

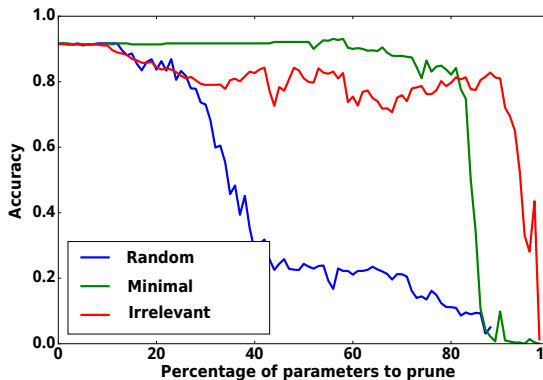
# **Bayesian selection of deep learning model structure**

Oleg Bakhteev, Vadim Strijov

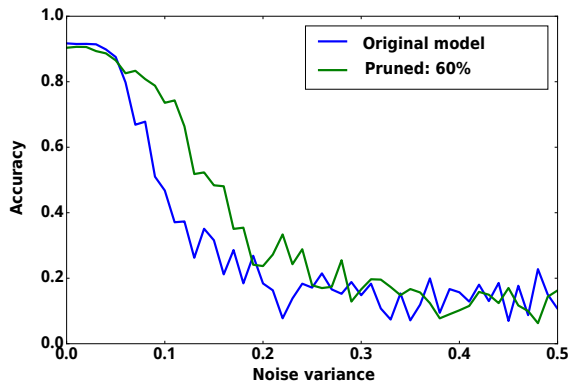
Moscow Institute of Physics and Technology  
July 2, 2021

# 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  $\Gamma$ .

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

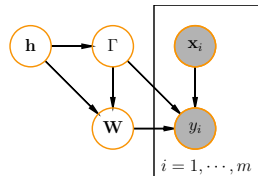
**Empirical model complexity estimations:**

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

# Prior distribution

## Definition

*Prior distribution* for parameters  $\mathbf{w}$  and structure  $\Gamma$  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,  $\lambda$  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.

# Model parameters optimization

## Log likelihood:

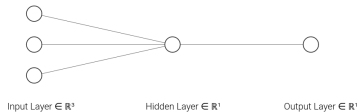
- $\log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{h}, \mathbf{f}) = \sum_{y, \mathbf{x}} \log p(\mathbf{f}[y]|\mathbf{x}, \mathbf{w}, \mathbf{h}, \mathbf{f})$  for classification problem, where  $\mathbf{f}[y]$  is the component  $y$  of the function  $y \in \mathcal{N}$ .
- $\log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{h}, \mathbf{f}) = (-\frac{1}{2}(\mathbf{f}(\mathbf{X}) - \mathbf{y})\mathbf{B}^{-1}(\mathbf{f}(\mathbf{X}) - \mathbf{y})^T) - \frac{1}{2}|\mathbf{B}| + C$ , for regression problem, where  $\mathbf{B}$  is a hyperparameter,  $\mathbf{B} \in \mathbf{h}$ ,  $C$  is a constant.

## Prior:

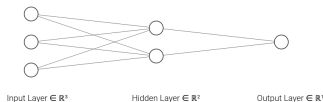
- $p(\mathbf{w}|\mathbf{h}) \sim \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1}) : \log p(\mathbf{w}|\mathbf{h}) = (-\frac{1}{2}(\mathbf{w}\mathbf{A}^{-1}\mathbf{w}^T)) - \frac{1}{2}|\mathbf{A}| + C$ , where  $\mathbf{A}$  is a hyperparameter,  $\mathbf{A} \in \mathbf{h}$ .

# Structure selection example

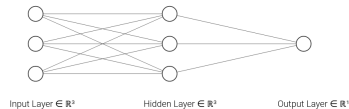
hidden layer dim = 1



hidden layer dim = 2



hidden layer dim = 3



All these models can be represented as  $\mathbf{f}(\mathbf{x}, \mathbf{w}) = \sigma \left( (\mathbf{w}^2)^T \sigma \left( (\mathbf{w}^1)^T \mathbf{x} \right) \right)$   
with similar shape of  $\mathbf{w}^1$  :  $\dim(\mathbf{w}^1) = 3 \times 3$ .

