

Interventions — November 2

Prof. Eric Wong

In this section we'll discuss interventions—ways in which you can adjust or edit your data or models to try to fix some of these problems.

- How can you intervene at the data level (balancing, source selection)?
- How can you intervene at the model level (editing, fine-tuning)

1 Data interventions

First, we'll look at how you can fix problems at the data level. One simple approach recently found to be competitive with state of the art is known as *data balancing*. Data balancing consists of rebalancing or reweighting your data to fix misbalanced proportions. One of the classic use-cases here is for worst-group performance (the same setting as group DRO from the robust learning module).

Data balancing. How do you balance the data? There are two main axes by which balancing can occur:

- What subsets do you balance? One could balance by classes, which does not require extra attribute labels. However, classes may not reflect inequalities between groups. Alternatively, one could balance by groups, but this requires group attribute information.
- How do you balance? One could perform a random subset to a fixed size so that all groups/classes have the same size. Alternatively, one could reweight the subsets according to their sizes so that all groups/classes have the same weight during training.

The difference between weighting classes and groups is clear—one is directly related to worst-group performance but requires more data annotations, while the other is only tangentially related but can be done for free. However, what is the difference between subsetting and reweighting?

To see this, consider the following example from Sagawa et al. (2020). Let $y \in \{-1, 1\}$ be the label and $a \in \{-1, 1\}$ is a spurious attribute. Then, generate some features $x = \begin{bmatrix} x_s \\ x_c \\ x_n \end{bmatrix} \in \mathbb{R}^{d+2}$ as follows:

- $p(x_s|y, a) = \mathcal{N}(a, \sigma_s^2) \in \mathbb{R}$, which is a spurious feature that has no effect on the true label y
- $p(x_c|y, a) = \mathcal{N}(y, \sigma_c^2) \in \mathbb{R}$, which is the core feature that is correlated with the label.

- $p(x_n|y, a) = \mathcal{N}(0, \sigma_n^2) \in \mathbb{R}^d$, which are noise features that are not correlated with anything.

Here, (y, a) determines a group—we can take $y = a$ to be the majority groups, and $y = -a$ to be the minority groups. The different variances control the amount of noise in each feature. if we take σ_c to be large and σ_s, σ_n to be small, then this encourages the use of spurious and noise features before the core feature. Overparameterized models can thus memorize individual examples with x_n , and otherwise use x_s or x_c to classify larger groups.

1. A vanilla ERM model will first rely on the spurious feature which separates the majority groups, and then memorizes the minority groups with the noise vectors. This results in poor accuracy on the minority groups.
2. Subsampling the majority groups decorrelates the spurious feature with the class label. The model thus does not prioritize the majority group, and uses the core feature as the most correlated feature first.
3. Reweighting can also have a similar effect to subsampling. However, repeated training on minority examples encourages their memorization with the noise features.

What exactly does it mean for this model to memorize a datapoint? We can decompose the weights on the noise features via the representer theorem as $w_n = \sum_i \alpha^{(i)} x_n^{(i)}$ and say that the model memorizes $x^{(i)}$ if $\alpha^{(i)}$ has large weight. Intuitively, since the noise vector is high dimensional, all noise vectors of different examples are nearly orthogonal. Thus, $w_n x_n^{(i)}$ will affect the prediction on $x^{(i)}$ but not other training or test points.

Why does a model favor the spurious feature over the core feature for majority groups? A separator that depends only on the spurious feature has a norm that scales only with the number of minority data points, since it only needs to memorize the minority groups. In contrast, a separator that depends on the core features has noisier examples, so a model that uses x_c will still need to memorize some amount of all points including some in the majority group. It can be shown that this results in a larger-norm separator (since it needs to memorize more points). Since inductive bias favors minimum-norm separators, models will use the spurious feature first.

There is a nice theorem from Sagawa et al. (2020) that proves that in the overparameterized regime (i.e. for large d), there exist parameters for this setting such that the error of a max margin classifier is lower bounded by $2/3$. On the other hand, for the underparameterized regime (i.e. for $d = 0$) the error is upper bounded.

Unadversarial examples If you can design the objects that your system needs to detect, then you can use adversarial examples techniques for non-nefarious purposes. Specifically, you can run an “adversarial attack” on your object to maximize the correct label, and then modify the object accordingly.

Generative modeling One other way to intervene on your data is to pass them through a generative model. For example, one can use generative models to convert data from one group to another group, or to simply generate more examples from a minority group. There are several things to watch out for here:

1. The generator needs to be diverse, and not suffer from mode collapse. This was a big problem for GANs, but may not be as big of an issue today.
2. Generative models typically needs lots of data. This is often not available for minority groups by definition. Often you can only use this approach when you already have sufficiently large data.
3. Generative models need to be high quality to be useful in this case. Otherwise, noisy or fuzzy generative models can hurt more than they help.

But if you can use a good quality generative model, then they can data augment your existing datasets to sometimes lead to big improvements.

2 References

Idrissi, Badr Youbi, et al. "Simple data balancing achieves competitive worst-group-accuracy." Conference on Causal Learning and Reasoning. PMLR, 2022.

Sagawa, Shiori, et al. "An investigation of why overparameterization exacerbates spurious correlations." International Conference on Machine Learning. PMLR, 2020.