

# Using GANs to Generate Synthetic Data

Joseph Despres and Yunus Shariff  
Michigan State University

## Introduction and Motivation

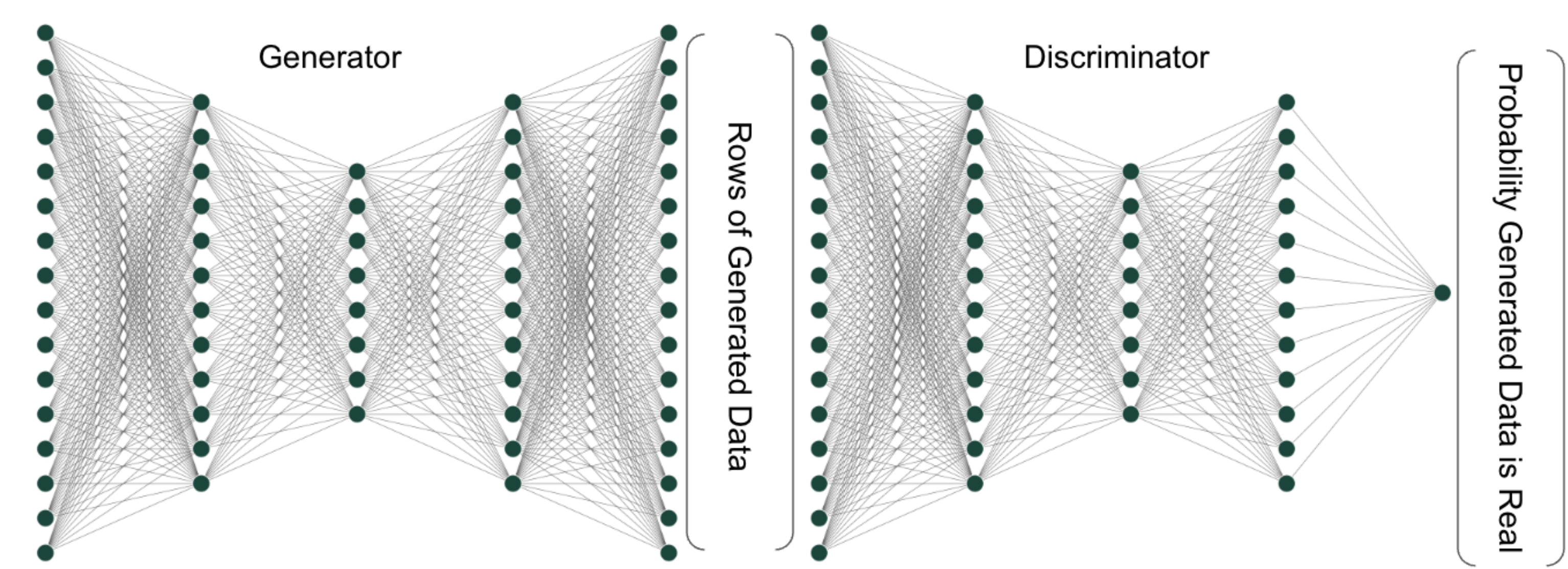
The aim of this project is to extend the Generative Adversarial Network beyond image generation into generating synthetic data. The objective of any GAN is to train a Neural Network to recognize a fake, while training an adjacent Neural Network to spot the fake, the model has converged to a desirable solution where fake data are indistinguishable from the real [1]. This area has shown substantial progress in the past several years today's GANs are known to generate astonishingly accurate generating images. Applications for this technology have yet to be fully realized. This project explores one of many potential use cases.

## Potential Use Cases

- Although data are more and more available in seeming abundance, there is still a need to generate more. First, big datasets are often convenient to collect and not controlled.
- **Use to generate additional datapoints in Controlled experiments**
- Since These data will be generated, confidentially is less of a concern. Therefore, these data could be used by peole training models outside of a particular organization.
- **Use to generate data to fit models where data are confidential**

## Method

The challenges with training a GAN to generate synthetic data is a row in a dataframe is several orders of magnitude smaller than a full image. An artificial nerual network has no difficulty fitting to it. The challenge is getting the right fit. This is a ballance between fitting the data so well it mimics the data seen and underfitting where it cannot generate. Therefor where the process could be generative yet,



