

Metadata Normalization

Mandy Lu, Qingyu Zhao, Jiequan Zhang, Kilian M. Pohl, Li Fei-Fei, Juan Carlos Niebles, Ehsan Adeli
Stanford University, Stanford, CA 94305

{mlu, eadeli}@stanford.edu

Abstract

Batch Normalization (BN) and its variants have delivered tremendous success in combating the covariate shift induced by the training step of deep learning methods. While these techniques normalize the feature distribution by standardizing with batch statistics, they do not correct the influence on features from extraneous variables or multiple distributions. Such extra variables, referred to as metadata here, may create bias or confounding effects (e.g., race when classifying gender from face images). We introduce the Metadata Normalization (MDN) layer, a new batch-level operation which can be used end-to-end within the training framework, to correct the influence of metadata on the feature distribution. MDN adopts a regression analysis technique traditionally used for preprocessing to remove (regress out) the metadata effects on model features during training. We utilize a metric based on distance correlation to quantify the distribution bias from the metadata and demonstrate that our method successfully removes metadata effects on four diverse settings: one synthetic, one 2D image, one video, and one 3D medical image dataset.

1. Introduction

Recent advances in fields such as computer vision, natural language processing, and medical imaging have been propelled by tremendous progress in deep learning [7]. These deep neural models owe their success to their large number of trainable parameters, which encode rich information from the data. However, since the learning process can be extremely unstable, much of the work is spent on carefully selecting a model through hyperparameter tuning, an integral part of approaches such as [31, 39]. To aid with model development, normalization techniques such as Batch Normalization (BN) [31] and Group Normalization (GN) [61] make the training process more robust and less sensitive to covariate or distribution shift.

BN and GN perform feature normalization by standardizing them solely using batch or group statistics (*i.e.*, mean and standard deviation). Although they have pushed the

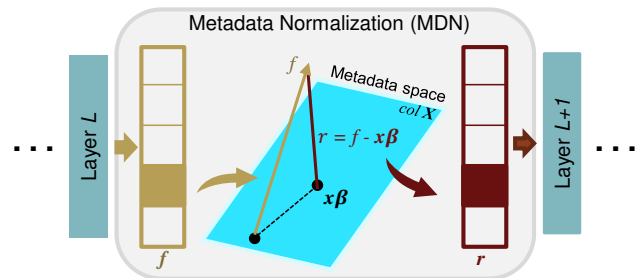


Figure 1. The proposed Metadata Normalization operation. MDN layer takes the learned features from the previous layer (f), analyzes the effects of the metadata on them, residualizes such effects, and outputs the distribution corrected features (r).

state-of-the-art forward, they do not handle extraneous dependencies in the data other than the input and output label variables. In many applications, confounders [57, 69] or protected variables [46] (sometimes referred to as bias variables [2]) may inject bias into the learning process and skew the distribution of the learned features. For instance, (i) when training a gender classification model from face images, an individual’s race (quantified by skin shade) has a crucial influence on prediction performance as shown in [11]; (ii) in video understanding, action recognition models are often driven by the scene [15, 30] instead of learning the harder movement-related action cues; (iii) for medical studies, patient demographic information or data acquisition site location (due to device and scanner differences) are variables that easily confound studies and present a troublesome challenge for the generalization of these studies to other datasets or everyday clinical usages [10, 68].

This additional information about training samples is often freely available in datasets (*e.g.*, in medical datasets as patient data) or can be extracted using off-the-shelf models (such as in [11, 15]). We refer to them as *metadata*, namely “data that provides information about other data” [59], an umbrella term for additional variables that provide information about training data but are not directly used as model input. The extraneous dependencies between the training data and metadata directly affect the distributions of the learned features; however, typical normalization op-

erations such as BN and GN operate agnostic to this extra information. Instead, current strategies to remove metadata effects include invariant feature learning [4, 42, 62] or domain adaptation [26, 55].

Traditional handcrafted and feature-based statistical methods often use intuitive approaches based on multivariate modeling to remove the effects of such metadata (referred to as study confounders in this setting). One such regression analysis method [41] builds a Generalized Linear Model (GLM) between the features and the metadata (see Fig. 1) to measure how much the feature variances are explained by the metadata versus the actual output (*i.e.*, ground-truth label) [1, 68]. The effects of the metadata can then be removed from the features by a technique referred to as “regressing out” the effects of the extraneous variables [1, 10, 45]. The application of this GLM-based method to deep end-to-end architectures has not yet been explored because it requires precomputed features to build the GLM and is traditionally performed on the dataset prior to training. Thus, this method is inapplicable to vision problems with pixel-level input and local spatial dependencies, which a GLM is unable to model. The key insight we use is that the later layers of a network represent high-level features with which we can build our GLM. In this paper, we extend this widely-explored and seminal regression analysis method by proposing a corresponding operation for deep learning architectures *within a network* to remove the metadata effects from the intermediate features of a network.

As illustrated in Fig. 1, we define a Metadata Normalization (MDN) layer which applies the aforementioned regression analysis as a normalization technique to remove the metadata effects from the features in a network¹. Our MDN operation projects each learned feature \mathbf{f} of the L^{th} layer to the subspace spanned by the metadata variables, denoted by \mathbf{X} , by creating a GLM [41, 43] $\mathbf{f} = \mathbf{X}\beta + \mathbf{r}$, where β is a learnable set of linear parameters, $\mathbf{X}\beta$ corresponds to the component in \mathbf{f} explained by the metadata, and \mathbf{r} is the residual component irrelevant to the metadata. The MDN layer removes the metadata-related components from the feature (Fig. 1) and regards the residual \mathbf{r} as the normalized feature impartial to metadata. We implement this operation in a (mini)batch iterative training setting.

As opposed to BN and its variants that aim at *normalizing the distribution* of the features throughout the training process, MDN focuses on *correcting the distribution* with respect to the chosen metadata variables. When employed in end-to-end settings, this enables deep learning architectures to remove the effects of confounders, protected variables, or biases during the training process. Moreover, the metadata will only correct the distributions if there are

¹Although metadata normalization was previously used to refer to the adjustment of metadata elements into standard formats [36], we redefine the term as an operation in deep learning architectures.

such distributions explained by the metadata. On the other hand, if the learned features are orthogonal to the metadata subspace (*i.e.*, features are not biased by the metadata variables), the β coefficients will be close to zero and hence will not alter the learning paradigm of the network.

In summary, our work makes the following primary contributions: (1) We propose the Metadata Normalization technique to correct the distribution of the data with respect to the additional, metadata, information; (2) We present a theoretical basis for removal of extraneous variable effects based on GLM and introduce a novel strategy to implement this operator in (mini)batch-level settings; (3) For the cases when output prediction variables are intrinsically correlated with the metadata, we outline a simple extension to MDN to ensure that only extraneous effects of the metadata are removed and not those that pertain to the actual output variables. Our implementation as a simple PyTorch layer module is available at <https://github.com/mlu355/MetadataNorm>. We show the effectiveness of MDN in four different experimental settings, including one synthetic, one image dataset for gender classification from face images, one video scene-invariant action recognition, and one multi-site medical image classification scenario.

