

Bayesian selection of deep learning model structure

Oleg Bakhteev

Supervisor: Prof. Vadim Strijov

Moscow Institute of Physics and Technology
November 21, 2019

Selection of deep learning model structure

Goal : to propose a method of selection of deep learning model structure.

Objectives:

- ① Proposal of suboptimal and optimal complexity criteria for deep learning models.
- ② Proposal of an algorithm suboptimal deep learning model selection and optimization of model parameters.

Investigated problems

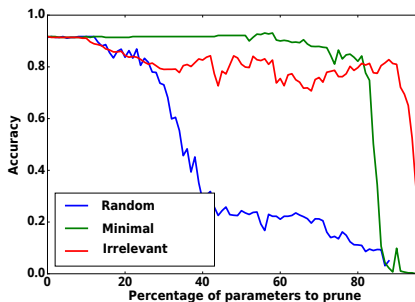
- ① Large number of parameters and hyperparameters, high computational complexity.
- ② Multiextremality and non-convexity of optimization.

Methods

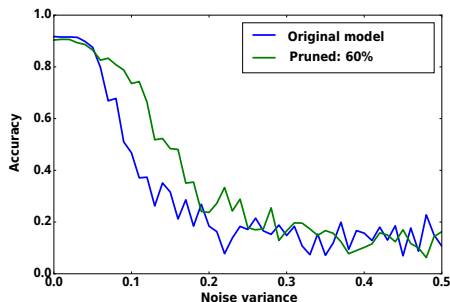
A deep learning model is considered as a multigraph. For the suboptimal model selection we use a composition of methods of automatic relevance determination and hyperparameter gradient optimization methods.

Model structure selection challenge

Data likelihood does not change with removing redundant parameters.



Redundancy of model parameters



Model robustness

Deep learning models have implicitly redundant complexity.

Deep learning model

Definition

Model $\mathbf{f}(\mathbf{w}, \mathbf{x})$ is a differentiable function with respect to parameters \mathbf{w} from the set of object descriptions into the set of labels:

$$\mathbf{f} : \mathbb{X} \times \mathbb{W} \rightarrow \mathbb{Y},$$

where \mathbb{W} is a space of parameters of model \mathbf{f} .

Main challenge of deep learning model selection is in large number of parameters of models. This disallows to use many classical approaches for the model and structure selection (AIC, BIC, cross-validation).

A model is defined by its parameters \mathbf{W} and structure Γ .

A **structure** defines a set of functional superpositions in the model. It is selected using statistical complexity criteria.

Empirical model complexity estimations:

- ① number of parameters;
- ② number of superpositions in the model.

Structure selection: one-layer network

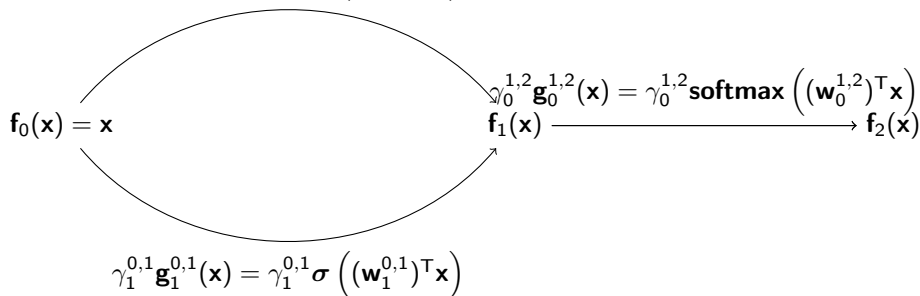
The model \mathbf{f} is defined by the **structure** $\Gamma = [\gamma^{0,1}, \gamma^{1,2}]$.

$$\text{Model: } \mathbf{f}(\mathbf{x}) = \text{softmax} \left((\mathbf{w}_0^{1,2})^\top \mathbf{f}_1(\mathbf{x}) \right), \quad \mathbf{f}(\mathbf{x}) : \mathbb{R}^n \rightarrow [0, 1]^{|Y|}, \quad \mathbf{x} \in \mathbb{R}^n.$$

$$\mathbf{f}_1(\mathbf{x}) = \gamma_0^{0,1} \mathbf{g}_0^{0,1}(\mathbf{x}) + \gamma_1^{0,1} \mathbf{g}_1^{0,1}(\mathbf{x}),$$

where $\mathbf{w} = [\mathbf{w}_0^{0,1}, \mathbf{w}_1^{0,1}, \mathbf{w}_0^{1,2}]^\top$ — parameter matrices, $\{\mathbf{g}_{0,1}^0, \mathbf{g}_{0,1}^1, \mathbf{g}_{1,2}^0\}$ — generalized-linear functions, alternatives of layers of the network.

$$\gamma_0^{0,1} \mathbf{g}_0^{0,1}(\mathbf{x}) = \gamma_0^{0,1} \sigma \left((\mathbf{w}_0^{0,1})^\top \mathbf{x} \right)$$



Deep learning model structure as a graph

Define:

