

INTERPRETABLE SEMANTIC VECTORS FROM A JOINT

MODEL OF BRAIN-AND TEXT-BASED MEANING

Fish et al



- VSM: vector space model → represent semantical meaning by assigning each word a point in high dimensional space. Ex vectors from corpus.
Distance between semantics maps into vectors' distance

IDEA: use it with brain activations to improve semantic models

DATASET: → human-annotated word semantics
→ brain-informed VSM (fMRI & MEG for 18 subjects → read 60 raves across 12 categories)

PREDICTION → infuse brain data to increase performance

MAP → map semantic concepts into the brain by jointly learning neural representation

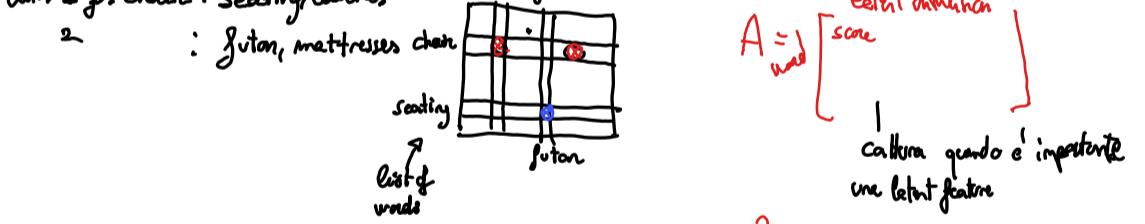
CORPUS → 16 billions words corpus data

METHODS:

$$\text{NNSE} \text{ (NON NEGATIVE SPARSE EMBEDDING)} \rightarrow \underset{A, D}{\operatorname{argmin}} \sum_{i=1}^w \|X_i - A_i x^D\|_F^2 + \lambda \|A\|_1$$

→ FIND BEST SCORING DIMENSION FOR EACH WORD

Eg: dim 1 for chair: seating, couches
2: futon, mattress, chair



JNSE → Adding brain activity

$$\underset{A, D, x^D}{\operatorname{argmin}} \sum_{i=1}^w \|X_i - A_i x^D\|_F^2$$

$$+ \sum_{i=1}^w \|Y_i - A_i x^D\|_F^2 + \lambda \|A\|_1$$

activity per word

partial paired data

mappa latent dimension nelle parole

PROS: handling partial paired data
merges different brain imaging

- give constraints in A and D
- get best solution possible
- we can specify
- choice of latent dimension

$$A = \begin{bmatrix} z \\ \vdots \\ z \end{bmatrix}$$

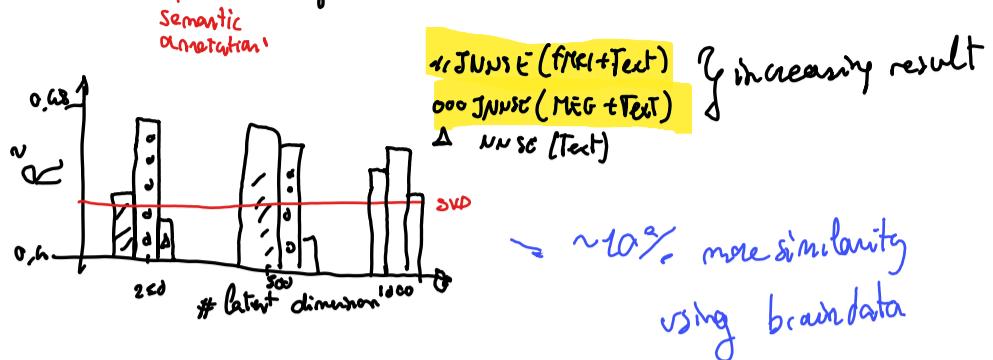
$$D = \begin{bmatrix} \text{Score} \\ \vdots \\ \text{Score} \end{bmatrix}$$

they want pre-defined # of dimensions

STEPS:

- ① try to predict behavioral measure using $U S M_{SA}$
- ② try to predict the 60×218 matrix from $60 \times A$ using $\delta N N S E, N N S E$

- matrix value must be sparse (for interpretability)
- we can exploit $D^{(b)}$ and $D^{(c)}$ to predict

