# Automated Machine Learning

**Lisheng Sun-Hosoya**

LIU Zhengying

Laboratoire en Recherche Informatique (LRI)

U. Paris-Sud / Inria / U. Paris Saclay
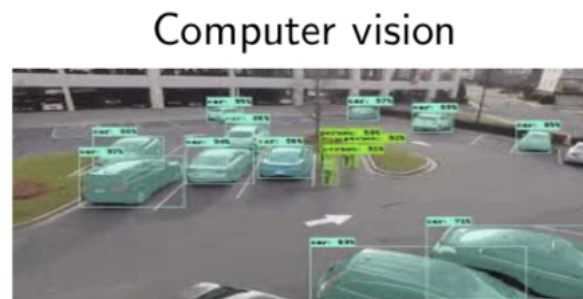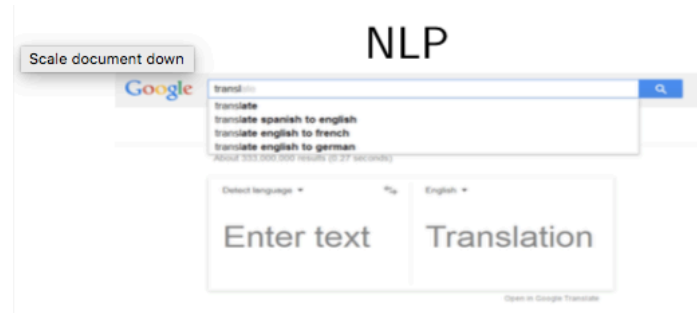
9 Mar 2021

# Contents

- AutoML: an intro

- AutoML methods
  with application to Deep Learning

- AutoML challenges
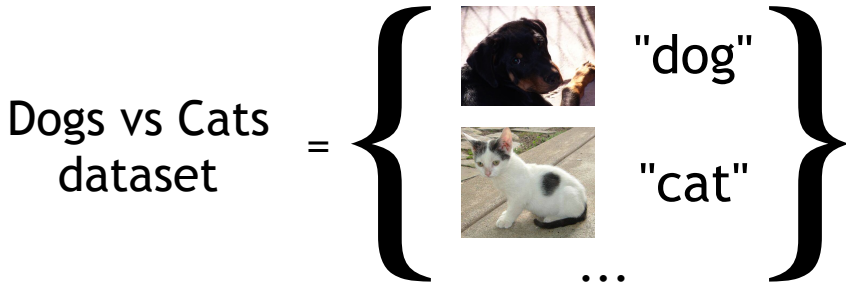
# AutoML: an intro

# Successes of Machine learning



NLP

Computer vision

Speech recognition

Games

… relies on **extensive** and **manual** tuning of algorithms and their hyperparameters
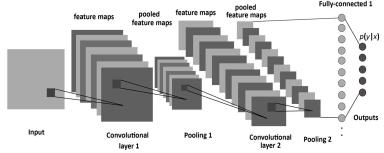
# Machine Learning

$$D = \{X_i, Y_i\} \xrightarrow{\beta_\lambda} \alpha_\theta \xRightarrow{D_{va}} P(\alpha_\theta)$$

(or $p_\theta(y|x)$)

performance
(e.g. accuracy)

Dogs vs Cats dataset $= \left\{ \begin{array}{l} \text{"dog"} \\ \text{"cat"} \\ \ldots \end{array} \right\} \longmapsto$ 

trained CNN

CIFAR-10 dataset $\longmapsto$ another trained CNN (for another $A$)

Iris dataset $\longmapsto$ trained SVM (for another $A$)

encoded by: hyperparameters $\lambda \in \Lambda$

hand-crafted by <u>ML experts</u>

Machine Learning algorithm:
Decision Tree, CNN, SVM, etc

$$D = \qquad \Rightarrow P(\alpha_\theta)$$

performance
(e.g. accuracy)

Dogs vs Cats
dataset

=

N  (for another $A$)

nother $A$)

encoded by: hyperparameters $\lambda \in \Lambda$

hand-crafted by ML experts

# Today's lecture

# The AutoML problem: definition

$$\max_{\gamma} \sum_{\substack{D_{tr}, D_{te} \\ \in \mathfrak{D}_{te}}} P(\hat{\alpha}; D_{te}) \qquad \text{where} \qquad \hat{\alpha} = \hat{\beta}(D_{tr}) \qquad \text{and} \qquad \hat{\beta} = \gamma(\mathfrak{D}_{tr})$$

supervised learning

reinforcement learning

learning to learn  ⇐  two layers of learning

$P(\hat{\alpha}; D_{te})$ may involve **time**  ⟹  computational efficiency: should be not only correct but also fast

initially we may have $\mathfrak{D}_{tr} = \varnothing$  ⟹  no prior experience BUT can be generated
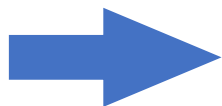
$$(D_{tr}, \beta_1, \alpha_1, P_1) , (D_{tr}, \beta_2, \alpha_2, P_2) , (D_{tr}, \beta_3, \alpha_3, P_3) , \dots$$

