**ampliasaúde**
observatório do período perinatal

# Exploration Trails:
# Data visualization methodology

The **Exploration Trails** are an interactive and customized tool that offers exploratory visualization, download, and sharing of maternal and neonatal health data using pollution data for 5,565 Brazilian municipalities over seven years (2012-2019).

**SUMMARY**

REALIZATION

UFRJ

LAB VIS

FIOCRUZ

FINANCING

BILL & MELINDA GATES foundation

CNPq

FAPERJ

SUPPORT

PCDaS

# THE OBJECTIVE OF THE AMPLIA SAÚDE (EXTENDED HEALTH) PROJECT

Maternal, neonatal, and children under-five mortality has been reduced in Brazil over the last few years, mainly by means of public health policies specific to the mother-infant axis within the Sustainable Development Goal (SDG) 3. A healthy pregnancy, birth care, sufficient breastfeeding, and timely vaccination favor children's growth and development. The fetus developing environment is essential for the child's onward life. It is mainly affected by social determinants and the health system (Figure 1), that is, demographic factors and the maternal socioeconomic conditions that make them more or less vulnerable, besides the access to specialized care, prenatal care, and vaccination. Environmental conditions consider only those in the immediate surroundings, such as access to drinking water, sewage treatment, and garbage collection. However, extreme weather conditions, environmental disasters, and pollution can directly affect human health as well as the availability and access to health services. Climate and environmental determinants are increasingly frequent, and their impact must be measured in maternal-neonatal health. Based on the National Health System (SUS) principle of comprehensiveness, the Amplia Saúde (Extended Health) project aimed to expand the analysis model of maternal and neonatal morbimortality indicators by integrating air pollution-related indicators.
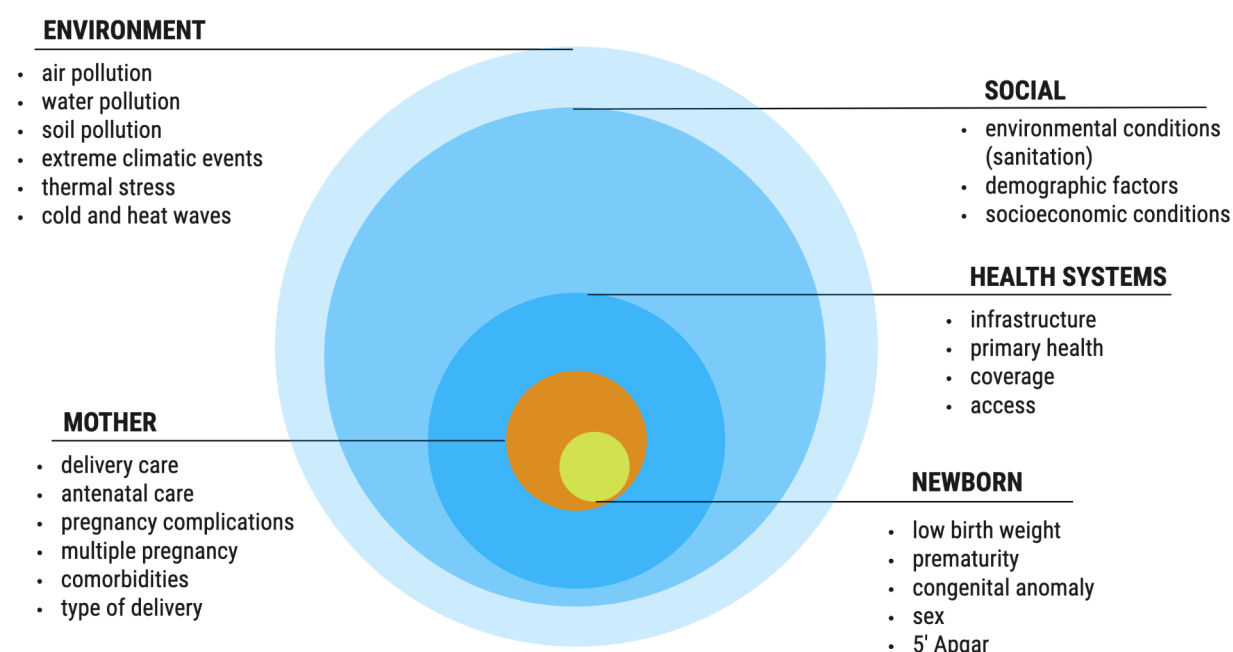


**ENVIRONMENT**
- air pollution
- water pollution
- soil pollution
- extreme climatic events
- thermal stress
- cold and heat waves