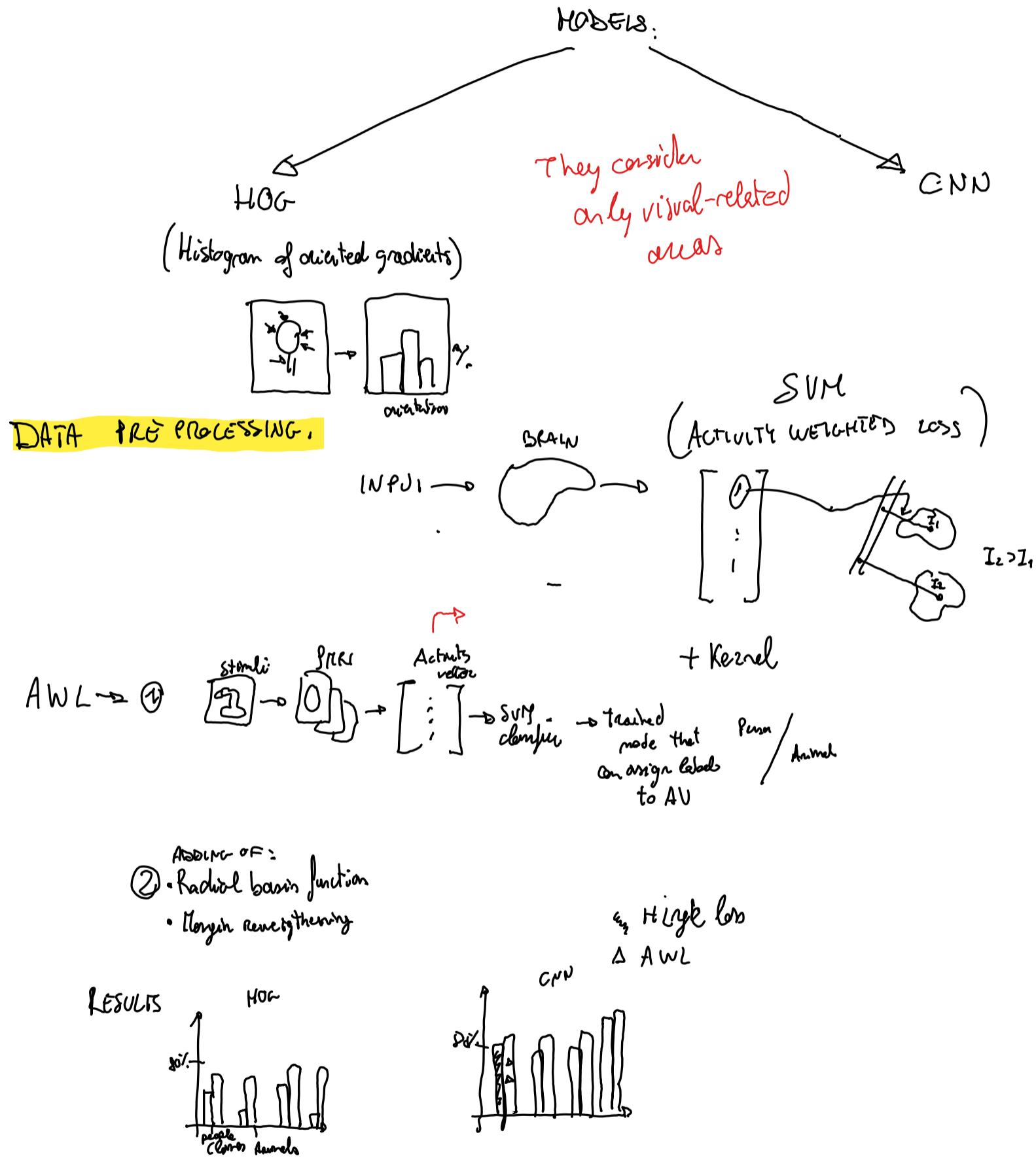


USING HUMAN BRAIN ACTIVITY TO GUIDE ML

Fang et al.



GOAL: improve representational geometry of ML methods



CONCLUSIONS:

ML can be improved by brain data

SURPRISE

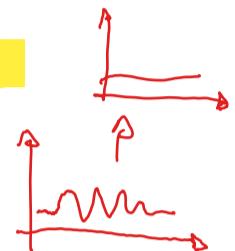
LEARNING AND PREDICTION (Henson)

SURPRISE-RESPONSE : depends on knowledge + low-level perception + accumulation of evidence, biases, noise

→ STUDY OF EXPECTATION: probe the state of a cognitive system prior to stimulus presentation and independent of stimuli-guided responses