# Structure selection: one-layer network

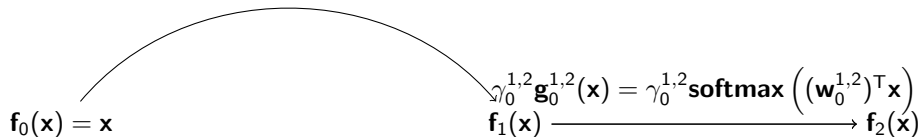
The model  $\mathbf{f}$  is defined by the **structure**  $\mathbf{\Gamma} = [\gamma^{0,1}, \gamma^{1,2}]$ .

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

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

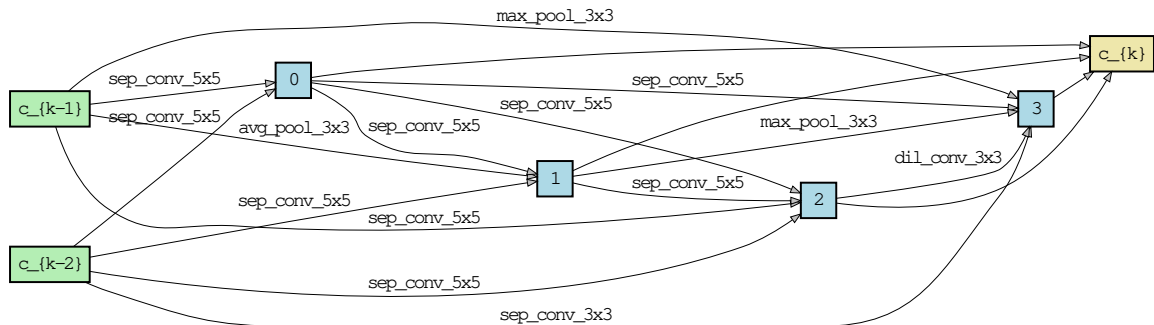
where  $\mathbf{w} = [\mathbf{w}_0^{0,1}, \mathbf{w}_1^{0,1}, \mathbf{w}_0^{1,2}]^\top$  — parameter matrices,  $\{\mathbf{g}_{0,1}^0, \mathbf{g}_{0,1}^1, \mathbf{g}_{1,2}^0\}$  — generalized-linear functions, alternatives of layers of the network.

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



$$\gamma_1^{0,1} \mathbf{g}_1^{0,1}(\mathbf{x}) = \gamma_1^{0,1} \sigma \left( (\mathbf{w}_1^{0,1})^\top \mathbf{x} \right)$$

# Neural architecture search example





# Structure selection: neural architecture search space

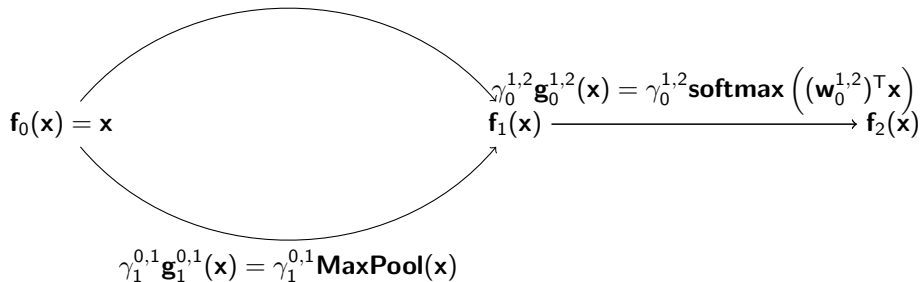
The model  $\mathbf{f}$  is defined by the **structure**  $\Gamma = [\gamma^{0,1}, \gamma^{1,2}]$ .

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

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

where  $\mathbf{w} = [\mathbf{w}_0^{0,1}, \mathbf{w}_0^{1,2}]^\top$  — parameter matrices,  $\mathbf{g}_{0,1}^0$  is a convolution,  $\mathbf{g}_{0,1}^1$  is a pooling operation,  $\mathbf{g}_{1,2}^0$  is a generalized-linear function.

$$\gamma_0^{0,1} \mathbf{g}_0^{0,1}(\mathbf{x}) = \gamma_0^{0,1} \text{Conv}(\mathbf{x}, \mathbf{w}_0^{0,1})$$



# 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

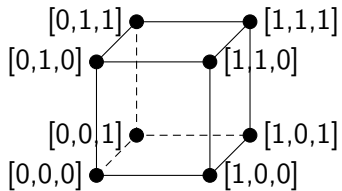
A *parametric set of models*  $\mathfrak{F}$  is a graph  $(V, E)$  with a set of primitive functions  $\{\mathbf{g}^{j,k}, (j, k) \in E\}$  and aggregation functions  $\{\mathbf{agg}_v, v \in V\}$ .

