

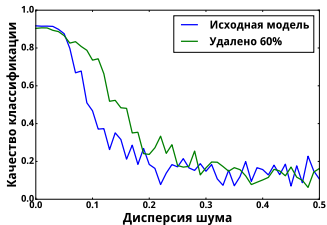
# Выбор структуры модели глубокого обучения

Бахтеев Олег

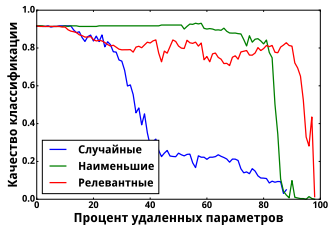
МФТИ

06.02.2019

# Сложность модели: зачем?



Устойчивость моделей при возмущении выборки



Качество классификации при удалении параметров

# Сложность модели: зачем?

Model	image size	# parameters	Multi-Adds	Top 1 Acc. (%)	Top 5 Acc. (%)
Inception V2 [29]	224×224	11.2 M	1.94 B	74.8	92.2
<b>NASNet-A (5 @ 1538)</b>	<b>299×299</b>	<b>10.9 M</b>	<b>2.35 B</b>	<b>78.6</b>	<b>94.2</b>
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
<b>NASNet-A (7 @ 1920)</b>	<b>299×299</b>	<b>22.6 M</b>	<b>4.93 B</b>	<b>80.8</b>	<b>95.3</b>
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
<b>SENet [25]</b>	<b>320×320</b>	<b>145.8 M</b>	<b>42.3 B</b>	<b>82.7</b>	<b>96.2</b>
<b>NASNet-A (6 @ 4032)</b>	<b>331×331</b>	<b>88.9 M</b>	<b>23.8 B</b>	<b>82.7</b>	<b>96.2</b>

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

# Глубокого обучение

## Определение

Моделью  $f(\mathbf{w}, \mathbf{x})$  назовем дифференцируемую по параметрам  $\mathbf{w}$  функцию из множества признаков описаний объекта во множество меток:

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

где  $\mathbb{W}$  — пространство параметров функции  $f$ .

**Особенность задачи** выбора модели *глубокого обучения* — значительное число параметров в моделях приводит к неприменимости классических методов оптимизации и выбора модели.

## Сложность модели:

- 1 количество параметров;
- 2 количество суперпозиций внутри модели.

# Принцип минимальной длины описания

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

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

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

$\mathbf{w}^*$  — оптимальные параметры модели.

$f_1$	$L(f_1)$	$L(w_1^* f_1)$	$L(\mathcal{D} w_1^*, f_1)$
$f_2$	$L(f_2)$	$L(w_2^* f_2)$	$L(\mathcal{D} w_2^*, f_2)$
$f_3$	$L(f_3)$	$L(w_3^* f_3)$	$L(\mathcal{D} w_3^*, f_3)$

# MDL и Колмогоровская сложность

**Колмогоровская сложность** — длина минимального кода для выборки на предварительно заданном языке.

## **Теорема инвариантности**

Для двух сводимых по Тьюрингу языков колмогоровская сложность отличается не более чем на константу, не зависящую от мощности выборки.

## **Отличия от MDL:**

- Колмогоровская сложность невычислима.
- Длина кода может зависеть от выбранного языка. Для небольших выборок теорема инвариантности не дает адекватных результатов.

# Байесовый подход к сложности

Правдоподобие модели (“Evidence”):

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

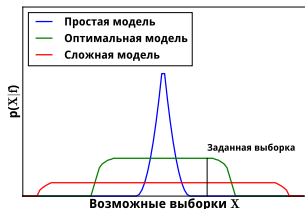
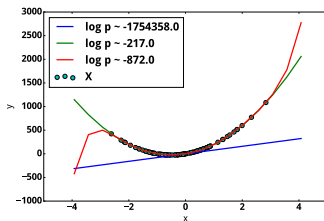


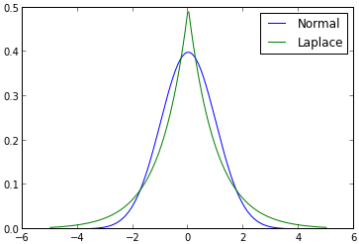
Схема выбора модели по правдоподобию



Пример: полиномы

# Evidence vs MDL

Evidence	MDL
Использует априорные знания	Независима от априорных знаний
Основывается на гипотезе о порождении выборки вне зависимости от их природы	Минимизирует длину описания выборки





# Оптимальность модели

## Определение

Пусть задано множество моделей  $M$ .