2. Related Work

Normalization in Deep Learning: Prior normalization techniques for neural models include Batch Normalization (BN) as a canonical example. BN has been a wildly effective mechanism to greatly speed up training and boost model convergence [31]. Operations such as Group Normalization (GN) [61] and Layer Normalization [6] are similar and apply the underlying concepts behind BN to smaller batch sizes and recurrent networks. Similar to BN and GN, MDN is applied at a dataset level implemented in (mini)batch-level settings, but it differs in how it shifts the distributions of the features (w.r.t. the metadata).

Statistical Methods for Regressing Out Confounders: Traditional feature-based statistical methods for removing confounders include stratification [19], techniques using Analysis of Variance (ANOVA) [45], and the use of multivariate modeling such as regression analysis [41] with the statistical GLM method outlined above [1, 68]. However, due to their ineffectiveness in dealing with pixel-level data and dependency on handcrafted features, vision-based tasks and end-to-end methods typically use other techniques to alleviate bias or the effects of study confounders. Common techniques include the use of data preprocessing techniques to remove dataset biases such as sampling bias [64] and label bias [32]. Other methods to remove bias from machine learning classifiers include the use of post-processing steps to enforce fairness on already trained, unfair classification models [23, 29, 68]. However, algorithms which decouple training from the fairness enforcement may lead to a subop-

timal fairness and accuracy trade-off [60]. Herein, we apply the ideas behind feature-based confounder removal through regression analysis as a batch-level module, which can be added synchronously to the training process.

Bias in Machine Learning: Bias in machine learning models is an increasingly scrutinized topic at the forefront of machine learning research. The prevalence of bias in large public datasets has been a cause for alarm due to their propagation or even amplification of bias in the models which use them. Recent examples of dataset bias in public image datasets such as ImageNet [63], IARPA Janus Benchmark A (IJB-A) face [35] and Adience [21] have shown that they are imbalanced with mainly light-skinned subjects and that models trained on them retain this bias in their predictions [11, 47]. Bias is prevalent in a wide range of disciplines, such as gender bias in natural language processing via word embeddings, representations, and algorithms [9, 16, 51] and medical domains [34] such as genomics [13] and Magnetic Resonance Imaging (MRI) analyses [3, 22], in which it is common for data to be skewed toward certain populations [25, 50]. Models trained on such biased settings produce biased predictions or can amplify existing biases. Many approaches have been developed to remove these adverse effects for both qualitative causes (*e.g.*, for fairness) and for quantitative causes (*e.g.*, improving the performance of a model by reducing its dependence on confounding effects).

Fair representation learning is an increasingly popular approach to learn debiased intermediate representations [65] that has been explored in numerous recent works [2, 5, 29, 37, 40, 58], with [33, 54] introducing methods to apply fair and invariant representation learning to continuous protected variables. Recently, adversarial learning has also become a popular area for exploration to mitigate bias in machine learning models [2, 66]. Our MDN paradigm can interpret this problem as correcting the feature distributions by treating the bias and protected variables as the study metadata, and hence has interesting applications for fair representation learning. In contrast with the domain adaptation and invariant feature learning frameworks, MDN is a layer which easily plugs into an end-to-end learning scheme and is also applicable to continuous protected variables.

3. Method

With a dataset including N training samples with prediction labels $(\mathbf{I}_i, \mathbf{y}_i)$ for $i \in \{1, \dots, N\}$, we train a neural network with trainable parameters Θ using a 2D or 3D backbone, depending on the application (2D for images, and 3D for videos or MRIs). MDN layer can be inserted between all convolutional and fully connected layers to correct the distribution of the learned features within the stochastic gradient descent (SGD) framework.

3.1. Metadata Normalization (MDN) via GLM

Let $x_i \in \mathbb{R}^K$ be a column vector storing the K -dimensional metadata of the i^{th} sample and $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top$ be the metadata matrix of all N training samples. Let $\mathbf{f} = [f_1, \dots, f_N] \in \mathbb{R}^N$ be a feature extracted at a certain layer of the network for the N samples. A general linear model associates the two variables by $\mathbf{f} = \mathbf{X}\beta + \mathbf{r}$, where β is an unknown set of linear parameters, $\mathbf{X}\beta$ corresponds to the component in \mathbf{y} explained by the metadata, and \mathbf{r} is the residual component irrelevant to the metadata. Therefore, the goal of the MDN layer is to remove the metadata-related components from the feature:

$$\mathbf{r} = \text{MDN}(\mathbf{f}; \mathbf{X}). \quad (1)$$

In this work, we use an ordinary least square estimator to solve the GLM so that the MDN layer can be reduced to a linear operator. Specifically, the optimal β is given by the closed-form solution

$$\beta = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{f}, \quad (2)$$

and the MDN layer can be written as

$$\mathbf{r} = \mathbf{f} - \mathbf{X}\beta = \mathbf{f} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{f} \quad (3)$$

$$= (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \mathbf{f} = (\mathbf{I} - \mathbf{P}) \mathbf{f} = \mathbf{R} \mathbf{f}. \quad (4)$$

Geometrically, \mathbf{P} is the projection matrix onto the linear subspace spanned by the metadata (column vectors of \mathbf{X}) [8]. The residualization matrix $\mathbf{R} \in \mathbb{R}^{N \times N}$ is the residual component orthogonal to the metadata subspace (Fig. 1).

3.2. Batch Learning

In a conventional GLM, both \mathbf{X} and \mathbf{R} are constant matrices defined with respect to all N training samples. This definition poses two challenges for batch stochastic gradient descent. To show this, let $\hat{\mathbf{X}} \in \mathbb{R}^{M \times K}$ and $\hat{\mathbf{f}} \in \mathbb{R}^M$ be the metadata matrix and feature associated with M training samples in a batch. In each iteration, we need to re-estimate the corresponding residualization matrix $\hat{\mathbf{R}} \in \mathbb{R}^{M \times M}$, which by Eq. (4) would require re-computing the matrix inverse $(\hat{\mathbf{X}}^\top \hat{\mathbf{X}})^{-1}$, a time-consuming task. Moreover, the GLM analysis generally results in sub-optimal estimation of β when few training samples are available (*i.e.*, $M \ll N$). To resolve these issues, we further explore the closed-form solution of Eq. (2), which can be re-written as

$$\beta = (\mathbf{X}^\top \mathbf{X})^{-1} \sum_{i=1}^N \mathbf{x}_i f_i \approx N \Sigma^{-1} \mathbb{E}[\mathbf{x} f], \quad (5)$$

where $\Sigma = \mathbf{X}^\top \mathbf{X}$ is a property solely of the metadata space independent of the learned features f and $\mathbb{E}[\mathbf{x} f] \approx \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i f_i$. We propose to pre-compute Σ^{-1} on all N

training samples to derive the most accurate characterization of the metadata space before training. During each training step, we compute the batch-level estimate of the expectation $\mathbb{E}[\mathbf{x}f] = \frac{1}{M} \hat{\mathbf{X}}^\top \hat{\mathbf{f}}$. Hence, the batch estimation of the residualization matrix is

$$\hat{\mathbf{R}} \approx \hat{\mathbf{f}} - \hat{\mathbf{X}}(N\Sigma^{-1}\mathbb{E}[\mathbf{x}f]) \quad (6)$$

$$\approx (\mathbf{I} - \frac{N}{M} \hat{\mathbf{X}} \Sigma^{-1} \hat{\mathbf{X}}^\top) \hat{\mathbf{f}}. \quad (7)$$

