

Creating fair machine learning models with Generative Adversarial Networks

Hamaad Musharaf Shah

License

Copyright 2020 Hamaad Musharaf Shah

Licensed under the Apache License, Version 2.0 (the “License”); you may not use this file except in compliance with the License. You may obtain a copy of the License at

<http://www.apache.org/licenses/LICENSE-2.0>

Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an “AS IS” BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.

Creating fair machine learning models with Generative Adversarial Networks

Author: Hamaad Shah

Disclaimer

This work is not necessarily a representation of any past or current employer of mine.

Introduction

- Ensuring fairness in machine learning is a key aspect especially in the financial services domain.
 - Think of a mortgage underwriting machine learning pipeline that perhaps might be sensitive to attributes such as gender and race without explicitly being trained on such features.
 - This can potentially cause a hidden cost of capital charge from an operational risk perspective if customers are not treated fairly.
 - Clearly a machine learning pipeline that makes unfair decisions given the same feature vector across multiple groups such as gender and race is unethical.
- To overcome this potential issue we discuss a method published in the 31st conference on Neural Information Processing Systems, i.e., NeurIPS 2017.
 - The title of that paper: “Learning to Pivot with Adversarial Networks”.
- I first came across the aforementioned NeurIPS paper in this blogpost.
 - <https://blog.godatadriven.com/fairness-in-ml>
- My lecture will discuss the aforementioned paper and Generative Adversarial Networks (GAN) in detail. It will also provide a R version of the code associated in the aforementioned blogpost.

Income Prediction

- We use the “Census Income” dataset to predict whether a person’s income is greater than \$50,000 per year. The link to the dataset is [here](#).

- <https://archive.ics.uci.edu/ml/datasets/Adult>
- The features used are denoted as below.
 - $x \sim P_X(x)$
- The sensitive attribute(s) are denoted as below.
 - $z \sim P_Z(z)$
- The target variable is denoted as below.
 - $y \sim P_Y(y)$

Classifier and Adversarial Model

- The classifier is a deep learner denoted as below.
 - $C(x|\theta_C) : x \rightarrow y$
- Note that we are not explicitly training the model on the sensitive attributes. However it is entirely possible that there might be some hidden bias within our features that perhaps act as proxies for the sensitive attributes.
- The sensitive attributes for this exercise are race and gender. These are binary variables.
- For each of these sensitive attributes we want to ensure the following holds.
 - $P(C(x|\theta_C)|z=1) = P(C(x|\theta_C)|z=0) \quad \forall x \sim P_X(x)$
- We introduce the concept of an adversarial deep learner model which takes in the predictions of the classifier in order to predict the sensitive attributes as below.
 - $R(C(x|\theta_C)|\theta_R) : C(x|\theta_C) \rightarrow z$
- In an ideal scenario the classifier will perform at its optimal level while simultaneously the adversarial model will be completely unable to predict sensitive attributes given the predictions from the optimal classifier. Therefore, $C(x|\theta_C^*)$ and Z will be independent random variables.
- We now introduce the GAN model in which the generator plays the role of the classifier and the discriminator plays the role of the adversarial model.

Generative Adversarial Networks

- The purpose of deep learning is to learn a representation of high dimensional and noisy data using a sequence of differentiable functions, i.e., geometric transformations, that can perhaps be used for supervised learning tasks among others.
- It has had great success in discriminative models while generative models have fared worse due to the limitations of explicit maximum likelihood estimation (MLE).
- Adversarial learning as presented in the GAN model aims to overcome these problems by using implicit MLE.

Generative Adversarial Networks

- There are 2 main components to a GAN, the generator and the discriminator, that play an adversarial game against each other.
- In doing so the generator learns how to create realistic synthetic samples from noise, i.e., the latent space z , while the discriminator learns how to distinguish between a real sample and a synthetic sample.
- The representation learnt by the discriminator can later on be used for other supervised learning tasks, i.e., automatic feature engineering or representation learning.

Generative Adversarial Networks

Generator

- Assume that we have a prior belief on where the latent space z lies: $p(z)$.
- Given a draw from this latent space the generator G , a deep learner parameterized by θ_G , outputs a synthetic sample.

$$- G(z|\theta_G) : z \rightarrow x_{\text{synthetic}}$$

Generative Adversarial Networks

Discriminator

- The discriminator D is another deep learner parameterized by θ_D and it aims to classify if a sample is real or synthetic, i.e., if a sample is from the real data distribution: $p_{\text{data}}(x)$
- Or the synthetic data distribution: $p_G(x)$
- Let us denote the discriminator D as follows.

$$- D(x|\theta_D) : x \rightarrow [0, 1]$$

- Here we assume that the positive examples are from the real data distribution while the negative examples are from the synthetic data distribution.

Generative Adversarial Networks

Game

- A GAN simultaneously trains the discriminator to correctly classify real and synthetic examples while training the generator to create synthetic examples such that the discriminator incorrectly classifies real and synthetic examples.
- This 2 player minimax game has the following objective function.

$$\min_{G(z|\theta_G)} \max_{D(x|\theta_D)} V(D(x|\theta_D), G(z|\theta_G)) = \mathbb{E}_{x \sim p_{\text{data}}(x)} \log D(x|\theta_D) + \mathbb{E}_{z \sim p(z)} \log (1 - D(G(z|\theta_G)|\theta_D))$$

Generative Adversarial Networks

Game

- Please note that the above expression is basically the objective function of the discriminator.

$$\mathbb{E}_{x \sim p_{\text{data}}(x)} \log D(x|\theta_D) + \mathbb{E}_{x \sim p_G(x)} \log (1 - D(x|\theta_D))$$

- It is clear from how the game has been set up that we are trying to obtain a solution θ_D for D such that it maximizes $V(D, G)$ while simultaneously we are trying to obtain a solution θ_G for G such that it minimizes $V(D, G)$.

