Overcoming set imbalance in data driven parameterization: A case study of gravity wave momentum transport

L. Minah Yang ¹, Edwin P. Gerber ¹

¹Center for Atmosphere Ocean Science, Courant Institute of Mathematical Sciences, New York University, 5 New York, New York, USA.

Key Points:

1

2

3

4

6

7

| 8 | ٠ | Unresolved geophysical processes often exhibit long tail distributions, which leads |
|----|---|---|
| 9 | | to imbalanced datasets for data-driven parameterizations. |
| 10 | • | Two strategies to overcome data imbalance are presented, where either the sam- |
| 11 | | pling or loss function is modified to better capture the tails. |
| 12 | • | Proof of concept is demonstrated by using a wind range metric to improve a ma- |
| 13 | | chine learning emulator of a physics based gravity wave parameterization. |

Corresponding author: L. Minah Yang, minah.yang@nyu.edu

14 Abstract

Machine learning for the parameterization of subgrid-scale processes in climate models 15 has been widely researched and adopted in a few models. A key challenge in develop-16 ing data-driven parameterization schemes is how to properly represent rare, but impor-17 tant events that occur in geoscience datasets. We investigate and develop strategies to 18 reduce errors caused by insufficient sampling in the rare data regime, under constraints 19 of no new data and no further expansion of model complexity. Resampling and impor-20 tance weighting strategies are constructed with user defined parameters that systemat-21 ically vary the sampling/weighting rates in a linear fashion and curb too much oversam-22 pling. Applying this new method to a case study of gravity wave momentum transport 23 reveals that the resampling strategy can successfully improve errors in the rare regime 24 at little to no loss in accuracy overall in the dataset. The success of the strategy, how-25 ever, depends on the complexity of the model. More complex models can overfit the tails 26 of the distribution when using non-optimal parameters of the resampling strategy. 27

²⁸ Plain Language Summary

Subgrid-scale parameterizations are a part of climate models that represent effects 29 of processes that cannot be directly modelled. In recent years, there have been many ef-30 forts to improve upon these parameterizations by applying machine learning techniques. 31 Since these methods rely heavily on the dataset they are learning from, it is important 32 33 to consider the frequency at which important events occur within the dataset because they are adept at learning frequent events at high accuracy but are prone to learning rare 34 but important events at low accuracy. To remedy this *data imbalance* problem, we de-35 veloped a resampling methodology that can be easily adjusted by tuning just two pa-36 rameters. We find that a right combination of those parameters can improve the accu-37 racy of an ML model at the rare event regime while keeping the accuracy high in the fre-38 quent regime. However, a "wrong" combination can actually increase the errors at the 39 rare event regime by overfitting to that regime. 40

41 **1** Introduction

Machine learning techniques have been used to develop data driven parameteriza-42 tion of un- or under-resolved processes in climate models, including a comprehensive rep-43 resentation of all missing terms, either at once (Brenowitz & Bretherton, 2019) or sep-44 arately (Yuval et al., 2021), or specific processes, including gravity wave momentum trans-45 port (Chantry et al., 2021; Espinosa et al., 2022) and radiative transfer (Ukkonen, 2022). 46 None of these attempts yielded a perfect sub-grid scale model, begging a general ques-47 tion: what can one do to improve a given data-driven parameterization? As these pro-48 cesses, and geoscience datasets more generally, are often high-dimensional and exhibit 49 long-tailed distributions, a common problem is to properly learn rare and extreme events. 50 This is particularly problematic if these extreme events have an outsized impact on the 51 climate, or become more prevalent in a changing climate. How can we capture impor-52 tant but rare events from the tail of the distribution as best as possible given the dataset 53 available to us? This is a data imbalance problem, and we propose strategies to com-54 bat it in this paper. 55

Set imbalance is a common challenge in machine learning (ML). In binary classification, the imbalanced dataset problem refers to a skewed distribution of the two target classes in a dataset. A naive learning algorithm will inherit an asymmetric class representation in the dataset, and will typically produce classifiers that predict the minority class with lower accuracy than it does for the majority class. These biased classifiers prove even more problematic when the minority class holds more importance or utility. As this combination of challenges is ubiquitous in real datasets, many methods that curb and minimize biases that stem from imbalanced datasets have been developed, as reviewed
 by He and Garcia (2009) and Krawczyk (2016).

⁶⁵ Data imbalance poses difficulties for ML tasks outside of binary classification. While ⁶⁶ it is straightforward to extend methods for treating imbalanced datasets for binary to ⁶⁷ multi-class classification, it has proven more difficult to extend this for regression tasks. ⁶⁸ Here, one seeks to learn a function g from a set of inputs \vec{x} to outputs \vec{y} where the ex-⁶⁹ ample pairs (\vec{x}, \vec{y}) is unevenly distributed. As with the classification problem, the task ⁷⁰ is particularly hard if we care especially about the behavior of g for rare pairs of (\vec{x}, \vec{y}) .

In this paper, we explore systematic methods for overcoming data imbalance in re-71 gression tasks, illustrating them with a case study of data driven parameterization grav-72 ity wave (GW) momentum transport. Gravity waves play an important role in forcing 73 the large scale atmospheric circulation, but their small scale makes them challenging to 74 properly represent directly. We seek a function g that maps vertical profiles of the re-75 solved wind, temperature, and GW source information within a column of an atmospheric 76 model: \vec{x} , to the profiles of the grid scale momentum tendency by unresolved gravity waves 77 associated with this large scale environment: \vec{y} . We assume limited resources, in that 78 one cannot simply increase the size of the dataset or complexity of our model q to over-79 come the problem: the goal is to work with the data and model one has on hand. 80

First steps have been taken towards deriving data-driven schemes for GWs by ex-81 ploring how well machine learning approaches can emulate existing, physics based pa-82 rameterizations (Chantry et al., 2021; Sun et al., 2023). Both studies found that data 83 imbalance was challenging, particularly for capturing the momentum forcing by grav-84 ity wave excited by orography. Not only are most grid cells of a GCM flat, but even where 85 there is topography, the waves themselves are highly intermittent. Here, we will focus 86 on non-orographic waves, but the method is general and an ad hoc version of it was used 87 by Sun et al. (2023) to emulate an orographic paramterization. More specifically, we build 88 on the work of Espinosa et al. (2022), who emulated a physics-based GW parameteri-89 zation (GWP) scheme (Alexander & Dunkerton, 1999) hereafter referred to as AD99, 90 with a deep neural network (DNN) architecture called WaveNet. We continue this in-91 vestigation to illustrate our approach for improving a generic ML methodology. Explor-92 ing our method in the context of emulation also allows us to explore the ability of a scheme 93 to generalize to different climates. 94

The strategy involves two distinct steps. First, one must identify the data imbalance. This requires "domain knowledge" of the problem, to identify key metric(s) that quantify rare cases where errors in the data-driven scheme limit its effectiveness. As detailed in Section 2, we establish a wind range metric to identify rare cases where WaveNet enmulator systematically fails. On top of being rare, these are cases where the physics of AD99 scheme become more non-local, and so more challenging to learn.

Once the data imbalance is identified, the second step is to treat it during model 101 training and implementation, as detailed in Section 3. We illustrate two strategies at the 102 learning stage, either to modify the sampling of training examples so that rarer cases are 103 better represented from the start, or to leave the distribution as it is, but adjust the loss 104 function to more strongly penalize mistakes on the rare cases. To construct a principled 105 method for this rebalancing, we borrow a concept from histogram equalization: a lin-106 ear interpolation of the original distribution to a more uniform distribution parameter-107 ized by a scalar t which can be varied from 0, where no change is made, to 1, where the 108 distribution is made completely uniform. The goal is to improve representation of the 109 rare cases without losing skill on the central part of the distribution or overfitting the 110 data in the tails, and the parameter t allows one to calibrate the degree of rebalancing. 111

These strategies assume that the ML model has enough complexity to learn the complex nonlinear behavior described by physics of g, but the data imbalance enables

the model to ignore rare samples and predominantly learn from the typical samples. As 114 we'll show in Section 4.2.1, overfitting can occur when the ML method is too complex 115 with respect to the amount of training data available. In addition to improving the train-116 ing of an ML scheme, one can mitigate data imbalance by applying a bias correction at 117 the inference stage. This involves computing the mean bias of the ML model as a func-118 tion of the relevant metric (the wind range in our case study of GWP emulation), and 119 subtracting the bias from the output. The remainder of the paper is structured as fol-120 lows. Section 2 illustrates how we identified data imbalance, Section 3 details modified 121 training and bias removal methods to overcome this imbalance. Our case study is pre-122 sented in Section 4. To demonstrate the generality of the method, we also introduce an 123 alternative ML strategy, an Encoder-Dense-Decorder (EDD). We use our approach to 124 improve both WaveNet and EDD. Furthermore, we illustrate how our approach can fail 125 when the complexity of the ML method exceeds the data available, leading to overfit-126 ting. Section 5 concludes our study and outlines possible future directions for this re-127 search. 128

¹²⁹ 2 Identifying data imbalance

A first step towards improving a data-driven parameterization – or more generally, 130 any data-driven task – is to identify potential imbalances in the training set. This pro-131 132 cess requires detailed knowledge of the application, as one is searching for metrics to quantify rare cases that are important for the performance of the task. The process is straight-133 forward in low dimensional data sets, i.e., if one needs to differentiate cats from dogs, 134 are the animals evenly distributed in the example data, but quickly becomes difficult in 135 high dimensional datasets. Here, we illustrate an example where the input data has 83 136 dimensions, but we seek one particular dimension that clearly identifies rare, but impor-137 tant, samples that need to be learned. 138

Our goal is to improve a data-driven emulator of the single column AD99 gravity wave parameterization, as implemented in the Model of an idealized Moist Atmosphere, MiMA (Garfinkel et al., 2020), following the work of Espinosa et al. (2022). We direct the reader to Alexander and Dunkerton (1999) for details on the parameterization and Espinosa et al. (2022) and Garfinkel et al. (2020) for details on the atmospheric model, but briefly review the most salient points here.

