

# **Hyperparameter optimization**

MIPT

2023

# Model selection: coherent inference

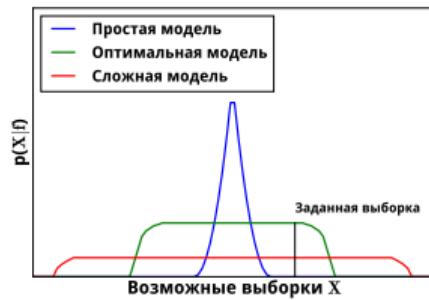
First level: select optimal parameters:

$$w = \arg \max \frac{p(\mathcal{D}|w)p(w|h)}{p(\mathcal{D}|h)},$$

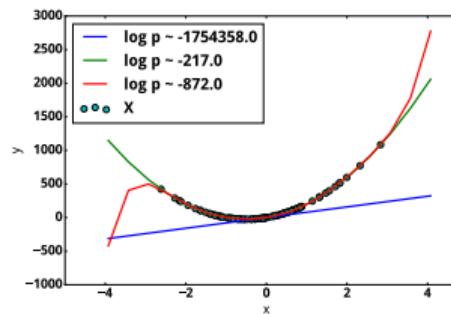
Second level: select optimal model (hyperparameters).

Evidence:

$$p(\mathcal{D}|h) = \int_w p(\mathcal{D}|w)p(w|h)dw.$$



Model selection scheme



Example: polynomials

# Hyperparameters

## Definition

*Prior* for parameters  $w$  and structure  $\Gamma$  of the model  $f$  is a distribution  $p(W, \Gamma|h) : \mathbb{W} \times \Gamma \times \mathbb{H} \rightarrow \mathbb{R}^+$ , where  $\mathbb{W}$  is a parameter space,  $\Gamma$  is a structure space.

## Definition

*Hyperparameters*  $h \in \mathbb{H}$  of the models are the parameters of  $p(w, \Gamma|h)$  (parameters of prior  $f$ ).

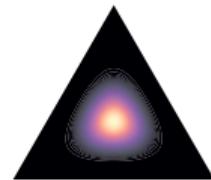
# Discrete distribution: relaxation



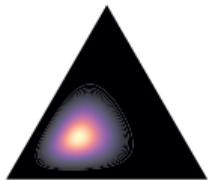
$\bar{\alpha} = [1, 1, 1]$ ,  $t = 0.9$



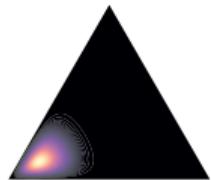
$t = 1.0$



$t = 10.0$



$\bar{\alpha} = [0.5, 0.25, 0.25]$ ,  $t = 30.0$



$[0.75, 0.125, 0.125]$



$[0.9, 0.05, 0.05]$

# Laplace approximation

Nonlinear case with  $m$  objects and  $n$  features:  $y \sim \mathcal{N}(f(X, w), \lambda^{-1})$ ,  $w \sim \mathcal{N}(0, A^{-1})$ .  
Write integral:

$$p(\mathfrak{D}|h) = p(y|X, A, \lambda) = \frac{\sqrt{\lambda \cdot |A|}}{\sqrt{(2\pi)^{m+n}}} \int_w \exp(-S(w)) dw.$$

Using Taylor series for  $S$ :

$$S(w) \approx S(\hat{w}) + \frac{1}{2} \Delta w^T H \Delta w$$

Integral reduces to the following expression:

$$\frac{\sqrt{\lambda \cdot |A|}}{\sqrt{(2\pi)^{m+n}}} S(\hat{w}) \int_w \exp\left(-\frac{1}{2} \Delta w^T H \Delta w\right) dw$$

The expression under integral corresponds to the unnormalized Gaussian PDF.

# Graves, 2011

Prior:  $p(w|\sigma) \sim \mathcal{N}(\mu, \sigma^2)$ .

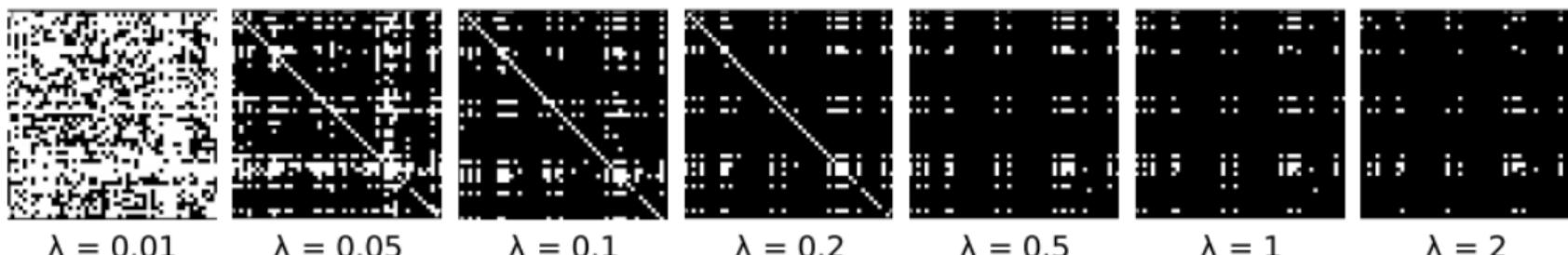
Variational inference:  $q(w) \sim \mathcal{N}(\mu_q, \sigma_q^2)$ .

Greedy optimization:

$$\mu = \hat{E}w, \quad \sigma = \hat{D}w.$$

Prune  $w_i$  using relative PDF:

$$\lambda = \frac{q(0)}{q(\mu_{i,q})} = \exp\left(-\frac{\mu_i^2}{2\sigma_i^2}\right).$$



# Problem statement

Let  $\theta \in \mathbb{R}^s$  be the set of all the optimized parameters (including variational parameters if needed).

