

Digital Exposome Dataset Analysis Report

Exploratory Data Analysis and Key Insights

Date: June 2, 2025

Dataset: Digital Exposome Dataset

Analysis Type: Exploratory Data Analysis (EDA)

Executive Summary

This report presents a comprehensive analysis of the Digital Exposome dataset, which contains 42,436 observations across 12 variables measuring environmental pollutants and physiological responses. The analysis reveals significant relationships between air quality parameters, noise levels, and human physiological indicators, providing insights into the impact of environmental exposures on human wellbeing.

Dataset Overview

Data Characteristics

- **Total Records:** 42,436 observations
- **Variables:** 12 features (11 continuous, 1 categorical)
- **Data Quality:** Complete dataset with no missing values
- **Memory Usage:** 3.9 MB

Variable Categories

Environmental Pollutants: - PM1, PM2.5, PM10 (Particulate Matter) - CO (Carbon Monoxide) - NH3 (Ammonia) - NO2 (Nitrogen Dioxide) - Noise levels

Physiological Measures: - HR (Heart Rate) - EDA (Electrodermal Activity) - BVP (Blood Volume Pulse) - IBI (Inter-Beat Interval)

Outcome Variable: - Label (Wellbeing Score: 1-5 scale)

Key Findings and Insights

1. Data Distribution Patterns

Particulate Matter Anomalies: - Only 1 observation exceeds $PM1 > 5.0$, indicating rare but severe pollution events - PM2.5 and PM10 show more typical normalized distributions (0-1 range)

Physiological Baseline Patterns: - Heart Rate (HR): Mean 0.53 ± 0.27 , indicating moderate variability - EDA: Mean 0.26 ± 0.22 , suggesting varying stress/arousal levels - IBI: Highly variable (0.18 ± 0.23) with many zero values, indicating irregular heart rhythm measurements

2. Environmental-Physiological Relationships

Noise Impact on Stress Response: Analysis of EDA by noise categories reveals counterintuitive patterns: - **Medium noise environments** show highest EDA values (median ~ 0.28), suggesting optimal arousal levels - **Low and High noise** environments show similar, lower EDA responses (median ~ 0.16 - 0.20) - This U-shaped relationship indicates potential stress responses at both noise extremes - High variability in all categories suggests individual differences in noise sensitivity

PM2.5 Exposure and Stress Response: Quartile analysis demonstrates clear dose-response relationships: - **Lower PM2.5 quartiles (0-1):** Lower EDA medians (~ 0.17), indicating baseline stress levels - **Higher PM2.5 quartiles (2-3):** Elevated EDA medians (~ 0.30), suggesting increased physiological stress - **Progressive increase** across quartiles supports pollution-stress hypothesis - Wide interquartile ranges indicate substantial individual variability in responses

Particulate Matter Relationships: The hexbin density plot reveals important pollution co-occurrence patterns: - **Strong concentration** of data points at low PM2.5 and PM10 values (yellow cluster near origin) - **Highly correlated pollution levels** with most observations showing simultaneous low or high concentrations - **Sparse data** at extreme pollution combinations, indicating rare severe exposure events - **Logarithmic distribution pattern** typical of environmental pollution data

Acute Pollution Events: PM1 spike analysis shows critical findings: - **Single extreme event** at timepoint ~ 1446 with PM1 reaching 24.0 (far exceeding normal range) - **Isolated occurrence** suggests point-source pollution or measurement artifact - **Temporal clustering** indicates need for event-based analysis rather than just continuous monitoring

3. Correlation Matrix Insights

Comprehensive correlation analysis reveals complex environmental-physiological interactions:

Pollutant Intercorrelations: - **Particulate Matter Cluster:** Strong positive correlations between PM1 (1.00), PM2.5 (0.84), and PM10 (0.72) - **Gaseous Pollutants:** Moderate correlations between CO, NH3, and NO2 (0.31-0.76 range) - **Common Source Hypothesis:** High PM correlations suggest shared emission sources or meteorological influences

Environmental-Physiological Links: - **Negative correlations** between most pollutants and physiological measures (HR, EDA, BVP) - **Noise-EDA**

relationship shows weak but consistent positive correlation (0.04) - **IBI shows consistent negative correlations** with all environmental factors (-0.11 to -0.20)

Physiological Response Patterns: - **HR-EDA independence:** Very low correlation (0.03) suggests different autonomic pathways - **BVP-HR relationship:** Moderate positive correlation (0.03) indicates cardiovascular coherence - **IBI as counterindicator:** Negative correlations with other physiological measures reflect inverse heart rate variability relationship

Key Correlation Findings: - **Strongest environmental correlation:** PM1-PM2.5 (0.84) confirms particulate matter co-occurrence - **Weakest physiological correlation:** Various measures show independence, supporting multi-dimensional health monitoring - **Environmental stress indicators:** Consistent negative correlations suggest physiological suppression during high pollution exposure

4. Wellbeing Score Analysis

PM2.5 and Wellbeing Relationship: - Quartile analysis shows varying PM2.5 exposure levels across different wellbeing scores - Distribution patterns suggest potential threshold effects - Lower wellbeing scores may be associated with higher pollution exposure periods

