

Scaling up AI Model Deployments with the Limited and Imperfectly Labeled Data

April 13, 2021 Denis Gudovskiy <u>denis.gudovskiy@us.panasonic.com</u> Panasonic Al Lab, Mountain View

Panasonic Al

Panasonic Business Fields



ADAS system for Car OEM Autonomous Commuter



Security



Air Conditioner Cold Chain



Imaging







Autonomous Ride Share Service



AI Development Process



Large scale multimodal dataset for living space

Real-time detection for self-driving mobility service on an embedded device



Home Action Genome Data Example

Home Action Genome (HOMAGE) [CVPR2021]: A large-scale multi-view video database of indoor daily activities

https://homeactiongenome.org/

This is a collaboration work with Stanford Vision and Learning Lab





Home Action Genome Data Example

HOMAGE encompasses the indoor actions and multimodal data



HOMAGE is a guest task of ActivityNet Challenge at CVPR2021 workshop



http://activity-net.org/challenges/2021/

HOMAGE workshop has been accepted to ICCV2021: http://campworkshop.org

Home Action Genome Data Example



Hic Sunt Dracones et AI - Here be Dragons and AI

The explored world – train data Scary dragons live in the test data



Hic Sunt Dracones et AI - Here be Dragons and AI

AI models minimize risk for the empirical train datasets... and empirical datasets are always limited, biased and noisy...



Hic Sunt Dracones et AI - Here be Dragons and AI

...but it is our job to make them better and, hence, our models more robust. Therefore, I encourage researchers to go beyond standard dataset setups!



Zoo of Relevant Methods and Industrial Applications

• Scaling AI applications to new domains is tough:

- $\,\circ\,$ Labeling is noisy and expensive
- $\,\circ\,$ Sometimes, labeling is not even possible
- $\,\circ\,$ Models are not robust to unseen/rare test data
- Safety-critical apps are the most vulnerable and leads to geo-fencing, remote operators etc.
- $\,\circ\,$ ML research community tries to catch up with:
 - Un/self/semi-supervised learning/pretraining
 - $\,\circ\,$ Active and zero/few-shot learning-assisted labeling
 - $\,\circ\,$ Augmentation methods optimize labeled datasets
 - $\circ~$ Domain adaptation: train-time and test-time
 - Generative modeling: syn2real, image-to-image translation etc.
- $\,\circ\,$ However, such methods usually are limited to:
 - Simple tasks: classification vs. 2D/3D object detection/tracking/segmentation, scene parsing
 - $\,\circ\,$ Simple data domains: sparse unimodal vs. dense multimodal
 - Standard public datasets: balanced train/test vs. biased data with label noise





Deep Active Learning for Biased Datasets via Fisher Kernel Self-Supervision

Denis Gudovskiy¹ Alec Hodgkinson¹ Takuya Yamaguchi² Sotaro Tsukizawa²

¹Panasonic AI Lab, Mountain View, CA ²Panasonic AI Solutions Center, Osaka Code: <u>github.com/gudovskiy/al-fk-self-supervision</u> Email: <u>denis.gudovskiy@us.panasonic.com</u>

Panasonic Al

Motivation

ODNNs require large annotated train datasets

Annotation process is costly and do not guarantee accuracy improvement





 \circ Challenges:

- How to select train examples for annotation?
- How to increase accuracy while minimizing annotation costs?
- How to utilize all unlabeled data?

