# METHODS FOR COMPARING REPRESENTATIONS OF ARTIFICIAL NEURAL NETWORKS

MARCO ZULLICH, PHD STUDENT @ DEPT. OF ENGINEERING AND ARCHITECTURE

UNIVERSITY OF TRIESTE

# WHY COMPARING REPRESENTATIONS OF ARTIFICIAL NEURAL NETWORKS?

- The dynamics behind the training process of SGD are still vaguely understood, as are the generalization capability and some properties of hidden representations

- It may be interesting, for instance, to compare layers during training

- On the other hand, given two identically structured ANNs trained on the same dataset with two different random seeds, do they learn similar representations or the representations are different even if the final performance is the same?

- Can we establish connections between the representations learned by ANNs and by their biological counterpart?

# COMPARING LAYERS OF ARTIFICIAL NEURAL NETWORKS (ANNs)

- What defines two ANNs as «similar»?

  - Their parameters (weights)

  - Their outputs

  - Their neurons

  *...other ideas?*

# COMPARING ANNs BY THEIR NEURONS

- What are their values given the input?

- How can we approximate their probability density (given the data manifold) over $\mathbb{R}$?

  - $p(a_{ij}|x)$

- «How does the neuron respond to the data manifold»?

$\rightarrow$ evaluate the network over a limited, yet «large enough» dataset of points

and collect the neurons activations.

«Monte-Carlo approximation»

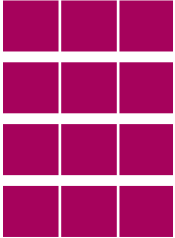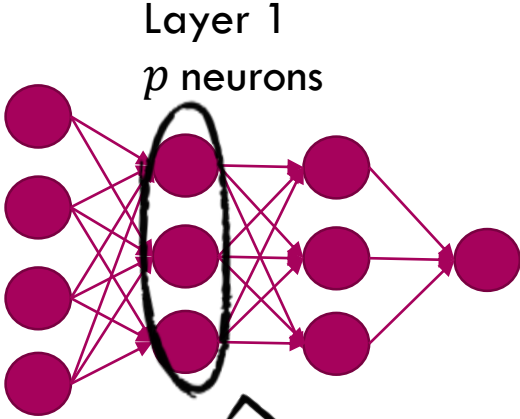# A BIT MORE FORMALLY



Layer 1
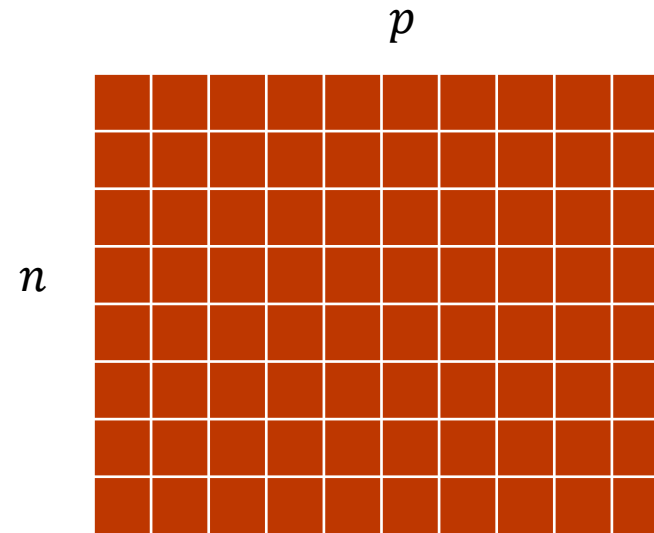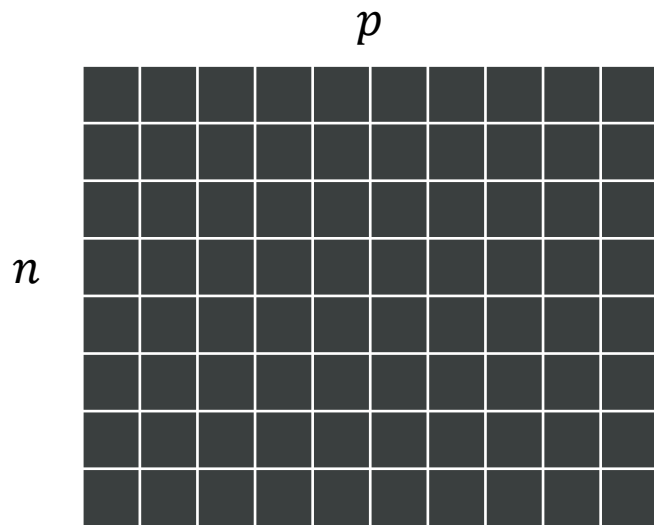$p$ neurons

