Automated Deep Learning Self-Service

Michael VACCARO (<u>michavaccaro@gmail.com</u>) 23/01/2021

Automated Deep Learning









- I. About the AutoDL challenge series
- II. The AutoDL Self-Service
- III. How does Deep Wisdom (winner) work?
- IV. Benchmarking the AutoDL solutions

I. About the AutoDL challenge series

AutoDL, a series of challenges

- Objective: provide "universal learning machines", which can learn and predict without any human intervention.
- Classification tasks: from image, video, speech, text or tabular domains, formatted in an uniform way.
- Series of challenges in 2019-2020 : AutoCV, AutoCV2, AutoNLP, AutoSpeech, AutoWSL and the final **AutoDL [1]** (combining all types of data).



AutoDL challenge design

- Five datasets from all domains in the final phase
- Time limit: 20min to initialize the model and 20min to run.
- Limited resources
- Baselines provided: Baseline0 (constant model), Baseline1 (Linear model), Baseline2 (CNN-based solution) and Baseline3 (combining winning solutions from previous domain-specific challenges).



AutoDL challenge design

• Metric : Area under the Learning Curve (ALC).



Each prediction is scored with the Normalized ROC AUC metric (NAUC).

NAUC = 2 AUC - 1

II. The AutoDL Self-Service

A Codalab* platform

The link: https://competitions.codalab.org/competitions/27082

Automated Deep Learning	Automated Deep Learnin Organized by michavaccaro - Current serve	Self-Service ne: Feb. 9, 2021, 1:20 p.m. UTC		
Self-Service	Current DeepWisdom	End Competition Ends		
	Dec. 1, 2020, 8 p.m. UTC	Never		
Learn the Details	Phases Participate Results			

Through this "Inverted" competition:

- Upload YOUR dataset
- Find automatically and train a deep learning model on your task
- Rate it on your test set
- Predict on unseen data afterwards

What is behind?

- Model used: *DeepWisdom**, winner of AutoDL** challenge
- Submissions run for 20 minutes on GPUs
- Metric used to "rate" submissions: ALC (based on predictions curve, scored by ROC

NAUC)



Each prediction is scored with the Normalized ROC AUC metric (NAUC).

NAUC = 2 AUC - 1



• Datasets can be **any classification task** formatted as 4D tensors (time, row, col,

channels) in a generic TFRecords format.

											• • • • • • •
٦	Га	b	ıla	ır —	→ (1, 1,	10	, 1)		Video \rightarrow (-1, 200, 300, 3)	Image → (1, 28, 28, 1)
13 50) 4(6	1	50	192485	5	1	1887	6		0123430789
1 40	2	9	1	38	159449	1 Na	N	0	4		1122151200
8 40	3	6	1	35	154641	1	1	0	4		0 (2 5 4 5 6 1 6 4
5 25	2	3	1	32	178649	1	1	0	1	III WE LEA WE LEA III III III III III III III III III I	1177151790
9 40	2	1	1	75	211013	1	1	0	2		0123700/07
9 40	3	1	1	17	386120	1	2	0	4		0122456786
6 50		0	1	26	16/200	5	1	0	1		0123430184
6 27	10	8	1	43	102074	1	3	0	4		11771111700
2 45	5	1	1	59	154950	1	1	0	6		01001010101
5 58	3 28	8	21	23	289293	1	1	0	11		A12345678A
5 50	30	0	1	67	212490	1	1	0	1		0123130181
6 60	48	8	1	59	73289	2	1	0	14		1177151700

• <u>Github: provided to convert your raw data into the right format.</u>

How to try?

- Log-in to Codalab and register to the "competition"
- Read the overview
- Try the interface with the starting kit:

Learn the Details	Phases	Participate Results					
Overview		Starting Kit					
Starting Kit		We provide one dataset of each domain: Tabular, Image, Video, Time Series and Text.					
Format my submission		Try to upload one now!					
Making predictions		Mini Starting Kit (1 dataset)	[Download]	128 MB			
About the Score		Starting Kit (5 datasets)	[Download]	2.7 GB			
Terms and Conditions		DeepWisdom's code	[Link]	(Github)			
		Format your own dataset	[Link]	(Github)			
Credits							

How to try?

- Format your own data in the AutoDL format
- Submit!

