

# Поиск эффективной архитектуры для решения задачи получения мультимодальных эмбеддингов для трех и более модальностей

Жаров Георгий

МФТИ

10 марта 2024 г.

## Цель работы

- Сравнить существующие методы получения мультимодальных эмбедингов методом ССА.
- Найти эффективную нейросетевую архитектуру для решения задачи получения мультимодальных эмбедингов для трех и более модальностей.

Даны данные из двух модальностей  $X_1 \in \mathbb{R}^{d_1}$  и  $X_2 \in \mathbb{R}^{d_2}$ , с ковариационными матрицами,  $\Sigma_{11}$  и  $\Sigma_{22}$ , соответственно, и кросс-ковариационной матрицей,  $\Sigma_{12}$

Рассматривается задача поиска линейных преобразований  $(u_1^*, u_2^*)$ , максимизирующих корреляцию между модальностями.

Итоговая задача оптимизации имеет вид:

$$\begin{aligned}(u_1^*, u_2^*) &= \arg \max_{u_1 \in \mathbb{R}^{d_1}, u_2 \in \mathbb{R}^{d_2}} \text{corr}(u_1^\top X_1, u_2^\top X_2) \\ &= \arg \max_{u_1 \in \mathbb{R}^{d_1}, u_2 \in \mathbb{R}^{d_2}} \frac{u_1^\top \Sigma_{12} u_2}{\sqrt{u_1^\top \Sigma_{11} u_1 u_2^\top \Sigma_{22} u_2}}\end{aligned}$$

Иначе

$$(u_1^*, u_2^*) = \arg \max_{u_1^\top \Sigma_{11} u_1 = u_2^\top \Sigma_{22} u_2 = 1} u_1^\top \Sigma_{12} u_2$$

# Deep CCA

Основное ограничение классического CCA — невозможность поиска нелинейного преобразования. Рассмотрим следующее расширение.

Пусть теперь  $f_1(X_1)$  и  $f_2(X_2)$  выходы двух нейросетей с весами  $W_1$  и  $W_2$  соответственно.

Тогда новая оптимизационная задача имеет вид:

$$(u_1^*, u_2^*, W_1^*, W_2^*) = \arg \max_{u_1, u_2} \text{corr}(u_1^\top f_1(X_1), u_2^\top f_2(X_2))$$

Еще одно расширение классического ССА.

Ставится задача нахождения некоторого общего представления  $G$  для  $J$  различных модальностей, где  $N$  — число объектов выборки,  $d_j$  — размерность  $j$ -ой модальности,  $r$  — размерность обучаемого представления, и  $X_j \in \mathbb{R}^{d_j \times N}$  — матрица данных для  $j$ -ой модальности.

Тогда итоговая задача имеет вид:

$$\begin{aligned} & \min_{U_j \in \mathbb{R}^{d_j \times r}, G \in \mathbb{R}^{r \times N}} \sum_{j=1}^J \|G - U_j^\top X_j\|_F^2 \\ & \text{subject to} \quad GG^\top = I_r \end{aligned}$$

## Deep GCCA

Пусть теперь  $f_i$  — сеть для  $i$ -ой модальности данных, а  $o_i$  — размерность ее выходного вектора. Тогда, по аналогии с GCCA можем поставить задачу для DGCCA:

$$\begin{aligned} & U_j \in \mathbb{R}^{o_j \times r}, G \in \mathbb{R}^{r \times N} \sum_{j=1}^J \|G - U_j^\top f_j(X_j)\|_F^2 \\ & \text{subject to} \quad GG^\top = I_r \end{aligned}$$

## Deep GCCA

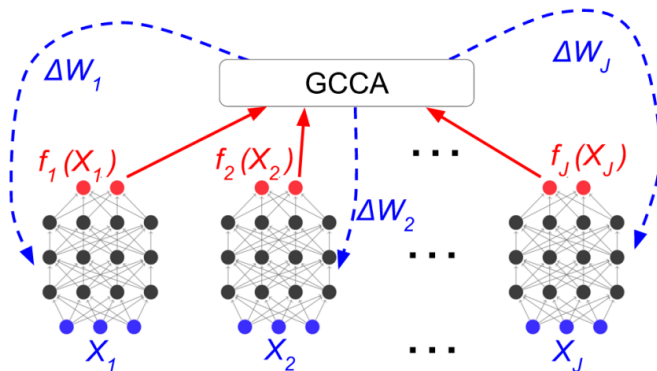
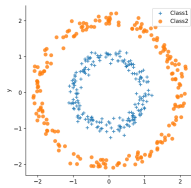


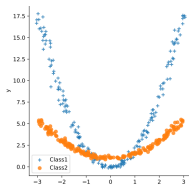
Рис.: Архитектура DGCCA



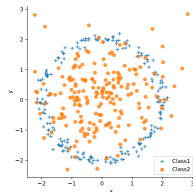
# Применение DGCCA



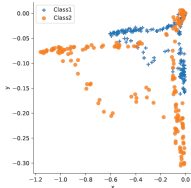
Mod 1 input



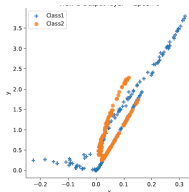
Mod 2 input



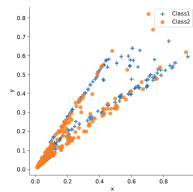
Mod 3 input



Mod 1 output

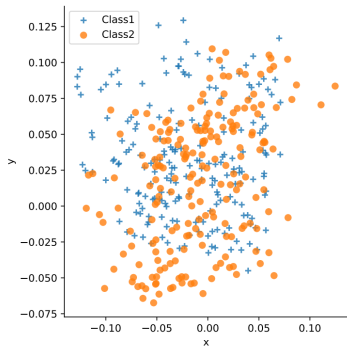


Mod 2 output

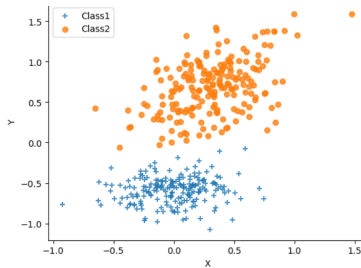


Mod 3 output

# Применение DGCCA



GCCA



DGCCA

Рис.: Итоговые мультимодальные эмбединги для GCCA и DGCCA

# Результаты

- Проведено исследование имеющихся методов получения мультимодальных эмбедингов.
- Проведены вычислительные эксперименты с нейросетевой архитектурой GCCA.
- Написан код для проведения вычислительного эксперимента на синтетических данных.