ANN

Dataset
$n$ images

Activation vectors
size $p$

Activation
matrix $n \times p$

Representation of layer 1

# COMPARING TWO SINGLE-DIMENSIONAL LAYERS



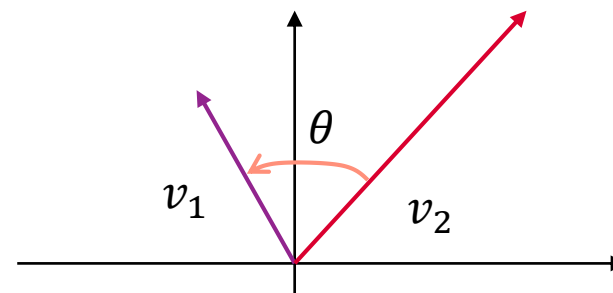What techinques do you know for comparing two matrices of the same size?

# COMPARING MATRICES WITH COSINE SIMILARITY

- The cosine similarity is a **similarity metric** between two vectors lying in the same space

- Is equivalent to the **cosine of the angle** between the two

$$v_1^T v_2 = \|v_1\| \cdot \|v_2\| \cdot \cos(\theta)$$

hence

$$\cos(\theta) = \frac{v_1^T v_2}{\|v_1\| \cdot \|v_2\|} = \frac{\langle v_1, v_2 \rangle}{\|v_1\| \cdot \|v_2\|}$$
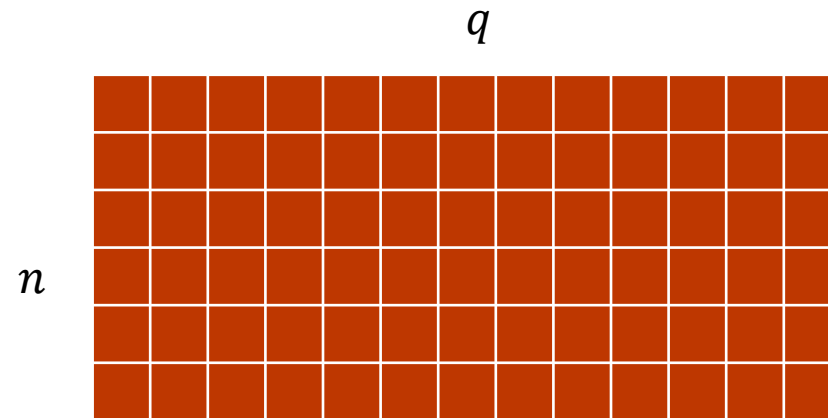
What is wrong with cosine similarity when comparing two flattened layers?

# THE REQUISITE OF ROTATION INVARIANCE

- One of the properties of the neural network is that, by carefully changing the order of the neurons (and their corresponding weights) in the hidden layers, we can obtain two networks behaving identically

- Ideally we would like a similarity metric to be invariant to that property

- If we view the **neurons as directions** within the **space of representations** the permutation can be seen as a **rotation** within this space and the requirement becomes **rotational invariance.**

# …AN ADDITIONAL REQUISITE

- In addition to the rotational invariance requisite, we would also like our metric to be more flexible and adapt to generic situations in which the two layers have different sizes

- In pratical terms…

$p$

$n$

$q$

$n$

$$p \lesseqqgtr q$$

# INTRODUCING CANONICAL CORRELATION ANALYSIS (CCA)

- CCA is a technique introduced in the '30s by the statistician Harold Hotelling

- Is a technique for **establishing connections between two generic sets of continuous variables** called **phenomena**

- $X = (X_1, \dots, X_p); \quad Y = (Y_1, \dots, Y_q), \quad p \gtreqless q$

- We operate by constructing a so-called **view** of these two **phenomena**

  - We consider $n$ statistical units (*individuals*)

