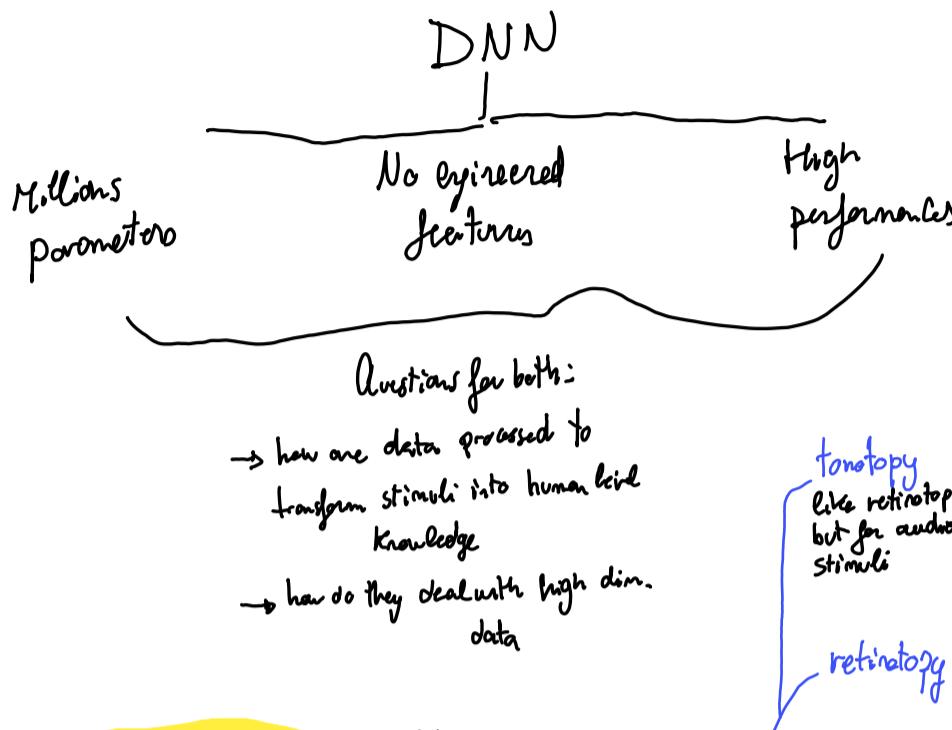


OVERVIEW OF ML APPROACHES

- BARRETT et al. → synergy between NN and artificial NN (DNN)



ANALOGIES

- RECEPTIVE FIELDS → relations between some areas and features processed there
- small receptive fields in first layers of visual cortex
 - receptive fields grows as information is transmitted to higher level areas of visual cortex
 - × this enables sensitivity to larger areas of space
 - × more complex features and abstraction
- Does this hold also for DNN?
• patterns
• pooling
- ↑ "concept cells"
- ABALATION: → studying brain lesions offer much information about potential function of brain areas (≈ REORGANIZATION)
- remove neuron with its outgoing weights
- we can always consider effect of deleting neurons on the output (structural pruning)
• NN which have good generalization are more robust to ablation
- "can be achieved with various tools" [?]
- DIMENSIONALITY REDUCTION:
- brain store infos. in a distributed way → redundancy
 - distributed manner dependencies
 - correlations of units means redundancy
- In DNN we have a lot of redundancy and correlation
- REPRESENTATIONAL GEOMETRIES
- how does a DNN represent informations across different layers in space and time
 - comparison of different states using canonical correlation analysis, identify few latent factors that maximise correlation
- PLS correlation
- SUPERVISED!
-
- feature
- object
- object
- similarity metrics
- Linear regression methods can predict neurobiological activation

SPATIAL METHODS: Spicer & Sanborn

place items on a multidimensional space and use their location to draw conclusions about categorization

→ CLASSIFICATION: locate items with hyperplane

↳ ASSUMPTION: there exist a statistical model with a center and some parameters that are learnable in training

Ex: GENERALIZED CONTEXT MODEL

$$P(C_j | i) = \frac{(\sum_{j \in i} S_{ij})^{\gamma}}{\sum_k (\sum_{k \in K} S_{ik})^{\gamma}} \rightarrow \text{No Softkey}$$

↑
prob. that i belongs
to G ↑
similarity of
 i to j ↑
deterministic
regulator
in simplicity = 1

LOGICAL METHODS

concepts are based on definition that is applied

the features of the object. → SEARCH FOR RULES } also probabilistic

THAT MAXIMIZE DISCRIMINATION BETWEEN STIMULI

↑
eg 1st order logic

ANNs: (artificial NN)

do not make assumptions about the representation involved, but offer an implementation method

Which model is the ~~best~~ a ~~nicer~~? Depends on the area
e.g.
of science you are working on → do you need underlying representations?