Пусть для каждой модели  $\mathbf{f}$  задано априорное распределение параметров:  $p(\mathbf{w}|\mathbf{h})$ , где  $\mathbf{h}$  — параметры априорного распределения.

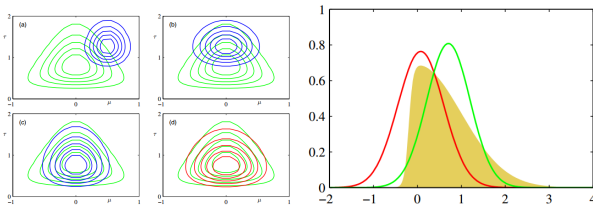
Модель  $\mathbf{f}$  назовем оптимальной среди моделей  $M$ , если достигается максимум интеграла:

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

# Вариационная оценка, ELBO

Вариационная оценка Evidence, Evidence lower bound — метод нахождения приближенного значения аналитически невычислимого распределения  $p(\mathbf{w}|\mathcal{D}, \mathbf{h})$  распределением  $q(\mathbf{w}) \in \mathbf{Q}$ . Получение вариационной нижней оценки обычно сводится к задаче минимизации

$$\text{KL}(q(\mathbf{w})||p(\mathbf{w}|\mathcal{D})) = - \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathbf{w}|\mathcal{D})}{q(\mathbf{w})} d\mathbf{w}.$$



Аппроксимация неизвестного распределения нормальным

Аппроксимация Лапласа (красная линия) и вариационная оценка (зеленая линия)

# Получение вариационной нижней оценки

$$\begin{aligned}\log p(\mathcal{D}|\mathbf{h}) &= \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathcal{D}, \mathbf{w}|\mathbf{h})}{q(\mathbf{w})} d\mathbf{w} + D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathcal{D}, \mathbf{h})) \geq \\ &\geq \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathcal{D}, \mathbf{w}|\mathbf{h})}{q(\mathbf{w})} d\mathbf{w} = \\ &= -D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{h})) + \int_{\mathbf{w}} q(\mathbf{w}) \log p(\mathcal{D}|\mathbf{w}, \mathbf{h}) d\mathbf{w},\end{aligned}$$

где

$$D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{h})) = - \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathbf{w}|\mathbf{h})}{q(\mathbf{w})} d\mathbf{w}.$$

## Определение

Модель  $\mathbf{f}$  назовем субоптимальной на множестве моделей  $M$ , если модель доставляет максимум нижней вариационной оценке:

$$\mathbf{f} = \arg \max_{\hat{\mathbf{f}} \in M} \max_{q \in Q} \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathbf{y}, \mathbf{w}|\mathcal{D}, \hat{\mathbf{f}})}{q(\mathbf{w})} d\mathbf{w}.$$

Максимизация вариационной нижней оценки

$$\int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathcal{D}, \mathbf{w} | \mathbf{h})}{q(\mathbf{w})} d\mathbf{w}$$

эквивалентна минимизации дивергенции между распределением  $q(\mathbf{w}) \in Q$  и апостериорным распределением параметров  $p(\mathbf{w} | \mathcal{D}, \mathbf{h})$ :

$$q = \operatorname{argmax}_{q \in Q} \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathcal{D}, \mathbf{w} | \mathbf{h})}{q(\mathbf{w})} d\mathbf{w} \Leftrightarrow q = \operatorname{argmin}_{q \in Q} D_{\text{KL}}(q(\mathbf{w}) || p(\mathbf{w} | \mathcal{D}, \mathbf{h})),$$