## Statement

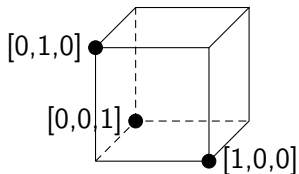
A function  $\mathbf{f} \in \mathfrak{F}$  is a model for each  $\gamma^{j,k} \in [0, 1]^{K^{j,k}}$ .

# Structure restrictions

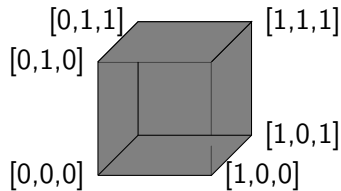
An example of restrictions for structure parameter  $\gamma$ ,  $|\gamma| = 3$ .



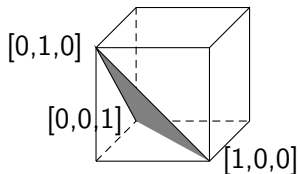
Cube vertices



Simplex vertices



Cube interior

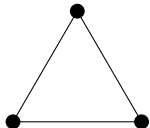


Simplex interior

# Prior distribution for the model structure

Every point in a simplex defines a model.

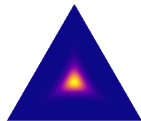
**Gumbel-Softmax distribution:**  $\boldsymbol{\Gamma} \sim \text{GS}(\mathbf{s}, \lambda_{\text{temp}})$



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

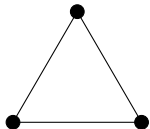


$$\lambda_{\text{temp}} = 0.995$$



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

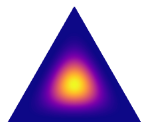
**Dirichlet distribution:**  $\boldsymbol{\Gamma} \sim \text{Dir}(\mathbf{s}, \lambda_{\text{temp}})$



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



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



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

# Bayesian model selection

- **parameters**

$\mathbf{w}_r^{j,k} \sim \mathcal{N}(0, (\mathbf{A}_r^{j,k})^{-1})$ ,  $\mathbf{A}_r^{j,k}$  is a diagonal matrix for the parameters of the primitive function  $\mathbf{g}_r^{j,k}$ ,

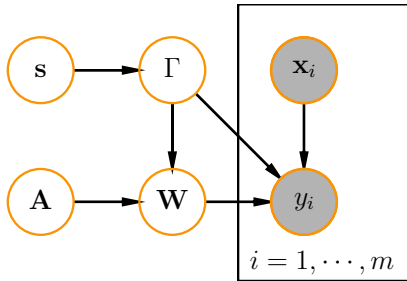
- **structure**

$\Gamma = \{\gamma^{j,k}, (j, k) \in E\}$ ,

$\gamma^{j,k} \sim \text{GS}(\mathbf{s}^{j,k}, \lambda_{\text{temp}})$ ,

- **hyperparameters**  $\mathbf{h} = [\text{diag}(\mathbf{A}), \mathbf{s}]$ ,

- **metaparameters**  $\lambda_{\text{temp}}$ .



# Evidence as a statistical complexity

**Minimum description length** for the model  $\mathbf{f}$ :

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

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

**Bayesian approach:**

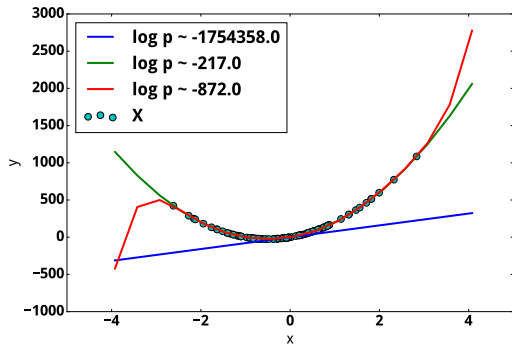
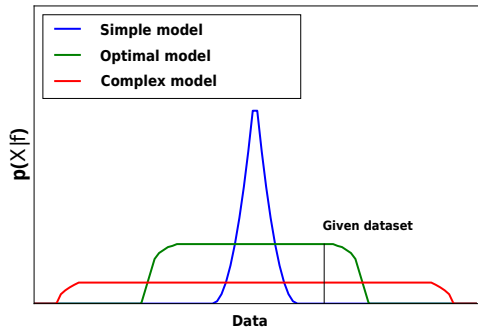
Obtain values of parameters  $\mathbf{w}$  with respect to **posterior distribution of parameters**:

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

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

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

# Evidence: example



# Evidence lower bound

The evidence is analytically intractable.

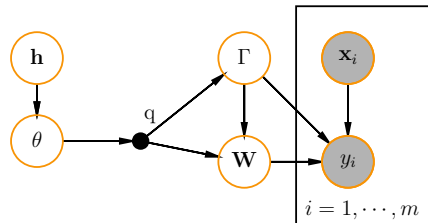
**Model evidence:**

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

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

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

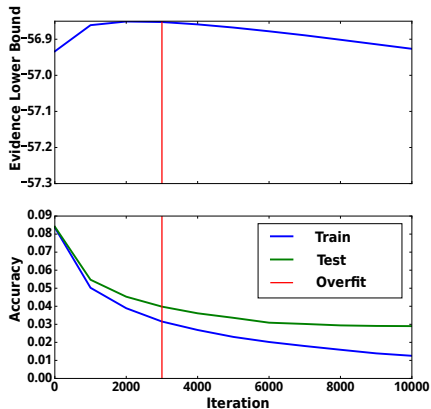
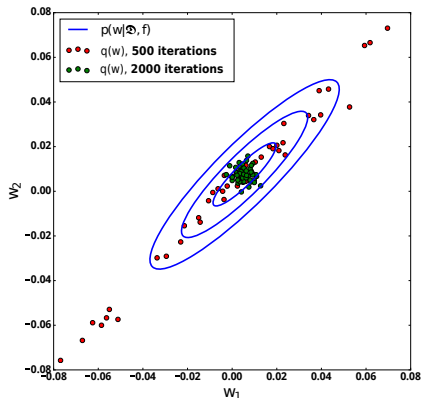


# Gradient descent as an evidence lower bound

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.



# Model selection problem

Define a variational distribution  $q = q_w q_r$  with parameters  $\theta$  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  $\theta^*$  is a solution for the following optimization:

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

# Proposed optimization problem

## Theorem [Bakhtreev, 2019]

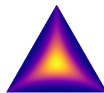
The following problem is generalizing:

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

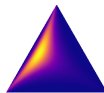
where

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

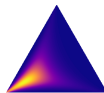
The proposed optimization generalized different optimization problems: maximum likelihood and evidence lower bound optimization, model complexity increase and decrease, exhaustive structure search.



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



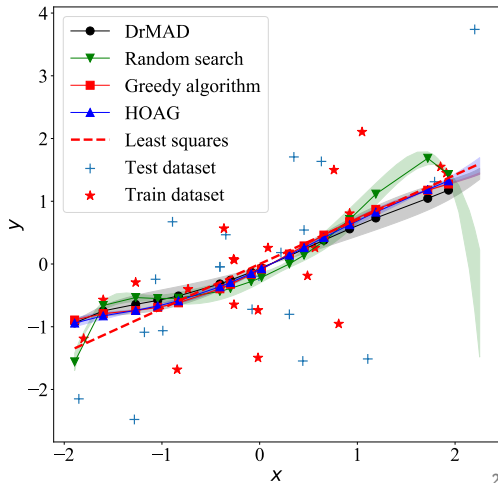
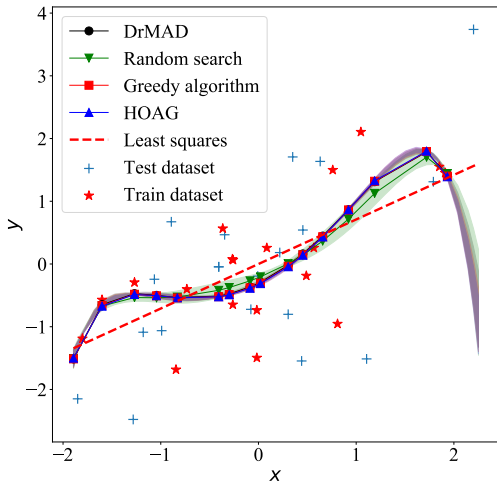
$$\lambda_{\text{struct}}^Q = [1; 0; 0].$$



