

Выбор структуры модели

Московский Физико-Технический Институт

2022

Выбор модели: связанный байесовский вывод

Первый уровень: выбираем оптимальные параметры:

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

Второй уровень: выбираем модель, доставляющую максимум обоснованности модели.

Обоснованность модели (“Evidence”):

$$p(\mathcal{D}|\mathbf{h}) = \int_{\mathbf{w}} p(\mathcal{D}|\mathbf{w})p(\mathbf{w}|\mathbf{h})d\mathbf{w}.$$

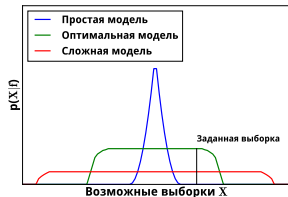
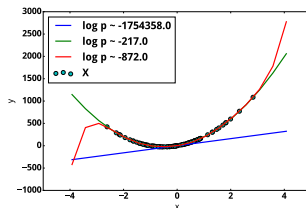


Схема выбора модели

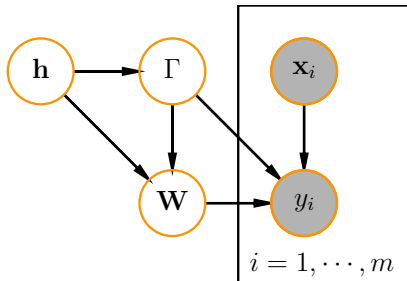


Пример: полиномы

Априорное распределение параметров

Определение

Априорным распределением параметров w и структуры Γ модели f назовем вероятностное распределение $p(W, \Gamma | h) : \mathbb{W} \times \mathbb{\Gamma} \times \mathbb{H} \rightarrow \mathbb{R}^+$, где \mathbb{W} — множество значений параметров модели, $\mathbb{\Gamma}$ — множество значений структуры модели.



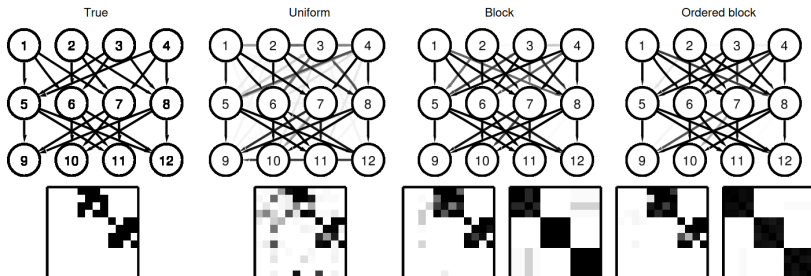
Определение

Гиперпараметрами $h \in \mathbb{H}$ модели назовем параметры распределения $p(w, \Gamma | h)$ (параметры распределения параметров модели f).

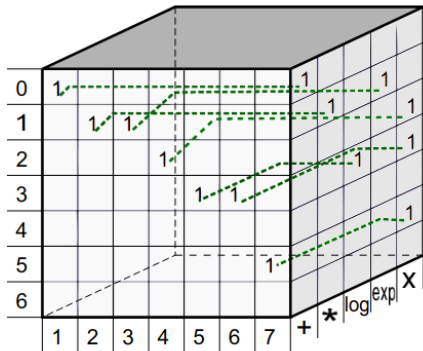
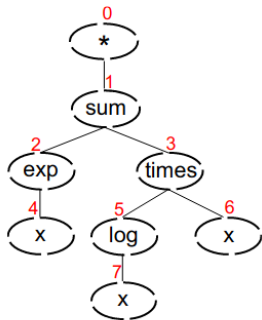
Модель f задается следующими величинами:

- **Параметры** $w \in \mathbb{W}$ задают суперпозиции f_v , из которых состоит модель f .
- **Структурные параметры** $\Gamma = \{\gamma^{j,k}\}_{(j,k) \in E} \in \mathbb{\Gamma}$ задают вклад суперпозиций f_v в модель f .
- **Гиперпараметры** $h \in \mathbb{H}$ задают распределение параметров и структурных параметров модели.
- **Метапараметры** $\lambda \in \mathbb{A}$ задают вид оптимизации модели.