$L(\theta, h)$  is a differential loss function  $f$ .

$Q(\theta, h)$  is a differential validation function.

The problem is to find optimal parameters  $\theta^*$  and hyperparameters  $h^*$  of the model that minimize

$$h^* = \arg \max_{h \in \mathbb{H}} Q(\theta^*(h), h),$$

$$\theta(h)^* = \arg \min_{\theta \in \mathbb{R}^s} L(\theta, h).$$

# Bayesian inference

Let  $\theta = [w]^T$ .

*First level:*

$$\theta^* = \arg \max(-L(\theta, h)) = p(w|X, y, h) = \frac{p(y|X, w)p(w|h)}{p(y|X, h)}.$$

*Second level:*

$$p(h|X, y) \propto p(y|X, h)p(h),$$

Considering  $p(h)$  improper flat prior we get the following expression:

$$Q(\theta, h) = p(y|X, h) = \int_{w \in \mathbb{R}^u} p(y|X, w)p(w|h) \rightarrow \max_{h \in \mathbb{H}}.$$

## Cross-validation

Split the dataset  $\mathcal{D}$  into  $k$  equal (maybe stratified) parts:

$$\mathcal{D} = \mathcal{D}_1 \sqcup \cdots \sqcup \mathcal{D}_k.$$

Optimize  $k$  models for each data part. Let  $\theta = [w_1, \dots, w_k]$ , where  $w_1, \dots, w_k$  are the model parameters for optimization  $k$ .

Let  $L$  be a loss function:

$$L(\theta, h) = -\frac{1}{k} \sum_{q=1}^k \left( \frac{k}{k-1} \log p(y \setminus y_q | X \setminus X_q, w_q) + \log p(w_q | h) \right). \quad (1)$$

Let  $Q$  be a validation loss:

$$Q(\theta, h) = \frac{1}{k} \sum_{q=1}^k k \log p(y_q | X_q, w_q).$$

# ELBO

Let  $L = -Q$ :

$$\log p(y|X, A) \geq \sum_{x,y} \log p(y|x, \hat{w}) - D_{KL}(q(w)||p(w|A)) = -L(\theta, A^{-1}) = Q(\theta, A^{-1}),$$

where  $q$  is a normal distribution with diagonal covariance matrix:

$$q \sim \mathcal{N}(\mu_q, A_q^{-1}),$$

$$D_{KL}(q(w)||p(w|f)) = \frac{1}{2}(\text{Tr}[AA_q^{-1}] + (\mu - \mu_q)^\top A(\mu - \mu_q) - u + \ln |A^{-1}| - \ln |A_q^{-1}|).$$

Use variational parameters of  $q$  as a vector of optimized parameters  $\theta$ :

$$\theta = [\alpha_1, \dots, \alpha_u, \mu_1, \dots, \mu_u].$$

# Evidence vs CV

Evidence estimation:

$$\log p(\mathcal{D}|f) = \log p(\mathcal{D}_1|f) + \log p(\mathcal{D}_2|\mathcal{D}_1, f) + \cdots + \log p(\mathcal{D}_n|\mathcal{D}_1, \dots, \mathcal{D}_{n-1}, f).$$

Leave-one-out estimation:

$$LOU = E \log p(\mathcal{D}_n|\mathcal{D}_1, \dots, \mathcal{D}_{n-1}, f).$$

Cross-validation uses expected values of the last term  $p(\mathcal{D}_n|\mathcal{D}_1, \dots, \mathcal{D}_{n-1}, f)$  as a complexity estimation.

Evidence considers **full** complexity.

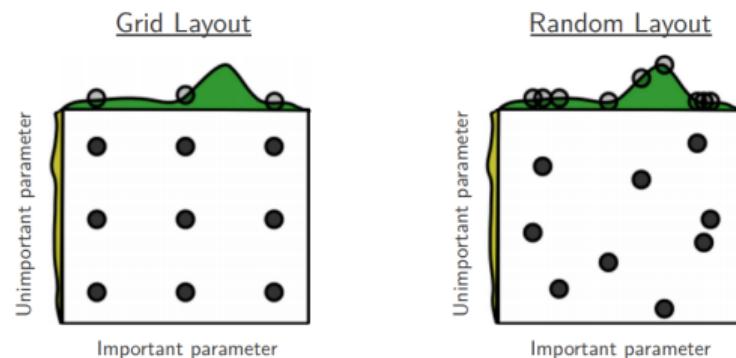
# Basic methods of hyperparameter optimization

Variants:

- Grid search;
- random search.

Both methods suffer from curse of dimensionality.

The random search can be more effective if the hyperparameter space is degenerate.



Bergstra et al., 2012

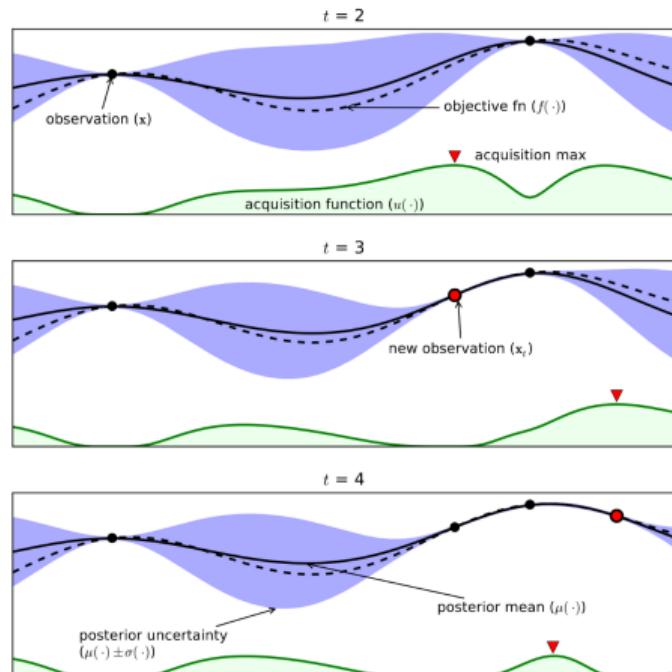
# Gaussian process

Idea: Model  $Q(\theta(h)^*, h)$  using Gaussian process depending on  $h$ .

