

# A Metric Learning Reality Check

Охрименко Дмитрий, 172

# Metric learning

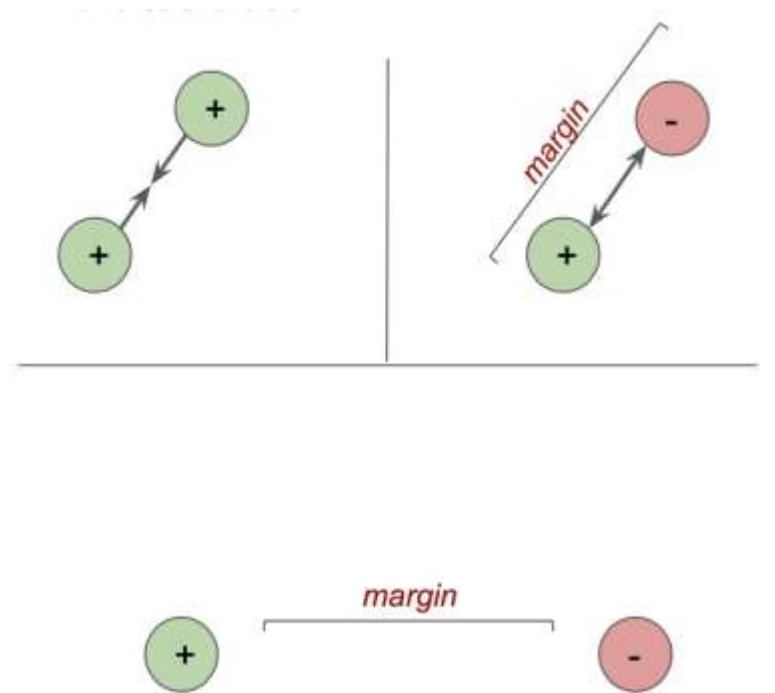
- Задача – определить похожесть двух объектов
- Пытаемся сопоставить данные с пространством эмбеддингов, где похожие данные находятся близко друг к другу, а разные данные далеко друг от друга.
- Можно использовать embedding losses и classificational losses

# Classificational losses

- Основаны на использовании весовой матрицы, где каждый столбец соответствует определенному классу
- Обучение состоит из матричного умножения весов на векторы эмбеддингов для получения логитов и последующего применения функции потерь к логитам.
- Есть множество вариаций (normalized softmax loss, ProxyNCA, SphereFace, CosFace, ArcFace, SoftTriple losses)

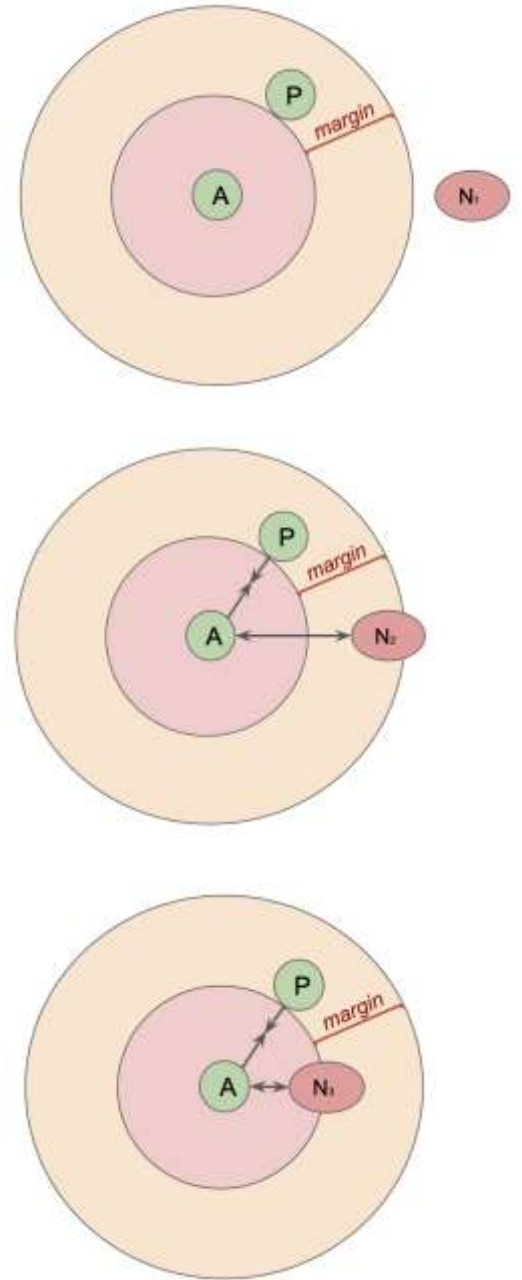
# Embedding losses

- Contrastive loss – одинаковые сдвигаем, только если они дальше  $m_{pos}$ , разные отталкиваем, только если они ближе  $m_{neg}$
- $L_{contrastive} = [d_p - m_{pos}]_+ + [m_{neg} - d_n]_+$
- $m_{pos}$  может быть равен 0 (как на картинке)