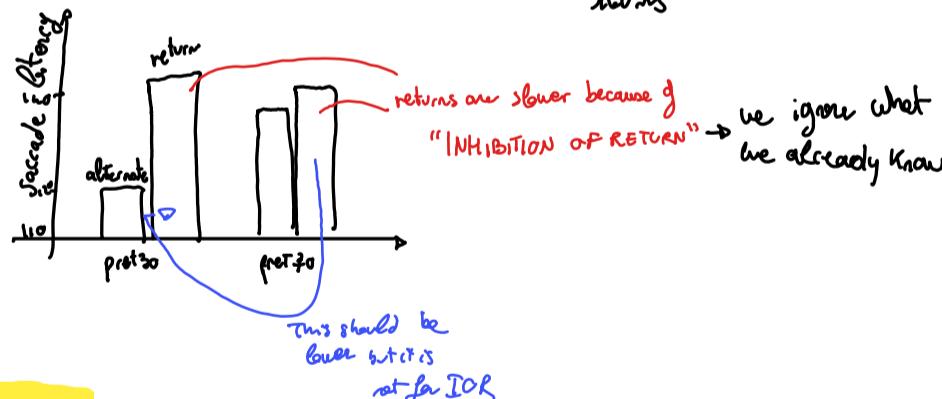
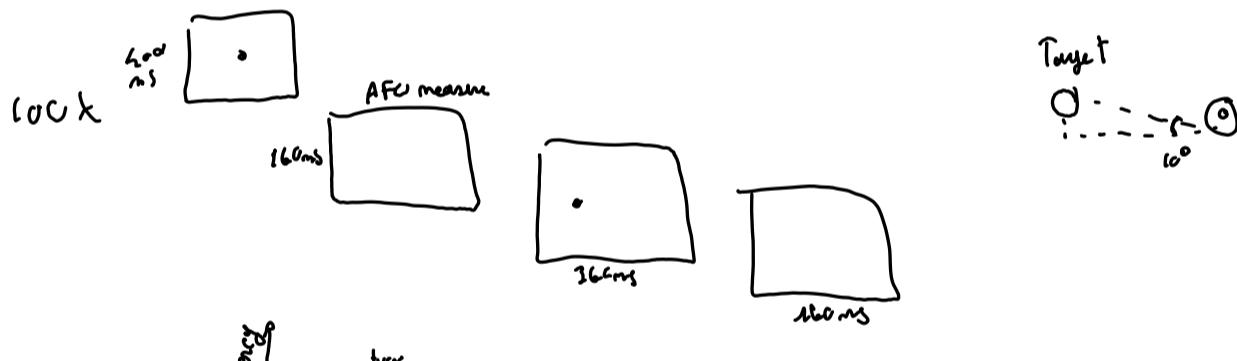
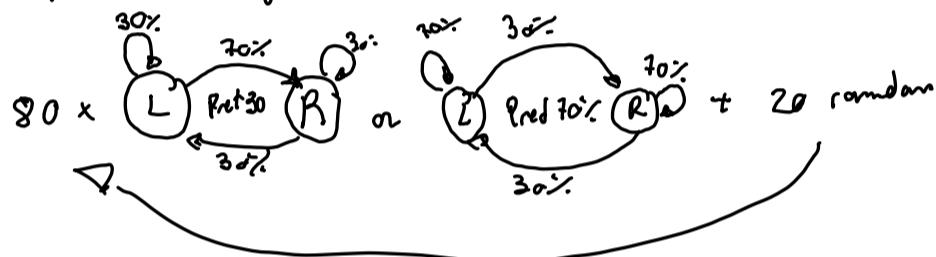
→ Previous works showed that presenting predictable stimuli have an impact on alpha waves
(after some time the brain signal decrease → LEARNING)

MEG



Dataset: N=21 volunteers

2 series → 100 samples with target position depending on Markov process (Pret 70, Pret 30)



Two important parameters:

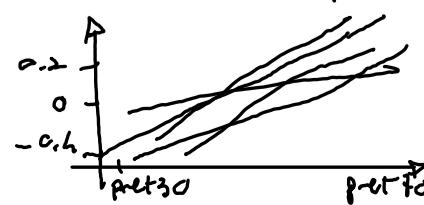
1. AGO (Anticipatory Gaze offset)

2. AFO (Anticipatory Fixation offset) → recording of gaze offset in relation to prior target

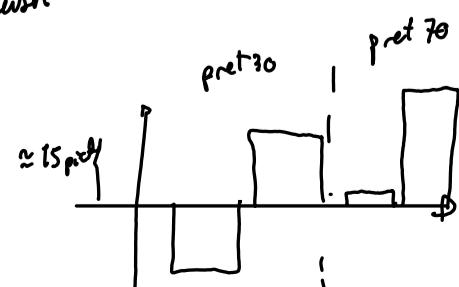
measured last 10ms of pretarget blank screen previous target

