

Выбор структуры модели глубокого обучения

Бахтеев Олег

МФТИ

20.11.2019

Резюме прошлых семинаров

Заданы:

- Вариационное распределение $q_{\mathbf{w}}(\mathbf{w}|\Gamma, \theta_{\mathbf{w}})$ с параметрами θ ;
- Априорное распределение $p(\mathbf{w}|\Gamma, \mathbf{h}, \lambda)$ с параметрами \mathbf{h} ;
- Функция потерь L и функция валидации Q .

Требуется: предложить метод выбора структуры модели Γ .

Вопросы:

- Как задать структуру модели?
- Как провести ее выбор?
- Какова вероятностная интерпретация структуры?

Automatic relevance determination

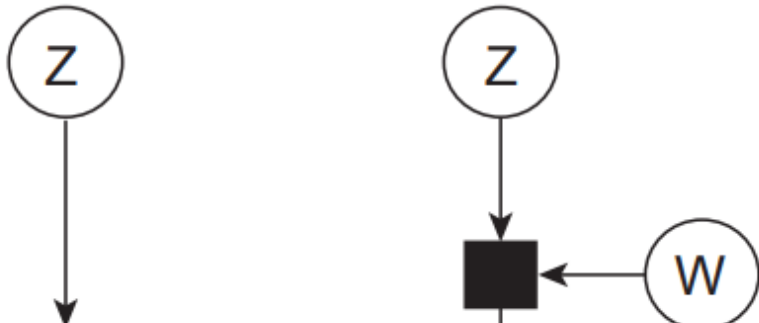
Пример: вариационный автокодировщик + ARD

VAE:

$$L = \int_{\mathbf{z}} p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}.$$

VAE + ARD:

$$L = \iint_{\mathbf{z}, \gamma} p(\mathbf{x}|\mathbf{z} \cdot \gamma)p(\mathbf{z})p(\gamma)d\mathbf{z}d\gamma.$$



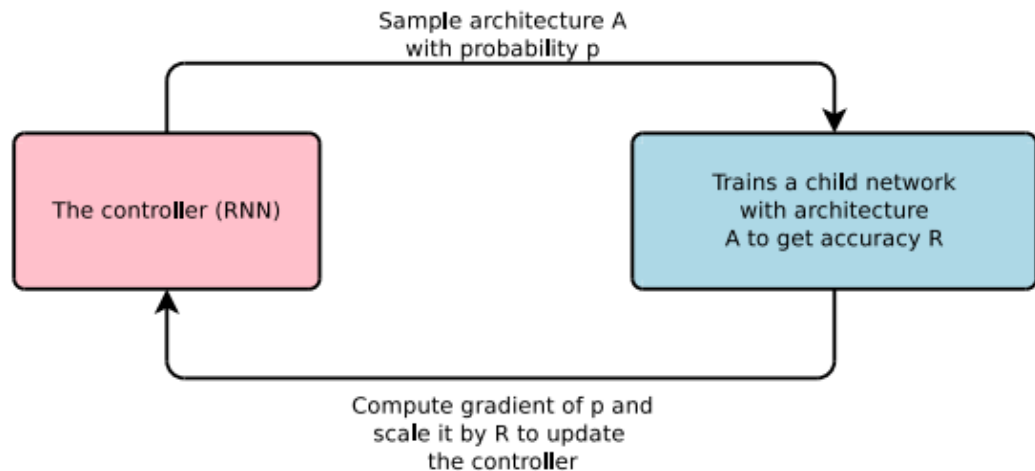
SpikeAndSlab

Гауссовый процесс для выбора структуры модели

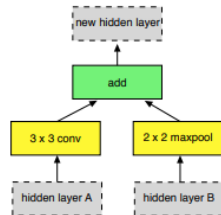
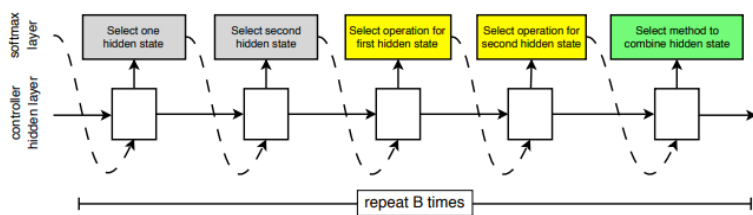
Индийский буфет, еще один пример

AdaNet

Neural Architecture Search



Neural Architecture Search



Neural Architecture Search: результаты

Model	image size	# parameters	Mult-Adds	Top 1 Acc. (%)	Top 5 Acc. (%)
Inception V2 [29]	224×224	11.2 M	1.94 B	74.8	92.2
NASNet-A (5 @ 1538)	299×299	10.9 M	2.35 B	78.6	94.2
Inception V3 [59]	299×299	23.8 M	5.72 B	78.0	93.9
Xception [9]	299×299	22.8 M	8.38 B	79.0	94.5
Inception ResNet V2 [57]	299×299	55.8 M	13.2 B	80.4	95.3
NASNet-A (7 @ 1920)	299×299	22.6 M	4.93 B	80.8	95.3
ResNeXt-101 (64 x 4d) [67]	320×320	83.6 M	31.5 B	80.9	95.6
PolyNet [68]	331×331	92 M	34.7 B	81.3	95.8
DPN-131 [8]	320×320	79.5 M	32.0 B	81.5	95.8
SENet [25]	320×320	145.8 M	42.3 B	82.7	96.2
NASNet-A (6 @ 4032)	331×331	88.9 M	23.8 B	82.7	96.2

Zoph et al., 2017. Сложность моделей отличается почти в два раза при одинаковом качестве.