**Table 7.1 Supervised learning illustration of the three-level formulation.** An algorithm's level is entirely determined by its type of *input* and *output*. For a given task, finding a good $\alpha$-level algorithm is the ultimate goal. $\gamma$-level algorithms exploit data from *all past experience*, in the form of a "meta-dataset", to allow us to select a better $\beta$-level algorithm, which in turn exploits the dataset of a given task to produce an $\alpha$-level algorithm by training.

| Level | Input | Output | Examples | Encoded by |
|---|---|---|---|---|
| $\alpha$-level | sample or example (e.g. an image) | prediction of label (e.g. 'dog' or 'cat') | heuristically hard-coded classifier or already trained classifier | parameters, hyper-parameters (if any) and meta-parameters (if any) |
| $\beta$-level | task/dataset (e.g. MNIST, CIFAR-10) | $\alpha$-level algorithm | learning algorithms (e.g. SVM, CNN); HPO algorithms (e.g. grid search cross-validation, SMAC [56], NAS [124]) | hyper-parameters and meta-parameters (if any) |
| $\gamma$-level | meta-dataset (e.g. OpenML [115]) | $\beta$-level algorithm | meta-learning algorithms (e.g. meta-learning part in Auto-sklearn [36]); **algorithms from this thesis.** | meta-parameters |

Image from L. Sun. Meta-Learning as a Markov Decision Process. Machine Learning, 2019

# AutoML: what's exciting?

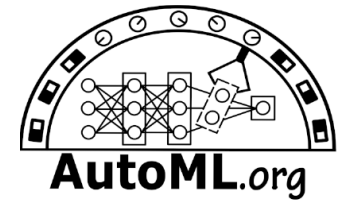- 100% autonomous
- Beat "no free lunch"
- Any time
- Any resource
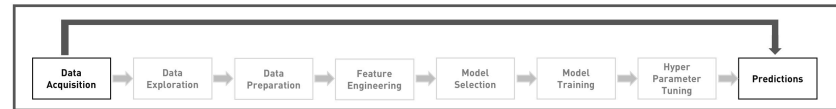
➡ AI for everyone

# AutoML: a trending topic



Google's AutoML







**Traditional Machine Learning Workflow**

Data Acquisition → Data Exploration → Data Preparation → Feature Engineering → Model Selection → Model Training → Hyper Parameter Tuning → Predictions

**AutoML Workflow**

Data Acquisition → Data Exploration → Data Preparation → Feature Engineering → Model Selection → Model Training → Hyper Parameter Tuning → Predictions

# AutoML methods

## with application to Deep Learning

# We'll focus on the simplest case

$\mathfrak{D}_{tr} = \varnothing$   (initially) and $\mathfrak{D}_{te} = \{(D_{tr}, D_{te})\}$   (single dataset)

$\Longrightarrow$ Hyperparameter Optimization

$\Longrightarrow$ single fixed training dataset: $D_{tr}$

$\Longrightarrow$ we only need to focus on $\beta_\lambda, \lambda \in \Lambda$

Reminder:

$$\max_{\gamma} \sum_{\substack{D_{tr}, D_{te} \\ \in \mathfrak{D}_{te}}} P(\hat{\hat{\alpha}}; D_{te}) \qquad \text{where } \hat{\hat{\alpha}} = \hat{\beta}(D_{tr}) \text{ and } \hat{\beta} = \gamma(\mathfrak{D}_{tr})$$

13

# Hyperparameter Optimization: a reformulation

an HPO algorithm aims to solve: $\max\limits_{\lambda \in \Lambda} P(\hat{\alpha}; D_{te})$ where $\hat{\alpha} = \beta_\lambda(D_{tr})$

unknown test score: $P(\hat{\alpha}; D_{te})$ $\Longrightarrow$ use an estimation (e.g. CV): $\hat{P}(\lambda)$

so usually the problem becomes

$$\max\limits_{\lambda \in \Lambda} \hat{P}(\lambda)$$

black-box optimization

expensive to compute

$\Longrightarrow$ surrogate model
(not discussed)

where

$$\hat{P} : \Lambda \to \mathbb{R}$$
$$\lambda \mapsto s = \hat{P}(\lambda) \approx P(\beta_\lambda(D_{tr}), D_{va})$$

is an estimation of the test score

Remark: some approaches optimize $\lambda$ and $\theta$ at the same time $\Longrightarrow$ bi-level optimization (ex. DARTS H. Liu et al., 2018)

$\beta_\lambda, \lambda \in \Lambda$ encodes an architecture A



Image adapted from: Automated Machine Learning - Methods, Systems, Challenges, Frank Hutter et. al, (2018) Springer.