↓
We can plot "fixation density"
pret 70, pret 30

negative AFO positive AFO



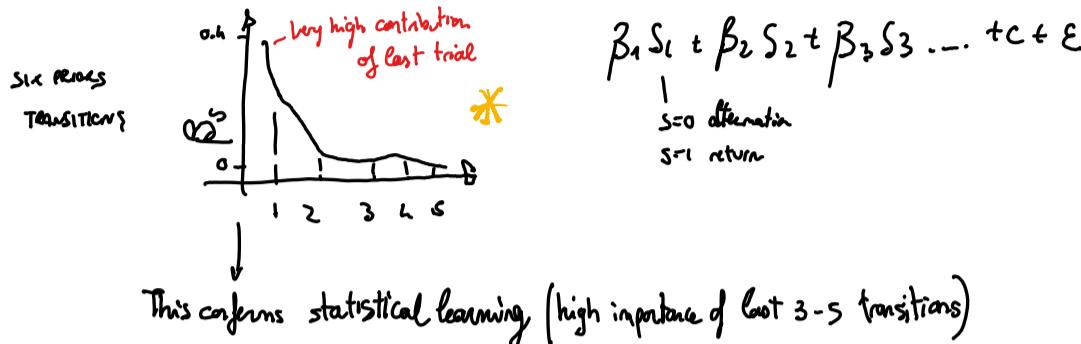
1 line per person



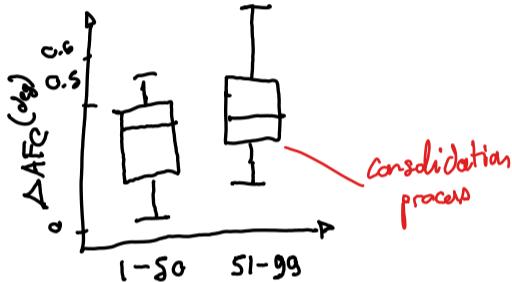
VALIDATION → INTERNAL (SPLIT-HALF) RELIABILITY: deriving two separate ΔAFO per participant: one for ODD and one for EVEN trials $S_{t+1} = C, S_t$

indicator of learning

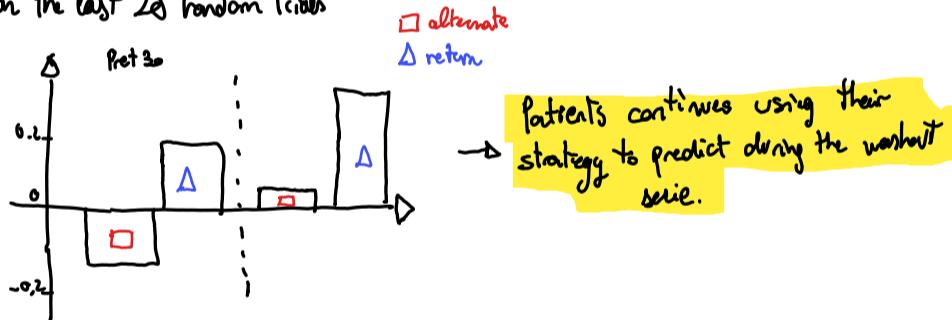
TIME SCALE OF LEARNING: we tend to do not predict the most likely outcome → DIRECTIONAL: we should always bet on the most likely bias on last results



TIME SCALES OF KNOWLEDGE CONSOLIDATION AND FORGETTING: verified looking at bets on the first split and the second split (S_{50}) (S_{99})



WASHOUT: Look at bets on the last 20 random trials



How CAN WE MODEL THIS BEHAVIOR?

• RIECKER-WAGNER learning model

"error driven learning model"
we can modify α according to

$$\begin{cases} P_{ret}(t+1) = P_{ret}(t) + \alpha(1 - P_{ret}(t)) & \text{after a return} \\ P_{ret}(t+1) = P_{ret}(t) - \alpha P_{ret}(t) & \text{after alternation} \end{cases}$$

α learning rate

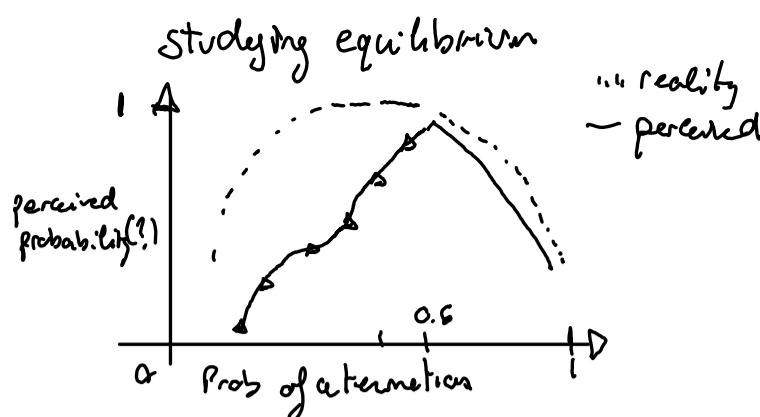
$$\Delta AFC(t+1) = K \cdot P_{ret}(t+1)$$