Generative Adversarial Networks

Game

- We do not simultaneously train D and G . We train them alternately: Train D and then train G while freezing D . We repeat this for a fixed number of steps.
- If the synthetic samples taken from the generator G are realistic then implicitly we have learnt the distribution $p_G(x)$. In other words, $p_G(x)$ can be seen as a good estimation of $p_{\text{data}}(x)$.
- The optimal solution will be as follows.

$$p_G(x) = p_{\text{data}}(x)$$

Generative Adversarial Networks

Game

- To show this let us find the optimal discriminator D^* given a generator G and sample x .

$$\begin{aligned}
 V(D, G) &= \mathbb{E}_{x \sim p_{\text{data}}(x)} \log D(x|\theta_D) + \mathbb{E}_{x \sim p_G(x)} \log (1 - D(x|\theta_D)) \\
 &= \int_x p_{\text{data}}(x) \log D(x|\theta_D) dx + \int_x p_G(x) \log (1 - D(x|\theta_D)) dx \\
 &= \int_x \underbrace{p_{\text{data}}(x) \log D(x|\theta_D) + p_G(x) \log (1 - D(x|\theta_D))}_{J(D(x|\theta_D))} dx
 \end{aligned}$$

Generative Adversarial Networks

Game

- Let us take a closer look at the discriminator's objective function for a sample x .

$$\begin{aligned}
 J(D(x|\theta_D)) &= p_{\text{data}}(x) \log D(x|\theta_D) + p_G(x) \log (1 - D(x|\theta_D)) \\
 \frac{\partial J(D(x|\theta_D))}{\partial D(x|\theta_D)} &= \frac{p_{\text{data}}(x)}{D(x|\theta_D)} - \frac{p_G(x)}{(1 - D(x|\theta_D))} \\
 0 &= \frac{p_{\text{data}}(x)}{D^*(x|\theta_{D^*})} - \frac{p_G(x)}{(1 - D^*(x|\theta_{D^*}))} \\
 p_{\text{data}}(x)(1 - D^*(x|\theta_{D^*})) &= p_G(x)D^*(x|\theta_{D^*}) \\
 p_{\text{data}}(x) - p_{\text{data}}(x)D^*(x|\theta_{D^*}) &= p_G(x)D^*(x|\theta_{D^*}) \\
 p_G(x)D^*(x|\theta_{D^*}) + p_{\text{data}}(x)D^*(x|\theta_{D^*}) &= p_{\text{data}}(x) \\
 D^*(x|\theta_{D^*}) &= \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_G(x)}
 \end{aligned}$$

Generative Adversarial Networks

Game

- We have found the optimal discriminator given a generator. Let us focus now on the generator's objective function which is essentially to minimize the discriminator's objective function.

$$\begin{aligned}
 J(G(x|\theta_G)) &= \mathbb{E}_{x \sim p_{\text{data}}(x)} \log D^*(x|\theta_{D^*}) + \mathbb{E}_{x \sim p_G(x)} \log (1 - D^*(x|\theta_{D^*})) \\
 &= \mathbb{E}_{x \sim p_{\text{data}}(x)} \log \left(\frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_G(x)} \right) + \mathbb{E}_{x \sim p_G(x)} \log \left(1 - \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_G(x)} \right) \\
 &= \mathbb{E}_{x \sim p_{\text{data}}(x)} \log \left(\frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_G(x)} \right) + \mathbb{E}_{x \sim p_G(x)} \log \left(\frac{p_G(x)}{p_{\text{data}}(x) + p_G(x)} \right) \\
 &= \int_x p_{\text{data}}(x) \log \left(\frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_G(x)} \right) dx + \int_x p_G(x) \log \left(\frac{p_G(x)}{p_{\text{data}}(x) + p_G(x)} \right) dx
 \end{aligned}$$

- We will note the Kullback–Leibler (KL) divergences in the above objective function for the generator:

$$D_{KL}(P||Q) = \int_x p(x) \log \left(\frac{p(x)}{q(x)} \right) dx$$

Generative Adversarial Networks

Game

- Recall the definition of a λ divergence.

$$D_\lambda(P||Q) = \lambda D_{KL}(P||\lambda P + (1 - \lambda)Q) + (1 - \lambda) D_{KL}(Q||\lambda P + (1 - \lambda)Q)$$

- If λ takes the value of 0.5 this is then called the Jensen-Shannon (JS) divergence. This divergence is symmetric and non-negative.

$$D_{JS}(P||Q) = 0.5 D_{KL} \left(P \middle| \middle| \frac{P+Q}{2} \right) + 0.5 D_{KL} \left(Q \middle| \middle| \frac{P+Q}{2} \right)$$

Generative Adversarial Networks

Game

- Keeping this in mind let us take a look again at the objective function of the generator.

$$\begin{aligned}
 J(G(x|\theta_G)) &= \int_x p_{\text{data}}(x) \log \left(\frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_G(x)} \right) dx + \int_x p_G(x) \log \left(\frac{p_G(x)}{p_{\text{data}}(x) + p_G(x)} \right) dx \\
 &= \int_x p_{\text{data}}(x) \log \left(\frac{2}{2} \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_G(x)} \right) dx + \int_x p_G(x) \log \left(\frac{2}{2} \frac{p_G(x)}{p_{\text{data}}(x) + p_G(x)} \right) dx \\
 &= \int_x p_{\text{data}}(x) \log \left(\frac{1}{2} \frac{1}{0.5} \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_G(x)} \right) dx + \int_x p_G(x) \log \left(\frac{1}{2} \frac{1}{0.5} \frac{p_G(x)}{p_{\text{data}}(x) + p_G(x)} \right) dx \\
 &= \int_x p_{\text{data}}(x) \left[\log(0.5) + \log \left(\frac{p_{\text{data}}(x)}{0.5(p_{\text{data}}(x) + p_G(x))} \right) \right] dx \\
 &\quad + \int_x p_G(x) \left[\log(0.5) + \log \left(\frac{p_G(x)}{0.5(p_{\text{data}}(x) + p_G(x))} \right) \right] dx
 \end{aligned}$$