т.к.

$$\log p(\mathcal{D} | \mathbf{h}) = \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathcal{D}, \mathbf{w} | \mathbf{h})}{q(\mathbf{w})} d\mathbf{w} + D_{\text{KL}}(q(\mathbf{w}) || p(\mathbf{w} | \mathcal{D}, \mathbf{h})) = \text{const.}$$

# Использование вариационной нижней оценки

**Для чего используют вариационный вывод?**

- получение оценок Evidence;
- получение оценок распределений моделей со скрытыми переменными (тематическое моделирование, снижение размерности).

**Зачем используют вариационный вывод?**

- сводит задачу нахождения апостериорной вероятности к методам оптимизации;
- проще масштабируется, чем аппроксимация Лапласа;
- проще в использовании, чем сэмплирующие методы.

**Вариационный вывод может давать сильно заниженную оценку.**

# ELBO: нормальное распределение

Пусть  $q \sim \mathcal{N}(\boldsymbol{\mu}_q, \mathbf{A}_q)$ .

Тогда вариационная оценка имеет вид:

$$\int_{\mathbf{w}} q(\mathbf{w}) \log p(\mathbf{Y}|\mathbf{X}, \mathbf{w}, \mathbf{h}) d\mathbf{w} - D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{h})) \simeq$$
$$\sum_{i=1}^m \log p(\mathbf{y}_i|\mathbf{x}_i, \hat{\mathbf{w}}) - D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{h})) \rightarrow \max_{\mathbf{A}_q, \boldsymbol{\mu}_q}, \quad \hat{\mathbf{w}} \sim q.$$

В случае, если априорное распределение параметров  $p(\mathbf{w}|\mathbf{h})$  является нормальным:

$$p(\mathbf{w}|\mathbf{h}) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{A}),$$

дивергенция  $D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{h}))$  вычисляется аналитически:

$$D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{h})) = \frac{1}{2} (\text{tr}(\mathbf{A}^{-1}\mathbf{A}_q) + (\boldsymbol{\mu} - \boldsymbol{\mu}_q)^T \mathbf{A}^{-1}(\boldsymbol{\mu} - \boldsymbol{\mu}_q) - n + \ln |\mathbf{A}| - \ln |\mathbf{A}_q|).$$

# ELBO: нормальное распределение

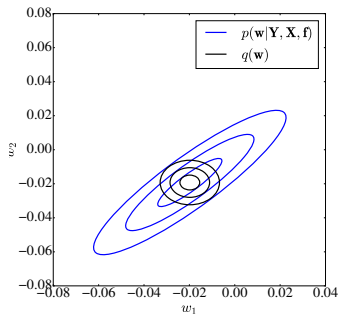
“Обычная” функция потерь:

$$L = \sum_{x,y \in \mathcal{D}} -\log p(y|x, \mathbf{w}) + \lambda \|\mathbf{w}\|_2^2.$$

Вариационный вывод при  
( $p(\mathbf{w}|\mathbf{h}) \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$ ):

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

Пример грубой аппроксимации нормальным диагональным распределением  $q$



# Оператор оптимизации, Maclaurin et. al, 2015

## Определение

Назовем оператором оптимизации алгоритм  $T$  выбора вектора параметров  $\mathbf{w}'$  по параметрам предыдущего шага  $\mathbf{w}$ :

$$\mathbf{w}' = T(\mathbf{w}).$$

## Определение

Пусть  $L$  — дифференцируемая функция потерь.