**SOCIAL**
- environmental conditions (sanitation)
- demographic factors
- socioeconomic conditions

**HEALTH SYSTEMS**
- infrastructure
- primary health
- coverage
- access

**MOTHER**
- delivery care
- antenatal care
- pregnancy complications
- multiple pregnancy
- comorbidities
- type of delivery

**NEWBORN**
- low birth weight
- prematurity
- congenital anomaly
- sex
- 5' Apgar

*Figure 1. The conceptual framework for analyzing neonatal morbimortality*

# THE TOOL'S SPECIFIC OBJECTIVE

→ Offer an **innovative visual analysis approach** of maternal and infant health data from the perinatal period, considering environmental factors (air pollution)

→ Create an **interactive tool that allows exploring** this relationship

# DATA SOURCES

The visualization tool uses data made openly available by the following government agencies:

→ DATASUS, Department of Informatics of the Unified Health System of Brazil - Ministry of Health

→ The "Queimadas" (Fires) Program of the National Institute for Space Research (INPE)

→ The Brazilian Institute of Geography and Statistics (IBGE)

Infant and fetal death surveillance has been mandatory in public and private health services of the Brazilian Unified Health System (SUS) since 2010 by decree of the Ministry of Health. Thus, we specifically adopted data from DATASUS Mortality Information System (SIM) and Live Births Information System (SINASC) from 2012 onwards. Besides mortality data, we also employed the Fetal Death Declaration (DOFET-SIM) databases, which were extracted, processed, and enriched by the PCDaS (Platform of Data Science applied to Health, the Institute of Scientific and Technological Communication and Information in Health from the Oswaldo Cruz Foundation-ICICT), using ETL (Extract, Transform, Load) methods.

We used INPE's Environmental Information System Integrated to Health (SISAM), developed through a partnership between INPE, PAHO/WHO (the Pan American Health Organization/World Health Organization), and FUNDEP (Foundation for Research Development). We used data on 2.5 micron air pollutants, (PM2.5) fine particulate matter, from 2012 to 2019, the last year available in the database. While other data, such as ozone concentration, carbon monoxide, and wind speed and direction, were available, they lacked the necessary consistency and were thus abandoned after the first tests. SISAM's data were downloaded by an automatic data collector (web scraping crawler). The database offers up to four daily PM2.5 measurements, although variations were identified in different municipalities. Thus, we opted for a weekly aggregation of pollution data

We employed the IBGE's 2010 Human Development Index (HDI), the last census carried out so far (2022), and the geolocation of the municipalities.

Only aggregated secondary data were used throughout the *Amplia Saúde* project, with no possibility of individual identification, including no common field that could link the events in the SINASC and SIM databases.

