# Using a Deep Autoencoder Model to Characterize Air Pollution

Nicolas Shannon[1], Bishnu Timalsena[2], Raul Chavez[3],
Vandana Nunna Lakshmi[2]

[1]Department of Computer Science, Loyola University Chicago, Chicago, IL, US
[2]Department of Computer Science, University North Texas, Denton, TX, US
[3]Department of Electrical and Computer Engineering, University of Texas - Rio Grande Valley,
Rio Grande Valley, TX, US

**Abstract**

Climate change is one of the most complex challenges humanity has ever faced. It is a multi-dimensional problem that poses economic, social, scientific, political, and moral questions that must be addressed. One of the most important ways of combating climate change is to research pollution mitigation strategies so that policy makers and government planners are better equipped to make informed decisions. We focused on the principle cause of global warming: the polluting gases responsible for the rapid global climate shift. We developed a technique to characterize cities by their levels of various polluting gases and particulate matter using an undercomplete autoencoder model. By creating a representation of pollution levels for numerous cities across the United States, we explored some of the more prominent features that characterize urban pollution. We compared two techniques for reducing the dimensionality of our input - Principle Component Analysis (PCA) and Undercomplete Autoencoders. The two first two hidden dimensions represented a large percentage of the explained variance in the model, so we clustered the cities in an effort to extract any useful features.

## 1   Introduction

Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis.

Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

## 2 Methodology

The main use of an undercomplete autoencoder is not in its ability to perfectly reconstruct the input, but in the useful features it is capable of distilling. In our case, we wanted to see what properties the autoencoder model could extract when given time series air pollution data for various cities throughout the United States. Using data from OpenWeather's Air Pollution API, we constructed data frames for eight polluting gases and particulates. We obtained the daily averages of each pollutant six and a half months for just over eighteen-thousand cities in the U.S. and other U.S. territories. After obtaining the requisite data, we cleaned and preprocessed the data by removing cities that had the same name and any entries with missing values. While this made the data quickly useable, a more ideal solution would be to add a state/country code to distinguish between cities and to impute any missing values. For more on this see the future improvements section. We compared the autoencoder model to Principle Component Analysis, a much older technique used for dimensionality reduction that has been a part of statistical literature since the early twentieth century.[2] The main distinction between autoencoders and PCA is that autoencoders can span a nonlinear subspace, whereas PCA learns the principle subspace. Autoencoders that employ nonlinar encoder and decoder functions can create a more powerful generalization than PCA. However, an autoencoder model that is allowed too much capacity is at risk of reconstructing the input "too perfectly" without extracting any useful information.[3]

As Goodfellow, Bengio, and Courville (2016) say, there's no reason... Jolliffe (2002) Agency (2009)

## References

Agency, U.S. Environmental Protection (2009). "Carbon Monoxide NAAQS: Scope and Methods Plan for Health Risk and Exposure Assessment". In:

https://www.epa.gov/sites/production/files/2020-07/documents/
2009_04_coscopeandmethodsplan.pdf.

Goodfellow, Ian, Yoshua Bengio, and Aaron Courville (2016). *Deep Learning*.
http://www.deeplearningbook.org. MIT Press.

Jolliffe, I.T. (2002). *Principle Component Analysis – 2nd Ed.* Library of Congress.