

# Математическое разложение оценки неопределенности для нейронных сетей.

Насыров Р.Р., Зайцев А.А.

МФТИ

21 декабря 2024 г.

# Содержание

- 1 Введение
  - Актуальность
  - Калибрация
  - Неопределенность
- 2 Постановка задачи
- 3 Решение задачи
  - Обозначения
  - Вывод формулы
  - Аппроксимация риска
- 4 Гипотезы
- 5 Вычислительный эксперимент
- 6 Выводы

# Введение

## Актуальность

- В чувствительных (здравоохранение, финансы) системах нужно оценивать надежность работы системы
- Нейросеть — один из компонентов системы
- Нужно оценивать неопределенность ответа нейросети

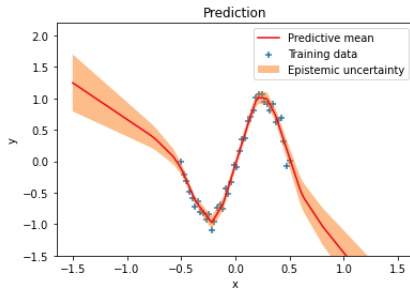


Рис.: Неопределенность.

# Введение

## Калибрация

Модель классификации  $h : \mathcal{X} \rightarrow \mathcal{Y} \times [0, 1]$  называется *скалиброванной*, если предсказываемые вероятности равны реальным вероятностям:

$$\mathbb{P}_{(X,Y) \sim \pi}(\hat{Y} = Y | \hat{P} = p) = p \quad \forall p \in [0, 1]$$

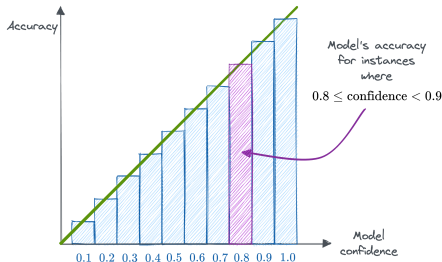


Рис.: Точность и уверенность.

# Введение

## Неопределенность

- *Алеаторная* неопределенность — связана с шумом в данных
- *Эпистемическая* неопределенность — связана с ограничениями модели и доступных данных

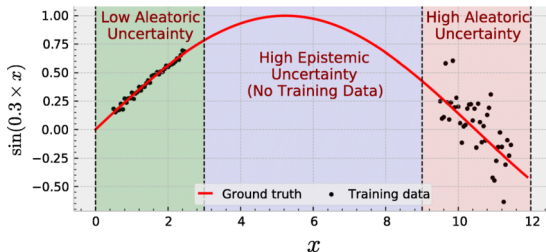


Рис.: 2 типа неопределенности.

[https://www.researchgate.net/figure/Illustration-of-epistemic-and-aleatoric-uncertainty\\_fig358723173](https://www.researchgate.net/figure/Illustration-of-epistemic-and-aleatoric-uncertainty_fig358723173)

**Цель исследования:** предложить метод обучения нейросетей, обеспечивающий высокую калибровку полученной модели без снижения качества классификации.

**Задача:** используя выведенное разложение полного риска модели, обеспечить лучшую скалиброванность моделей, чем при тренировке с обычным риском.

$$R_{Tot} = R_{Tot}^{OHE} + \langle G'(\hat{\eta}), ohe(y) - \eta \rangle$$

## Постановка задачи

- Решается задача классификации.
- Обозначим

$$D = \{(X_i, Y_i)\}_{i=1}^N; \quad X_i \in \mathbb{R}^d, \quad Y_i \in \mathcal{Y} = \{1, \dots, K\}$$

— i.i.d. данные из истинного распределения

$$P(x, y) = P(x)P(y|x)$$

.

- Обозначим истинное условное распределение и модельное как

$$\eta(x) = \mathbb{P}(Y|X = x); \quad \hat{\eta}(x)_\theta = \hat{\eta}(x) = \mathbb{P}(Y|X = x, \theta)$$

соответственно.

- Обозначим лосс функцию

$$l : \Delta^K \times \mathcal{Y} \rightarrow \mathbb{R}$$

## Решение задачи. Обозначения

- Proper scoring rule — такое семейство лосс функций

$l(P, y) : \mathcal{P}_K \times Y \rightarrow \mathbb{R}$ , для которых выполнено:

$$\int l(P, y) dQ(y) = l(P, Q) \geq l(Q, Q) \quad \forall P, Q \in \Delta^K$$

- Под некоторыми условиями на  $l$ , любой такой лосс представим в виде

$$l(\eta, y) = \langle G'(\eta), \eta \rangle - G'_y(\eta) - G(\eta)$$

для какой-то выпуклой функции

$$G : \mathcal{P}_K \rightarrow \mathbb{R}$$

- Введем обозначения полного (total) риска и частного (OHE) риска:

$$R_{Tot} = \langle l(\hat{\eta}), \eta \rangle$$

$$R_{Tot}^{OHE} = l(\hat{\eta}, y) = \langle l(\hat{\eta}), ohe(y) \rangle$$



## Решение задачи. Вывод формулы

Распишем формулу ОНЕ риска для proper scoring rule:

$$R_{Tot}^{OHE} = R_{Tot} + \langle l(\hat{\eta}), ohe(y) - \eta \rangle$$

$$\begin{aligned} \langle l(\hat{\eta}), ohe(y) - \eta \rangle &= \langle (\langle G'(\hat{\eta}), \hat{\eta} \rangle - G(\hat{\eta})) \cdot \mathbb{1} - G'(\hat{\eta}), ohe(y) - \eta \rangle = \\ &= (\langle G'(\hat{\eta}), \hat{\eta} \rangle - G(\hat{\eta})) \cdot \langle \mathbb{1}, ohe(y) - \eta \rangle - \langle G'(\hat{\eta}), ohe(y) - \eta \rangle = \\ &\quad - \langle G'(\hat{\eta}), ohe(y) - \eta \rangle \end{aligned}$$