3 ingredients in HPO (NAS):

• Search space

• Search strategy

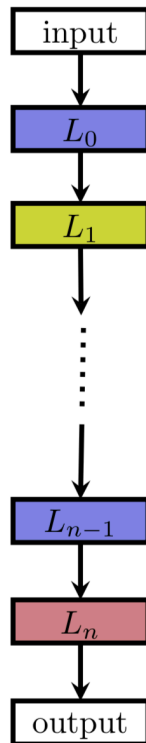• Performance estimation strategy

# Search Space (for DL)

$\beta_\lambda, \lambda \in \Lambda$ : architecture, optimizer, regularization, etc
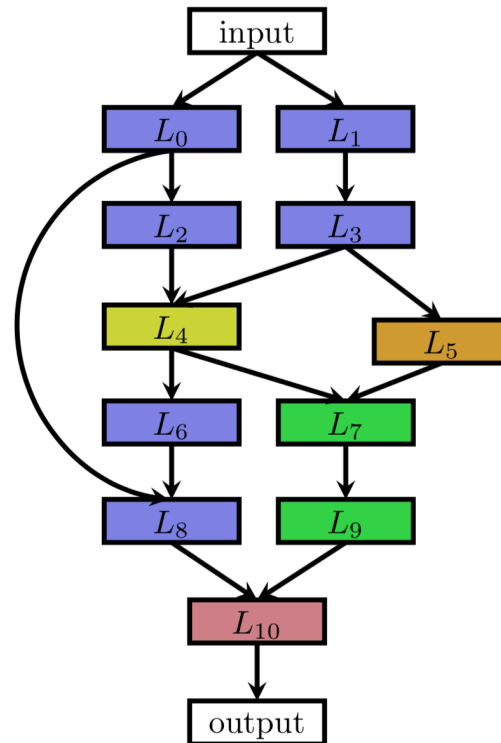
**chain-structured (feed-forward)**
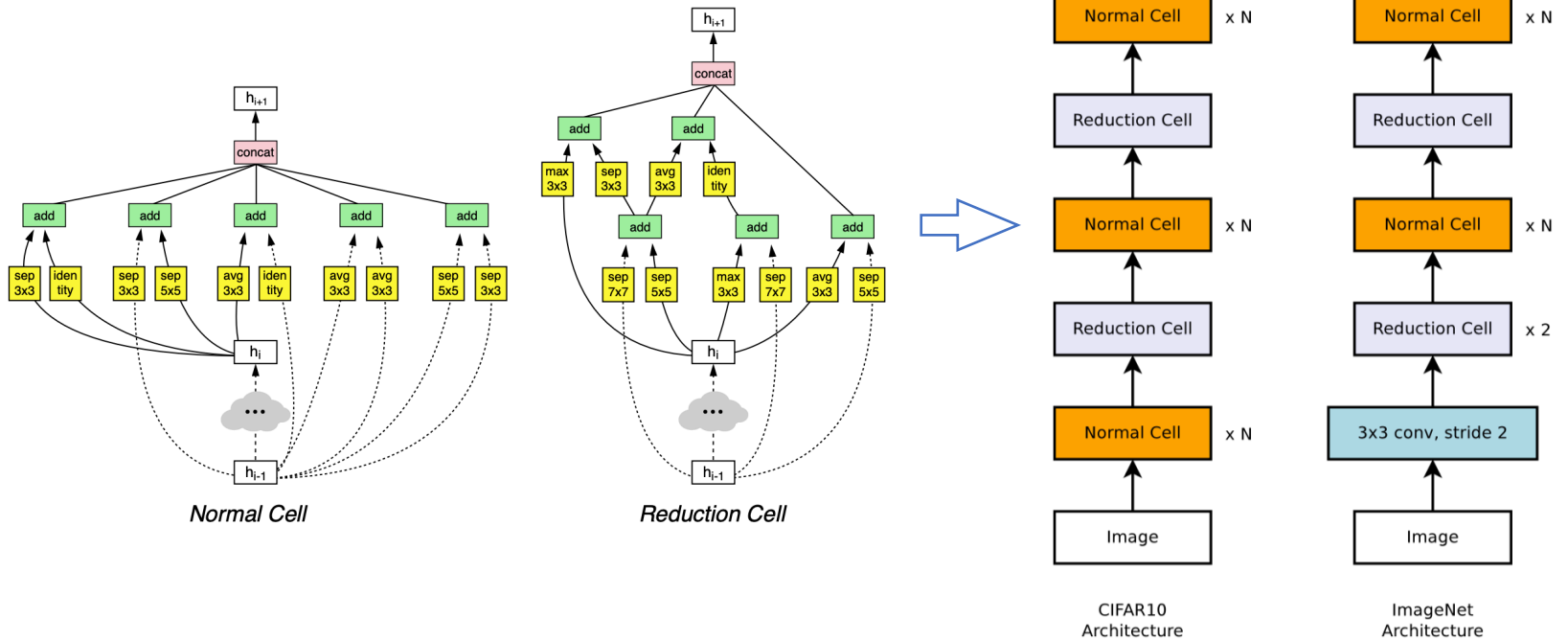
$$A = L_n \circ L_{n-1} \circ \dots \circ L_0$$

$$L_i^{in} = L_{i-1}^{out}$$

**multi-branch**

$$L_i^{in} = g_i(L_{i-1}^{out}, \dots, L_0^{out})$$



Different layer types are visualized by different colors.

Automated Machine Learning - Methods, Systems, Challenges, Frank Hutter et. al, (2018) Springer.