We have rederived our residual solution from Eq. (4) with the addition of the scaling constant $\frac{N}{M}$ and $\hat{\mathbf{X}}^\top \hat{\mathbf{X}}^{-1}$ replaced by our precomputed Σ^{-1} .

3.3. Evaluation

During training, we store aggregated batch-level statistics to use during evaluation, when we may not have a large enough batch to form a reliable GLM solution. Observe from Eq. (5) that β is re-estimated in each training batch because the features \mathbf{f} are updated after the batch. Since the testing process does not update the model, the underlying association between features and metadata is fixed, so the β estimated from the training stage can be used to perform the metadata residualization. To ensure that the β estimation can accurately encode the GLM coefficients associated with the entire training set and to avoid oscillation from random sampling of batches), we update β at each iteration using a momentum model [52]:

$$\beta^k = \eta\beta^k + (1 - \eta)\beta^{k-1}, \quad (8)$$

where $k \in \{1, \dots, \tau\}$ is the batch index and η is the momentum constant. During testing, we no longer solve for β and instead use the estimate β_τ from the last training batch

$$\text{MDN}(\mathbf{f}; \mathbf{X}) = \mathbf{f} - \mathbf{X}\beta_\tau. \quad (9)$$

The batch-level GLM solution will approach the optimal group-level solution with increasing batch size, as larger batches produce a better estimate for the dataset-level GLM solution during both training and evaluation.

3.4. Collinearity between Metadata and Labels

In more complicated scenarios where confounding effects occur, the metadata not only affects the training input but also correlates with the prediction label. In this case, we need to remove the direct association between \mathbf{f} and \mathbf{X} while preserving the indirect association created via \mathbf{y} [69]. We control for the effect of \mathbf{y} by reformulating the GLM as

$$\mathbf{f} = \mathbf{X}\beta_X + \mathbf{y}\beta_y + \mathbf{r} = \tilde{\mathbf{X}}\tilde{\beta} + \mathbf{r}, \quad (10)$$

where \mathbf{y} is prediction labels vector of the training samples, $\tilde{\mathbf{X}}$ is the horizontal concatenation of $[\mathbf{X}, \mathbf{y}]$, and $\tilde{\beta}$ is the

vertical concatenation of $[\beta_X; \beta_y]$. This multiple regression formulation allows us to separately model the variance within the features explained by the metadata and by the labels, so that we only remove the metadata-related variance from the features. To perform MDN in this scenario, we first estimate the composite $\tilde{\beta}$ in a similar way as in Eq. (5) for each batch during training

$$\tilde{\beta} \approx N\tilde{\Sigma}^{-1}\mathbb{E}[[\mathbf{x}f; \mathbf{y}f]], \quad (11)$$

where $\tilde{\Sigma}$ is the covariance matrix of $\tilde{\mathbf{X}}$ estimated on the whole training population, and the expectation \mathbb{E} is computed on the batch level. Next, unlike the previous MDN implementation in Eq. (7), the residualization is now only performed with respect to β_X

$$\text{MDN}(\mathbf{f}; \mathbf{X}) = \mathbf{f} - \mathbf{X}\beta_X. \quad (12)$$

Controlling for the labels when fitting features to the metadata preserves the components informative to prediction in our residual and thus in the ensuing features of the network.

4. Experiments

We test our method on a variety of datasets covering a diverse array of settings, including both categorical and continuous metadata variables, binary and multi-class classification, and multi-label. For all experiments, our baseline is a vanilla convolutional neural network (2D or 3D CNN), to which we add the proposed MDN to assess its influence on the model learning process. We show that adding MDN to a model can result in improved or comparable prediction accuracy while reducing model dependence on the metadata. The collinearity of the metadata with the labels is handled by adopting the MDN implementation in Section 3.3. We test our method by (1) adding MDN to solely the final fully-connected layers of the network (which we refer to as MDN-FC) and (2) adding MDN to the convolutional layers in addition to the final linear layers (MDN-Conv). For comparison, we add other normalization layers such as Batch Normalization (BN) and Group Normalization (GN) to all convolutional layers of the baseline.

Computational efficiency is one of the strengths of our method, as there are no learnable parameters due to the closed form solution for batch learning in Section 3.2. Therefore, memory cost is negligible and training time comparable to models without MDN. We used NVIDIA GTX 1080 Ti with 11GB VRAM for image experiments and TI-TAN RTX with 24GB VRAM for video experiments.

4.1. Metrics

For each of the experiments, we use the squared distance correlation (dcor²) [53] between our model features and the metadata variables as the primary quantitative metric for measuring the magnitude of the metadata effect on

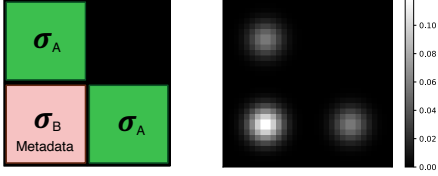


Figure 2. Synthetic data generation process with a generated training sample. The main diagonal Gaussian σ_A differentiates the two groups while the off-diagonal Gaussian σ_B serves as the metadata.

the features. Unlike univariate linear correlation, dcor^2 measures non-linear dependency between high-dimensional variables. The lower the dcor^2 , the less the learned features (and hence the model) are affected by the extraneous variables, with $\text{dcor}^2 = 0$ indicating statistical independence. The goal is to minimize dcor^2 and reduce the dependence between the features and metadata variables. We additionally compute balanced accuracy (bAcc) for all experiments to measure prediction performance. With regard to individual experiments, we compute the following additional metrics: (1) for the GS-PPB experiment, we compute accuracy per shade; (2) for the HVU experiment, we compute mean average precision (mAP) for the action classification task.