CATEGORICAL PERCEPTION

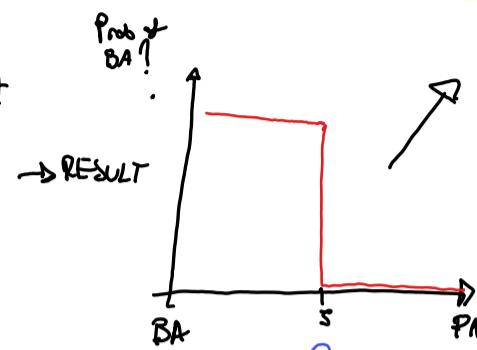
The phenomenon in which people perceive stimuli from different categories as more different from each other than stimuli from within the same category example: sound

EXAMPLE

stimuli
100 - 8000 Hz



VOT: voice onset time of consonant
similarity
judgements/
generalization/
confusion



we are very good to distinguish
4 to 6 but we are
unable to distinguish
2 to 2

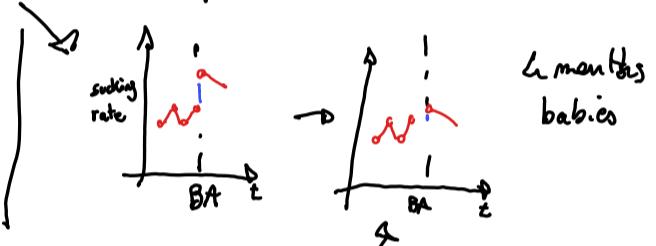
CATEGORICAL PERCEPTION

If there is no category boundary, then
the two stimuli are the same for us

} we are fast to learn that category

EVIDENCES

sordida
CHANGE DEAFNESS: participants were repeated words in a stream - Halfway the voice changed but just 40% of participants noticed it



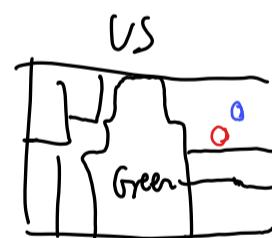
habituation

LINGUISTIC CATEGORICAL PERCEPTION (Colors)

more easy to distinguish blue than green than shades even if they have same distance

↓
are discriminative better when we introduce categories

Borinno tribes distinguish better some colors because they give a particular name to some colors which we usually consider just shades



↓
better discrimination

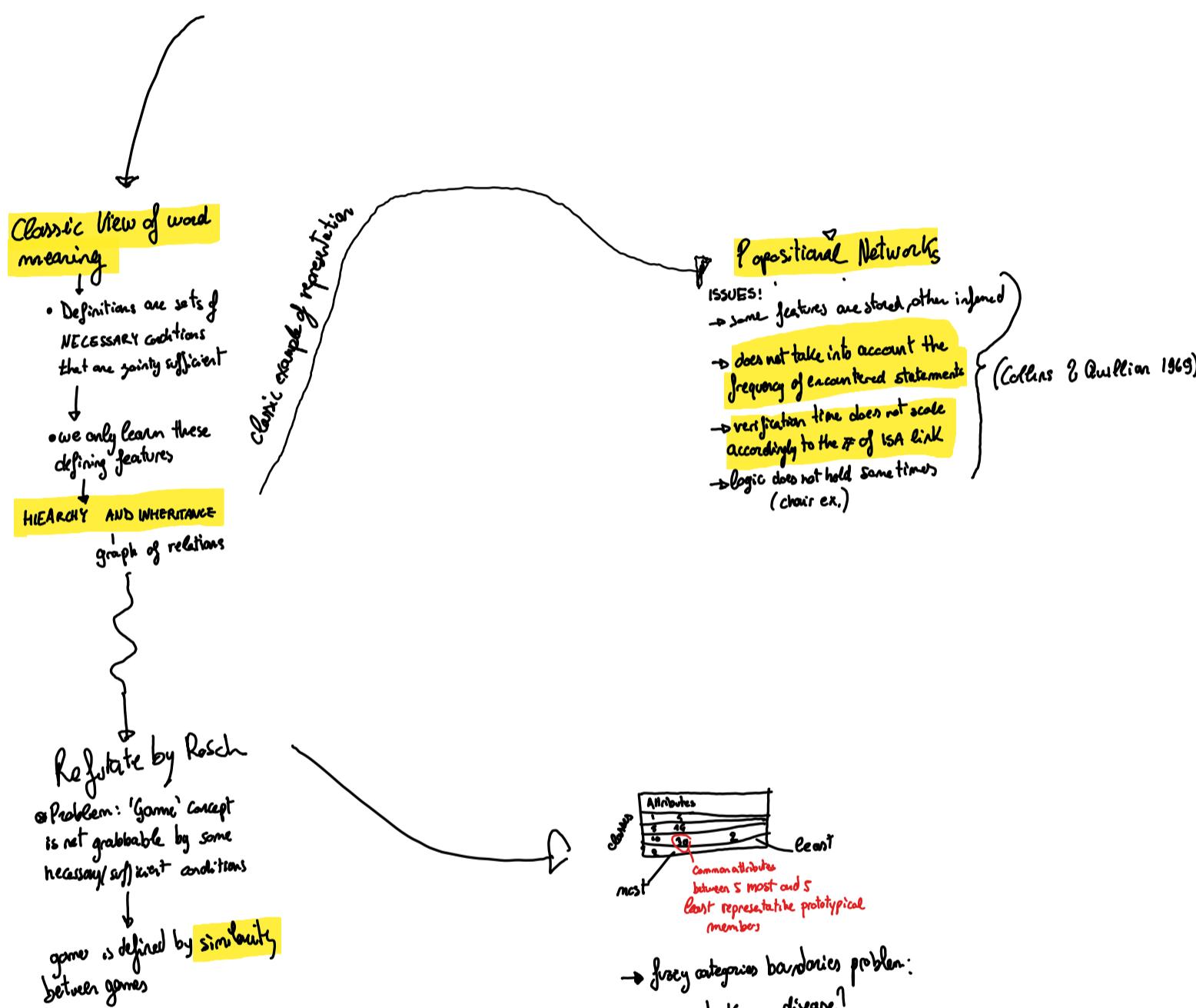
CONCEPTUAL STRUCTURE

- Typical approach → assume that words are associated with concepts representations