# Search Space (for DL)

observation: some approaches only use some
building blocks (sub-modules): ResNes, Inception, ...



## "NASNet search space" only uses two building blocks

Zoph B, Vasudevan V, Shlens J, Le QV. Learning Transferable Architectures for Scalable Image Recognition. *CVPR2018*

# Search Strategy

# Grid Search (exhaustive search)

$$\Lambda = \Lambda_1 \times \Lambda_2 \text{ with } \Lambda_1 = \{1,2,3,4\} \text{ and } \Lambda_2 = \{0.001, 0.001, 0.1, 1\}$$

# neurons in hidden layer          learning rate

try every possible combination in

$$\Lambda = \Lambda_1 \times \Lambda_2$$

evaluate it and return argmax in the end

curse of dimensionality!

# Random Search

$$\Lambda = \Lambda_1 \times \Lambda_2 \text{ with } \Lambda_1 = \{1,2,3,4\} \text{ and } \Lambda_2 = \{0.001, 0.001, 0.1, 1\}$$

# neurons in hidden layer      learning rate

Randomly sample certain number of combinations in

$$\Lambda = \Lambda_1 \times \Lambda_2$$

evaluate it and return argmax in the end
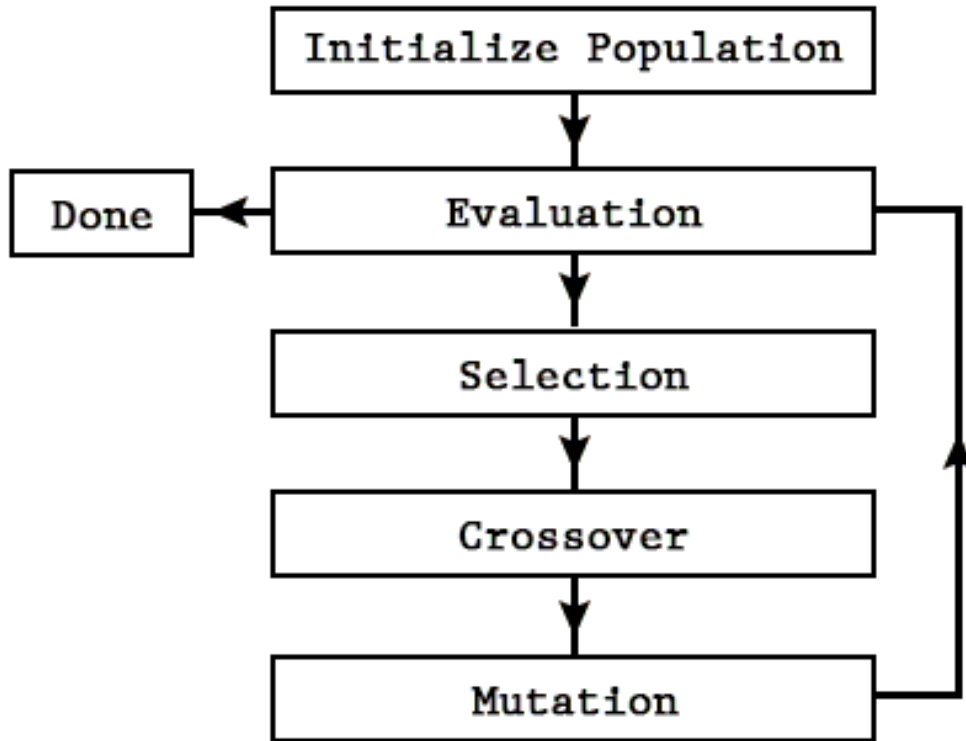
# Grid Search and Random Search

two model-free black-box optimization methods



RS tends to perform better than GS when some HP are more important than others

Random Search provides already a strong HPO baseline (surprisingly...?)

Bergstra J, Bengio Y. **Random Search for Hyper-Parameter Optimization**. *JMLR*. 2012

# Evolutionary Algorithms

Population-based derivative-free optimization methods



Optimize w.r.t a population (a set of points) or a distribution instead of one single point

Often encode an individual by "chromosome"

Explore new points by mutation or crossover

Select individuals by fitness

Just some vocabulary...but the idea is simple

Easy to parallelize

similar to: genetic algorithms, evolutionary strategies, particle swarm optimization

22

# Evolutionary Algorithm: an example

Real E, Moore S, Selle A, et al. **Large-Scale Evolution of Image Classifiers**. *ICML2017*

1000 individuals

fitness: accuracy on validation dataset

pair-wise competition
(select two individuals and kill the weaker one)

the winner gets to reproduce and mutate

massively-parallel
(due to huge computation cost)

chromosome (DNA): tensor graph

begins from single layer individuals

possible mutations:

- ALTER-LEARNING-RATE
- IDENTITY
- RESET-WEIGHTS
- INSERT-CONVOLUTION
- REMOVE-CONVOLUTION.
- ALTER-STRIDE
- ALTER-NUMBER-OF-CHANNELS
- FILTER-SIZE
- INSERT-ONE-TO-ONE
- ADD-SKIP
- REMOVE-SKIP