Automated Deep Learning	Automated Deep	Current server time: Dec. 7, 2020, 3:39 p.m. UTC			
Self-Service	Current DeepWisdom Dec. 1, 2020, 8 p.m. UTC	End Competition Ends Never			
Learn the Details	Phases Participate	Results			
Submit / View Resu	Under DeepWise	lom			
	Phase d Submit da	escription atasets to Deep Wisdom			
	Max subr Max subr	nissions per day: 5 nissions total: 100			
	Click the	Submit button to upload a new submission.			
	Optiona	ally add more information about this submission			
	Submit				

Get the results

Your submission







0.0012304 0.00041728 0.000270869 0.00016897 3.00002258 0.00022784 0.00022696 0.9002564 0.0001697 0.002327 0.000282824 0.0007167 0.0101102 0.040168 0.001494 0.010162 0.011541 0.000267 0.00078987 0.960258 0.00071612 0.0041503 0.002823 0.00059116 0.00026824 0.00078754 0.00051477 7.0005000 0.0007895 0.960258 0.00071637 0.001253 0.001161 0.000258 0.00026916 0.0002695 0.0005755 0.00050841 0.002875 0.0005147 0.0005120 0.001253 0.001111 0.001253 0.0001261 0.0002647 0.0005157 0.00050495 0.0015649 0.00056247 0.0002550 0.001577 0.0005130 0.00027650 0.0005757 0.005134 0.000504975 0.96134 0.00056247 0.00056247 0.0005250 0.001567 0.0005230 0.001567 0.001563 0.0015675

Training and testing



• Advantage : self-service (no ML knowledge needed, no computational resources needed)

Next:

- Getting users
- Adding other models
- Benchmarking?

III. How does Deep Wisdom work?



A multi-modal solution adapted to every classification tasks (single-label or multilabel)



Meta-learning framework to learn domain-specific workflows

Use offline datasets (public datasets...)

Image domain

Transfer learning:

- In the early learning phases (t<0.2): Resnet-18 (architecture already used in **Baseline3**)
- Resnet-9 then.



The resnets are pre-trained on **ImageNet [2] dataset**. Batch Normalization and Bias parameters are reinitialized.

Fast autoaugmentation [3] is applied on later training phases to improve last AUCs.

Image domain



Higher NAUC



Transfer learning:

MC3 [4] (mixed convolutions) instead of Resnet-18 (used in **Baseline3**) pretrained on Kinetics.



First two layers are freezed. Bias and linear weights are reinitialized.

Frame extraction: extracting few frames to accelerate predictions (3D Conv being slower).

Weighting ensemble strategy: predict on 3-,10-,12-frames extracted data with MC3 and combine them.

Video domain



21

Speech domain

Resulting of Meta-training on offline datasets.

Model selection: Logistic Regression and ThinResNet34 [5] pretrained on VoxCeleb2 [6].

First layers are freezed.

Workflow optimization:

- Multilabel/unilabel specificity
- Softmax \rightarrow Sigmoid for last activation, different losses
- Skip validation in the beginning (no impact on ALC)

Best-last-predictions ensemble strategy

Module	Input Spectrogram $(257 \times T \times 1)$	Output Size
Thin ResNet	conv2d, 7 × 7, 64	$257 \times T \times 64$
	max pool, 2×2 , stride $(2, 2)$	$128 \times T/2 \times 64$
	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$128 \times T/2 \times 96$
	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$64 \times T/4 \times 128$
	$\begin{array}{c} {\rm conv}, 1 \times 1, 128 \\ {\rm conv}, 3 \times 3, 128 \\ {\rm conv}, 1 \times 1, 256 \end{array} \times 3$	$32 \times T/8 \times 256$
	$\begin{array}{c} {\rm conv}, 1\times 1, 256\\ {\rm conv}, 3\times 3, 256\\ {\rm conv}, 1\times 1, 512 \end{array} \times 3 \end{array}$	$16 \times T/16 \times 512$
	max pool, 3×1 , stride $(2, 2)$	$7 \times T/32 \times 512$
	conv2d, 7 × 1, 512	$1 \times T/32 \times 512$

Text domain

Various preprocessing methods: word frequency filtering, word segmentation, sentences truncating.

Features engineering :

- word / char level features, tokenization with metadata as guideline.
- word embedding : Fast text embedding [7] (pre-trained or reinitialized)

Models :

- include TextCNN, RCNN, GRU, GRU with attention and SVM.
- using datasets metadata to find training parameters

Weighted ensembling : returns weighted top 20 models (for NAUC).

Tabular domain

Use optimized **boosting frameworks** such as: LightGBM, XGBoost, CatBoost and DNN.