As in Espinosa et al. (2022), we use an integration of MiMA at triangular trunca-145 tion T42 resolution (corresponding to a $\approx 3^{\circ}$ grid) with model parameters configured 146 to produce a realistic representation of northern hemisphere climate by Garfinkel et al. 147 (2020). The model is integrated for 60 years, and after discarding the first 20 years' data 148 as spin-up, we use years 21-30 for the training and years 56-60 for the validation set. Out-149 put from the model is saved 4 times a day, yielding over 1.1×10^9 samples, where each 150 sample consists of vertical profiles of winds and temperature (the inputs), one for each 151 column on a 128×64 longitude-latitude grid, and the parameterized gravity wave ten-152 dency as the output. For simplicity, we focus only on the zonal (East-West) gravity wave 153 tendencies. 154

AD99 is a multi-wave GW parameterization that adheres closely to the scheme es-155 tablished by (Lindzen, 1981), which assumes the conservation of wave action flux and 156 wave-mean flow interactions under linear theory. The scheme determines GW momen-157 tum transport by launching a spectrum of non-interacting, monochromatic waves. Ther-158 modynamic breaking criteria determine when each wave breaks and deposits its momen-159 tum into the mean flow: waves tend to break when they near a critical level, where the 160 speed of the large scale winds equals that of the GW, or when their amplitude becomes 161 sufficiently large to overturn. This latter criteria is favored at upper levels where den-162 sity decays. Additional criteria account for waves that would be filtered out at the source 163 level (the nominal tropopause) or reflected downward. Important for our application, 164



Figure 1. Left: Two zonal wind profiles sampled near the South Pole at different times in the control integration; Middle: The physics based (AD99) computation of gravity wave momentum deposition (GWD) associated with these two profiles in the left panel; Right: The GWD output by the WaveNet emulator of AD99 for the same input profiles.

momentum carried by waves that do not break before reaching the model top are deposited in the upper levels of the column, thereby preventing a leak of momentum through
the model top (Shaw et al., 2009). A key simplification of the scheme is that the source
spectrum is only a function of latitude, meant to capture a simple background of waves
generated by convection, frontegenis, and orography.

Physical intuition can be garnered from Figure 1, which shows two example wind 170 profiles from an integration of the MiMA in the left panel, and the momentum tendency 171 computed by AD99 in the center. The scheme also uses the temperature profile (not shown) 172 to determine when convective overturning will lead to GW breaking, but winds are the 173 most important for prediction. The blue profile exhibits a more typical case; we will de-174 fine 'typical' precisely below. Critical line wave breaking leads to deposition of easterly 175 momentum in easterly shear zones, e.g., near 100 hPa, and conversely westerly momen-176 tum in westerly shear zones, e.g., near 1 hPa. The orange profile demonstrates a less typ-177 ical case with easterly flow in the troposphere below strong westerly shear throughout 178 the atmospheric column. Westerly waves are filtered out by easterly winds at the source 179 level (hence no westerly forcing), but the easterly half of the spectrum never experience 180 a critical level. The scheme thus deposits them all near the model top. 181

The right panel of Figure 1 provides anecdotal evidence that the WaveNet emulator does a reasonable job of capturing the momentum tendencies from the more typical blue profile case, but fails rather spectacularly with the orange profile. As detailed by Connelly and Gerber (n.d.), WaveNet is good at capturing critical level behavior, but struggles to capture non-local effects on the momentum tendencies, both the impact of source level filtering and integrated behavior, where an absence of easterly shear allows waves to reach the top.

We hypothesize that WaveNet's emulation of AD99 in MiMA suffers from data imbalance, in that gravity wave breaking is most often associated with local critical levels. WaveNet learns this relationship well. Cases where the momentum forcing depends on non-local behavior (e.g., when surface level filtering or low level critical levels remove

part of the spectrum low in the atmosphere, or when a lack of critical levels leads to mo-193 mentum deposition near the model top) are more seldom seen, and so tend to be poorly 194 captured the data-driven scheme. The challenge is to translate this physical intuition into 195 an objective metric to identify the rarer cases dominated by non-local effects. The in-196 put space is 83 dimensional (zonal wind \vec{u} and temperature \vec{T} at 40 levels each, plus sur-197 face pressure, latitude, and longitude), but we want a single metric to sort the data. Af-198 ter significant trial and error we developed a simple "wind range" metric that captures 199 many of these rare cases. 200

The wind shear is a crucial quantity in computing GW forcing on the mean flow. Large shear at any given level favors wave breaking, as GWs over a wider range of phase speeds will experience a critical level. Profiles with large shear, particularly at lower levels, tend to exhibit non-local behavior, as the GW spectrum is rapidly depleted, rending upper level critical levels moot. (This is to say, a second shear zone will not be associated with GW breaking because waves have already broken below.) In addition, strong shear in one direction can lead to cases like that exhibited in Figure 1, where the momentum conservation criterion leads to momentum tendencies near the model top, even if individual waves wouldn't otherwise break there. An admittedly crude proxy metric we consider to represent the overall presence of shear is the wind range, the total span of winds throughout the atmospheric column. Formally,

wind range =
$$\left(\max_{i=1,\cdots,\mathtt{nlev}} u_i\right) - \left(\min_{i=1,\cdots,\mathtt{nlev}} u_i\right).$$
 (1)

The wind metric is illustrated by the arrows in the left panel of Figure 1. It suggests that WaveNet may struggle when the wind range is large (the orange profile). While this metric was motivated by the physical argument that these high shear cases are more challenging to learn due to non-local effects, Figure 2 shows that these high wind range cases are rare as well.

The wind range exhibits the two key features of data imbalance. First, the input 206 data exhibits a long tailed distribution with respect to the wind range, and second the 207 ML based emulator systematically struggles with the tail of this distribution. This is most 208 clearly illustrated in Figure 2, which shows the distribution of errors for different val-209 ues of the wind shear. The spread of error increases superlinearly with respect to wind 210 range. For profiles with a wind spread of 50 m/s, at the mode of the distribution, the 211 error is the prediction of the drag is less than 5 m/s/day for over 90% of cases. For pro-212 files with range of 100 m/s, the error rates are only modestly worse, 85% of profiles ex-213 hibit an error less than 5 m/s/day. With a wind range of 150 or 200 m/s, however, only 214 70 and 30% of the profiles, respectively, can be predicted with an error of less than 5 m/s. 215 Error rates at the 90 percentile are associated with 16 and 28 m/s/day, respectively, a 216 full three to five times worse for cases at the mode of the distribution. 217

Figure 2 motivates another, even simpler approach of addressing data imbalance: 218 bias removal. The high absolute error rates for rare profiles with large wind range are 219 in part associated with systematic mean biases in the prediction (not shown). In gen-220 eral, a well-trained ML scheme will have no bias in the overall mean, but it can system-221 atically under and over-predict profiles with respect to metrics like the wind range. For 222 example, it may trivially under-predict the GW tendencies over the main part of the dis-223 tribution, but massively over-predict the tendencies at the tail. As discussed in Section 224 3.3 one can remove these biases at the time of inference. 225

For the remainder of the paper, we use the wind range metric, and the data imbalance it reveals, to improve the training and implementation of WaveNet and a related ML scheme. These methods are generic, and ready to apply once a user has identified the metric to quantify the imbalance. The better one can sort prediction errors in a high dimensional dataset along a single (or at least a small number of) dimension(s), however, the better one is positioned to use these strategies to improve the scheme.



Figure 2. Bottom panel shows the histogram of the dataset where each sample is represented by its zonal wind range Eq. (1). Frequency is the number of samples in a bin relative to the total number of samples. Top left: For each of the 100 equal-width bins of the histogram, we show 5th to 95th absolute error percentiles at 5-percentile increments. Thus we can view the error spread as a function of wind range. Due to noisy error statistics for samples with wind range >200 m/s, we exclude those samples in the analysis in the following sections. Top right: The error percentiles for a select few bins show that larger errors are incurred more often as the wind range increases.

²³² **3** Treating data imbalance

Our goal is to help the data driven scheme perform better on the tails of the distribution *without* decreasing performance over the main part of the distribution. This makes the typical balancing act between "bias" and "variance" that one seeks with any machine learning task more challenging. Good performance requires a scheme that both learns the training data well (has low bias) and works equally well on new data (has low variance). By this, we mean that the skill is uniform for different samples from the underlying distribution, so it generalizes well to new inputs it has not seen before.

A large bias is associated with under-fitting, where the method lacks enough train-240 ing data and/or expressivity to capture the relationships, while a large variance is as-241 sociated with over-fitting, where the ML scheme uses "noise" (unimportant features) in 242 the training data to reduce the bias. This is a case of having too much expressivity rel-243 ative to the amount of data. The expressivity of a ML scheme is related to its complex-244 ity (roughly, the flexibility it has to identify relationships between inputs and outputs, 245 which is a function of both the method and the number of free parameters it is given). 246 For our application, we are given some ML scheme of fixed complexity (i.e., WaveNet). 247 We must ensure there is still enough training data in the center of the distribution to 248 avoid under-fitting it, and not too much emphasis on the tails to cause over-fitting. 249

Learning from unbalanced datasets is challenging. For example, consider a dataset where 99% of the dataset is class A and the remaining 1% is class B. A binary classifier that always predicts class A can still be considered very good under a seemingly innocent metric such as average accuracy, defined as

average accuracy $\equiv \frac{\text{#correctly labeled samples}}{\text{#of total samples}}$,

with a value of 0.99, although it completely fails to learn the characteristics of class B.
Methods to remedy difficulties attributed to imbalanced datasets for classification are
far and plenty (He & Garcia, 2009; Johnson & Khoshgoftaar, 2019), and are used in a
variety of applications including object detection (Oksuz et al., 2021).