Pros:

- Flexibility.
- Probabilistic model, cheaper than exhaustive search.

Cons: cubic complexity,  $O(|\mathbb{H}|^3)$ .



Shahriari et. al, 2016. GP example.

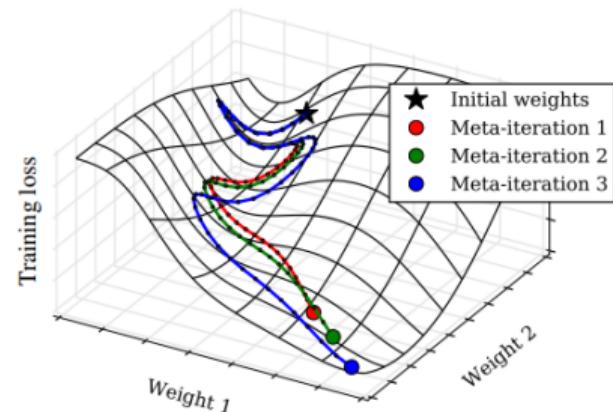
# Gradient methods

**Idea:** Optimize hyperparameters using the full parameter optimization trajectory.

**Pros:**

- Hyperparameter optimization will consider the features of the parameter optimization.
- The complexity is linear on the number of hyperparameters.

**Cons:** the computational cost is very expensive.



MacLaurin et. al, 2015. Example.

# Problem statement: gradient-based optimization

## Definition

A SGD operator  $T$  is an operator providing  $\eta$  steps of optimization:

$$\hat{\theta} = T \circ T \circ \cdots \circ T(\theta_0, h) = T^\eta(\theta_0, h), \quad (2)$$

where

$$T(\theta, h) = \theta - \lambda \nabla L(\theta, h)|_{\hat{\mathcal{D}}},$$

$\lambda$  is a learning rate,  $\theta_0$  is an initial value of  $\theta$ ,  $\hat{\mathcal{D}}$  is a random subset of  $\mathcal{D}$ .

Rewrite the optimization problem:

$$h^* = \arg \max_{h \in \mathbb{H}} Q(T^\eta(\theta_0, h)).$$

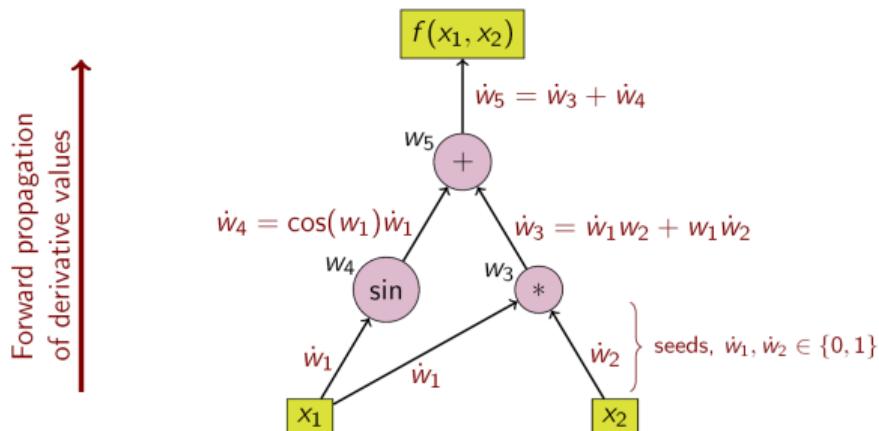
# Forward-mode differentiation

Main idea:

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial w_{n-1}} \frac{\partial w_{n-1}}{\partial x} = \frac{\partial y}{\partial w_{n-1}} \left( \frac{\partial w_{n-1}}{\partial w_{n-2}} \frac{\partial w_{n-2}}{x} \partial x \right) = \dots$$

Example (wiki):

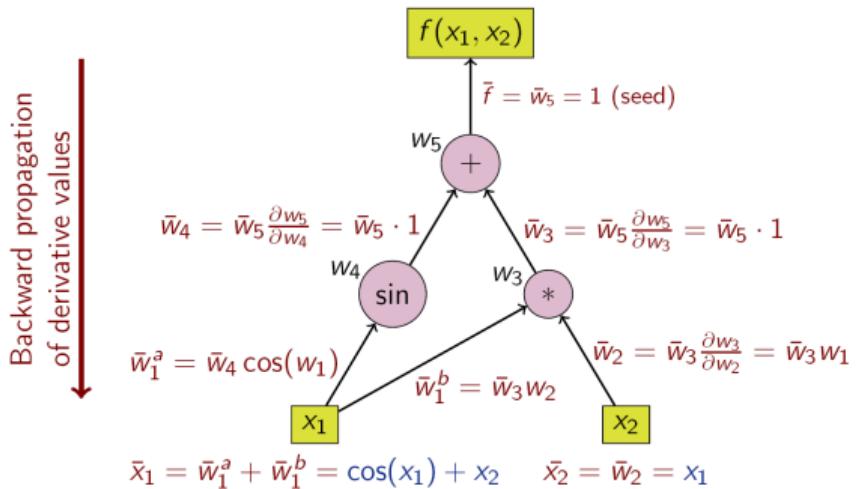
$$x_1 x_2 + \sin(x_1)$$



# Reverse-mode differentiation

Idea:

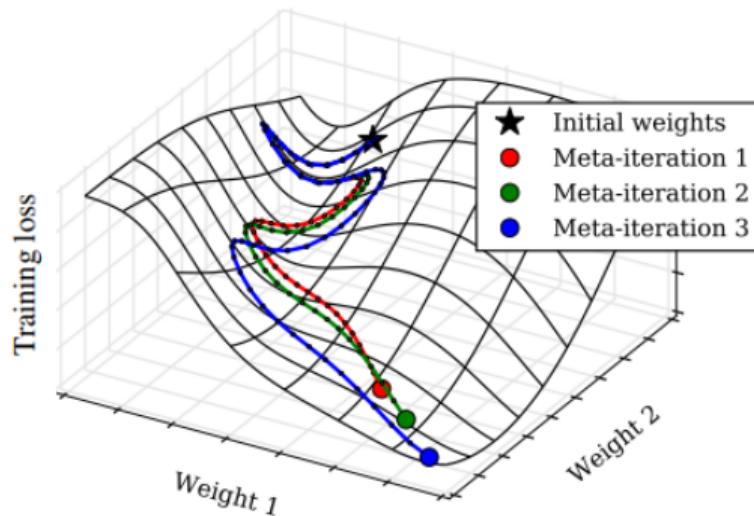
$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial w_1} \frac{\partial w_1}{\partial x} = \left( \frac{\partial y}{\partial w_2} \frac{\partial w_2}{\partial w_1} \right) \frac{\partial w_1}{\partial x} = \dots$$



# RMAD, Maclaurin et. al, 2015

- ① Run  $\eta$  steps of momentum optimization  $\gamma$ :  
 $\theta = T(\theta_0, h)$ .
- ② Let  $\hat{\nabla}h = \nabla_h Q(\theta, h)$ .
- ③ Let  $dv = 0$ .
- ④ For  $\tau = \eta \dots 1$  repeat:
  - ⑤ calculate  $\theta^{\tau-1}$ .
  - ⑥ Calculate gradient on the step  $\tau - 1$ , using RMD.

RMAD algorithm is based on the Reverse-mode differentiation.



# DrMAD

DrMad is a simplified RMAD. Idea: the parameter optimization trajectory is linear.

- |   |  |
|---|--|
| <ul style="list-style-type: none"><li>① Run <math>\eta</math> steps of memontum optimization <math>\gamma</math>:<br/><math>\theta = T(\theta_0, h)</math>.</li><li>② Let <math>\hat{\nabla}h = \nabla_h Q(\theta, h)</math>.</li><li>③ Let <math>dv = 0</math>.</li><li>④ For <math>\tau = \eta \dots 1</math> repeat:<ul style="list-style-type: none"><li>⑤ Calculate <math>\theta^{\tau-1}</math>.</li><li>⑥ Calculate gradient on the step <math>\tau - 1</math>, using RMD.</li></ul></li></ul> | <ul style="list-style-type: none"><li>① Run <math>\eta</math> steps of memontum optimization <math>\gamma</math>:<br/><math>\theta = T(\theta_0, h)</math>.</li><li>② Let <math>\hat{\nabla}h = \nabla_h Q(\theta, h)</math>.</li><li>③ Let <math>dv = 0</math>.</li><li>④ For <math>\tau = \eta \dots 1</math> repeat:<ul style="list-style-type: none"><li>⑤ <math>\theta^{\tau-1} = \theta_0 + \frac{\tau-1}{\eta} \theta^\eta</math>.</li><li>⑥ Calculate gradient on the step <math>\tau - 1</math>, using RMD.</li></ul></li></ul> |
|---|--|

# Optimization closed form

## Statement (Pedregosa, 2016)

Let  $L$  be a differential function such that all the stationary points of  $L$  are global minima. Let the Hessian  $H$  of  $L$  be invertible in each stationary point.

Then

$$\nabla_h Q(T(\theta_0, h), h) = \nabla_h Q(\theta^\eta, h) - \nabla_h \nabla_\theta L(\theta^\eta, h)^T H^{-1} \nabla_\theta Q(\theta^\eta, h).$$

# Greedy optimization of hyperparameters

On each step of  $\theta$ :

$$h' = h - \lambda_h \nabla_h Q(T(\theta, h), h) = h - \lambda_h \nabla_h Q(\theta - \lambda \nabla L(\theta, h), h),$$

where  $\lambda_h$  is a learning rate of hyperparameter optimization.

- The method can be considered as a simplified RMAD that uses only one step of parameter optimization.
- The method can be considered as an approximation of the closed form expression for  $H^{-1} \sim I$ .
- Complexity:  $O(|\theta| \cdot |h|)$
- Can be simplified using finite differences, see DARTS,  $O(|\theta| + |h|)$

# HOAG

Approximation of closed-form formula:

$$\nabla_h Q(\theta^\eta, h) - \nabla_h \nabla_\theta L(\theta^\eta, h)^T H^{-1} \nabla_\theta Q(\theta^\eta, h).$$