Generative Adversarial Networks

Game

$$\begin{aligned}
 &= \log \left(\frac{1}{4} \right) + \int_x p_{\text{data}}(x) \left[\log \left(\frac{p_{\text{data}}(x)}{0.5(p_{\text{data}}(x) + p_G(x))} \right) \right] dx \\
 &\quad + \int_x p_G(x) \left[\log \left(\frac{p_G(x)}{0.5(p_{\text{data}}(x) + p_G(x))} \right) \right] dx \\
 &= -\log(4) + D_{KL} \left(p_{\text{data}}(x) \left\| \frac{p_{\text{data}}(x) + p_G(x)}{2} \right\| \right) + D_{KL} \left(p_G(x) \left\| \frac{p_{\text{data}}(x) + p_G(x)}{2} \right\| \right) \\
 &= -\log(4) + 2 \left(0.5 D_{KL} \left(p_{\text{data}}(x) \left\| \frac{p_{\text{data}}(x) + p_G(x)}{2} \right\| \right) + 0.5 D_{KL} \left(p_G(x) \left\| \frac{p_{\text{data}}(x) + p_G(x)}{2} \right\| \right) \right) \\
 &= -\log(4) + 2 D_{JS}(p_{\text{data}}(x) || p_G(x))
 \end{aligned}$$

Generative Adversarial Networks

Game

- It is clear from the objective function of the generator above that the global minimum value attained is $-\log(4)$ which occurs when the following holds.

$$p_G(x) = p_{\text{data}}(x)$$

- When the above holds the Jensen-Shannon divergence, i.e., $D_{JS}(p_{\text{data}}(x) || p_G(x))$, will be zero. Hence we have shown that the optimal solution is as follows.

$$p_G(x) = p_{\text{data}}(x)$$

Generative Adversarial Networks

Game

- Note that the above implies that the optimal value for the discriminator is as follows.

$$\begin{aligned}
 D^*(x|\theta_{D^*}) &= \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_G(x)} \\
 &= \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_{\text{data}}(x)} \\
 &= \frac{1}{2}
 \end{aligned}$$

- The optimal discriminator cannot distinguish between real and synthetic data and this is precisely what we need for the adversarial model to do as well: To not be able to distinguish between sensitive attributes.

Classifier and Adversarial Model

- Returning back to our problem at hand, we use adversarial training to find a classifier such that its predictions cannot be used to predict for sensitive attributes by the adversarial model.
- This 2 player minimax game has the following objective function. Note that \tilde{y} are predictions from the classifier.

$$\begin{aligned}
& \max_{C(x|\theta_C)} \min_{R(C(x|\theta_C)|\theta_R)} V(C(x|\theta_C), R(C(x|\theta_C)|\theta_R)) \\
&= \underbrace{\mathbb{E}_{x \sim P_X(x)} \mathbb{E}_{y \sim P_{Y|X}(y|x)} \log C(x|\theta_C)}_{\text{Classifier maximizes this term}} - \underbrace{\mathbb{E}_{x \sim P_X(x)} \mathbb{E}_{\tilde{y} \sim P_{\tilde{Y}|X}(\tilde{y}|x)} \mathbb{E}_{z \sim P_{Z|\tilde{Y}}(z|\tilde{y})} \log R(C(x|\theta_C)|\theta_R)}_{\text{Classifier minimizes this term and the adversarial maximizes this term}}
\end{aligned}$$

Classifier and Adversarial Model

- We note that the classifier is aiming to maximize it's likelihood while minimizing that of the adversarial while the adversarial is aiming to maximize it's own likelihood.
- The solution to this 2 player game theoretically, albeit not empirically as we shall observe, results in a classifier that is both optimal and fair.
- For the optimal classifier $C(x|\theta_C^*)$ to be fair as well it has to be shown to be independent from the sensitive attributes Z .
- Essentially this means that the adversarial model $R(C(x|\theta_C^*)|\theta_R)$ is not able to predict the sensitive attributes, i.e., it essentially becomes an uninformative prior on the sensitive attributes as opposed to an informative posterior since the predictions, i.e., new information, from the classifier do not help in updating this uninformative prior, i.e., $P(Z)$.

Classifier and Adversarial Model

- In light of the above information, the objective function of the 2 player game is bounded above as follows.
- The upper bound holds when $C(x|\theta_C^*) \perp\!\!\!\perp Z$.
- This basically means that $R(C(x|\theta_C^*)|\theta_R)$ cannot be used to predict for the sensitive attributes, i.e., we force the adversarial model towards the uninformative prior $P(Z)$.

$$\begin{aligned}
& \max_{C(x|\theta_C)} \min_{R(C(x|\theta_C)|\theta_R)} V(C(x|\theta_C), R(C(x|\theta_C)|\theta_R)) \\
&= \underbrace{\mathbb{E}_{x \sim P_X(x)} \mathbb{E}_{y \sim P_{Y|X}(y|x)} \log C(x|\theta_C)}_{\text{Classifier maximizes this term}} - \underbrace{\mathbb{E}_{x \sim P_X(x)} \mathbb{E}_{\tilde{y} \sim P_{\tilde{Y}|X}(\tilde{y}|x)} \mathbb{E}_{z \sim P_{Z|\tilde{Y}}(z|\tilde{y})} \log R(C(x|\theta_C)|\theta_R)}_{\text{Classifier minimizes this term and the adversarial maximizes this term}} \\
&\leq \mathbb{E}_{x \sim P_X(x)} \mathbb{E}_{y \sim P_{Y|X}(y|x)} \log C(x|\theta_C^*) - \mathbb{E}_{z \sim P_Z(z)} \log R(C(x|\theta_C^*)|\theta_R)
\end{aligned}$$

- The optimal classifier $C(x|\theta_C^*)$ maximizes it's own likelihood and minimizes the likelihood of the adversarial such that the upper bound holds, i.e., $R(C(x|\theta_C^*)|\theta_R)$ is forced to move towards the uninformative prior $P(Z)$. Hence the classifier is theoretically both optimal and fair.
- Empirically the upper bound is not strict therefore a trade off has to be made between the classifier being optimal and fair. We shall see this in the following experiment using the "Census Income" data.