These methods can be broadly categorized into data-level, algorithm-level, and the 254 hybrid of those two. Data-level methods manipulate the distribution of the training data 255 distribution: such as undersampling from the majority class and oversampling from the 256 minority class (Chawla et al., 2004), or generating synthetic samples of the minority class 257 (Chawla et al., 2002) through randomly weighted linear combinations of samples. Algorithm-258 level methods adjust the learning algorithm to increase/decrease the impact of samples 259 from minority/majority class. The latter case falls under cost-sensitive learning as it is 260 implemented by imbuing a cost or penalty term in the learning process (Krawczyk, 2016; 261 Elkan, 2001). 262

Although many methods for treating data imbalance are established for classifi-263 cation, extending them for regression is nontrivial. There have been some efforts on this 264 front as done by Torgo et al. (2015); Ding et al. (2019); and Rudy and Sapsis (2023). Torgo 265 et al. (2015) extends the Synthetic Minority Oversampling TEchnique (SMOTE; (Chawla 266 et al., 2002)) to regression by assuming near linearity of the model being learned, Rudy 267 and Sapsis (2023) extends relative entropy based loss functions from scalar outputs to 268 low dimensional vector outputs, and Ding et al. (2019) proposes a new loss function and 269 a model design that memorizes extreme events for time series applications. Some short-270 comings of these solutions are that they are incompatible with nonlinear problems and 271 difficult to implement in applications with high dimensional datasets. 272

We prepare two methods to address data imbalance in regression tasks. Both meth-273 ods require first identifying a metric along which the high-dimensional dataset yields a 274 long-tailed distribution; in our case, the wind range. We project our high-dimensional 275 dataset to the low-dimensional space identified by the metric. Section 3.1 shows how his-276 togram equalization can be applied to transform unbalanced distribution to one more 277 uniform. This idea is closely related to transportation theory (optimal transport), which 278 is the study of allocation of resources with a constraint of cost appended to the trans-279 portation of those resources. Since we merely intend to modify the data distribution en-280 countered by the training algorithm, rather than to transform the data itself, we drop 281 the transportation cost constraint. In Section 3.2, we describe the data rebalance method, 282 which extends the ideas of over/undersampling methods to treating data imbalance for 283 regression tasks by applying linear transformations to the probability distribution func-284 tion (PDF) of the dataset. Finally, we describe mean bias removal in Section 3.3. 285

3.1 Histogram equalization

Histogram equalization is an image processing method that adjusts the contrast of an image by changing the shape of the histogram of the intensity values, and is the simplest optimal transport method for 1D data. The extent to which the shape of the histogram is modified is parameterized by $t \in [0, 1]$ where t = 0 yields the original histogram, and t = 1 a target histogram. By equalization, we aim for a target distribution that is uniform, with an equal number of pixels in each intensity bin. Figure 3 shows an example of this applied to a grayscale image where each sample has a value in [0, 1] which represents a greyscale value between black and white. The original histogram (t = 0) has the majority of pixels in the moderate intensity region, and very few pixels are close to minimum and maximum intensities. As the parameter t increases to 1, the distribution is flattened in the peak region and elevated in the extreme regions. Lighter pixels are made lighter and darker pixels are made darker, qualitatively yielding images with greater contrast as t increases.



Figure 3. An example of histogram equalization performed for image processing with t ranging from 0 to 1. The original image corresponds to t=0. As t increases, moderate saturation pixels are pushed towards their nearest extremes. At t=1, the pixels are distributed almost uniformly.

Let us describe this procedure in more detail. Let x_i denote the intensity of the *i*th pixel of an $m \times m$ image, and let permutation σ be defined such that $\{x_{\sigma(j)}\}_{j=1}^{m^2}$ are sorted in increasing order,

$$x_{\sigma(1)} \leq \ldots \leq x_{\sigma(m^2)}.$$

Assign $\{y_j\}_{j=1}^{m^2}$ to the cumulative distribution function (CDF) of the target distribution. This corresponds to m^2 equispaced, ordered nodes from 0 to 1 since the target is the uniform distribution for histogram equalization:

$$y_j = (j-1)/(m^2-1), \ j = 1, \cdots, m^2.$$

In general, the CDF of any desired target distribution suffices as the values of y_j 's. Then, the new intensity value for the i^{th} node is given by

$$z_i := (1-t)x_i + ty_{\sigma^{-1}(i)}.$$
(2)

Here is a numerical example of applying this to a 2×2 image. The original image is given by pixels

$$\begin{bmatrix} x_1 & x_2 \\ x_3 & x_4 \end{bmatrix} = \begin{bmatrix} 0.60 & 0.52 \\ 0.25 & 0.44 \end{bmatrix}$$

The sorting permutation is $\sigma = [3, 4, 2, 1]$ for a row-wise uncoiling of the matrix, and the target values are $y_1 = 0$, $y_2 = 1/3$, $y_3 = 2/3$, $y_4 = 1$. Thus, the transformation yields

$$\begin{bmatrix} y_{\sigma^{-1}(1)=4} & y_{\sigma^{-1}(2)=3} \\ y_{\sigma^{-1}(3)=1} & y_{\sigma^{-1}(4)=2} \end{bmatrix} = \begin{bmatrix} 1 & 2/3 \\ 0 & 1/3 \end{bmatrix}$$

for t = 1, and the general formula for any $t \in [0, 1]$ is given by

$$(1-t)\begin{bmatrix} x_1 & x_2\\ x_3 & x_4 \end{bmatrix} + t \begin{bmatrix} y_4 & y_3\\ y_1 & y_2 \end{bmatrix} = (1-t)\begin{bmatrix} 0.60 & 0.52\\ 0.25 & 0.44 \end{bmatrix} + t \begin{bmatrix} 1 & 2/3\\ 0 & 1/3 \end{bmatrix}.$$

301 3.2 Data Rebalancing

Our goal is to change the distribution of the training dataset while taking full use 302 of the available data and without generating synthetic data. Histogram equalization for 303 image processing achieves the reshaping of the dataset distribution by transforming the 304 values of the sample from x_i to z_i as shown in Eq. (2). Doing so may move a sample from 305 one histogram bin to another, thereby changing the histogram directly. Our method uses 306 the linear mapping from the original to the new intensity values described in Eq. (2), but 307 apply the mapping to the PDF instead. The newly assigned probability may increase 308 or decrease a sample's contribution to the training process. We describe the method in 309 detail below, and propose two implementations of the method in Sections 3.2.1 and 3.2.2, 310 respectively. 311

Let $H^{(0)}$ be the histogram of the training dataset X_{training} with N bins,

$$\{[b_0, b_1), \ldots, [b_{N-1}, b_N]\}.$$

The count of samples in the *n*th bin, $[b_{n-1}, b_n)$ is $h_n^{(0)}$, and the ideal count of the samples in the *n*th bin in the ideal histogram is $h_n^{(1)}$. Here, the ideal histogram is uniform with N equal width bins, so $h_n^{(f)} = M/N$ for all $n = 1, \dots, N$ for a dataset with M samples. The new count of the *n*th bin for parameter t is then:

$$h_n^{(t)} = (1-t)h_n^{(0)} + th_n^{(1)}.$$
(3)

Since the *n*th bin originally represented $h_n^{(0)}/M$ of the training set and now we want it to represent $h_n^{(t)}/M$ of the training set, the ratio between the two determines the resampling rate in the *n*th bin.

$$\alpha_n^{(t)} := \begin{cases} h_n^{(t)} / h_n^{(0)} = (1-t) + t h_n^{(f)} / h_n^{(0)} &, h_n^{(0)} > 0\\ 0 &, h_n^{(0)} = 0 \end{cases}$$
(4)

These ratios determine the new sampling rates for the training data. We found in practice that fairly low *t*-values still yielded very large α ratios at bins belonging to the extreme tail of the distribution. To avoid unreasonable resampling rates being assigned to rare data points, we bound the ratios by the maximum repeat parameter as shown in Eq. (5),

$$\tilde{\alpha}_{n}^{(t)} := \begin{cases} \min\{\alpha_{n}^{(t)}, \max_\texttt{repeat}\} &, \ h_{n}^{(0)} > 0\\ 0 &, \ h_{n}^{(0)} = 0. \end{cases}$$
(5)

Thus, the final resampling rate, $\tilde{\alpha}_n^t$, is determined by three decisions: 1) choice of histogram bins; 2) t, the linear mapping parameter; and 3) the maximum repeat parameter. The resampling strategy is no longer a simple bilinear interpolation between the original $(h^{(0)})$ and desired $(h^{(1)})$ histograms due the maximum value of the resampling rate. The counts for the bins of the new, resampled histogram for some $t \in [0, 1]$ and max_repeat is,

$$\tilde{h}_{n}^{(t)} = \tilde{\alpha}_{n}^{(t)} h_{n}^{(0)}.$$
(6)

The process is easier to visualize than spell out: Figure 4 shows the original his-312 togram, $h^{(0)}$, plotted in foreground with the new histograms, $\tilde{h}^{(t)}$, with increasing val-313 ues of t for each panel, as well as three different values for max_repeat in each panel. The 314 impact of max_repeat is seen most clearly in the bottom three panels. The zonal wind 315 range at which the lower values of max_repeat diverge from the highest value is depen-316 dent on t as expected. The number of histogram bins was kept constant here. It gov-317 erns how finely one resolves the distribution. One could also allow the width of the bins 318 to vary, say to more ifnely capture the center vs. the tails. 319



Figure 4. Each of the panels correspond to *t* values ranging from 0.05 to 0.60. The 3 lines for each panel represent the impact of maxrepeat parameter values 10, 100, and 500. The original histogram is shown filled in as a basis for comparison.