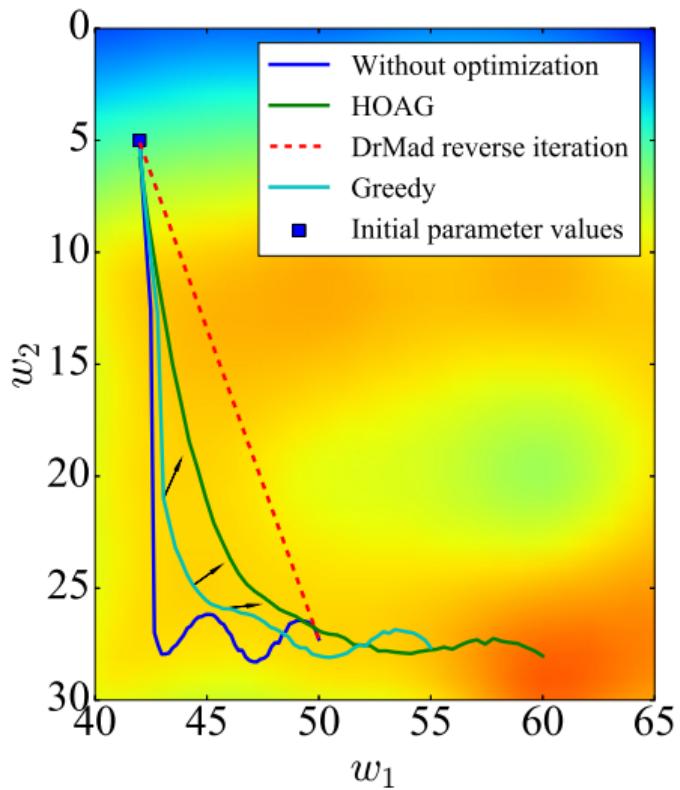
- ① run  $\eta$  optimization steps:  $\theta = T(\theta_0, h)$ .
- ② Solve linear system for  $\lambda$ :  $H(\theta)\lambda = \nabla_\theta Q(\theta, h)$ .
- ③ Calculate approximate value of hyperparameter gradients  
$$\hat{\nabla}_h Q = \nabla_h Q(\theta, h) - \nabla_{\theta,h} L(\theta, h)^T \lambda.$$

Update formula:

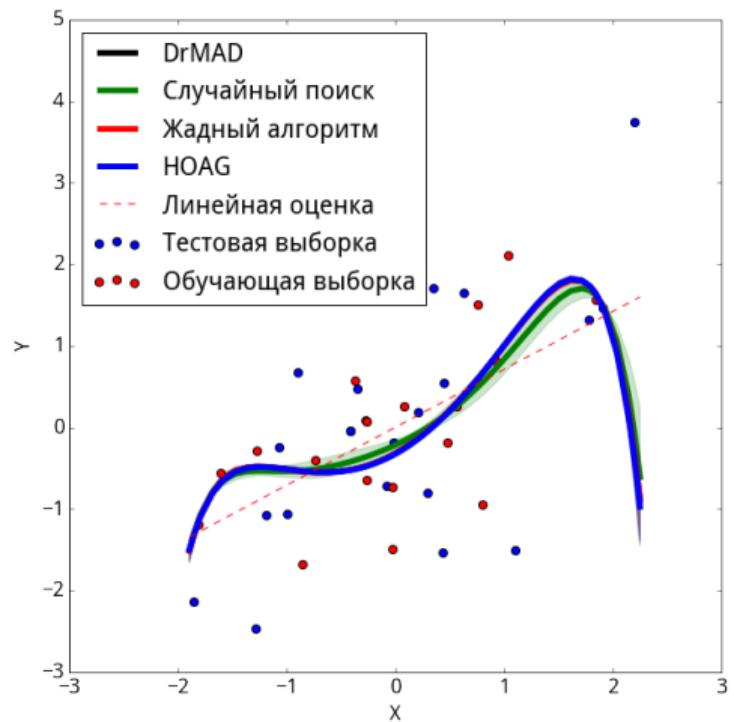
$$h' = h - \gamma_h \hat{\nabla}_h Q.$$

# Algorithm comparison

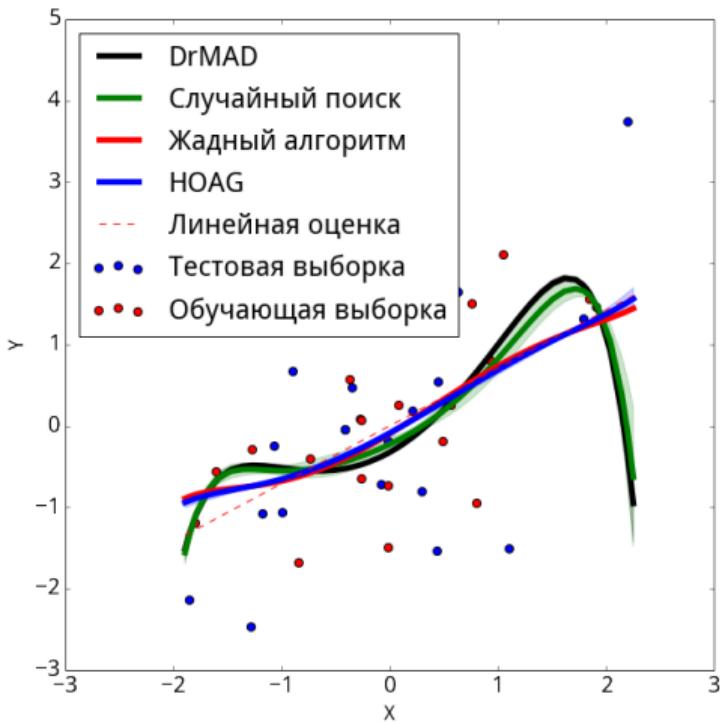
Algorithm	+	-
Random search	Easy to implement	Curse of dimensionality
Greedy optimization	Easy to implement, can be run inside common gradient optimization loop	Non-optimality.
HOAG	Fast convergence.	Results quality depends on the linear system solution precision $H(\theta)\lambda = \nabla_{\theta} Q(\theta, h)$ .
DrMAD	Considers optimization features. Can be used for metaparameter optimization.	Not very robust. Strong assumptions about linearity of the optimization trajectory.