# Embedding losses

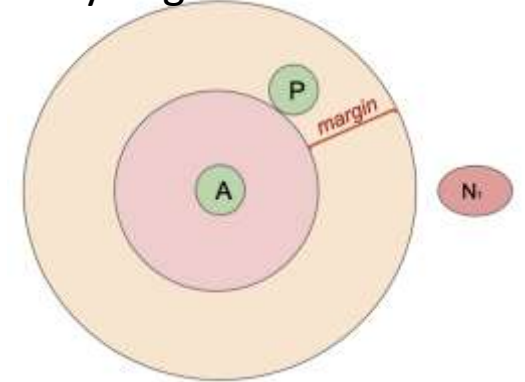
- Triplet loss – A (anchor) похож на Positive больше, чем на Negative. Цель – сделать расстояние AP меньше чем AN на выбранный margin.
- $L_{triplet} = [d_{ap} - d_{an} + m]_+$



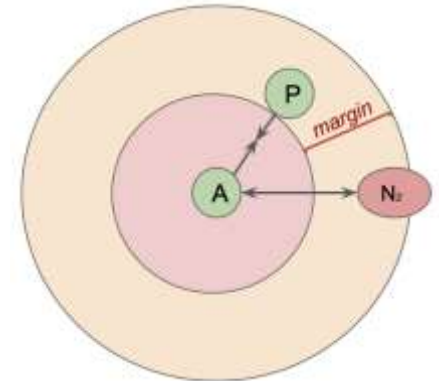
# Pair and triplet mining

- Offline mining – выбираем тройки до конструирования батчей
- Online mining – находим тройки в каждом рандомном батче
- Easy negatives – не влияют на обучение
- Hard negatives – вызывают обучение на плохих данных
- Semi-hard negatives – остается только этот вариант

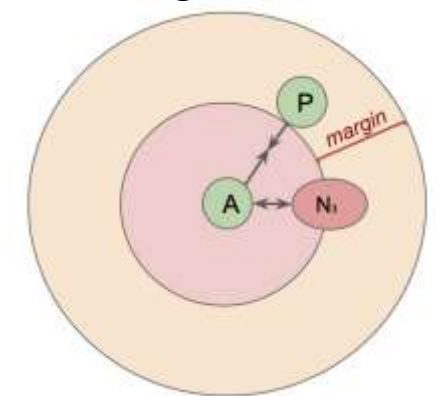
Easy negatives



Semi-hard negatives



Hard negatives



# Проблемы в существующих статьях

- Нечестные сравнения
- Неточность используемых метрик
- Переобучение на тестовой выборке

# Нечестные сравнения

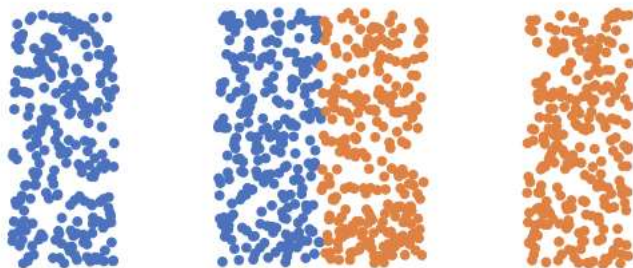
- Чтобы заявить, что новый алгоритм превосходит существующие методы, важно сохранять как можно больше параметров постоянными, однако в большинстве существующих работ этого не происходит
- Простые способы повысить точность - обновить сетевую архитектуру или использовать более сложные способы увеличения изображения



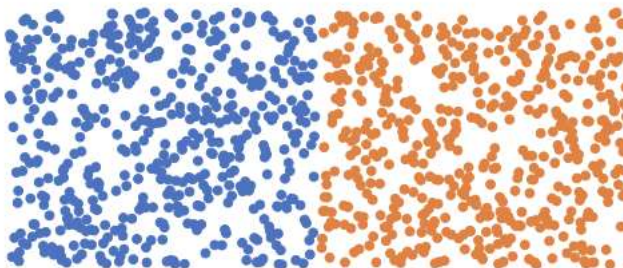
# Неточность используемых метрик

- Обычно используются Recall@K, Normalized Mutual Information (NMI) и F1 score
- Три картинки снизу демонстрируют различные разделения на классы, однако показания метрик очень схожи

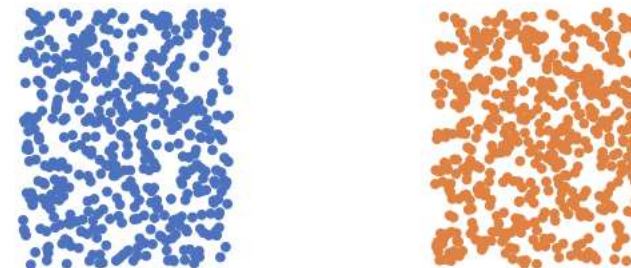
NMI: 95.6% F1: 100% R@1: 99%,  
R-Precision: 77.4% MAP@R: 71.4%



