



Automatic Machine Learning (AutoML)

Jie Tang

Tsinghua University

June 5, 2019

Overview

- 1 Modern Hyperparameter Optimization
- 2 Neural Architecture Search
- 3 Meta-learning
- 4 Conclusions

Successes of Deep Learning

Images & Video

flickr
Google
YouTube



Product
Recommendation
amazon



Text & Language



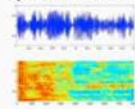
REUTERS
AP Associated Press

Relational Data/
Social Network

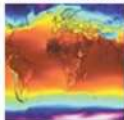
facebook
twitter



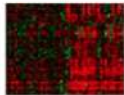
Speech & Audio



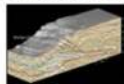
Climate Change



Gene Expression

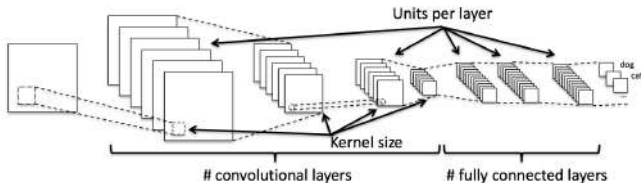


Geological Data



One Problem of Deep Learning

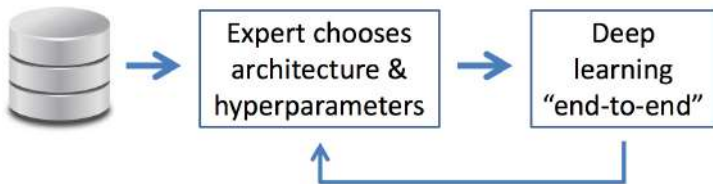
- Performance is very **sensitive** to **many hyperparameters**
 - Architectural hyperparameters



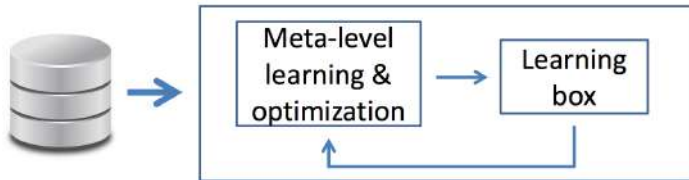
- Optimization algorithm, learning rates, momentum, batch normalization, batch sizes, dropout rates, weight decay, data augmentation, ...
 - **Easily 20-50 design decisions**
- A highly trained team of human experts is necessary: **data scientists** + **domain experts**

Deep Learning and AutoML

Current deep learning practice



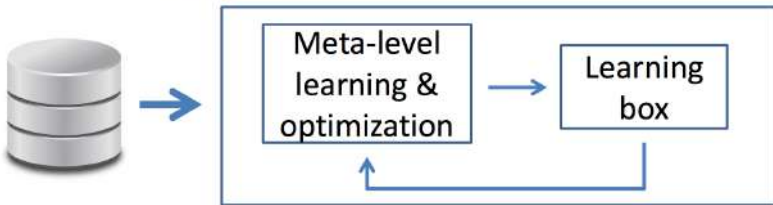
AutoML: true end-to-end learning



Learning box is not restricted to deep learning

- Traditional **machine learning pipeline**:
 - Clean & preprocess the data
 - Select / engineer better features
 - Select a model family
 - Set the hyperparameters
 - Construct ensembles of models
 - ...

AutoML: true end-to-end learning



Outline

- 1 Modern Hyperparameter Optimization
- 2 Neural Architecture Search
- 3 Meta-learning
- 4 Conclusions

Hyperparameter Optimization

Definition

Let

- $\lambda \in \Lambda$ be the hyperparameters of a ML algorithm A
- $\mathcal{L}(A_\lambda, D_{\text{train}}, D_{\text{valid}})$ denotes the loss of A , using hyperparameters λ trained on D_{train} and evaluated on D_{valid}

The **hyperparameter optimization (HPO)** problem is to find a hyperparameter configuration λ^* that minimizes this loss:

$$\lambda^* \in \arg \min_{\lambda \in \Lambda} \mathcal{L}(A_\lambda, D_{\text{train}}, D_{\text{valid}})$$

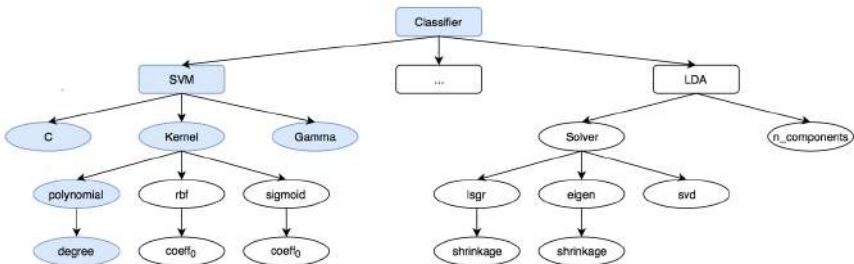
Types of Hyperparameters

- Continuous
 - Example: learning rate
- Integer
 - Example 1: #units in NN
 - Example 2: #neighbors in k-nearest neighbors
- Categorical
 - Finite domain, unordered
 - Example 1: algorithm $A \in \{\text{SVM, RF, NN}\}$
 - Example 2: activation function $\sigma \in \{\text{ReLU, sigmoid, tanh}\}$
 - Example 3: operator $\in \{\text{conv3x3, max pool, } \dots \}$
 - Example 4: the splitting criterion used for decision trees
 - Special case: binary

Conditional hyperparameters

- **Conditional hyperparameters** B are only active if other hyperparameters A are set a certain way
 - Example 1:
 - A = choice of optimizer (Adam or SGD)
 - B = Adam's momentum hyperparameter (only active if A =Adam)
 - Example 2:
 - A = type of layer k (convolution, max pooling, fully connected, ...)
 - B = conv. kernel size of that layer (only active if A = convolution)
 - Example 3:
 - A = choice of classifier (RF or SVM)
 - B = SVM's kernel parameter (only active if A = SVM)

Conditional Hyperparameters Example



AutoML as Hyperparameter Optimization

Definition: Combined Algorithm Selection and Hyperparameter Optimization (CASH)

Let

- $\mathcal{A} = \{A^{(1)}, \dots, A^{(n)}\}$ be a set of algorithms
- $\Lambda^{(i)}$ denote the hyperparameter space of $A^{(i)}$, for $i = 1, \dots, n$
- $\mathcal{L}(A_{\lambda}^{(i)}, D_{train}, D_{valid})$ denote the loss of $A^{(i)}$, using $\lambda \in \Lambda^{(i)}$ trained on D_{train} and evaluated on D_{valid} .