# Experiment: polynomials

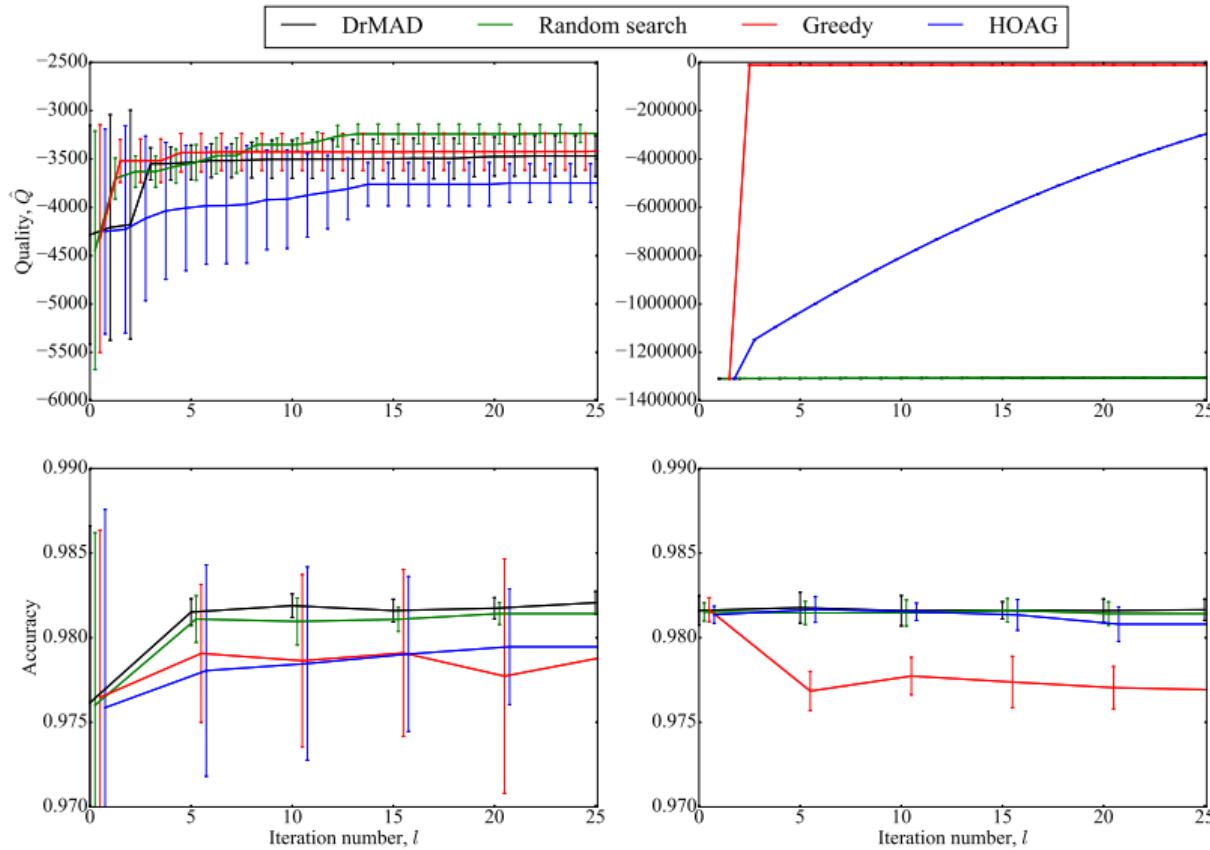


CV



Evidence

# Experiment: MNIST

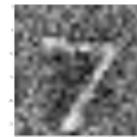


# Experiment: MNIST

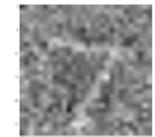
Adding Gaussian noise  $\mathcal{N}(0, \sigma^2 I)$ :