4.2. Synthetic Experiments

The synthetic experiments are constructed as a binary classification task on a dataset of random synthetically generated images comprised of two groups of data, each containing 1000 images of resolution 32×32 pixels. Each image is generated by 4 Gaussians, the magnitude of which is controlled by parameters σ_A (controlling quadrants II and IV) and σ_B (controlling quadrant III). Images from Group 1 are generated by sampling σ_A and σ_B from a uniform distribution $\mathcal{U}(1, 4)$, while images from Group 2 are generated with stronger intensities from $\mathcal{U}(3, 6)$ (see Figure 2). The difference in σ_A between the two groups is associated with the true discriminative cues that should be learned by a classifier, whereas σ_B is a metadata variable. Therefore, an unbiased model which is agnostic to σ_B should predict the group label purely based on the two diagonal Gaussians without depending on the off-diagonal Gaussian. The overlapping sampling range of σ_A between the two groups leads to a theoretical maximum accuracy of 83.33%.

Our baseline is a simple CNN with 2 convolution/ReLU stacks followed by 2 fully-connected layers of dimension (18432, 84, 1) with Sigmoid activation. We observe the effect of adding MDN to various layers of the baseline: the first (MDN-FC) has one MDN layer applied to the first fully-connected layer, and the second (MDN-Conv) additionally applies MDN to the convolutional layers.

Table 1 shows the results of 100 runs of each model over batch sizes 200, 1000, and 2000 with 95% confidence intervals (CIs). Our baseline achieves 94.1% training accu-

Table 1. Comparison of models on the Synthetic Dataset over 100 runs with 95% CIs for dcor^2 (lower is better) and bAcc (closer to 83.3% is better). Note that the theoretical maximum accuracy of an unbiased model is 83.3%, so significantly higher values indicate that the model is “cheating” by using the metadata.

Model	Batch	dcor^2	bAcc
Baseline [28]	200	0.399 ± 0.014	94.1 ± 0.0
	1000	0.464 ± 0.004	94.1 ± 0.0
	2000	0.479 ± 0.005	94.1 ± 0.0
BN [31]	200	0.331 ± 0.003	93.2 ± 0.1
	1000	0.289 ± 0.004	93.5 ± 0.1
	2000	0.273 ± 0.004	93.6 ± 0.1
GN [61]	200	0.368 ± 0.010	94.0 ± 0.1
	1000	0.399 ± 0.009	94.0 ± 0.1
	2000	0.435 ± 0.009	94.0 ± 0.1
MDN-FC	200	0.189 ± 0.010	90.7 ± 0.1
	1000	0.043 ± 0.008	86.7 ± 0.7
	2000	0.028 ± 0.012	82.4 ± 1.2
MDN-Conv	200	0.181 ± 0.019	89.5 ± 0.8
	1000	0.017 ± 0.007	82.8 ± 0.4
	2000	0.003 ± 0.000	83.4 ± 0.1

acy, significantly higher than the theoretical maximum accuracy of 83.3%, so it must be falsely leveraging the metadata information for prediction. Similarly, BN and GN produce accuracies around 93% and 94%, confirming that neither technique corrects for the distribution shift caused by the metadata. On the other hand, MDN-Conv and MDN-FC produce accuracies much closer to the theoretical unbiased optimum. As batch size increases, both MDN-FC and MDN-Conv decrease in accuracy until hitting the theoretical optimum, which suggests they have completely removed their dependence on the metadata without removing components of features which aid in prediction. MDN-FC reaches the max accuracy at batch size 2000 with 82.4% and MDN-Conv reaches the max accuracy more quickly, with 82.8% at batch 1000. We measure dcor^2 between the features in the first FC layer and the metadata variable for samples from each group separately (Fig. 3) and then record the average in Table 1. We observe that dcor^2 for both MDN models is significantly lower than the lowest baseline dcor^2 of 0.399 (Table 1). The dcor^2 decreases as batch size increases; *e.g.*, when using a batch size of 2000, the correlation drops to virtually 0, indicating an exact independence between network features and the metadata variable σ_B . These results corroborate our expectation from section 3.2 that the batch-level GLM solution will approach the optimal group-level solution with large batches.

MDN-Conv has superior results to MDN-FC in terms of lower dcor^2 and accuracy closer to 83.33%. This suggests that forming a sequence of linear models by applying MDN to successive layers (after each convolution) may gradually remove nonlinear effects between features and metadata variables. Figure 4 shows the tSNE visualization of features extracted from the baseline and from MDN-Conv

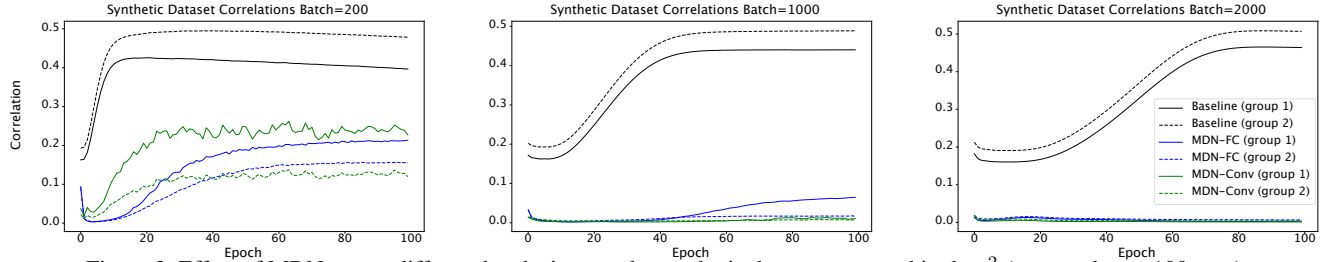


Figure 3. Effect of MDN across different batch sizes on the synthetic dataset measured in $dcor^2$ (averaged over 100 runs).

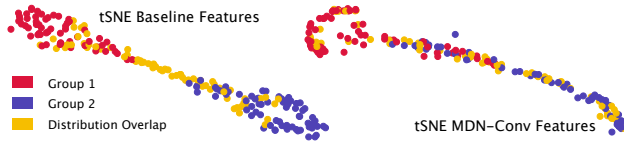


Figure 4. tSNE of features extracted from baseline and MDN-Conv on the synthetic dataset. Samples from the overlapping region $\mathcal{U}(3, 4)$ are separable in baseline but not in MDN-Conv.

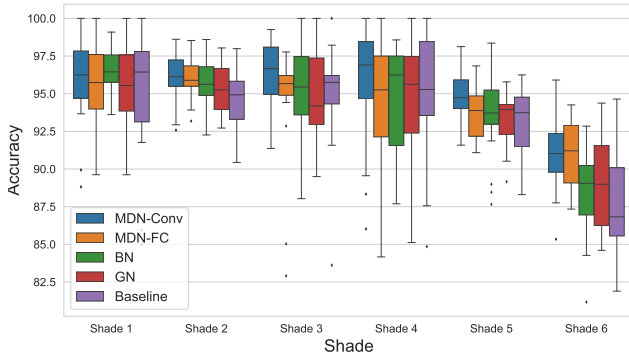


Figure 5. Accuracy by shade on the GS-PPB dataset for different methods pretrained on ImageNet.

with 3 groups: samples with kernels sampled from the overlapping $\mathcal{U}(3, 4)$ region, which should not be separable without using the metadata, and samples separable into Group 1 and Group 2 using only σ_A (with kernels in $\mathcal{U}(1, 3)$ and $\mathcal{U}(4, 6)$). Our tSNE plot shows that the overlapping region is separable in the baseline features but not in MDN-Conv.