Пример: байесовские сети



Пример: прогнозирование ранжирующей функции



$$f = \exp(x) + (\log x)x$$

Трехиндексная матрица связей Z_f дерева Γ_f

- вершины дерева пронумерованы;
- первые два индекса — номера вершин в ребре;
- третий индекс — выбранная элементарная функция на конце ребра.

Optimal Brain Damage

Рассматривается задача удаления неинформативных параметров.

Идея метода: Разложим функцию потерь в ряд Тейлора в окрестности максимума θ^* :

$$L(\theta^* + \Delta\theta) - L(\theta^*) = -\frac{1}{2}\theta^T H \theta + o(\|\Delta\theta\|^3),$$

где H — гессиан функции $-L$.

Для простоты вычисления будем полагать гессиан диагональным. Задача удаления параметров сводится к рассмотрению задач условной оптимизации вида:

$$L(\theta^* + \Delta\theta) \rightarrow \max$$

при

$$\theta_i^* + \Delta\theta_i = 0.$$

Показатель информативности параметра:

$$\frac{\theta_i^2}{2[H^{-1}]_{i,i}}.$$

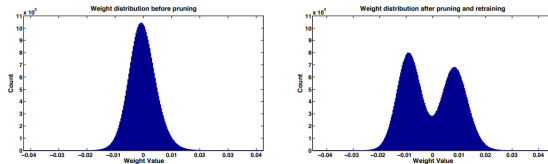
Learning both Weights and Connections for Efficient Neural Networks

Идея подхода:

- 1 Оптимизируем модель;
- 2 Удаляем наименьшие по модулю параметры;
- 3 Запускаем оптимизацию заново.

Почти очевидные факты, которые подтверждаются в статье:

- L_2 лучше для прунинга, чем L_1 в случае, если после прунинга идет оптимизация.
- Оптимизацию лучше производить из предыдущего оптимума, чем из случайной точки.
- После прунинга распределение параметров становится мультимодальным.

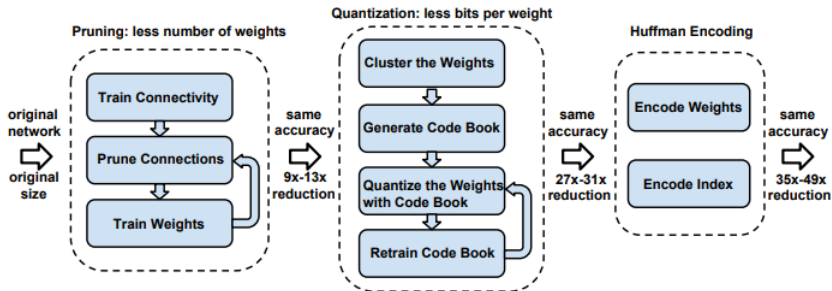


Deep Compression

Идея подхода:

- 1 Удаляем ненужные параметры модели, аналогично предыдущем подходу.
- 2 Кластеризуем параметры (K-means на каждом слое).
- 3 Производим повторную оптимизацию на центроидах.
- 4 Кодлируем индексы параметров с использованием кодов Хаффмана.

Результат: уменьшение размеров модели в 40 раз, ускорение в 3 раза.



Graves, 2011

$$\text{MDL}(f, \mathcal{D}) = L(f) + L(\mathcal{D}|f),$$

где f — модель, \mathcal{D} — выборка, L — длина описания в битах.

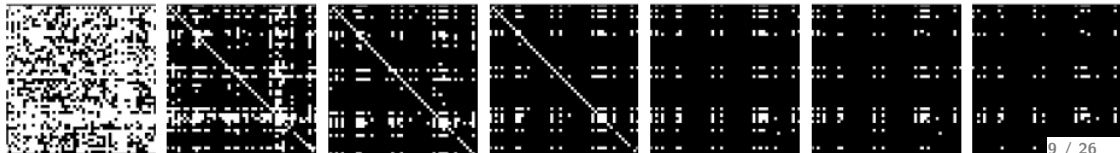
$$\text{MDL}(f, \mathcal{D}) \sim L(f) + L(w^*|f) + L(\mathcal{D}|w^*, f),$$

w^* — оптимальные параметры модели.

$$L = \sum_{x,y} \log p(y|x, \hat{w}) + \frac{1}{2} (\text{tr}(A_q) + \mu_q^T A^{-1} \mu_q - \ln |A_q|).$$