A weighted ensemble strategy is applied:

The train set $D = (x_1, y_1, ..., x_n, y_n)$ is divided in **K**-folds, which are themselves divided in 3 sets (one for training, one for hyperparameter search and one for weights computing) :

 $\{(D_{train}^{(1)}, D_{valid1}^{(1)}, D_{valid2}^{(1)}), ..., (D_{train}^{(k)}, D_{valid1}^{(k)}, D_{valid2}^{(k)})\}$

Ensemble learning strategy

Ensemble learning : a stacking/blending strategy

Algorithm 2: AutoEnsemble **Require:** F: Model Space, D: Data Space, M: metric **Require:** k: the number of data folds, m: the number of models, stratified splits D into k folds $\{(D_{train}^{(1)}, D_{valid1}^{(1)}, D_{valid2}^{(1)}), ..., (D_{train}^{(k)}, D_{valid1}^{(k)}, D_{valid2}^{(k)})\}$ for $i \in [1, k]$ do for $j \in [1, m]$ do $f_{\lambda^*}^j(\theta^*) = \operatorname*{arg\,min}_{\theta,\lambda}(L(f_{\lambda}^j(\theta), D_{train}^{(i)}, D_{valid1}^{(i)}))$ $score_{ij} = M(f_{\lambda^*}^j(D_{valid2}|\theta^*))$ $preds_{ij} = f_{\lambda*}^{j}(D_{test}|\theta^{*})$ end for end for $w_{ij} = rac{score_{ij}}{\sum\limits_{i=1}^{k}\sum\limits_{j=1}^{m}(score_{ij})}$ $preds = \sum_{k=1}^{k} \sum_{j=1}^{m} (w_{ij} * preds_{ij})$



Modality	Non-neural	Neural	Pretrained	Generalization/ Acceleration
Image	_	ResNet9 ResNet18	ImageNet	re-initializing
Video	-	MC3	Kinetics	re-initializing freezing
Speech	LR	ThinResnet34	VoxCeleb2	freezing
Text	SVM	TextCNN/RCNN GRU/GRUA	-	-
Tabular	LightGBM/Xgboost Catboost	DNN	-	_

IV. Benchmarking the AutoDL solutions



13 models

- 9 participants (DeepWisdom, DeepBlueAI, PasaNJU, AutoMLFreiburg, Inspur_AutoDL, Frozenmad, TeamZhaw, Surromind, Kon)
- 4 baselines (Baseline1: Linear NN, Baseline2: 3DCNN, Baseline3: All winner solution, Pre-trained Inception)

66 datasets

- 17 image datasets
- 10 video datasets
- 16 text datasets
- 16 time series datasets
- 7 tabular datasets

Deep Wisdom results



Worst datasets:

- Ideal (avg ALC: 0.7737)
- Kitsune (avg ALC: 0.2093)
- Ray (avg ALC: 0.2555)
- Viktor (avg ALC: 0.2763)

Deep Wisdom v. Baseline3





Time comparison with Baseline3

	Image	Video	Time	Text	Tabular	Overall
Average duration Deep Wisdom (s)	571.1211	937.5797	900.0690	577.6960	870.2734	742.2048
Average duration Baseline3 (s)	461.5072	488.4777	1067.8111	195.9854	305.7676	536.8543
Average first prediction duration Deep Wisdom (s)	10.2858	17.2410	12.0937	9.4392	0.8963	10.5709
Average first prediction duration Baseline3 (s)	11.9154	20.2322	28.8647	97.0922	13.6596	37.3415

ALC benchmark on each domain



Last NAUC benchmark on each domain





Automated Deep Learning

Try it out!

Deep Wisdom : multimodal, no unified AutoDL framework emerged.

References

 Zhengying Liu, Isabelle Guyon, Julio Jacques Junior, Meysam Madadi, Sergio Escalera, Adrien Pavao, Hugo Jair Escalante, Wei-Wei Tu, Zhen Xu, and Sebastien Treguer. AutoCV Challenge Design and Baseline Results. In *CAp 2019 - Conférence sur l'Apprentissage Automatique*, 2019.
 Deng J, Dong W, Socher R, Li L-J, Li K, Fei-Fei L. Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. 2009. p. 248–55.
 Sungbin Lim, Ildoo Kim, Taesup Kim, Chiheon Kim and Sungwoong Kim. Fast AutoAugment. 2019. ArXiv eprint: 1905.00397

[4] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition. In CVPR, pages 6450–6459, 2018.

[5] W. Xie, A. Nagrani, J. S. Chung, and A. Zisserman. Utterance-level aggregation for speaker recognition in the wild. In *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5791–5795, 2019.

[6] Joon Son Chung, Arsha Nagrani and Andrew Zisserman. VoxCeleb2: Deep Speaker Recognition. In *Interspeech 2018*. ISCA, 2018.

[7] Piotr Bojanowski, Edouard Grave, Armand Joulin, Tomas Mikolov. Enriching Word Vectors with Subword Information. arXiv preprint arXiv:1607.04606. 2016.