

Meta-Learning for AutoML

ACDL 2021, part 2

Joaquin Vanschoren (TU Eindhoven)

Cover art: MC Escher, Ustwo Games



Recap: What can we learn to learn?



1. architectures / pipelines (hyperparameters, structures)

Focus of this part

From hand-designed to learned learning algorithms ... to Al-generating algorithms?

See part 2





Task: distribution of samples q(x) outputs y, loss $\mathcal{L}(x,y)$

Learner: model parameters ϕ , Neural architectures, pipelines, other hyperparameters, ...

When the new task is quite similar, keep λ , (meta-)learn the model parameters φ

When the new task is quite different, (meta-)learn the hyper-parameters λ

Note: we can also learn λ and ϕ at the same time (bilevel optimization)



Automatic Machine Learning (AutoML)

Manual machine learning



AutoML: build models in a *data-driven, intelligent, purposeful* way



Hutter et al. 2019





Cleaning, preprocessing, feature selection/engineering features, model selection, hyperparameter tuning, adapting to concept drift,...





Cleaning, preprocessing, feature selection/engineering features, model selection, hyperparameter tuning, adapting to concept drift,...





Cleaning, preprocessing, feature selection/engineering features, model selection, hyperparameter tuning, adapting to concept drift,...





Cleaning, preprocessing, feature selection/engineering features, model selection, hyperparameter tuning, adapting to concept drift,...

Figure source: Nick Gillian



AutoML example: Neural Architecture Search

Architecture:

- Type of operators
- Size of layers
- Filter sizes
- Skip connections
- Pre-trained layers
- Transformers

Optimization:

. . .

- Gradient descent hyperparameters
- Regularization



Figure source: Elsken et al., 2018



AutoML + meta-learning



Hutter et al. 2019

Meta-learn how to design architectures/pipelines and tune hyper parameters Human data scientists also learn from experience



Meta-learning for AutoML: how?

hyperparameters = architecture + hyperparameters

Learning hyperparameter priors



Warm starting (what works on *similar* tasks?)



Meta-models (learn how to build models/components)



Vanschoren 2018



Observation: current AutoML strongly depends on learned priors





• Most successful pipelines have a similar structure





- Most successful pipelines have a similar structure
- Fix architecture, encode all choices as extra hyperparameters
 - Architecture search becomes hyperparameter optimization



- + smaller search space
- you can't learn entirely new architectures



- Most successful pipelines have a similar structure
- Fix architecture, encode all choices as extra hyperparameters
 - Architecture search becomes hyperparameter optimization



- + smaller search space
- you can't learn entirely new architectures

autosklearn Feurer et al. 2015 autoWEKA Thornton et al. 2013 hyperopt-sklearn *Komer et al. 2014* AutoGluon-Tabular *Erickson et al. 2020*





- Most successful pipelines have a similar structure
- Fix architecture, encode all choices as extra hyperparameters
 - Architecture search becomes hyperparameter optimization



- + smaller search space
- you can't learn entirely new architectures

autosklearn Feurer et al. 2015 autoWEKA Thornton et al. 2013 hyperopt-sklearn *Komer et al. 2014* AutoGluon-Tabular *Erickson et al. 2020*







Elsken et al. 2019

Parameterized Graph

- Choose:
- branching
- skip connections
- types of layers
- hyperparameters of layers

+ more flexible - much harder to search





Successful deep networks often have repeated motifs (cells) e.g. Inception v4:



Figure source: Szegedy et al 2016



Compositionality: learn hierarchical building blocks to simplify the task

Cell search space

- learn parameterized building blocks (cells)
- stack cells together in macroarchitecture
- + smaller search space + cells can be learned on a small dataset & transferred to a larger dataset
- strong domain priors, doesn't generalize well



Google NASNet Zoph et al 2018

Figure source: Elsken et al., 2019



Compositionality: learn hierarchical building blocks to simplify the task

Cell search space

- learn parameterized building blocks (cells)
- stack cells together in macroarchitecture
- + smaller search space + cells can be learned on a small dataset & transferred to a larger dataset
- strong domain priors, doesn't generalize well

Can we meta-learn hierarchies / components that generalize better?

Google NASNet Zoph et al 2018



Figure source: Elsken et al., 2019



- Cell construction with RL-based search (SOTA ImageNet, 450 GPUs, 3-4 days):
 - Select existing layers (hidden states, e.g. cell input) H_i to build on
 - Add operation (e.g. 3x3conv) on H_i
 - Combine into new hidden state (e.g. concat, add,...)
 - Iterate over *B* blocks



NASNet, Zoph et al 2018



Cell construction with neuro-evolution (SOTA ImageNet)





normal cell

reduction cell

Figure source: Real et al., 2019



Cell construction with multi-fidelity random search!