Итого, получим:

$$R_{Tot} = R_{Tot}^{OHE} + \langle G'(\hat{\eta}), ohe(y) - \eta \rangle$$

# Решение задачи. Аппроксимация риска

Имеем:

$$R_{Tot} = R_{Tot}^{OHE} + \langle G'(\hat{\eta}), ohe(y) - \eta \rangle$$

$$s = ohe(y) - \eta$$

Изучим свойства  $s$ :

- ❶  $\langle \mathbb{1}, s \rangle = 0$ .
- ❷  $\exists! i : s_i \geq 0 \wedge \forall j \neq i : s_j \leq 0$ .  $s_i \in [-1, 0], s_j \in [0, 1]$ .
- ❸ Нет градиентов  $s$  по  $\hat{\eta}$ .
- ❹  $\mathbb{E}_{y \sim \eta(y|x)} s = 0$ , т.к.  $\mathbb{E}_{y \sim \eta(y|x)} ohe(y) = \eta$ .

# Решение задачи. Аппроксимация риска

3 варианта выбора  $s$ , удовлетворяющих свойствам выше:

① (Равномерное)

$$s_k = \begin{cases} \varepsilon(1 - \frac{1}{K}), & \text{if } y = k \\ -\frac{\varepsilon}{K}, & \text{иначе} \end{cases}$$

② (Априорное)

$$s_k = \begin{cases} \varepsilon(1 - \frac{N_k}{N}), & \text{if } y = k \\ -\varepsilon \frac{N_k}{N}, & \text{иначе} \end{cases}$$

③ (Предсказанное)

$$s_k = \begin{cases} \varepsilon(1 - SG(\hat{\eta}_k)), & \text{if } y = k \\ -\varepsilon SG(\hat{\eta}_k), & \text{иначе} \end{cases}$$

# Гипотезы

Гипотезы к проверке:

- ➊ Addition term помогает построить лучше откалиброванную модель. Таким образом, с этим членом  $ECE$  должен быть меньше.
- ➋ Accuracy не сильно снижается при введении addition term. Таким образом, модель остается точной.
- ➌  $ECE$  модели является выпуклой функцией эпохи/точности: сначала она уменьшается, затем увеличивается. Это соответствует следующему утверждению: в начале обучения ONE risk направлен на калибровку, а к концу он подталкивают модель к тому, чтобы быть менее калиброванной и более точной.

# Результаты экспериментов. CE loss



Рис.: Cross Entropy Loss and Uniform addition

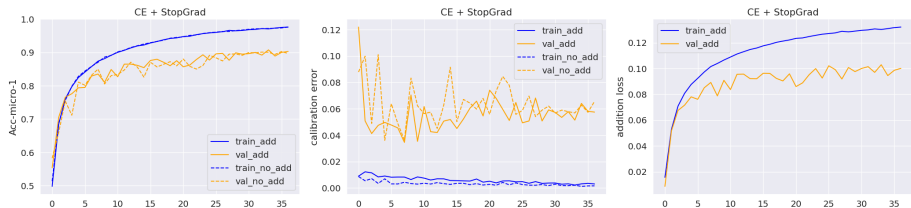


Рис.: Cross Entropy Loss and StopGrad addition

# Результаты экспериментов. Brier loss



Рис.: Brier Loss and Uniform addition

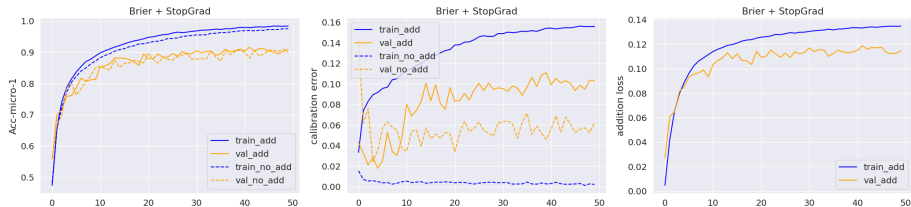


Рис.: Brier Loss and StopGrad addition

# Выводы

- Гипотеза №2 (про Accuracy) подтвердилась
- Гипотеза №3 (про выпуклость ESE) подтвердилась частично, больше экспериментов нужно провести.
- Гипотеза №1 (про калибрационные свойства addition term) под вопросом.

## Следующие шаги

- Провести модельный эксперимент для экспериментального подтверждения гипотез
- Сформулировать и доказать теорему о выпуклости ЕСЕ
- Подтвердить гипотезы с помощью построения доверительных интервалов на картиночных данных.



# Литература

- ① Nikita Durasov et al. “Zigzag: Universal sampling-free uncertainty estimation through two-step inference”. In: arXiv preprint arXiv:2211.11435
- ② Chuan Guo et al. “On calibration of modern neural networks”. In: International conference on machine learning. PMLR. 2017, pp. 1321–1330.
- ③ Nikita Kotelevskii and Maxim Panov. “Predictive Uncertainty Quantification via Risk Decompositions for Strictly Proper Scoring Rules”. In: arXiv preprint arXiv:2402.10727
- ④ Jeremy Nixon et al. “Measuring Calibration in Deep Learning.” In: CVPR workshops. Vol. 2. 7. 2019.
- ⑤ Cheng Wang. “Calibration in deep learning: A survey of the state-of-the-art”. In: arXiv preprint arXiv:2308.01222 (2023).