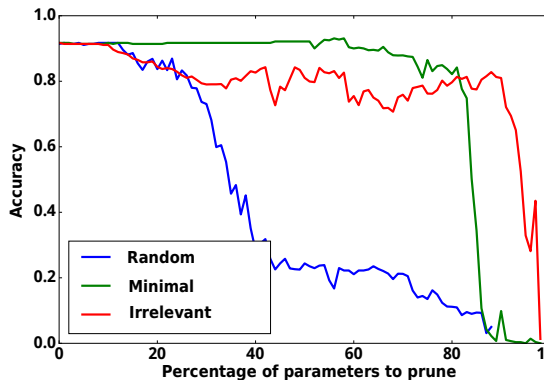# Bayesian selection of
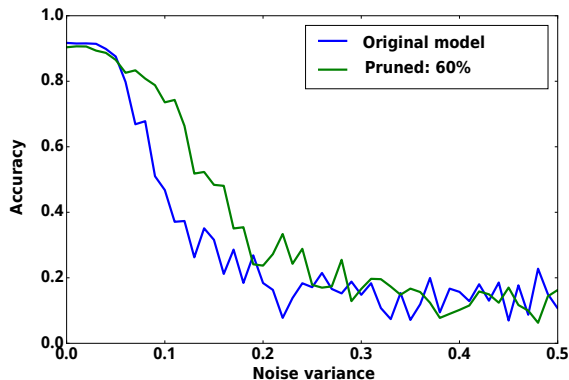# deep learning model structure

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

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

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

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

A model is defined by its parameters W and structure Γ.
A **structure** defines a set of functional superpositions in the model. It is selected using statistical complexity criteria.
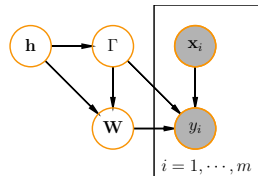**Empirical model complexity estimations:**

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

# Prior distribution



**Definition**

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

**Definition**

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

A model f is defined by:

- **Parameters** $w \in \mathbb{W}$ that define superpositions $f_v$ in the model f.
- **Structure** $\Gamma = \{\gamma^{j,k}\}_{(j,k) \in E} \in \mathbb{\Gamma}$ that define the contribution of all the superpositions $f_v$ into f.
- **Hyperparameters** $h \in \mathbb{H}$ that define the prior distribution.
- **Metaparameters** $\lambda \in \mathbb{\Lambda}$ that define the optimization function.
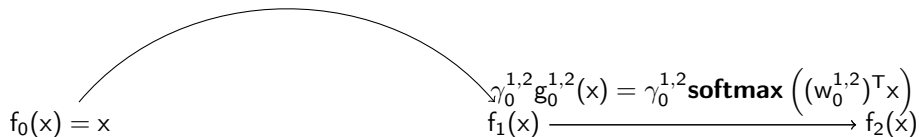
## Structure selection: one-layer network

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

$$\text{Model: } f(x) = \textbf{softmax}\left((w_0^{1,2})^\mathsf{T} f_1(x)\right), \quad f(x) : \mathbb{R}^n \to [0,1]^{|\mathbb{Y}|}, \quad x \in \mathbb{R}^n.$$

$$f_1(x) = \gamma_0^{0,1} g_0^{0,1}(x) + \gamma_1^{0,1} g_1^{0,1}(x),$$

where $w = [w_0^{0,1}, w_1^{0,1}, w_0^{1,2}]^\mathsf{T}$ — parameter matrices, $\{g_{0,1}^0, g_{0,1}^1, g_{1,2}^0\}$ — generalized-linear functions, alternatives of layers of the network.

$$\gamma_0^{0,1} g_0^{0,1}(x) = \gamma_0^{0,1} \boldsymbol{\sigma}\left((w_0^{0,1})^\mathsf{T} x\right)$$

$$f_0(x) = x$$

$$\gamma_0^{1,2} g_0^{1,2}(x) = \gamma_0^{1,2}\textbf{softmax}\left((w_0^{1,2})^\mathsf{T} x\right)$$
$$f_1(x) \longrightarrow f_2(x)$$

$$\gamma_1^{0,1} g_1^{0,1}(x) = \gamma_1^{0,1} \boldsymbol{\sigma}\left((w_1^{0,1})^\mathsf{T} x\right)$$