(a) Normal Cell

Li & Talwalkar 2019 <u>Yu et al. 2019</u> Real et al. 2019

If you constrain the search space enough, you can get SOTA results with random search!



(b) Reduction Cell

Convolutional Cells on CIFAR-10 Benchmark: Best architecture found by random search with weight-sharing.

Figure source: Li & Talwalkar., 2019





Manual priors: Weight sharing

Weight-agnostic neural networks

- ALL weights are shared
- Only evolve the architecture? ullet
 - Minimal description length \bullet
 - Baldwin effect?



Gaier & Ha, 2019







Figure source: Gaier & HA., 2019





Learning hyperparameter priors



Learn hyperparameter importance

Functional ANOVA⁷ lacksquare

- Select hyperparameters that cause lacksquarevariance in the evaluations.
- Useful to speed up black-box \bullet optimization techniques



¹ van Rijn & Hutter 2018

ResNets for image classification

Figure source: van Rijn & Hutter, 2018





Learn defaults + hyperparameter importance

• **Tunability** ^{1,2,3} *Learn* good defaults, measure importance as improvement via tuning

function max_features m = 0.16*p $m = p^0.74$ $m = 1.15^sqrt(p)$ m = sqrt(p) gamma m = 0.00574*p m = 1/p m = 0.006	
$max_features m = 0.16*p m = p^0.74 m = 1.15^sqrt(p) m = sqrt(p) gamma m = 0.00574*p m = 1/p m = 0.006$	function
m = 0.16*p $m = p^0.74$ $m = 1.15^sqrt(p)$ m = sqrt(p) gamma m = 0.00574*p m = 1/p m = 0.006	max_features
$m = p \ ^0.74$ $m = 1.15 \ ^sqrt(p)$ m = sqrt(p) gamma $m = 0.00574^*p$ m = 1/p m = 0.006	m = 0.16*p
$m = 1.15^{sqrt(p)}$ m = sqrt(p) gamma $m = 0.00574^{*}p$ m = 1/p m = 0.006	$m = p \ ^0.74$
m = sqrt(p) gamma m = 0.00574*p m = 1/p m = 0.006	$m = 1.15^sqrt(p)$
gamma m = $0.00574*p$ m = $1/p$ m = 0.006	m = sqrt(p)
m = 0.00574*p m = 1/p m = 0.006	gamma
m = 1/p m = 0.006	m = 0.00574*p
m = 0.006	m = 1/p
	m = 0.006



Learned defaults

¹ *Probst et al. 2018*

² Weerts et al. 2018

³ <u>van Rijn et al. 2018</u>



Bayesian Optimization (interlude)

- Start with a few (random) hyperparameter configurations
- Fit a *surrogate model* to predict other configurations
- Probabilistic regression: mean μ and standard deviation σ (blue band)
- Use an *acquisition function* to trade off exploration and exploitation, e.g. Expected • Improvement (EI)
- Sample for the best configuration under that function



Mockus, 1974



Bayesian Optimization

- Repeat until some stopping • criterion:
 - Fixed budget
 - Convergence ullet
 - El threshold •
- Theoretical guarantees

Srinivas et al. 2010, Freitas et al. 2012, Kawaguchi et al. 2016

- Also works for non-convex, ulletnoisy data
- Used in AlphaGo



Figure source: <u>Shahriari 2016</u>



Learn basis expansions for hyperparameters

- Hyperparameters can interact in very non-linear ways
- Use a neural net to learn a suitable basis expansion $\Phi_z(\lambda)$ for all tasks
- You can use Bayesian linear models, transfers info on configuration space

Learn basis expansion on lots of data (e.g. OpenML)







Surrogate model transfer

- If task j is *similar* to the new task, its surrogate model S_i will likely transfer well
- Sum up all S_i predictions, weighted by task similarity (as in active testing)¹
- Build combined Gaussian process, weighted by current performance on new task²



¹ Wistuba et al. 2018 ² Feurer et al. 2018



Surrogate model transfer

- Store surrogate model S_{ij} for every pair of task i and algorithm j
- Simpler surrogates, better transfer
- Learn weighted ensemble -> significant speed up in optimization





Warm starting (what works on *similar* tasks?)



How to measure task similarity?