Fairness Metric

- To assess a classifier's performance as optimal we can use the AUROC metric.
- To assess a classifier's performance as fair we can use the p-rule metric. This is defined as follows.

$$- \min \left(\frac{\mu(1_{C(x|\theta_C^*) > \tau} | z=1)}{\mu(1_{C(x|\theta_C^*) > \tau} | z=0)}, \frac{\mu(1_{C(x|\theta_C^*) > \tau} | z=0)}{\mu(1_{C(x|\theta_C^*) > \tau} | z=1)} \right) \geq \frac{p}{100}$$

- The intuition above is that at least $p\%$ of the time the optimal classifier should give the same predicted outcome label under different levels of the sensitive attribute(s), i.e., think gender and race.
- The optimal classifier is deemed fair $p\%$ of the time if it's predicted outcome label does not deviate across gender or race.
- Ideally one would want a high AUROC and a high p-rule. In the experiment on “Census Income” we observe that there is a trade-off to be made between having an optimal classifier and a fair classifier.
- Let me emphatically emphasize that trading off some degree of classifier performance to generate a sufficiently fair classifier, say that has a p-rule score of 80%, is absolutely worth it. This represents a clear opportunity for machine learning to remove the hidden bias in data and make fair and simultaneously optimal decisions.

R Code

Load up R libraries.

```
library(package = fastDummies)
library(package = data.table)
library(package = dplyr)

##
## Attaching package: 'dplyr'
##
## The following objects are masked from 'package:data.table':
##
##   between, first, last
##
## The following objects are masked from 'package:stats':
##
##   filter, lag
##
## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union
library(package = keras)
library(package = hmeasure)
library(package = ggplot2)
library(package = scales)
```

Load up data, do some basic preprocessing and create the set of features, sensitive attributes and target variable.

```
data <- fread(
  input = "/home/hamaad/Downloads/adult.data",
  col.names = c(
    "age",
    "workclass",
    "fnlwgt",
    "education",
    "education_num",
    "marital_status",
    "occupation",
    "relationship",
    "race",
    "sex",
    "capital_gain",
    "capital_loss",
    "hours_per_week",
    "country",
    "target"
  ),
  na.strings = "?"
)

data <- data[which(data$race %in% c("White", "Black"))]

sensitive.attrs <- c("race", "sex")
```



```

data[, race := ifelse(race == "White", 1, 0)]
data[, sex := ifelse(sex == "Male", 1, 0)]

# 0 is black and 1 is white.
# 0 is female and 1 is male.
Z <- data[, sensitive.attrs, with = FALSE]

y <- data[, target := ifelse(target == ">50K", 1, 0)][, target]

X <- data[, ":(target = NULL,
  race = NULL,
  sex = NULL)]

is.missing.ind <- is.na(X)
X[is.missing.ind] <- "Unknown"

X <- dummy_cols(.data = X, select_columns = c(
  "workclass",
  "education",
  "marital_status",
  "occupation",
  "relationship",
  "country"
), remove_first_dummy = TRUE)
X <- X[, ":(workclass = NULL,
  education = NULL,
  marital_status = NULL,
  occupation = NULL,
  relationship = NULL,
  country = NULL)]

```

Split the data into train and test sets.

```

tr.ind <- sample(
  x = 1:dim(X)[1],
  size = as.integer(0.5 * dim(X)[1]),
  replace = FALSE
)
ts.ind <- setdiff(
  x = 1:dim(X)[1],
  y = tr.ind
)

```

Standardise the data.

```

mu <- apply(
  X = X[tr.ind],
  MARGIN = 2,
  FUN = mean
)

```

```

std <- apply(
  X = X[tr.ind],
  MARGIN = 2,
  FUN = sd
) + .Machine$double.eps

standardised.X.tr <- sweep(
  x = X[tr.ind],
  MARGIN = 2,
  STATS = mu,
  FUN = "-"
)
standardised.X.tr <- sweep(
  x = standardised.X.tr,
  MARGIN = 2,
  STATS = std,
  FUN = "/"
)
standardised.X.ts <- sweep(
  x = X[ts.ind],
  MARGIN = 2,
  STATS = mu,
  FUN = "-"
)
standardised.X.ts <- sweep(
  x = standardised.X.ts,
  MARGIN = 2,
  STATS = std,
  FUN = "/"
)

```

Create the classifier.

```

classifier.input.layer <- layer_input(
  shape = c(dim(standardised.X.tr)[2]),
  name = "classifier.input"
)

classifier.hidden.layers <- classifier.input.layer %>%
  layer_dense(
    units = 32,
    activation = "relu",
    name = "classifier.hidden.layer.1"
  ) %>%
  layer_dropout(
    rate = 0.2,
    name = "classifier.dropout.layer.1"
  ) %>%
  layer_dense(
    units = 32,
    activation = "relu",
    name = "classifier.hidden.layer.2"
  ) %>%