# 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 $\mathbf{g}^{j,k} = [g_0^{j,k}, \ldots, g_{K^{j,k}}^{j,k}]$ with length of $K^{j,k}$;
3. for each vertex $v \in V$: a differentiable aggregation function $\mathbf{agg}_v$.
4. a function $\mathbf{f} = \mathbf{f}_{|V|-1}$ :

$$f_v(w, x) = \mathbf{agg}_v \left( \{ \langle \gamma^{j,k}, g^{j,k} \rangle \circ f_j(x) | j \in \mathrm{Adj}(v_k) \} \right), v \in \{1, \ldots, |V| - 1\}, \quad f_0(x) = x \qquad (1)$$

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 $\{g^{j,k}, (j, k) \in E\}$ and aggregation functions $\{\mathbf{agg}_v, v \in V\}$.
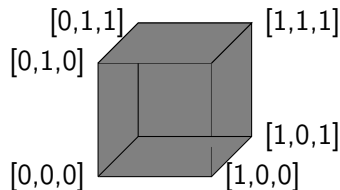
---

**Statement**

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

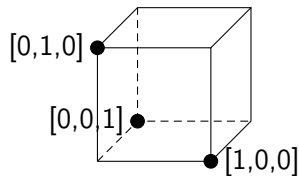## Structure restrictions

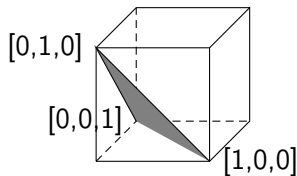An example of restrictions for structure parameter $\gamma$, $|\gamma| = 3$.



Cube vertices
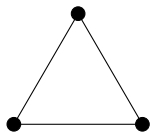


Cube interior



Simplex vertices



Simplex interior

# Prior distribution for the model structure
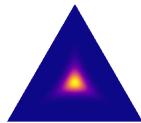
Every point in a simplex defines a model.

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



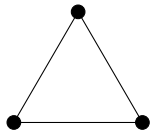$\lambda_{\text{temp}} \to 0$          $\lambda_{\text{temp}} = 0.995$          $\lambda_{\text{temp}} = 5.0$
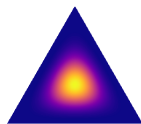
**Dirichlet distribution:** $\Gamma \sim Dir(s, \lambda_{\text{temp}})$



$\lambda_{\text{temp}} \to 0$          $\lambda_{\text{temp}} = 0.995$          $\lambda_{\text{temp}} = 5.0$

# Bayesian model selection

- **parameters**
  $w_r^{j,k} \sim \mathcal{N}(0, (\gamma_r^{j,k})^2 (A_r^{j,k})^{-1})$, $A_r^{j,k}$ is a diagonal matrix for the parameters of the primitive function $g_r^{j,k}$,
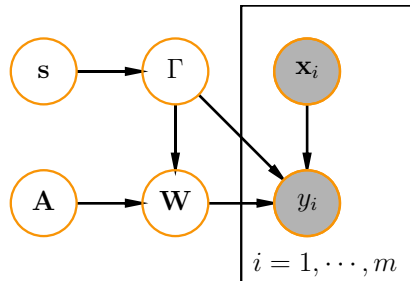
- **structure**
  $\Gamma = \{\gamma^{j,k}, (j,k) \in E\}$,
  $\gamma^{j,k} \sim GS(s^{j,k}, \lambda_{\text{temp}})$,

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

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

# Evidence as a statistical complexity

**Minimum description length** for the model f:

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

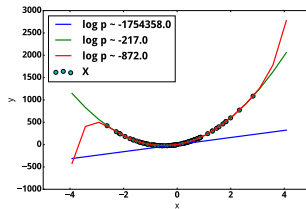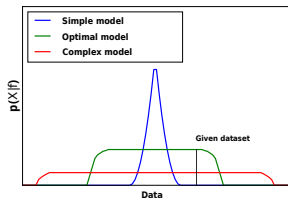where $\delta\mathfrak{D}$ is an information transmission precision.