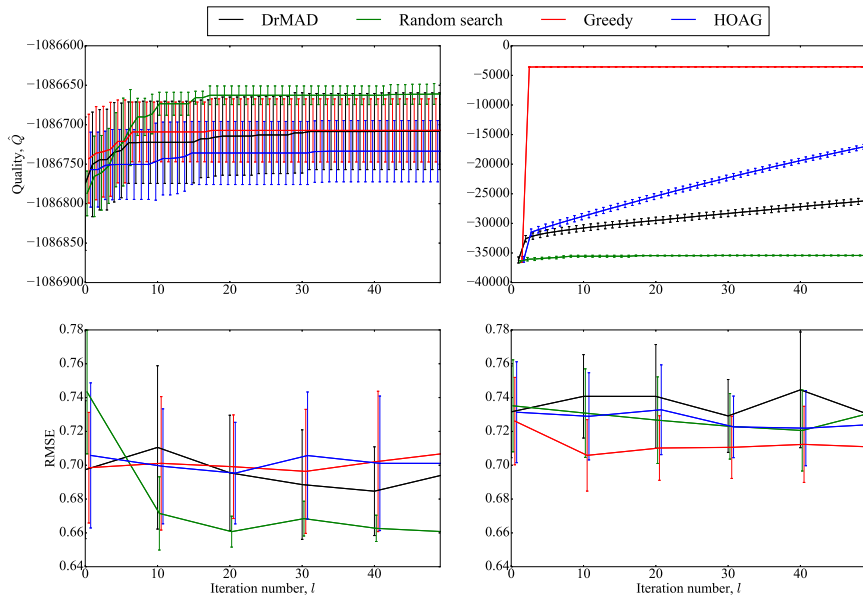
$$\lambda_{\text{struct}}^Q = [1; 1; 0].$$

# Hyperparameter optimization: example

Toy example: polynomial regression with potential overfitting.



# Experiments: WISDM



# Toy example

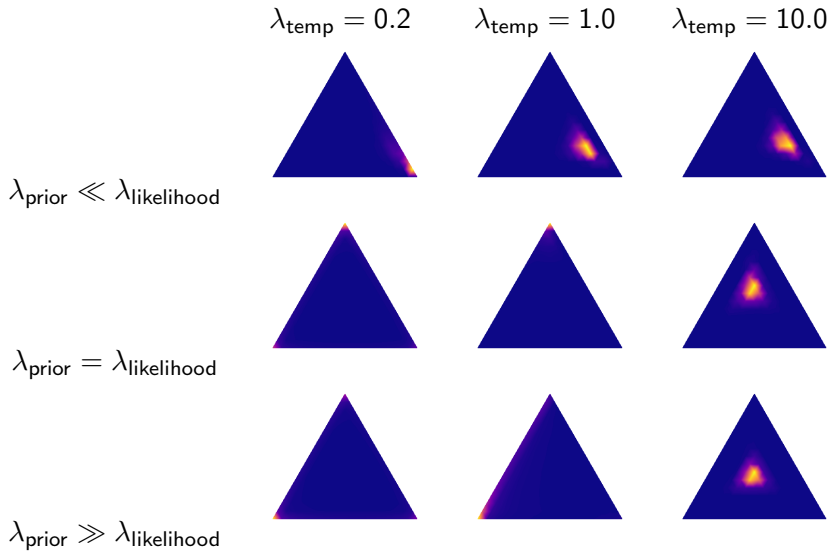
A model  $\mathbf{f}$  is an ensemble of 3 models:

- ①  $\mathbf{g}_0^{0,1} = \tanh(wx);$
- ②  $\mathbf{g}_1^{0,1} = \tanh(\mathbf{w}^T[x, x^2, \dots, x^{10}]);$
- ③  $\mathbf{g}_2^{0,1} = w.$