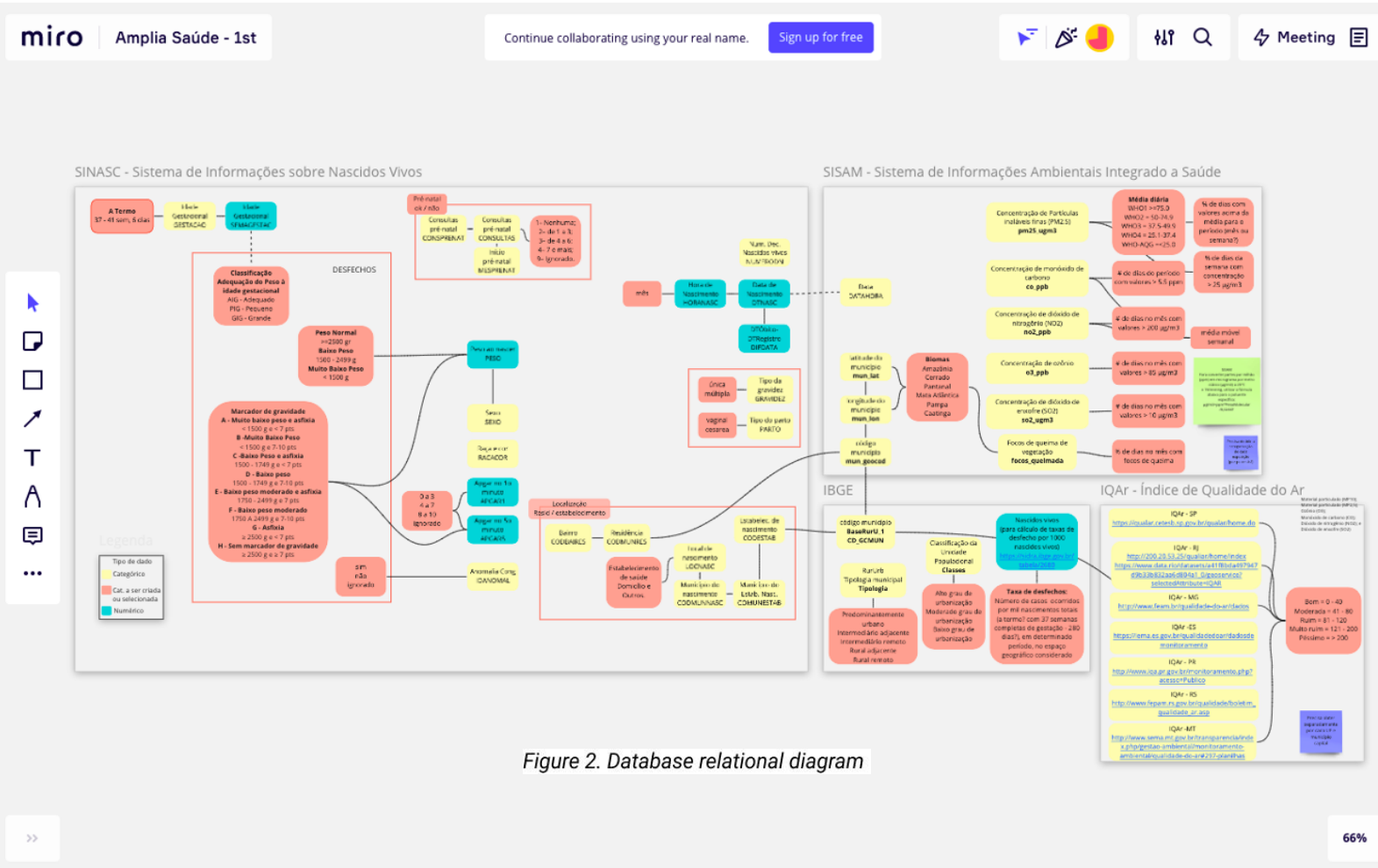# RELEVANT INFORMATION ON DATA RELATIONSHIP

The data tool developed within the *Amplia Saúde* project comprises a set of exploratory visualizations that combine a vast perinatal data volume with socioeconomic, geographic, and air pollution incidence attributes. The exploratory visualizations can identify data patterns, offer insights, and formulate or test hypotheses.

Since they gather and combine heterogeneous and large-scale data (Big Data) retrieved from unlinked databases, these visualizations are not indicated to seek direct cause-and-effect relationships between air pollution and health outcomes for newborns, requiring further studies that substantiate the identified correlations.

# THE PROJECT'S ASSUMPTIONS

As it addresses large volumes of data collected over several years in different formats from diverse locations, we had to conduct a preliminary survey of the available databases and verify their quality, besides possible aggregation and interaction, creation of indicators, or use of existing indicators. In this initial stage, a diagram was created on the Miro platform (Figure 2), where data from sources related to maternal-neonatal health and environmental pollution were listed. Subsequently, some of the selected databases were not employed, as was the case with the National Registry of Health Facilities (CNES) and the Hospital and Outpatient Information System (SIH).

The Tool provides exploratory, on-demand, and detailed information. It is a layer aimed at those seeking to perform specific queries, formulate or test hypotheses, and make data-driven decisions. Therefore, we had to create a customized architecture for the tool (Figure 3). We adopted an architecture where the data is stored in and retrieved from a cloud server to generate interactive visualizations. However, the visualizations themselves are built by code running in the user's browser. An additional layer was introduced in the system architecture (labeled "middleware" in Figure 3). It allows compressing the data sent over the network from the server to the client and faster access to data to work around performance issues typical of large data visualization systems.
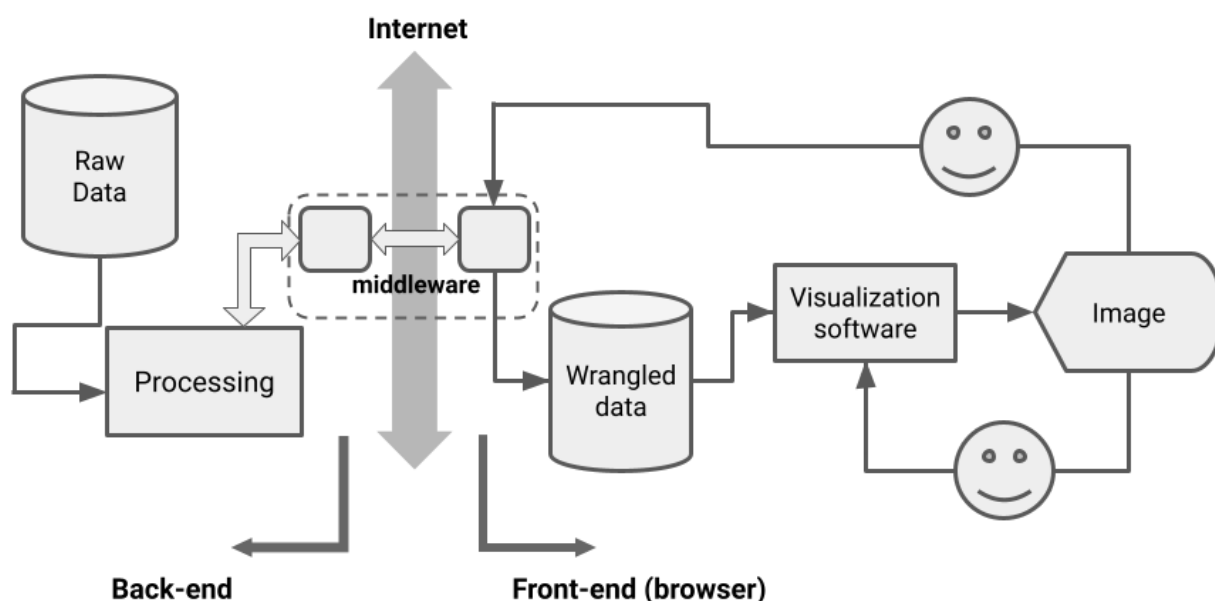