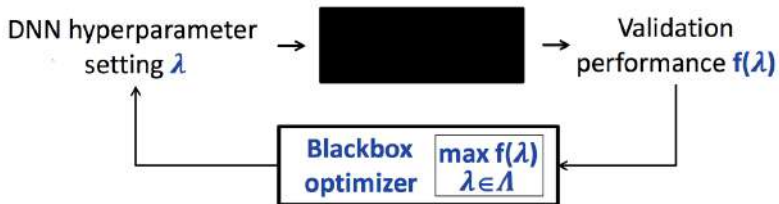
The **Combined Algorithm Selection and Hyperparameter Optimization (CASH)** problem is to find a combination of algorithm $A^* = A^{(i)}$ and hyperparameter configuration $\lambda^* \in \Lambda^{(i)}$ that minimizes this loss:

$$A_{\lambda^*}^* \in \arg \min_{A^{(i)} \in \mathcal{A}, \lambda \in \Lambda^{(i)}} \mathcal{L}(A_{\lambda}^{(i)}, D_{train}, D_{valid})$$

- $\text{CASH}^1 = \text{HPO} + \text{choice of algorithm}$

¹Chris Thornton, et al. Auto-WEKA: Combined Selection and Hyperparameter Optimization of Classification Algorithms. In KDD 2013.

Blackbox Hyperparameter Optimization



- The blackbox function is expensive to evaluate
- sample efficiency is important

Grid Search

- Each continuous hyperparameter is discretized into k equidistant values
- For categorical hyperparameters each value is used
- Cartesian product of the discretized hyperparameters

$$\Lambda_{GS} = \lambda_{1:k_1}^{(1)} \times \lambda_{1:k_2}^{(2)} \times \dots \times \lambda_{1:k_n}^{(n)}$$

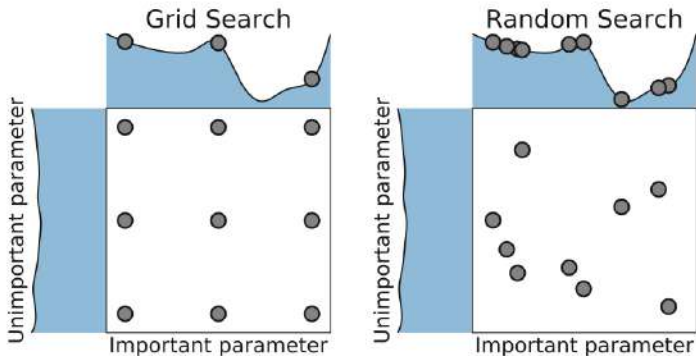
- Curse of dimensionality
- Does not exploit knowledge of well performing regions
 - Coarse grid + Finer grid

Random Search

- Converge faster than grid search
- Easier parallelization
- Flexible resource allocation
- Random search is a useful **baseline**
- Does not exploit knowledge of well performing regions
- Still very expensive

Grid Search and Random Search

- Random search works better than grid search when some hyperparameters are much more important than others



Bayesian Optimization

- An **iterative** algorithm
 - Fit a **probabilistic model** (e.g., **Gaussian Process**) to the function evaluations $\langle \lambda, f(\lambda) \rangle$
 - **Acquisition function** determines the utility of different candidate points, trading off **exploration** and **exploitation**
 - expected improvement (EI)

$$\mathbb{E}[\mathbb{I}(\lambda)] = \mathbb{E}[\max(f_{\min} - y, 0)]$$

- Upper confidence bound (UCB)

$$a_{\text{UCB}}(\lambda; \beta) = \mu(\lambda) - \beta\sigma(\lambda)$$

- ...

- Popular since Mockus[1974]
 - Sample-efficient
 - Works when objective is nonconvex, noisy, has unknown derivatives, etc
 - Recent results [Srinivas et al, 2010; Bull 2011; de Freitas et al, 2016; Kawaguchi et al, 2016]

Illustration of Bayesian optimization

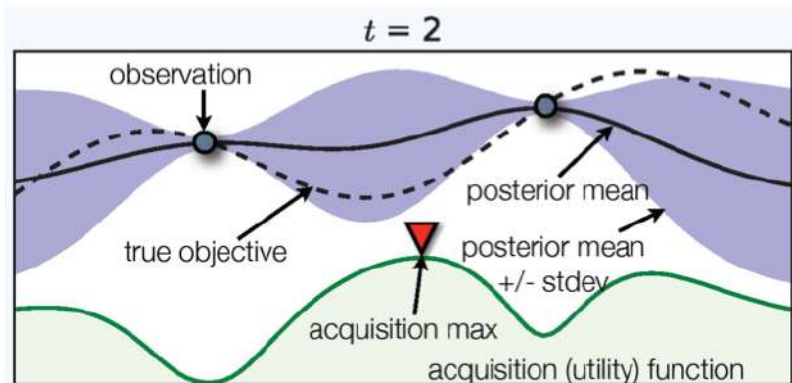


Illustration of Bayesian optimization

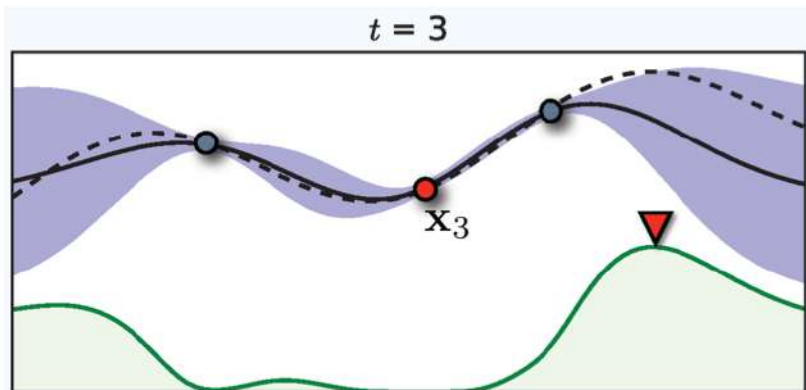
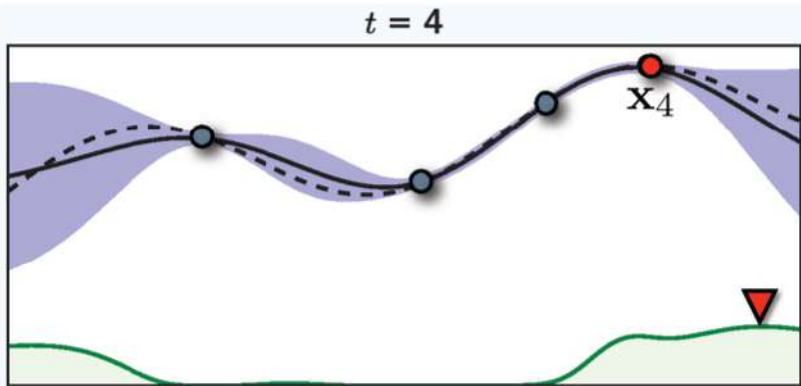


Illustration of Bayesian optimization



Example: Bayesian Optimization in AlphaGo

- “During the development of AlphaGo, **its many hyperparameters were tuned with Bayesian optimization multiple times.**”
- “This automatic tuning process resulted in **substantial improvements** in playing strength. For example, prior to the match with Lee Sedol, we tuned the latest AlphaGo agent and this improved its win-rate from 50% to 66.5% in self-play games. This tuned version was deployed in the final match.”
- “Of course, since we tuned AlphaGo many times during its development cycle, the **compounded contribution was even higher than this percentage.**”