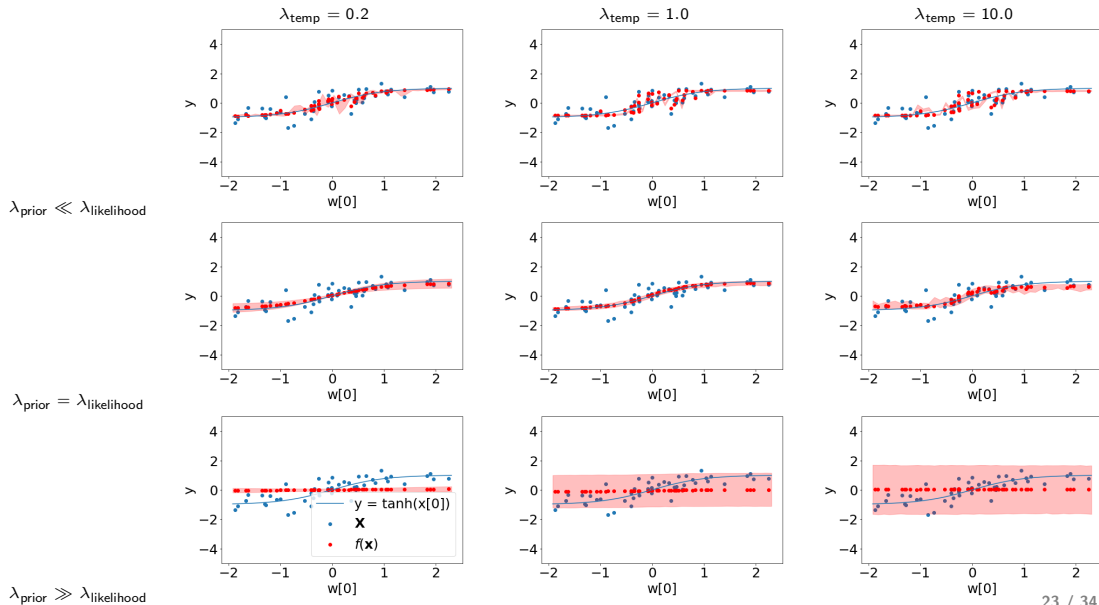
The optimization ran with three regimes:

- $\lambda_{\text{prior}} \ll \lambda_{\text{likelihood}};$
- $\lambda_{\text{prior}} = \lambda_{\text{likelihood}};$
- $\lambda_{\text{prior}} \gg \lambda_{\text{likelihood}};$

# Toy dataset: structures



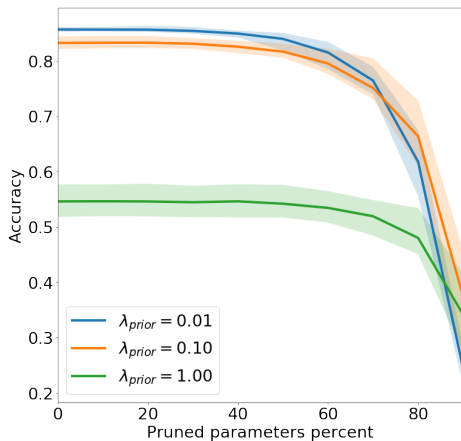
# Toy dataset: prediction performance





## Example

$\lambda_{\text{prior}}$  controls the importance of the prior distribution. With its increasing the model complexity decreases.



Grebenkova, Bakhteev, Strijov. Hypernetworks for deep model complexity control, 2021. (*work in progress*).

# Current challenge

- Can we control the model complexity at the inference step?
- Can we select robust architecture? What properties should it have?

# Model complexity control

## Hypernetworks

A hypernetwork is a mapping from a set of variables responsible for the properties of a desired model to a set of its parameters.

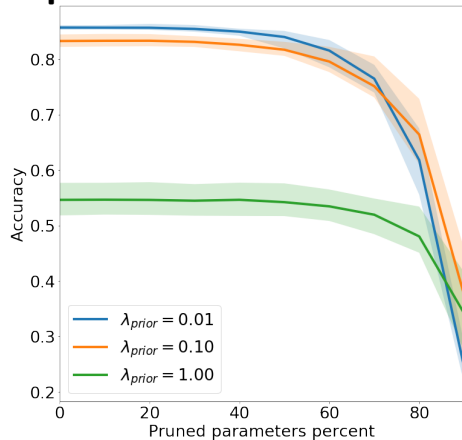
Optimize the model with hypernetworks in the following optimization procedure:

$$\mathbb{E}_{\lambda \sim P(\lambda)} (\log p(\mathcal{D} | \mathbf{w}(\lambda))) - \lambda D_{\text{KL}}(q(\mathbf{w}(\lambda)) || p(\mathbf{w})) \rightarrow \max.$$