# Evolutionary Algorithm: an example



Real E, Moore S, Selle A, et al. Large-Scale Evolution of Image Classifiers. *ICML2017*

# Bayesian Optimization

$$\max_{\lambda \in \Lambda} \hat{P}(\lambda) \qquad \text{with } \hat{P} : \Lambda \to \mathbb{R}$$
$$\lambda \mapsto s$$

Original idea:

$\lambda$ and $s = \hat{P}(\lambda)$ follow prior distributions $p(\lambda), p(s \mid \lambda)$

we choose next point to evaluate by maximizing an acquisition function (active learning-like)

we gain more information and update $p(\lambda)$ and $p(s \mid \lambda)$ (or $p(s, \lambda)$)

repeat until convergence



Automated Machine Learning - Methods, Systems, Challenges, Frank Hutter et. al, (2018) Springer.

# Bayesian Optimization (cont'd)

$$\max_{\lambda \in \Lambda} \hat{P}(\lambda) \qquad \text{with } \hat{P} : \Lambda \to \mathbb{R}$$
$$\lambda \mapsto s$$

usual acquisition function:
**Expected Improvement (EI)**

$$a_{EI}(\lambda \,|\, D_n) = \mathbb{E}[\max(\hat{P}(\lambda) - s_{\max}, 0)]$$

usual prior model:
**Gaussian Process (GP)**

but state-of-the-art tends to use tree-based classifier such as **Random Forest** to model

$$\hat{P}(\lambda) \text{ (or } p(s\,|\,\lambda) \text{ )}$$

(thus not so Bayesian anymore…), see Auto-sklearn

# Bayesian Optimization: an example

Swersky K, Snoek J, Ada;s RP. **Freeze-Thaw Bayesian Optimization**. 2014

Intuition:
Maintains a set of "frozen" (partially completed but not being actively trained) models and uses an information-theoretic criterion to determine which ones to "thaw" and continue training

Use Bayesian Optimization for:
- learning curve prediction —> offers quick evaluations
- HP space modeling



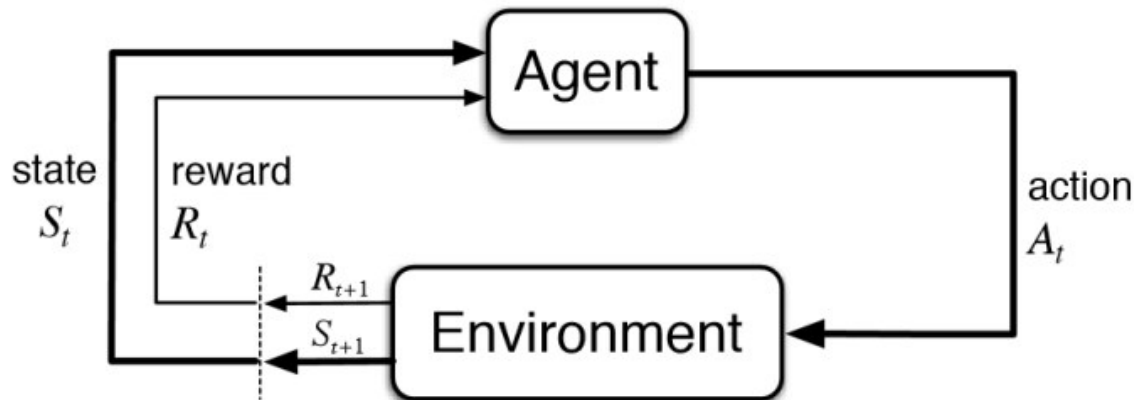(a) Graphical Model

(b) Training curve predictions
$$p(f_{t+1}|\mathcal{D}_{1:t})$$

(c) Asymptotic GP
$$p(f_{new}|\mathcal{D}_{1:t}, x_{new})$$

use notation $f : x \mapsto y$ instead of $\hat{P} : \lambda \mapsto s$

# Reinforcement Learning

A reminder:



State space: $S$  
Action space: $A$

Transition model: $\mathscr{P}^a_{ss'} = p(s'|s,a) : S \times A \times S \to [0,1]$  
Reward: $\mathscr{R}^a_{ss'} : S \times A \times S \to \mathbb{R}$