**Bayesian approach:**

Obtain values of parameters w with respect to **posterior distribution of parameters**:

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

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

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

# Evidence lower bound
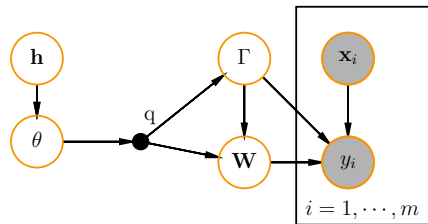
The evidence is analytically intractable.

**Model evidence:**

$$p(y|X, h, \lambda) = \iint\limits_{w, \Gamma} p(y|X, w, \Gamma)p(w, \Gamma|h, \lambda)dwd\Gamma.$$

**Definition**

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

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


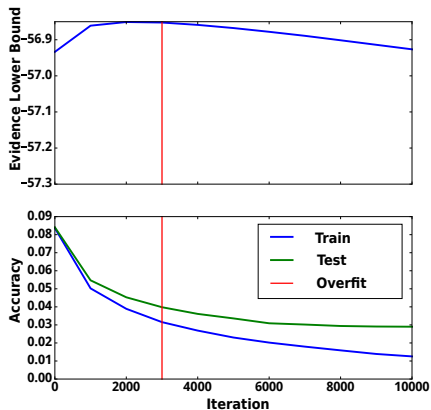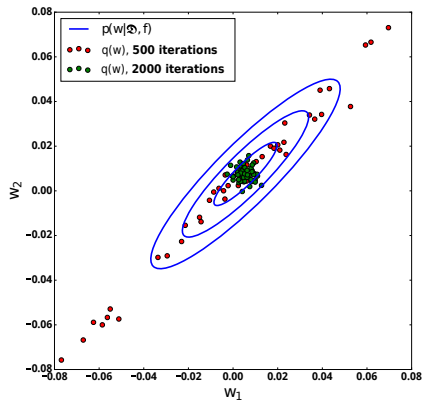
Lower bound of $\log p(y|X, h, \lambda)$:

$$\log p(y|X, h, \lambda) \geq E_q \log p(y|X, w, \Gamma) - D_{KL}(q(w, \Gamma)||p(w, \Gamma|h, \lambda)).$$

# Gradient descent as an evidence lower bound

Empirical distribtuion 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 $q = q_w q_\Gamma$ with parameters $\boldsymbol{\theta}$ that approximates posterior distribution $p(w, \Gamma | X, y, h, f)$.

**Definition**

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

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

The *model selection problem* f is a level optimization:

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

where $\boldsymbol{\theta}^*$ is a solution for the following optimization:

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

# Proposed optimization problem

**Theorem [Bakhtreev, 2019]**

The following problem is generalizing:

$$h^* = \arg\max_h Q =$$

$$= \lambda^Q_{\text{likelihood}} E_{q(w,\Gamma|\theta^*)} \log p(y|X, w, \Gamma, h, \lambda) -$$

$$- ^{\text{prior}}_Q D_{KL}\big(q(w, \Gamma|\theta^*)||p(w, \Gamma|h, \lambda)\big) -$$

$$- \sum_{p' \in \mathfrak{P}, \lambda \in \lambda^{\text{struct}}_Q} \lambda D_{KL}(p(\Gamma|h, \lambda)|p') + \log p(h|\lambda),$$

where

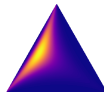$$\theta^* = \arg\max_\theta L = E_q \log p(y|X, w, \Gamma, h, \lambda)$$

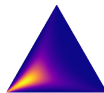$$- ^{\text{prior}}_L D_{KL}\big(q^*(w, \Gamma)||p(w, \Gamma|h, \lambda)\big).$$

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



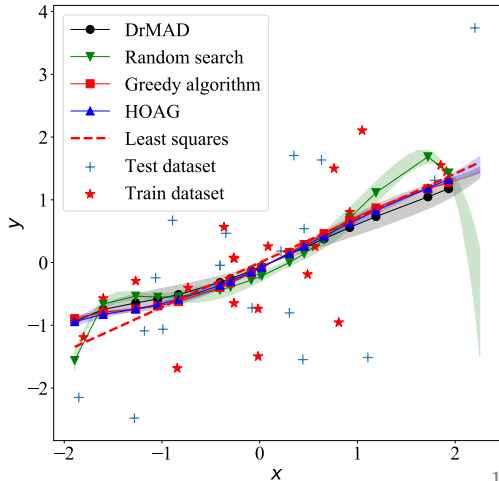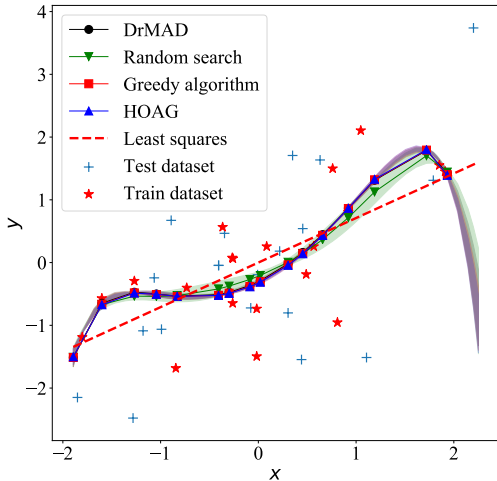$\lambda^Q_{\text{struct}} = [0; 0; 0].$