Neural Architecture Search: постановка задачи

TODO

DARTS

DARTS

Графовое представление модели глубокого обучения

Заданы:

- 1 ациклический граф (V, E) ;
- 2 для каждого ребра $(j, k) \in E$: вектор базовых дифференцируемых функций $\mathbf{g}^{j,k} = [\mathbf{g}_0^{j,k}, \dots, \mathbf{g}_{K^{j,k}}^{j,k}]$ мощности $K^{j,k}$;
- 3 для каждой вершины $v \in V$: дифференцируемая функция агрегации \mathbf{agg}_v .
- 4 Функция $\mathbf{f} = \mathbf{f}_{|V|-1}$, задаваемая по правилу

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

и являющаяся функцией из признакового пространства \mathbb{X} в пространство меток \mathbb{Y} при значениях векторов, $\gamma^{j,k} \in [0, 1]^{K^{j,k}}$.

Определение

Граф (V, E) со множеством векторов базовых функций $\{\mathbf{g}^{j,k}, (j, k) \in E\}$ и функций агрегаций $\{\mathbf{agg}_v, v \in V\}$ назовем *параметрическим семейством моделей* \mathfrak{F} .

Утверждение

Для любого значения $\gamma^{j,k} \in [0, 1]^{K^{j,k}}$ функция $\mathbf{f} \in \mathfrak{F}$ является моделью.

Выбор структуры: двуслойная нейросеть

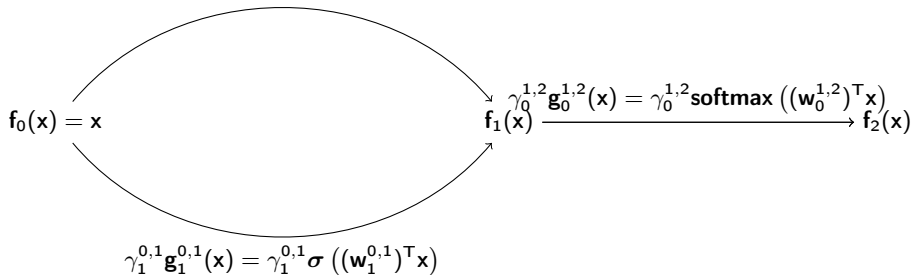
Модель f задана структурой $\Gamma = [\gamma^{0,1}, \gamma^{1,2}]$.

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

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

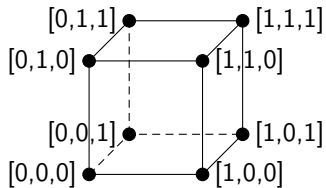
где $w = [w_0^{0,1}, w_1^{0,1}, w_0^{1,2}]^T$ — матрицы параметров, $\{g_0^0, g_0^1, g_1^0\}$ — обобщенно-линейные функции скрытых слоев нейросети.

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

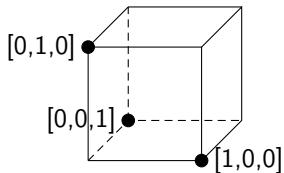


Ограничения на структурные параметры

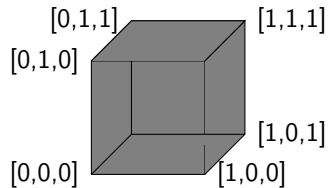
Примеры ограничений для одного структурного параметра γ , $|\gamma| = 3$.



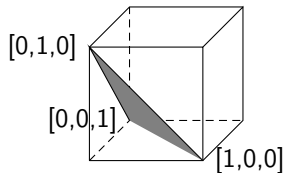
На вершинах куба



На вершинах симплекса



Внутри куба



Внутри симплекса

Репараметризация

Определение

Случайную величину ψ с распределением q с параметрами θ_ψ назовем репараметризованной через случайную величину ε , чье распределение не зависит от параметров θ_ψ , если:

$$\psi = g(\varepsilon, \theta_\psi)$$

где g — некоторая непрерывная функция.

Пример

$$E_{q_w(w|\Gamma, \theta_w)} \log p(y|X, w, \Gamma) = \int_w \log p(y|X, w, \Gamma) q_w(w|\Gamma, \theta_w) dw.$$

Продифференцируем по параметрам θ_w :

$$\nabla_{\theta_w} E_{q_w(w|\Gamma, \theta_w)} \log p(y|X, w, \Gamma) = \int_w \log p(y|X, w, \Gamma) \nabla_{\theta_w} q_w(w|\Gamma, \theta_w) dw.$$

Пусть возможна репараметризация: $w = g(\varepsilon, \theta_w)$. Тогда:

$$\begin{aligned} \nabla_{\theta_w} E_{q(w, \Gamma|\theta)} \log p(y|X, w, \Gamma) &= \nabla_{\theta_w} E_\varepsilon \log p(y|X, g(\varepsilon), \Gamma) = \\ &= \int_\varepsilon \nabla_{\theta_w} \log p(y|X, g(\varepsilon), \Gamma) p(\varepsilon) d\varepsilon = E_\varepsilon \nabla_{\theta_w} \log p(y|X, g(\varepsilon), \Gamma). \end{aligned}$$

Reparametrization

Logit-Normal

Gumbel-Softmax

Proposed Method

Proposed Method

Properties

Examples

Используемые материалы