NMI: 100% F1: 100% R@1: 99.8%  
R-Precision: 83.3% MAP@R: 77.9%



NMI: 100% F1: 100% R@1: 100%,  
R-Precision: 99.8% MAP@R: 99.8%



# Переобучение на тестовой выборке

- В большинстве статей данные делятся на тестовые и обучающие 50/50 и гиперпараметры настраиваются в зависимости от поведения на тестовых данных
- Это может привести к переобучению на тестовых данных, что является большой проблемой

# Предлагаемые решения

- Честные сравнения и воспроизводимость
- Информативные метрики
- Подбор гиперпараметров через кросс-валидацию

# Честные сравнения и воспроизводимость

- Использовалась предоубченная на ImageNet BN-Inception сеть с размером выходных эмбеддингов = 128
- Размер батча = 32
- Батчи создаются случайной выборкой  $C$  классов, а затем случайным разделением  $M$  картинок на  $C$  классов
- Для embedding losses  $C=8$ ,  $M=4$ ; для classification losses  $C=32$ ,  $M=1$
- Во время обучения изображения увеличиваются с использованием стратегии обрезки со случайным изменением размера.
- Все параметры сети оптимизированы с помощью RMSprop с  $lr = 1e-6$
- Эмбеддинги L2 нормализуются перед подсчетом лосса и во время оценки

# Информативная метрика

- $R\text{-precision} = \frac{r}{R}$ ,  $R$  – ближайших к запросу элементов,  $r$  – относящиеся к тому же классу
- $R\text{-precision}$  не учитывает ранжирование правильных выборов, поэтому используем  $MAP@R$  – Mean Average Precision at  $R$ :

$$MAP@R = \frac{1}{R} \sum_{i=1}^R P(i)$$
$$P(i) = \begin{cases} precision_i, & \text{если } i\text{-е извлечение верно} \\ 0, & \text{иначе} \end{cases}$$

# Информативная метрика

- MAP@R информативнее R@1
- Может быть вычислен прямо из пространства эмбедингов
- Вознаграждает хорошо сгруппированные пространства эмбедингов
- Легкая в понимании
- Стабильнее, чем R@1

$$MAP@R = \frac{1}{R} \sum_{i=1}^R P(i)$$
$$P(i) = \begin{cases} precision_i, & \text{если } i\text{-е извлечение верно} \\ 0, & \text{иначе} \end{cases}$$

# Информативная метрика

- MAP@R информативнее R@1
- Может быть вычислен прямо из пространства эмбедингов
- Вознаграждает хорошо сгруппированные пространства эмбедингов
- Легкая в понимании
- Стабильнее, чем R@1

| Retrieval results                                 | Recall@1 | R-Precision | MAP@R |
|---|----------|-------------|-------|
| 10 results, of which only the 1st is correct      | 100      | 10          | 10    |
| 10 results, of which the 1st and 10th are correct | 100      | 20          | 12    |
| 10 results, of which the 1st and 2nd are correct  | 100      | 20          | 20    |
| 10 results, of which all 10 are correct           | 100      | 100         | 100   |

# Подбор гиперпараметров через кросс-валидацию

- Запускается 50 итераций байесовской проверки, каждая состоит из 4-х кратной перекрестной проверки
- В результате наборы для обучения и проверки всегда не пересекаются с классами, поэтому оптимизация производительности набора проверки должна быть хорошим показателем точности для задач с открытым набором
- Гиперпараметры оптимизированы для максимизации средней точности проверки. Для получения лучших гиперпараметров загружается контрольная точка наивысшей точности для каждого раздела обучающего набора, вычисляются его эмбединги для тестового набора и нормализуются L2 нормой. Точность вычисляется объединенным и раздельным способами
- Результаты слабо подвержены начальному шуму, так как проводится 10 тренировочных прогонов с использованием лучших гиперпараметров и выбирается среднее



# Эксперименты

- Провели эксперименты с 13 лоссами на трех датасетах: CUB200, Cars196 и Stanford Online Products (SOP)

| Method                                       | Year | Loss type      |
|--|------|----------------|
| Contrastive [16]                             | 2006 | Embedding      |
| Triplet [63]                                 | 2006 | Embedding      |
| NT-Xent [50,38,6]                            | 2016 | Embedding      |
| ProxyNCA [35]                                | 2017 | Classification |
| Margin [65]                                  | 2017 | Embedding      |
| Margin / class [65]                          | 2017 | Embedding      |
| Normalized Softmax (N. Softmax) [58,31,72]   | 2017 | Classification |
| CosFace [57,59]                              | 2018 | Classification |
| ArcFace [11]                                 | 2019 | Classification |
| FastAP [3]                                   | 2019 | Embedding      |
| Signal to Noise Ratio Contrastive (SNR) [70] | 2019 | Embedding      |
| MultiSimilarity (MS) [62]                    | 2019 | Embedding      |
| MS+Miner [62]                                | 2019 | Embedding      |
| SoftTriple [41]                              | 2019 | Classification |

