

## **Motivation Statement**

Estimating smoke impact on respiratory illness-caused mortality in Fresno, CA.

The linkages between increased particulate pollution due to wildfires and increased mortality in California are well-documented (Chen et al., Gwon et. al). This analysis will use a predictive estimate for the impact of smoke caused by wildfires within 650 miles of Fresno. Data on mortality caused by respiratory illness from the California Department of Health and Human Services will then be added into the predictive estimate to quantify future impacts of smoke from a public health perspective. This project will present a refinement of a smoke estimate predictive model by using non-periodic time series forecasting techniques such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) (Tensorflow). We hope to generate correlations between the refined predictive estimate for smoke impact and changes in Fresno's mortality rates caused by respiratory conditions. In turn, calling for the future study of causal connections between the two variables.

## **Data To Be Used**

*Data:* Death Profiles by Zip Code (1989-2022)

*Source:* Hosted on the California Health and Human Services Open Data Portal, provided by the California Department of Public Health

*Summary:*

“This dataset contains counts of deaths for California residents by ZIP Code based on information entered on death certificates. Final counts are derived from static data and include out-of-state deaths of California residents. The data tables include deaths of residents of California by ZIP Code of residence (by residence). The data are reported as totals, as well as stratified by age and gender. Deaths due to all causes (ALL) and selected underlying cause of death categories are provided. See temporal coverage for more information on which combinations are available for which years.”

*Dataset:*

[CSV Files](#) of:

- 2019-2022 Final Deaths by Year by ZIP Code
  - “This data table reports the annual number of deaths of California residents by the ZIP Code of residence regardless of where the death occurred (by residence). The cause of death categories are based solely on the underlying cause of death as coded by the International Classification of Diseases.” (CalHHS)

- 2009-2018 Final Deaths by Year by ZIP CodeCSV Popular
- 1999-2008 Final Deaths by Year by ZIP CodeCSV Popular
- 1989-1998 Final Deaths by Year by ZIP CodeCSV Popular
- Data Dictionary - Deaths by ZIP CodeCSV Popular

*License/Terms of Use:* Open for public use contingent on following the [CalHHS terms of use](#)

---

*Data:* ICD-9 and ICD-10 Insurance Codes for Respiratory Illnesses

*Source:* Wikipedia

*Summary:* Unified and standardized codes attributed to “diseases, signs and symptoms, abnormal findings, complaints, social circumstances, and external causes of injury or diseases”, used for medical or insurance purposes (Wikipedia).

*Dataset:* [ICD-9](#) [460-519] and [ICD-10](#) [J00-J99] codes

*License/Terms of Use:* Open for public use under the [Creative Commons Attribution-ShareAlike 4.0 International License](#), additional terms may apply.

### **Model To Be Used**

*Model:* The initial predictive model was a Holt-Winters Exponential Smoothing Time-Series Forecasting Model. This model is good for periodic and noisy data, but has limitations in predicting non-periodic outcomes. Therefore, this extension will seek to refine this predictive model by switching to predicting smoke estimates with Convolutional Neural Networks (CNNs) provided by Tensorflow. Time will also be dedicated to explore the success of Recurrent Neural Networks (RNNs) on this question.

The foundational theory and practice of these networks comes from a Jupyter Notebook named [Time Series Forecasting](#), by an unnamed team of Tensorflow authors. This model is distributed under the Apache License, Version 2.0 (Tensorflow).

*Smoke Estimate Model Adaptations:* This extension will adapt the initial smoke estimate by creating a smoke IMPACT estimate on mortality. This will be a new variable measuring not smoke, but smoke impact on human-centered mortality rates from respiratory illnesses in Fresno, CA.

The correlation between increases in smoke estimate and increases in mortality will inherently carry some lag time. Therefore, a large part of updating my model will be finding the lag time, in

years, which maximizes the correlation between my initial smoke estimate and Fresno's mortality rate due to respiratory diseases.