- ① acyclic graph (V, E) ;
- ② for each edge $(j, k) \in E$: a vector primitive differentiable functions $\mathbf{g}^{j,k} = [\mathbf{g}_0^{j,k}, \dots, \mathbf{g}_{K^{j,k}}^{j,k}]$ with length of $K^{j,k}$;
- ③ for each vertex $v \in V$: a differentiable aggregation function \mathbf{agg}_v .
- ④ a function $\mathbf{f} = \mathbf{f}_{|V|-1}$:

$$\mathbf{f}_v(\mathbf{w}, \mathbf{x}) = \mathbf{agg}_v \left(\{ \langle \gamma^{j,k}, \mathbf{g}^{j,k} \rangle \circ \mathbf{f}_j(\mathbf{x}) \mid j \in \text{Adj}(v_k) \} \right), v \in \{1, \dots, |V|-1\}, \quad \mathbf{f}_0(\mathbf{x}) = \mathbf{x} \quad (1)$$

that is a function from \mathbb{X} into a set of labels \mathbb{Y} for any value of $\gamma^{j,k} \in [0, 1]^{K^{j,k}}$.

Definition

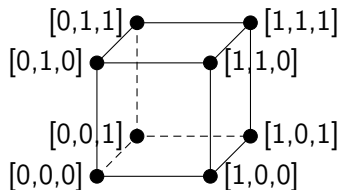
A *parametric set of models* \mathfrak{F} is a graph (V, E) with a set of primitive functions $\{\mathbf{g}^{j,k}, (j, k) \in E\}$ and aggregation functions $\{\mathbf{agg}_v, v \in V\}$.

Statement

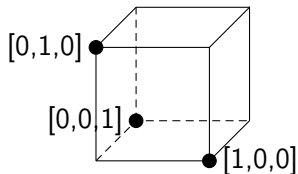
A function $\mathbf{f} \in \mathfrak{F}$ is a model for each $\gamma^{j,k} \in [0, 1]^{K^{j,k}}$.

Structure restrictions

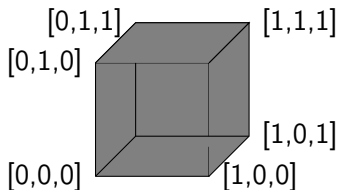
An example of restrictions for structure parameter γ , $|\gamma| = 3$.



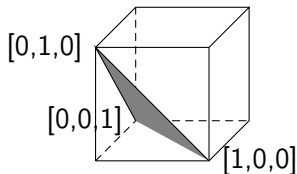
Cube vertices



Simplex vertices



Cube interior



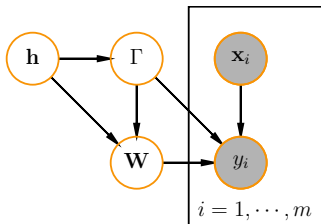
Simplex interior

Prior distribution

Definition

Prior distribution for parameters \mathbf{w} and structure Γ of model \mathbf{f} is a distribution

$p(\mathbf{W}, \Gamma | \mathbf{h}, \lambda) : \mathbb{W} \times \mathbb{\Gamma} \times \mathbb{H} \rightarrow \mathbb{R}^+$, where \mathbb{W} is a parameter space, $\mathbb{\Gamma}$ is a structure space, λ is a vector of metaparameters.



Definition

Hyperparameters $\mathbf{h} \in \mathbb{H}$ are the parameters of prior distribution $p(\mathbf{w}, \Gamma | \mathbf{h}, \mathbf{f})$ (parameters of the distribution of the parameters and structure of model \mathbf{f}).

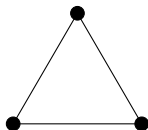
A model \mathbf{f} is defined by:

- **Parameters** $\mathbf{w} \in \mathbb{W}$ that define superpositions \mathbf{f}_v in the model \mathbf{f} .
- **Structure** $\Gamma = \{\gamma^{j,k}\}_{(j,k) \in E} \in \mathbb{\Gamma}$ that define the contribution of all the superpositions \mathbf{f}_v into \mathbf{f} .
- **Hyperparameters** $\mathbf{h} \in \mathbb{H}$ that define the prior distribution.
- **Metaparameters** $\lambda \in \mathbb{A}$ that define the optimization function.

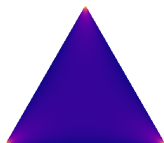
Prior distribution for the model structure

Every point in a simplex defines a model.

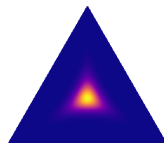
Gumbel-Softmax distribution: $\Gamma \sim \text{GS}(\mathbf{s}, \lambda_{\text{temp}})$



$$\lambda_{\text{temp}} \rightarrow 0$$

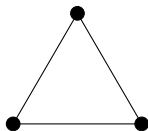


$$\lambda_{\text{temp}} = 0.995$$

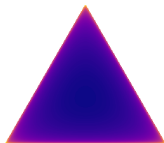


$$\lambda_{\text{temp}} = 5.0$$

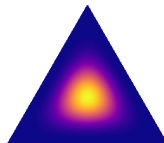
Dirichlet distribution: $\Gamma \sim \text{Dir}(\mathbf{s}, \lambda_{\text{temp}})$



$$\lambda_{\text{temp}} \rightarrow 0$$



$$\lambda_{\text{temp}} = 0.995$$

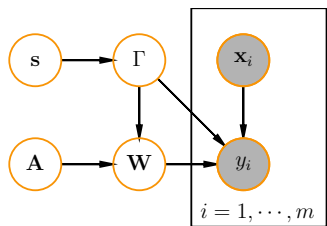


$$\lambda_{\text{temp}} = 5.0$$

Bayesian model selection

Base model:

- parameters
 $\mathbf{w} \sim \mathcal{N}(0, \alpha^{-1})$,
- hyperparameters
 $\mathbf{h} = [\alpha]$.