  - We evaluate these individuals over our phenomena

# CCA {2}

- Now, we may store these views in two matrices, $M_X, M_Y$

$M_X$ $p$

$n$

$M_Y$ $q$

$n$

# C C A {3}

- CCA acts by applying a **linear transformation** to two **unknown vectors** $w_X \in \mathbb{R}^p, w_y \in \mathbb{R}^q$

- The linear transformation is the one **implied** by $M_X$ and $M_Y$

- $M_X w_X = z_X \in \mathbb{R}^n; \quad M_Y w_Y = z_Y \in \mathbb{R}^n$

The constraint over $w_X, w_y$ is that the corresponding $z_X, z_y$
1) are unit vectors
2) have maximum Pearson correlation.
Geometrically, Pearson correaltion = cosine of enclosing angle $\theta$ (*)



$$z_X = M_X w_X$$
$$z_Y = M_Y w_Y$$

$w_X$    $M_X$    $\theta$

$w_Y$    $M_Y$

Does this ring a bell?

(*) for 0-mean random variables/representations

# CCA {4}

- Call $\rho \triangleq \cos \theta \rightarrow$ CANONICAL CORRELATION (CC)
- We now wish to obtain another set $w_X^{(2)}, w_Y^{(2)}$
- Such that the corresponding $z_X^{(2)}, z_Y^{(2)}$ respect all previous properties
- and they're orthogonal to $z_X, z_y$ respectively
- We will have the corresponding $\rho^{(2)} \leq \rho$
- We can find **an iterative method** which produces a **decreasing sequence of CCs** $\left(\rho^{(1)}, \dots, \rho^{(\min(p,q))}\right)$

# CCA {5}

- We can summarize all this in matrix notation

- $W_X \in \mathbb{R}^{p \times \min(p,q)}, Z_X \in \mathbb{R}^{n \times \min(p,q)}$ (analogous for $W_Y, Z_Y$)

- $M_X W_X = Z_X$ (analogous for $W_Y, Z_Y, M_Y$)

- The rows of $Z_X$ hold the s.c. **canonical variables**

- $Z_X, Z_Y$ are **orthonormal bases of the space** $\mathbb{R}^n$

- Pearson correlation between $z_X^{(i)}, z_Y^{(i)}$ is maximum

- $\mathrm{P} = \left( \rho^{(i)} \right)_i$ is trivially obtained as the row-wise cosine similarity between $Z_X, Z_Y$

# C C A {6}

But how can we obtain $W_X, W_Y$?

| $\Sigma_{XX}$ | $\text{VAR}(M_X)$ |
|---|---|
| $\Sigma_{YY}$ | $\text{VAR}(M_Y)$ |
| $\Sigma_{XY}$ | $COV(M_X, M_Y)$ |

Singular Value Decomposition (SVD)

$$\Sigma_{XX}^{-1/2} \Sigma_{XY} \Sigma_{YY}^{-1/2} = USV$$

Left singular vectors
$\Sigma_{XX}^{1/2} u_i$ yields the corresponding $w_X^{(i)}$

Right singular vectors
$\Sigma_{YY}^{1/2} v_i$ yields the corresponding $w_Y^{(i)}$

Singular values
$s_i$ corresponds to the canonical correlation $\rho^{(i)}$

# MEAN CCA SIMILARITY

| $\boldsymbol{\rho^{(1)}}$ | $\boldsymbol{\rho^{(2)}}$ | $\boldsymbol{\rho^{(3)}}$ | ... | $\boldsymbol{\rho^{(r)}}$ |

$r = \min(p, q)$

These coefficient convey an information on **relatedness** between $M_X$ and $M_Y$.

$$\mathrm{CCA}_{\mathrm{sim}}(M_X, M_Y) \triangleq \frac{\sum_{j=1}^{r} \rho^{(j)}}{r}$$

Mean CCA similarity

# RECAP ON CCA

Study connections between two phenomena $X, Y$ of **possibly different sizes**

Build views by means of $n$ observations over $X, Y \rightarrow M_X, M_Y$

Linearly project the columns of $M_X, M_Y$ in orthonormal bases of size $\mathbb{R}^n \rightarrow Z_X, Z_Y$

The projection maximizes correlation $\rho$ between rows of these two bases

$\rho$'s can be easily obtained *one-shot* via SVD applied on variance-covariance matrices of $M_X, M_Y$

