

Контрастное обучение в задачах компьютерного зрения для повышения интерпретируемости модели

Андрей Семёнов¹²⁴ Владимир Иванов³ Александр Безносовых¹²⁴

¹MIPT ²Yandex Research ³Innopolis University ⁴MMO Laboratory

Май 2024

Sparse Concept Bottleneck Models: Gumbel Tricks in Contrastive Learning

<https://arxiv.org/abs/2404.03323>

ICML 2024 submit

- 1 Постановка задачи
- 2 Сопутствующие работы
- 3 Наш вклад
- 4 Результаты
 - Тяжелые шумы в NLP и CV
 - Наш фреймворк для CBM
 - Результаты Sparse-CBM, ℓ_1 -CBM, Contrastive-CBM
 - Сравнение с другими результатами
- 5 Ссылки на статьи из презентации

Concept Bottleneck модели

$y \in \mathbb{R}, x \in \mathbb{R}^d, c \in \mathbb{R}^k, \{(x^{(i)}, y^{(i)}, c^{(i)})\}_{i=1}^n,$

$g : \mathbb{R}^d \rightarrow \mathbb{R}^k, F : \mathbb{R}^k \rightarrow \mathbb{R}, (F, g) - \text{информационный буттлнек}$

- Independent bottleneck $\hat{F} = \operatorname{argmin}_F \sum_i \mathcal{L}_Y(F(c^{(i)}), y^{(i)})$

$$\hat{g} = \operatorname{argmin}_g \sum_{i,j} \mathcal{L}_{C_j}(g_j(x^{(i)}))$$

- Sequential bottleneck $\hat{F} = \operatorname{argmin}_F \sum_i \mathcal{L}_Y(F(\hat{g}(x^{(i)})), y^{(i)})$

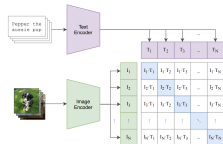
- Joint bottleneck

$$\hat{F}, \hat{g} = \operatorname{argmin}_{F, g} \sum_i [\lambda_1 \mathcal{L}_Y(F(g(x^{(i)}))) + \lambda_2 \sum_j \mathcal{L}_{C_j}(g_j(x^{(i)}))]$$

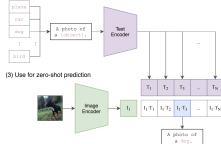
- Standard model $\hat{F}, \hat{g} = \operatorname{argmin}_{F, g} \sum_i \mathcal{L}_Y(F(g(x^{(i)})), y^{(i)})$

Сопутствующие работы

(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

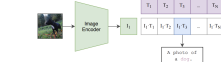


Рис.: CLIP

Label-free
CBM

Step 1: Generate and
filter concept set

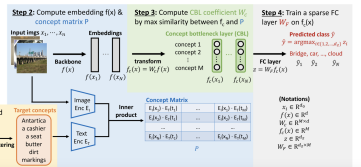


Рис.: Label-free

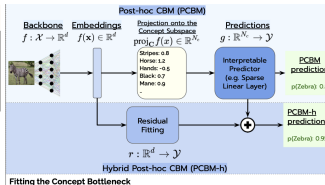
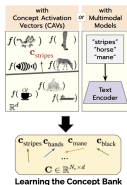


Рис.: Post-hoc CBM

- Мы предложили два новых варианта архитектуры и алгоритм для обучения Concept Bottleneck моделей (CBM).
- Мы формально описываем алгоритм для повышения точности CLIP – Concept Matrix Search (CMS), и в то же время делаем модель более интерпретируемой. Мы также приводим анализ латентного пространства CLIP.
- Получили неожиданный результат о полезности внутренних разреженных слоев в CBM.
- Исследовали задачу распределенной оптимизации, предложили способы устранения распределений тяжелых шумов в совместной задаче NLP и CV.