↓
 → Problem of polysemy → we need context
 (Murphy) → Klein and Murphy: we can store multiple meaning
 for the same word

→ Problem of homophones → need context
 • homonyms → need context
 (bank, bank)

* anomia → a form of aphasia, difficulty
 to retrieve lexical form
 of a concept



Typicality (?)

Categories have central and peripheral members

Attributes	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100
classes	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100
most	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100
least	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100

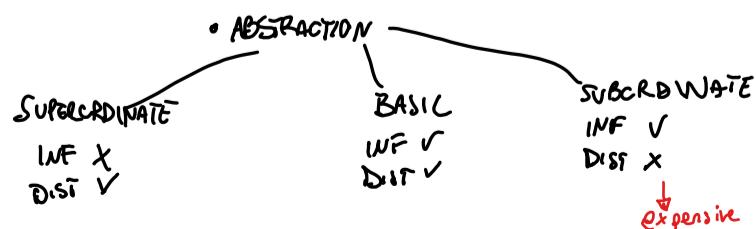
→ fuzzy categories boundaries problem:

is stroke a disease?

→ insecurity, fast changes in opinions

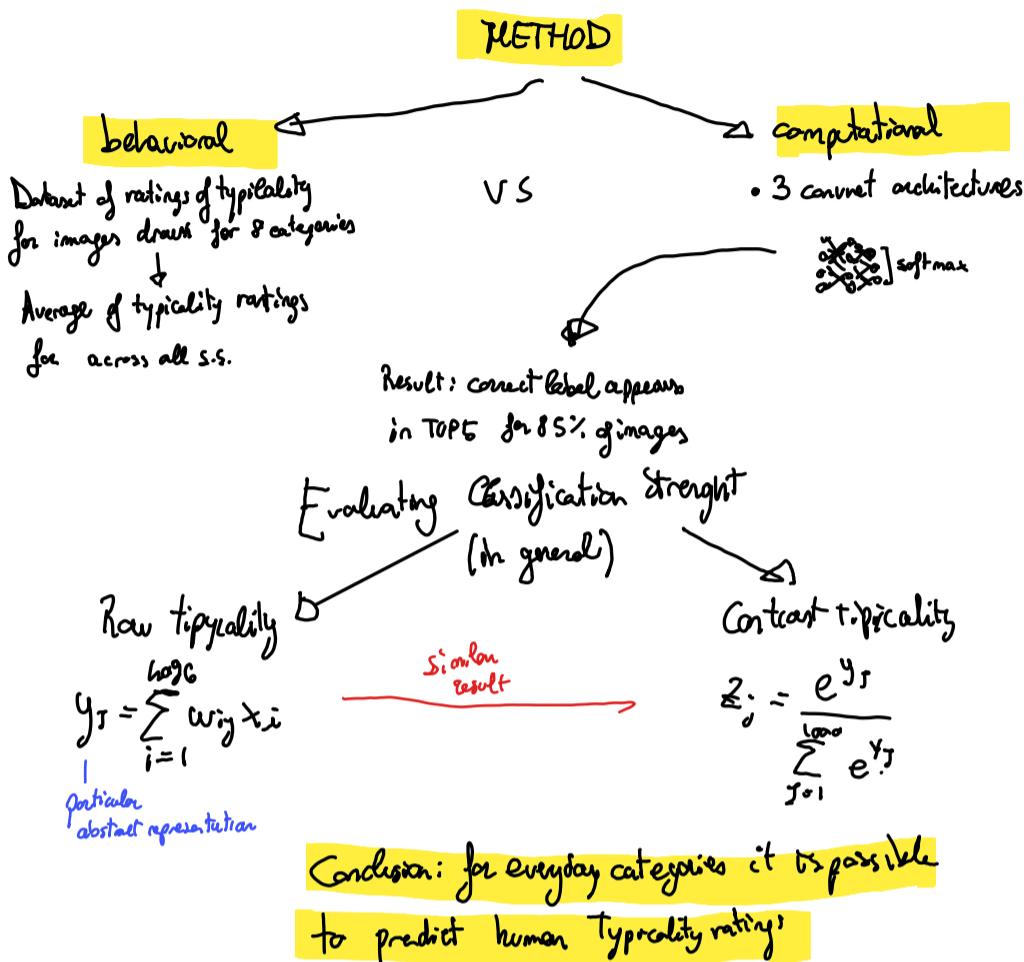
Why do we choose a particular word? (cat)

- level of information (why adding less information?)
- distinctiveness (what makes it different)

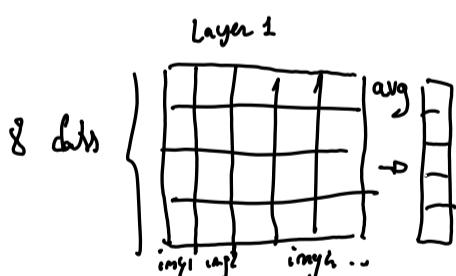


MODELING CONCEPTUAL ORGANIZATION (Lake et al.)