Прунинг параметра w_i определяется относительной плотностью:

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



Порождение моделей: пример

Adams et al., 2010:

- Порождаются глубокие сети доверия (Deep belief networks)
- структура модели Γ — последовательность матриц инцидентности для каждого слоя
- Порождение через Монте-Карло с использованием процесса индийского буффета в качестве априорного с параметрами α, β
- Интерпретация параметров: ширина и разреженность структуры



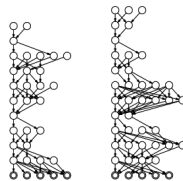
(a) $\alpha = 1, \beta = 1$



(b) $\alpha = \frac{1}{2}, \beta = 1$



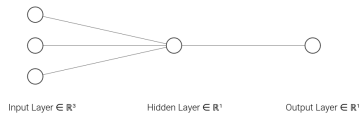
(c) $\alpha = 1, \beta = 2$



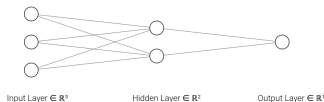
(d) $\alpha = \frac{3}{2}, \beta = 1$

Structure selection example

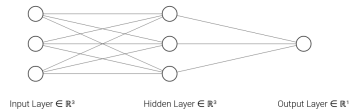
hidden layer dim = 1



hidden layer dim = 2



hidden layer dim = 3



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

Structure selection: one-layer network

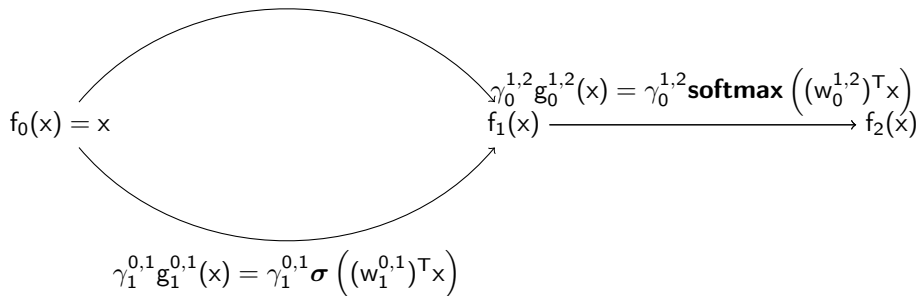
The model f is defined by the **structure** $\Gamma = [\gamma^{0,1}, \gamma^{1,2}]$.

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

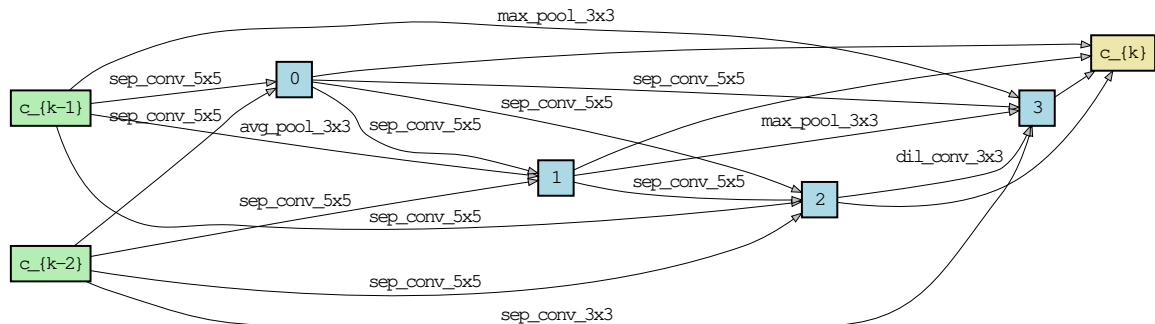
$$f_1(x) = \gamma_0^{0,1} g_0^{0,1}(x) + \gamma_1^{0,1} g_1^{0,1}(x),$$