K scaling factor transforming internal probability to overt behavior

e.g.: $P_{ret}(t)=0.6, \alpha=0.1$

ret $\rightarrow P_{ret}(t+1)=0.64$
alt $\rightarrow P_{ret}(t+1)=0.56$

Problem: for $ret > 0.5$, people will start predicting returns because prob. exceed chance while in reality we mis-estimate chance \rightarrow equilibrium at $P_{ret}=0.6$

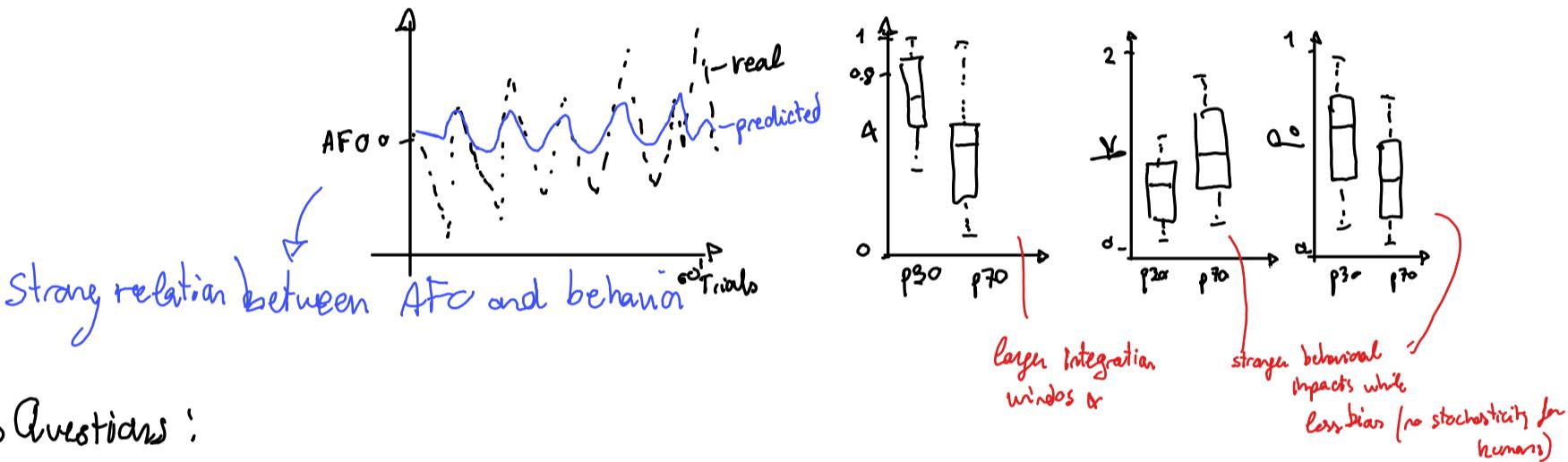


We can now connect the Rescorla-Wagner model with this information

$$\left\{ \begin{array}{l} P_{\text{ret}}(t+1) = P_{\text{ret}}(t) + \alpha (1 - P_{\text{ret}}(t)) \quad \text{after return} \\ P_{\text{ret}}(t+1) = P_{\text{ret}}(t) - \alpha P_{\text{ret}}(t) \quad \text{after alternation} \\ \text{AFC}(t+1) = K \cdot (P_{\text{ret}}(t+1) - P_0) \quad \begin{array}{l} \text{probability equilibrium} \\ \text{point} \\ \text{(bias)} \end{array} \end{array} \right.$$

Now that we have a good model \Rightarrow TRAIN

Leave one out ($19/20$ participants)



Questions:

\rightarrow If AFC measure prediction:

Larger AFC \rightarrow faster return saccade and slower alternation saccade
 ≥ 0
 verified by the experiment

\rightarrow Which measure captures more info about predictions? AFC or SL

We can answer considering **MUTUAL INFORMATION (MI)** \rightarrow consider the variation of uncertainty knowing another variable

$$I(x,w) = \sum_{x \in X} \sum_{w \in W} p(x,w) \log \left(\frac{p(x,w)}{p(x)p(w)} \right)$$

Result: AFC conveyed about twice as much information about the statistical process wrt SL

CONCLUSIONS! \rightarrow Learning can be studied quantifying observable behavior over time

o macro-scale properties of the env. can be separated to micro-scale occurrences in recent past \rightarrow Both impact behavior

o learning translates into surprise signals \rightarrow slow reaction to surprising events



o Learning produces anticipation/prediction and we can model it with error-driven learning

o Anticipatory signals can contain more inf. about the env. than stimulus related responses

SURPRISE AND NEGATION

Requirements:

- Beliefs

- Attention to the world

- Appropriate contrast between event and beliefs

- suddenly, unexpectedly experience

- Novelty? → Def: not being previously experienced or encountered

Ekman & Davidson (1994) → surprise is an emotion that arises from mismatch between expectation and what is observed

not beliefs but outcomes
of beliefs about the world

"VIOLATION OF EXPECTATION"

→ We can consider as a formal account of S as the "probability of the encountered event"

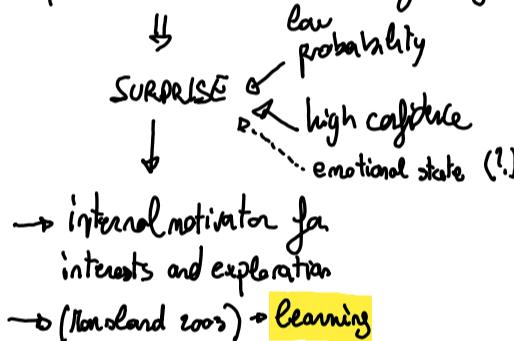


as Frequency distribution



"Surprise is proportional to the expected probability of an encountered event"

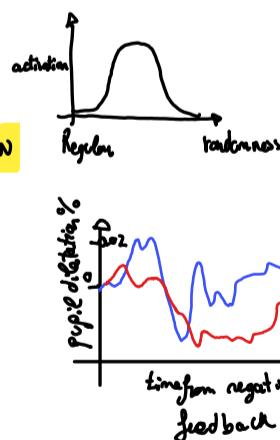
Problem: Surprise does not capture confidence → determined by # of evidences



USE OF SURPRISE IN LEARNING → Learning is proportional to surprise

investigation of neural correlates of responses related to surprise

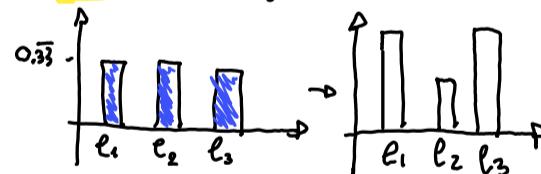
SURPRISE IN THE BRAIN



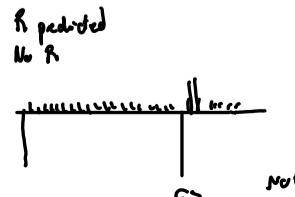
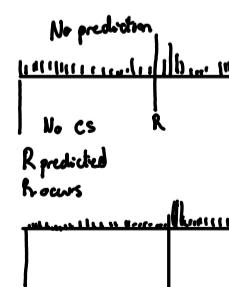
Stimulus → Reward paradigms

(CS = stimulus learned to be followed by reward)

Shannon theory



Repetition priming: brain activity goes down when we have repetition



MULTIPLE SURPRISES AT DIFFERENT SCALES

- Evidence that there is both sensory specific surprise systems
general systems
- Some respond to surprise in a stimulus
- others respond to overall surprise in input stream
- These responses may track avg. surprise linearly or non-linearly (inverse v)
- They can show increased activity for EEG - low-Shannon entropy streams (many repetitions)

SURPRISE IN LANGUAGE (DISCUSSED)

Use EEG to measure neuroelectrical activity at a very fine temporal scale

words that are more "surprising" in context produce stronger response within 0.5 sec



CONTRAST MODELS:

→ improbability and unexpectedness don't really mean the same

likelihood	depends on brain processes
------------	----------------------------

ES of roulette

→ 2S/7S% vs 2S/27/38%

ES of balls extracted from a bag

whether or not you predict a different outcome seems to play an important role

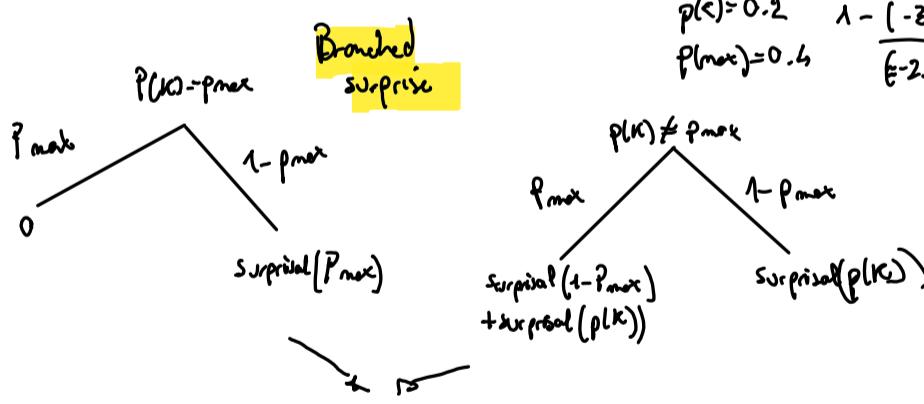
$$V_{Placedo}(k) = p_{max} - p(k)$$

extremeness of p_{max}

$$V_{bounded}(k) = 1 - \frac{\log_2(1-p(k))}{\log_2(1-p_{max})} \rightarrow \text{surprise is bigger when } p_{max} \text{ is higher}$$

$$p(k=0.2) \\ p_{max}=0.8 \rightarrow 1 - \frac{(-3)}{(-1)} = -2$$

$$p(\epsilon)=0.2 \\ p(max)=0.4 \quad 1 - \frac{(-3)}{(-2.8)} \approx 0.8 \quad ??$$



a psychological model explains experienced surprise as a function of outcomes and P_{max}