Active Learning

• Active learning (AL) aims to select only relevant train examples for annotation





 \odot How active learning works:

- \checkmark A **representation** *z* is pooled from the task model
- ✓ AL acquisition function $\mathcal{R}(z)$ is calculated
- ✓ A b^{th} batch \mathbb{N}^{b} of examples is selected for annotation using $\mathcal{R}(z)$
- ✓ Steps are repeated upon reaching target accuracy

Problem Statement for Biased Datasets

• **Prior methods** assume: $Q_x^{\text{test}} = Q_x$ and only train data is accessed • Then, a trained classifier **misses on underrepresented** test instances



Case #1: autonomous vehicle in a rare traffic situation
 ◆Vehicle may have an accident with unfamiliar objects
 Case #2: face recognition with gender and race biases
 ◆Photo cannot be recognized for very rare types of faces
 Our idea: explicitly minimize val/train distribution shift R_{opt} = arg min D_{KL}(Q̂_x^v || Q̂_x)

Our AL Method via Self-Supervised Fisher Kernel

• Main approach: density matching in feature space using Fisher kernel



З

 \circ Initially, we pretrain classifier using unsupervised learning [1]

○ Next, our algorithm repeats the following steps:

- 1) Pool compact features $m{z}$ and gradients $m{g}$ for val/unlabeled data
- 2) Estimate practical Fisher kernel (PFK) between validation and unlabeled data

3) Annotate and add selected examples to train data that maximize PFK References:

[1] Gidaris et al. <u>Unsupervised representation learning by predicting image rotations</u>. In ICLR18

Details of Self-Supervised Fisher Kernel



• We pool **multi-scale** features $z_i \in \mathbb{R}^L$ and gradients $g_i \in \mathbb{R}^L$ from the task DNN • Pseudo-labels for gradients are **estimated** using $S = \hat{p}(y|z)$ metric • Practical Fisher kernel is a **tractable** similarity matrix for DNNs:

$$\boldsymbol{R}_{\boldsymbol{z},\boldsymbol{g}} = \boldsymbol{R}_{\boldsymbol{z}} \circ \boldsymbol{R}_{\boldsymbol{g}} = \sum_{j} \left(\left(\boldsymbol{Z}_{v}^{j} \right)^{T} \boldsymbol{Z}^{j} \circ \left(\boldsymbol{G}_{v}^{j} \right)^{T} \boldsymbol{G}^{j} \right)$$

Examples similar to clustered misclassified validation data are added to train dataset

Complexity Estimates

 \odot Comparison in terms of forward and backward DNN passes

- > M size of validation dataset
- > N size of all unlabeled train data
- $> N^{b}$ current labeled train data size
- $\gg \acute{N}^{b}$ current unlabeled data size
- $\succ K$ number of stochastic samples
- \geq *I* number of train epochs
- $\succ E$ number of ensembles

Method	AL	Train
Uncert. (varR) [2]	$K \acute{N}^b$	$2IN^b$
Ensembl. uncert. [3]	$EK\acute{N^b}$	$2EIN^{b}$
VAAL [4]	$\acute{N}^b + 2NI_{\mathrm{VAE,D}}$	$2IN^b$
PCC: R_{z}	$M + \acute{N^b}$	$2IN^b$
PFK: $oldsymbol{R}_{oldsymbol{z},oldsymbol{g}}$ (ours)	$2(M + \acute{N^b})$	$2IN^b$

• Our AL speedup compared to [2,3]
$$\frac{EK\dot{N}^b}{2(M+\dot{N}^b)} \approx \frac{EK}{2}$$
, since $\dot{M}^b \gg M$

 \odot Typically, complexity of our method is at least 10× lower

References:

[2] Gal et al. <u>Deep Bayesian active learning with image data</u>. In ICML17

[3] Beluch et al. <u>The power of ensembles for active learning in image classification</u>. In CVPR18

[4] Sinha et al. <u>Variational adversarial active learning</u>. In ICCV19

ImageNet w/o Class Imbalance

- ResNet-18 model
- $\,\circ\,$ Practically unrealistic data setup
- Our AL method with pseudo-labels:
 - \circ 1.5% accuracy \uparrow
 - \circ 20% labels \checkmark
 - $_{\odot}$ 16× speedup \uparrow
- \odot Theoretical limit with true labels:
 - 7% accuracy ↑
 70% labels ↓