320

3.2.1 Implementation I: Direct sampling

State-of-the-art optimization methods for deep neural networks rely on incremen-321 tal, iterative updates of the model weights. They are incremental in that each update 322 is based only a subset of the training dataset called a *batch*, and iterative in that the train-323 ing dataset is passed through the optimization method many times before the model weights 324 converge to an acceptably optimal state. An *epoch* is a measure of unit for the progress 325 of the training of a model defined by a single pass over the training dataset, for which 326 each sample in the training dataset processed exactly once. Since our strategy changes 327 the contribution of each sample to the training algorithm based on where in the data dis-328 tribution the sample belongs, some samples will be seen more often than others. There-329 fore, we modify the definition of an epoch to mean a single-pass over a resampled sub-330 set of the dataset. We outline the procedure for resampling in context of a general NN 331 training algorithm, which is written as a pseudoalgorithm (Algorithm 1) in Appendix 332 А. 333

First, compute resampling rates $\tilde{\alpha}_n^{(t)}$ for each bin using Eq. (5). Next, for each bin labelled by $n = 1, \dots, N$, resample and collect the indices of the chosen samples. If $\tilde{\alpha}_n^{(t)} <$ 1, then it is straightforward to sample from the *n*th bin with probability $\tilde{\alpha}_n^{(t)}$ by randomly choosing a subset of the bin of size $\tilde{h}_n^{(t)}$ without replacement. Another method is to sample from the uniform distribution $h_n^{(0)}$ times and keep the indices that correspond to sampled values less than $\tilde{\alpha}_n^{(t)}$. For both methods, the selected indices are recorded. On the other hand, if $\tilde{\alpha}_n^{(t)} > 1$, then include every sample from this bin floor($\tilde{\alpha}_n^{(t)}$) times, and then sample with probability $\tilde{\alpha}_n^{(t)} - \text{floor}(\tilde{\alpha}_n^{(t)})$. Following good practice, the collected indices from all N bins should be combined, shuffled, and separated into batches. These batches should then be fed to the training algorithm, which will update the NN model weights once for each batch.

Once all of the batches are processed and if further training is needed, repeat the 345 resampling step to select another realization of the new data distribution. Note that that 346 every iteration of resampling is done without replacement, but samples may be repeated 347 from one iteration to the next. It is straightforward to include an additional step to re-348 sample at the next iteration without replacement by keeping track of which samples and 349 how many times those had been picked in previous iterations. When sampling without 350 replacement is implemented across epochs, all of the samples to be seen by the training 351 algorithm at least once after ceiling $\left(\left(\min_{n} \tilde{\alpha}_{n}^{(t)}\right)^{-1}\right)$ epochs. We include a pseudoal-352 gorithm for the resampling method in Algorithm 2 in Appendix A. 353

354

3.2.2 Implementation II: Weighted Loss Function

An alternative implementation of our approach is to modify the loss function to 355 account for disparity in the distribution. Success in training deep NNs are attributed to 356 efficient back-propagation, a method of updating model weights with the goal of min-357 imizing a loss computed from a batch of samples. Since loss functions are typically de-358 fined for a single pair of the target and the predicted value, the loss over a batch of sam-359 ples is an average of the loss function values for each of the samples in that batch. This 360 implies that every sample in the batch has equal importance in updating the model weights. 361 Our resampling strategy aims to modify the data distribution to lend importance to some 362 samples and reduce impact from other samples. We propose using a weighted average 363 in the accumulation of loss function values of a batch, where the weight for each sam-364 ple corresponds to the resampling rate of the bin the sample belongs to. For a sample 365 indexed by i that belongs to bin n, the weight is determined by parameters t and \max -366 **repeat** via Eq. (5): $w_i \equiv \tilde{\alpha}_n^{(t)}$. The weights can be computed for the entire training dataset 367 prior to any training and passed to the training loop to compute a weighted average of 368 the loss function for each batch, as shown in Section 3.2.2. 369

$$\operatorname{Loss}_{\operatorname{avg}}(\{y_i\}_{i=1}^{\operatorname{batch size}}, \{\hat{y}_i\}_{i=1}^{\operatorname{batch size}}) = \frac{1}{\operatorname{batch size}} \sum_{i=1}^{\operatorname{batch size}} \operatorname{Loss}(y_i, \hat{y}_i)$$
(7)

$$\text{Loss}_{\text{weighted avg}}(\{y_i\}_{i=1}^{\text{batch size}}, \{\hat{y}_i\}_{i=1}^{\text{batch size}}) = \frac{1}{\text{batch size}} \sum_{i=1}^{\text{batch size}} w_i \text{Loss}(y_i, \hat{y}_i).$$
(8)

370

375

3.2.3 Maximum repeat: Fail-safe against overfitting

The maximum repeat parameter, Eq. (5), puts a threshold on the oversampling rate to prevent overfitting. This allows us to fine tune treatment of the data imbalance by relaxing the computed resampling rates of bins with high α ratios, which typically occur at the the extreme tail of the distribution.

3.3 Bias removal

In addition to the resampling method, we propose a correction method to be employed at time of inference to further enhance the quality of the ML model. This tactic applies a first-order correction to remedy the bias of a trained model, where the bias is computed along the metric used to identify the data imbalance. There are a couple of ways to compute the bias. Consider a dataset of M samples that were binned into Nbins where \mathcal{B}_n is the set of indices of samples that belong to the *n*th bin. The output variable has dimension d, and we denote the target and predicted variable of the *i*th sample by

$$\vec{y_i} = \begin{bmatrix} y_{i,1} \\ \vdots \\ y_{i,k} \end{bmatrix}, \vec{\hat{y}_i} = \begin{bmatrix} \hat{y}_{i,1} \\ \vdots \\ \hat{y}_{i,k} \end{bmatrix}$$

where $\hat{\cdot}$ is used to denote the ML predictions. The mean error profile for the entire dataset can be computed by

mean error profile =
$$M^{-1} \sum_{i=1}^{M} \vec{y}_i - \hat{\vec{y}}_i$$

For a well trained scheme, the mean error profile should be close to a vector of zeros. Similarly, we can compute the mean error profile can be computed for each bin,

mean error profile for bin
$$n = \left\{ |\mathcal{B}_n|^{-1} \sum_{i \in \mathcal{B}_n} \vec{y}_i - \hat{y}_i \right\}_{n=1}^N.$$
 (9)

Large errors in bins of the tails can be balanced by smaller errors in the fat pail of the distribution. At inference, we simply determine the bin the sample belongs to and subtract the appropriate mean bias profile.

³⁷⁹ 4 Case study: Data-driven GWP emulation

Section 4.1 describes two model architectures we use to test our method: WaveNet from Espinosa et al. (2022) and a convolutional NN encoder-dense-decoder (EDD). Both implementations of the data rebalancing, with varying t parameters, are applied during training on the same MiMA dataset. Offline results are presented in Section 4.2, and the emulators with the best offline results are tested online in Section 4.3. Here, online refers to replacing AD99 within MiMA integrations with our trained ML emulators.

4.1 Model Architectures

386

We include a short summary of WaveNet here, and refer readers to Espinosa et al. (2022) for a full description. WaveNet takes in a concatenation of all of the input variables and applies several dense layers that split into pressure level-specific "branches". The branches themselves are also dense layers that output GWD values for a specific pressure level of the MiMA vertical grid, and do not communicate with one another.

The EDD architecture uses 1D convolutional layers in the encoder and decoder sec-392 tions and dense layers in the middle section. This structure is imposed to encourage the 393 model to learn local interactions in the encoder section via convolutions while downsam-394 pling layers compress the outputs. This combination of convolutional layers followed by 395 downsampling is commonly used in autoencoders, which can serve as a nonlinear dimen-396 sion reduction technique that extract essential information. The middle dense section 397 allows the processing of global relations and the decoder section reassembles the verti-398 cal profile of the zonal gravity wave drag with transposed convolutions and upsampling. 399 Additional details are included in Appendix B. 400

The hyperparameters for these architectures, listed in Table 1, include the num-401 ber and width of the dense layers, the number of (transposed) convolution layers and the 402 size and number of filters for each of these (transposed) convolution layers. Some degrees 403 of freedom were removed by restricting the encoder and decoder halves to be as sym-404 metric as possible, while accounting for the fact that the encoder receives multiple chan-405 nels and the decoder outputs a single channel. For the remaining degrees of freedom, we 406 used RayTune (see Liaw et al. (2018)) to thoroughly tune the hyperparameters. We con-407 trast two sizes for each architecture: a smaller network of approximately 350,000 pa-408 rameters; and a larger network of approximately 700,000 parameters. Espinosa et al. 409 (2022) found that large networks yielded better offline skill than their smaller counter-410 parts, but at the expense of additional computational costs. We present two metrics that 411 are closely related to the mean squared error (MSE), the loss function used during train-412 ing. 413

 Table 1.
 Number of trainable parameters in section of each model architecture. The EDD is comprised of 3 sections: encoder, dense, decoder; WaveNet is comprised of 2 sections: shared layers and 33 branches for the top 33 pressure levels.

| Model Type/Size | Convolutional Layers | Dense Layers | # Layers per section |
|------------------------|----------------------|----------------------|----------------------|
| Small EDD Large EDD | $26,237 \\ 50,337$ | $328,800 \\ 650,800$ | $3/3 \\ 3/3$ |
| Model Type/Size | Shared Layers | Branched Layers | # Layers per section |
| Small WaveNet | 10,368 | 342,177 | 1/3 |
| Large WaveNet | 14,904 | 704,385 | 1/3 |

The absolute norm error (AE) is defined as,

absolute norm $\operatorname{error}(y, \hat{y}) = \|y - \hat{y}\|_2$,

and was shown, for instance in Fig. 2. We also consider the relative norm error (RE) expressed as,

relative norm
$$\operatorname{error}(y, \hat{y}) = \frac{\|y - \hat{y}\|_2}{\|y\|_2}.$$

The relative norm error scales the norm of the error by the magnitude of the target vector, and ensures that the trend of the error norms are not simply proportional to the trend of the target vector norms.