## Experiment

### Generate and Label Data

| Minimal Example                                   |            |            |            |            |         |  |
|---|------------|------------|------------|------------|---------|--|
| Generate fake data, append to real, and train GAN |            |            |            |            |         |  |
| rdelay  | age        | rasp3      | td         | expdd      | is_real |  |
| -0.1526527  | -1.2121831 | -0.4472136 | -0.4923660 | 1.4494064  | 1       |  |
| -0.6383658  | 0.2020305  | 0.6708204  | -0.0820610 | -0.0607479 | 1       |  |
| 1.6514244   | -0.8081220 | -1.5652476 | -0.3282440 | -0.1689550 | 1       |  |
| -0.9159161  | 1.2121831  | 0.6708204  | -0.8206099 | -1.3592337 | 1       |  |
| 0.0555101   | 0.6060915  | 0.6708204  | 1.7232809  | 0.1395303  | 1       |  |
| 0.2316489   | -0.2192497 | 0.5248911  | 0.0797766  | -0.1968897 | 0       |  |
| -0.0481826  | 1.5192747  | -0.0045090 | 0.9744888  | -2.8427019 | 0       |  |
| -0.3158222  | 1.5258328  | -0.9236113 | 0.0106245  | -0.1436559 | 0       |  |
| -0.2907037  | -0.1617303 | 0.9790730  | -1.5406935 | 0.7550172  | 0       |  |
| -1.1800350  | -0.1189155 | 1.3039990  | 1.1063480  | 2.0710810  | 0       |  |

| Minimal Example                    |            |            |            |            |                     |  |
|------------------------------------|------------|------------|------------|------------|---------------------|--|
| Estimated probability data is real |            |            |            |            |                     |  |
| rdelay                             | age        | rasp3      | td         | expdd      | probability_is_real |  |
| -0.1526527                         | -1.2121831 | -0.4472136 | -0.4923660 | 1.4494064  | 0.9698882           |  |
| -0.6383658                         | 0.2020305  | 0.6708204  | -0.0820610 | -0.0607479 | 0.9697784           |  |
| 1.6514244                          | -0.8081220 | -1.5652476 | -0.3282440 | -0.1689550 | 0.9674411           |  |
| -0.9159161                         | 1.2121831  | 0.6708204  | -0.8206099 | -1.3592337 | 0.9669423           |  |
| 0.0555101                          | 0.6060915  | 0.6708204  | 1.7232809  | 0.1395303  | 0.9698892           |  |
| -0.4375128                         | 0.5214567  | -0.1713268 | 1.0914937  | -0.3610044 | 0.3533136           |  |
| 0.1673733                          | -0.6796552 | -0.9850925 | -0.955354  | -0.3855441 | 0.3856075           |  |
| 1.0546583                          | 0.4201501  | 0.6421450  | 1.1982469  | -0.1962640 | 0.3217127           |  |
| -0.0803883                         | -1.9345402 | 0.2203373  | -0.3808446 | 0.7414579  | 0.3097427           |  |
| -0.2065684                         | -0.7166563 | 0.8721342  | 2.2525326  | 0.7977623  | 0.3449469           |  |

| Minimal Example                    |            |            |            |            |                     |  |
|------------------------------------|------------|------------|------------|------------|---------------------|--|
| Estimated probability data is real |            |            |            |            |                     |  |
| rdelay                             | age        | rasp3      | td         | expdd      | probability_is_real |  |
| -0.1526527                         | -1.2121831 | -0.4472136 | -0.4923660 | 1.4494064  | 0.9751343           |  |
| -0.6383658                         | 0.2020305  | 0.6708204  | -0.0820610 | -0.0607479 | 0.9634157           |  |
| 1.6514244                          | -0.8081220 | -1.5652476 | -0.3282440 | -0.1689550 | 0.9576347           |  |
| -0.9159161                         | 1.2121831  | 0.6708204  | -0.8206099 | -1.3592337 | 0.9692259           |  |
| 0.0555101                          | 0.6060915  | 0.6708204  | 1.7232809  | 0.1395303  | 0.9813991           |  |
| -1.0309191                         | 1.1386368  | 1.0842576  | -0.7508808 | 0.3402721  | 0.5540787           |  |
| -0.3814221                         | 0.8675030  | -0.2946905 | -0.1045899 | -1.4483958 | 0.7900096           |  |
| -1.7107864                         | 1.4801052  | -1.3323751 | -0.2920365 | -2.9327913 | 0.6587892           |  |
| -0.2353828                         | 0.5057158  | 0.5928036  | 0.1671478  | 0.9925682  | 0.7344067           |  |
| -1.4030078                         | -0.1905782 | 0.0841914  | -0.6393145 | 1.3587482  | 0.7860073           |  |