```

```

layer_dropout(
  rate = 0.2,
  name = "classifier.dropout.layer.2"
) %>%
layer_dense(
  units = 32,
  activation = "relu",
  name = "classifier.hidden.layer.3"
) %>%
layer_dropout(
  rate = 0.2,
  name = "classifier.dropout.layer.3"
)

classifier.output.layer <- classifier.hidden.layers %>%
  layer_dense(
    units = 1,
    activation = "sigmoid",
    name = "classifier.output"
  )

classifier.model <- keras_model(
  inputs = classifier.input.layer,
  outputs = classifier.output.layer
)

classifier.model %>% compile(
  optimizer = "adam",
  loss = "binary_crossentropy"
)

summary(classifier.model)

```

```

## Model: "functional_1"
## -----
## Layer (type)                Output Shape          Param #
## -----
## classifier.input (InputLayer)  [(None, 94)]          0
## -----
## classifier.hidden.layer.1 (Dense)  (None, 32)           3040
## -----
## classifier.dropout.layer.1 (Dropout) (None, 32)            0
## -----
## classifier.hidden.layer.2 (Dense)  (None, 32)           1056
## -----
## classifier.dropout.layer.2 (Dropout) (None, 32)            0
## -----
## classifier.hidden.layer.3 (Dense)  (None, 32)           1056
## -----
## classifier.dropout.layer.3 (Dropout) (None, 32)            0
## -----
## classifier.output (Dense)         (None, 1)             33
## -----
## Total params: 5,185

```

```
## Trainable params: 5,185
## Non-trainable params: 0
## -----
```

Create the adversarial model.

```
adversarial.input.layer <- layer_input(
  shape = c(1),
  name = "adversarial.input"
)

adversarial.hidden.layers <- adversarial.input.layer %>%
  layer_dense(
    units = 32,
    activation = "relu",
    name = "adversarial.hidden.layer.1"
  ) %>%
  layer_dense(
    units = 32,
    activation = "relu",
    name = "adversarial.hidden.layer.2"
  ) %>%
  layer_dense(
    units = 32,
    activation = "relu",
    name = "adversarial.hidden.layer.3"
  )

race.layer <- adversarial.hidden.layers %>%
  layer_dense(
    units = 1,
    activation = "sigmoid",
    name = "race.layer.output"
  )

sex.layer <- adversarial.hidden.layers %>%
  layer_dense(
    units = 1,
    activation = "sigmoid",
    name = "sex.layer.output"
  )

adversarial.model <- keras_model(
  inputs = adversarial.input.layer,
  outputs = c(
    race.layer,
    sex.layer
  )
)

summary(adversarial.model)
```

```
## Model: "functional_3"
```

```
## -----
## Layer (type)           Output Shape      Param #   Connected to
```

```

## =====
## adversarial.input (InputL [(None, 1)])      0
## -----
## adversarial.hidden.layer. (None, 32)      64      adversarial.input[0][0]
## -----
## adversarial.hidden.layer. (None, 32)      1056     adversarial.hidden.layer.1[
## -----
## adversarial.hidden.layer. (None, 32)      1056     adversarial.hidden.layer.2[
## -----
## race.layer.output (Dense) (None, 1)      33      adversarial.hidden.layer.3[
## -----
## sex.layer.output (Dense) (None, 1)      33      adversarial.hidden.layer.3[
## =====
## Total params: 2,242
## Trainable params: 2,242
## Non-trainable params: 0
## -----

```

We freeze the classifier weights and unfreeze the adversarial model weights in order to train the adversarial model.

```

unfrozen.adversarial.model <- keras_model(
  inputs = classifier.input.layer,
  outputs = adversarial.model(object = classifier.output.layer)
)
freeze_weights(
  object = classifier.model,
  from = "classifier.input",
  to = "classifier.output"
)

unfreeze_weights(
  object = adversarial.model,
  from = "adversarial.input",
  to = "race.layer.output"
)
unfreeze_weights(
  object = adversarial.model,
  from = "adversarial.input",
  to = "sex.layer.output"
)

unfrozen.adversarial.model %>% compile(
  optimizer = "adam",
  loss = c(
    "binary_crossentropy",
    "binary_crossentropy"
  )
)

summary(unfrozen.adversarial.model)

## Model: "functional_5"
## -----

```

| ## Layer (type) | Output Shape | Param # |
|---|------------------------|---------|
| ## classifier.input (InputLayer) | [(None, 94)] | 0 |
| ## classifier.hidden.layer.1 (Dense) | (None, 32) | 3040 |
| ## classifier.dropout.layer.1 (Dropout) | (None, 32) | 0 |
| ## classifier.hidden.layer.2 (Dense) | (None, 32) | 1056 |
| ## classifier.dropout.layer.2 (Dropout) | (None, 32) | 0 |
| ## classifier.hidden.layer.3 (Dense) | (None, 32) | 1056 |
| ## classifier.dropout.layer.3 (Dropout) | (None, 32) | 0 |
| ## classifier.output (Dense) | (None, 1) | 33 |
| ## functional_3 (Functional) | [(None, 1), (None, 1)] | 2242 |
| ## ===== | | |
| ## Total params: 7,427 | | |
| ## Trainable params: 2,242 | | |
| ## Non-trainable params: 5,185 | | |
| ## ===== | | |