AutoML Challenges for Bayesian Optimization

- Problems for standard Gaussian Process (GP) approach:
 - scale cubically in the number of data points
 - poor scalability to high dimensions
 - Mixed continuous/discrete hyperparameters
 - Conditional hyperparameters
- Simple solution used in **SMAC** framework²: random forests

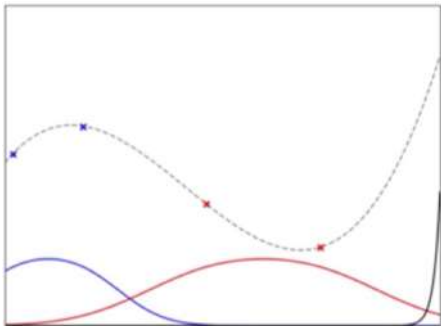
²Frank Hutter, Holger H. Hoos, Kevin Leyton-Brown. Sequential Model-Based Optimization for General Algorithm Configuration. In: Coello C.A.C. (eds) Learning and Intelligent Optimization. LION 2011. Lecture Notes in Computer Science, vol 6683. Springer, Berlin, Heidelberg

Bayesian Optimization with Neural Networks

- The simplest way: NN as a **feature extractor** to preprocess inputs and then use the outputs of the **final hidden layer** as basis functions for **Bayesian linear regression**. [Snoek et al, ICML 2015]
- **Fully Bayesian** neural network trained with stochastic gradient Hamiltonian Monte Carlo. [Springenberg et al, NIPS 2016]
- A **variational auto-encoder** can be used to embed complex inputs into a real-valued vector such that a regular Gaussian process can handle it. [Xiaoyu Lu et al, ICML 2018]
- ...

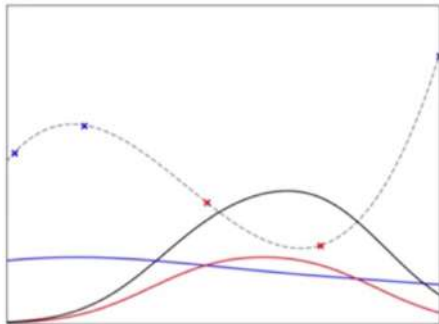
Tree of Parzen Estimators (TPE)

- Non-parametric KDEs for $p(\lambda \text{ is good})$ and $p(\lambda \text{ is bad})$, rather than $p(y|\lambda)$
- Acquisition function
 - $p(\lambda \text{ is good})/p(\lambda \text{ is bad})$
 - Equivalent to expected improvement
- Pros:
 - Efficient: $O(N^*d)$
 - Parallelizable
 - Robust
- Cons:
 - Less sample-efficient than GPs



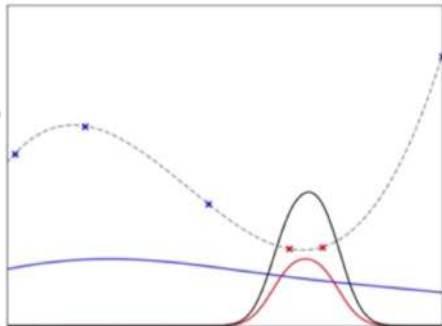
Tree of Parzen Estimators (TPE)

- Non-parametric KDEs for $p(\lambda \text{ is good})$ and $p(\lambda \text{ is bad})$, rather than $p(y|\lambda)$
- Acquisition function
 - $p(\lambda \text{ is good})/p(\lambda \text{ is bad})$
 - Equivalent to expected improvement
- Pros:
 - Efficient: $O(N^*d)$
 - Parallelizable
 - Robust
- Cons:
 - Less sample-efficient than GPs



Tree of Parzen Estimators (TPE)

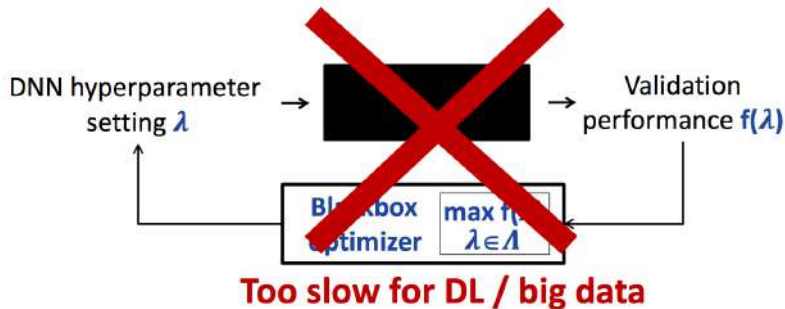
- Non-parametric KDEs for $p(\lambda \text{ is good})$ and $p(\lambda \text{ is bad})$, rather than $p(y|\lambda)$
- Acquisition function
 - $p(\lambda \text{ is good})/p(\lambda \text{ is bad})$
 - Equivalent to expected improvement
- Pros:
 - Efficient: $O(N^*d)$
 - Parallelizable
 - Robust
- Cons:
 - Less sample-efficient than GPs



Population-based methods

- population-based methods
 - maintain a **population**, i.e., a set of configurations
 - local perturbations (so-called **mutations**) and combinations of different members (so-called **crossover**) to obtain a new generation of better configurations
- **genetic algorithms, evolutionary algorithms, particle swarm optimization...**
- covariance matrix adaption evolutionary strategy (**CMA-ES**)
 - samples configurations from a multivariate Gaussian whose mean and covariance are updated in each generation based on the success of the populations individuals.
 - **dominating** the Black-Box Optimization Benchmarking (BBOB) challenge

Beyond Blackbox Hyperparameter Optimization



Hyperparameter Gradient Descent

- Formulation as bilevel optimization problem

$$\begin{aligned} & \min_{\lambda} \mathcal{L}_{val}(w^*(\lambda), \lambda) \\ \text{s.t. } & w^*(\lambda) = \operatorname{argmin}_w \mathcal{L}_{train}(w, \lambda) \end{aligned}$$

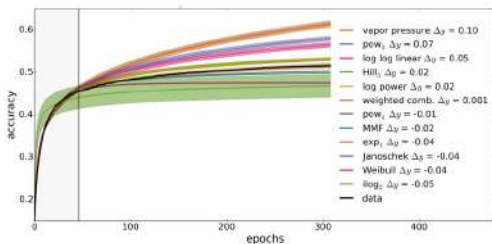
- Derive through the entire optimization process [MacLaurin et al, ICML 2015]
- Interleave optimization steps [Luketina et al, ICML 2016]

Hyperparameter gradient step w.r.t. $\nabla_{\lambda} \mathcal{L}_{val}$

Parameter gradient step w.r.t. $\nabla_w \mathcal{L}_{train}$

Probabilistic Extrapolation of Learning Curves

- Humans have one advantage: when they evaluate a poor hyperparameter setting they can quickly detect (after a few steps of SGD) and terminate the corresponding evaluation to save time
- Mimic the **early termination** of bad runs using a probabilistic model that extrapolates the performance from the first part of a **learning curve**
- **Speed up** automatic hyperparameter optimization
- **Parametric** learning curve models [Domhan et al, IJCAI 2015]