Figure 5. Baseline (t=0) absolute and relative error norms of two sizes of WaveNet and EDD are shown. The errors are shown as a function of wind range in the validation set, as in Figure 2, which showed results only from the large large WaveNet model.

416

Figure 5 shows the absolute and relative errors of the four models with no resampling strategy; this establishes a baseline for comparison with our resampling strategies. The validation set (data not observed in training) errors are averaged for each bin of the zonal wind range. We have dropped data points whose zonal wind range are greater than 421 200 m/s, as the errors here are too noisy for robust analysis. We show the relative er-422 ror on the right panel to highlight how all four variants learn the peak portion of the dis-423 tribution (refer to the histogram in Fig. 2) best, but fail at the tail (>125m/s). A rel-424 ative norm error of 1 (100% relative error) implies that the magnitude of the error is as 425 large as the target profile itself: a scheme predicting zero drag all the time would sat-426 isfy this condition. This suggests that the schemes are doing a pretty awful job for wind 427 range above 125 m/s; predicting nothing at all would be more accurate.

We observe that the EDD models outperform the WaveNet models, albeit the er-428 rors are of similar magnitudes. Overall, the disparity in errors is more significant between 429 the model architectures than between network sizes. Despite having approximately the 430 same number of learnable parameters as their EDD counterparts, the WaveNet models 431 have not acquired as much skill given identical training conditions; the number of learn-432 able parameters is not all in all when it comes to model complexity. The larger variants 433 of both architectures yield smaller average absolute errors than the smaller variants. The 434 disparity grows slightly for larger zonal wind ranges, though this slight lead of the larger 435 models falters for zonal wind ranges greater than ≈ 125 m/s for the relative error. 436

437

4.2 Data Resampling and Offline Results

⁴³⁸ Of the three tunable parameters of the resampling strategy, we study the impact ⁴³⁹ of tuning t. The maximum repeat parameter and resolution of the histogram were set ⁴⁴⁰ at 100 maximum repeats and 100 equal-width bins after an initial survey. We investi-⁴⁴¹ gated values of t = 0.05, 0.10, 0.15, 0.20, 0.40, 0.60 following intuition that t closer to ⁴⁴² 1 is likely more damaging than helpful given the shape of the distribution of our dataset. ⁴⁴³ Figure 4 shows the new shape of the data distribution of the 6 configurations on the teal ⁴⁴⁴ (medium-width) lines with the original distribution shaded in green in the background.

Figures 6 to 9 show the baseline error (t = 0, shown in Fig. 5) in black lines, and 445 the deviation of the error relative to this baseline for t > 0 in colors ranging from brown 446 to yellow. In all instances we see very little, if any, loss of accuracy in the peak region 447 (a wind range of roughly 10 to 100 m/s). We have achieved one criterion for success: re-448 sampling, either directly or through a weighted loss function, does not damage perfor-449 mance for typical inputs. Now the harder part: does resampling improve performance 450 in the tail, from 100 to 200 m/s in our wind metric? Here we found success in most cases, 451 though not uniformly. We acknowledge the failure first. In our best baseline network, 452 the large EDD, direct oversampling led to overfitting. In all other cases, however, we were 453 able to successfully reduce error in the tail. 454

455

4.2.1 Overfitting vs Underfitting

As we feared, the resampling strategy can encourage overfitting of the tail in a data 456 driven scheme with sufficient complexity. Figure 6 shows the result of training the large 457 EDD model. The left panel shows the direct sampling implementation (Algorithm 2). 458 For the direct sampling implementation, samples with wind range greater than 125 m/s 459 in the training set suggest impressive gains when compared to the baseline error, albeit 460 with no clear correlation with the t parameter. This improvement, however, fails to gen-461 eralize to samples unseen during training: the mean absolute error of the validation set 462 is larger than that of the baseline error. We observe that larger t corresponds to larger 463 growth in error, suggesting that the trained models suffer from overfitting triggered by 464 the inflation of samples in the moderate tail region. 465

Typically, overfitting is diagnosed during training when validation error stops improving (or even start to get worse) while training error further improves. While we only show the errors at the end of training, it is clear from the design of this experiment that the resampling strategy resulted in models that learned the noise at the tail rather than



Figure 6. This figure show errors from the the baseline (t=0, black) model and the errors of models trained with the sampling strategy (t>0, colors brighten as t increases) implemented with direct sampling on the large EDD model architecture. The top and middle rows show the errors on the training and validation sets, and the bottom row shows the histogram of the dataset with respect to the zonal wind range between 0 and 200 m/s.

learning an intrinsic principle tied to the tail. A potential cause for overfitting is larger
model complexity (number of trainable parameters) relative to the complexity of the pattern being learned, which then leads to the model learning noise associated with the specific instance of the training set. We suspect that oversampling of the tail combined with
the large network size created a learning environment in which the EDD had the capacity to learn noise in the tail.

The right panel of Fig. 6 shows the experiment results with the weighted loss im-476 plementation (Algorithm 3). Unlike the direct sampling implementation, we observe about 477 the same magnitude of improvement in the tail for the training set and the validation 478 set. The upper-middle range t-values (0.15, 0.20, 0.40) exhibit no improvements from 479 the baseline in the validation set, but the extreme t-values (0.05, 0.10, 0.60) all show slight 480 improvements. Since the only difference between the left and the right panels is in the 481 implementation details of the resampling strategy, this suggests that the weighted loss 482 function implementation may be less amenable to overfitting than the direct sampling 483 method. We further analyze the comparison between the two implementation methods 484 in Section 4.2.2. 485

Next, we repeat the experiment in the previous section for the large WaveNet ar-486 chitecture, which has a comparable number of tunable parameters for both implemen-487 tation methods, and show the result in Fig. 7. We observe that the validation set errors 488 at the tail are smaller than the baseline error for most t values, and there is no signif-489 icant change to the errors at the peak. Unlike the example in Fig. 6, these large networks 490 did not overfit to the samples at the tail of the training set relative to the baseline er-491 ror. If network size is a potential cause for overfitting in the direct sampling large EDD 492 case, why do we not see similar results in the large WaveNet cases? We speculate that 493 the baseline WaveNet model was underfitting and there was more room for improvement 494



Figure 7. Both columns show errors in the same fashion as Fig. 6. Left column shows errors for the large WaveNet instances with direct sampling implementation, and right column shows errors for the large WaveNet instances with weighted loss function sampling implementation.

to be garnered from applying the resampling strategy. If the baseline EDD model was
not underfitting, then the resampling strategy could not reduce the approximation error (bias) much more than was already achieved by the baseline model, and all there was
left to learn were noisy traits unique to the training set.

With the exception of the overfitting case, the resampling strategy successfully reduces underfitting at the tail without penalty in the peak, thereby reducing the bias overall. In the next section, we show further evidence of success of the resampling strategy and compare the two implementation methods.

503

4.2.2 Sampling strategy comparison: weighted sampling vs weighted loss

We now compare the two implementations (Algs. 2 and 3) on the small EDD mod-504 els. Figure 8 shows the baseline errors and the deviations from the baseline errors as we 505 vary t over the training and the validation sets. Figure 8 reveals improvements in the 506 tail, albeit modest, with little to no damage in the peak. The notable exceptions occur 507 at t = 0.20 and t = 0.40 for the weighted loss implementation, where there are almost 508 no change if not a decline in performance on the tail. These occur in both the training 509 and validation set, however, and therefore are not likely an issue of overfitting. Outside 510 of those exceptions, improvements occur for a wider range of the distribution, with larger 511 magnitudes of improvement in the training set than in the validation set as expected. 512 The weighted loss experiment (right plots of Fig. 8) shows a slightly larger disparity be-513 tween the training and validation set errors than the direct sampling experiment; the train-514 ing set errors show larger improvements with the weighted loss implementation than di-515 rect sampling, but the validation errors are comparable between the two implementa-516 tions. With direct sampling, all t values except for t = 0.60 still yield improvement in 517 error in the moderate tail region. 518

Next, we discuss the experiment results for the small WaveNet model. As shown in Fig. 9, the difference between direct sampling and weighted loss are less pronounced



Figure 8. Both columns show errors in the same fashion as Fig. 6. Left column shows errors for the small EDD instances with direct sampling implementation, and right column shows errors for the small EDD instances with weighted loss function sampling implementation.

than in the EDD model. Also, the errors of the training set and the validation set are much closer than in the experiments for the small EDD models. The largest difference between the implementation methods for the small WaveNet models is in which t values are the most optimal. The direct sampling method is optimized for the smallest and largest t values, whereas the weighted loss method prefers moderate t values ($t \approx 0.15$).