Градиентный клиппинг

SGD:

$$x^{k+1} = x^k - \gamma \nabla f(x^k, \xi^k) \quad (1)$$

Clipped-SGD:

$$x^{k+1} = x^k - \gamma \text{clip}(\nabla f(x^k, \xi^k), \lambda) \quad (2)$$

- $\text{clip}(x, \lambda) = \min\{1, \lambda/\|x\|\}x$
- $\mathbb{E}_{\xi^k}[\text{clip}(\nabla f(x^k, \lambda))] \neq \nabla f(x^k)$
- При $\beta_1 = 0$ Adam можно интерпретировать как Clipped-SGD с "адаптивным" λ

Тяжелый шум: теоретическая справка

Случайный вектор X имеет распределение с легкими хвостами если:

$$\mathbb{P}\{\|X - \mathbb{E}[X]\|\} \leq 2 \exp\left(-\frac{b^2}{2\sigma^2}\right) \forall b > 0, \quad (3)$$

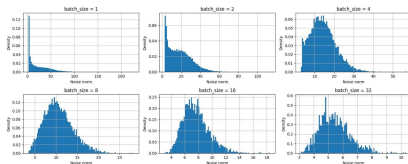
что эквивалентно

$$\mathbb{E}\left[\exp\left(\frac{\|X - \mathbb{E}[X]\|^2}{\sigma^2}\right)\right] \leq \exp 1 \quad (4)$$

Для задач с шумом в виде тяжелых хвостов в теории используют предположение:

$$\forall \alpha \in (1, 2] : \mathbb{E}[\|X - \mathbb{E}[X]\|^\alpha] \leq \sigma^\alpha \quad (5)$$

Тяжелый шум: экспериментальное подтверждение



Нормы градиентов, CIFAR10.

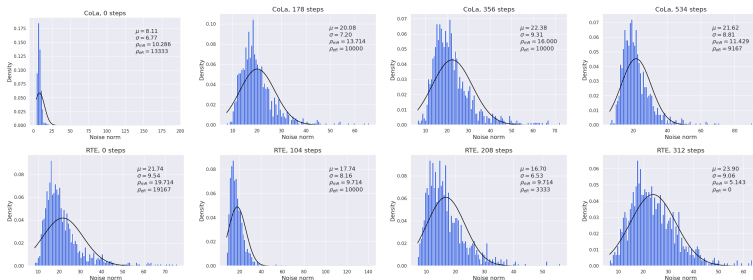
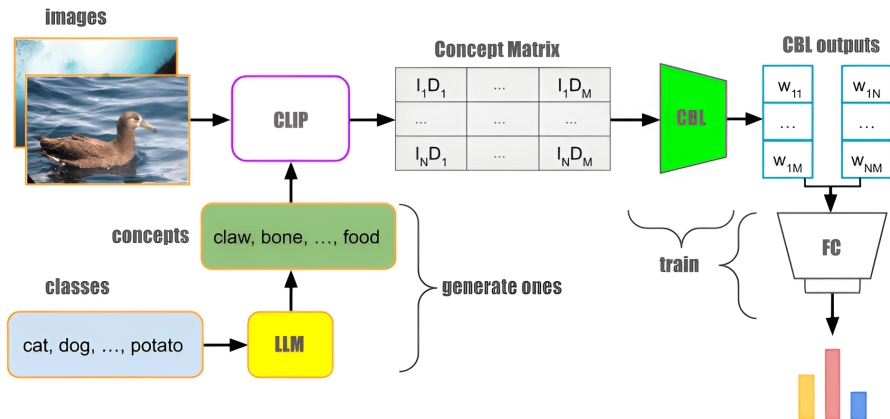


Рис.: Оценка градиентного шума для Adam на датасетах CoLa (первая строка) и RTE (вторая строка).

Наша архитектура для CBM

Concept Bottleneck Model framework



Наш CBM фреймворк.