Figure 2. Database relational diagram

*Figure 3. Data/tool system architecture*

We had to make several design-related decisions in producing visualizations due to the large data volume that can be displayed in multiple ways according to choices made by the user. For instance, we limited to four the number of municipalities that can be plotted simultaneously in the Trails visualization. Similarly, when comparing municipalities with extreme pollution values, the air pollution graphs of municipalities with low pollution levels may show no variation, because the axis scale was fixed to cover the  wide range of values.

Five design assumptions were determined due to the several challenges observed throughout data structuring and while creating the visualizations, which will be detailed below.

### DATA AGGREGATION AND TEMPORAL RELATIONSHIP APPROXIMATION

Maternal-neonatal health data and pollution data have different temporal granularities. While the former are usually organized and analyzed by month or year, the latter are measured several times daily. We had to create a single aggregation mode to allow the simultaneous visual analysis of the two datasets (health and air pollution). Thus, all data were aggregated by epidemiological week and year, considering their use in the visualization groups of the exploratory trails and the map/scatter plot.

### SELECTION OF SINGLETON BIRTHS

We decided to exclude those born from multiple pregnancies from the project database to allow a basis for comparing pregnancy periods and because fetal morbimortality increases in pregnancies with two or more fetuses in-uterus.

### POLLUTION AND HUMAN HEALTH RATES

We adopted the benchmark values determined by CONAMA (National Council for the Environment) on air quality in daily measurements of particulate matter 25 μ (PM2.5) to visually indicate the pollution levels unsuitable for human health, which in turn were based on the WHO recommendations of daily (15μg/m3) and annual (5μg/m3) thresholds without prejudice to health (*World Health Organization. Air Quality Guidelines - Update 2021. Copenhagen, Denmark: WHO Regional Office for Europe*).

### REORGANIZING DATA BY CONCEPTION RATE

The main challenge to visualizing and analyzing the possible effects of pollution on pregnancies and maternal-neonatal health lies in the secondary data-imposed limitation. In other words, we cannot monitor pregnancies and births individually. Women are pregnant at different gestational periods at any

one time, and pregnancies are of varying lengths. Thus, the same air pollution peak period can occur in women in early pregnancy while others are in the middle or at the end of their pregnancy. Furthermore, the period of most susceptibility for the fetus is known to occur between the 22nd and 38th gestation weeks. However, we cannot observe the effect over different gestation periods when analyzing the exposure to pollution over a set of pregnancies. From these considerations, we decided to adopt an estimated conception date, which was calculated by subtracting the length of gestation in weeks from the date of birth. In the visualization, we used the term "window of susceptibility" to demarcate this period between the 22nd and 38th weeks.

## HIGHLIGHTED DATA

The SIM and SINASC databases have more than one hundred attributes when considering the original attributes and those added later at enrichment. However, birth weight is a critical variable since this information is crucial for calculating several indicators, including gestational age birthweight adequacy. After establishing the relevance of the weight attribute, only cases with no missing birthweight were included. Gestational age in weeks, necessary for the calculations, was included in the SINASC database in 2009 but found to be consistent only from 2012 onwards. This date determined the onset of the period covered by the visualizations.

# DEFINITION AND INDICATORS

The tool comprises two sets of visualizations: **Map** and **Trails**.