SUMMARY:

- Surprise ≠ Novelty (expectation and deviation in context)
- If surprise = probability → ignore confidence and the role of alternative expectation
- Contrast Models → combine all together but are difficult to test
- Humans can code for surprise of a stimuli but also the avg. surprise of an entire series of stimuli (independent aspects of the environment)

Finalmente



LANGUAGE COMPREHENSION AND INTEGRATION

- we predict the future using learned association and recent past
- Question: do we store language associations? And how do we use them for comprehension?

- ① From experience **memorize structures** (words or types of phrases)
- ② we learn **abstract rules** that specify constraints on grammaticality
- ③ we acquire **statistical knowledge** on co-occurrence of language entities

Behavioral psychology: S-R model for language → word is stimulus

Language comprehension: the learned sequences of adjacent elements are internally represented as automatically characterizing sentence as it is encountered

BEHAVIORISM: Knowledge of language is knowing which words follow another word

} NO EXPLANATION
OF MEANING

THE SR APPROACH (which word comes after another) Horse → races ↗ resolve just ① and ②

- { 1. The → X → Y-es (article, noun, verb)
2. The → Xes → Y- (def^{def}, noun, verb)

①

Acceptance of some level of **abstraction** → discourse structures → sentence → longer phrase → combination of shorter ones → single words

e.g.: "NP₁ V + ed NP₂ ↔ NP₂ was V + ed by NP₁" ↗
| |
Kernel Dative sentence
declarative constructions
transformed into sentence
Transformation → what we learn

②

TWO PARALLEL IDEAS

+ other views SIGN, ASSIGN
Zebra = horse + stripes

SUBSTITUTION VIEW

Signs (words) substituted for the objects in the world and both produce a similar reaction

DISPOSITION VIEW

Language produces a tendency (disposition) towards certain behaviors

CHOMSKY'S CHALLENGE TO S-R LEARNING (ABSTRACT RULES) 1959

Multiple weakness:

1. inability to explain comprehension of sentences not heard previously
2. inability to handle hierarchical structure and long-distance dependencies



TRANSFORMATIONAL GRAMMAR MODEL

→ GRAMMAR RULE DRIVEN COMPOSITIONALITY

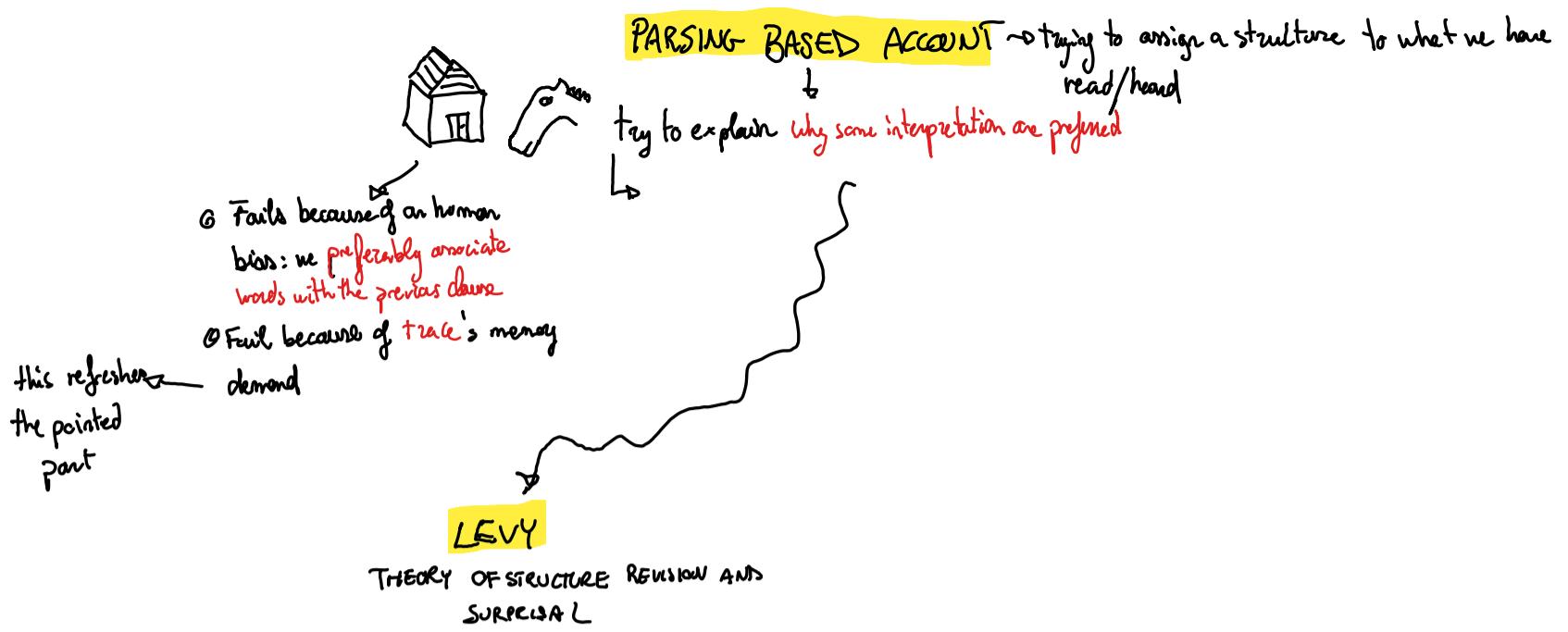
acoustic representation improves when words are not independent