# Эксперименты

**Table 4.** Accuracy on CUB200

|              | Concatenated (512-dim)             |                                    |                                    | Separated (128-dim)                |                                    |                                    |
|--------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|
|              | P@1                                | RP                                 | MAP@R                              | P@1                                | RP                                 | MAP@R                              |
| Pretrained   | 51.05                              | 24.85                              | 14.21                              | 50.54                              | 25.12                              | 14.53                              |
| Contrastive  | <b>68.13 <math>\pm</math> 0.31</b> | 37.24 $\pm$ 0.28                   | 26.53 $\pm$ 0.29                   | 59.73 $\pm$ 0.40                   | 31.98 $\pm$ 0.29                   | 21.18 $\pm$ 0.28                   |
| Triplet      | 64.24 $\pm$ 0.26                   | 34.55 $\pm$ 0.24                   | 23.69 $\pm$ 0.23                   | 55.76 $\pm$ 0.27                   | 29.55 $\pm$ 0.16                   | 18.75 $\pm$ 0.15                   |
| NT-Xent      | 66.61 $\pm$ 0.29                   | 35.96 $\pm$ 0.21                   | 25.09 $\pm$ 0.22                   | 58.12 $\pm$ 0.23                   | 30.81 $\pm$ 0.17                   | 19.87 $\pm$ 0.16                   |
| ProxyNCA     | 65.69 $\pm$ 0.43                   | 35.14 $\pm$ 0.26                   | 24.21 $\pm$ 0.27                   | 57.88 $\pm$ 0.30                   | 30.16 $\pm$ 0.22                   | 19.32 $\pm$ 0.21                   |
| Margin       | 63.60 $\pm$ 0.48                   | 33.94 $\pm$ 0.27                   | 23.09 $\pm$ 0.27                   | 54.78 $\pm$ 0.30                   | 28.86 $\pm$ 0.18                   | 18.11 $\pm$ 0.17                   |
| Margin/class | 64.37 $\pm$ 0.18                   | 34.59 $\pm$ 0.16                   | 23.71 $\pm$ 0.16                   | 55.56 $\pm$ 0.16                   | 29.32 $\pm$ 0.15                   | 18.51 $\pm$ 0.13                   |
| N. Softmax   | 65.65 $\pm$ 0.30                   | 35.99 $\pm$ 0.15                   | 25.25 $\pm$ 0.13                   | 58.75 $\pm$ 0.19                   | 31.75 $\pm$ 0.12                   | 20.96 $\pm$ 0.11                   |
| CosFace      | 67.32 $\pm$ 0.32                   | <b>37.49 <math>\pm</math> 0.21</b> | <b>26.70 <math>\pm</math> 0.23</b> | 59.63 $\pm$ 0.36                   | 31.99 $\pm$ 0.22                   | 21.21 $\pm$ 0.22                   |
| ArcFace      | 67.50 $\pm$ 0.25                   | 37.31 $\pm$ 0.21                   | 26.45 $\pm$ 0.20                   | <b>60.17 <math>\pm</math> 0.32</b> | <b>32.37 <math>\pm</math> 0.17</b> | <b>21.49 <math>\pm</math> 0.16</b> |
| FastAP       | 63.17 $\pm$ 0.34                   | 34.20 $\pm$ 0.20                   | 23.53 $\pm$ 0.20                   | 55.58 $\pm$ 0.31                   | 29.72 $\pm$ 0.16                   | 19.09 $\pm$ 0.16                   |
| SNR          | 66.44 $\pm$ 0.56                   | 36.56 $\pm$ 0.34                   | 25.75 $\pm$ 0.36                   | 58.06 $\pm$ 0.39                   | 31.21 $\pm$ 0.28                   | 20.43 $\pm$ 0.28                   |
| MS           | 65.04 $\pm$ 0.28                   | 35.40 $\pm$ 0.12                   | 24.70 $\pm$ 0.13                   | 57.60 $\pm$ 0.24                   | 30.84 $\pm$ 0.13                   | 20.15 $\pm$ 0.14                   |
| MS+Miner     | 67.73 $\pm$ 0.18                   | 37.37 $\pm$ 0.19                   | 26.52 $\pm$ 0.18                   | 59.41 $\pm$ 0.30                   | 31.93 $\pm$ 0.15                   | 21.01 $\pm$ 0.14                   |
| SoftTriple   | 67.27 $\pm$ 0.39                   | 37.34 $\pm$ 0.19                   | 26.51 $\pm$ 0.20                   | 59.94 $\pm$ 0.33                   | 32.12 $\pm$ 0.14                   | 21.31 $\pm$ 0.14                   |