Оператором градиентного спуска назовем следующий оператор:

$$T(\mathbf{w}) = \mathbf{w} - \beta \nabla L(\mathbf{w}, \mathbf{y}, \mathcal{D}).$$



# Градиентный спуск для оценки правдоподобия

Рассмотрим максимизацию совместного распределения параметров:

$$L = -\log p(\mathcal{D}, \mathbf{w}|\mathbf{h}) = - \sum_{\mathcal{D} \in \mathcal{D}} \log p(\mathcal{D}|\mathbf{w}, \mathbf{h})p(\mathbf{w}|\mathbf{h})$$

Проведем оптимизацию нейросети из  $r$  различных начальных приближений  $\mathbf{w}_1, \dots, \mathbf{w}_r$  с использованием градиентного спуска:

$$\mathbf{w}' = T(\mathbf{w}).$$

Векторы параметров  $\mathbf{w}_1, \dots, \mathbf{w}_r$  соответствуют некоторому скрытому распределению  $q(\mathbf{w})$ .

# Энтропия

Формулу вариационной оценки можно переписать с использованием энтропии:

$$\log p(\mathcal{D}|\mathbf{f}) \geq \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathcal{D}, \mathbf{w}|\mathbf{h})}{q(\mathbf{w})} d\mathbf{w} = \\ E_{q(\mathbf{w})}[\log p(\mathcal{D}, \mathbf{w}|\mathbf{h})] + S(q(\mathbf{w})),$$

где  $S(q(\mathbf{w}))$  — энтропия:

$$S(q(\mathbf{w})) = - \int_{\mathbf{w}} q(\mathbf{w}) \log q(\mathbf{w}) d\mathbf{w}.$$

# Градиентный спуск для оценки правдоподобия

При достаточно малой длине шага оптимизации  $\beta$  разность энтропии на различных шагах оптимизации вычисляется как:

$$S(q'(\mathbf{w})) - S(q(\mathbf{w})) \simeq \frac{1}{r} \sum_{g=1}^r (-\beta \text{Tr}[\mathbf{H}(\mathbf{w}'^g)] - \beta^2 \text{Tr}[\mathbf{H}(\mathbf{w}'^g)\mathbf{H}(\mathbf{w}'^g)]).$$

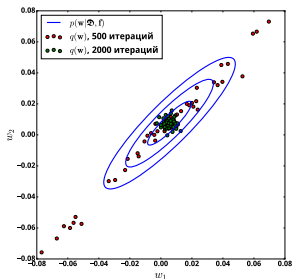
Итоговая оценка на шаге оптимизации  $\tau$ :

$$\begin{aligned} \log \hat{p}(\mathbf{Y}|\mathcal{D}, \mathbf{h}) &\sim \frac{1}{r} \sum_{g=1}^r L(\mathbf{w}_\tau^g, \mathcal{D}, \mathbf{Y}) + S(q^0(\mathbf{w})) + \\ &+ \frac{1}{r} \sum_{b=1}^{\tau} \sum_{g=1}^r (-\beta \text{Tr}[\mathbf{H}(\mathbf{w}_b^g)] - \beta^2 \text{Tr}[\mathbf{H}(\mathbf{w}_b^g)\mathbf{H}(\mathbf{w}_b^g)]), \end{aligned}$$

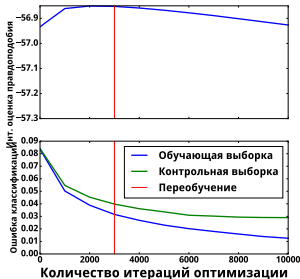
$\mathbf{w}_b^g$  — вектор параметров старта  $g$  на шаге  $b$ ,  $S(q^0(\mathbf{w}))$  — начальная энтропия.

# Переобучение, Maclaurin et. al, 2015

Градиентный спуск не минимизирует дивергенцию  $KL(q(\mathbf{w})||p(\mathbf{w}|\mathcal{D}, \mathbf{h}))$ . При приближении к моде распределения снижается оценка Evidence, что интерпретируется как переобучение модели.



Схождение распределения к моде



Оценка начала переобучения

# Задача оптимизации гиперпараметров

Задана дифференцируемая по параметрам модель, приближающая зависимую переменную  $y$ :

$$f(\mathbf{w}, \mathbf{x}) : \mathbb{W} \times \mathbb{X} \rightarrow \mathbb{Y}, \quad \mathbf{w} \in \mathbb{W}$$

Пусть  $\theta \in \mathbb{R}^u$  — вариационные параметры распределения.