4.3. Gender Prediction Using the GS-PPB Dataset

The next experiment is gender prediction on the face images in the Gender Shades Pilot Parliaments Benchmark (GS-PPB) dataset [11]. GS-PPB contains 1,253 facial images of 561 female and 692 male subjects, each labeled with a shade on the Fitzpatrick six-point labeling system [24] from type 1 (lighter) to type 6 (darker). Face detection is used to crop the images to ensure that our classification relies solely on facial features [27]. It has been shown that models pre-trained on large public image datasets such as ImageNet amplify the dependency between shade and target labels due to dataset imbalance [11, 63]. A large dis-

Table 2. Per-shade bAcc and per-class $dcor^2$ on GS-PPB. Results are averaged over 5 runs of 5-fold CV with 95% Confident Intervals. Best results are bolded and second best are underlined.

Shade	MDN-Conv	MDN-FC	Baseline [28]	BN [31]	GN [61]
1	96.2±1.1	95.9±1.2	95.2±1.5	96.6±0.6	95.6±1.5
2	96.2±0.6	<u>96.0±0.6</u>	94.5±0.8	95.7±0.7	95.6±0.9
3	96.4±0.8	94.5±1.7	94.9±1.2	<u>95.2±1.1</u>	95.1±1.4
4	95.7±1.6	94.3±1.8	<u>95.0±1.6</u>	94.9±1.4	94.6±2.1
5	95.1±0.7	<u>93.8±0.8</u>	93.3±0.8	<u>93.8±1.0</u>	93.3±1.0
6	91.0±0.9	<u>90.9±1.0</u>	87.6±1.4	<u>88.5±1.0</u>	89.0±1.6
Avg.	95.1±0.4	<u>94.2±0.5</u>	93.4±0.7	94.1±0.4	93.8±0.5
$dcor^2$ F	0.06±0.01	<u>0.07±0.01</u>	0.14±0.02	0.22±0.02	0.17±0.01
$dcor^2$ M	0.08±0.01	<u>0.09±0.01</u>	0.05±0.01	0.10±0.01	0.10±0.01
Avg.	0.07±0.01	<u>0.08±0.01</u>	0.10±0.01	0.16±0.01	0.13±0.01

crepancy in classification accuracy of such pre-trained models has been observed between lighter and darker shades, with lowest accuracy in shades 5 and 6. In this experiment, we aim to reduce the shade bias in a baseline VGG16 backbone model [48] pre-trained on ImageNet [18] (chosen for its known dataset bias to shade [63]) by fine-tuning on the GS-PPB dataset using MDN with shade as the metadata variable. In our VGG16 baseline, we replace the final FC layers with a simple predictor of two FC layers. We test MDN by applying it to the first FC layer (MDN-FC) and additionally to the last convolutional layer (MDN-Conv).

Table 2 shows prediction results across five runs of 5-fold cross-validation. Per shade accuracies are further visualized in the box plot Figure 5. MDN-Conv and MDN-FC both achieve higher accuracies on the darker shades 5 and 6 with comparable or higher performance for other shades, correcting for the bias in the baseline VGG16 pretrained on ImageNet. Both MDN models also achieve the highest average bAcc and lowest correlations, with MDN-Conv obtaining the highest average bAcc of 95.1% and the lowest $dcor^2$ on both classes (F:0.06, M:0.08, Avg:0.07). BN and GN operations increase accuracy when applied to the baseline, which is expected, as these operations have been shown to improve model stability and performance. However, when compared with the MDN models, they produce higher correlation and less robust results MDN, with smaller average accuracy and a much higher $dcor^2$ for females than for males, indicating that they are more heavily leveraging if a person is female vs male. This difference is greatly reduced by MDN, which has $dcor^2$ agnostic to gender. This shows that in settings

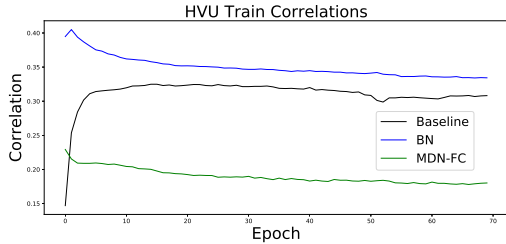


Figure 6. Comparison of train $dcor^2$ for MDN-FC, Baseline and BN on HVU pretrained on the Kinetics-700 Dataset

where features are heavily impacted by metadata, MDN can successfully correct the feature distribution to improve results over normalization methods that only perform standardization over using batch or group statistics.

4.4. Action Recognition Using the HVU Dataset

The Holistic Video Understanding (HVU) dataset [20] is a large scale dataset that contains 572k video clips of 882 different human actions. In addition to the action labels, the videos are annotated with labels of other categories including 282 scene labels. We use the original split from the paper, with 481k videos in the training set and 31k videos in the validation set. Our task is action recognition with scene as our metadata, aiming to reduce the direct reliance of our model on scene. Action recognition architectures are often biased by the background scene because videos of the same action are captured in similar scenes [15, 30]. These architectures may capture the easier scene cues rather than the harder-to-understand movement cues that define the action in time, which can reduce generalizability to unseen cases.

We use a 3D-ResNet-18 [28] architecture pretrained on the Kinetics-700 dataset [12] as our BN model and our baseline is the same with no normalization (BN layers removed). MDN is added to the final FC layer of the baseline before the output layer (MDN-FC). The metadata variables are one-hot vector encodings of the top 50 scenes by occurrence. The model is fine-tuned on the HVU dataset until validation accuracy converges, at around 70 epochs. Figure 6 and Table 3 show that while the baseline model increases in $dcor^2$ during training, MDN and BN decrease, even though BN has higher $dcor^2$. MDN displays the lowest $dcor^2$ by far of 0.182, so we have clearly succeeded at reducing the dependence between our model features and scene. However, both MDN and BN experience slightly lower mAP than the baseline without normalization. This is not surprising since scenes may provide information about the action itself, separating cases where the model directly uses scene for prediction and indirectly does so with action as an intermediate dependency. Thus, removing scene dependence may hurt model performance by either preventing it from “cheating” by directly using the scene for prediction, or removing useful components of the features.

Table 3. Performance of MDN and baseline on the HVU dataset.