# Эксперименты

**Table 5.** Accuracy on Cars196

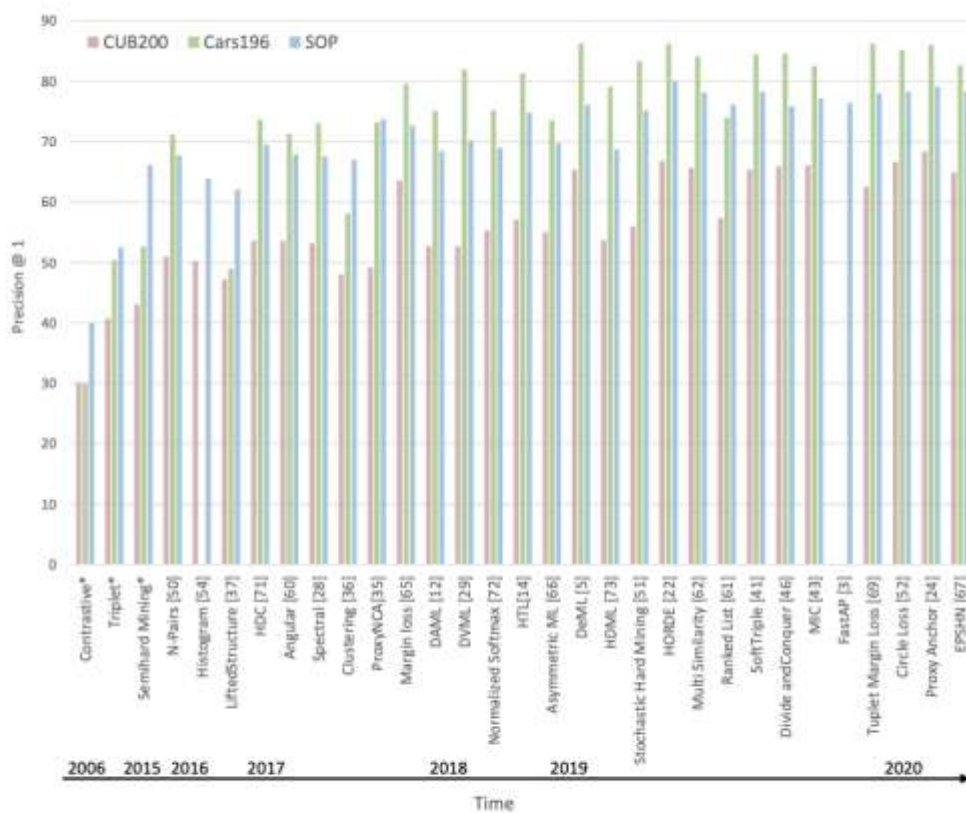
|                | Concatenated (512-dim)             |                                    |                                    | Separated (128-dim)                |                                    |                                    |
|----------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|
|                | P@1                                | RP                                 | MAP@R                              | P@1                                | RP                                 | MAP@R                              |
| Pretrained     | 46.89                              | 13.77                              | 5.91                               | 43.27                              | 13.37                              | 5.64                               |
| Contrastive    | 81.78 $\pm$ 0.43                   | 35.11 $\pm$ 0.45                   | 24.89 $\pm$ 0.50                   | 69.80 $\pm$ 0.38                   | 27.78 $\pm$ 0.34                   | 17.24 $\pm$ 0.35                   |
| Triplet        | 79.13 $\pm$ 0.42                   | 33.71 $\pm$ 0.45                   | 23.02 $\pm$ 0.51                   | 65.68 $\pm$ 0.58                   | 26.67 $\pm$ 0.36                   | 15.82 $\pm$ 0.36                   |
| NT-Xent        | 80.99 $\pm$ 0.54                   | 34.96 $\pm$ 0.38                   | 24.40 $\pm$ 0.41                   | 68.16 $\pm$ 0.36                   | 27.66 $\pm$ 0.23                   | 16.78 $\pm$ 0.24                   |
| ProxyNCA       | 83.56 $\pm$ 0.27                   | 35.62 $\pm$ 0.28                   | 25.38 $\pm$ 0.31                   | 73.46 $\pm$ 0.23                   | 28.90 $\pm$ 0.22                   | 18.29 $\pm$ 0.22                   |
| Margin         | 81.16 $\pm$ 0.50                   | 34.82 $\pm$ 0.31                   | 24.21 $\pm$ 0.34                   | 68.24 $\pm$ 0.35                   | 27.25 $\pm$ 0.19                   | 16.40 $\pm$ 0.20                   |