$L(\theta | \mathbf{h}, \mathbf{X}, \mathbf{y})$  — дифференцируемая функция потерь по которой производится оптимизация функции  $f$ .

$Q(\mathbf{h} | \theta, \mathbf{X}, \mathbf{y})$  — дифференцируемая функция определяющая итоговое качество модели  $f$ .

Требуется найти параметры  $\theta^*$  и гиперпараметры  $\mathbf{h}$  модели, доставляющие минимум следующему функционалу:

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

$$\theta^*(\mathbf{h}) = \arg \min_{\theta \in \mathbb{R}^u} L(\theta | \mathbf{h}, \mathbf{X}, \mathbf{y}),$$

# Байесовский вывод

*Первый уровень:*

$$\theta^* = \arg \max(-L) = p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \mathbf{h}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w}|\mathbf{h})}{p(\mathbf{y}|\mathbf{X}, \mathbf{h})}.$$

*Второй уровень:*

$$p(\mathbf{h}|\mathbf{X}, \mathbf{y}) \propto p(\mathbf{y}|\mathbf{X}, \mathbf{h})p(\mathbf{h}).$$

Полагая распределение параметров  $p(\mathbf{h})$  равномерным на некоторой большой окрестности, получим задачу оптимизации гиперпараметров:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{h}) = \int_{\mathbf{w} \in \mathbb{R}^u} p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w}|\mathbf{h}) = -Q \rightarrow \max_{[\alpha_1, \dots, \alpha_n] \in \mathbb{R}^n}. \quad (1)$$

## Другие примеры $L, Q$

- Кросс-валидация ( $L$  — ошибка на обучении,  $Q$  — на контроле);
- вариационная оценка ( $L = Q = \text{ELBO}$ ).

# Формальная постановка задачи: градиентная оптимизация

## Определение

Пусть задан оператор  $T$ , проводящий  $\eta$  шагов оптимизации по функции  $L$ :

$$\theta^* = T \circ T \circ \dots \circ T(\theta^0, \mathbf{h}) = T^\eta(L, \theta_0, \mathbf{h}), \quad (2)$$

где  $\beta$  — длина шага градиентного спуска,  $\theta^0$  — начальное значение параметров  $\theta$ .

Перепишем итоговую задачу оптимизации:

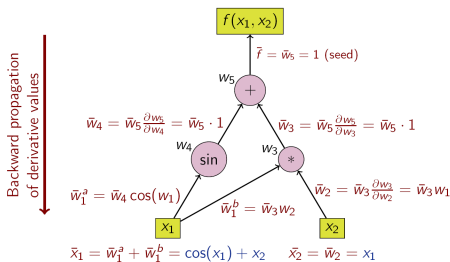
$$\mathbf{h}^* = \arg \min_{\mathbf{h} \in \mathbb{R}^b} Q(T^\eta(L, \theta_0, \mathbf{h})).$$