We freeze the adversarial weights and unfreeze the classifier model weights in order to train the classifier model.

```
frozen.adversarial.model <- keras_model(
  inputs = classifier.input.layer,
  outputs = c(
    classifier.output.layer,
    adversarial.model(object = classifier.output.layer)
  )
)

unfreeze_weights(
  object = classifier.model,
  from = "classifier.input",
  to = "classifier.output"
)

freeze_weights(
  object = adversarial.model,
  from = "adversarial.input",
  to = "race.layer.output"
)

freeze_weights(
  object = adversarial.model,
  from = "adversarial.input",
  to = "sex.layer.output"
)
```

```
summary(frozen.adversarial.model)
```

```
## Model: "functional_7"
## -----
## Layer (type)                Output Shape          Param #
## =====
## classifier.input (InputLayer)  [(None, 94)]          0
## -----
## classifier.hidden.layer.1 (Dense)  (None, 32)           3040
## -----
## classifier.dropout.layer.1 (Dropout) (None, 32)            0
## -----
## classifier.hidden.layer.2 (Dense)  (None, 32)           1056
## -----
## classifier.dropout.layer.2 (Dropout) (None, 32)            0
## -----
## classifier.hidden.layer.3 (Dense)  (None, 32)           1056
## -----
## classifier.dropout.layer.3 (Dropout) (None, 32)            0
## -----
## classifier.output (Dense)          (None, 1)             33
## -----
## functional_3 (Functional)          [(None, 1), (None, 1)] 2242
## =====
## Total params: 7,427
## Trainable params: 5,185
## Non-trainable params: 2,242
## -----
```

```
frozen.adversarial.model %>% compile(
  optimizer = "adam", metrics = NULL,
  loss_weights = c(1, -100, -50),
  loss = c(
    "binary_crossentropy",
    "binary_crossentropy",
    "binary_crossentropy"
  )
)
```

Pre-train the classifier.

```
unfreeze_weights(
  object = classifier.model,
  from = "classifier.input",
  to = "classifier.output"
)

classifier.model %>% fit(
  x = as.matrix(standardised.X.tr),
  y = y[tr.ind],
  epochs = 25,
  batch_size = 100,
  verbose = 1,
)
```

```

  shuffle = TRUE
)

```

Pre-train the adversarial model.

```

freeze_weights(
  object = classifier.model,
  from = "classifier.input",
  to = "classifier.output"
)

unfreeze_weights(
  object = adversarial.model,
  from = "adversarial.input",
  to = "race.layer.output"
)

unfreeze_weights(
  object = adversarial.model,
  from = "adversarial.input",
  to = "sex.layer.output"
)

unfrozen.adversarial.model %>% fit(
  x = as.matrix(standardised.X.tr),
  y = list(
    as.matrix(Z[tr.ind, "race", with = FALSE]),
    as.matrix(Z[tr.ind, "sex", with = FALSE])
  ),
  epochs = 25,
  batch_size = 100,
  verbose = 1,
  shuffle = TRUE
)

```

Make some basic plots and check the AUC and P-rule measures to check whether the model is optimal and fair.

```

pred.test <- classifier.model %>% predict(x = as.matrix(standardised.X.ts))

plot.out <- data.table(
  pred.test,
  Z[ts.ind]
)

colnames(plot.out) <- c(
  "predictions",
  sensitive.attrs
)

thrs <- 0.5

y.z.1 <- ifelse(test = plot.out$predictions[which(plot.out$race == 1)] > thrs,
  yes = 1.0,

```



```

    no = 0.0
  )
  y.z.0 <- ifelse(test = plot.out$predictions[which(plot.out$race == 0)] > thrs,
    yes = 1.0,
    no = 0.0
  )
  odds <- mean(x = y.z.1) / mean(x = y.z.0)
  out.race <- min(x = c(odds, 1 / odds))

  y.z.1 <- ifelse(test = plot.out$predictions[which(plot.out$sex == 1)] > thrs,
    yes = 1.0,
    no = 0.0
  )
  y.z.0 <- ifelse(test = plot.out$predictions[which(plot.out$sex == 0)] > thrs,
    yes = 1.0,
    no = 0.0
  )
  odds <- mean(x = y.z.1) / mean(x = y.z.0)
  out.sex <- min(x = c(odds, 1.0 / odds))

output.plot <- ggplot(data = plot.out) +
  geom_density(aes(
    x = predictions,
    fill = factor(race)
  ),
  alpha = 0.5
) +
  ggtitle(paste(
    "Predictions for income being higher than $50,000 per year\nGrouped by race\nP-rule:",
    percent(out.race),
    "\nAUROC:",
    percent(HMeasure(
      true.class = y[ts.ind],
      scores = pred.test
    )$metrics$AUC)
  )) +
  scale_fill_discrete(
    name = "Race",
    labels = c("Black", "White")
  ) +
  scale_x_continuous(labels = percent) +
  xlab(label = "Classifier predictions") +
  ylab(label = "Density")
ggsave(
  filename = "/home/hamaad/Projects/fair_ml_R/pre_train_result_race.png",
  plot = output.plot
)

```

Saving 6.5 x 4.5 in image

```

output.plot <- ggplot(data = plot.out) +
  geom_density(aes(
    x = predictions,
    fill = factor(sex)
  )

```

```

),
alpha = 0.5
) +
ggtitle(paste(
  "Predictions for income being higher than $50,000 per year\nGrouped by gender\nP-rule:",
  percent(out.sex),
  "\nAUROC:",
  percent(HMeasure(
    true.class = y[ts.ind],
    scores = pred.test
  )$metrics$AUC)
)) +
scale_fill_discrete(
  name = "Gender",
  labels = c("Female", "Male")
) +
scale_x_continuous(labels = percent) +
xlab(label = "Classifier predictions") +
ylab(label = "Density")
ggsave(
  filename = "/home/hamaad/Projects/fair_ml_R/pre_train_result_gender.png",
  plot = output.plot
)