# SVCCA [1]

- A technique thought explicitly for applying CCA to compare ANN layers

- Assumption: *most of the neurons within a layer contribute close to nothing to the variance of its representation*



Operate SVD for dimensionality reduction

$$M_X = U \; S \; V^T$$

$$s_{jj}^2 = \text{VAR}(M_X)_j$$
$$s_{jj} \geq s_{kk}, j > k$$

Operate CCA and obtain similarity

# SIMILARITY OF REPRESENTATIONS USING KERNELS

Linear kernel over a representation $M_X$

To *measure* the dissimilarity between these two matrices, we can use the inner product of their vectorized version

Note: $\text{COV}(A,B) = \frac{AB^T}{n-1} \Rightarrow \|\text{COV}(A,B)\|^2 = \frac{\|AB^T\|^2}{(n-1)^2}$

$$L_X = M_X M_X^T \in \mathbb{R}^{n \times n}$$

$$L_Y = M_Y M_Y^T \in \mathbb{R}^{n \times n}$$

$$\langle \text{vec}(L_X), \text{vec}(L_Y) \rangle = \text{tr}(L_X L_Y)$$

$$= \left\| M_Y^T M_X \right\|_{\text{Fr}}^2$$

$$= (n-1)^2 \left\| \text{COV}(M_X^T, M_Y^T) \right\|^2$$

$$\left\| \text{COV}(M_X^T, M_Y^T) \right\|^2 = \frac{\text{tr}(L_X L_Y)}{(n-1)^2}$$

# HILBERT SCHMIDT INDEPENDENCE CRITERION (HSIC)

$$\left\|\mathrm{COV}(M_X^T, M_Y^T)\right\|^2 = \frac{\mathrm{tr}(L_X L_Y)}{(n-1)^2}$$

HSIC for linear kernels

$$\mathrm{HSIC}\left(\widetilde{K_X}, \widetilde{K_Y}\right) = \frac{\mathrm{tr}(\widetilde{K_X}\widetilde{K_Y})}{(n-1)^2}$$

$\widetilde{K_X} = K_X(I - n^{-1}\mathbf{1}\mathbf{1}^T)$ kernels centered w.r.t. row and column means

HSIC is a statistic for determining whether two generic sets of variables are independent

HSIC $\to$ 0 stochastic independence
HSIC $\to$ 1 stochastic dependence

# CENTERED KERNEL ALIGNMENT (CKA)

- Normalization of HSIC

$$\text{CKA}(\widetilde{K_X}, \widetilde{K_Y}) = \frac{\text{HSIC}(\widetilde{K_X}, \widetilde{K_Y})}{\sqrt{\text{HSIC}(\widetilde{K_X}, \widetilde{K_X})\text{HSIC}(\widetilde{K_Y}, \widetilde{K_Y})}} = \frac{\langle \widetilde{K_X}, \widetilde{K_Y} \rangle}{\|\widetilde{K_X}\|\|\widetilde{K_Y}\|} \in [0,1]$$

- Linear CKA

$$\text{CKA}_{\text{lin}}(M_X, M_Y) = \frac{\|M_Y^T M_X\|^2}{\|M_X^T M_X\|\|M_Y^T M_Y\|} \in [0,1]$$

# CHOICE OF KERNELS

- **Radial Basis Function** (RBF) kernel:

- $\kappa(x, y) = \exp\left(\frac{-\|x-y\|^2}{2\sigma^2}\right)$

- [2] cites no substantial difference between using RBF kernel with $\sigma \in [0.2, 0.6]$ w.r.t. a linear kernel

- On the other hand, [3] refers that CKA with RBF with *very small* sigma may be a better choice for a more accurate similarity metric, but more on that on the next keynote.