Even though we saw that the loss function sampling avoided overfitting for the large 527 EDD experiment, we do not see a similar advantage of the loss function implementation 528 over the direct sampling implementation in the small EDD, small WaveNet, and large 529 WaveNet experiments. However, we do see modest improvements in the tail for mod-530 els trained with the resampling strategy for the majority of t values for those three ex-531 periments, although there is no clear trend of which t values are optimal. Future exper-532 iments that may reveal tighter trends, include studying the sensitivity of learning algo-533 rithm, and increasing the density of t values. 534

535

4.3 Bias Removal and Online Results

We conclude our case study with a brief discussion of how our modified data-driven 536 parameterizations perform when coupled "online" with the MiMA atmospheric model. 537 An important evaluation of a new parameterization scheme is conducted by computing 538 statistics from long-time integrations where the scheme is coupled with the model, as op-539 posed to the "offline" metrics we showed in the previous section. Online coupling is a 540 more challenging task, as errors in the GWP can lead to biases in the large scale flow, 541 forcing the scheme to make inferences in regimes it has not yet seen, which often leads 542 to instability (Brenowitz et al., 2020). 543

To test a selection of our trained ML emulators, we follow Espinosa et al. (2022), coupling them with MiMA for 40-year integrations after 20 years of model spin-up. The simulations with the data-driven emulators can then be compared against the control



Figure 9. Both columns show errors in the same fashion as Figs. 6 to 8. Left column shows errors for the small WaveNet instances with direct sampling implementation, and right column shows errors for the small WaveNet instances with weighted loss function sampling implementation.

integration with the "true" gravity wave forcing provided by the AD99 physics based pa-547 rameterization. Coupling also allowed us to implement the bias correction, which can 548 be implemented independently or in addition to the rebalancing strategies. To summa-549 rize quickly, the new data driven parameterizations successfully couple with the model, 550 producing climatological statistics (mean and variability) that were consistent with the 551 original model. Differences between the model with the baseline schemes and our re-balanced 552 versions, however were not statistically significant. It is likely that a longer integration 553 could eventually reveal significant differences, but an improvement that requires a cen-554 tury or more to observe is of modest utility. We conclude that while re-balancing the data 555 did improve performance based on the wind metric, this bias was either not critical to 556 performance of the parameterization in the model, or we have not sufficiently improved 557 the tails to see a significant effect. 558

For completeness, we show a few results here, focusing on the coupled model's abil-559 ity to generate the Quasi-Biennial Oscillation (QBO), a vacillation of easterly and west-560 erly jets in the tropical stratosphere over a period of approximately 28 months. We high-561 light this metric because the QBO is in large part driven by gravity wave momentum 562 transport. This emergent behavior on a time scale of years, generated from gravity waves 563 that operate on time scales of hours, is viewed as critical test of gravity wave parame-564 terizations (Richter et al., 2022; Anstey et al., 2022; Bushell et al., 2022). An important 565 difference between the online runs in this manuscript and that of (Espinosa et al., 2022) 566 is in the model parameters of MiMA that generated the training data. We employed pa-567 rameters that were optimized for simulation of the Northern hemisphere (Garfinkel et 568 al., 2020), not the QBO. Thus the oscillation is the control integration had a period of approximately 35 months, not 28 months, as shown in Figure 10. Capturing the right 570 period of the QBO is generally achieved by tuning the GWP, as was done in Garfinkel 571 et al. (2022). 572

| Emulator Description | Transition Time |
|---|------------------|
| control | 35.01 ± 2.46 |
| small EDD, $t = 0$ | 38.37 ± 6.59 |
| small EDD, $t = 0$, bias removed | 37.01 ± 7.70 |
| small EDD, $t = 0.05$, direct sampling | 37.05 ± 2.98 |
| small EDD, $t = 0.05$, direct sampling, bias removed | 38.12 ± 3.69 |
| small EDD, $t = 0.05$, weighted loss | 39.66 ± 8.97 |
| small EDD, $t = 0.05$, weighted loss, bias removed | 36.42 ± 5.47 |

Table 2.

We show results with the smaller EDD models, as the rebalancing strategies ex-573 hibited the largest offline improvement. Table 2 lists the QBO period for the baseline 574 575 model (t = 0) and the various combination of resampling strategy and bias removal. The QBO period was computed using the Transition Time (TT) method of Richter et 576 al. (2020). First, the zonal wind was averaged zonally in the tropical region (latitudes 577 between $5^{\circ}S$ and $5^{\circ}N$), as shown in Figure 10. Then the intervals between QBO phase 578 changes are defined as times when the signs of zonal mean zonal wind reversal near 10 579 hPa (denoted by the plus signs). The resulting mean and the standard error of those val-580 ues give us a proxy for a confidence interval of the QBO period. A robust implementa-581 tion of the TT method requires smoothing the field with 15 to 30 day windows, to avoid 582 double counting small deviations around transitions. 583

The baseline model exhibited a slightly longer QBO period of 38 months, though 584 40 years of simulation was insufficient to establish whether this bias is statistically sig-585 nificant. We found that all of our modified data-driven approaches exhibited shorter QBO 586 periods, an improvement relative to the baseline, but still biased long relative to the con-587 trol. The best performing model is highlighted in Figure 10, but as quantified in Table 2, 588 these integrations are not long enough to establish whether these differences are statis-589 tically significant. As noted above, this could be due to the fact that the QBO bias is 590 unrelated to errors in the rare cases highlighted by the wind metric, or that our correc-591 tion is insufficiently large to make a dent. It highlights the importance of domain knowl-592 edge to identify the key metric(s) of data imbalance that matter for the problem of in-593 terest. 594

595 5 Conclusions and Future Directions

With the growing prevalence data-driven methods being used for various tasks in 596 modeling earth system models, it is crucial to properly learn from geoscience datasets. 597 We address what one can do to improve a data-driven parameterization given that there 598 is no additional data to learn from, nor computational capacity to allow for a larger, more 599 complex model. In other words, this is the typical scenario for modeling various subgrid-600 scale mechanisms in climate models. In particular, we proposed two strategies to com-601 bat data imbalance with the goal of improving data-driven models, and applied it to a 602 case study of improving a data-driven GWP model. 603

Both methods rely on first identifying a metric or a projection that yields reveal an imbalance in the available dataset that has an inherent significance to the physical process being modelled. This process is unique to each application and requires expert scientific knowledge of the modelled process, and doubles as a dimension reduction step that allows the practitioners to view the original high-dimensional dataset in a new context. Ideally, this new context should illuminate the differences between frequent (and therefore easy to model) instances from rare (and difficult to model) instances. Despite



Figure 10. Both plots show the zonal mean zonal wind averaged over years 20-40 in latitudes between 5°S and 5°N. The crosses indicate the times where QBO phase changes are detected by the TT method. Left: Control run with AD99; Right: Best emulator (small EDD with resampling strategy via weighted loss and t = 0.05 and bias removal).

resulting from the same physical mechanisms, these two types of instances occupy al-611 most two distinct regimes due to the natural variability in the model system. A neces-612 sary complicating factor is that these two *classes* are not sharply partitioned like dis-613 crete distributions, but rather can be viewed as the peak and the tail of a continuous dis-614 tribution. In our case study, we chose wind range of a model column as the appropri-615 ate metric for our physical process, gravity waves. This choice stemmed from the obser-616 vation that wind range can crudlely approximate shear, an important quantity in deter-617 mining the level at which GWs break. 618

Data rebalancing can be achieved in two ways. In the first method, we use the dis-619 tribution of the dataset along the identified metric to systematically undersample from the peak and oversample from the tail. Our motivation to undersample from the peak 621 is from the intuition that these samples are over-represented relative to the variability 622 they cover within the dataset, resulting in trained models that may overfit to this region. 623 On the other hand, oversampling the tail is justified by the exact inverse logic: these rare 624 samples are undervalued in their influence over training models. In the second method, 625 the sampling is left unchanged, but the loss function is weighted by the same ratio to 626 increase the penalty on the under-represented class and reduce it for over-represented 627 class. Both data rebalancing strategies generate a new distribution/weighting function 628 on the training dataset with a linear interpolation of the original distribution to a de-629 sired distribution (i.e., uniform distribution) parameterized by $t \in [0, 1]$, much like his-630 togram equalization. We add in an additional parameter to prevent too much oversam-631 pling/weighting in the tail, by the name of maximum repeat. The implementation of these 632 methods requires discretizing the continuous distribution into discrete bins; the choice 633 of the histogram is also an important choice. The methods are implemented by either 634 providing to the learning algorithm a subsample of the dataset that realizes the new dis-635 tribution, or using the new weights in the loss function such that the new distribution 636 is implicitly represented. 637

In our case study, we found that data rebalancing successfully reduces the errors in the moderate tail region while maintaining approximately the same error levels in the peak under most scenarios. In the exception case, data resampling increased the generalization error in the tail, which we attribute to the large size of the ML model. Too large of a model complexity can cause a model to learn noise rather than pattern in the dataset, a phenomenon exacerbated by oversampling in the tail. Unfortunately, we do not observe a clear advantage of the direct sampling implementation over the weighted
loss implementation, nor an unambiguous indication of how to choose the method parameters. Further studies are needed to address these issues on well-understood datasets:
the dataset used for our case study is likely not the best tool for developing intuition for
this method.