→ Notion of typicality → Can DL systems imitate cognitive models?
 ↳ changes categorization speed

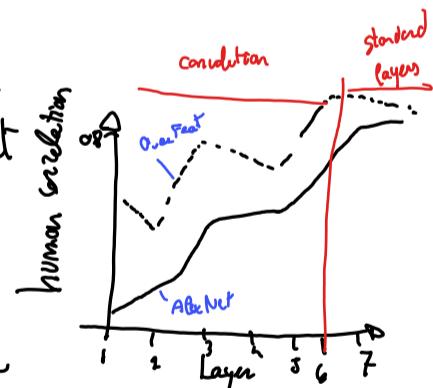


→ More in depth: examination within layers: → pass dataset to NN
 → get activation vectors



→ average by rows over all activation matrices to get typicality response of a class

→ compute it with single image response



Words features as models of

Mitchell et. al 2008

Cognition

- How do humans represent concepts in the brain?
- Which features (brain functions)

Previously: regression on fMRI on words sampled from different categories
 → predict activation for specific word

→ Problem of high # of features using encoder (they use just the 25 already known to be significant)

fMRI $\xrightarrow{\text{predict}}$ word

$\xrightarrow{\text{pure decoder}}$

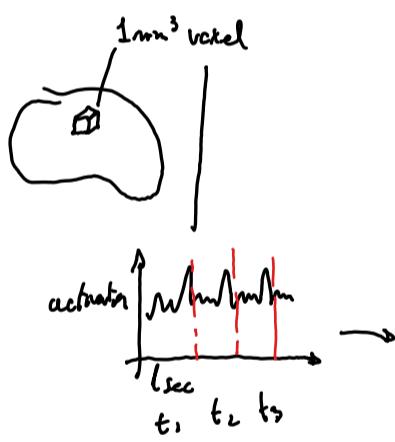
GENERATIVE ENCODER MODEL

- ① encode words into semantic vectors

that should capture meaning.

$$\begin{matrix} \text{- useful} \\ \text{- grey} \\ \text{- pet} \\ \dots \\ \vdots \end{matrix} \xrightarrow{\text{25 features}} \begin{bmatrix} 1 \\ 0 \\ 0 \\ \dots \\ 1 \end{bmatrix}$$

extracted from corpus counting co-occurrences
 25 verbs
 celery: eat $\rightarrow 0.83$
 taste $\rightarrow 0.3661$



- ② single voxel analysis: collect activation for each voxel for each semantic feature

collect fMRI data		
t_1	t_2	t_3
x_1	x_2	x_3
y_1	y_2	y_3
z_1	z_2	z_3
1	\bar{x}	1
:	58	:

\rightarrow f voxel

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_{58} \end{bmatrix}, E = \begin{bmatrix} e_1 \\ \vdots \\ e_{58} \end{bmatrix}, \beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_{25} \end{bmatrix}, X = \begin{bmatrix} 1 & x_{11} & \dots & x_{125} \\ \vdots & \vdots & & \vdots \\ 1 & x_{51} & \dots & x_{525} \end{bmatrix}$$

semantic features

activation of a single voxel for each of the 58 words

celery

sector

$$\begin{bmatrix} 1 & 0.8368 & 0.3661 & \dots \end{bmatrix}$$

from Corpus Statistics

③ predict neural fMRI activation of every voxel location as a weighted sum of neural activations contributed by each of the semantic features (58/60)

We want to predict 60th word using the other 58

↓
Just multiply the semantic features vector of that word by betas and add bias

$$\text{bias} + 0.21 \begin{bmatrix} \text{eat} \\ 0 \end{bmatrix} + 0.35 \begin{bmatrix} \text{toe} \\ 0 \end{bmatrix}$$