Proposed model:

- parameters
 $\mathbf{w}_r^{j,k} \sim \mathcal{N}(0, (\gamma_r^{j,k})^2 (\mathbf{A}_r^{j,k})^{-1})$, $\mathbf{A}_r^{j,k}$ is a diagonal matrix for the parameters of the primitive function $\mathbf{g}_r^{j,k}$,
 $(\mathbf{A}_r^{j,k})^{-1} \sim \text{inv-gamma}(\lambda_1, \lambda_2)$,
- structure
 $\Gamma = \{\gamma^{j,k}, (j, k) \in E\}$,
 $\gamma^{j,k} \sim \text{GS}(\mathbf{s}^{j,k}, \lambda_{\text{temp}})$,
- hyperparameters $\mathbf{h} = [\text{diag}(\mathbf{A}), \mathbf{s}]$,
- metaparameters $\lambda_1, \lambda_2, \lambda_{\text{temp}}$.

Evidence as a statistical complexity

Minimum description length for the model f :

$$\text{MDL}(\mathbf{y}, f) = -\log p(\mathbf{h}|f) - \log p(\hat{\mathbf{w}}|\mathbf{h}, f) - \log (p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{w}}, f)\delta\mathcal{D}),$$

where $\delta\mathcal{D}$ is an information transmission precision.

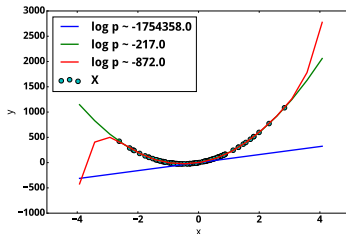
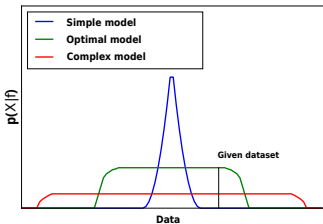
Bayesian approach:

Obtain values of parameters \mathbf{w} with respect to **posterior distribution of parameters**:

$$L = \log p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \mathbf{h}, \lambda) \propto \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{h}, \lambda) + \log p(\mathbf{w}|\mathbf{h}, \lambda).$$

Hyperparameters are optimized using **posterior distribution of hyperparameters**:

$$Q = \log p(f|\mathbf{X}, \mathbf{y}) \propto \log p(\mathbf{h}|f) + \log \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \lambda) p(\mathbf{w}|\mathbf{h}, \lambda) d\mathbf{w}.$$



Evidence lower bound

The evidence is analytically intractable.

Model evidence:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{h}, \lambda) = \int \int_{\mathbf{w}, \Gamma} p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \Gamma) p(\mathbf{w}, \Gamma|\mathbf{h}, \lambda) d\mathbf{w} d\Gamma.$$

Definition

Variational parameters of the model $\theta \in \Theta$ are the parameters of the distribution q that approximates posterior distribution $p(\mathbf{w}, \Gamma|\mathbf{X}, \mathbf{y}, \mathbf{h}, \lambda)$:

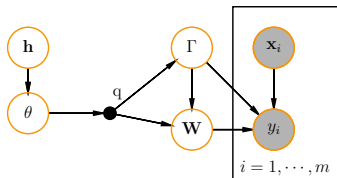
$$q \approx \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \Gamma) p(\mathbf{w}, \Gamma|\mathbf{h}, \lambda)}{\iint_{\mathbf{w}', \Gamma'} p(\mathbf{y}|\mathbf{X}, \mathbf{w}', \Gamma') p(\mathbf{w}', \Gamma'|\mathbf{h}, \lambda) d\mathbf{w}' d\Gamma'}.$$

Lower bound of $\log p(\mathbf{y}|\mathbf{X}, \mathbf{h}, \lambda)$:

$$\log p(\mathbf{y}|\mathbf{X}, \mathbf{h}, \lambda) \geq \mathbb{E}_q \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \Gamma) - D_{\text{KL}}(q(\mathbf{w}, \Gamma) || p(\mathbf{w}, \Gamma|\mathbf{h}, \lambda)).$$

The lower bound equals to evidence when

$$D_{\text{KL}}(q(\mathbf{w}, \Gamma) || p(\mathbf{w}, \Gamma|\mathbf{y}, \mathbf{X}, \mathbf{h}, \lambda)) = 0.$$



Model selection problem

Define a variational distribution $q = q_w q_r$ with parameters θ that approximates posterior distribution $p(\mathbf{w}, \Gamma | \mathbf{X}, \mathbf{y}, \mathbf{h}, \mathbf{f})$.

Definition

Loss function $L(\theta | \mathbf{y}, \mathbf{X}, \mathbf{h}, \lambda)$ is a differentiable function interpreted as a performance of the model on the train dataset.

Validation function $Q(\mathbf{h} | \mathbf{y}, \mathbf{X}, \theta, \lambda)$ is a differentiable function interpreted as a general performance of the model.