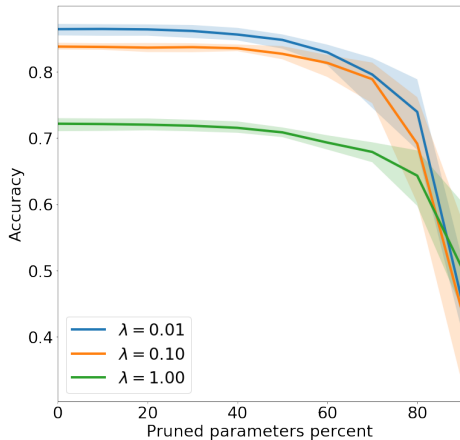
## Theorem, Grebenkova 2021

The hypernetwork approximates not only deep learning model's performance, but also it's statistical propoerties.

## Example: CIFAR-10



CNN



CNN with hypernetwork

Grebenkova, Bakhteev, Strijov. Hypernetworks for deep model complexity control, 2021. (*work in progress*).

# Architecture complexity control

The hypernetworks can approximate not only the model parameter  $\mathbf{w}$ , but also structural parameters  $\gamma$ .

## Baseline: DARTS

A model architecture is a directed graph with non-linear operations  $\mathbf{f}^{(i,j)}$  that are induced by basic functions  $\mathbf{g}^{(i,j)}$  with weights obtained by softmax function application:

$$\mathbf{f}^{(i,j)}(\mathbf{x}) = \langle \text{softmax}(\gamma^{(i,j)}), \mathbf{g}^{(i,j)}(\mathbf{x}) \rangle$$

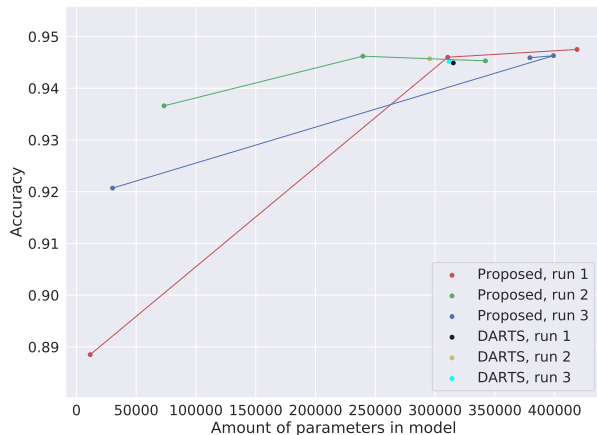
## Our proposal

To use a mapping  $\gamma(\lambda_n)$  instead of constant structural parameters  $\gamma(\lambda_n)$ , where  $\lambda_n$  is a regularization term for the loss function:

$$E_{\lambda_n} \left( \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \Gamma(\lambda_n)) + \lambda_n \sum_{(i,j)} \left\langle \text{softmax} \left( \frac{\gamma(\lambda_n)^{(i,j)}}{\lambda_{\text{temp}}} \right), \mathbf{n}(\mathbf{g}^{(i,j)}) \right\rangle \right),$$

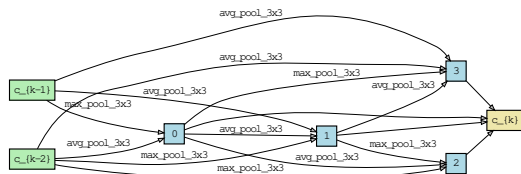
where  $\mathbf{n}(\mathbf{g}^{(i,j)})$  is a vector of amount of parameters for all the basic functions  $\mathbf{g}$ .

# Example: Fashion-MNIST

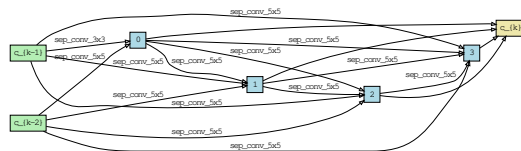


Yakovlev, Grebenkova, Bakhteev, Strijov. Automated architecture search with model complexity control, 2021. (*work in progress*).