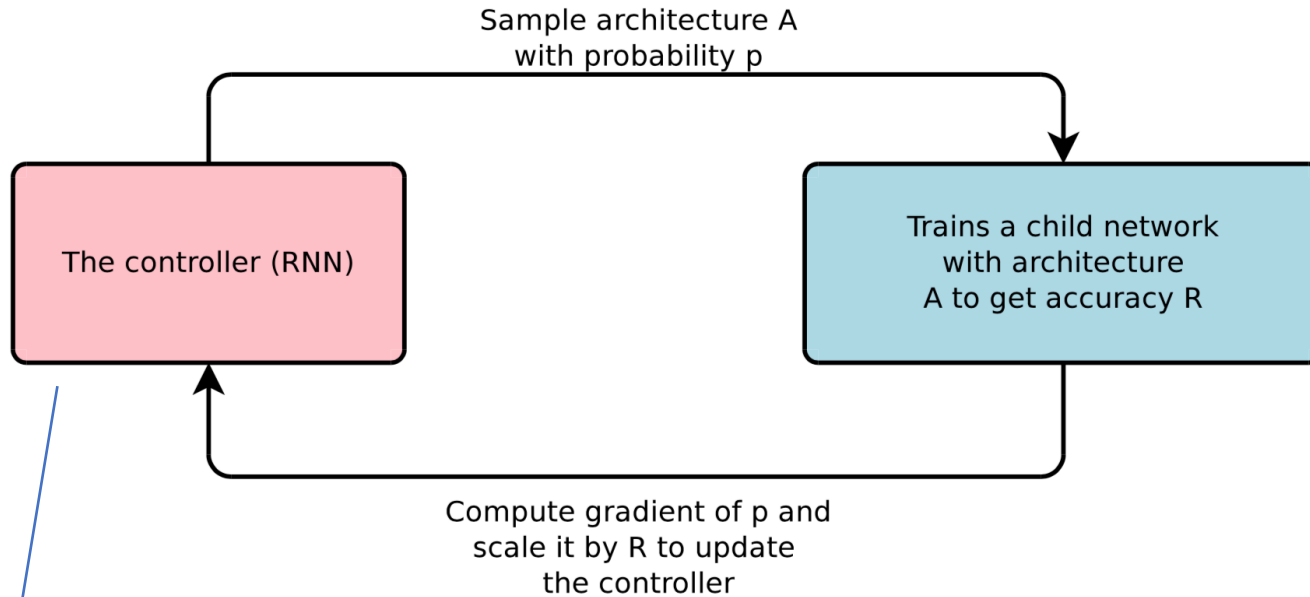Goal: Learn a policy: $\pi(s,a) = p(a|s) : S \times A \to [0,1]$

that maximizes the (discounted) expected return

$$\mathbb{E}_\pi \left[ \sum_{t=1}^{T} \gamma^t r_t \right]$$

with $T \in [0, +\infty], \gamma \in [0,1]$ and $s_0, a_0, r_1, s_1, a_1, r_2, s_2, a_2, \ldots$ the agent's trajectory

Sutton, Richard S., and Andrew G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

# Reinforcement Learning: an example

Zoph B, Le QV. **Neural Architecture Search with Reinforcement Learning**. ICLR 2017



Objective:

$$J(\theta_c) = E_{P(a_{1:T};\theta_c)}[R]$$

REINCFORCE rule:

$$\nabla_{\theta_c} J(\theta_c) = \sum_{t=1}^{T} E_{P(a_{1:T};\theta_c)} \left[ \nabla_{\theta_c} \log P(a_t|a_{(t-1):1};\theta_c)R \right]$$

an estimation:

$$\frac{1}{m} \sum_{k=1}^{m} \sum_{t=1}^{T} \nabla_{\theta_c} \log P(a_t|a_{(t-1):1};\theta_c)R_k$$

# Summary

| Method | Type | How to take next action | Update/Learn |
|---|---|---|---|
| **Grid Search** | model-free | loop over all choices (Cartesian product) | take max |
| **Random Search** | model-free | totally random | take max |
| **Bayesian Optimization** | sequential-based | maximizes acquisition function | update surrogate model |
| **Evolutionary Algorithms** | population-based | each individual randomly mutates | eliminate the weakest (with least fitness) |
| **Reinforcement Learning** | mixed/can be very general | according to learned policy | policy gradient method |
| **Differentiable Methods** | gradient-based | follow (negative) gradient | gradient descent |

There is learning in EVERY method

Is there exploration-exploitation trade-off in each method?

How do we do benchmarking and fairly evaluate these methods?

⟹ AutoDL challenge!!!

# Some other AutoML methods

Transfer Learning

Meta-learning

Ensemble methods
(competition winners)

embedded methods*: bi-level optimization methods
(related to transfer learning)

filter methods*: narrowing down the model space,
without training the learning machine
(related to meta-learning)

* Guyon I, Bennett K, Cawley G, et al. Design of the 2015 ChaLearn AutoML challenge. *IJCNN 2015*

# From one to multiple datasets: meta-learning

Given:

- ▶ Algorithms $j = 1, \ldots, m$
- ▶ PAST datasets $i = 1, \ldots, n - 1$
- ▶ a NEW dataset $n$

**Meta-dataset**: $S$ where $S(i, j)$ = score of algo. $j$ applied on dataset $i$.

Find

$$\mathrm{argmax}_{j=1,\ldots,m} \, S(n, j)$$

I.e. We want to learn some transferable knowledge across datasets (a meta-learning model $\gamma$), to solve a new dataset better and faster.

* Sun-Hosoya. Meta-learning as a Markov decision process. *2019*

# Meta-Learning: 1st trial with Auto-sklearn

Feurer M, Klein A, Eggensperger K, Springenberg JT, Blum M, Hutter F. **Efficient and Robust Automated Machine Learning**. 2015

Intuition:
Warm start the BO with meta-learning techniques, ensemble the top models.



Figure 1: Our improved approach to AutoML. We add two components to Bayesian hyperparameter optimization of an ML framework: meta-learning for initializing the Bayesian optimizer and automated ensemble construction from configurations evaluated during optimization.

