S9372

Automatic Model Tuning Using Amazon SageMaker

Cyrus M Vahid Principal Evangelist – AWS AI Labs CyrusMV@amazon.com







Hyperparameters

Search Based HPO

Bayesian HPO

Amazon SageMaker AMT







Hyperparameters





What is a Hyperparameter

- **Hyperparameter** = algorithm parameter
- It affects how an algorithm behaves during model learning process
- "Any decision an algorithm author can't make for you"





Examples of Hyperparameters

Model:

Number of layers: 1, 2, 3, ... Activation functions: Sigmoid, tanh, RELU, ...



Optimization:

Method: SGD, Adam, AdaGrad, ... Learning Rate: 0.01 to 2

Data:

Batch Size: 8, 16, 32 ... Augmentation: Resize, Normalize, Color Jitter, ...





Optimizing Loss Function for a Model



...



Optimizing Loss Function for a Model





© 2019, Amazon Web Services, Inc. or its affiliates. All rights reserved

•••



Optimizing Loss Function for a Model

Optimization: min $f(x|\theta)$, where $\theta = model params$





- 'lbsqd'

•••

- num_layers = 15
- batch_size = 2000

Hyper Parameters

init=orthogonal(scale=31.5)

dropout = .2, .5, 0, .1, ..., .2



Blackbox Optimization

- We aim to minimize the objective function J.
- We have no knowledge of what the objective function is.
- We do have access to the aradients of the objective function









Model-Free HPO





Grid Search







Curse of Dimensionality

- In grid search the user specifies a finite set of values for each hyperparameter and then .
- Each hyperparameter increases degree of freedom and results in combinatorial explosion.

$$\mathbf{P} = \mathbf{C}_{\mathbf{P}} = \prod_{i=1}^{m} \mathbf{c}_{\mathbf{p}_{i}} = \mathbf{C}_{\mathbf{p}_{1}} \times \mathbf{C}_{\mathbf{p}_{2}} \times \cdots \times \mathbf{C}_{\mathbf{p}_{n}}$$





Random Search

- Random search samples configurations at random until a certain budget B for exhausted.
- The number of evaluations for Nhyperparameters is only B^{1/N}vs. B for Grid S



Bergstra, Bengio. 2012. "Random Search for Hyper-Parameter Optimization"



or the search is

earch.



Model-Based HPO for Automatic Model Training





Model Optimization

 The aim of model optimization to to find a set of model θ*, that returns the most accurate results for a given data

Optimize model-params: θ

 $\theta^* = \arg\min_{\theta \in \Theta} f(\lambda, \theta, \mathcal{D}_t \mathcal{D}_v)$



© 2019, Amazon Web Services, Inc. or its affiliates. All rights reserved.

parameters, taset.



Hyperparameter Optimization







Bayesian Optimization

- Keeps track of previous evaluations and infers expected b
- It is Bayesian in a sense that the surrogate model model probability distribution to make predictions about the po $P(Y|X) \propto P(Y|X)P(Y)$
- An acquisition function selects next points to be evaluate expected improvement.
 F[I(2)] = F[max(f_min - Y, 0)]
- As per iteration improves our beliefs about the objective applying iterative learning.





ehaviour.

uses prior sterior.

d. Commonly

function by



Bayesian Optimization

1: for n = 1, 2, 3, ... do2: select new X_{n+1} by optimizing acquisition function α $X_{n+1} = \arg \max \alpha(x; D)$

- 3: query objective function to obtain y_{n+1}
- 4: augment data $D_{n+1} = \{D_n, (x_{n+1}, y_{n+1})\}$
- 5: update the surrogate model

https://ieeexplore.ieee.org/document/7352306





Gaussian Process







Distribution Function



Gaussian Process – continued...

- Each distribution corresponds to a set of hyperparameter $\lambda_i \in \Lambda = \prod_{i=1}^n \Lambda_i$
- A Gaussian process is fully specified by a mean $\mu(\lambda)$ and a function $k(\lambda, \lambda')$.

* Co-variance functions encode all assumptions about the form of function that we are modelling. In general, covariance represents some form of distance or similarity.



a covariance



Gaussian Process for Model of Model Loss

۲











Gaussian Process for Model of Model Loss

۲



$$f(X_{t_1}), f(X_{t_2}), \dots, f(X_{t_n}) \sim \mathcal{N}(\mu, Z)$$









Gaussian Process for Model of Model Loss

۲











Covariant Function

- The mean function is usually constant in Bayesian optimization
- The quality of the Gaussian process is solely dependent on the quality of the covariant function.
- A common choice is the Matren 5/2 kernel, with its hyperparameters integrated out by Markov Chain Monte Carlo







Acquisition Function: Probability of Improvement







Acquisition Function: Expected Improvement







Using Acquisition Function

- Expected improvement
 [maximining the dashed line] has
 two components:
 - One is dependent on $-\mu$ [solid line]
 - The other dependent on uncertainty or variance $k(\lambda, \lambda')$ [blue line]
- There fore we maximize the acquisition function wherever:
 - Mean, μ , is low, or
 - Uncertainty, $k(\lambda, \lambda')$, is high.



Parallelism through Thompson Sampling

synTS

https://www.cs.cmu.edu/~kkandasa/misc/automl-slides.pdf

HPO using in Amazon SageMaker AMT

Automatic Model Tuning - workflow

Automatic Model Tuning - architecture

Fully

Input Mapping

Create a unit hypercube from all input hyperparameters

- Continuous e.g. learning rate
- Integral e.g. number of epochs
- Categorical e.g. type of activation

Sigmoid, RELU

$[0,2] \rightarrow [0,1]$ $\dots, 50\} \rightarrow [0,1]$ $tanh \} \rightarrow [0,1]^3$

DEMO!

Thank you!

Cyrus Vahid cyrusmv@amazon.com