The *model selection problem* \mathbf{f} is a level optimization:

$$\mathbf{h}^* = \arg \max_{\mathbf{h} \in \mathbb{H}} Q(\mathbf{h} | \mathbf{y}, \mathbf{X}, \theta^*, \lambda),$$

where θ^* is a solution for the following optimization:

$$\theta^* = \arg \max_{\theta \in \mathbb{U}} L(\theta | \mathbf{y}, \mathbf{X}, \mathbf{h}, \lambda).$$

Generalizing optimization problem

The model selection problem \mathbf{h}^*, θ^* is a generalizing problem on the compact $U_\theta \times U_h \times U_\lambda \subset \mathbb{R}^u \times \mathbb{H} \times \mathbb{A}$, if the following conditions are met:

- ① For each parameter, hyperparameter and metaparameters its domain is not empty and not a point.
- ② For each $\mathbf{h} \in U_h$ and each $\lambda \in U_\lambda$ the solution θ^* is uniquely defined.
- ③ **Continuance:** L, Q are continuous with respect to metaparameters.
- ④ **optimal structure exhaustive search:** there is a constant $K_3 > 0$ and a value for the metaparameters λ such that for all pairs of local optima $\mathbf{h}_1, \mathbf{h}_2$ of Q with metaparameters λ such that

$$D_{\text{KL}}(p(\Gamma|\mathbf{h}_1, \lambda)|p(\Gamma|\mathbf{h}_1, \lambda)) > K_3, D_{\text{KL}}(p(\Gamma|\mathbf{h}_1, \lambda)|p(\Gamma|\mathbf{h}_2, \lambda)) > K_3,$$

$$Q(\mathbf{h}_1|\lambda) > Q(\mathbf{h}_2|\lambda),$$

there exists another value of metaparameters $\lambda' \neq \lambda$ that

- ① the correspondence between optimal variational parameters and hyperparameters $\theta^*(\mathbf{h}_1), \theta^*(\mathbf{h}_2)$ remains for λ' ,
- ② the following inequality is satisfied: $Q(\mathbf{h}_1|\lambda') < Q(\mathbf{h}_2|\lambda')$.

Generalizing optimization problem

The model selection problem \mathbf{h}^*, θ^* is generalizing on the compact $U_\theta \times U_h \times U_\lambda \subset \mathbb{R}^u \times \mathbb{H} \times \Lambda$, if the following conditions are met:

- ⑤ **Likelihood maximization:** there is a metaparameter value $\lambda \in U_\lambda$ and $K_1 \in \mathbb{R}_+$ such that for each pair of hyperparameter vectors $\mathbf{h}_1, \mathbf{h}_2 \in U_h$, $Q(\mathbf{h}_1) - Q(\mathbf{h}_2) > K_1$ the following inequality is satisfied :
$$\mathbb{E}_q \log p(\mathbf{y}|\mathbf{X}, \theta^*(\mathbf{h}_1), \lambda_{\text{temp}}, \mathbf{f}) > \log \mathbb{E}_q p(\mathbf{y}|\mathbf{X}, \theta^*(\mathbf{h}_2), \lambda_{\text{temp}}, \mathbf{f}).$$
- ⑥ **Complexity minimization:** there is a metaparameter value $\lambda \in U_\lambda$ and $K_2 \in \mathbb{R}_+$ such that for each pair of hyperparameter vectors $\mathbf{h}_1, \mathbf{h}_2 \in U_h$, $Q(\mathbf{h}_1) - Q(\mathbf{h}_2) > K_2$, $\mathbb{E}_q \log p(\mathbf{y}|\theta_1, \lambda_{\text{temp}}, \mathbf{f}) = \log \mathbb{E}_q p(\mathbf{y}|\theta_2, \lambda_{\text{temp}}, \mathbf{f})$, the complexity of the first model is less than the second one.
- ⑦ **Evidence lower bound optimization:** there is a metaparameter value λ , such that the optimization is equivalent to the evidence lower bound optimization:
$$\mathbf{h}^* \propto \arg \max \log \mathbb{E}_{q(\mathbf{w}, \Gamma|\theta)} p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \Gamma) - D_{\text{KL}}(q(\mathbf{w}, \Gamma|\theta) || p(\mathbf{w}, \Gamma|\mathbf{h}, \lambda)) + \log p(\mathbf{h}|\lambda),$$

$$\theta^* = \arg \min D_{\text{KL}}(q|p(\mathbf{w}, \Gamma|\mathbf{y}, \mathbf{X}, \mathbf{h}, \lambda)).$$

Model selection problem analysis

Theorem [Bakhteev, 2019]

The following problems are not generalizing:

- ① maximum likelihood criterion: $\max_{\theta} E_q \log p(\mathbf{y}|\mathbf{X}, \theta, \mathbf{h}, \lambda);$
- ② maximum posterior probability criterion: $\max_{\theta} E_q \log p(\mathbf{y}|\mathbf{X}, \theta, \mathbf{f}) p(\theta|\mathbf{h}, \lambda);$
- ③ evidence lower bound maximization:
 $\max_{\mathbf{h}} \max_{\theta} E_q \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \Gamma) - D_{\text{KL}}(p(\mathbf{w}, \Gamma|\mathbf{h}, \lambda) || q(\mathbf{w}, \Gamma|\theta)) + \log p(\mathbf{h}|\mathbf{f});$
- ④ cross-validation: $\max_{\mathbf{h}} E_q \log p(\mathbf{y}_{\text{valid}}|\mathbf{X}_{\text{valid}}, \theta^*),$
 $\theta^* = \arg \max_{\theta} E_q \log p(\mathbf{y}_{\text{train}}|\mathbf{X}_{\text{train}}, \mathbf{h}, \lambda) p(\theta|\mathbf{h}).$
- ⑤ AIC: $\max_{\theta} E_q \log p(\mathbf{y}|\mathbf{X}, \theta, \lambda_{\text{temp}}, \mathbf{f}) - |\theta_i : D_{\text{KL}}(q(w_i)|p(w_i|\Gamma, \mathbf{h}, \lambda) < \lambda|;$
- ⑥ BIC:
 $\max_{\theta} E_q \log p(\mathbf{y}|\mathbf{X}, \theta, \lambda_{\text{temp}}, \mathbf{f}) - \frac{1}{2} \log(|\mathbb{W}| |\theta_i : D_{\text{KL}}(q(w_i)|p(w_i|\Gamma, \mathbf{h}, \lambda) < \lambda|;$
- ⑦ structure exhaustive search:
 $\max_{\Gamma'} \max_{\theta} E_q \log p(\mathbf{y}|\mathbf{X}, \theta, \lambda_{\text{temp}}, \mathbf{f}) \mathbb{I}(q(\Gamma\Gamma = p')),$ where p' is a distribution on a structure (metaparameter).

Proposed optimization problem

Theorem [Bakhtreev, 2019]

The following problem is generalizing:

$$\begin{aligned} \mathbf{h}^* &= \arg \max_{\mathbf{h}} Q = \\ &= \lambda_{\text{likelihood}}^Q \mathbb{E}_{q(\mathbf{w}, \Gamma | \theta^*)} \log p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \Gamma, \mathbf{h}, \lambda) - \\ &\quad - \lambda_{\text{prior}}^Q D_{KL}(q(\mathbf{w}, \Gamma | \theta^*) || p(\mathbf{w}, \Gamma | \mathbf{h}, \lambda)) - \\ &\quad - \sum_{p' \in \mathfrak{P}, \lambda \in \lambda_{\text{Q}}^{\text{struct}}} \lambda D_{KL}(\Gamma | p') + \log p(\mathbf{h} | \lambda), \end{aligned}$$

where

$$\begin{aligned} \theta^* &= \arg \max_{\theta} L = \mathbb{E}_q \log p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \Gamma, \mathbf{h}, \lambda) \\ &\quad - \lambda_{\text{prior}}^Q D_{KL}(q^*(\mathbf{w}, \Gamma) || p(\mathbf{w}, \Gamma | \mathbf{h}, \lambda)). \end{aligned}$$

The proposed optimization generalized different optimization problems: maximum likelihood and evidence lower bound optimization, model complexity increase and decrease, exhaustive structure search.



$$\lambda_{\text{struct}}^Q = [0; 0; 0].$$



$$\lambda_{\text{struct}}^Q = [1; 0; 0].$$



$$\lambda_{\text{struct}}^Q = [1; 1; 0].$$

Bayesian interpretation of the proposed optimization

Theorem, [Bakhteev, 2018]

Define a set of variational distribution $q(\theta)$.

Let $\lambda_{\text{likelihood}}^L = \lambda_{\text{prior}}^L = \lambda_{\text{prior}}^Q = 1$, $\lambda_{\text{struct}}^Q = 0$. Then:

- 1 Solution of the proposed optimization problem obtains a maximum posterior distribution for the hyperparameters with evidence lower bound approximation:
$$\log \hat{p}(\mathbf{y}|\mathbf{X}, \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f}) + \log p(\mathbf{h}|\mathbf{f}) \rightarrow \max_{\mathbf{h}}.$$

- 2 Variational distribution q for the solution approximates posterior distribution $p(\mathbf{w}, \Gamma|\mathbf{y}, \mathbf{X}, \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f})$ in the best way:
$$D_{\text{KL}}(q||p(\mathbf{w}, \Gamma|\mathbf{y}, \mathbf{X}, \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f})) \rightarrow \min_{\theta}.$$

Let q be decomposed into two distributions for parameters \mathbf{w} and structure Γ of the model \mathbf{f} :

$$q = q_{\mathbf{w}} q_{\Gamma}, q_{\Gamma} \approx p(\Gamma|\mathbf{y}, \mathbf{X}, \mathbf{h}, \mathbf{f}), q_{\mathbf{w}} \approx p(\mathbf{w}|\Gamma, \mathbf{y}, \mathbf{X}, \mathbf{h}, \mathbf{f}).$$

If there are values for the variational parameters such that $q(\mathbf{w}) = p(\mathbf{w}|\Gamma, \mathbf{h}, \lambda)$, $q(\Gamma) = p(\Gamma|\mathbf{h}, \lambda)$, then the solution of optimization of L is equal to these values.

Optimization operator

Definition

An *optimization operator* T is an estimation of the new vector of parameters θ' using the previous one θ .