Результаты: Общая схема для наших CBM

- i -эмбединг картинки, d -эмбединг концепта, $\mathcal{D} = (x, t, l)$ – датасет.
- CLIP: $\psi(x, t) = (\langle i, d_1 \rangle, \dots, \langle i, d_{|D|} \rangle)^\top \in \mathbb{R}^{|D|}$

$$\min_{W_{\text{CBL}}} \mathbb{E}_{(x, t, l) \sim \mathcal{D}} [\mathcal{L}_{\text{CBL}}(W_{\text{CBL}} \psi(x, t))], \quad (6)$$

$$\min_{W_{\text{F}}} \mathbb{E}_{(x, t, l) \sim \mathcal{D}} [\mathcal{L}_{\text{CE}}(W_{\text{F}} W_{\text{CBL}} \psi(x, t), l)]. \quad (7)$$

Contrastive-CBM:

$$-\frac{1}{2|B|} \sum_{k=1}^{|B|} \left(\log \frac{e^{\alpha \langle w_k, \varphi_k \rangle}}{\sum_{j=1}^{|B|} e^{\alpha \langle w_k, \varphi_j \rangle}} + \log \frac{e^{\alpha \langle w_k, \varphi_k \rangle}}{\sum_{j=1}^{|B|} e^{\alpha \langle w_j, \varphi_k \rangle}} \right). \quad (8)$$

Sparse-CBM:

$$z = 1 \left(\arg \max_k [g_k + \log \pi_k] \right), \quad (9)$$

$$- \frac{1}{2|B|} \sum_{k=1}^{|B|} \left(\log \frac{e^{(\log(\alpha \langle w_k, \varphi_k \rangle) + g_k)/\tau}}{\sum_{j=1}^{|B|} e^{(\log(\alpha \langle w_k, \varphi_j \rangle) + g_j)/\tau}} \right. \\ \left. + \log \frac{e^{(\log(\alpha \langle w_k, \varphi_k \rangle) + g_k)/\tau}}{\sum_{j=1}^{|B|} e^{(\log(\alpha \langle w_j, \varphi_k \rangle) + g_j)/\tau}} \right). \quad (10)$$

$u \in \text{Uniform}(0, 1)$, $g = -\log \log u \in \text{Gumbel}(0, 1)$

ℓ_1 -CBM:

$$\min_{W_{\text{CBL}}} \mathbb{E}_{(x,t,l) \sim \mathcal{D}} \left[\mathcal{L}_{\text{CE}}(W_{\text{F}} W_{\text{CBL}} \psi(x, t), l) + \frac{\lambda}{|D|} \Omega(W_{\text{CBL}}) \right], \quad (11)$$

где $\Omega(W_{\text{CBL}})$ соответствует регуляризатору. Мы используем:

$$\Omega(W_{\text{CBL}}) = \|W_{\text{CBL}}\|_1 \quad (12)$$

Результаты: Concept Matrix Search алгоритм

- 1: **Input:** Batch of image embeddings $I_{|B|}$, labels, all classes C and concepts D embeddings.
- 2: Build $\mathcal{V} \in \mathbb{R}^{|B| \times |D|}$, $\mathcal{T} \in \mathbb{R}^{|C| \times |D|}$ matrices, store \mathcal{T} .
- 3: **for** $k = 0, 1, 2, \dots, |B| - 1$ **do**
- 4: **for** $m = 0, 1, 2, \dots, |C| - 1$ **do**
- 5: Compute and store $\cos(\mathcal{V}_{k,\cdot}^\top, \mathcal{T}_{m,\cdot}^\top)$
- 6: **end for**
- 7: Find $m_{\max} = \max_m \cos(\mathcal{V}_{k,\cdot}^\top, \mathcal{T}_{m,\cdot}^\top)$
- 8: **if** $\text{label}(k) = m_{\max}$ **then**
- 9: the hypothesis has been proven, increase Accuracy
- 10: **else**
- 11: the hypothesis has been disproved
- 12: **end if**
- 13: **end for**
- 14: **return** Mean accuracy

Сравнение наших методов с предшествующими

Таблица: Сравнение перформанса Bottleneck моделей на основных датасетах. Мы наблюдаем превосходство Sparse-CBM над другими архитектурами на CIFAR10, CIFAR100 и CUB200 датасетах.