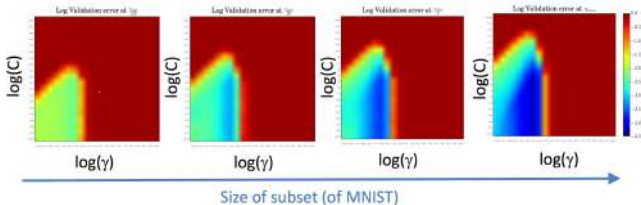
Reference name	Formula
vapor pressure	$\exp(a + \frac{b}{x} + c \log(x))$
pow ₃	$c - ax^{-\alpha}$
log log linear	$\log(a \log(x) + b)$
Hill ₃	$\frac{y_{max} x^\eta}{\kappa^\eta + x^\eta}$
log power	$\frac{a}{1 + (\frac{x}{c})^b}$
pow ₄	$c - (ax + b)^{-\alpha}$
MMF	$\alpha - \frac{\alpha - \beta}{1 + (\kappa x)^\delta}$
exp ₄	$c - e^{-ax^\alpha + \beta}$
Janoschek	$\alpha - (\alpha - \beta)e^{-\kappa x^\delta}$
Weibull	$\alpha - (\alpha - \beta)e^{-(\kappa x)^\delta}$
ilog ₂	$c - \frac{\alpha}{\log x}$

Multi-Fidelity Optimization

- Use cheap approximations of the blackbox, performance on which correlates with the blackbox, e.g.
 - Subsets of the data
 - Fewer epochs of iterative training algorithms (e.g., SGD)
 - Shorter MCMC chains in Bayesian deep learning
 - Fewer trials in deep reinforcement learning
 - Downsampled images in object recognition

Multi-fidelity Optimization

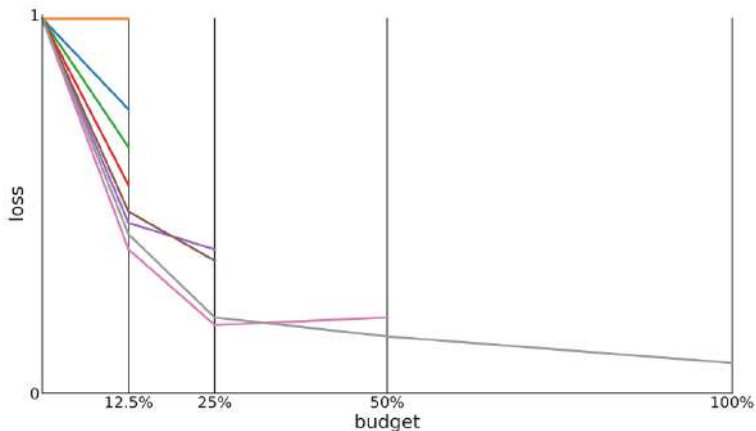
- Make use of cheap low-fidelity evaluations
 - E.g., subsets of the data (here: SVM on MNIST)



- Many cheap evaluations on small subsets
- Few expensive evaluations on the full data
- **Up to 1000x speedups** [Klein et al, AISTATS 2017]

Successive Halving (SH)

- For a given initial budget, query all algorithms for that budget; then, remove the **half** that performed worst, **double** the budget and successively repeat until only a single algorithm is left.



Hyperband

- SH suffers from **budget-vs-number of configurations** trade off
 - try many configurations and only assign a small budget to each
 - may prematurely terminate good configurations
 - try only a few and assign them a larger budget.
 - may run poor configurations too long and thereby wasting resources

Algorithm 1: HYPERBAND algorithm for hyperparameter optimization.

```
input      :  $R, \eta$  (default  $\eta = 3$ )
initialization:  $s_{\max} = \lfloor \log_{\eta}(R) \rfloor, B = (s_{\max} + 1)R$ 
1 for  $s \in \{s_{\max}, s_{\max} - 1, \dots, 0\}$  do
2    $n = \lceil \frac{B}{R} \frac{\eta^s}{(s+1)} \rceil, \quad r = R\eta^{-s}$ 
   // begin SUCCESSIVEHALVING with  $(n, r)$  inner loop
3    $T = \text{get\_hyperparameter\_configuration}(n)$ 
4   for  $i \in \{0, \dots, s\}$  do
5      $n_i = \lfloor n\eta^{-i} \rfloor$ 
6      $r_i = r\eta^i$ 
7      $L = \{\text{run\_then\_return\_val\_loss}(t, r_i) : t \in T\}$ 
8      $T = \text{top\_k}(T, L, \lfloor n_i/\eta \rfloor)$ 
9   end
10 end
11 return Configuration with the smallest intermediate loss seen so far.
```

Hyperband

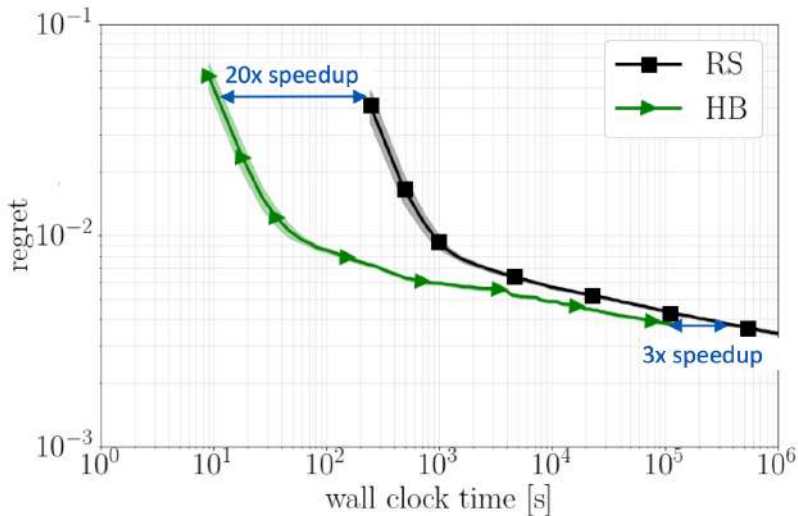
- Hyperband
 - the **outer loop** iterates over different values of **n** and **r** (lines 1-2)
 - the **inner loop** invokes **Successive Halving** for fixed values of **n** and **r** (lines 3-9)

i	$s = 4$		$s = 3$		$s = 2$		$s = 1$		$s = 0$	
	n_i	r_i	n_i	r_i	n_i	r_i	n_i	r_i	n_i	r_i
0	81	1	27	3	9	9	6	27	5	81
1	27	3	9	9	3	27	2	81		
2	9	9	3	27	1	81				
3	3	27	1	81						
4	1	81								