$\sigma = 0.1$

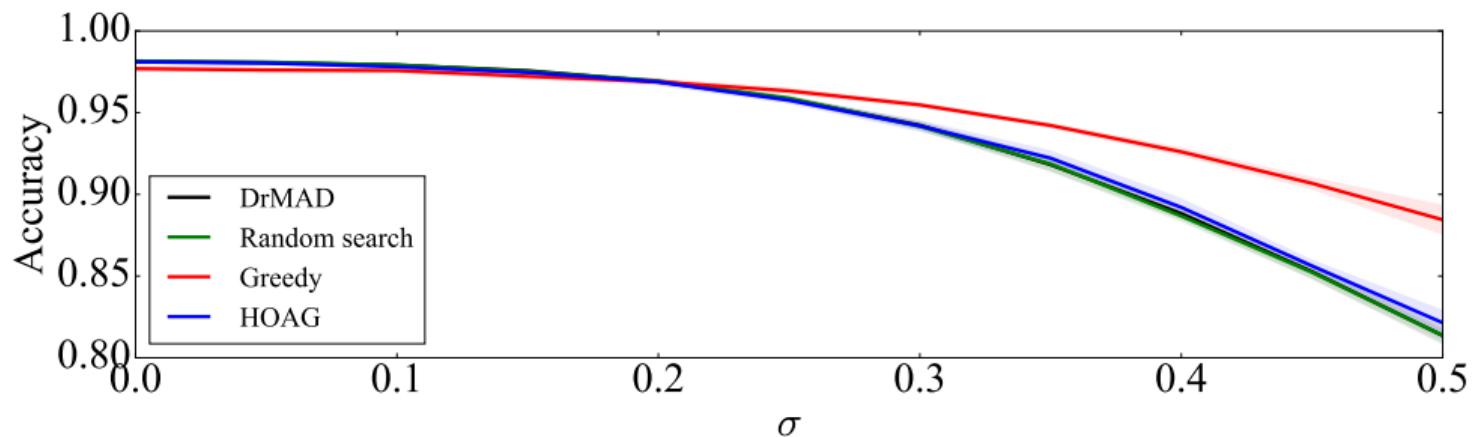


$\sigma = 0.25$

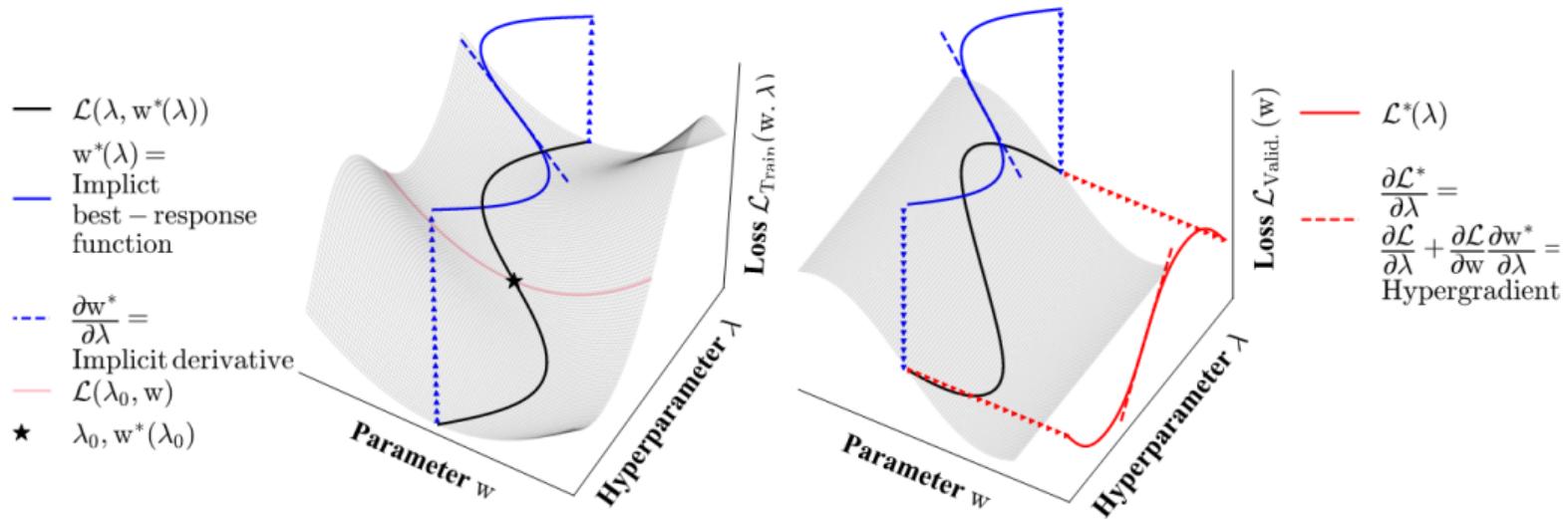


$\sigma = 0.5$

Без шума



# Optimizing Millions of Hyperparameters by Implicit Differentiation



# Optimizing Millions of Hyperparameters by Implicit Differentiation

**Theorem 1** (Cauchy, Implicit Function Theorem). *If for some  $(\lambda', \mathbf{w}')$ ,  $\frac{\partial \mathcal{L}_T}{\partial \mathbf{w}}|_{\lambda', \mathbf{w}'} = 0$  and regularity conditions are satisfied, then surrounding  $(\lambda', \mathbf{w}')$  there is a function  $\mathbf{w}^*(\lambda)$  s.t.  $\frac{\partial \mathcal{L}_T}{\partial \mathbf{w}}|_{\lambda, \mathbf{w}^*(\lambda)} = 0$  and we have:*

$$\frac{\partial \mathbf{w}^*}{\partial \lambda}|_{\lambda'} = - \underbrace{\left[ \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \mathbf{w}} \right]^{-1}}_{\text{training Hessian}} \times \underbrace{\frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \lambda}}_{\text{training mixed partials}} \Big|_{\lambda', \mathbf{w}^*(\lambda')} \quad (\text{IFT})$$

$$\left[ \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \mathbf{w}} \right]^{-1} = \lim_{i \rightarrow \infty} \sum_{j=0}^i \left[ I - \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \mathbf{w}} \right]^j$$

