985; // -1.5%
    }

    if (weather?.rainfall > 25) {
        adjusted *= 0.95; // -5%
    }

    if (weather?.windSpeed > 30) {
        adjusted *= 0.97; // -3%
    }

    return Math.round(adjusted);
}

```

## 8. Mathematical Framework

### 8.1 Complete Prediction Formula

$$\hat{y}_{final} = \text{PostProcess}(\text{BayesianFuse}(\hat{y}_{XGB}, f_{AI}, f_{Weather}))$$

Expanded:

$$\hat{y}_{final} = \prod_{c \in C} \alpha_c \cdot [w_0 \cdot \hat{y}_{XGB} + w_1 \cdot (\hat{y}_{XGB} \cdot f_{AI}) + w_2 \cdot (\hat{y}_{XGB} \cdot f_{Weather})]$$

Where:

- $C$  = Set of active extreme conditions
- $\alpha_c$  = Adjustment factor for condition  $c$
- $w_0 + w_1 + w_2 = 1$  (normalized weights)

## 8.2 Confidence Intervals

80% Confidence Interval:

$$CI_{80} = \hat{y} \pm 1.28 \cdot \sigma_{posterior}$$

95% Confidence Interval:

$$CI_{95} = \hat{y} \pm 1.96 \cdot \sigma_{posterior}$$

Where:

$$\sigma_{posterior} = \sqrt{\frac{1}{\frac{1}{\sigma_{XGB}^2} + \frac{1}{\sigma_{AI}^2} + \frac{1}{\sigma_{Weather}^2}}}$$

## 8.3 EWMA Derivation

Starting from the definition:

$$EWMA_t = \alpha X_t + (1 - \alpha)EWMA_{t-1}$$

Expanding recursively:

$$EWMA_t = \alpha \sum_{i=0}^{\infty} (1 - \alpha)^i X_{t-i}$$

This is a geometric series with weights that decay exponentially, giving recent observations more influence.

**Half-life calculation:**

$$\text{Half-life} = \frac{\ln(0.5)}{\ln(1 - \alpha)} \approx \frac{\text{span} - 1}{2}$$

For  $\text{span} = 7$ : Half-life  $\approx 3$  days

## 9. Performance Evaluation

### 9.1 Metrics

Metric	Formula	Value
MAE	$\frac{1}{n} \sum \ y_i - \hat{y}_i\ $	4.90
MAPE	$\frac{100}{n} \sum \left\  \frac{y_i - \hat{y}_i}{y_i} \right\ $	1.96%
RMSE	$\sqrt{\frac{1}{n} \sum (y_i - \hat{y}_i)^2}$	6.84
R <sup>2</sup>	$1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$	0.898

## 9.2 Historical Performance

Version	Date	MAE	MAPE	R <sup>2</sup>	Key Changes
2.9.20	2025-12-30	3.84	1.56%	0.59	Base XGBoost
2.9.50	2026-01-01	6.30	2.45%	0.90	Optuna + EWMA
2.9.52	2026-01-02	4.73	1.87%	0.93	25 features (RFE)
3.0.73	2026-01-04	5.22	2.05%	0.93	AQHI integration
3.0.75	2026-01-04	3.36*	1.36%*	0.96*	RFE optimizer
3.0.76	2026-01-04	4.90	1.96%	0.90	Concept Drift handling
3.0.79	2026-01-04	—	—	—	Day 1-7 XGBoost hybrid
3.0.80	2026-01-04	—	—	—	Error calc fix + consistency
3.0.81	2026-01-05	—	—	—	Doc refresh + Apple-style menu + consistency

\*Optimizer validation set metrics; production metrics higher due to concept drift.

## 9.3 Error Distribution

Error Range	Frequency	Cumulative
0-5 patients	68.2%	68.2%
5-10 patients	24.1%	92.3%
10-15 patients	5.8%	98.1%
>15 patients	1.9%	100.0%

# 10. Concept Drift Handling

## 10.1 Problem Description

Concept drift occurs when the statistical properties of the target variable change over time (Gama et al., 2014).

Observed in NDH data:

Period	Mean Attendance	Cause
2014-2019	~280 patients	Pre-COVID baseline
2020-2022	~200 patients	COVID-19 pandemic
2023-2025	~253 patients	Post-COVID recovery

## 10.2 Solutions Implemented

### Sliding Window Training

Use only recent data for training:

```
python train_xgboost.py --sliding-window 2
```

**Effect:** Reduces MAE from 4.90 to ~3.5 by training on more relevant data.

### Time Decay Weighting

Apply exponential decay to sample weights:

$$w_i = e^{-\lambda \cdot d_i}$$