BOHB: Bayesian Optimization & Hyperband

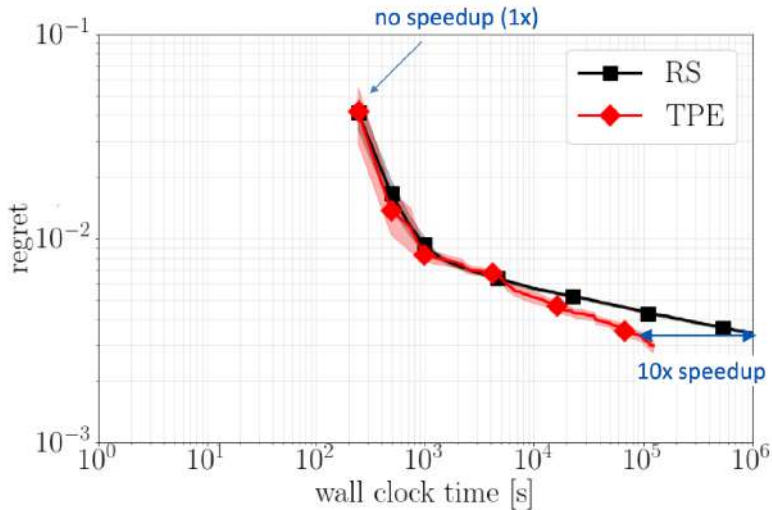
- Combining the best of both worlds in BOHB
 - Bayesian optimization
 - for choosing the configuration to evaluate
 - **strong final performance** (good performance in the long run by replacing HyperBands random search by Bayesian optimization)
 - Hyperband
 - for deciding how to allocate budgets
 - **strong anytime performance** (quick improvements in the beginning by using low fidelities in HyperBand)

Hyperband vs. Random Search



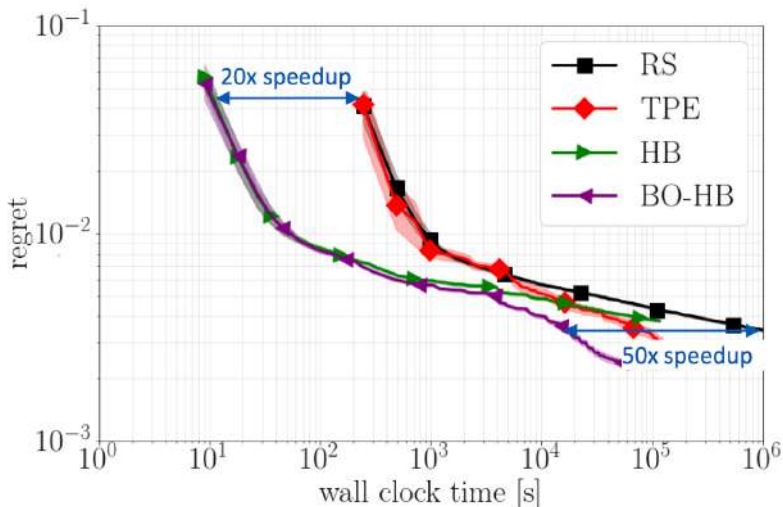
Biggest advantage: much improved **anytime** performance

Bayesian Optimization vs. Random Search



Biggest advantage: much improved **final** performance

Combining Bayesian Optimization & Hyperband



Best of both worlds: strong **anytime** and **final** performance

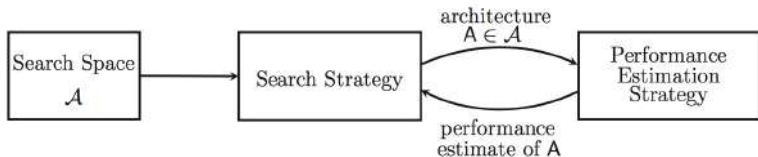
HPO Tools

- If you have access to multiple fidelities
 - BOHB
 - Combines the advantages of TPE and Hyperband
- If you do not have access to multiple fidelities
 - Low-dim, continuous: Gaussian Process-based BO (e.g., Spearmint)
 - High-dim, categorical, conditional: SMAC or TPE
 - CMA-ES
- Open-source AutoML tools based on HPO: Auto-WEKA, Hyperopt-sklearn, Auto-sklearn, TPOT, H2O AutoML...

Outline

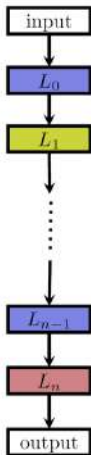
- 1 Modern Hyperparameter Optimization
- 2 Neural Architecture Search**
- 3 Meta-learning
- 4 Conclusions

Neural Architecture Search

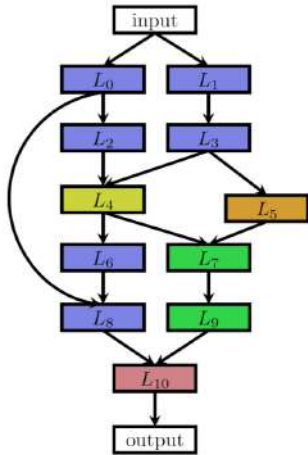


- A **search strategy** selects an architecture A from a predefined **search space** \mathcal{A} . The architecture is passed to a **performance estimation strategy**, which returns the estimated performance of A to the search strategy.

Basic Neural Architecture Search Spaces

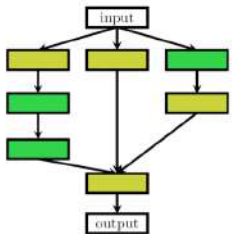
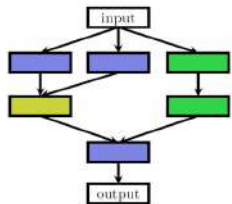


Chain-structured space
(different colours:
different layer types)

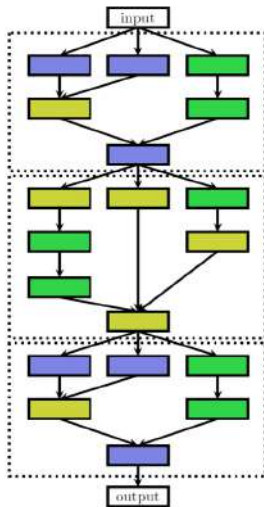
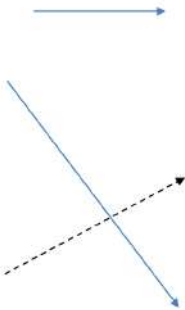


More complex space
with multiple branches
and skip connections

Cell Search Spaces



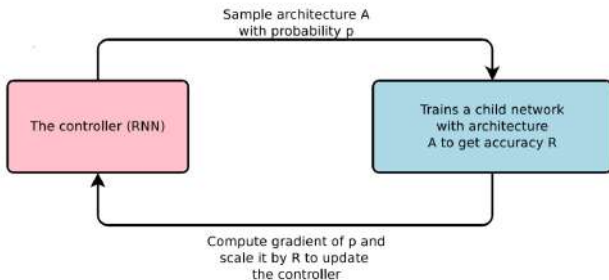
Two possible cells



Architecture composed of stacking together individual cells

Reinforcement Learning

- NAS became a mainstream research topic in the machine learning community after NAS with Reinforcement Learning [Zoph& Le, ICLR 2017]
 - State-of-the-art results for CIFAR-10, Penn Treebank
 - Large computational demands
 - 800 GPUs for 28 days, 12,800 architectures evaluated



- Different RL approaches differ in how they represent the agent's **policy** and how they **optimize** it