# Example



Simple CNN cell architecture



Complex CNN cell architecture

Yakovlev, Grebenkova, Bakhteev, Strijov. Automated architecture search with model complexity control, 2021.  
(*work in progress*).

# Robust architecture

- Robustness to noise in data
  - ▶ Random noise
  - ▶ Adversarial attacks
- Robustness to model modification
  - ▶ Parameters modification
  - ▶ Structure modification



# Experiments: MNIST

Correct hyperparameter optimization leads to the model robustness under noise adjustment:  $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ .



Original images



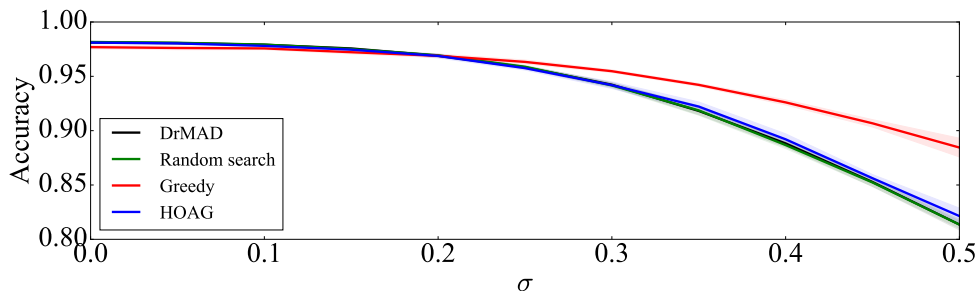
$\sigma = 0.1$



$\sigma = 0.25$

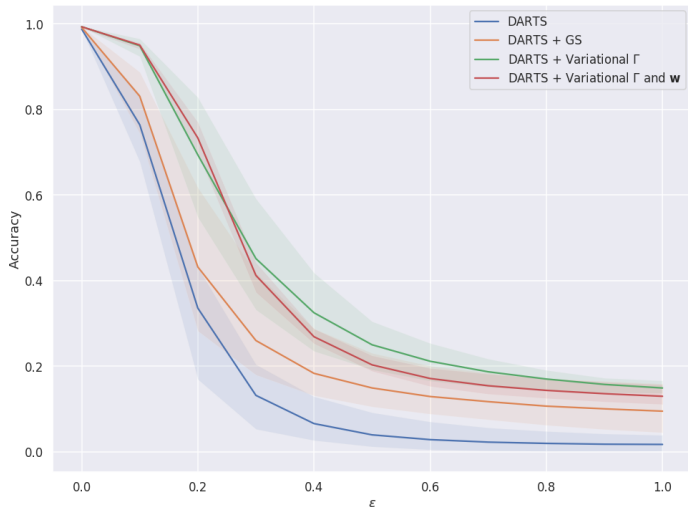


$\sigma = 0.5$

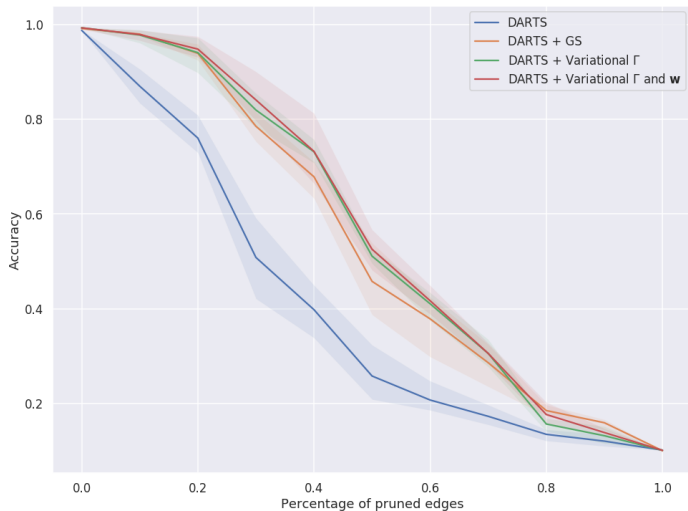


# Robustness to adversarial attacks

FGSM-method:  $\hat{\mathbf{x}} = \mathbf{x} + \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} \log p(y|\mathbf{x}, \mathbf{w}, \Gamma, \mathbf{f}))$ .



# Robustness to structure pruning



Simple CNN cell architecture

# References