**Meta-learning** [Brazdil et al., 2009]:
- characterize the dataset using meta-features,
- Initialize BP with config. That performed well on old similar dataset

**BO subroutine**: SMAC [Hutter et al. 2011]:
- Random Forest prior
- Expected improvement acquisition
- 1 fold quick evaluation

# Model-Agnostic Meta-Learning [Finn et al. 2017]

- ▶ Assumption: a single learning algorithm (NN)
- ▶ Setting: Given a distribution of datasets noted $\mathbf{D}$; with $\omega_i$ the optimal model for $D_i$

MAML finds a generally good solution:

$$\omega = \operatorname{argmax} \sum_{D_i} s_{D_i}(\omega - \alpha \nabla_\omega s_{D_i})$$

This solution is used as starting point for the new pb.



$\omega$

——— Meta-learning

············· Adaptation

$\omega$    parameter vector being meta-learned

$\omega_i^*$    optimal parameter vector for task I

$\nabla S_1$   $\nabla S_2$   $\nabla S_3$

$\omega_1^*$   $\omega_2^*$   $\omega_3^*$

34

# AutoML challenges

# The AutoML challenge (Guyon et al., 2015-2016)



**image**　　　　**medical**　　　　**speech**　　　　**marketing**

**Task variabilities:**

- classification / regression
- various scoring functions
- various time budget
- etc.

Goal: Find a process to identify the best $\beta_\lambda$ for each task

[1]: Design of the 2015 ChaLearn AutoML challenge, Guyon et al., 2015
[2]: Lessons learned from the AutoML challenge, Sun-Hosoya, Guyon and Sebag, 2018

# After the AutoML challenge series

http://automl.chalearn.org/

# AutoDL

# AutoDL challenge 2019-2020



- IMAGE
- VIDEO
- SPEECH
- TEXT
- TABULAR
- Multi-label tasks

Liu Z, Xu Z, Rajaa S, Madadi M. Towards Automated Deep Learning: Analysis of the AutoDL challenge series 2019. To appear in *NeurIPSCD2019* in Proceedings of Machine Learning Research (PMLR) 2019:10.

**(1) Raw data** from **5 modalities**: Image, Video, Speech, Text, Tabular.

**(2) Fixed time budget. Any-time learning** (ALC metric). **Blind testing**.

**(3) Starting kit, sample "public" data** and **baselines** provided.

**(4) Fixed computational resources.**

**(5)** Using **Deep Learning** was **NOT imposed.**

# Neural architectures used in the winning approaches



Base predictor / architecture
10 responses

| Architecture name | # Parameters | Domains | Teams |
|---|---|---|---|
| ResNet-18, ResNet-9 (He et al 2015) | 11.4M, 5.7M | image, video | Kakaobrain, DeepWisdom, automl_freiburg |
| MC3 (Du Tran et al CVPR 2018) | **32.8M** | video | DeepWisdom |
| EfficientNet-(b0, b1, b2) (M. Tan and Q. Le. 2019) | 5.3M, 7.8M, 9.2M | image, video | DeepWisdom, automl_freiburg |
| MobileNetV2 (M. Sandler et al 2019) | 3.4M | image, video | team_zhaw, DeepBlueAI |
| TextCNN | variable | text | Upwind_flys, DeepWisdom |
| Fast RCNN (Ross Girshick) | | text | DeepWisdom |
| LSTM, BiLSTM (Hochreiter, Schmidhuber 1997) | 0.2M-1M | text, speech | frozenmad, PASA_NJU |
| GRU, BiGRU, (Kyunghyun Cho et al 2014) GRU with Attention | 0.1M-1M | text, speech | DeepBlueAI, DeepWisdom |
| BERT-like (Tiny-BERT(X.Jiao et al)) | <110M | text | frozenmad, upwind_flys |
| DNN | <1M | tabular | DeepWisdom |

# AutoML techniques vs domains

| Approach | image | video | speech | text | tabular |
|---|---|---|---|---|---|
| Meta-learning | Offline meta-training transferred with AutoFolio [25] based on meta-features (*automl freiburg*) Offline meta-training generating solution agents, searching for optimal sub-operators in predefined sub-spaces, based on dataset meta-data. (*DeepWisdom*) MAML-like method [17] (*team zhaw*) | | | | |
| Preprocessing | Image cropping and data augmentation (*PASANJU*), fast autoaugment (*DeepBlueAI*) | Sub-sampling keeping 1/6 frames and adaptive image size (*DeepBlueAI*) Adaptive image size | MFCC, Mel Spectrogram, STFT | root features extractions with stemmer, meaningless words filtering (*DeepBlueAI*) | Numerical and Categorical data detection and encoding |
| Hyperparameter Optimization | Offline with BOHB [26] (Bayesian Optimization and Multi-armed Bandit) (*automl freiburg*) Sequential Model-Based Optimization for General Algorithm Configuration (SMAC) (*automl freiburg*) | | | | Baysien Optimization (*PASANJU*) HyperOpt [27] (*Inspur AutoDL*) |
| Transfer learning | Pre-trained on ImageNet [28] (all teams except *Kon*) | Pre-trained on ImageNet [28] (all top-8 teams except *Kon*) MC3 model pretrained on Kinetics (*DeepWisdom*) | ThinResnet34 pre-trained on VoxCeleb2 (*DeepWisdom*) | FastText pre-trained on Common Crawl (*frozenmad*) | |
| Ensemble learning | Adaptive Ensemble Learning (ensemble latest 2 to 5 predictions) (*DeepBlueAI*) | Ensemble Selection [29] (top 5 validation predictions are fused) (*DeepBlueAI*); Ensemble models sampling 3, 10, 12 frames (*DeepBlueA*) | last best predictions ensemble strategy (*DeepWisdom*) averaging 5 best overall and best of each model: LR, CNN, CNN+GRU (*DeepBlueA*) | Weighted Ensemble over 20 best models [29] (*DeepWisdom*) | LightGBM ensemble with bagging method [30] (*DeepBlueAI*), Stacking and blending (*DeepWisdom*) |