## MAP

The Map visualization offers several alternatives for comparing municipalities, filtering, and allows the selection of up to four municipalities, considering the in-depth analysis to be conducted later with the Trails visualization. On the map itself, we can observe the size of the population of each municipality and the mean annual value of a primary variable chosen among six possible ones, facilitating the selection based on geographic

positioning. This page also contains a scatter plot to analyze up to three variables from all available municipalities per the desired filtering. A bar chart showing the variation of the primary variable over the years for up to four municipalities is also available.

KEY INDICATORS:
→ Pollution
→ Mean annual pollution
→ Low birthweight
→ Neonatal Mortality
→ Perinatal Mortality
→ Infant Mortality

## TRAILS

In the Trails, we can investigate the municipalities selected on the Map and optionally modify this choice. The Trails comprise three investigative sets: Live Births, Birth Characteristics, and Perinatal and Infant Mortality. All have a line graph with the data of air pollutants in parallel with other graphs related explicitly to the neonatal health data addressed. Neonatal health line graphs will vary by attributes selected in the main menu or the comparison menu.

## LIVE BIRTHS

KEY INDICATORS:
→ Pollutant
→ Birth weight
→ (Birthweight) Adequacy to gestational age
→ Total live births

ATTRIBUTES:
→ **Birth weight**: Normal, Low weight, Extremely low weight, Unknown
→ **Gestational age**: Post-term, Late term, Full term, Early term, Moderate preterm, Very preterm, Extremely preterm, Unknown
→ **Gestational age adequacy**: LGA, AGA, SGA, Unknown
→ **Sex**: Female, Male, Unspecified
→ **Congenital anomaly**: No, Yes
→ **Mother's ethnicity/skin color**: White, Black, Brown, Other, Unknown
→ **Mother's Schooling**: Higher Education, High School, Elementary School, Incomplete Elementary School, Illiterate, Unknown
→ **Mother's Marital Status**: Single, Married or Common-Law Marriage, Widowed, Separated or Divorced, Unknown

- → **Mother's age**: 10-14, 15-19, 20-29, 30-39, 40-49, 50+, Unknown
- → **Prenatal visits**: 7+, 4-6, 1-3, 0, Unknown
- → **Delivery type**: Vaginal, Cesarean, Unspecified
- → **Induced labor**: Yes, No
- → **Prenatal control Adequacy**: Not performed, Inadequate, Intermediate, Adequate, More than adequate, Unknown

## BIRTHS CHARACTERISTICS

KEY INDICATORS:
- → Pollutant
- → Birth weight
- → Robson groups
- → Total live births

ATTRIBUTES:
- → **Birth weight**: Normal, Low weight, Extremely low weight, Unknown
- → **Gestational age**: Post-term, Late term, Full term, Early term, Moderate preterm, Very preterm, Extremely preterm, Unknown
- → **Adequacy (of birthweight) to Gestational age:** LGA (large for gestational age), AGA (adequate for gestational age), SGA (small for gestational age), Unknown
- → **Sex**: Female, Male, Unspecified
- → **Congenital anomaly**: No, Yes
- → **Mother's ethnicity/skin color**: White, Black, Brown, Other, Unknown
- → **Mother's Schooling**: Higher Education, High School, Elementary School, Incomplete Elementary School, Illiterate, Unknown
- → **Mother's Marital Status**: Single, Married or Common-Law Marriage, Widowed, Separated or Divorced, Unknown
- → **Mother's age**: 10-14, 15-19, 20-29, 30-39, 40-49, 50+, Unknown
- → **Apgar 5th Minute:** Good vitality, Moderate asphyxia, Severe asphyxia, Unknow
- → **Delivery type**: Vaginal, Cesarean, Unspecified
- → **Induced labor**: Yes, No
- → **Prenatal Adequacy**: Not performed, Inadequate, Intermediate, Adequate, More than adequate, Unknown

## PERINATAL AND INFANT MORTALITY

KEY INDICATORS:
- → Pollutant
- → Fetal deaths
- → Deaths versus delivery
- → Total deaths

ATTRIBUTES:
- → **Birth weight**: Normal, Low weight, Extremely low weight, Unknown
- → **Sex**: Female, Male, Unspecified
- → **Mother's ethnicity/skin color**: White, Black, Brown, Other, Unknown
- → **Mother's Schooling**: Higher Education, High School, Elementary School, Incomplete Elementary School, Illiterate, Unknown
- → **Mother's age**: 10-14, 15-19, 20-29, 30-39, 40-49, 50+, Unknow
- → **Delivery type**: Vaginal, Cesarean, Unspecified