# Optimizing Millions of Hyperparameters by Implicit Differentiation

David Duvenaud, Tom Salimbeni, David Hernández-Lobato

Method	Steps	Eval.	Hypergradient Approximation	
Exact IFT	$\infty$	$\mathbf{w}^*(\boldsymbol{\lambda})$	$\frac{\partial \mathcal{L}_V}{\partial \boldsymbol{\lambda}} - \frac{\partial \mathcal{L}_V}{\partial \mathbf{w}} \times \left[ \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \mathbf{w}} \right]^{-1} \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \boldsymbol{\lambda}} \Big _{\mathbf{w}^*(\boldsymbol{\lambda})}$	
Unrolled Diff. [Maclaurin et al., 2015]	$i$	$\mathbf{w}_0$	$\frac{\partial \mathcal{L}_V}{\partial \boldsymbol{\lambda}} - \frac{\partial \mathcal{L}_V}{\partial \mathbf{w}} \times \sum_{j \leq i} \left[ \prod_{k < j} I - \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \mathbf{w}} \Big _{\mathbf{w}_{i-k}} \right] \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \boldsymbol{\lambda}} \Big _{\mathbf{w}_{i-j}}$	
$L$ -Step Truncated Unrolled Diff.	$i$	$\mathbf{w}_L$	$\frac{\partial \mathcal{L}_V}{\partial \boldsymbol{\lambda}} - \frac{\partial \mathcal{L}_V}{\partial \mathbf{w}} \times \sum_{L \leq j \leq i} \left[ \prod_{k < j} I - \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \mathbf{w}} \Big _{\mathbf{w}_{i-k}} \right] \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \boldsymbol{\lambda}} \Big _{\mathbf{w}_{i-j}}$	
Larsen et al. [1996]	$\infty$	$\widehat{\mathbf{w}}^*(\boldsymbol{\lambda})$	$\frac{\partial \mathcal{L}_V}{\partial \boldsymbol{\lambda}} - \frac{\partial \mathcal{L}_V}{\partial \mathbf{w}} \times \left[ \frac{\partial \mathcal{L}_T \partial \mathcal{L}_T^T}{\partial \mathbf{w} \partial \mathbf{w}} \right]^{-1} \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \boldsymbol{\lambda}} \Big _{\widehat{\mathbf{w}}^*(\boldsymbol{\lambda})}$	
Bengio [2000]	$\infty$	$\widehat{\mathbf{w}}^*(\boldsymbol{\lambda})$	$\frac{\partial \mathcal{L}_V}{\partial \boldsymbol{\lambda}} - \frac{\partial \mathcal{L}_V}{\partial \mathbf{w}} \times \left[ \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \mathbf{w}} \right]^{-1} \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \boldsymbol{\lambda}} \Big _{\widehat{\mathbf{w}}^*(\boldsymbol{\lambda})}$	
$T_1 - T_2$ [27]	1	$\widehat{\mathbf{w}}^*(\boldsymbol{\lambda})$	$\frac{\partial \mathcal{L}_V}{\partial \boldsymbol{\lambda}} - \frac{\partial \mathcal{L}_V}{\partial \mathbf{w}} \times [I]^{-1} \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \boldsymbol{\lambda}} \Big _{\widehat{\mathbf{w}}^*(\boldsymbol{\lambda})}$	
<b>Ours</b>	$i$	$\widehat{\mathbf{w}}^*(\boldsymbol{\lambda})$	$\frac{\partial \mathcal{L}_V}{\partial \boldsymbol{\lambda}} - \frac{\partial \mathcal{L}_V}{\partial \mathbf{w}} \times \left( \sum_{j \leq i} \left[ I - \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \mathbf{w}} \right]^j \right) \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \boldsymbol{\lambda}} \Big _{\widehat{\mathbf{w}}^*(\boldsymbol{\lambda})}$	
Conjugate Gradient (CG) $\approx$	-	$\widehat{\mathbf{w}}^*(\boldsymbol{\lambda})$	$\frac{\partial \mathcal{L}_V}{\partial \boldsymbol{\lambda}} - \left( \arg \min_{\mathbf{x}} \ \mathbf{x} \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \mathbf{w}} - \frac{\partial \mathcal{L}_V}{\partial \mathbf{w}}\  \right) \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \boldsymbol{\lambda}} \Big _{\widehat{\mathbf{w}}^*(\boldsymbol{\lambda})}$	
Hypernetwork	-	-	$\frac{\partial \mathcal{L}_V}{\partial \boldsymbol{\lambda}} + \frac{\partial \mathcal{L}_V}{\partial \mathbf{w}} \times \frac{\partial \mathbf{w}_\phi^*}{\partial \boldsymbol{\lambda}}$ where $\mathbf{w}_\phi^*(\boldsymbol{\lambda}) = \arg \min_\phi \mathcal{L}_T(\boldsymbol{\lambda}, \mathbf{w}_\phi(\boldsymbol{\lambda}))$	
Bayesian Optimization	-	-	$\frac{\partial \mathcal{L}_V^*}{\partial \boldsymbol{\lambda}}$ where $\mathcal{L}_V^* \sim \text{Gaussian-Process}(\{\boldsymbol{\lambda}_i, \mathcal{L}_V(\boldsymbol{\lambda}_i, \mathbf{w}^*(\boldsymbol{\lambda}_i))\})$	

# References

- Bishop C. M. Pattern recognition //Machine learning. – 2006. – T. 128. – №. 9.
- Bakhteev O. Y., Strijov V. V. Comprehensive analysis of gradient-based hyperparameter optimization algorithms //Annals of Operations Research. – 2020. – T. 289. – №. 1. – C. 51-65.
- Graves A. Practical variational inference for neural networks //Advances in neural information processing systems. – 2011. – T. 24.
- Bergstra et al., Random Search for Hyper-Parameter Optimization, 2012
- Dougal Maclaurin et. al, Gradient-based Hyperparameter Optimization through Reversible Learning, 2015
- Jelena Luketina et. al, Scalable Gradient-Based Tuning of Continuous Regularization Hyperparameters, 2016
- Jie Fu et. al, DrMAD: Distilling Reverse-Mode Automatic Differentiation for Optimizing Hyperparameters of Deep Neural Networks, 2016
- Fabian Pedregosa, Hyperparameter optimization with approximate gradient, 2016
- Bobak Shahriari et. al, Taking the Human Out of the Loop: A Review of Bayesian Optimization, 2016
- Liu H., Simonyan K., Yang Y. Darts: Differentiable architecture search //arXiv preprint arXiv:1806.09055. – 2018.
- Lorraine, J., Vicol, P., & Duvenaud, D. (2020, June). Optimizing millions of hyperparameters by implicit differentiation. In International Conference on Artificial Intelligence and Statistics (pp. 1540-1552). PMLR. (see also Victor Pankratov's slides, BMM-2021)