= 

$$\text{prediction} = \sum_{i=1}^{35} f_i(w) c_i$$

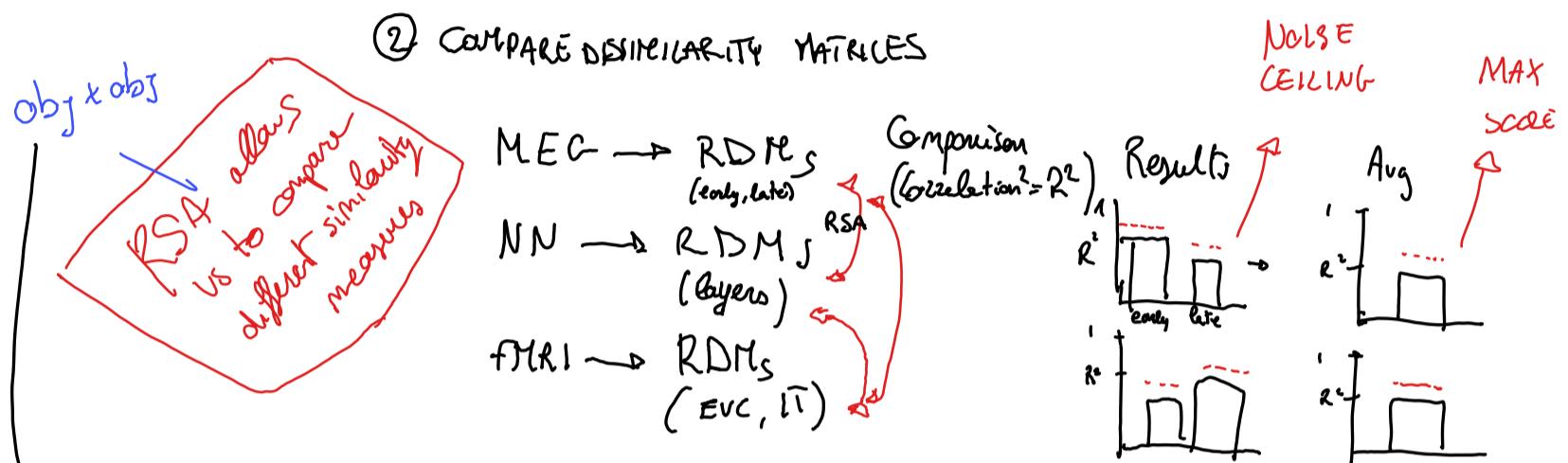
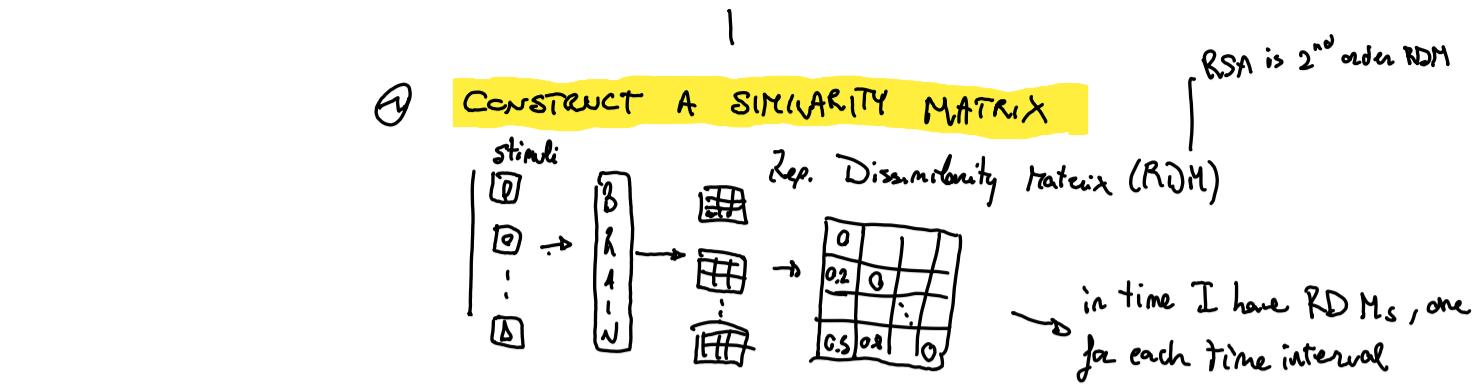
↳ test it for each word using leave-one-out t-test

↓
Mean accuracy 0.79

④ Examine importance: → activation of areas of the brain that are known to be related to some semantics (eg. food) → most activating words are in those semantics
 → random basis set is crucial → random sets always do worst (0.6 vs 0.79)

STUDYING BRAIN REPRESENTATION VIA SIMILARITY SPACES

Kriegeskorte 2008
RSA



obviously this lead to a loss of dimension if # features > # obj

$$\text{e.g. } \text{obj} \times \text{obj} \rightarrow \text{voxels} \times \text{obj}$$