# ON THE *DESIRABLE* AND *INDESIRABLE* INVARIANCES OF METRICS

- Invariance to orthogonal transformations was already discussed before

  Good

- **Invariance to isotropic scaling** (arbitrary scaling of the features)

  - $\text{SIM}(M_X, M_Y) = \text{SIM}(\alpha M_X, \beta M_y), \quad \alpha, \beta \in \mathbb{R}^+$

    Good

**Invariance to invertible linear transformations**   Bad

This invariance poses problems when $n < p$, as for full-rank matrices $A, B$, the similarity $SIM(A, C) = SIM(B, C)$ [2]

The training process is sensible w.r.t. invertible linear transforms. Just think of BATCH NORMALIZATION
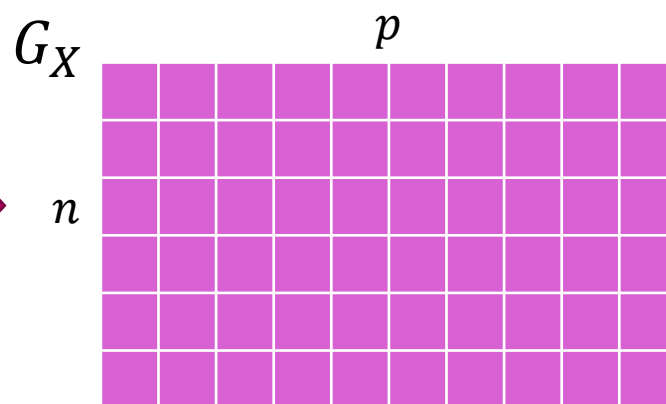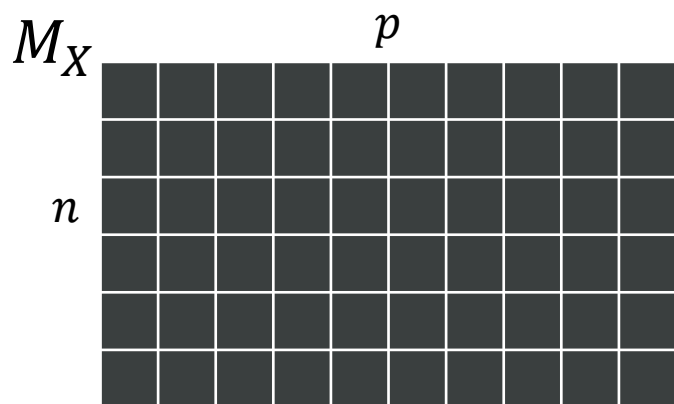
# SUMMARY OF METRICS INVARIANCES

| Similarity Index | Formula | Invariant to | | |
|---|---|---|---|---|
| | | **Invertible Linear Transform** | **Orthogonal Transform** | **Isotropic Scaling** |
| Linear Reg. ($R^2_{\mathrm{LR}}$) | $\|Q_Y^{\mathrm{T}} X\|_{\mathrm{F}}^2 / \|X\|_{\mathrm{F}}^2$ | $Y$ only | ✓ | ✓ |
| CCA ($R^2_{\mathrm{CCA}}$) | $\|Q_Y^{\mathrm{T}} Q_X\|_{\mathrm{F}}^2 / p_1$ | ✓ | ✓ | ✓ |
| CCA ($\bar{\rho}_{\mathrm{CCA}}$) | $\|Q_Y^{\mathrm{T}} Q_X\|_* / p_1$ | ✓ | ✓ | ✓ |
| SVCCA ($R^2_{\mathrm{SVCCA}}$) | $\|(U_Y T_Y)^{\mathrm{T}} U_X T_X\|_{\mathrm{F}}^2 / \min(\|T_X\|_{\mathrm{F}}^2, \|T_Y\|_{\mathrm{F}}^2)$ | If same subspace kept | ✓ | ✓ |
| SVCCA ($\bar{\rho}_{\mathrm{SVCCA}}$) | $\|(U_Y T_Y)^{\mathrm{T}} U_X T_X\|_* / \min(\|T_X\|_{\mathrm{F}}^2, \|T_Y\|_{\mathrm{F}}^2)$ | If same subspace kept | ✓ | ✓ |
| PWCCA | $\sum_{i=1}^{p_1} \alpha_i \rho_i / \|\alpha\|_1, \alpha_i = \sum_j |\langle \mathbf{h}_i, \mathbf{x}_j \rangle|$ | ✗ | ✗ | ✓ |
| Linear HSIC | $\|Y^{\mathrm{T}} X\|_{\mathrm{F}}^2 / (n-1)^2$ | ✗ | ✓ | ✗ |
| Linear CKA | $\|Y^{\mathrm{T}} X\|_{\mathrm{F}}^2 / (\|X^{\mathrm{T}} X\|_{\mathrm{F}} \|Y^{\mathrm{T}} Y\|_{\mathrm{F}})$ | ✗ | ✓ | ✓ |
| RBF CKA | $\mathrm{tr}(KHLH) / \sqrt{\mathrm{tr}(KHKH)\mathrm{tr}(LHLH)}$ | ✗ | ✓ | ✓* |