| Margin / class | 80.04 $\pm$ 0.61                   | 33.78 $\pm$ 0.51                   | 23.11 $\pm$ 0.55                   | 67.54 $\pm$ 0.60                   | 26.68 $\pm$ 0.40                   | 15.88 $\pm$ 0.39                   |
| N. Softmax     | 83.16 $\pm$ 0.25                   | 36.20 $\pm$ 0.26                   | 26.00 $\pm$ 0.30                   | 72.55 $\pm$ 0.18                   | 29.35 $\pm$ 0.20                   | 18.73 $\pm$ 0.20                   |
| CosFace        | <b>85.52 <math>\pm</math> 0.24</b> | 37.32 $\pm$ 0.28                   | 27.57 $\pm$ 0.30                   | <b>74.67 <math>\pm</math> 0.20</b> | 29.01 $\pm$ 0.11                   | 18.80 $\pm$ 0.12                   |
| ArcFace        | 85.44 $\pm$ 0.28                   | 37.02 $\pm$ 0.29                   | 27.22 $\pm$ 0.30                   | 72.10 $\pm$ 0.37                   | 27.29 $\pm$ 0.17                   | 17.11 $\pm$ 0.18                   |
| FastAP         | 78.45 $\pm$ 0.52                   | 33.61 $\pm$ 0.54                   | 23.14 $\pm$ 0.56                   | 65.08 $\pm$ 0.36                   | 26.59 $\pm$ 0.36                   | 15.94 $\pm$ 0.34                   |
| SNR            | 82.02 $\pm$ 0.48                   | 35.22 $\pm$ 0.43                   | 25.03 $\pm$ 0.48                   | 69.69 $\pm$ 0.46                   | 27.55 $\pm$ 0.25                   | 17.13 $\pm$ 0.26                   |
| MS             | 85.14 $\pm$ 0.29                   | <b>38.09 <math>\pm</math> 0.19</b> | <b>28.07 <math>\pm</math> 0.22</b> | 73.77 $\pm$ 0.19                   | <b>29.92 <math>\pm</math> 0.16</b> | <b>19.32 <math>\pm</math> 0.18</b> |
| MS+Miner       | 83.67 $\pm$ 0.34                   | 37.08 $\pm$ 0.31                   | 27.01 $\pm$ 0.35                   | 71.80 $\pm$ 0.22                   | 29.44 $\pm$ 0.21                   | 18.86 $\pm$ 0.20                   |
| SoftTriple     | 84.49 $\pm$ 0.26                   | 37.03 $\pm$ 0.21                   | 27.08 $\pm$ 0.21                   | 73.69 $\pm$ 0.21                   | 29.29 $\pm$ 0.16                   | 18.89 $\pm$ 0.16                   |