Neuroevolution

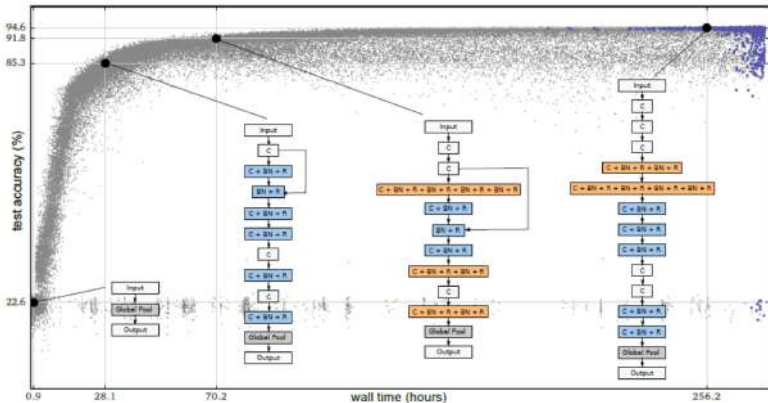
- **Neuroevolution**: use evolutionary algorithms for optimizing the neural architecture (already since the 1989³)
 - Optimize both architecture and weights with evolutionary methods
 - Use **gradient-based** methods for optimizing weights and solely use evolutionary algorithms for optimizing the neural architecture
 - scale to neural architectures with **millions** of weights for supervised learning tasks

³Miller, G., Todd, P., Hedge, S.: Designing neural networks using genetic algorithms. In: 3rd International Conference on Genetic Algorithms (ICGA89) (1989)

Neuroevolution

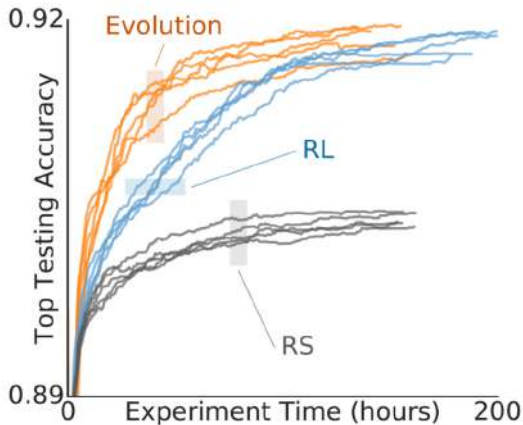
- Neuroevolution algorithms
 - a **population** of models, i.e., a set of (possibly trained) networks
 - in every evolution step, at least one model from the population is sampled and serves as a parent to generate offsprings by applying mutations to it.
 - **mutation**: local operation: adding or removing a layer, altering the hyperparameters of a layer, adding skip connections, altering training hyperparameters...
 - After training the offsprings, their fitness (e.g., performance on a validation set) is evaluated and they are added to the population
- Neuro-evolutionary methods differ in how they sample parents, update populations, and generate offsprings.

Neuroevolution



Comparison of evolution, RL and random search

- comparing RL, evolution, and random search (RS)
 - RL and evolution perform equally well in terms of final test accuracy
 - Evolution has better anytime performance and finds smaller models

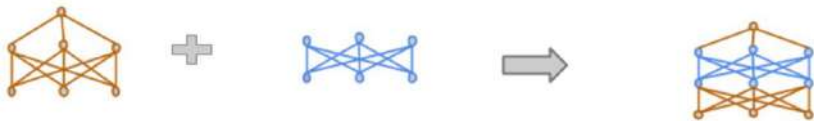


Bayesian Optimization

- Joint optimization of a vision architecture with 238 hyperparameters with TPE [Bergstra et al, ICML 2013]
- Auto-Net
 - Joint architecture and hyperparameter search with SMAC
 - First Auto-DL system to win a competition dataset against human experts [Mendoza et al, AutoML 2016]
- Kernels for GP-based NAS
 - Arc kernel [Swersky et al, BayesOpt2013]
 - NASBOT [Kandasamy et al, NIPS 2018]
- Sequential model-based optimization
 - PNAS [Liu et al, ECCV 2018]

Network morphisms

- Network morphisms
 - Change the network structure, but not the modelled function
 - for every input the network yields the same output as before applying the network morphism
- Allow efficient moves in architecture space
- Deeper, wider



Network morphisms

Definition

Network morphism Type I. Let $f_i^{w_i}(x)$ be some part of a NN $f^w(x)$, e.g., a layer or a subnetwork. We replace $f_i^{w_i}$ by

$$\tilde{f}_i^{\tilde{w}_i}(x) = Af_i^{w_i}(x) + b$$

The [network morphism equation](#) obviously holds for $A = \mathbf{1}, b = \mathbf{0}$.

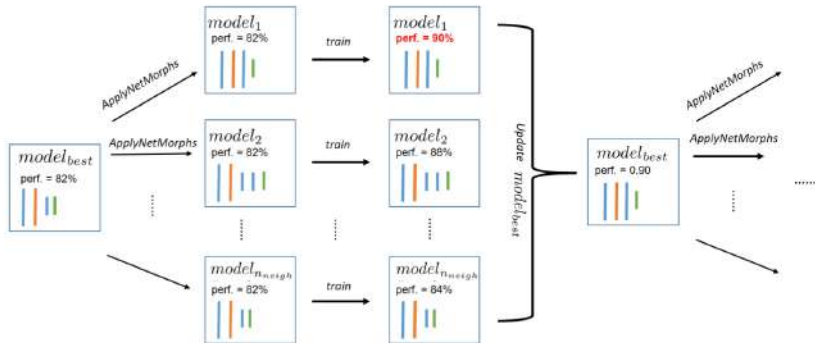
Definition

Network morphism Type II. Assume $f_i^{w_i}$ has the form $f_i^{w_i}(x) = Ah^{w_h}(x) + b$ for an arbitrary function h . We replace $f_i^{w_i}$, $w_i = (w_h, A, b)$ by

$$\tilde{f}_i^{\tilde{w}_i}(x) = \left(A \quad \tilde{A} \right) \begin{pmatrix} h^{w_h}(x) \\ \tilde{h}^{w_{\tilde{h}}}(x) \end{pmatrix} + b$$

The [network morphism equation](#) can trivially be satisfied by setting $\tilde{A} = \mathbf{0}$.

Weight inheritance & network morphisms



→ enables efficient architecture search

Outline

- 1 Modern Hyperparameter Optimization
- 2 Neural Architecture Search
- 3 Meta-learning**
- 4 Conclusions

Meta-learning

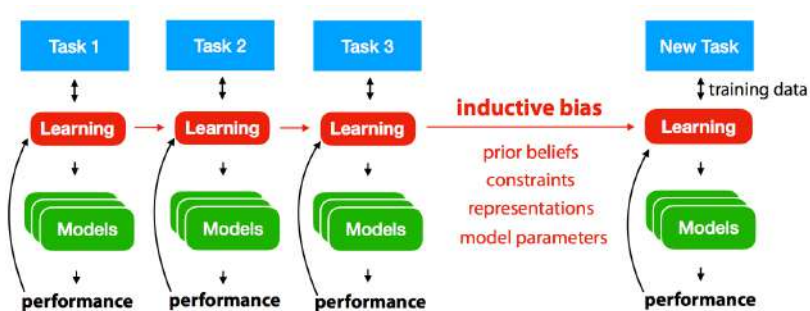
- Given a new unknown ML task, ML methods usually start from scratch to build an ML pipeline
- Meta-learning is the science of **learning to learn**
- Based on the observation of various configurations on **previous ML tasks**, meta-learning builds a model to construct promising configurations for a new unknown ML task leading to **faster convergence** with **less trial and error**

