Aaron Klein

#### **Effects of Changing the Learning Rate**



Cast the problem as an optimization problem:

$$\mathbf{x}^{\star} = \operatorname{arg\,min}_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$$

- **X** : denotes all hyperparameters
- **f(x)** : training and validation of the machine learning algorithm with hyperparameters **X**

## Comparison

Hyperparameter optimization:

- very expensive (hours, days)
- no gradient information (at least not easy to compute)
- noisy
- non-convex

#### SGD optimization:

- single function evaluation are cheap (seconds)
- gradient information
- noisy
- non-convex

### **Configuration Space**

	Name	Range	Default	log scale	e Type	Conditional
	batch size	[32, 4096]	32	$\checkmark$	float	-
Network hyperparame-	number of updates	[50, 2500]	200	$\checkmark$	int	-
	number of layers	[1,6]	1	-	int	-
ters	learning rate	$[10^{-6}, 1.0]$	$10^{-2}$	$\checkmark$	float	-
	$L_2$ regularization	$[10^{-7}, 10^{-2}]$	$10^{-4}$	$\checkmark$	float	-
	dropout output layer	[0.0, 0.99]	0.5	$\checkmark$	float	-
	solver type	{SGD, Momentum, Adam, Adadelta, Adagrad, smorm, Nesterov }	smorm3s	20	cat	-
	lr-policy	{Fixed, Inv, Exp, Step}	fixed	(L)	$\operatorname{cat}$	<u>~</u>
Conditioned on solver type	$\beta_1$	$[10^{-4}, 10^{-1}]$	$10^{-1}$	~	float	√
	$\beta_2$	$[10^{-4}, 10^{-1}]$	$10^{-1}$	$\checkmark$	float	$\checkmark$
	ρ	[0.05, 0.99]	0.95	$\checkmark$	float	$\checkmark$
	momentum	[0.3, 0.999]	0.9	$\checkmark$	float	$\checkmark$
Conditioned on lr-policy	$\gamma$	$[10^{-3}, 10^{-1}]$	$10^{-2}$	~	float	✓
	$\dot{k}$	[0.0, 1.0]	0.5	_	float	$\checkmark$
	8	[2,20]	2	( <del>-</del> 1)	$\operatorname{int}$	$\checkmark$
Per-layer hy- perparameters	activation-type	{Sigmoid, TanH, ScaledTanH, ELU, ReLU, Leaky, Linear}	ReLU	-	cat	√
	number of units	[64, 4096]	128	$\checkmark$	int	$\checkmark$
	dropout in layer	[0.0, 0.99]	0.5	-	float	$\checkmark$
	weight initialization	{Constant, Normal, Uniform, Glorot-Uniform, Glorot-Normal, He-Normal, He-Uniform, Orthogonal, Sparse}	He-Normal	l -	$\operatorname{cat}$	$\checkmark$
	std. normal init.	$[10^{-7}, 0.1]$	0.0005	(20)	float	$\checkmark$

different ways to tackle this optimization problem:

- manual tuning
- grid search
- random search
- Bayesian optimization
- evolutionary algorithms
- reinforcement learning
- hypergradient descent

different ways to tackle this optimization problem:

- manual tuning
- grid search
- random search
- Bayesian optimization
- evolutionary algorithms
- reinforcement learning
- hypergradient descent

### **Example: SVM with 2 hyperparameters**



#### **Grid Search**

- select a finite set of values for each hyperparameter
- 2. evaluate the Cartesian product of all hyperparameters



### **Grid Search**

#### • Pros:

- easy to implement
- parallelizable

#### • Cons:

- $\circ \quad \ \ {\rm might\ never\ find\ the\ true\ optimum}$
- $\circ$  number of function evaluations grows exponentially with the number of dimensions

### **Random Search**

- 1. Specify a distribution over the input space, e g. uniform
- draw and evaluate configurations sampled from this distribution until a predefined stopping criterion has been reached



### **Random Search**

#### • Pros:

- $\circ$  easy to implement
- parallelizable
- $\circ$  in theory it always finds the optimum
- Cons:
  - can't make use of previous observations and often needs too long to approach the optimum

#### **Exercise: Incumbent Trajectory**



## Conclusions

- Hyperparameters are important and need to be set correctly
- Instead of manually tuning them, we can cast it as an optimization problem
- Random Search works better than Grid Search
- In the AutoML track we will learn how to make hyperparameter optimization more efficient