Stochastic gradient descent operator :

$$\begin{aligned}\hat{\theta} &= T \circ T \circ \dots \circ T(\theta_0, \mathbf{h}) = T^\eta(\theta_0, \mathbf{h}), \quad \text{где } T(\theta, \mathbf{h}) = \\ &= \theta - \lambda_{lr} \nabla (-L(\theta, \mathbf{h})|_{\hat{\mathcal{D}}}),\end{aligned}$$

λ_{lr} is a learning rate, θ_0 is an initial state for θ , $\hat{\mathcal{D}}$ is a random subsample of the dataset \mathcal{D} .

Reformulate the optimization problem:

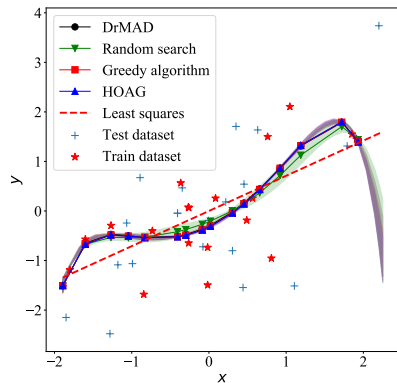
$$\mathbf{h}' = T^\eta(Q, \mathbf{h}, T^\eta(L, \theta_0, \mathbf{h})).$$

Theorem, [Bakhteev, 2019]

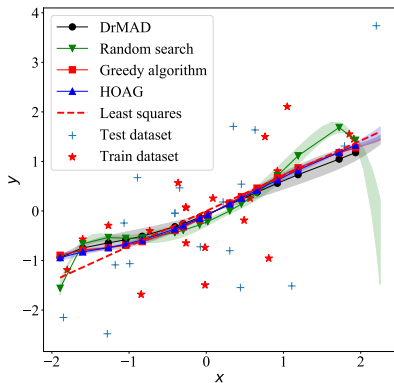
Let $\frac{\lambda_{\text{prior}}^Q}{\lambda_{\text{likelihood}}^Q} = \lambda_{\text{prior}}^L$. Then the proposed optimization is an one-level optimization.

Hyperparameter optimization: example

The hyperparameter gradient-based optimization methods were investigated.

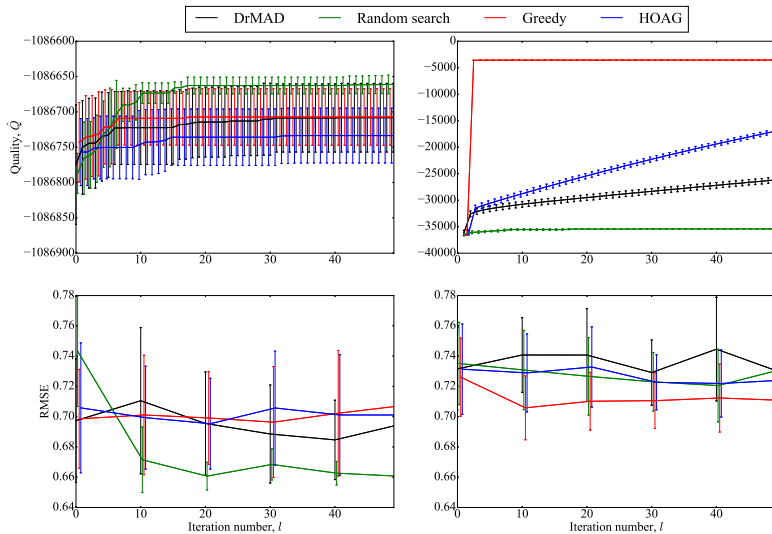


Cross-validation

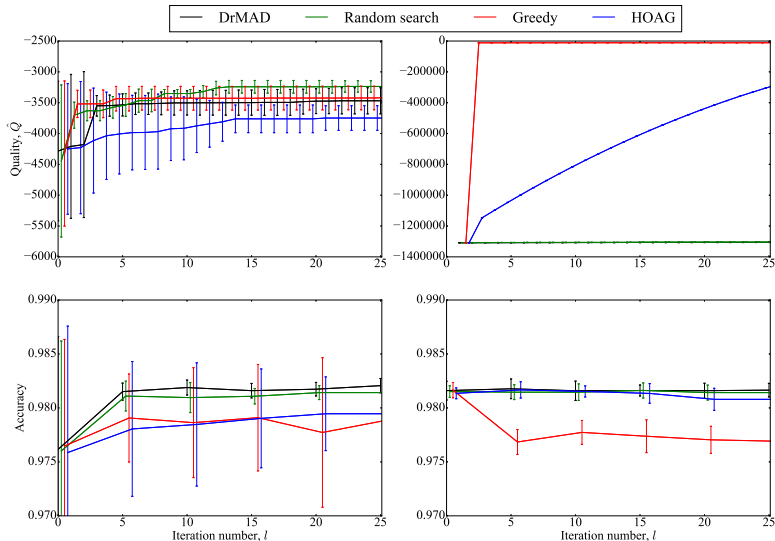


Evidence lower bound

Experiments: WISDM



Experiments: MNIST



Experiments: MNIST

Noise adjustment $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$:



Original images



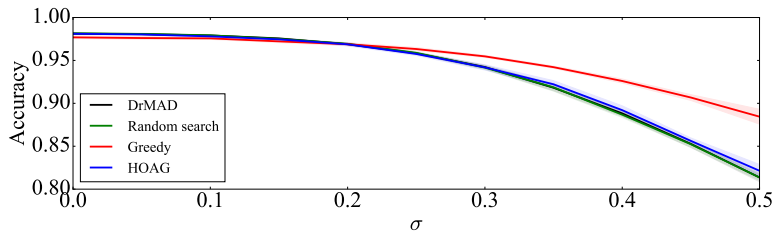
$\sigma = 0.1$



$\sigma = 0.25$



$\sigma = 0.5$



Evidence lower bound using multi-start

$$\log p(\mathbf{y}|\mathbf{X}, \mathbf{h}, \mathbf{f}) \geq \mathbb{E}_{q(\mathbf{w})} \log p(\mathbf{y}, \mathbf{w}|\mathbf{X}, \mathbf{h}, \mathbf{f}) - \mathbb{E}_{q_{\mathbf{w}}}(-\log(q_{\mathbf{w}})).$$

Theorem [Bakhteev, 2016]

Let L be a loss function with continuously-differentiable gradient with Lipschitz constant C .

Let $\boldsymbol{\theta} = [\mathbf{w}^1, \dots, \mathbf{w}^k]$ be a vector of initial states of multiple model optimizations, λ_{lr} is a learning rate.

Then the difference of differentiable entropies for the optimization step can be estimated:

$$\mathbb{E}_{q_{\mathbf{w}}^{\tau}}(-\log(q_{\mathbf{w}}^{\tau})) - \mathbb{E}_{q_{\mathbf{w}}^{\tau-1}}(-\log(q_{\mathbf{w}}^{\tau-1})) \approx \frac{1}{k} \sum_{r=1}^k (\lambda_{lr} \text{Tr}[\mathbf{H}(\mathbf{w}^r)] - \lambda_{lr}^2 \text{Tr}[\mathbf{H}(\mathbf{w}^r)\mathbf{H}(\mathbf{w}^r)]),$$

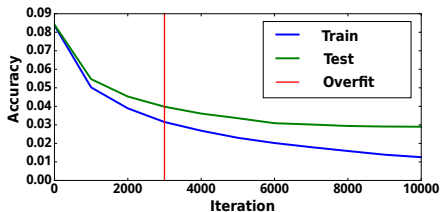
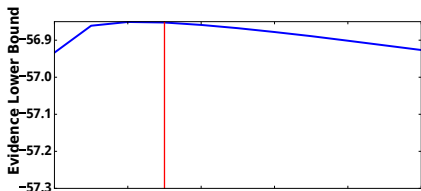
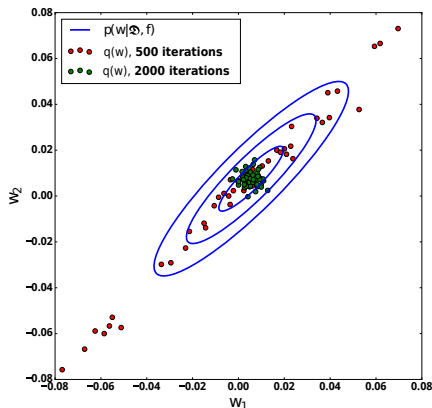
where \mathbf{H} is a Hessian of the negative loss function $-L$, $q_{\mathbf{w}}^{\tau}$ is a distribution q at the iteration τ .

Gradient descent as an evidence lower bound

Empirical distribution of the optimized model parameters is a variational distribution.

Gradient descent does not optimize evidence lower bound.

Evidence lower bound decrease is a signal of overfitting.



Proposed optimization analysis

Theorem, [Bakhteev, 2018]

Let $\lambda_{\text{prior}}^L > 0, m \gg 0, \frac{m}{\lambda_{\text{prior}}^L} \in \mathbb{N}$. Then optimization of

$$L = E_q \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f}) - \lambda_{\text{prior}}^L D_{\text{KL}}(q||p(\mathbf{w}, \mathbf{\Gamma}|\mathbf{h}, \lambda_{\text{temp}}, \mathbf{f})))$$

is a minimization of $E_{\hat{\mathbf{X}}, \hat{\mathbf{y}} \sim p(\mathbf{X}, \mathbf{y})} D_{\text{KL}}(q||p(\mathbf{w}, \mathbf{\Gamma}|\hat{\mathbf{X}}, \hat{\mathbf{y}}, \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f})))$, where $\hat{\mathbf{X}}, \hat{\mathbf{y}}$ is a random sample of size $\frac{m}{\lambda_{\text{prior}}^L}$.

Definition

Parametric complexity of the model is a minimal divergence:

$$C_p = \min_h D_{\text{KL}}(q||p(\mathbf{w}, \mathbf{\Gamma}|\mathbf{h}, \lambda_{\text{temp}}, \mathbf{f}))).$$

Theorem, [Bakhteev, 2018]

Let $\lambda_{\text{struct}}^Q = 0$. Let $\theta_1, \theta_2, \mathbf{h}_1, \mathbf{h}_2$ are the optimization solutions for different metaparameter values $\lambda_{\text{prior}_1}^Q, \lambda_{\text{prior}_2}^Q, \lambda_{\text{prior}_1}^Q > \lambda_{\text{prior}_2}^Q$ on a compact U . Let function $Q(\mathbf{h}|\mathbf{y}, \mathbf{X}, \theta, \lambda)$ be concave on U for $\lambda_{\text{prior}_2}^Q$. Then:

$$C_p(\theta_1|U_{\mathbf{h}}, \lambda_1) - C_p(\theta_2|U_{\mathbf{h}}, \lambda_2) < \frac{\lambda_{\text{prior}}^L}{\lambda_{\text{prior}_2}^Q} (\lambda_{\text{prior}_2}^Q - \lambda_{\text{prior}}^L) C,$$

where C is a constant.