Meta-learning v.s. Multi-task learning v.s. Ensemble learning

- **Multi-task learning** learns **multiple related tasks simultaneously**
- **Ensemble learning** builds **multiple models** on the **same task**
- They do not in themselves involve learning from **prior experience on other tasks**

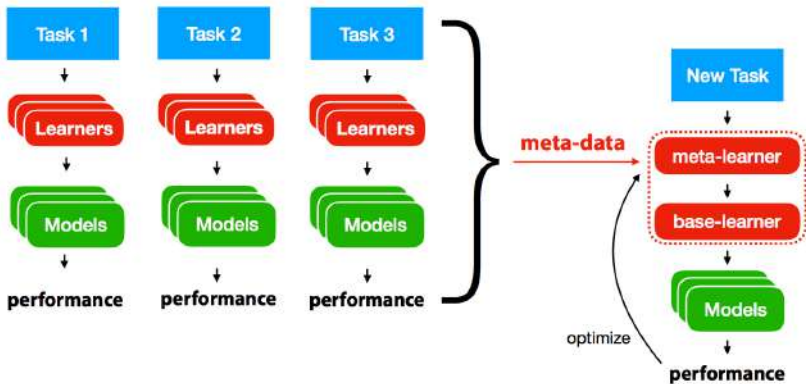
Learning to learn

- **Inductive bias:** all assumptions added to the training data to learn effectively
- If prior tasks are similar, we can transfer prior knowledge to new tasks
- if not it may actually harm learning



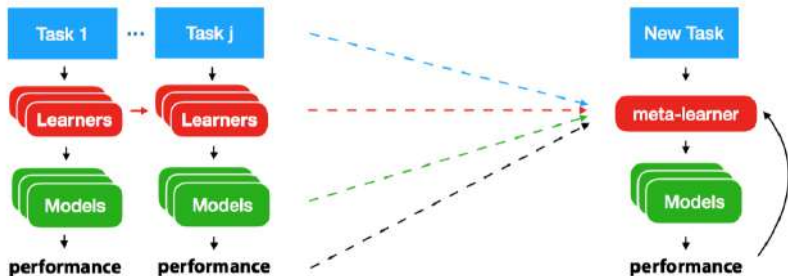
Meta-learning

- Collect meta-data about learning episodes and learn from them
- Meta-learner learns a (base-)learning algorithm, end-to-end

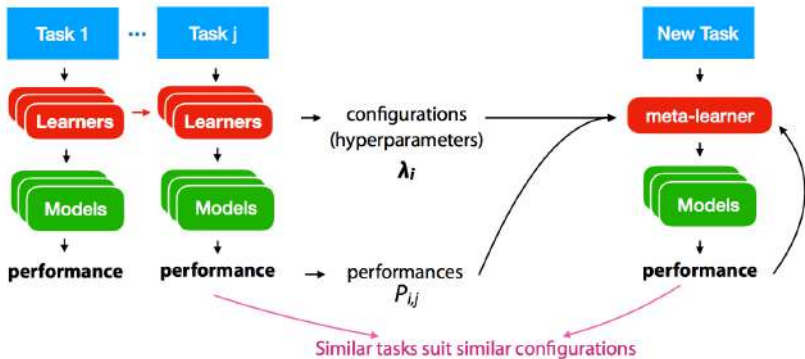


Three approaches

- Learning from Model Evaluations
- Learning from Task Properties
- Learning from Prior Models

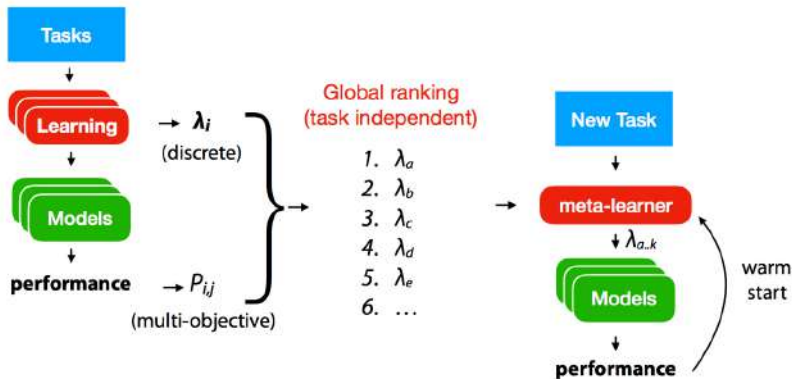


Learning from Model Evaluations



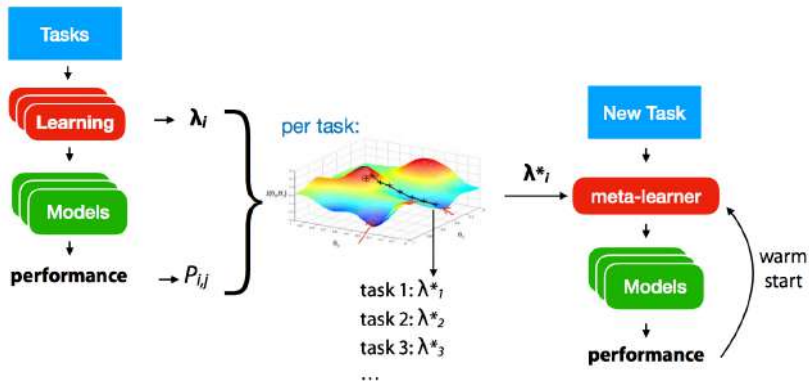
Top-K recommendation

- Build a global (multi-objective) ranking, recommend the top-K
- Requires fixed selection of candidate configurations
- Can be used as a warm start for optimization techniques



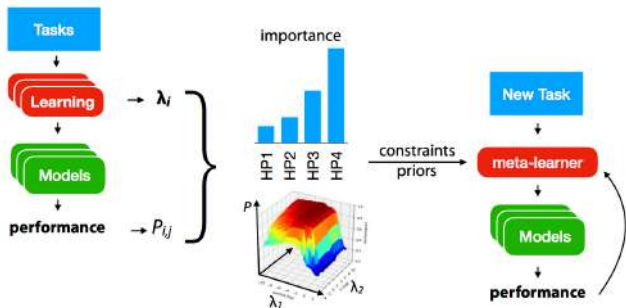
Warm-starting with plugin estimators

- What if prior configurations are not optimal?
- Per task, fit a differentiable plugin estimator on all evaluated configurations
- Do gradient descent to find optimized configurations, recommend those