ImageNet with 100× Class Imbalance

- \circ ResNet-18 model
- \odot Practical data setup case
- Imbalance = {500 random classes}/{500 other classes} images
- Our AL method with pseudo-labels:
 2% accuracy ↑
 42% labels ↓
 16× speedup ↑
- \odot Theoretical limit with true labels:
 - \circ 6% accuracy \uparrow
 - \circ 90% labels \checkmark



SVHN with 100× Imbalance

- \circ ResNet-10 DNN model
- \circ Imbalance = {0...4}/{5...9} images
- \circ Our AL method with pseudo-labels:
 - $_{\odot}$ 10% accuracy \uparrow
 - \circ 40% labels \downarrow
 - \circ 32× speedup \uparrow
- \odot Theoretical limit with true labels:
 - \circ 17% accuracy \uparrow
 - \circ 80% labels \downarrow



Confusion Matrices for 100× Imbalanced MNIST

Imbalance = $\{0...4\}/\{5...9\}$ images = $100 \times$ at AL iteration b = 3

a) Prior varR (uncertainty-based) method: **36%** accuracy for digits {5,8,9}

- b) Our method with estimated pseudolabels: 75% accuracy for digits {5,8,9}
- c) Our method with all true labels: 89% theoretical limit for PFK



T-SNE Clustering: MNIST 100× Imbalance

Imbalance = $\{0...4\}/\{5...9\}$ images = $100 \times$ at AL iteration b = 3Balls are **misclassified** and dots are **correctly** classified

- a) Prior varR [2] method: **36%** accuracy for digits {5,8,9}
- b) Our method with estimated pseudolabels: **75%** accuracy for digits {5,8,9}
- c) Our method with all true labels: 89% theoretical limit for PFK





To appear in CVPR2021



AutoDO: Robust AutoAugment for Biased Data with Label Noise via Scalable Probabilistic Implicit Differentiation

Denis Gudovskiy¹ Luca Rigazio¹ Shun Ishizaka² Kazuki Kozuka² Sotaro Tsukizawa²

¹Panasonic AI Lab, Mountain View, CA ²Panasonic Technology Division, Osaka Code: <u>github.com/gudovskiy/autodo</u> Email: <u>denis.gudovskiy@us.panasonic.com</u>

Panasonic Al

Labeled Data Optimization

After AL-assisted labeling, efficiency of labeled data can be improved
 Data pre-processing can virtually increase train dataset and improve generalization



 \odot Labeled data optimization steps:

✓ Select **pre-processing functions** (*e.g.* augmentations) to increase data variability

- ✓ Set processing hyperparameters or learn them automatically (AutoAugment/AutoDO)
- ✓ Train AI model with the virtually expanded train dataset to improve generalization
- ✓ Continue optimizing data processing functions and their hyperparameters

Labeled Data Optimization

Typical data pre-processing functions: augmentations, noise, mixup, CutMix [1]
 Manual selection is time-consuming and requires domain knowledge

• Recent methods aim to search or learn pre-processing automatically: AutoAugment [2]

Learnable pre-processing can be achieved using automatic hyperparameter optimization

Image	ResNet-50	Mixup [48]	Cutout [3]	CutMix
Label	Dog 1.0	Dog 0.5 Cat 0.5	Dog 1.0	Dog 0.6 Cat 0.4
ImageNet	76.3	77.4	77.1	78.6
Cls (%)	(+0.0)	(+1.1)	(+0.8)	(+2.3)
ImageNet	46.3	45.8	46.7	47.3
Loc (%)	(+0.0)	(-0.5)	(+0.4)	(+1.0)
Pascal VOC	75.6	73.9	75.1	76.7
Det (mAP)	(+0.0)	(-1.7)	(-0.5)	(+1.1)



References:

[1] Yun et al. CutMix: Regularization Strategy to Train Strong Classifiers. In ICCV19