Proposed optimization analysis

Definition

Relative variational density is a ratio:

$$\rho(w|\Gamma, \theta_w, h, \lambda) = \frac{q_w(\text{mode } p(w|\Gamma, h, \lambda))}{q_w(\text{mode } q_w)}.$$

Theorem, [Bakhteev, 2018]

Given $U_h \subset \mathbb{H}$, $U_{\theta_w} \subset \Theta_w$, $U_{\theta_r} \subset \Theta_r$, variational and prior distributions $q_w(w|\Gamma, \theta_w)$, $p(w|\Gamma, h, \lambda)$ are absolutely continuous and unimodal U_θ with equality of mode and mean. Let mode and mean of prior distribution be independent on the hyperparameters h and the structure Γ .

Given a infinite sequence $\theta[1], \theta[2], \dots, \theta[i], \dots \in U_\theta$ such that $\lim_{i \rightarrow \infty} C_p(\theta[i]|U_h, \lambda) = 0$. Then

$$\lim_{i \rightarrow \infty} E_{q_r(r|\theta_r[i])} \rho(w|\Gamma, \theta_w[i], h[i], \lambda)^{-1} = 1, h[i] = \arg \min D_{KL}(q(w, \Gamma|\theta_i) || p(w, \Gamma|h, \lambda)).$$

Main results

The following results were proposed:

- ① method of Bayesian selection of suboptimal structure;
- ② optimal and suboptimal complexity criteria;
- ③ deep learning model graph description;
- ④ generalizing function that includes other methods of model selection:
 - ▶ evidence lower bound;
 - ▶ sequential complexity increase;
 - ▶ sequential complexity decrease;
 - ▶ structure exhaustive search;
- ⑤ method of evidence lower bound optimization based on multistart model optimization;
- ⑥ algorithm of optimization hyperparameters, structure and parameters for deep learning model.
- ⑦ The properties of the proposed optimization were investigated and comprehensively analyzed.

Publications

Main publications

- 1 Bakhteev, O., Kuznetsova, R., Romanov, A. and Khritankov, A. A monolingual approach to detection of text reuse in Russian-English collection // In 2015 Artificial Intelligence and Natural Language and Information Extraction, Social Media and Web Search FRUCT Conference (AINL-ISMW FRUCT) (pp. 3-10). IEEE.
- 2 Бахтеев О.Ю., Попова М.С., Стрижов В.В. Системы и средства глубокого обучения в задачах классификации. // Системы и средства информатики. 2016. № 26.2. С. 4-22.
- 3 Romanov, A., Kuznetsova, R., Bakhteev, O. and Khritankov, A. Machine-Translated Text Detection in a Collection of Russian Scientific Papers. // Computational Linguistics and Intellectual Technologies. 2016.
- 4 Bakhteev, O. and Khazov, A., 2017. Author Masking using Sequence-to-Sequence Models // In CLEF (Working Notes). 2017.
- 5 Бахтеев О.Ю., Стрижов В.В. Выбор моделей глубокого обучения субоптимальной сложности. // Автоматика и телемеханика. 2018. №8. С. 129-147.
- 6 Огальцов А.В., Бахтеев О.Ю. Автоматическое извлечение метаданных из научных PDF-документов. // Информатика и её применения. 2018.
- 7 Смердов А.Н., Бахтеев О.Ю., Стрижов В.В. Выбор оптимальной модели рекуррентной сети в задачах поиска парафраза. // Информатика и ее применения. 2019.
- 8 Грабовой А.В., Бахтеев О.Ю., Стрижов В.В. Определение релевантности параметров нейросети. // Информатика и её применения. 2019.
- 9 Bakhteev O., Strijov V. Comprehensive analysis of gradient-based hyperparameter optimization algorithms // Annals of Operations Research. 2019.

Conference talks

- 1 "Восстановление панельной матрицы и ранжирующей модели в разнородных шкалах", Всероссийская конференция «57-я научная конференция МФТИ», 2014.
- 2 "Выбор модели глубокого обучения субоптимальной сложности с использованием вариационной оценки правдоподобия", Международная конференция «Интеллектуализация обработки информации», 2016.
- 3 "Градиентные методы оптимизации гиперпараметров моделей глубокого обучения", Всероссийская конференция «Математические методы распознавания образов ММРО», 2017.
- 4 "Детектирование переводных заимствований в текстах научных статей из журналов, входящих в РИНЦ", Всероссийская конференция «Математические методы распознавания образов ММРО», 2017.
- 5 "Байесовский выбор наиболее правдоподобной структуры модели глубокого обучения", Международная конференция «Интеллектуализация обработки информации», 2018.