

Quality-Diversity Generative Sampling for Learning with Synthetic Data Allen Chang¹, Matthew C. Fontaine¹, Serena Booth², Maja J. Matarić¹, Stefanos Nikolaidis¹ ¹ University of Southern California² Massachusetts Institute of Technology

Introduction

Synthetic datasets can supplement real datasets to improve training.

However, synthetic training datasets should be balanced to prevent a transfer of bias to a downstream task.

Biased Real Data



Biased Synthetic Data

Contributions:

 \rightarrow Use quality-diversity algorithms to simultaneously balance data features while creating synthetic training datasets

 \rightarrow Improve fairness without harming accuracy when the data are used to pretrain downstream facial recognition classifiers

Results (1) More Uniformity in QDGS Distributions



...achieve more uniformity than random sampling. **QDGS** datasets



(2) > Fairness Increase, and (3) \approx Accuracy Increase

Density Map 30	Pretraining	Dark-skinned (RFW)	Light-skinned (RFW)	ACC (LFW AgeDB)*
	Train: CASIA	(Imbalanced)	8	8 /
15	None	88.08 ± 0.07	94.05 ± 0.06	94.52 ± 0.03
	Rand15	88.17 ± 0.08	94.41 ± 0.07	94.83 ± 0.03
	Rand50	88.69 ± 0.10	94.67 ± 0.07	94.89 ± 0.03
0.0 0.1	QD15 (Ours)	88.25 ± 0.09	94.42 ± 0.06	94.83 ± 0.03
on.	QD50 (Ours)	88.94 ± 0.07	94.62 ± 0.10	94.92 ± 0.02
Density Map 30	Train: BUPT (Balanced)		
	None	96.22 ± 0.03	97.58 ± 0.05	96.27 ± 0.02
15	Rand15	96.22 ± 0.05	97.72 ± 0.05	96.35 ± 0.02
	Rand50	96.17 ± 0.09	97.82 ± 0.04	96.31 ± 0.03
	QD15 (Ours)	96.29 ± 0.06	97.75 ± 0.03	96.32 ± 0.03
	QD50 (Ours)	96.32 ± 0.06	97.80 ± 0.04	96.30 ± 0.02
on.				

c is the averaged accuracy of LFW, CFPFP, CPLFW, CALFW, and AgeDB

... are helpful for balanced and imbalanced training sets.

Prior Work:

Single-factor diversity optimization

QDGS:

Quality-diversity (QD) optimization for multiple measures of diversity optimization + text prompts to describe desired attributes

5. Optimize via objective

 $oldsymbol{c} \sim \mathcal{N}(oldsymbol{\mu}, \, \Sigma)$ 6. Adapt and sample Gaussian

 $oldsymbol{ heta} + oldsymbol{c} oldsymbol{
abla}$

Conclusions

Synthetic datasets created from generative models can be helpful to train classifiers, especially when they are balanced.

QDGS uses QD optimization and text prompts to control for which attributes should be made more balanced in these datasets.

QDGS + facial recognition classifiers \rightarrow more fairness without decrease in accuracy!