↓ what is the keypoint?

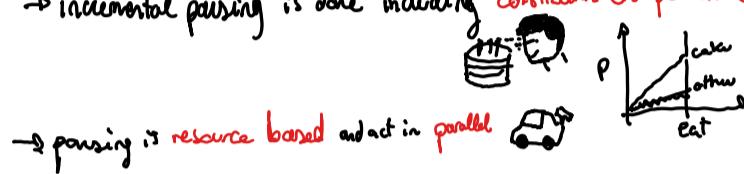
SENTENCE MEANING IS CAPTURED BY ITS STRUCTURE

↓
DERIVATIONAL THEORY OF COMPLEXITY → abandoned because of criticism
of method





→ incremental parsing is done including constraints on potential continuations



→ parsing is resource based and act in parallel

→ includes the role of expectation which consists on ranking possible continuations based on experience

① distributions change ($w(1\dots w-1)$ vs $w(1\dots w)$) can be studied with Kullback Liebler divergence

② Levy found out that surprise $\approx -\log(P(\text{word}|\text{context}))$

→ we don't need to predict but just a MODEL for each word in time

PCFG

→ incoming words sharper next-constituent expectations and thereby decreases surprise at the final word itself

A phrase is difficult is related just to how expensive is to integrate it into prior context

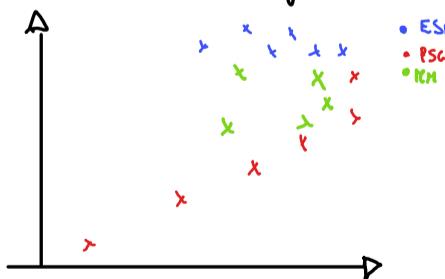
FRANK AND BoD's

→ Difficulty = prob. without syntax

3 Test → PSG (phrase structure grammar) → They tried 4 combinations
 ↓
 ↗ Markov Process (N-gram)
 ↗ Echo-state Network

Test how surprising is the truth
 ↓
 → surprise → bad accuracy

Result → based on → PSYCHOBOLICAL ACCURACY reading time
 ↗ LINGUISTIC ACCURACY from corpus



→ DATASET: Brain activity of listeners that listen to a story

→ METHOD: • Describe how surprising was each word

• Compute vector-space representation

$$\text{surprisal}(t) = -\log_{10} P(W_t | W_1, \dots, W_{t-1}) \rightarrow \text{set window to 10}$$

meaning → adjacency of a word with 1 of 1000 anchors (like Mitchell)

$$v_i = p(c_i | w) / p(c_i)$$

$$\text{Semantic weighted surprise} = P(W_t | W_{t-1}, \dots, W_{t-10}) \sum_i w_i v_i p(c_i)$$

semantic vector \vec{s}_t , $t=67$
 anchor value (\cdot) / scaling factor
 from corpus
 2nd order Markov-based surprise
 semantic weight using anchors

Regression Model

→ Predict brain activation given word and context



Compute semantic weighted surprise

outcome → good prediction

Test 1 → remove probability → result is bad

Keeping only
semantic part

Test 2 → remove conditional prob. → acceptable
Keeping a raw probability → result

→ To conclude: might be difficult to disentangle different
Surprise models as they may be highly correlated
SWS and LS were correlated at $r=0.77$