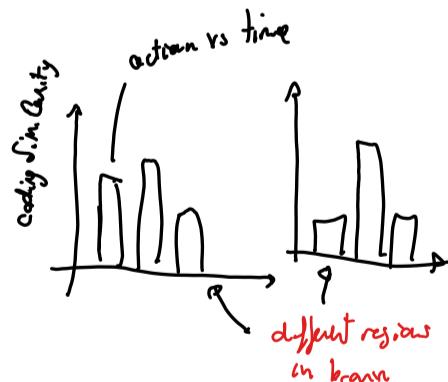
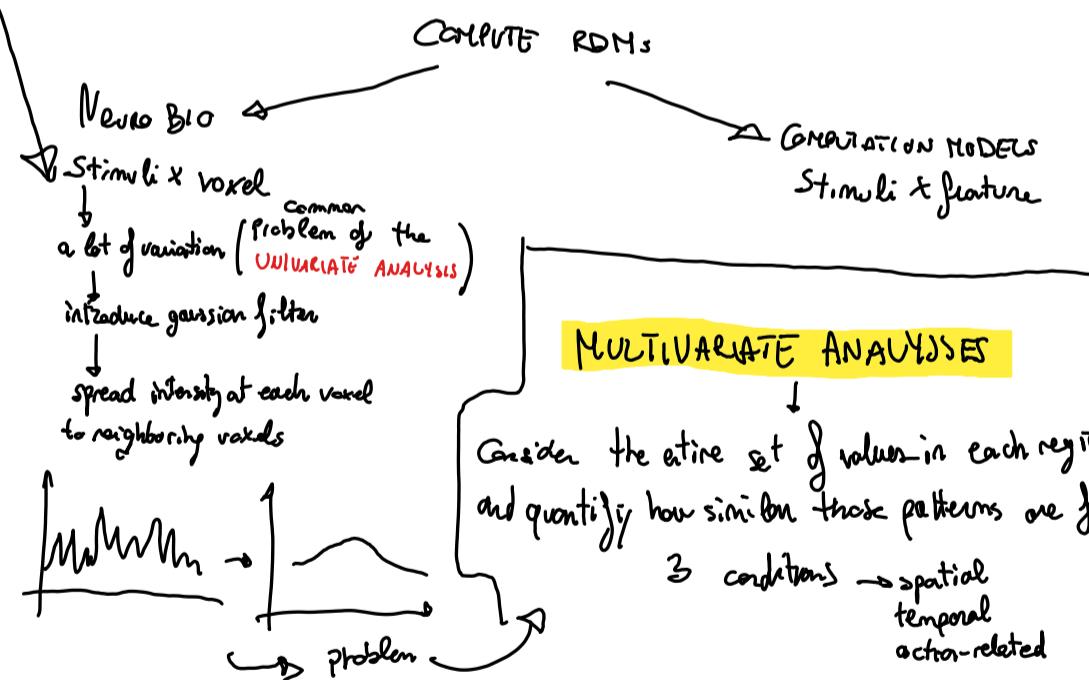
Noise Ceiling \rightarrow maximum value given the noise \rightarrow comes from stochasticity of human behavior

\hookrightarrow can be measured repeating some experiment and computing correlation

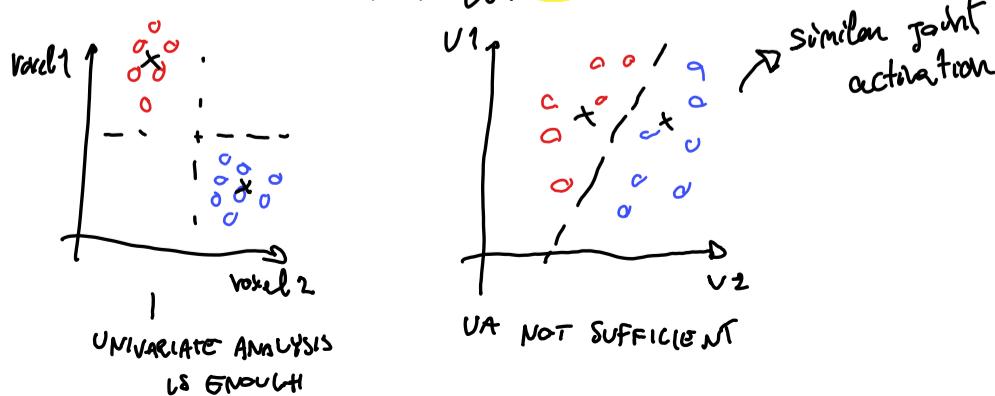
Why IS THIS RELEVANT?

\hookrightarrow We can check if different areas react to the same stimuli.

test hypothesis of area where informations are processed
compare goodness of same brain areas for same tasks.



MULTI-VOXEL PATTERN ANALYSIS



EVALUATING AND IMPROVING

Joshua C. Peterson

THE CORRESPONDENCE BETWEEN DNN → RSA

AND HUMAN REPRESENTATIONS → similarity judgements

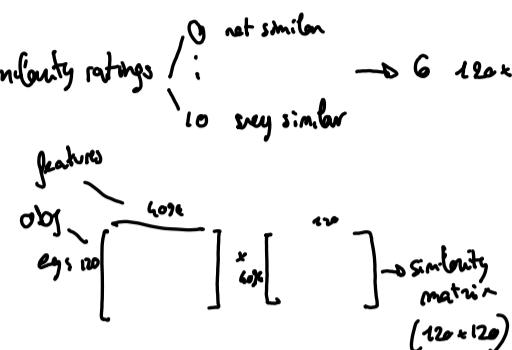
DATASET: 6 categories, 120 images each

for HUMAN → Amazon Mechanical Turk → pairwise similarity ratings → 6 120×120 matrices

$$S = F \cdot F^T$$

human
Similarity
RDM
(0,1)

feature matrix
(dot product)