An important subfield of mortality studies in demography is the study of excess mortality. This process involves proving causality: proving that some political regime, pandemic, or other large-scale change in a population caused an excess of deaths, compared to if that large-scale change hadn't occurred. This requires computing the counterfactual mortality rate, the rate of deaths in a population if that large-scale change hadn't occurred. The reason excess mortality will not be studied in this extension is that it is beyond the scope of this assignment; this is a human-centered data science assignment, not a demography paper. Ultimately, this extension will aim to provide evidence suggesting that research establishing such causality would be scientifically impactful in the California Central Valley.

### **Unknowns and Dependencies**

Due to the techniques used and author's understanding of demographic methods, this analysis will be inherently restricted to building a refined smoke impact score informed by correlations with mortality rates. This integration of mortality rates will also seek to identify the approximate lag in years between high-wildfire smoke years and spikes in mortality. Furthermore, there exist inherent limitations to this extension's ability to optimally tune the CNN, due to the timeline of this project. All forecasted smoke impact estimates will therefore be documented as such.

The data contains all relevant zip codes for Fresno, CA, but not all years from the USGS Wildfire dataset. This extension will be dependent on a narrower timeframe, from 1989 - 2019.

Furthermore, this analysis will be informed by hospital-reported mortality rates, as this is the primary source for the California Department of Public Health. This extension has a tradeoff between the availability of data capturing all of Fresno's deaths, informing their death rates, and the precision of the smoke estimate's impact on mortality; it will be assumed that not all of Fresno's deaths during this time period have been fully recorded.

### Timeline to Completion

Sunday	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday
<b>27</b>	<b>28</b>	<b>29</b>	<b>30</b> Collect Data	<b>31</b>	<b>1</b>	<b>2</b> Build Initial Model
<b>3</b>	<b>4</b> <i>Class</i>	<b>5</b>	<b>6</b> Improve Mortality Rate Integration	<b>7</b>	<b>8</b>	<b>9</b>
<b>10</b>	<b>11</b> <i>Class</i> Finalize Mortality Rate-Informed Smoke Impact Score	<b>12</b>	<b>13</b>	<b>14</b>	<b>15</b> First Iteration Model Tuning	<b>16</b> Document Process
<b>17</b>	<b>18</b> <i>Class</i>	<b>19</b> Second Iteration Model Tuning	<b>20</b> Document Process	<b>21</b>	<b>22</b> Third Iteration Model Tuning	<b>23</b> Document Process
<b>24</b> Final Model Tuning Done	<b>25</b> <i>Class</i>	<b>26</b> Visualizations Complete	<b>27</b> Project Part 3 - Presentation	<b>28</b>	<b>29</b>	<b>30</b> Document Process
Sunday	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday
<b>1</b> Finalize Documentation	<b>2</b> <i>Class</i> Project Part 5 - Presentation Feedback	<b>3</b>	<b>4</b> Project Part 4 - Final Repo	<b>5</b>	<b>6</b>	<b>7</b>

## **Bibliography**

California Department of Public Health. California Comprehensive Master Death File (Static), 2014-2022. Compiled by Center for Health Statistics and Informatics.

California Department of Public Health. Death Statistical Master File (Static), 1989-2013. Compiled by Center for Health Statistics and Informatics.

Chen, Gongbo, et al. "Mortality risk attributable to wildfire-related PM<sub>2.5</sub> pollution: A global time series study in 749 locations." *The Lancet Planetary Health*, vol. 5, no. 9, Sept. 2021, [https://doi.org/10.1016/s2542-5196\(21\)00200-x](https://doi.org/10.1016/s2542-5196(21)00200-x).

Gwon, Y.; Ji, Y.; Bell, J.E.; Abadi, A.M.; Berman, J.D.; Rau, A.; Leeper, R.D.; Rennie, J. The Association between Drought Exposure and Respiratory-Related Mortality in the United States from 2000 to 2018. *Int. J. Environ. Res. Public Health* 2023, 20, 6076. <https://doi.org/10.3390/ijerph20126076>

The Tensorflow Authors, (2024) Time Series Forecasting. Tensorflow. [https://www.tensorflow.org/tutorials/structured\\_data/time\\_series](https://www.tensorflow.org/tutorials/structured_data/time_series)

United States Mortality DataBase. University of California, Berkeley (USA). Available at [www.usa.mortality.org](http://www.usa.mortality.org) (data downloaded on 31 Oct. 2024).