- Hand-designed (statistical) meta-features that describe (tabular) datasets ¹
- Task2Vec: task embedding for image data ²
- Optimal transport: similarity measure based on comparing probability distributions ³ Metadata embedding based on textual dataset description ⁴
- \bullet lacksquare
- Dataset2Vec: compares batches of datasets ⁵
- Distribution-based invariant deep networks ⁶

- Vanschoren 2018
- ² Achille et al. 2019
- ³ Alvarez-Melis et al. 2020
 - ⁴ Drori et al. 2019
 - ⁵ Jooma et al. 2020
 - ⁶ de Bie et al. 2020





Warm-starting with kNN

- Find k most similar tasks, warm-start search with best λ_i
 - Auto-sklearn: Bayesian optimization (SMAC) ullet
 - Meta-learning yield better models, faster
 - Winner of AutoML Challenges •







Probabilistic Matrix Factorization

Ρ

- Collaborative filtering: configurations λ_i are `rated' by tasks t_j
- Learn latent representation for tasks T and configurations λ
- Use meta-features to warm-start on new task
- Returns probabilistic predictions for Bayesian optitmization

Fusi et al. 2017



DARTS: Differentiable NAS

- Fixed (one-shot) structure, learn which operators to use
- Give all operators a weight α_i
- Optimize α_i and model weights ω_i using bilevel optimization
 - approximate $\omega_i^*(\alpha_i)$ by adapting ω_i after every training step



One-shot model

Liu et al. 2018



Warm-started DARTS

- Warm-start DARTS with architectures that worked well on similar problems
- Slightly better performance, but much faster (5x)





(b) Normal cell found on *aircraft* task

Grobelnik and Vanschoren, 2021

(a) FLOWER_V3 transfer architecture for normal cell

(c) Normal cell found on *birds* task

Meta-models (learn how to build models/components)



Algorithm selection models

- *Learn* direct mapping between meta-features and $P_{i,i}$
 - Zero-shot meta-models: predict best λ_i given meta-features ¹ lacksquare

$$m_j \rightarrow \text{meta-learner} \rightarrow \lambda_{best}$$

Ranking models: return ranking $\lambda_{1..k}^2$

$$m_j \rightarrow \text{meta-learner} \rightarrow \lambda_{1..k}$$

Predict which algorithms / configurations to consider / tune³

$$m_j \rightarrow M_j \rightarrow \Lambda$$

Predict performance / runtime for given Θ_i and task⁴

$$m_{j,}\lambda_i \longrightarrow \text{meta-learner} \longrightarrow P_{ij}$$

• Can be integrated in larger AutoML systems: warm start, guide search,...

¹ Brazdil et al. 2009, Lemke et al. 2015 ² Sun and Pfahringer 2013, Pinto et al. 2017 ³ Sanders and C. Giraud-Carrier 2017



Learning model components

- - E.g. Swish can outperform ReLU

Swish:
$$\frac{x}{1 + e^{-\beta x}}$$

- Learn optimizers: RL-based search of space of likely useful update rules ²
 - E.g. PowerSign can outperform Adam, RMPprop

PowerSign : $e^{sign(g)sign(m)}g$

Learn acquisition functions for Bayesian optimization³

Learn nonlinearities: RL-based search of space of likely useful activation functions ¹



g: gradient, m:moving average



Figure source: Ramachandran et al., 2017 (top), Bello et al. 2017 (bottom)





Monte Carlo Tree Search + reinforcement learning

- Self-play:
 - Game actions: insert, delete, replace components in a pipeline ullet
 - Monte Carlo Tree Search builds pipelines given action probabilities
 - With grammar to avoid invalid pipelines
 - Neural network (LSTM) Predicts pipeline performance (can be pre-trained on prior datasets) \bullet



MOSAIC [*Rakotoarison et al. 2019*] AlphaD3M [Drori et al. 2019]

self play training examples actual pipeline evaluations

Figure source: Drori et al., 2019



Neural Architecture Transfer learning

- Warm-start a deep RL controller based on prior tasks
- Much faster than single-task equivalent



Wong et al. 2018

Figure source: Wong et al., 2018





Meta-Reinforcement Learning for NAS

- Train an agent how to build a neural net, across tasks
- Should transfer but also adapt to new tasks

[0, 0, 0, 0, 0][0, 0, 0, 0, 0][1, 1, 1, 0, 0][2, 2, 2, 1, 0][3, 1, 3, 0, 0][4, 3, 2, 3, 0][5, 1, 5, 0, 0][6, 2, 2, 5, 0][7, 5, 0, 2, 4][8, 7, 0, 0, 0]

Actions: add/remove certain layers in certain locations





Meta-Reinforcement Learning for NAS

Results on increasingly difficult tasks:

- a few tasks
- Policy entropy shows learning/re-• learning

omniglot



vgg_flower

dtd









MetaNAS: MAML + Neural Architecture Search

- Combines gradient based meta-learning (REPTILE) with NAS During meta-train, it optimizes the meta-architecture (DARTS weights) along with the meta-parameters (initial weights) Θ
- During meta-test, the architecture can be adapted to the novel task through gradient descent



Figure source: Elsken et al., 2020





Thank you

Image source: Pesah et al. 2018