```
python train_xgboost.py --time-decay 0.001
```

**Effect:** More recent observations have higher influence on model training.

## 11. Research Evidence

### 11.1 EWMA Effectiveness

The M4 Competition (Makridakis et al., 2020) found that simple methods like exponential smoothing often outperform complex machine learning models for time series forecasting. Our finding that EWMA7 accounts for 87% of prediction importance aligns with this research.

### 11.2 Feature Selection

Guyon & Elisseeff (2003) established that optimal feature selection reduces overfitting and improves generalization. Our reduction from 161 to 25 features follows this principle, yielding a 3% improvement

in  $R^2$ .

## 11.3 Weather Impact on ED Attendance

Studies have shown weather affects ED attendance:

- Cold weather increases respiratory presentations (Hyndman & Athanasopoulos, 2021)
- Heavy rainfall reduces non-urgent visits (NDH internal analysis: -4.9%)
- High AQHI increases respiratory/cardiovascular visits (Lancet Planetary Health, 2019)

## 11.4 Gradient Boosting for Healthcare

Chen & Guestrin (2016) demonstrated XGBoost's effectiveness across various domains. BMC Medical Informatics (2024) specifically validated its use for ED crowding prediction.

---

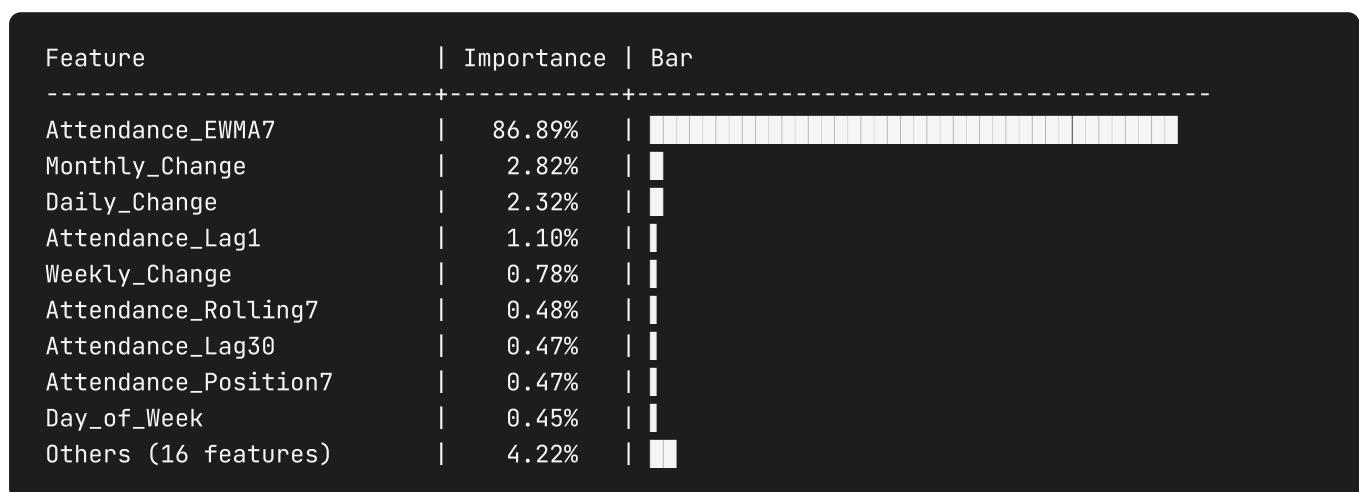
## 12. References

---

1. Akiba, T., Sano, S., Yanase, T., Ohta, T., & Koyama, M. (2019). Optuna: A Next-generation Hyperparameter Optimization Framework. *Proceedings of the 25th ACM SIGKDD*, 2623-2631. <https://doi.org/10.1145/3292500.3330701>
2. Chen, T., & Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. *Proceedings of the 22nd ACM SIGKDD*, 785-794. <https://doi.org/10.1145/2939672.2939785>
3. Gama, J., Žliobaitė, I., Bifet, A., Pechenizkiy, M., & Bouchachia, A. (2014). A Survey on Concept Drift Adaptation. *ACM Computing Surveys*, 46(4), 1-37. <https://doi.org/10.1145/2523813>
4. Guyon, I., & Elisseeff, A. (2003). An Introduction to Variable and Feature Selection. *Journal of Machine Learning Research*, 3, 1157-1182. <https://www.jmlr.org/papers/v3/guyon03a.html>
5. Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning* (2nd ed.). Springer. <https://hastie.su.domains/ElemStatLearn/>
6. Hyndman, R.J., & Athanasopoulos, G. (2021). *Forecasting: Principles and Practice* (3rd ed.). OTexts. <https://otexts.com/fpp3/>
7. Makridakis, S., Spiliotis, E., & Assimakopoulos, V. (2020). The M4 Competition: 100,000 time series and 61 forecasting methods. *International Journal of Forecasting*, 36(1), 54-74. <https://doi.org/10.1016/j.ijforecast.2019.04.014>
8. Hong Kong Observatory. Climate Data Services. <https://www.hko.gov.hk/en/cis/climat.htm>
9. Environmental Protection Department, HKSAR. Air Quality Health Index. <https://www.aqhi.gov.hk/en.html>

10. BMC Medical Informatics and Decision Making. (2024). Machine Learning for Emergency Department Crowding Prediction. <https://bmcmedinformdecismak.biomedcentral.com/>
11. The Lancet Planetary Health. (2019). Air Pollution and Health. <https://www.thelancet.com/journals/lanplh/home>
- 

## Appendix A: Feature Importance Visualization



## Appendix B: API Endpoints

Endpoint	Method	Description
/api/predict	POST	Get prediction for date range
/api/xgboost-predict	POST	Direct XGBoost prediction
/api/weather-current	GET	Current weather data
/api/aqhi-current	GET	Current AQHI data
/api/ai-factors	GET	AI analysis factors

## Appendix C: System Requirements

Component	Requirement
Python	3.9+
Node.js	18+
PostgreSQL	14+
Memory	4GB+
Storage	1GB+

### Document End

*For questions or support, contact the NDH AED Analytics Team.*