Mean bias removal, an additional approach to fix errors with data imbalance, corrects the extant bias in a fully trained data-driven model as a function of the data imbalancerevealing metric This is a first-order correction as it assumes that the mean bias profile of the trained model evolves meaningfully across this metric. The main source of error for this method is generalization error as the mean bias profiles of the training set may not be representative of the instances available at time of inference.

In conclusion, data rebalancing and bias removal show modest improvements in pro-655 ducing data-driven models less inclined to mirror the imbalance apparent in the dataset. The lack of overwhelming evidence of the success of these methods can be attributed to 657 several factors. First, our research did not investigate how to choose the projection used 658 to identify data imbalance, a crucial component to both data rebalancing and mean bias 659 removal. Thus, it may be that the wind range metric is not the most ideal projection 660 for the dataset used in our case study, or that any 1D projection is too simple to cap-661 ture the data imbalance for this dataset. Second, our assumptions on how the data im-662 balance impacts the training of the data-driven models may be overly simplistic, espe-663 cially in its treatment of the tail. We view samples from the tail as in need of a greater significance in training the ML model. However, a more pressing issue at the tail may 665 be that the dataset available to us does not cover the variability inherent to that region. 666 If so, any oversampling does not increase coverage in this region but instead lead to over-667 fitting. We attempt to curb this by introducing the maximum repeat parameter, but this 668 introduces another parameter to be tuned in the rebalancing method. Scarcity of rare 669 (and extreme) phenomena in datasets is a common challenge in geoscience datasets that 670 may be alleviated by rare event sampling, and beyond the scope of the methods presented 671 in this paper. 672

673 6 Open Research

674

6.1 Data Availability

All neural networks used in this manuscript were (re-)written in PyTorch (Paszke 675 et al., 2019). The WaveNet implementation in PyTorch exactly followed the descriptions 676 in (Espinosa et al., 2022). Model of an Idealized Moist Atmosphere (MiMA) (Jucker & 677 Gerber, 2017; Garfinkel et al., 2020) is maintained at https://github.com/mjucker/MiMA 678 and available at https://doi.org/10.5281/zenodo.3984605. The model code, forpy cou-679 pling code, trained ANNs, run parameters, and modified configuration for MiMA are avail-680 able at https://github.com/yangminah/GWPRebalance. The coupling library, forpy, 681 developed and maintained by Elias Rabel is well documented and available at https:// 682 github.com/ylikx/forpy. 683

684 Acknowledgments

This work was supported by the U.S. National Science Foundation through award OAC-2004572 and Schmidt Sciences, as part of the Virtual Earth System Research Institute (VESRI). The manuscript benefited greatly from conversations with Joan Alexander and Pedram Hassanzadeh. We also thank the NYU High Performance Computing center, where the model integrations were performed.

690 References

| 691 | Alexander, M. J., & Dunkerton, T. J. (1999, December). A Spectral Parame- |
|-----|---|
| 692 | terization of Mean-Flow Forcing due to Breaking Gravity Waves. Journal |
| 693 | of the Atmospheric Sciences, 56(24), 4167–4182. Retrieved from http:// |
| 694 | journals.ametsoc.org/doi/10.1175/1520-0469(1999)056%3C4167:ASPOMF% |
| 695 | 3E2.0.C0;2 doi: 10.1175/1520-0469(1999)056(4167:ASPOMF)2.0.CO:2 |
| 696 | Anstev, J. A., Osprev, S. M., Alexander, J., Baldwin, M. P., Butchart, N., Grav, L., |
| 697 | Richter, J. H. (2022). Impacts, processes and projections of the quasi- |
| 608 | biennial oscillation Nature Reviews Earth & Environment 3(9) 588–603 |
| 600 | Brenowitz N D Beucler T Pritchard M & Bretherton C S (2020) Interpret- |
| 700 | ing and stabilizing machine-learning parametrizations of convection <i>Journal of</i> |
| 700 | the Atmospheric Sciences $\gamma\gamma(12)$ $A357-A375$ |
| 701 | Brenowitz N D & Bretherton C S (2010) Spatially Extended Tests of a Neural |
| 702 | Notwork Parametrization Trained by Coarse Craining Lawrad of Advances |
| 703 | in Modeling Farth Systems $11(8)$ 2728 2744 Batrioved 2023 05 26 from |
| 704 | https://onlinelibrory.viley.com/doi/obg/10_1020/2010MS001711 |
| 705 | mutps://onlinelibrary.wiley.com/doi/abs/10.1029/2019MS001711 (|
| 706 | eprint: https://oninenbrary.wney.com/doi/pdi/ $10.1029/2019M5001711$ doi: 10.1020/2010MC001711 |
| 707 | 10.1029/2019MS001711 |
| 708 | Bushell, A., Anstey, J., Butchart, N., Kawatani, Y., Osprey, S., Richter, J., oth- |
| 709 | ers (2022). Evaluation of the quasi-biennial oscillation in global climate models |
| 710 | for the sparc qbo-initiative. Quarterly Journal of the Royal Meteorological |
| 711 | Society, 148(744), 1459–1489. |
| 712 | Chantry, M., Hatfield, S., Dueben, P., Polichtchouk, I., & Palmer, T. (2021). |
| 713 | Machine Learning Emulation of Gravity Wave Drag in Numerical |
| 714 | Weather Forecasting. Journal of Advances in Modeling Earth Sys- |
| 715 | <i>tems</i> , 13(7), e2021MS002477. Retrieved 2023-06-08, from https:// |
| 716 | onlinelibrary.wiley.com/doi/abs/10.1029/2021MS002477 (_eprint: |
| 717 | https://onlinelibrary.wiley.com/doi/pdf/10.1029/2021MS002477) doi: |
| 718 | 10.1029/2021 MS002477 |
| 719 | Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). Smote: |
| 720 | synthetic minority over-sampling technique. Journal of artificial intelligence re- |
| 721 | search, 16, 321–357. |
| 722 | Chawla, N. V., Japkowicz, N., & Kotcz, A. (2004). Special issue on learning from |
| 723 | imbalanced data sets. ACM SIGKDD explorations newsletter, $6(1)$, 1–6. |
| 724 | Connelly, D. S., & Gerber, E. P. (n.d.). Regression forest approaches to gravity wave |
| 725 | parameterization for climate projection. |
| 726 | Ding, D., Zhang, M., Pan, X., Yang, M., & He, X. (2019, July). Modeling Ex- |
| 727 | treme Events in Time Series Prediction. In <i>Proceedings of the 25th ACM</i> |
| 728 | SIGKDD International Conference on Knowledge Discovery & Data Min- |
| 729 | ing (pp. 1114–1122). Anchorage AK USA: ACM. Retrieved 2023-05- |
| 730 | 07, from https://dl.acm.org/doi/10.1145/3292500.3330896 doi: |
| 731 | 10.1145/3292500.3330896 |
| 732 | Elkan, C. (2001). The foundations of cost-sensitive learning. In International joint |
| 733 | conference on artificial intelligence (Vol. 17, pp. 973–978). |
| 734 | Espinosa, Z. I., Sheshadri, A., Cain, G. R., Gerber, E. P., & DallaSanta, K. J. (2022, |
| 735 | April). Machine Learning Gravity Wave Parameterization Generalizes to Cap- |
| 736 | ture the QBO and Response to Increased CO2. <i>Geophysical Research Letters</i> . |
| 737 | 49(8), (Publisher: John Wiley and Sons Inc) doi: 10.1029/2022GL098174 |
| 738 | Garfinkel, C. L. Gerber, F. P., Shamir, O., Bao, J., Jucker, M., White, L. & Paldor, |
| 739 | N. (2022). A abo cookbook: Sensitivity of the quasi-biennial oscillation to res- |
| 740 | olution, resolved waves, and parameterized gravity waves. <i>Journal of Advances</i> |
| 741 | in Modeling Earth Systems 1/(3) e2021MS002568 |
| 742 | |
| 142 | Gartinkel C. I. White I. Gerber E. P. Jucker M. & Erez M. (2020). The build- |
| 742 | Garfinkel, C. I., White, I., Gerber, E. P., Jucker, M., & Erez, M. (2020). The build- ing blocks of northern hemisphere wintertime stationary waves. <i>Journal of Cli</i> |
| 743 | Garfinkel, C. I., White, I., Gerber, E. P., Jucker, M., & Erez, M. (2020). The build- ing blocks of northern hemisphere wintertime stationary waves. <i>Journal of Cli-</i> mate 22(13) 5611-5633 |