Model	$dcor^2$	mAP (%)
BN (3D-ResNet-18) [31]	0.335	39.7
Baseline (3D-ResNet-18 - BN) [28]	0.307	42.4
MDN-FC (3D-ResNet-18 + MDN)	0.182	40.3

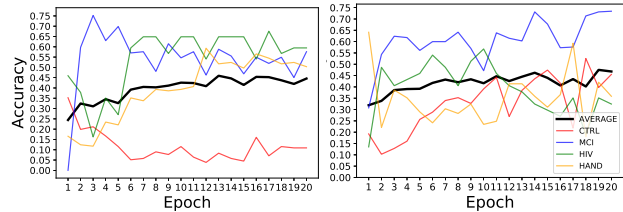


Figure 7. Prediction accuracy of each cohort and the balanced accuracy over the 4 cohorts on the testing folds (averaged over the 5 folds) versus training iterations.

We have demonstrated that MDN successfully removes direct reliance of our model on scene, but further exploration is needed to interpret the effect on model performance. This video action recognition experiment is also particularly challenging due to the large dimensionality of video input. Thus, it is difficult to produce the large batch sizes needed to best estimate the GLM parameters. This is a common problem in video recognition tasks for batch-level operations such as BN. Several prior works [14, 17, 49] have proposed work-around solutions by calculating aggregated gradients from several batches to virtually increase the batch size or aggregated batch-level statistics (for BN and for our case the GLM parameters). This is an area that requires further study and we anticipate that implementing such strategies may improve results further.

4.5. Classification of Multi-Site Medical Data

The last experiment is diagnostic disease classification of 4 cohorts of participants based on their T1-weighted 3D MRI scans. The 4 cohorts are healthy control (CTRL) subjects, subjects that show Mild Cognitive Impairment

Table 4. Multi-site multi-label disease diagnosis classification based on 3D MRIs and accuracy scores with respect to the target (UCSF) dataset. Recall rate (accuracy) for each cohort, bAcc of all cohorts, standard deviation of per-cohort recall, and distance correlation between the learned features and dataset labels.

Class	# of Subjects		Baseline	BN [31]	GN [61]	MDN
	Total	UCSF				
CTRL	460	156	8.3%	9.6%	14.1%	41.7%
HIV	112	37	63.0%	62.2%	66.9%	32.4%
MCI	732	335	37.8%	54.6%	51.4%	73.7%
HAND	145	145	42.1%	55.9%	52.4%	42.0%
Overall bAcc			37.8%	45.6%	46.1%	47.5%
Standard Deviation			22.5%	25.2%	22.5%	18.1%
$dcor^2$			0.26	0.30	0.34	0.06

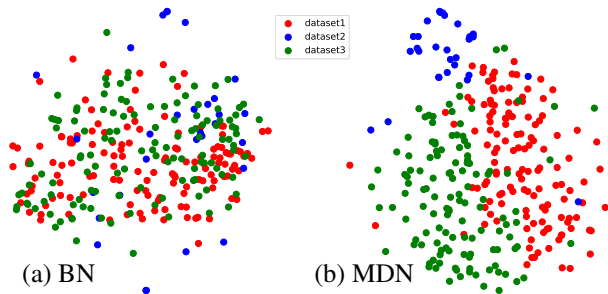


Figure 8. tSNE visualization of the features extracted from the convolutional layers with BN or with MDN. Each data point is color-coded based on the dataset label.

(MCI), those diagnosed with Human Immunodeficiency Virus (HIV) infection, and subjects with HIV-Associated Neurocognitive Disorder (HAND). Since HAND is a comorbid condition that combines the characteristics of HIV and MCI, the classification is formulated as a multi-label binary classification problem, where we predict for each subject whether 1) the subject has MCI diagnosis; and 2) the subject is HIV-positive. The HAND patients are positive for both labels and the CTRLs are negative for both.

The T1-weighted MR images used in this study were collected at the Memory and Aging Center, University of California - San Francisco (UCSF; PI: Dr. Valcour) [67] shown in Table 4. Since the number of subjects is relatively small, especially for the HIV cohort ($N = 37$), we augment the training dataset with MRI scans collected by the Neuroscience Program, SRI International (PI: Dr. Pfefferbaum), consisting of 75 CTRLs and 75 HIV-positive subjects [1], and by the public Alzheimer’s Disease Neuroimaging Initiative (ADNI1) [44], which contributed an additional 229 CTRLs and 397 MCI subjects. To perform classification on such multi-site data, the source of the data (dataset label) becomes the metadata, which is parameterized here by one-hot encodings. Medical imaging datasets acquired in multiple sites with different scanning protocols is a core challenge for machine learning algorithms in medicine, [38, 56], as different scanning protocols lead to different image formations. Differing class formations across sites (as in this experiment) creates a simple undesirable cue for the model to leverage during prediction as a confounder. We corrected the features distribution by deeming the acquisition site as our metadata variable and employing our MDN operation.

The baseline classification model consists of a feature extractor and a classifier. We designed the feature extractor as 4 stacks of $3 \times 3 \times 3$ convolution/ReLU/BN/max-pooling layers with dimension (16, 32, 64, 32). The classifier consists of 2 fully connected layers with dimension (2048, 128, 16). We construct each batch by sampling 10 subjects from each cohort of each dataset (with replacement). The model accuracy is evaluated by 5-fold cross-validation with re-

spect to the 4 cohorts of UCSF, which is the primary goal of the experiment. We train the model for 20 epochs until the bAcc on the testing folds (averaged over the 5 testing folds) converges. We then rerun the experiment by replacing the BN layers in the baseline model with MDN layers.

We observe from Table 4 and Fig. 7 that MDN improves bACC for the multi-label prediction compared to the baselines with BN and GN. The baseline models exhibit highly imbalanced predictive power among the 4 cohorts reflected by the large discrepancy in per-cohort recall (std of 25.2% and 22.5%). This is in part because the 3 datasets represent distinct cohort constructions (*e.g.*, ADNI only contains MCI, but no HIV) so the multi-domain feature distribution is likely to bias the discriminative cues related to the neurological disorders. This, however, is not the case for MDN, which successfully reduces the accuracy discrepancy among cohorts (std of 18.1%). The CTRL group received an especially high increase in bAcc, from 9.6% in BN and 14.1% in GN to 41.7% in MDN.

The reduced dataset bias is also evident in the distance correlation analysis, which examines the dependency between the features extracted from the convolutional layers and the dataset label. Table 4 records the average $dcor^2$ derived over the 5 testing folds, and MDN achieves a significantly lower metric than the baseline model. Qualitatively, we randomly select a testing fold and use t-SNE to project the features learned by the two models into a 2D space and color-code the data point by their dataset label (Fig. 8). The features are clearly clustered by dataset assignment for the baseline model, whereas this adverse clustering effect is significantly reduced after MDN.