```

Saving 6.5 x 4.5 in image

Initial model performance

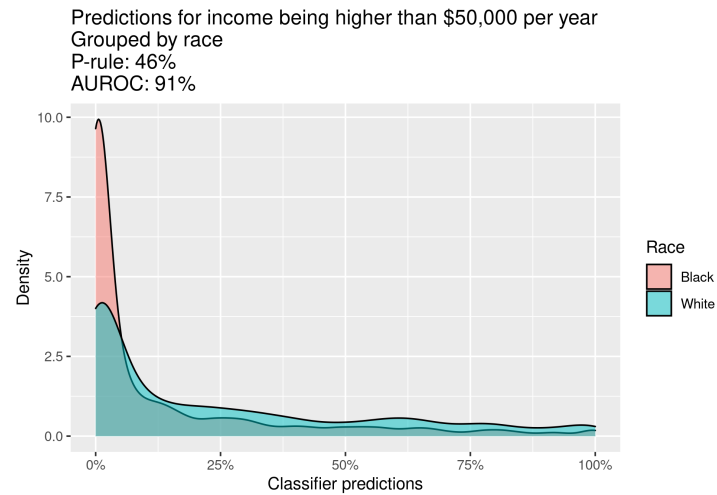


Figure 1: Initial model performance with regards to optimality and fairness grouped by race

Initial model performance

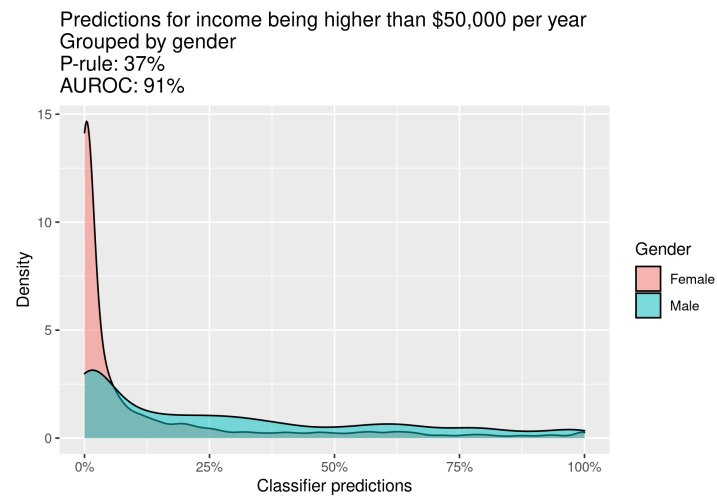


Figure 2: Initial model performance with regards to optimality and fairness grouped by gender

Train the classifier and the adversarial model alternately.

```
for (i in 1:500) {  
  freeze_weights(  
    object = classifier.model,  
    from = "classifier.input",  
    to = "classifier.output"  
  )  
  unfreeze_weights(  
    object = adversarial.model,  
    from = "adversarial.input",
```

```

    to = "race.layer.output"
  )
  unfreeze_weights(
    object = adversarial.model,
    from = "adversarial.input",
    to = "sex.layer.output"
  )
  batch.ind <- sample(x = 1:length(tr.ind), size = 100, replace = FALSE)
  unfrozen.adversarial.model %>% train_on_batch(
    x = as.matrix(standardised.X.tr)[batch.ind, ],
    y = list(
      as.matrix(Z[tr.ind, "race", with = FALSE])[batch.ind, ],
      as.matrix(Z[tr.ind, "sex", with = FALSE])[batch.ind, ]
    )
  )
  pred.test <- classifier.model %>% predict(as.matrix(standardised.X.ts))

  plot.out <- data.table(pred.test, Z[ts.ind])
  colnames(plot.out) <- c(
    "predictions",
    sensitive.attrs
  )

  y.z.1 <- ifelse(test = plot.out$predictions[which(plot.out$race == 1)] > thr,
    yes = 1.0,
    no = 0.0
  )
  y.z.0 <- ifelse(test = plot.out$predictions[which(plot.out$race == 0)] > thr,
    yes = 1.0,
    no = 0.0
  )
  odds <- mean(x = y.z.1) / mean(x = y.z.0)
  out.race <- min(x = c(odds, 1 / odds))

  y.z.1 <- ifelse(test = plot.out$predictions[which(plot.out$sex == 1)] > thr,
    yes = 1.0,
    no = 0.0
  )
  y.z.0 <- ifelse(test = plot.out$predictions[which(plot.out$sex == 0)] > thr,
    yes = 1.0,
    no = 0.0
  )
  odds <- mean(x = y.z.1) / mean(x = y.z.0)
  out.sex <- min(x = c(odds, 1.0 / odds))

  unfreeze_weights(
    object = classifier.model,
    from = "classifier.input",
    to = "classifier.output"
  )
  freeze_weights(
    object = adversarial.model,
    from = "adversarial.input",

```

```

    to = "race.layer.output"
  )
  freeze_weights(
    object = adversarial.model,
    from = "adversarial.input",
    to = "sex.layer.output"
  )

  batch.ind <- sample(
    x = 1:length(tr.ind),
    size = 100,
    replace = FALSE
  )
  frozen.adversarial.model %>% train_on_batch(
    x = as.matrix(standardised.X.tr)[batch.ind, ],
    y = list(
      as.matrix(data.table(y)[tr.ind])[batch.ind, ],
      as.matrix(Z[tr.ind, "race", with = FALSE])[batch.ind, ],
      as.matrix(Z[tr.ind, "sex", with = FALSE])[batch.ind, ]
    )
  )
}

```

Final results.

```

output.plot <- ggplot(data = plot.out) +
  geom_density(aes(
    x = predictions,
    fill = factor(race)
  ),
  alpha = 0.5
) +
  ggtitle(paste(
    "Predictions for income being higher than $50,000 per year\nGrouped by race\nP-rule:",
    percent(out.race),
    "\nAUROC:",
    percent(HMeasure(true.class = y[ts.ind], scores = pred.test)$metrics$AUC)
  )) +
  scale_fill_discrete(
    name = "Race",
    labels = c("Black", "White")
  ) +
  scale_x_continuous(labels = percent) +
  xlab(label = "Classifier predictions") +
  ylab(label = "Density")
ggsave(
  filename = "/home/hamaad/Projects/fair_ml_R/post_train_result_race.png",
  plot = output.plot
)