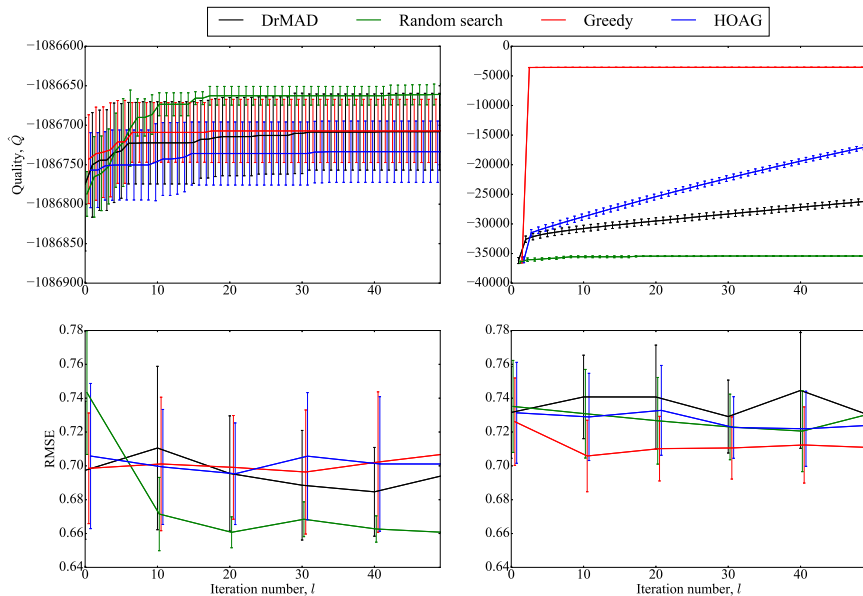
$\lambda^Q_{\text{struct}} = [1; 0; 0].$



$\lambda^Q_{\text{struct}} = [1; 1; 0].$

# Hyperparameter optimization: example

Toy example: polynomial regression with potential overfitting.
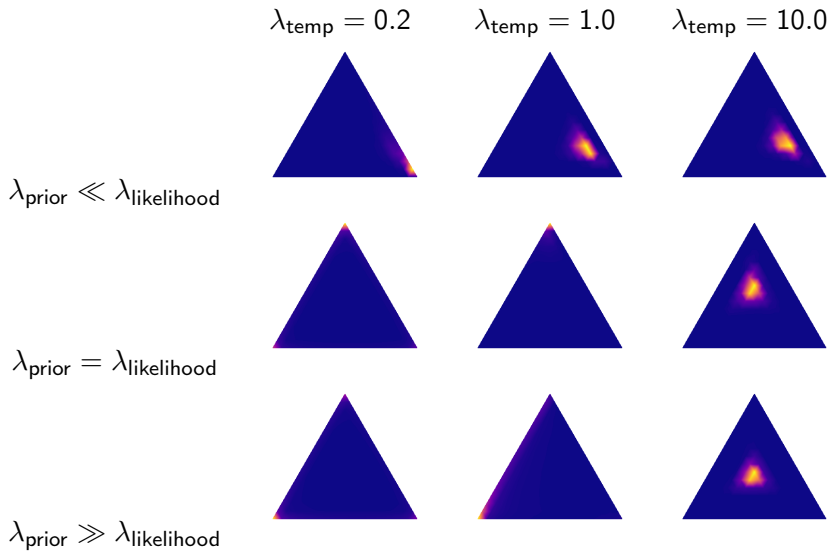
# Experiments: WISDM

## Toy example

A model f is an ensebmle of 3 models:

1. $g_0^{0,1} = tanh(wx)$;
2. $g_1^{0,1} = tanh(w^T[x, x^2, \ldots, x^{10}])$;
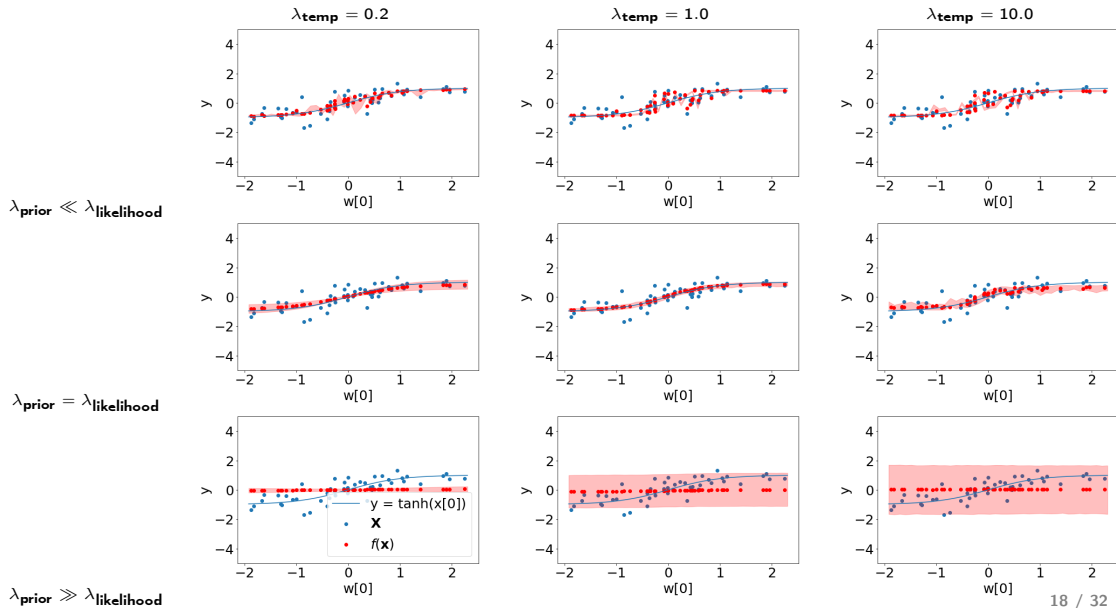3. $g_2^{0,1} = w$.

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



$\lambda_{\text{temp}} = 0.2$  $\lambda_{\text{temp}} = 1.0$  $\lambda_{\text{temp}} = 10.0$

$\lambda_{\text{prior}} \ll \lambda_{\text{likelihood}}$

$\lambda_{\text{prior}} = \lambda_{\text{likelihood}}$

$\lambda_{\text{prior}} \gg \lambda_{\text{likelihood}}$

# Toy dataset: prediction performance
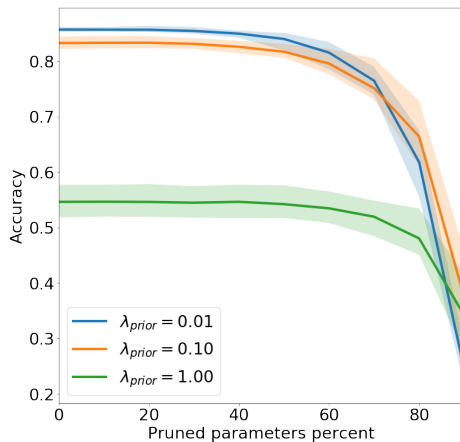
# Example

$\lambda_{\text{prior}}$ 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 archirecture?
- Can we train not a single model, but an ensemble of models? Can we control their deiversity?

# 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.
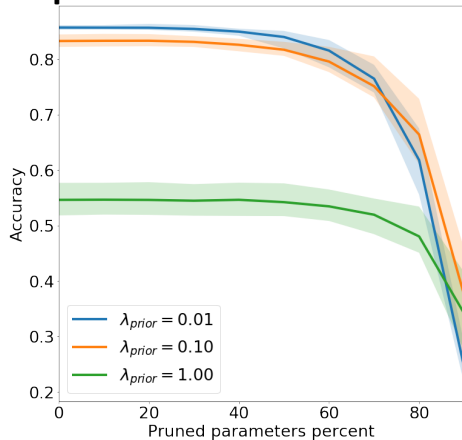
Optimize the model with hypernetworks in the following optimization procedure:

$$\mathbb{E}_{\lambda \sim P(\lambda)}(\log p(\mathfrak{D}|\mathsf{w}(\lambda)) - \lambda D_{\mathsf{KL}}(q(\mathsf{w}(\lambda))||p(\mathsf{w})) \to \max.$$
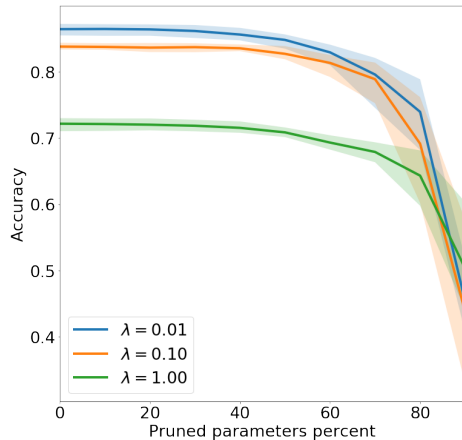
**Theorem, Grebenkova 2021**

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

# Example



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 w, but also structural parameters $\gamma$.