5. Conclusion

We presented a novel normalization operation for deep learning architectures, denoted by Metadata Normalization (MDN). This operation, used in end-to-end settings with any architecture, removes undesired extraneous relations between the learned features of a model and the chosen metadata. MDN extends traditional statistical methods to deep learning architectures as a network layer that corrects the distribution shift of model features from metadata, differing from BN and its variants that only standardize the features. Therefore, it can effectively combat bias in deep learning models as well as remove the effects of study confounders in medical studies. Our results on four diverse datasets have shown that MDN is successful at removing the dependence of the learned features of a model on metadata variables while maintaining or improving performance.

Acknowledgements. This study was partially supported by NIH Grants (AA017347, MH113406, and MH098759), Schmidt Futures Gift, and Stanford Institute for Human-Centered AI (HAI) AWS Cloud Credit.

References

- [1] Ehsan Adeli, Dongjin Kwon, Qingyu Zhao, Adolf Pfefferbaum, Natalie M Zahr, Edith V Sullivan, and Kilian M Pohl. Chained regularization for identifying brain patterns specific to hiv infection. *Neuroimage*, 183:425–437, 2018.
- [2] Ehsan Adeli*, Qingyu Zhao*, Adolf Pfefferbaum, Edith V Sullivan, Li Fei-Fei, Juan Carlos Niebles, and Kilian M Pohl. Representation learning with statistical independence to mitigate bias. *arXiv preprint arXiv:1910.03676*, 2019.
- [3] Ehsan Adeli*, Qingyu Zhao*, Natalie M Zahr, Aimee Goldstone, Adolf Pfefferbaum, Edith V Sullivan, and Kilian M Pohl. Deep learning identifies morphological determinants of sex differences in the pre-adolescent brain. *NeuroImage*, 223:117293, 2020.
- [4] Kei Akuzawa, Yusuke Iwasawa, and Yutaka Matsuo. Adversarial invariant feature learning with accuracy constraint for domain generalization. *arXiv preprint arXiv:1904.12543*, 2019.
- [5] Mohsan Alvi, Andrew Zisserman, and Christoffer Nellåker. Turning a blind eye: Explicit removal of biases and variation from deep neural network embeddings. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 0–0, 2018.
- [6] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [7] Yasaman Bahri, Jonathan Kadmon, Jeffrey Pennington, Sam S Schoenholz, Jascha Sohl-Dickstein, and Surya Ganguli. Statistical mechanics of deep learning. *Annual Review of Condensed Matter Physics*, 2020.
- [8] Alexander Basilevsky. Dover, 2005.
- [9] Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Advances in neural information processing systems*, pages 4349–4357, 2016.
- [10] M Alan Brookhart, Til Stürmer, Robert J Glynn, Jeremy Rassen, and Sebastian Schneeweiss. Confounding control in healthcare database research: challenges and potential approaches. *Medical care*, 48(6 0):S114, 2010.
- [11] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91, 2018.
- [12] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6299–6308, 2017.
- [13] Ruth Chadwick. Gender and the human genome. *Mens Sana Monographs*, 7(1):10, 2009.
- [14] Tianyi Chen, Georgios Giannakis, Tao Sun, and Wotao Yin. Lag: Lazily aggregated gradient for communication-efficient distributed learning. In *Advances in Neural Information Processing Systems*, pages 5050–5060, 2018.
- [15] Jinwoo Choi, Chen Gao, Joseph Messou, and Jia-Bin Huang. Why can’t i dance in the mall? learning to mitigate scene bias in action recognition. 12 2019.
- [16] Marta R Costa-jussà. An analysis of gender bias studies in natural language processing. *Nature Machine Intelligence*, pages 1–2, 2019.
- [17] Alexandre Défossez and Francis Bach. Adabatch: Efficient gradient aggregation rules for sequential and parallel stochastic gradient methods. *arXiv preprint arXiv:1711.01761*, 2017.
- [18] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [19] NA Diamantidis, Dimitris Karlis, and Emmanouel A Giakoumakis. Unsupervised stratification of cross-validation for accuracy estimation. *Artificial Intelligence*, 116(1-2):1–16, 2000.
- [20] Ali Diba, Mohsen Fayyaz, Vivek Sharma, Manohar Paluri, Jürgen Gall, Rainer Stiefelhagen, and Luc Van Gool. Large scale holistic video understanding. In *European Conference on Computer Vision*, pages 593–610. Springer, 2020.
- [21] Eran Eidinger, Roeen Enbar, and Tal Hassner. Age and gender estimation of unfiltered faces. *IEEE Transactions on Information Forensics and Security*, 9(12):2170–2179, 2014.
- [22] Lise Eliot. Neurosexism: the myth that men and women have different brains. *Nature*, 566(7745):453–455, 2019.
- [23] Michael Feldman. *Computational fairness: Preventing machine-learned discrimination*. PhD thesis, 2015.
- [24] Thomas B Fitzpatrick. The validity and practicality of sun-reactive skin types i through vi. *Archives of dermatology*, 124(6):869–871, 1988.
- [25] Anna Fry, Thomas J Littlejohns, Cathie Sudlow, Nicola Doherty, Ligia Adamska, Tim Sprosen, Rory Collins, and Naomi E Allen. Comparison of sociodemographic and health-related characteristics of uk biobank participants with those of the general population. *American journal of epidemiology*, 186(9):1026–1034, 2017.
- [26] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *The Journal of Machine Learning Research*, 17(1):2096–2030, 2016.
- [27] Adam Geitgey. Face recognition. *Adam Geitgey*, 3, 2017.
- [28] Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. Learning spatio-temporal features with 3d residual networks for action recognition. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 3154–3160, 2017.
- [29] Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. In *Advances in neural information processing systems*, pages 3315–3323, 2016.
- [30] De-An Huang, Vignesh Ramanathan, Dhruv Kumar Mahajan, Lorenzo Torresani, Manohar Paluri, Li Fei-Fei, and Juan Carlos Niebles. What makes a video a video: Analyzing temporal information in video understanding models and datasets. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7366–7375, 2018.
- [31] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [32] Heinrich Jiang and Ofir Nachum. Identifying and correcting label bias in machine learning. In *International Conference on Artificial Intelligence and Statistics*, pages 702–712,