# Эксперименты

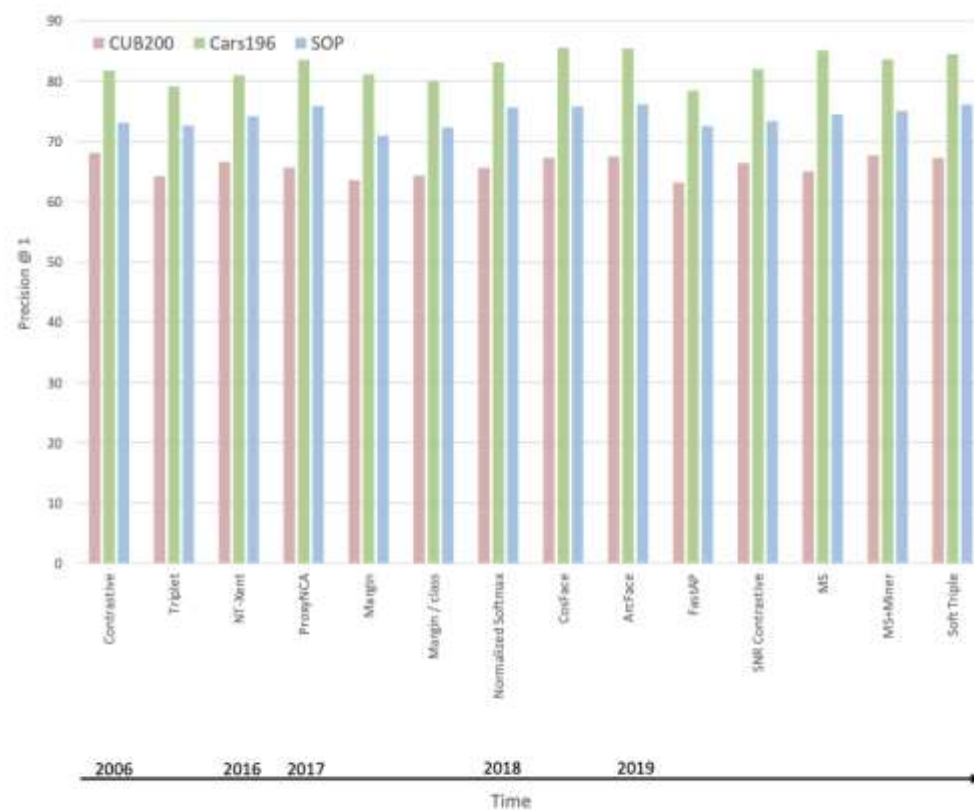
**Table 6.** Accuracy on SOP

|                | Concatenated (512-dim)             |                                    |                                    | Separated (128-dim)                |                                    |                                    |
|----------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|
|                | P@1                                | RP                                 | MAP@R                              | P@1                                | RP                                 | MAP@R                              |
| Pretrained     | 50.71                              | 25.97                              | 23.44                              | 47.25                              | 23.84                              | 21.36                              |
| Contrastive    | $73.12 \pm 0.20$                   | $47.29 \pm 0.24$                   | $44.39 \pm 0.24$                   | $69.34 \pm 0.26$                   | $43.41 \pm 0.28$                   | $40.37 \pm 0.28$                   |
| Triplet        | $72.65 \pm 0.28$                   | $46.46 \pm 0.38$                   | $43.37 \pm 0.37$                   | $67.33 \pm 0.34$                   | $40.94 \pm 0.39$                   | $37.70 \pm 0.38$                   |
| NT-Xent        | $74.22 \pm 0.22$                   | $48.35 \pm 0.26$                   | $45.31 \pm 0.25$                   | $69.88 \pm 0.19$                   | $43.51 \pm 0.21$                   | $40.31 \pm 0.20$                   |
| ProxyNCA       | $75.89 \pm 0.17$                   | $50.10 \pm 0.22$                   | $47.22 \pm 0.21$                   | $71.30 \pm 0.20$                   | $44.71 \pm 0.21$                   | $41.74 \pm 0.21$                   |
| Margin         | $70.99 \pm 0.36$                   | $44.94 \pm 0.43$                   | $41.82 \pm 0.43$                   | $65.78 \pm 0.34$                   | $39.71 \pm 0.40$                   | $36.47 \pm 0.39$                   |
| Margin / class | $72.36 \pm 0.30$                   | $46.41 \pm 0.40$                   | $43.32 \pm 0.41$                   | $67.56 \pm 0.42$                   | $41.37 \pm 0.48$                   | $38.15 \pm 0.49$                   |
| N. Softmax     | $75.67 \pm 0.17$                   | $50.01 \pm 0.22$                   | $47.13 \pm 0.22$                   | <b><math>71.65 \pm 0.14</math></b> | <b><math>45.32 \pm 0.17</math></b> | <b><math>42.35 \pm 0.16</math></b> |
| CosFace        | $75.79 \pm 0.14$                   | $49.77 \pm 0.19$                   | $46.92 \pm 0.19$                   | $70.71 \pm 0.19$                   | $43.56 \pm 0.21$                   | $40.69 \pm 0.21$                   |
| ArcFace        | <b><math>76.20 \pm 0.27</math></b> | <b><math>50.27 \pm 0.38</math></b> | <b><math>47.41 \pm 0.40</math></b> | $70.88 \pm 1.51$                   | $44.00 \pm 1.26$                   | $41.11 \pm 1.22$                   |
| FastAP         | $72.59 \pm 0.26$                   | $46.60 \pm 0.29$                   | $43.57 \pm 0.28$                   | $68.13 \pm 0.25$                   | $42.06 \pm 0.25$                   | $38.88 \pm 0.25$                   |
| SNR            | $73.40 \pm 0.09$                   | $47.43 \pm 0.13$                   | $44.54 \pm 0.13$                   | $69.45 \pm 0.10$                   | $43.34 \pm 0.12$                   | $40.31 \pm 0.12$                   |
| MS             | $74.50 \pm 0.24$                   | $48.77 \pm 0.32$                   | $45.79 \pm 0.32$                   | $70.43 \pm 0.33$                   | $44.25 \pm 0.38$                   | $41.15 \pm 0.38$                   |
| MS+Miner       | $75.09 \pm 0.17$                   | $49.51 \pm 0.20$                   | $46.55 \pm 0.20$                   | $71.25 \pm 0.15$                   | $45.19 \pm 0.16$                   | $42.10 \pm 0.16$                   |
| SoftTriple     | $76.12 \pm 0.17$                   | $50.21 \pm 0.18$                   | $47.35 \pm 0.19$                   | $70.88 \pm 0.20$                   | $43.83 \pm 0.20$                   | $40.92 \pm 0.20$                   |