# Teams vs domains

| Team | image | video | speech | text | tabular |
|---|---|---|---|---|---|
| DeepWisdom | [**ResNet-18** and ResNet-9 models] [pretrained on ImageNet] | [MC3 model] [pretrained on Kinetics] | [fewshot learning ] [LR, ThinRestnet34 models] [pretrained on VoxCeleb2] | [fewshot learning] [task difficulty and similarity evaluation for model selection] [SVM, **TextCNN**,[fewshot learning] RCNN, GRU, GRU with Attention] | [**LightGBM**, Xgboost, Catboost, DNN models] [no pretrained] |
| DeepBlueAI | [data augmentation with Fast AutoAugment] [**ResNet-18** model] | [subsampling keeping 1/6 frames] [Fusion of 2 best models ] | [iterative data loader (7, 28, 66, 90%)] [MFCC and Mel Spectrogram preprocessing] [LR, CNN, CNN+GRU models] | [Samples truncation and meaningless words filtering] [Fasttext, **TextCNN**, BiGRU models] [Ensemble with restrictive linear model] | [3 **LightGBM** models] [Ensemble with Bagging] |
| PASA NJU | **ResNet-18** and SeResnext50; preprocessing: shape standardization and image flip (data augmentation) | **ResNet-18** and SeResnext50; preprocessing: shape standardization and image flip (data augmentation) | [data truncation(2.5s to 22.5s)][LSTM, **VggVox** ResNet with pretrained weights of DeepWis- dom(AutoSpeech2019) ThinRestnet34?] | [data truncation(300 to 1600 words)][TF-IDF and word embedding] | [iterative data loading] [Non Neural Nets models] [models complexity increasing over time] [Baysien Optimization of hyperparameters] |
| frozenmad | [images resized under 128x128] [progressive data loading increasing over time and epochs] [**ResNet-18** model] [pretrained on ImageNet] | [Successive frames difference as input of the model] [pretrained **ResNet-18** with RNN models] | [progressive data loading in 3 steps 0.01, 0.4, 0.7] [time length adjustment with repeating and clipping] [STFT and Mel Spectrogram preprocessing] [LR, LightGBM, **VggVox** models] | [TF-IDF and BERT tokenizers] [ SVM, RandomForest , CNN, tinyBERT ] | [progressive data loading] [no preprocessing] [Vanilla Decision Tree, RandomForest, Gradient Boosting models applied sequentially over time] |

# Lessons learned from the AutoDL challenge

(1) The winning methods are capable of generalizing on new unseen datasets  =>  Potential universal AutoML solutions

(2) Domain-dependent approaches are dominant
=> No universal workflows, mostly hand-tuned meta-learning

(3) We cannot afford to run expensive NAS for every new task
=> Need transferability of learned architectures

(4) Beating Baseline 3 by using "true" meta-learning is hard
=> Need more meta-train datasets (public datasets)

# MetaDL challenge

| | Input | Output | Comp. Ex. |
|---|---|---|---|
| **Alpha level:**<br>predict() in sklearn,<br>a **classifier** | $x$<br>example/sample<br>(e.g. an image) | $y$<br>labels | **Code Jam**<br>**LeetCode** |
| **Beta level:**<br>fit() in sklearn,<br>a **learning algo.** | $T$<br>ML task<br>(dataset) | $\alpha$<br>alpha-level algo | **(Auto)ML**<br>**challenges**<br>**& AutoDL** |
| **Gamma level:**<br>meta_fit() on a<br>**meta-dataset** | $\mathfrak{D}$<br>Meta-dataset | $\beta$<br>beta-level algo | **MetaDL** |

**Check and stay tuned** https://metalearning.chalearn.org/

# Conclusion

# Take-home messages

AutoML problem can be formulated in 3 levels:
$$\alpha \leftarrow \beta \leftarrow \gamma$$

Domain specific AutoML solution generalizes

Hand-crafted gamma-level learning
=> Cross-domain meta-learning yet to be studied

Any-time learning aspect to be studied further

# Stay tuned! autodl.chalearn.org