He, H., & Garcia, E. A. (2009, September). Learning from imbalanced data. IEEE 745 Transactions on Knowledge and Data Engineering, 21(9), 1263–1284. doi: 10 746 .1109/TKDE.2008.239 747 Johnson, J. M., & Khoshgoftaar, T. M. (2019, December). Survey on deep learning 748 with class imbalance. Journal of Big Data, 6(1). (Publisher: SpringerOpen) 749 doi: 10.1186/s40537-019-0192-5 750 Jucker, M., & Gerber, E. P. (2017, September). Untangling the annual cycle of the 751 tropical tropopause layer with an idealized moist model. Journal of Climate. 752 30(18), 7339-7358.(Publisher: American Meteorological Society) doi: 10 753 .1175/JCLI-D-17-0127.1 754 Krawczyk, B. (2016, November). Learning from imbalanced data: open challenges 755 and future directions. Progress in Artificial Intelligence, 5(4), 221–232. (Pub-756 lisher: Springer Verlag) doi: 10.1007/s13748-016-0094-0 757 Liaw, R., Liang, E., Nishihara, R., Moritz, P., Gonzalez, J. E., & Stoica, I. (2018).758 Tune: A research platform for distributed model selection and training. arXiv759 preprint arXiv:1807.05118. 760 Lindzen, R. S. (1981). Turbulence and stress owing to gravity wave and tidal break-761 down. Journal of Geophysical Research, 86(C10), 9707. (Publisher: American 762 Geophysical Union (AGU)) doi: 10.1029/jc086ic10p09707 763 Oksuz, K., Cam, B. C., Kalkan, S., & Akbas, E. (2021, October). Imbalance 764 Problems in Object Detection: A Review. IEEE Transactions on Pattern 765 Analysis and Machine Intelligence, 43(10), 3388–3415. (Conference Name: 766 IEEE Transactions on Pattern Analysis and Machine Intelligence) doi: 767 10.1109/TPAMI.2020.2981890 768 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., ... others 769 (2019). Pytorch: An imperative style, high-performance deep learning library. 770 Advances in neural information processing systems, 32. 771 Richter, J. H., Anstey, J. A., Butchart, N., Kawatani, Y., Meehl, G. A., Osprey, S., 772 & Simpson, I. R. (2020).Progress in simulating the quasi-biennial oscilla-773 tion in cmip models. Journal of Geophysical Research: Atmospheres, 125(8), 774 e2019JD032362. 775 Richter, J. H., Butchart, N., Kawatani, Y., Bushell, A. C., Holt, L., Serva, F., ... 776 others (2022).Response of the quasi-biennial oscillation to a warming cli-777 Quarterly Journal of the Royal Meteorological mate in global climate models. 778 Society, 148(744), 1490–1518. 779 Rudy, S. H., & Sapsis, T. P. (2023, January). Output-weighted and relative en-780 tropy loss functions for deep learning precursors of extreme events. Physica 781 D: Nonlinear Phenomena, 443, 133570. Retrieved 2023-05-23, from https:// 782 www.sciencedirect.com/science/article/pii/S0167278922002743 doi: 783 10.1016/j.physd.2022.133570 784 Shaw, T. A., Sigmond, M., Shepherd, T. G., & Scinocca, J. F. (2009).Sensitiv-785 ity of simulated climate to conservation of momentum in gravity wave drag 786 parameterization. Journal of climate, 22(10), 2726–2742. 787 Sun, Y. Q., Hassanzadeh, P., Alexander, M. J., & Kruse, C. G. (2023).Quanti-788 fying 3D Gravity Wave Drag in a Library of Tropical Convection-Permitting 789 Simulations for Data-Driven Parameterizations. Journal of Advances in Mod-790 eling Earth Systems, 15(5), e2022MS003585. Retrieved 2023-06-12, from 791 https://onlinelibrary.wiley.com/doi/abs/10.1029/2022MS003585 (_-792 eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1029/2022MS003585) doi: 793 10.1029/2022MS003585 794 Torgo, L., Branco, P., Ribeiro, R. P., & Pfahringer, B. (2015). Resampling strategies 795 for regression. Expert Systems, 32(3), 465-476. doi: 10.1111/exsy.12081 796 Ukkonen. P. (2022).Exploring Pathways to More Accurate Machine Learning 797 Emulation of Atmospheric Radiative Transfer. Journal of Advances in Mod-798 eling Earth Systems, 14(4), e2021MS002875. Retrieved 2023-05-25, from 799

| 800 | https://onlinelibrary.wiley.com/doi/abs/10.1029/2021MS002875 (| |
|-----|---|-----|
| 801 | eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1029/2021MS002875) dc | oi: |
| 802 | 10.1029/2021 MS002875 | |
| 803 | Yuval, J., O'Gorman, P. A., & Hill, C. N. (2021, March). Use of Neural Ne | t- |
| 804 | works for Stable, Accurate and Physically Consistent Parameterization of | |
| 805 | Subgrid Atmospheric Processes With Good Performance at Reduced Pre- | |
| 806 | cision. Geophysical Research Letters, 48(6), e2020GL091363. Retrieve | ed |
| 807 | 2023-05-26, from https://onlinelibrary.wiley.com/doi/abs/10.1029/ | |
| 808 | 2020GL091363 (arXiv: 2010.09947 Publisher: Blackwell Publishing Ltd) do | oi: |
| 809 | 10.1029/2020GL091363 | |
| | | |

⁸¹⁰ Appendix A Formal Algorithm Details

Algorithm 1 shows an example of how to incorporate the direct sampling implementation of the resampling strategy within the framework of any stochastic gradient descent-type learning algorithm that processes batches of training samples at a time. Next, Algorithms 2 and 3 show the direct sampling and weighted loss sampling implementations in detail. Algorithm 1 can easily be modified to use Algorithm 3, where the computed weights are passed into the loss function in the optimization step in line 6, and lines 1, 3, and 4 can be omitted.

| Algorithm 1: Training structure. | | |
|---|--|--|
| Input: \mathcal{X} , Training set; $\hat{\varphi}$, machine learning model; $\{C_n^{(1)}\}_{n=1}^N$, counts of bins of | | |
| ideal histogram; t, linear parameter; max_repeat, maximum repeat | | |
| parameter.r. | | |
| 1 $\{I_n^{(0)}\}_{n=1}^N \leftarrow 	ext{Bin } \mathcal{X} 	ext{ into } N 	ext{ bins.}$ // $I_n^{(0)}$ is the list of indices in the | | |
| nth bin. | | |
| 2 while $\hat{\varphi}$ needs further improvement do | | |
| <pre>// This while-block encompasses a pass over the training set.</pre> | | |
| $3 I^{(t)} \leftarrow \texttt{resample}(\{\mathtt{I}_{\mathtt{n}}^{(0)}\}_{\mathtt{n=1}}^{\mathtt{N}}, \{\mathtt{C}_{\mathtt{n}}^{(0)}\}_{\mathtt{n=1}}^{\mathtt{N}}, \mathtt{t}, \mathtt{max_repeat})$ | | |
| 4 Shuffle $I^{(t)}$ and divide it into B batches $(I^{(t)} = \bigcup_{b=1}^{B} I_b)$. | | |
| 5 for $b=1:B$ do | | |
| 6 Coptimize $\hat{\varphi}$ over $\mathcal{X}[I_b]$. | | |
| 7 return $\hat{\varphi}$ // Trained model | | |

817

Appendix B Architecture Details

We process each of the input features separately with the 1D convolutions. To achieve 819 this, we horizontally stack the features (vertical profiles of zonal wind, U, meridional wind, 820 V, vertical wind, ω , temperature, T as "channels"), resulting in a 2D input shape of **nlev** 821 $\times 4$. (Note that the nomenclature of channels originates from Red Green Blue (RGB) chan-822 nels in image processing.) Additional information such as longitude, latitude, and sur-823 face pressure are concatenated to the flattened output of the encoder. The resulting 1D 824 array is pushed through dense layers intended to represent global relations. Finally, the 825 output from the dense section is reshaped to be processed via transposed convolutions 826 and upsampling layers in the decoder. 827

$\label{eq:algorithm} \textbf{Algorithm 2:} I^{(t)} \gets \texttt{resample}(\{\textbf{I}_n^{(0)}\}_{n=1}^{\mathbb{N}}, \{\textbf{C}_n^{(1)}\}_{n=1}^{\mathbb{N}}, \texttt{t})$

Input: $\{I_n^{(0)}\}_{n=1}^N$, binned indices; $\{C_n^{(1)}\}_{n=1}^N$, counts of bins of ideal histogram; t, linear parameter; max_repeat, maximum repeat parameter. // $I_n^{(0)}$ is the list of indices in the $n{\rm th}$ bin. 1 $I^{(t)} \leftarrow []$ // $I^{(t)}$ is an empty list. $\mathbf{2} \ \mathbf{for} \ n = 1: N \ \mathbf{do}$ Compute $\alpha_n^{(t)}$. // Use Eqs. (3) to (5). 3 $l \leftarrow c_n^{(t)}$ if $\alpha_n^{(t)} \ge 1$ then 4 $\mathbf{5}$ Append $I_n^{(0)}$ floor($lpha_n^{(t)}$) times to $I^{(t)}$. 6 $l \leftarrow c_n^{(t)} - \left(\texttt{count}(\mathbf{I}_{\mathbf{n}}^{(0)}) \times \texttt{floor}(\alpha_{\mathbf{n}}^{(t)}) \right).$ 7 // l is now an integer that satisfies $0 \le l \le \operatorname{count}(\operatorname{I}_n^{(0)})$. Append a random subset of $I_n^{(0)}$ with length l picked without replacement to 8 $I^{(t)}$. 9 return $I^{(t)}$ // New indices.

| $ \textbf{Algorithm 3:} J \gets \texttt{weights}(\{\texttt{I}_n^{(0)}\}_{n=1}^{\mathbb{N}}, \{\texttt{C}_n^{(0)}\}_{n=1}^{\mathbb{N}}, \\$ | t,M) | |
|---|--|--|
| Input: $\{I_n^{(0)}\}_{n=1}^N$, binned indices; $\{C_n^{(1)}\}_{n=1}^N$, counts of bins of ideal histogram; t, linear parameter; M, the size of dataset. | | |
| // $I_n^{(0)}$ is the list of indices in the n th | h bin. | |
| 1 $J \leftarrow \texttt{zeros}(\mathtt{M})$ | // $I^{(t)}$ is an empty list. | |
| 2 for $n = 1 : N$ do | | |
| 3 Compute $\alpha_n^{(t)}$. | // Use Eqs. (3) and (4). | |
| $4 \left[\begin{array}{c} J[I_n^{(t)}] = \alpha_n^{(t)} \end{array} \right]$ | | |
| 5 return J | <pre>// New weights for samples.</pre> | |