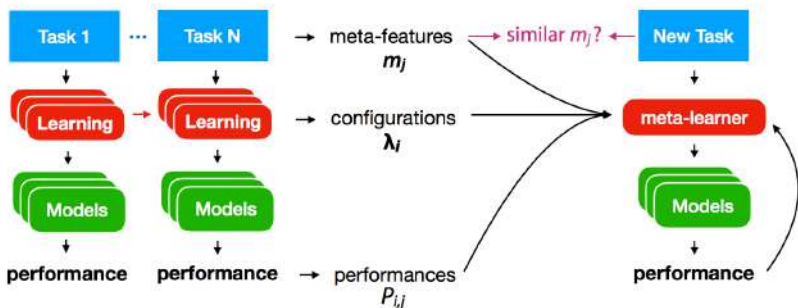
Configuration space design

- Prior evaluations can also be used to learn a better configuration space Θ^*
 - speed up the search as more relevant regions of the configuration space are explored
- **Functional ANOVA**: hyperparameters are important if they explain most of variance
- **Tunability**: learn an optimal hyperparameter, and define hyperparameter importance as the performance gain by tuning



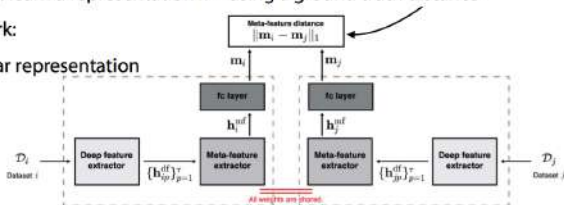
Learning from Task Properties

- Another rich source of meta-data are characterizations (meta-features) of the task at hand



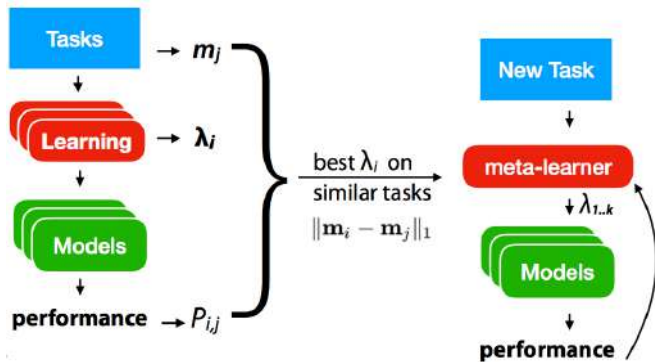
Meta-Features

- **Hand-crafted (interpretable) meta-features¹**
 - **Number of** instances, features, classes, missing values, outliers,...
 - **Statistical:** skewness, kurtosis, correlation, covariance, sparsity, variance,...
 - **Information-theoretic:** class entropy, mutual information, noise-signal ratio,...
 - **Model-based:** properties of simple models trained on the task
 - **Landmarkers:** performance of fast algorithms trained on the task
 - Domain specific task properties
- **Learning a joint task representation**
 - Deep metric learning: learn a representation h^{mf} using a ground truth distance²
 - With Siamese Network:
 - Similar task, similar representation

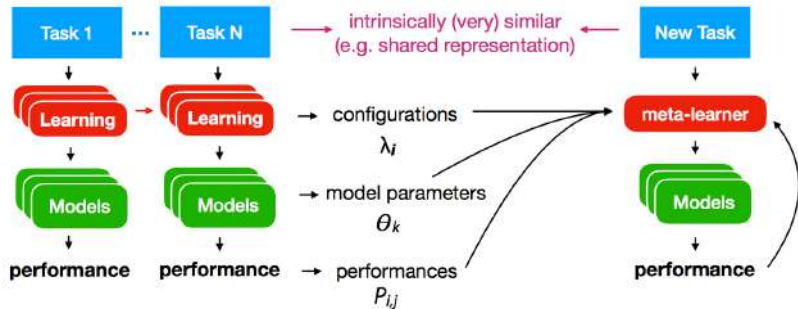


Warm-starting from similar tasks

- Find k most similar tasks, warm-start search with best λ_i
- Collaborative filtering: configurations λ_i are “related” by tasks t_j

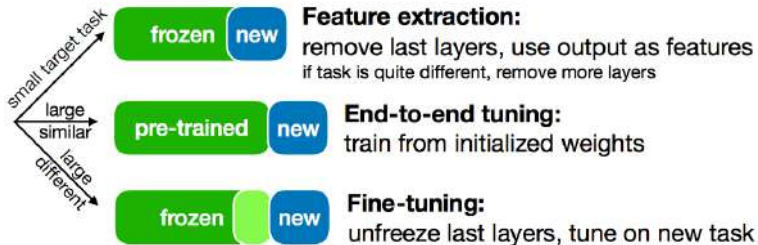


Learning from Prior Models



Transfer Learning

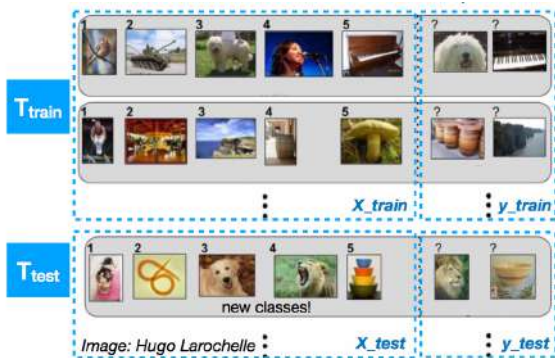
- Select source tasks, transfer trained models to similar target task
- Use as starting point for tuning, or freeze certain aspects
- Reinforcement learning: start policy search from prior policy
- Neural networks: both structure and weights can be transferred
 - Large image datasets (e.g. ImageNet)
 - Large text corpora (e.g. Wikipedia)
- Fails if tasks are not similar enough



Few-shot learning

- Learn how to learn from few examples (given similar tasks)
- Meta-learner must learn how to train a base-learner based on prior experience
- Parameterize base-learner model and learn the parameters

$$\text{cost}(\theta_i) = \frac{1}{|\mathcal{T}_{\text{test}}|} \sum_{t \in \mathcal{T}_{\text{test}}} \text{loss}(\theta_i, t)$$



Few-shot learning: approaches

- Existing algorithm as meta-learner:
 - LSTM + gradient descent
 - Learn Θ_{init} + gradient descent
 - KNN-like: Memory + similarity
 - Learn embedding + classifier
 - ...
- Black-box meta-learner:
 - Neural Turing machine (with memory)
 - Neural attentive learner
 - ...

Model-agnostic meta-learning⁴

Algorithm 1 Model-Agnostic Meta-Learning

Require: $p(\mathcal{T})$: distribution over tasks

Require: α, β : step size hyperparameters

- 1: randomly initialize θ
 - 2: **while** not done **do**
 - 3: Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
 - 4: **for all** \mathcal{T}_i **do**
 - 5: Evaluate $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ with respect to K examples
 - 6: Compute adapted parameters with gradient descent: $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
 - 7: **end for**
 - 8: Update $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$
 - 9: **end while**
-

⁴Finn, Chelsea, Pieter Abbeel, and Sergey Levine. 2017. Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks. International Conference on Machine Learning, 112635.

Outline

- 1 Modern Hyperparameter Optimization
- 2 Neural Architecture Search
- 3 Meta-learning
- 4 Conclusions**

AutoML: Further Benefits and Concerns

- Democratization of data science :)
- We directly have a strong baseline :)
- Reducing the tedious part of our work, freeing time to focus on problems humans do best (creativity, interpretation,...) :)
- People will use it without understanding anything :(

Thanks.

HP: <http://keg.cs.tsinghua.edu.cn/jietang/>
Email: jietang@tsinghua.edu.cn