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

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



Neural architecture search example



Structure selection: neural architecture search space

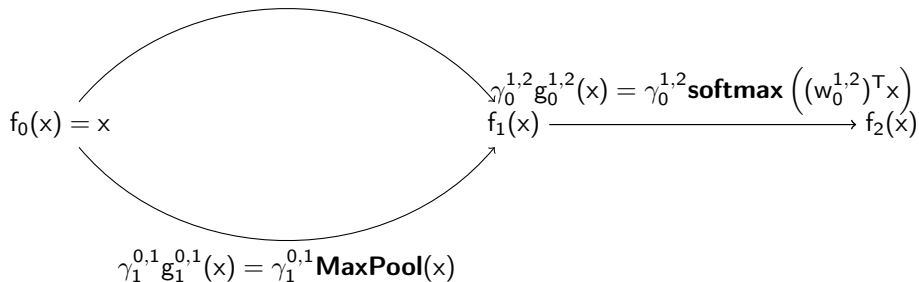
The model f is defined by the **structure** $\Gamma = [\gamma^{0,1}, \gamma^{1,2}]$.

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

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

where $w = [w_0^{0,1}, w_0^{1,2}]^T$ — parameter matrices, $g_{0,1}^0$ is a convolution, $g_{0,1}^1$ is a pooling operation, $g_{1,2}^0$ is a generalized-linear function.

$$\gamma_0^{0,1} g_0^{0,1}(x) = \gamma_0^{0,1} \mathbf{Conv}(x, w_0^{0,1})$$



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

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

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

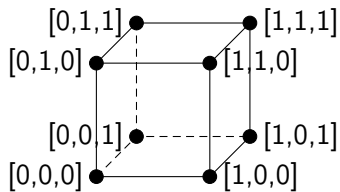
A *parametric set of models* \mathfrak{F} is a graph (V, E) with a set of primitive functions $\{g^{j,k}, (j, k) \in E\}$ and aggregation functions $\{\mathbf{agg}_v, v \in V\}$.

Statement

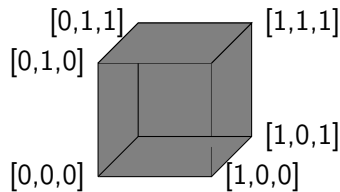
A function $f \in \mathfrak{F}$ is a model for each $\gamma^{j,k} \in [0, 1]^{K^{j,k}}$.

Structure restrictions

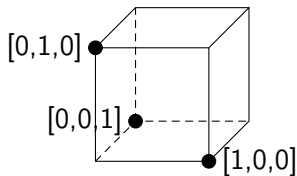
An example of restrictions for structure parameter γ , $|\gamma| = 3$.



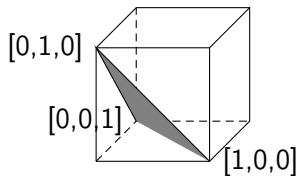
Cube vertices



Cube interior



Simplex vertices

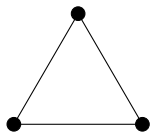


Simplex interior

Prior distribution for the model structure

Every point in a simplex defines a model.

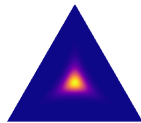
Gumbel-Softmax distribution: $\boldsymbol{\Gamma} \sim \text{GS}(\mathbf{s}, \lambda_{\text{temp}})$



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

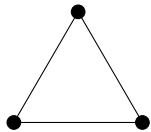


$$\lambda_{\text{temp}} = 0.995$$



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

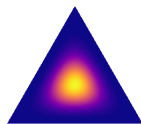
Dirichlet distribution: $\boldsymbol{\Gamma} \sim \text{Dir}(\mathbf{s}, \lambda_{\text{temp}})$



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



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



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

Neural Architecture Search: постановка задачи

\mathbf{w} — параметры модели, оптимизируемые при заданной структуре.

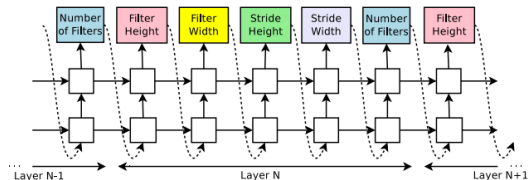
Γ — структура модели, задается контроллером, должна доставлять максимум валидации.

$$\Gamma^* = \arg \max Q(\mathbf{w}^*, \Gamma),$$

$$\mathbf{w}^* = \arg \max L(\mathbf{w}, \Gamma).$$

Neural Architecture Search with Reinforcement Learning

Структура выбирается контроллером. В цикле выбора структуры производится полная оптимизация параметров модели.