[2] Lim et al. <u>Fast AutoAugment</u> In NeurIPS19

Automatic Dataset Optimization (AutoDO)



[1] Gudovskiy et al. <u>AutoDO: robust AutoAugment for biased data with label noise</u>. To appear in CVPR21
 [2] Lorraine et al. <u>Optimizing millions of hyperparameters by implicit differentiation</u>. In AISTATS20

Automatic Dataset Optimization (AutoDO)

Main focus: optimization of realistic train datasets with data biases and noisy labels
 Train data distribution should be adjusted to test data in actual systems
 AutoDO can flexibly change train distribution due to per-point model

 \odot This prevents overfitting to train data and increases generalization!

The first construction of the first construction for the first construction of the first constr						
Alg./IR-NR	1-0.0	100-0.0	1-0.1	100-0.1		
Baseline	$3.6_{\pm 0.10}$	13.6 ± 0.69	5.3 ± 0.27	20.0±1.92		
RAA [6]	$2.7{\scriptstyle\pm0.04}$	$10.9{\scriptstyle \pm 0.66}$	$3.4_{\pm 0.11}$	13.6 ± 0.96		
FAA [17]	2.8 ± 0.02	$11.5{\scriptstyle\pm0.32}$	$3.7{\scriptstyle\pm0.08}$	15.3 ± 1.07		
DADA [16]	$2.9{\scriptstyle \pm 0.03}$	$12.2{\scriptstyle\pm0.54}$	$4.1{\scriptstyle \pm 0.13}$	$16.5{\scriptstyle \pm 1.51}$		
$oldsymbol{\lambda}^{A_{ ext{SHA}}}$ (ours)	$2.8_{\pm 0.10}$	12.6±1.53	$3.0_{\pm 0.17}$	$13.7_{\pm 0.77}$		
$\boldsymbol{\lambda}^{A}$ (ours)	$2.7 \scriptstyle \pm 0.09$	10.2 ± 0.50	3.0 ± 0.07	$12.3{\scriptstyle \pm 0.80}$		
$oldsymbol{\lambda}^{A,W}$ (ours)	2.8 ± 0.04	6.1 ± 0.22	2.8 ± 0.07	$8.1{\scriptstyle \pm 0.14}$		
$\boldsymbol{\lambda}^{A,W,S}$ (ours)	$2.5{\scriptstyle \pm 0.04}$	$5.3{\scriptstyle \pm 0.21}$	$2.6{\scriptstyle \pm 0.05}$	6.3 ± 0.57		

Table 1 WRNet28-10 SVHN top-1 test error rate $\mu_{\perp} = -\infty$

References:

[6] Cubuk et al. <u>RandAugment</u>. In CVPRW19
[17] Lim et al. <u>Fast AutoAugment</u>. In NeurIPS19
[16] Li et al. <u>DADA</u>. In ECCV20



Automatic Dataset Optimization (AutoDO)

 \circ WideResNet28-10 model

 \circ SVHN dataset

- 100× class imbalance
- 10% noisy labels
- \circ Confusion matrices (top)
- \circ T-SNE of embeddings (bottom)
 - a) Standard augmentations
 - b) Fast AutoAugment
 - c) AutoDO
- AutoDO improves generalization when train dataset is biased, and labels are noisy



Conclusions

○ I presented recent Panasonic AI-related projects:

- ✓ Autonomous ride share service (big thanks to BDD)
- \checkmark Our new challenging multimodal Home Action Genome Dataset
- ✓ Recent research about robust active learning and AutoAgument
- I encourage research community to:
 - ✓ Consider a realistic imperfect data collection and labeling process
 - \checkmark Introduce some data bias when develop your models
 - ✓ Research new methods to overcome limitations of empirical datasets

• Main takeaways:

- ✓ Methods to deal with the imperfect data not only save corp money...
- ✓ They are good for publications¹ and discovering all unknown yet dragons!



Questions?

Panasonic Al