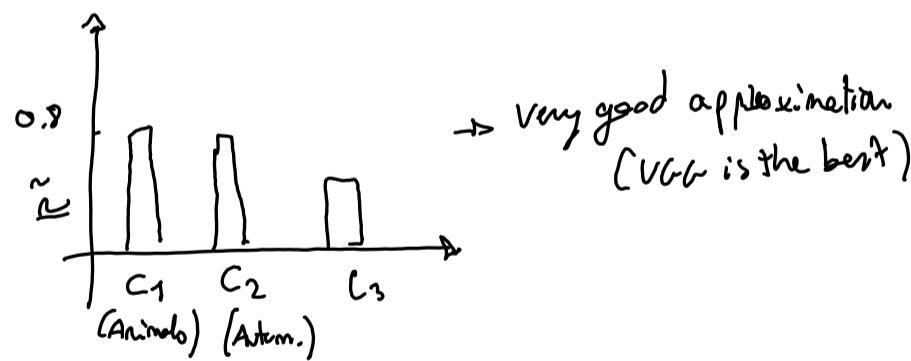
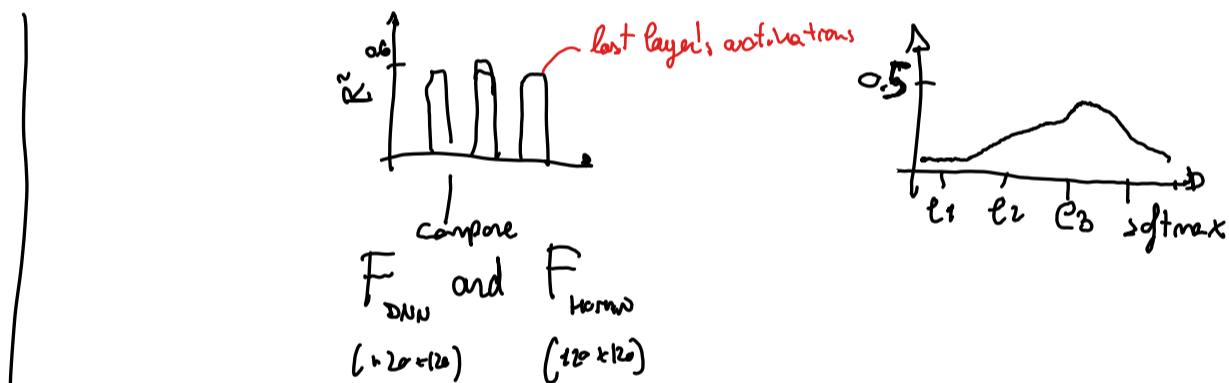
for DNN → consider activation values

↳ used

$$\begin{matrix} & & & \text{neurons in layer} \\ 1 & \left[\begin{matrix} a_{11}^1 & \dots & a_{1n}^1 \\ \vdots & & \vdots \\ m & \left[\begin{matrix} a_{m1}^1 & \dots & a_{mn}^1 \end{matrix} \right] \end{matrix} \right] \end{matrix}$$

Multidimensional
feature representation → matrix of activations
for each image

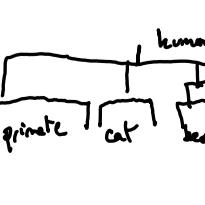
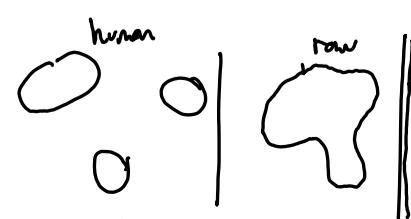
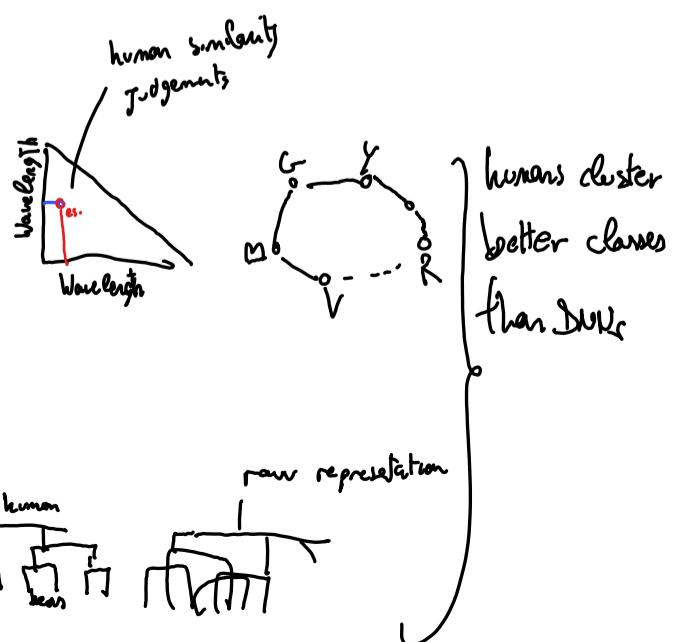
COMPARISON → Compute Pearson correlation between human and DNN's similarity matrix



OTHER METHODS:

Used for better
understanding how
DNNs succeed and
fail in representing
psychological
concepts

① **NON-METRIC MULTIDIMENSIONAL SCALING:**
(convert similarity into spatial representation)