- 2020.
- [33] James E Johndrow, Kristian Lum, et al. An algorithm for removing sensitive information: application to race-independent recidivism prediction. *The Annals of Applied Statistics*, 13(1):189–220, 2019.
- [34] Amit Kaushal, Russ Altman, and Curt Langlotz. Geographic distribution of us cohorts used to train deep learning algorithms. *Jama*, 324(12):1212–1213, 2020.
- [35] Brendan F Klare, Ben Klein, Emma Taborsky, Austin Blanton, Jordan Cheney, Kristen Allen, Patrick Grother, Alan Mah, and Anil K Jain. Pushing the frontiers of unconstrained face detection and recognition: Iarpa janus benchmark a. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1931–1939, 2015.
- [36] Jason Koh, Dezhi Hong, Rajesh Gupta, Kamin Whitehouse, Hongning Wang, and Yuvraj Agarwal. Plaster: An integration, benchmark, and development framework for metadata normalization methods. In *Proceedings of the 5th Conference on Systems for Built Environments*, pages 1–10, 2018.
- [37] Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, and Richard Zemel. The variational fair autoencoder. *arXiv preprint arXiv:1511.00830*, 2015.
- [38] Qiongmin Ma, Tianhao Zhang, Marcus V Zanetti, Hui Shen, Theodore D Satterthwaite, Daniel H Wolf, Raquel E Gur, Yong Fan, Dewen Hu, Geraldo F Busatto, et al. Classification of multi-site mr images in the presence of heterogeneity using multi-task learning. *NeuroImage: Clinical*, 19:476–486, 2018.
- [39] Dougal Maclaurin, David Duvenaud, and Ryan Adams. Gradient-based hyperparameter optimization through reversible learning. In *International Conference on Machine Learning*, pages 2113–2122, 2015.
- [40] David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. Learning adversarially fair and transferable representations. *arXiv preprint arXiv:1802.06309*, 2018.
- [41] Roseanne McNamee. Regression modelling and other methods to control confounding. *Occupational and environmental medicine*, 62(7):500–506, 2005.
- [42] Daniel Moyer, Shuyang Gao, Rob Brekelmans, Aram Galstyan, and Greg Ver Steeg. Invariant representations without adversarial training. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 9084–9093, 2018.
- [43] J. Neter, M. H. Kutner, C. J. Nachtsheim, and W. Wasserman. volume 4. McGraw-Hill/Irwin, Chicago, 1996.
- [44] Ronald Carl Petersen, PS Aisen, Laurel A Beckett, MC Donohue, AC Gamst, Danielle J Harvey, CR Jack, WJ Jagust, LM Shaw, AW Toga, et al. Alzheimer’s disease neuroimaging initiative (adni): clinical characterization. *Neurology*, 74(3):201–209, 2010.
- [45] Mohamad Amin Pourhoseingholi, Ahmad Reza Baghestani, and Mohsen Vahedi. How to control confounding effects by statistical analysis. *Gastroenterology and hepatology from bed to bench*, 5(2):79, 2012.
- [46] Babak Salimi, Luke Rodriguez, Bill Howe, and Dan Suci. Interventional fairness: Causal database repair for algorithmic fairness. In *Proceedings of the 2019 International Conference on Management of Data*, pages 793–810, 2019.
- [47] Shreya Shankar, Yoni Halpern, Eric Breck, James Atwood, Jimbo Wilson, and D Sculley. No classification without representation: Assessing geodiversity issues in open data sets for the developing world. *arXiv preprint arXiv:1711.08536*, 2017.
- [48] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [49] Saurabh Singh and Shankar Krishnan. Filter response normalization layer: Eliminating batch dependence in the training of deep neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11237–11246, 2020.
- [50] Kraig R Stevenson, Joseph D Coolon, and Patricia J Witkopp. Sources of bias in measures of allele-specific expression derived from rna-seq data aligned to a single reference genome. *BMC genomics*, 14(1):536, 2013.
- [51] Tony Sun, Andrew Gaut, Shirllyn Tang, Yuxin Huang, Mai ElSherief, Jieyu Zhao, Diba Mirza, Elizabeth Belding, Kai-Wei Chang, and William Yang Wang. Mitigating gender bias in natural language processing: Literature review. *arXiv preprint arXiv:1906.08976*, 2019.
- [52] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *International conference on machine learning*, pages 1139–1147. PMLR, 2013.
- [53] Gábor J Székely, Maria L Rizzo, Nail K Bakirov, et al. Measuring and testing dependence by correlation of distances. *The annals of statistics*, 35(6):2769–2794, 2007.
- [54] Zilong Tan, Samuel Yeom, Matt Fredrikson, and Ameet Talwalkar. Learning fair representations for kernel models. In *International Conference on Artificial Intelligence and Statistics*, pages 155–166. PMLR, 2020.
- [55] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7167–7176, 2017.
- [56] Christian Wachinger, Anna Rieckmann, and Sebastian Pölsterl. Detect and correct bias in multi-site neuroimaging datasets. *arXiv preprint arXiv:2002.05049*, 2020.
- [57] Haohan Wang, Zhenglin Wu, and Eric P Xing. Removing confounding factors associated weights in deep neural networks improves the prediction accuracy for healthcare applications. In *Pac Symp Biocomput.*, pages 54–65. World Scientific, 2019.
- [58] Tianlu Wang, Jieyu Zhao, Mark Yatskar, Kai-Wei Chang, and Vicente Ordonez. Balanced datasets are not enough: Estimating and mitigating gender bias in deep image representations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5310–5319, 2019.
- [59] Merriam Webster. *Metadata* <https://www.merriam-webster.com/dictionary/metadata>, 2020 (accessed Nov 14, 2020).
- [60] Blake Woodworth, Suriya Gunasekar, Mesrob I Ohannesian, and Nathan Srebro. Learning non-discriminatory predictors. *arXiv preprint arXiv:1702.06081*, 2017.
- [61] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European conference on computer vision (ECCV)*, pages 3–19, 2018.

- [62] Qizhe Xie, Zihang Dai, Yulun Du, Eduard Hovy, and Graham Neubig. Controllable invariance through adversarial feature learning. In *Advances in Neural Information Processing Systems*, pages 585–596, 2017.
- [63] Kaiyu Yang, Klint Qinami, Li Fei-Fei, Jia Deng, and Olga Russakovsky. Towards fairer datasets: Filtering and balancing the distribution of the people subtree in the imagenet hierarchy. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 547–558, 2020.
- [64] Bianca Zadrozny. Learning and evaluating classifiers under sample selection bias. In *Proceedings of the twenty-first international conference on Machine learning*, page 114, 2004.
- [65] Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. Learning fair representations. In *International Conference on Machine Learning*, pages 325–333, 2013.
- [66] Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. Mitigating unwanted biases with adversarial learning. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, pages 335–340, 2018.
- [67] Yong Zhang, Dongjin Kwon, Pardis Esmacili-Firidouni, Adolf Pfefferbaum, Edith V Sullivan, Harold Javitz, Victor Valcour, and Kilian M Pohl. Extracting patterns of morphometry distinguishing hiv associated neurodegeneration from mild cognitive impairment via group cardinality constrained classification. *Human brain mapping*, 37(12):4523–4538, 2016.
- [68] Qingyu Zhao*, Ehsan Adeli*, Adolf Pfefferbaum, Edith V Sullivan, and Kilian M Pohl. Confounder-aware visualization of convnets. *arXiv preprint arXiv:1907.12727*, 2019.
- [69] Qingyu Zhao*, Ehsan Adeli*, and Kilian M Pohl. Training confounder-free deep learning models for medical applications. *Nature Communications*, 2020 In Press.