## Baseline: DARTS

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

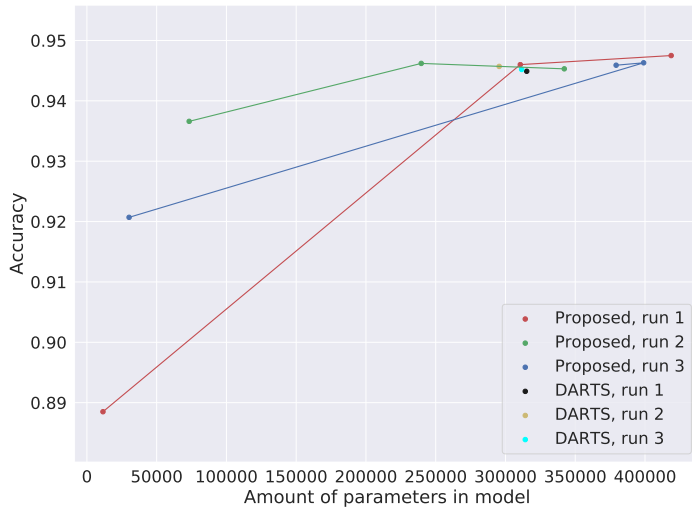$$f^{(i,j)}(x) = \langle \text{softmax}(\gamma^{(i,j)}), g^{(i,j)}(x) \rangle$$

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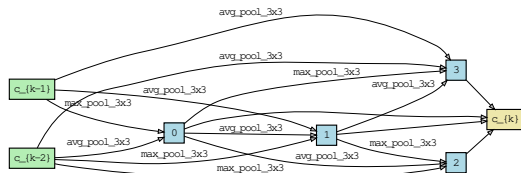
$$E_{\lambda_n} \left( \log p(y|X, w, \Gamma(\lambda_n)) + \lambda_n \sum_{(i,j)} \langle \text{softmax}\left( \frac{\gamma(\lambda_n)^{(i,j)}}{\lambda_{\text{temp}}} \right), n(g^{(i,j)}) \rangle \right),$$

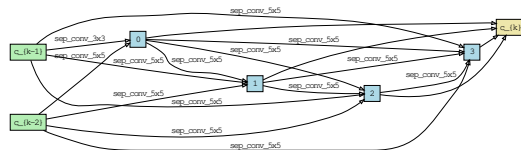where $n(g^{(i,j)})$ is a vector of amount of parameters for all the basic functions g.

# Example: Fashion-MNIST

# Example



Simple CNN cell architecture

Complex CNN cell architeuctre

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

Robustness and architecture ensembling

Intro + references

# Experiments: MNIST

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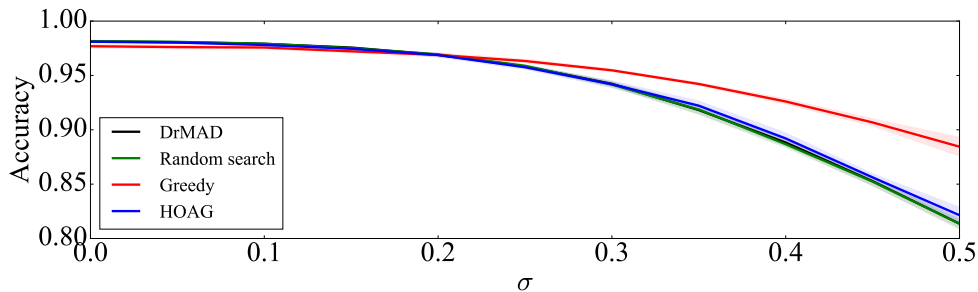


Original images      $\sigma = 0.1$      $\sigma = 0.25$      $\sigma = 0.5$

Plots on MNIST

Ensemblings

Plots on MNIST

# References