```

Saving 6.5 x 4.5 in image

```

output.plot <- ggplot(data = plot.out) +
  geom_density(aes(
    x = predictions,
    fill = factor(sex)
  ),
  alpha = 0.5
) +
  ggtitle(paste(
    "Predictions for income being higher than $50,000 per year\nGrouped by gender\nP-rule:",
    percent(out.sex),
    "\nAUROC:",
    percent(HMeasure(true.class = y[ts.ind], scores = pred.test)$metrics$AUC)
  )) +
  scale_fill_discrete(
    name = "Gender",
    labels = c("Female", "Male")
  ) +
  scale_x_continuous(labels = percent) +
  xlab(label = "Classifier predictions") +
  ylab(label = "Density")
ggsave(
  filename = "/home/hamaad/Projects/fair_ml_R/post_train_result_gender.png",
  plot = output.plot
)

```

```
## Saving 6.5 x 4.5 in image
```

Final model performance

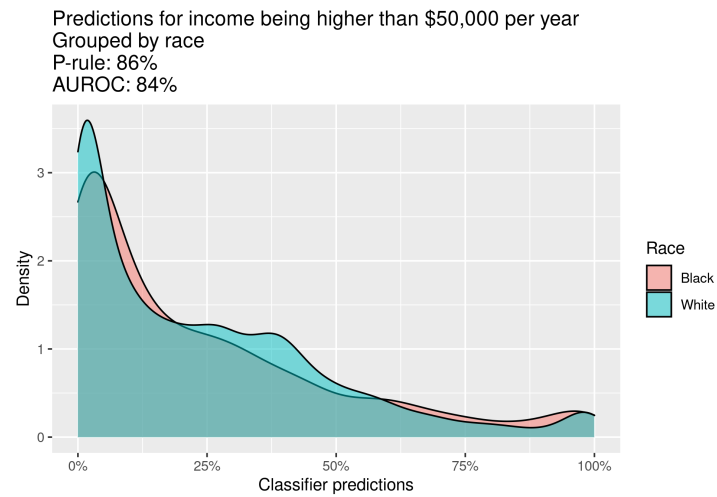


Figure 3: Final model performance with regards to optimality and fairness grouped by race

Final model performance

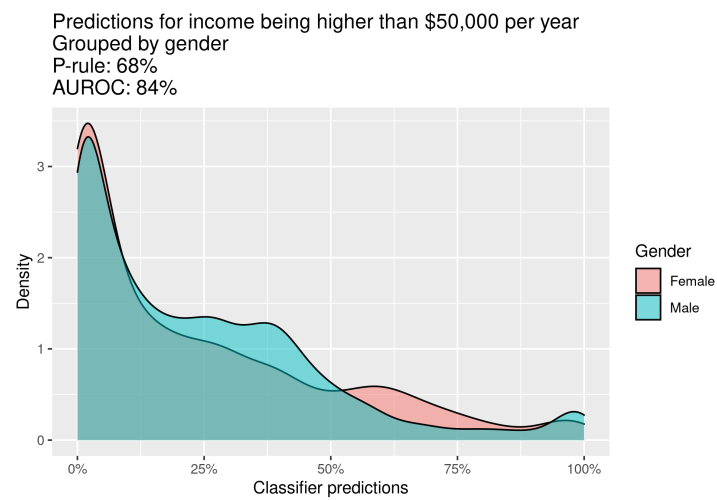


Figure 4: Final model performance with regards to optimality and fairness grouped by gender

Conclusion

- We have shown how to use adversarial training, inspired by the GAN model, in order to provide a trade off between an optimal classifier and a fair classifier, i.e., one that is not sensitive to attributes such as race and gender.
- The results shows that machine learning might actually be used to train fair decision making pipelines such that the data's hidden bias is not reflected in this decision making, even though sensitive attributes might have been explicitly removed.
- This is a powerful and interesting way of using the idea of adversarial training derived from the GAN model for a goal that is highly desirable from an ethics perspective.

References

1. Goodfellow, I., Bengio, Y. and Courville A. (2016). Deep Learning (MIT Press).
2. Geron, A. (2017). Hands-On Machine Learning with Scikit-Learn & Tensorflow (O'Reilly).
3. Radford, A., Luke, M. and Chintala, S. (2015). Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks (<https://arxiv.org/abs/1511.06434>).
4. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y. (2014). Generative Adversarial Networks (<https://arxiv.org/abs/1406.2661>).
5. Chollet, F., Allaire, J. J. (2018). Deep Learning with R (Manning).
6. <https://blog.godatadriven.com/fairness-in-ml>
7. Louppe, G., Kagan, M., Kranmer, K. (2017). Learning to Pivot with Adversarial Networks (<https://papers.nips.cc/paper/6699-learning-to-pivot-with-adversarial-networks.pdf>)
8. Zafar, M. B., Valera, I., Rodriguez, M. G., Gummadi, K. P. (2017). Fairness Constraints: Mechanisms for Fair Classification (<https://arxiv.org/pdf/1507.05259.pdf>)
9. <https://towardsdatascience.com/automatic-feature-engineering-using-generative-adversarial-networks-8e24b3c16bf3>

Blog on Medium and Code on GitHub

- On GitHub:
 - https://github.com/hamaadshah/autoencoders_tensorflow
 - https://github.com/hamaadshah/gan_tensorflow
 - https://github.com/hamaadshah/market_risk_gan_tensorflow
 - https://github.com/hamaadshah/autoencoders_pytorch
 - https://github.com/hamaadshah/gan_deeplearning4j
 - https://github.com/hamaadshah/fair_ml_R
- On Medium - search for my article titles below:
 - Automatic feature engineering using deep learning and Bayesian inference.
 - Automatic feature engineering using Generative Adversarial Networks.
 - Using Bidirectional Generative Adversarial Networks to estimate Value-at-Risk for Market Risk Management.