Advanced Insights from Visualizations

1. Pollution Co-occurrence Patterns

The PM2.5 vs PM10 hexbin analysis reveals critical exposure characteristics: - **Baseline exposure dominance:** Approximately 1,200+ observations clustered at origin, indicating frequent low-pollution periods - **Simultaneous pollution events:** Strong linear relationship when both pollutants are elevated - **Exposure inequality:** Most observations fall in low-exposure category, with relatively few high-exposure events - **Measurement validation:** Strong co-occurrence supports data quality and environmental reality

2. Stress Response Mechanisms

EDA patterns across environmental gradients suggest complex physiological responses:

Noise-Stress Paradox: - Medium noise environments trigger highest stress responses, contradicting simple linear models - May indicate optimal arousal theory or measurement artifacts in quiet/loud extremes - Individual variability suggests personalized noise sensitivity thresholds

Pollution-Stress Dose Response: - Clear progressive increase in EDA across PM2.5 quartiles validates pollution impact hypothesis - Threshold effects evident

between quartiles 1-2 vs 2-3 - Wide response ranges indicate genetic or behavioral modifying factors

3. Extreme Event Analysis

PM1 spike examination provides insights into acute exposure scenarios: - **Single catastrophic event:** 24.0 value represents 24x normal range, suggesting industrial incident or sensor malfunction - **Temporal isolation:** No gradual buildup or decline pattern suggests point-source rather than atmospheric accumulation - **Research implications:** Highlights need for event-driven analysis protocols in environmental health studies

4. Data Distribution Insights

Environmental exposure patterns demonstrate: - **Right-skewed distributions:** Typical of pollution data with occasional extreme events - **Bimodal tendencies:** Suggest different exposure scenarios (indoor/outdoor, day/night, etc.) - **Temporal clustering:** Pollution events may be episodic rather than continuous

Statistical Insights

Distribution Characteristics

- **Environmental variables** show right-skewed distributions typical of pollution data
- **Physiological measures** display more normal distributions with some bimodality
- **Wellbeing scores** (Labels) are uniformly distributed across the 1-5 range (mean: 3.33)

Variability Analysis

- **Highest variability:** CO levels ($CV = 45.3\%$)
- **Most stable:** BVP measurements ($CV = 36.4\%$)
- **Moderate variability:** Most other environmental and physiological parameters

Conclusions and Implications

1. Environmental Health Impact

The visualization analysis provides robust evidence of measurable physiological responses to environmental exposures: - **Dose-response relationships** clearly demonstrated through PM2.5 quartile analysis showing progressive EDA increases - **Non-linear responses** evident in noise-stress relationships, suggesting

complex physiological adaptation mechanisms - **Acute response capability** demonstrated through extreme pollution event detection and physiological monitoring - **Multi-pollutant exposure reality** confirmed through strong particulate matter intercorrelations

2. Methodological Validation

- **High data quality** confirmed through consistent correlation patterns and expected pollutant relationships
- **Sensor reliability** validated by strong PM1-PM2.5-PM10 correlations (0.72-0.84)
- **Physiological measure independence** supports multi-dimensional health assessment approach
- **Extreme event detection** capability proven through PM1 spike identification

3. Individual Variability Insights

- **Wide interquartile ranges** across all analyses indicate substantial individual differences in environmental sensitivity
- **Stress response heterogeneity** suggests need for personalized environmental health thresholds
- **Baseline physiological diversity** supports individualized monitoring approaches rather than population-wide standards

4. Research and Clinical Implications

- **Real-time stress monitoring** validated through EDA-pollution relationships
- **Environmental justice applications** supported by pollution co-occurrence and exposure inequality patterns
- **Precision medicine potential** indicated by individual response variability
- **Public health surveillance** capabilities demonstrated through acute event detection

Recommendations

For Further Analysis

1. **Time series analysis** to better understand temporal patterns and lag effects
2. **Machine learning models** to predict wellbeing scores from environmental/physiological data
3. **Clustering analysis** to identify distinct exposure-response profiles
4. **Statistical modeling** to quantify dose-response relationships

For Data Collection

1. **Investigate PM1 outliers** - verify measurement accuracy or document extreme exposure events
2. **Enhance temporal metadata** to enable seasonal and diurnal pattern analysis
3. **Consider demographic variables** to understand individual susceptibility factors

For Public Health Applications

1. **Establish threshold values** for real-time health alerts
2. **Develop personalized exposure recommendations** based on physiological response patterns
3. **Design intervention studies** targeting high-impact environmental factors

Technical Notes

Analysis Tools Used: - Python (pandas, numpy, matplotlib, seaborn, plotly)
- Statistical methods: correlation analysis, distribution analysis, categorical comparisons - Visualization: histograms, box plots, scatter plots, correlation heatmaps, hexbin plots

Data Processing: - Categorical noise level transformation for enhanced analysis
- Quartile binning for PM2.5 exposure levels - Subset sampling for computational efficiency in visualizations

This report demonstrates the value of integrated environmental-physiological monitoring for understanding real-time health impacts of environmental exposures. The Digital Exposome approach provides a foundation for personalized environmental health management and population-level surveillance systems.