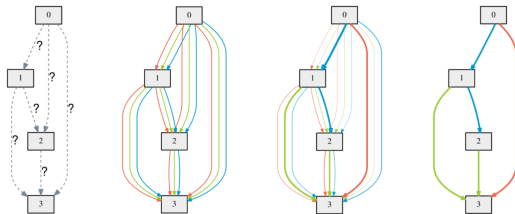


DARTS

Модель — мультиграф, где ребра $[g^e]$ соответствуют подмоделям, а вершины $f_v(x)$ — результату действия подмоделей на выборку.

Результат применения подмоделей:

$$f_v = \langle \gamma, \text{softmax}([g^e(x)]) \rangle.$$



DARTS

Задача оптимизации:

$$\Gamma^* = \arg \max Q(\mathbf{w}^*, \Gamma),$$

$$\mathbf{w}^* = \arg \max L(\mathbf{w}, \Gamma).$$

Оптимизация структуры производится жадным градиентным методом:

$$\nabla_{\Gamma} Q(\mathbf{w}', \Gamma) = \lambda_L \nabla_{\Gamma, \mathbf{w}} L(\mathbf{w}, \Gamma) \nabla_{\mathbf{w}} Q(\Gamma, \mathbf{w}').$$

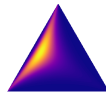
Перебор структур

Процесс перебора можно осуществить путем добавления регуляризации на структуру:

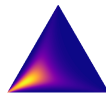
$$\lambda_1 \text{KL}(\mathbf{\Gamma}|\mathbf{\Gamma}_1) + \lambda_2 \text{KL}(\mathbf{\Gamma}|\mathbf{\Gamma}_2) + \dots$$



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



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



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

Эксплуатационные критерии качества для выбора структуры

- Количество параметров в элементах структуры
- Количество вершин в структуре
- Количество ребер в структуре
- Сложность вычисления подфункции

FBNet

$$\min_{\Gamma} \min_w L \cdot \lambda_1 \log^{\lambda_2} \text{LAT}(\Gamma),$$

где LAT — функция аппаратной задержки операций в структуре, измеренная для целевого железа.

FBNet

Model	#Parameters	#FLOPs	Latency on iPhone X	Latency on Samsung S8	Top-1 acc (%)
FBNet-iPhoneX	4.47M	322M	19.84 ms (target)	23.33 ms	73.20
FBNet-S8	4.43M	293M	27.53 ms	22.12 ms (target)	73.27

Table 5. FBNets searched for different devices.

Литература

- Bishop C. M., Nasrabadi N. M. Pattern recognition and machine learning. – New York : springer, 2006. – Т. 4. – №. 4. – С. 738.
- Бахтеев О. Ю. Байесовский выбор субоптимальной структуры модели глубокого обучения, диссертация
- Mansinghka V. et al. Structured priors for structure learning //arXiv preprint arXiv:1206.6852. – 2012.
- Варфоломеева А. А. Методы структурного обучения в задаче обобщения структур прогностических моделей, магистерская диссертация
- LeCun Y., Denker J., Solla S. Optimal brain damage //Advances in neural information processing systems. – 1989. – Т. 2.
- Han S. et al. Learning both weights and connections for efficient neural network //Advances in neural information processing systems. – 2015. – Т. 28.
- Graves A. Practical variational inference for neural networks //Advances in neural information processing systems. – 2011. – Т. 24.
- Han S., Mao H., Dally W. J. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding //arXiv preprint arXiv:1510.00149. – 2015.
- Adams R. P., Wallach H., Ghahramani Z. Learning the structure of deep sparse graphical models //Proceedings of the thirteenth international conference on artificial intelligence and statistics. – JMLR Workshop and Conference Proceedings, 2010. – С. 1-8.
- Jang E., Gu S., Poole B. Categorical reparameterization with gumbel-softmax //arXiv preprint arXiv:1611.01144. – 2016.
- Zoph B., Le Q. V. Neural architecture search with reinforcement learning //arXiv preprint arXiv:1611.01578. – 2016.
- Liu H., Simonyan K., Yang Y. Darts: Differentiable architecture search //arXiv preprint arXiv:1806.09055. – 2018.
- Wu B. et al. Fbnet: Hardware-aware efficient convnet design via differentiable neural architecture search //Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. – 2019. – С. 10734-10742.