MODEL	CIFAR10	CIFAR100	IMAGENET	CUB200	PLACES
SPARSE-CBM	91.17%	74.88%	71.61%	80.02%	41.34%
ℓ_1 -CBM	85.11%	73.24%	71.02%	74.91%	40.87%
CONTR-CBM	84.75%	68.46%	70.22%	67.04%	40.22%
[2]	86.40%	65.13%	71.95%	74.31%	43.68%
[3]	83.34%	57.20%	62.57%	63.92%	39.66%
[4]	87.90%	69.10%	70.40%	71.80%	39.43%
PROBING	96.12%	80.03%	83.90%	79.29%	48.33%

Сравнение Concept Matrix Search алгоритма с предшествующими методами

Таблица: Сравнение CMS и "DescriptionCLS"[5] на основных датасетах.

METHOD	CIFAR10	CIFAR100	IMAGENET	CUB200	PLACES
CMS	85.03%	62.95%	77.82%	65.17%	39.43%
[5]	81.61%	68.32%	75.00%	63.46%	40.55%
ZERO-SHOT	81.79%	52.84%	76.20%	62.63%	41.12%

Визуализация пространства эмбедингов

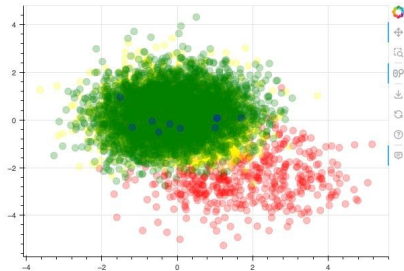


Рис.: CIFAR10 t-SNE. Зелёные точки – проекция эмбедингов концептов, синие – проекция эмбедингов классов, красные – картинки, и желтые – случайные слова.

Визуализация CBM

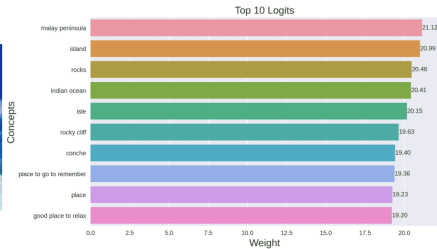


Рис.: Концепты, извлекаемые с помощью CLIP.

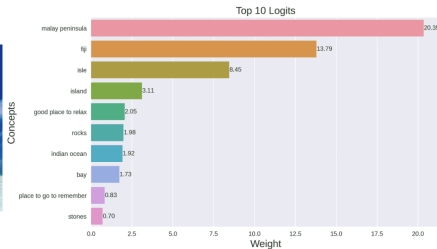


Рис.: Концепты, извлекаемые с помощью Sparse-CBM.

Визуализация CBM

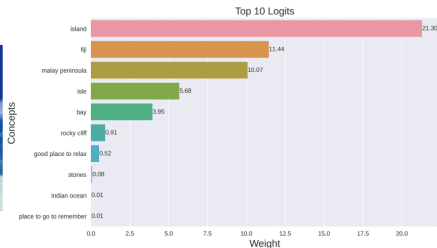


Рис.: Концепты, извлекаемые с помощью ℓ_1 -CBM.

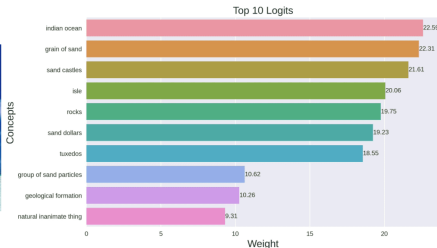


Рис.: Концепты, извлекаемые с помощью Contrastive-CBM.

- [1] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021.
- [2] Tuomas Oikarinen, Subhro Das, Lam M. Nguyen, and Tsui-Wei Weng. Label-free concept bottleneck models. In The Eleventh International Conference on Learning Representations, 2023.
- [3] Yuksekgonul, M., Wang, M., and Zou, J. Post-hoc concept bottleneck models, 2023.
- [4] Yang, Y., Panagopoulou, A., Zhou, S., Jin, D., Callison-Burch, C., and Yatskar, M. Language in a bottle: Language model guided concept bottlenecks for interpretable image classification, 2023.
- [5] Menon, S. and Vondrick, C. Visual classification via description from large language models. ICLR, 2023.