### Generate Data and Retrain

| Minimal Example                    |            |            |            |            |                     |  |
|------------------------------------|------------|------------|------------|------------|---------------------|--|
| Estimated probability data is real |            |            |            |            |                     |  |
| rdelay                             | age        | rasp3      | td         | expdd      | probability_is_real |  |
| -0.1526527                         | -1.2121831 | -0.4472136 | -0.4923660 | 1.4494064  | 0.9679135           |  |
| -0.6383658                         | 0.2020305  | 0.6708204  | -0.0820610 | -0.0607479 | 0.9685082           |  |
| 1.6514244                          | -0.8081220 | -1.5652476 | -0.3282440 | -0.1689550 | 0.9603543           |  |
| -0.9159161                         | 1.2121831  | 0.6708204  | -0.8206099 | -1.3592337 | 0.9690194           |  |
| 0.0555101                          | 0.6060915  | 0.6708204  | 1.7232809  | 0.1395303  | 0.9998083           |  |
| -0.2772837                         | 0.9719187  | -0.4466996 | -2.3410539 | 0.7101667  | 0.0204801           |  |
| 1.0260006                          | 1.2748769  | 0.9795976  | -1.0115168 | 0.4987356  | 0.0017582           |  |
| 3.5675369                          | -0.4570953 | 0.8062306  | -3.0644811 | 0.3857205  | 0.0007435           |  |
| 0.3801538                          | 0.5359162  | -2.8023581 | -0.5935911 | 0.9851182  | 0.0223810           |  |
| -0.5490505                         | 0.2701278  | 0.4339166  | -0.1718862 | -0.6262546 | 0.0491006           |  |

| Minimal Example                    |            |            |            |            |                     |  |
|------------------------------------|------------|------------|------------|------------|---------------------|--|
| Estimated probability data is real |            |            |            |            |                     |  |
| rdelay                             | age        | rasp3      | td         | expdd      | probability_is_real |  |
| -0.1526527                         | -1.2121831 | -0.4472136 | -0.4923660 | 1.4494064  | 0.9637378           |  |
| -0.6383658                         | 0.2020305  | 0.6708204  | -0.0820610 | -0.0607479 | 0.9623097           |  |
| 1.6514244                          | -0.8081220 | -1.5652476 | -0.3282440 | -0.1689550 | 0.9564374           |  |
| -0.9159161                         | 1.2121831  | 0.6708204  | -0.8206099 | -1.3592337 | 0.9594558           |  |
| 0.0555101                          | 0.6060915  | 0.6708204  | 1.7232809  | 0.1395303  | 0.9841065           |  |
| 0.6132995                          | 0.2966554  | 0.1725356  | 0.0663309  | 0.7656293  | 0.2198298           |  |
| 1.1766507                          | 0.2963944  | 1.1329624  | -0.6495373 | -0.0093986 | 0.1423325           |  |
| -1.4541459                         | 0.8453177  | -1.8864023 | -0.9861465 | -0.0857276 | 0.1255824           |  |
| 0.8717983                          | 1.0047278  | 0.2884087  | 0.4371844  | 0.7561790  | 0.1219416           |  |
| -1.2912063                         | -0.3895090 | -0.7358911 | 1.0391212  | -2.7888273 | 0.1337178           |  |

| Minimal Example  |            |            |            |            |                     |  |
|--|------------|------------|------------|------------|---------------------|--|
| Train GAN until generated data is indistinguishable from real data |            |            |            |            |                     |  |
| rdelay   | age        | rasp3      | td         | expdd      | probability_is_real |  |
| -0.1526527   | -1.2121831 | -0.4472136 | -0.4923660 | 1.4494064  | 0.9957382           |  |
| -0.6383658   | 0.2020305  | 0.6708204  | -0.0820610 | -0.0607479 | 0.9979365           |  |
| 1.6514244  | -0.8081220 | -1.5652476 | -0.3282440 | -0.1689550 | 0.9969307           |  |
| -0.9159161   | 1.2121831  | 0.6708204  | -0.8206099 | -1.3592337 | 0.9901561           |  |
| 0.0555101  | 0.6060915  | 0.6708204  | 1.7232809  | 0.1395303  | 0.9921287           |  |
| -1.0858977   | 0.1258421  | 0.7208802  | -0.6415371 | 1.2507096  | 0.9987849           |  |
| 1.3405701  | -1.3444258 | -0.2103912 | 1.3688782  | -0.5387795 | 0.9971853           |  |
| 0.4695888  | 1.2875312  | 1.2613931  | 2.2261421  | -0.9526827 | 0.9960572           |  |
| -0.4728902   | 0.5750984  | 1.2177754  | 0.1072202  | 0.0881042  | 0.9911092           |  |
| -0.0267091   | -0.2755770 | -0.9962036 | -0.1694681 | 0.0481810  | 0.9979257           |  |

This project we use New York City Taxi from 2016. Primarily because of well collected. data from the N [2]

## Results

## References

- [1] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.
- [2] New York City Taxi and Limousine Commission.