# Paper vs Reality

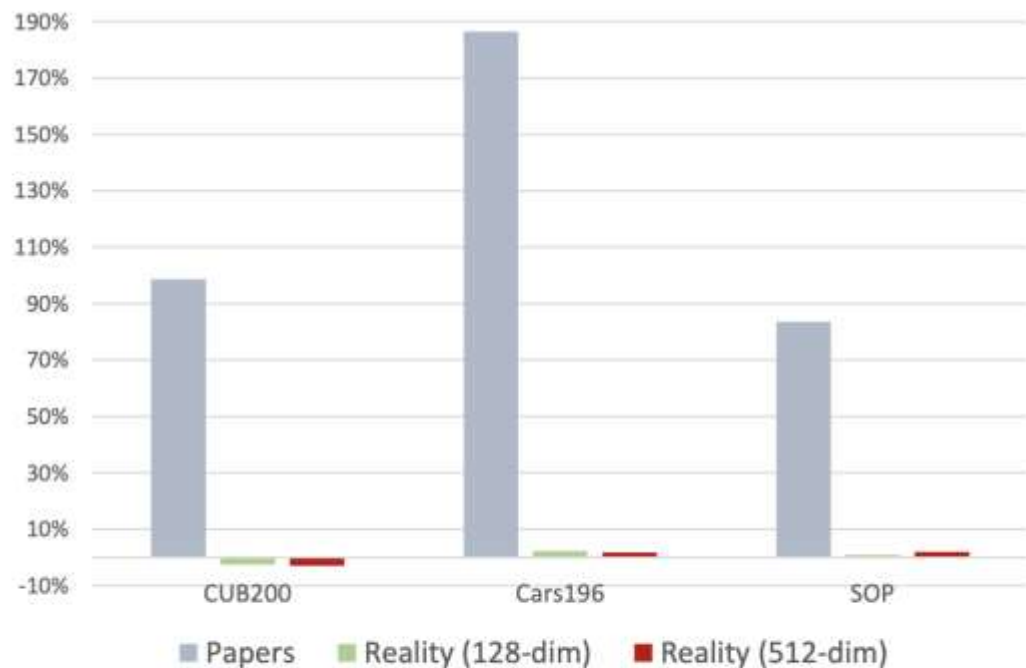


(a) The trend according to papers

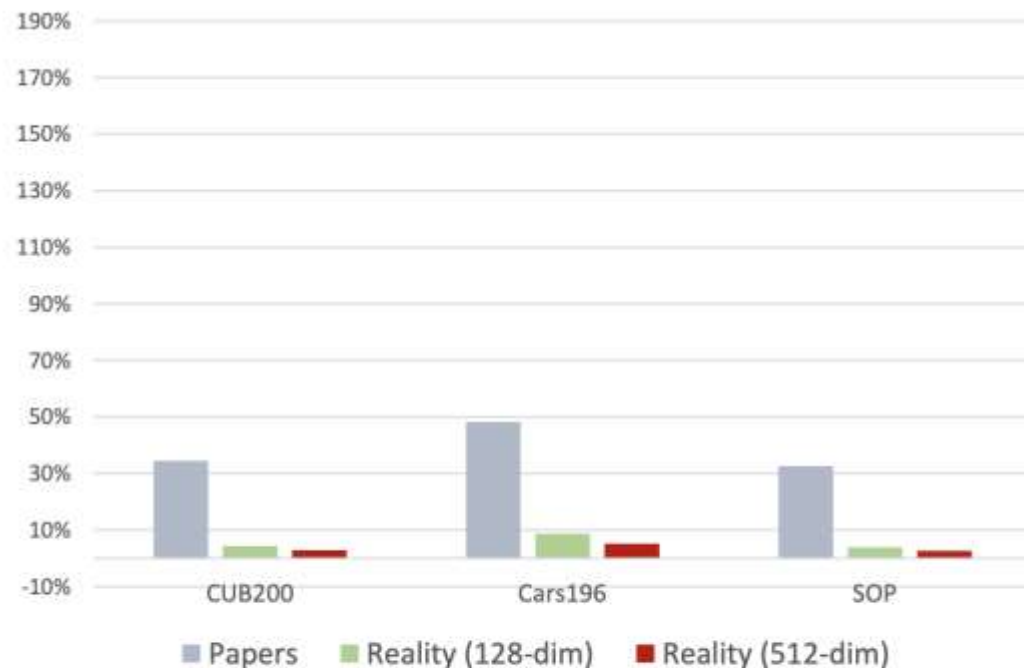


(b) The trend according to reality

# Paper vs Reality



(a) Relative improvement over the contrastive loss



(b) Relative improvement over the triplet loss

# Paper vs Reality

- Обнаружилось, что в статьях резко преувеличены улучшения по сравнению с contrastive и triplet лоссами. Это происходит из-за чрезвычайно низкой точности, связанной с этими потерями.
- При правильной реализации мы получим, что методы 2006 и 2019 работают +- одинаково. Другими словами, алгоритмы метрического обучения не добились впечатляющего прогресса, которого, как они утверждают, добились.

# Заключение

- Обнаружили некоторые недостатки в современной литературе по Metric Learning:
- Несправедливые сравнения, вызванные изменениями в сетевой архитектуре, размерах встраиваемых файлов, способах увеличения изображений и оптимизаторах.
- Использование метрик, которые либо вводят в заблуждение, либо не дают полной картины пространства эмбедингов.
- Обучение без валидационной выборки, то есть с обратной связью только с тестовым набором.



# Заключение

- Затем провели эксперименты с исправлением этих проблем и обнаружили, что современные функции потерь работают немного лучше, а иногда и наравне с классическими методами. Это резко контрастирует с утверждениями, сделанными в статьях.
- В будущем можно будет изучить взаимосвязь между оптимальными гиперпараметрами и комбинациями наборов данных / архитектуры, а также причины, по которым разные потери работают одинаково.

# Вопросы

- 1. В чем идея triplet loss и в чем он превосходит более простые версии, например, contrastive loss?
- 2. В чем проблема использования обычных метрик в задачах metric learning? Какую метрику предлагают авторы статьи и чем она лучше? (С формулой)
- 3. Как авторы предлагают подбирать гиперпараметры для модели?