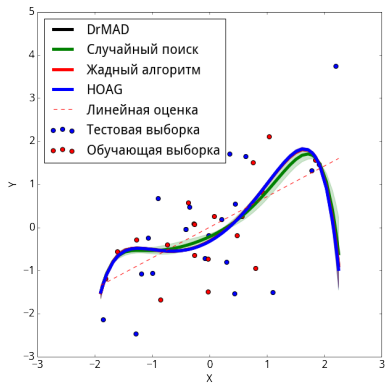
# RMAD, Maclaurin et. al, 2015

- 1 Провести  $\eta$  шагов оптимизации:  
 $\theta = T(\theta_0, \mathbf{A}^{-1})$ .
- 2 Положим  $\hat{\nabla} \mathbf{A}^{-1} = \nabla_{\mathbf{A}}^{-1} Q(\theta, \mathbf{A}^{-1})$ .
- 3 Положим  $d\mathbf{v} = \mathbf{0}$ .
- 4 Для  $\tau = \eta \dots 1$  повторить:
- 5  $\theta^{\tau-1} = \theta^{\tau} - \gamma \mathbf{v}^{\tau}$ .
- 6  $\mathbf{v}^{\tau-1} = \mathbf{v}^{\tau} + \gamma \hat{\nabla} \theta$ .
- 7  $d\mathbf{v} = \gamma \hat{\nabla} \theta$ .
- 8  $\hat{\nabla} \mathbf{A}^{-1} = \hat{\nabla} \mathbf{A}^{-1} - d\mathbf{v} \nabla_{\mathbf{A}^{-1}} \nabla_{\theta} Q$ .
- 9  $\hat{\nabla} \theta = \hat{\nabla} \theta - d\mathbf{v} \nabla_{\theta} \nabla_{\theta} Q$ .

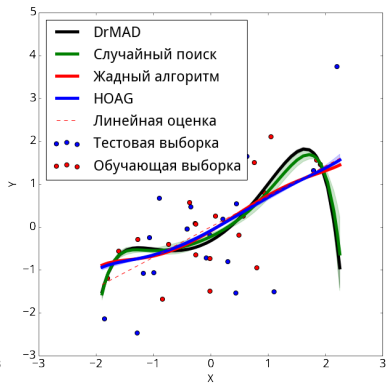
Алгоритм RMAD основывается на Reverse-mode differentiation.



# Эксперименты: полиномы



Кросс-валидация



Evidence

# Задача выбора структуры модели

Однослойная нейросеть:

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

$$\mathbf{f}_1(\mathbf{x}) = \gamma_{0,1}^1 \mathbf{g}_{0,1}^1(\mathbf{x}) + \dots + \gamma_{0,1}^K \mathbf{g}_{0,1}^K(\mathbf{x}) = \gamma_{0,1}^1 \sigma(\mathbf{W}_1^T \mathbf{x}) + \dots + \gamma_{0,1}^K \sigma(\mathbf{W}_K^T \mathbf{x}),$$

где  $\mathbf{W} = [\mathbf{W}_0, \mathbf{W}_1, \dots, \mathbf{W}_K]^T$  — матрицы параметров,  $\{\mathbf{g}_{0,1}^i\}_{i=1}^K$  — базовые функции скрытого слоя нейросети.

**Структурные параметры:**  $\Gamma = [\gamma_{0,1}]$ .

**Структура модели** задается вершиной  $K$ -мерного симплекса.

# Задача выбора структуры модели: два скрытых слоя

Двухслойная нейросеть:

$$\mathbf{f}(x) = \text{softmax}(\mathbf{W}^T \mathbf{f}_2(x)), \quad f(x) : \mathbb{R}^n \rightarrow [0, 1]^{|Y|}, \quad x \in \mathbb{R}^n.$$

$$\mathbf{f}_2(x) = \gamma_{1,2}^1 \mathbf{g}_{1,2}^1(\mathbf{f}_1(x)) + \dots + \gamma_{1,2}^K \mathbf{g}_{1,2}^K(\mathbf{f}_1(x)) = \gamma_{1,2}^1 \sigma(\mathbf{W}_{K+1}^T \mathbf{f}_1(x)) + \dots + \gamma_{1,2}^K \sigma(\mathbf{W}_{2K}^T \mathbf{f}_1(x))$$

$$\mathbf{f}_1(x) = \gamma_{0,1}^1 \mathbf{g}_{0,1}^1(x) + \dots + \gamma_{0,1}^K \mathbf{g}_{0,1}^K(x) = \gamma_{0,1}^1 \sigma(\mathbf{W}_1^T x) + \dots + \gamma_{0,1}^K \sigma(\mathbf{W}_K^T x),$$

где  $\mathbf{W} = [\mathbf{W}_0, \mathbf{W}_1, \dots, \mathbf{W}_{2K}]^T$  — матрицы параметров,  $\{\mathbf{g}_{0,1}^i, \mathbf{g}_{1,2}^i\}_{i=1}^K$  — базовые функции скрытых слоев нейросети.