*Table from [2].*

# AUGMENTING CKA WITH INFORMATION ON GRADIENTS

[4] proposes an augmentation of CKA by incorporating information on gradients for the given layer(s)



$$K_X = K_{M_X} \odot K_{G_X}$$

$$K_Y = K_{M_Y} \odot K_{G_Y}$$

$$\text{CKA}_{Gr.}\left(\widetilde{K_X}, \widetilde{K_Y}\right) = \frac{\left\langle \widetilde{K_X}, \widetilde{K_Y} \right\rangle}{\left\| \widetilde{K_X} \right\| \left\| \widetilde{K_Y} \right\|}$$

# WHY CKA AND CCA MIGHT BE WRONG

**Classical statistics:**
FIXED features (*variables*)
VARIABLE datapoints (*observations*)
$p$ fixed; $n \to \infty$

**High dimensional statistics [7]:**
VARIABLE features (*variables*)
VARIABLE datapoints (*observations*)
$p \to \infty; n \to \infty$

In Deep Learning, the focus is either on depth and width. In wide ANNs, we're essentially increasing $p$ which may be way larger than $n$.
Consider moreover that both CCA and CKA are essentially based on the concept of COVARIANCE.
Maybe, we might want to detach from a classical statistical view and go to more "*non-parametric*" techniques to obtain metrics.

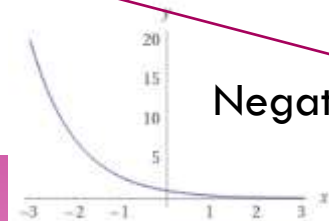Most of the results we know and use everyday in research adhere to this paradigm.
**e.g. law of large numbers**

In high dimensional statistics, some results from classical statistics do not hold or may present abnormally large errors

$$P\left(\left\|\Sigma - \hat{\Sigma}\right\| \geq \delta\right) \leq 2p \exp\left(\frac{-n\delta^2}{2b(\|\Sigma\| + \delta)}\right), b \in \mathbb{R}, \delta \to 0^+$$

What is this?

Linear increase

Negative exponential decay

# IMD

CRAZY STUFF!

- IMD is a metric recently proposed at ICLR 2020 [6]

- Compares generic data manifolds, unaligned and different in dimension

- Underlying theory is absurdly difficult

- If you really want to have a go at it

  - https://github.com/xgfs/imd

- Focus is on generative models and language models, but it should work fine on simple MLPs as well

# RECAPPING

**Comparing hidden representations produced by MLPs is a difficult task**

*Unaligned representations:* Neuron $i$ in ANN 1 might not be the same as neuron $i$ in ANN2

*Different dimensionalities:* Representations may not be composed by the same number of neurons

*Metrics such as CKA and (SV)CCA overcome these hurdles*

*Curse of dimensionality:* Human intuition fails when the number of variables (*neurons*) in the representation is very high

*High dimensional statistics:* Regular statistical results may fail when $p \rightarrow inf$. It is shown that covariance is tricky in that scenario

# REFERENCES

1. Raghu, Maithra, et al. "Svcca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability." *NeurIPS, 2018*.

2. Kornblith, Simon, et al. "Similarity of neural network representations revisited." *International Conference on Machine Learning*. PMLR, 2019.

3. Doimo, Diego, et al. "Hierarchical nucleation in deep neural networks." *NeurIPS, 2020*.

4. Tang, Shuai, et al. "Similarity of neural networks with gradients." *arXiv preprint arXiv:2003.11498* (2020).

5. Zullich, Marco, et al. "Investigating Similarity Metrics for Convolutional Neural Networks in the Case of Unstructured Pruning". ICPRAM, 2020.

6. Tsitsulin, Anton, et al. "The Shape of Data: Intrinsic Distance for Data Distributions." ICLR, 2020.

7. Wainwright, Martin J. *High-dimensional statistics: A non-asymptotic viewpoint*. Vol. 48. Cambridge University Press, 2019.