→ Try to improve alignment of DNNs representations

$$S = FWF^T$$

$$\left| \begin{array}{c} \\ \\ \end{array} \right| \left[\begin{array}{c} a \\ B_y \\ C \\ 0 \\ \vdots \\ \epsilon \end{array} \right]$$

$$S_{ij} = \sum_k w_k f_{ik} f_{jk}$$

two images feature

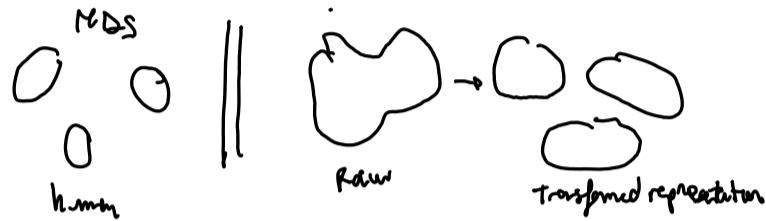
$$Y_i = \beta_i X$$

DNN's learn the correct basis set of features but with incorrect scale of saliency → reweighting improves a lot accuracy

OVERFITTING PROBLEM

$\frac{120+160}{2} \approx 7000$ samples against 60k parameters
 ↓
 → Feature selection
 → regularized L2
 $(w^T w)$

→ GOOD RESULTS → but we are weighting dot product \Rightarrow NOT LINEAR MAPPING!



Conclusion: reweighting standard AI approaches is sufficient to obtain a good approximation



They then asked themselves if they can use the same W for different categories

IMPROVING ISOMORPHISM

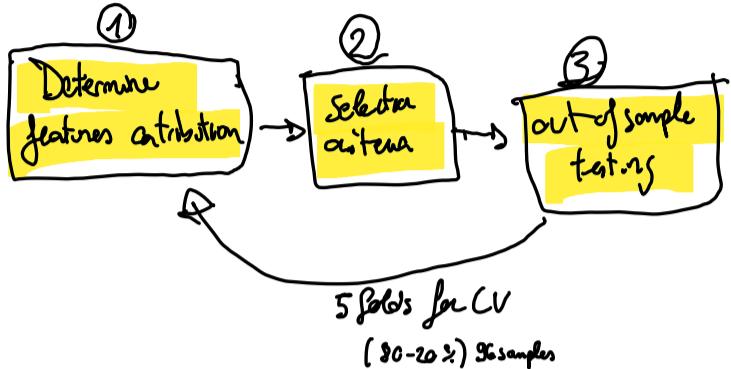
(Pruning vs Reneging)

IDFA: DNNs do acquire the correct saliency of the but features but they're diluted by irrelevant ones

[es: considering every word for the task]

HYPOTHESES: DNN's features embeddings approximate HSS better than 2nd order isomorphisms (Pearson Correlation)

- feature pruning by ranking can improve the fit between DNNs and human



① Done by removing each feature once at the time and retrain \rightarrow 4696 features \rightarrow RANKING
Comparing score with a baseline

② Add one feature at the time from scratch based on ranking and compute RDM

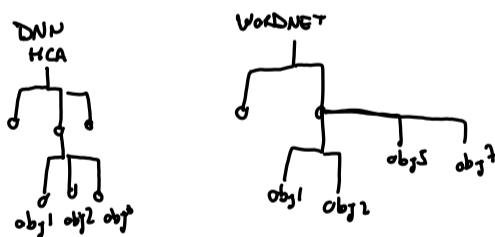
③ Prune the DNN based on features \rightarrow check activation in pruned layers

Check with different architecture \rightarrow PRUNED outperforms the other architectures (Lasso, Sim-DR, PAC-18...) \downarrow
 Problems

\rightarrow Pruning with DNN reflects ranking of pairwise similarities
but does not directly address the hierarchical structure

\downarrow
sol: Compare the H. structure latent in DNNs to WordNet

- ① Produce clusters from DNNs.
- ② Define a Neighbors for each leaf
- ③ Jaccard distance (fitness of neighbors)



$$IoU = \frac{A \cap B}{A \cup B} = \frac{(1,2)}{(1,2,3,5,7)}$$

how the net is capable of learning
better hierarchy

Another test: Prune from brain similarity rather than human similarity

Dataset: 122 images with brain activity for each area (fMRI) \rightarrow combined in 8 brain areas

Result \rightarrow Pruned DNN performs a lot better than raw

Why?

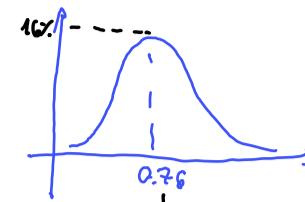
Hu et. al

- Overfitting \rightarrow STUDIED WITH THE NOTION OF APOZ

APoZ: Average percentage of zeros of a single neuron \rightarrow % of 0 activations of that neuron after the ReLU mapping

$$\frac{\sum_k^N \sum_j^M f(O_{c,j}^{(i)}(k) = 0)}{N \cdot M}$$

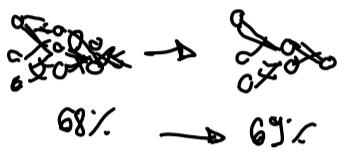
samples features



Consider the entire curve!!

ABORTION
↓
test
↓
fine tuning
↓
test

find the useless nodes and prune them!



\rightarrow then retrain (fine tuning)

\rightarrow They then add one feature at the times according to their APoZ

H to L \rightarrow very good R²

L to H \rightarrow good R² \rightarrow This means that informations are redundant, so some neurons

fire with low freq. bc others do the same task

\rightarrow Then they try to remove H to L

Prune is good

even to approximate human SJs

