

# Lawrence Lin

925-918-2478 | [lawrencedlin@gmail.com](mailto:lawrencedlin@gmail.com) | [linkedin.com/in/llinnn](https://www.linkedin.com/in/llinnn) | [github.com/lawrencedlin](https://github.com/lawrencedlin)

## EDUCATION

---

### University of San Francisco

*M.S. Data Science*

San Francisco, CA

*July 2021 - August 2022 (Expected)*

### University of California, Santa Barbara

*B.S. Statistics*

Santa Barbara, CA

*August 2017 - June 2021*

**Courses:** Advanced Machine Learning, Bayesian Statistics, Distributed Computing, Time Series, Stochastic Processes

## EXPERIENCE

---

### Data Science Intern

*Walmart Labs*

November 2021 - Present

*Sunnyvale, CA*

- Analyzed customers' seasonal purchase behavior for festivals using Apache Spark and Seaborn
- Worked on feature engineering and data cleaning using Apache Hadoop
- Developed a Next-item Sequential Recommendation Neural Network Model in TensorFlow
- Deployed model on Google Dataproc
- Created validation metrics to assess seasonal inference of model

### Research Assistant

*Sansum Diabetes Research Institute*

January 2021 - June 2021

*Santa Barbara, CA*

- Geographically visualized diabetes severity by zip code using Folium
- Modeled HbA1c with regression models achieving an  $R^2$  of 0.77
- Authored and presented weekly written reports of insights to SDRI researchers

### Gretler Fellow

*University of California, Santa Barbara*

September 2019 - June 2020

*Santa Barbara, CA*

- Web-scraped thousands of pages of state legislative data using BeautifulSoup
- Cleaned, processed, and validated data using Pandas

## PROJECTS

---

### Predict Implicit Ratings (Kaggle) | *Pytorch, FastAi, NumPy, Pandas*

- Developed Matrix Factorization model in PyTorch
- Implemented and trained Tabular Neural Network model in FastAi using cyclical learning rate
- Created various negative sampling algorithms for models
- Achieved 1<sup>st</sup> place on Kaggle leaderboard with loss of 0.40320

### Feature Importance | *Scikit-Learn, NumPy, Pandas*

- Implemented Spearman correlation, PCA, permutation importance, and drop-column importance from scratch
- Visualized the cross-validated Gradient-Boosted Tree's  $R^2$  trained on k most important features
- Implemented automatic forward feature selection algorithm using permutation importance
- Calculated variance and empirical p-value of feature importances using bootstrap samples

## TECHNICAL SKILLS

---

**Languages:** Python, R, C++, SQL (Postgres), NoSQL (Mongo) HTML/CSS, Bash

**Frameworks:** Hadoop Ecosystem (HDFS, YARN, Spark, SparkMLlib, HiveQL) TensorFlow, PyTorch, FastAI, Scikit-Learn, Statsmodels, Scipy, Numpy, Pandas, Matplotlib, Seaborn, Flask, BeautifulSoup, Selenium

**Developer Tools:** Git, Docker, Google Cloud Platform, Amazon Web Services, MongoDB