Структурные параметры:  $\Gamma = [\gamma_{0,1}, \gamma_{1,2}]$ .

Структура модели задается вершинами **двух**  $K$ -мерных симплексов.

# Графовое представление модели глубокого обучения

## Определение

Задан граф  $(V, E)$ . Для каждого ребра  $(j, k) \in E$  определен вектор базовых функций мощности  $K^{j,k}$ :

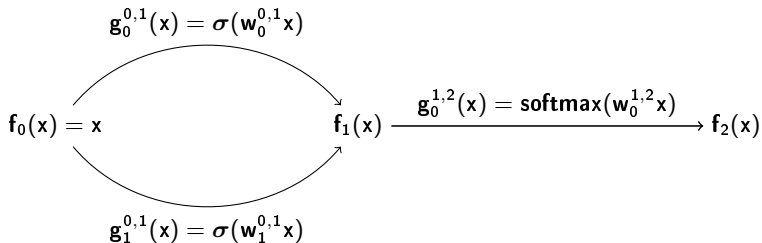
$$\mathbf{g}^{j,k} = [\mathbf{g}_0^{j,k}, \dots, \mathbf{g}_{K^{j,k}}^{j,k}]$$

. Пусть для каждой вершины  $v \in V$  определена функция агрегации  $\mathbf{agg}_v$ . Граф  $(V, E)$  в совокупности со множеством векторов базовых функций  $\{\mathbf{g}^{j,k}, (j, k) \in E\}$  и множеством функций агрегаций  $\{\mathbf{agg}_v, v \in V\}$  называется *параметрическим семейством моделей*  $\mathfrak{F}$ , если функция, задаваемая как

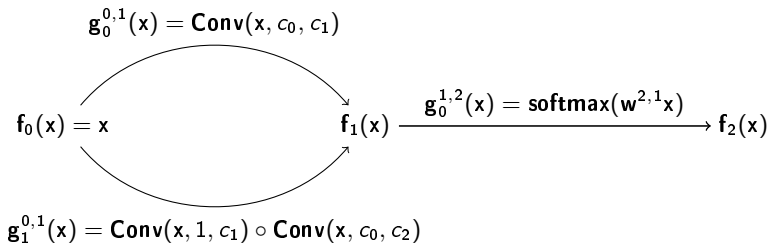
$$\mathbf{f}_k(\mathbf{x}) = \mathbf{agg}_k (\{ \langle \gamma^{j,k}, \mathbf{g}^{j,k} \rangle (\mathbf{f}_j(\mathbf{x})) \mid j \in \text{Adj}(v_k) \}), \quad \mathbf{f}_0(\mathbf{x}) = \mathbf{x} \quad (3)$$

является моделью при любых значениях векторов,  $\gamma^{j,k} \in [0, 1]^{K^{j,k}}$ .

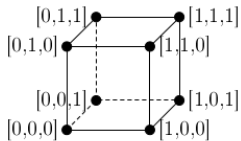
# Пример: двуслойная нейросеть



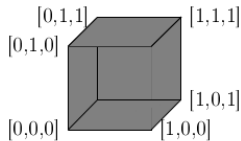
# Пример: сверточная сеть



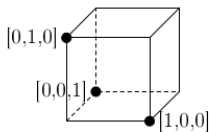
# Ограничения на структурные параметры



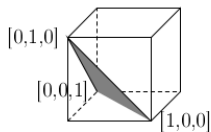
(a)



(б)



(в)



(г)



# Статистические критерии качества модели

**Параметрическая сложность** — наименьшая дивергенция между априорным распределением параметров и апостериорным распределением параметров:

$$C_{\text{param}} = \min_{\mathbf{h}} D_{\text{KL}}(p(\mathbf{W}, \Gamma | \mathbf{y}, \mathbf{X}) || p(\mathbf{W}, \Gamma | \mathbf{h})).$$

**Структурная сложность модели** — энтропия апостериорного распределения структуры модели:

$$C_{\text{struct}} = -E_p \log p(\Gamma | \mathbf{y}, \mathbf{X}).$$