

Connectionist and Statistical Language Processing

Lecture 7: Learning Linguistic Structure in Simple Recurrent Networks



Matthew W Crocker

Computerlinguistik
Universität des Saarlandes

Reading: J Elman (1991). Distributed Representations, simple recurrent networks, and grammatical structure. *Machine Learning*.
J Elman (1993). Learning and development in neural networks: the importance of starting small. *Cognition*, **48**:71-99.

Summary of Elman 1990

- Some problems change their nature when expressed as temporally:
 - E.g. sequential XOR developed frequency sensitive units
 - Time varying error signal can be a clue to temporal structure:
 - Lower error in prediction suggests structure exists
 - Increased sequential dependencies don't result in worse performance:
 - Longer, more variable sequences were successfully learned
 - Also, the network was able to make partial predictions (e.g. "consonant")
 - The representation of time and memory is task dependent:
 - Networks intermix immediate task, with performing a task over time
 - No explicit representation of time: rather "processing in context"
 - Memory is bound up inextricably with the processing mechanisms
 - Representation need not be flat, atomistic or unstructured:
 - Sequential inputs give rise to "hierarchical" internal representations
- "SRNs can discover rich representations implicit in many tasks, including structure which unfolds over time"**

Challenges for a connectionist account

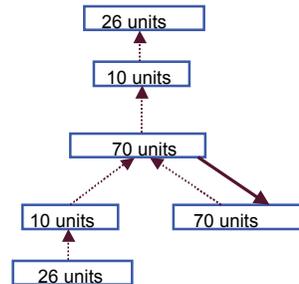
- What is the nature of the linguistic representations?
 - ❑ Localist representations seem too limited (fixed and simplistic)
 - ❑ Distributed are poorly understood, but greater capacity, can be learned
- How can complex structural relationships such as constituency be represented? Consider “noun” versus “subject” versus “role”:
 - ❑ The boy broke the *window*
 - ❑ The rock broke the *window*
 - ❑ The *window* broke
- How can the “open-ended” nature of language be accommodated by a fixed resource system?
 - ❑ Especially problematic for localist representations
- In a famous article, Fodor & Pylyshyn argue that connectionist models:
 - ❑ Cannot encode for the fully compositional structure/nature of language
 - ❑ Cannot provide for the open-ended generative capacity

Learning Linguistic Structure

- Construct a language, generated by a grammar which enforces diverse linguistic constraints:
 - ❑ Subcategorisation
 - ❑ Recursive embedding
 - ❑ Long-distance dependencies
- Training the network:
 - ❑ Prediction task
 - ❑ Structure of the training data is necessary
- Assess the performance:
 - ❑ Evaluation of predictions (as in Elman 1990), not RMS error
 - ❑ Cluster analysis? Only really informs us of the similarity of words, not the dynamics of processing
 - ❑ Principle component analysis: permits us to investigate the role of specific hidden units

Learning Constituency: Elman (1991)

- So far, we have seen how SRNs can find structure in sequences
- How can complex structural relationships such as constituency be represented?
- The Stimuli:
 - Lexicon of 23 items
 - Encoded orthogonally, in 26 bit vector
- Grammar:
 - q $S \rightarrow NP VP \text{ " . "}$
 - q $NP \rightarrow PropN \mid N \mid N RC$
 - q $VP \rightarrow V (NP)$
 - q $RC \rightarrow \text{who NP VP} \mid \text{who VP (NP)}$
 - q $N \rightarrow \text{boy} \mid \text{girl} \mid \text{cat} \mid \text{dog} \mid \text{boys} \mid \text{girls} \mid \text{cats} \mid \text{dogs}$
 - q $PropN \rightarrow \text{John} \mid \text{Mary}$
 - q $V \rightarrow \text{chase} \mid \text{feed} \mid \text{see} \mid \text{hear} \mid \text{walk} \mid \text{live} \mid \text{chases} \mid \text{feeds} \mid \text{sees} \mid \text{hears} \mid \text{walks} \mid \text{lives}$
 - Number agreement, verb argument patterns



Training

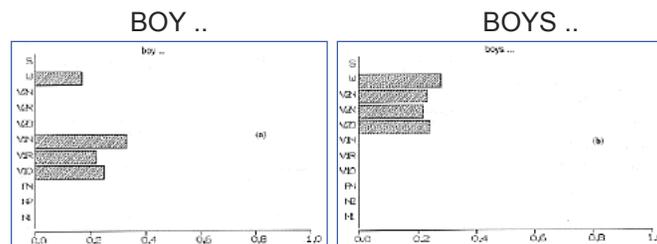
- Verb subcategorization
 - Transitives: *hit, feed*
 - Optional transitives: *see, hear*
 - Intransitives: *walk, live*
- Interaction with relative clauses:
 - q *Dog who chases cat sees girl*
 - q *Dog who cat chases sees girl*
 - Agreement can span arbitrary distance
 - Subcategorization doesn't always hold (superficially)
- Recursion: *Boys who girls who dogs chase see hear*
- Viable sentences: where should end of sentence occur?
 - *Boys see (.) dogs (.) who see (.) girls (.) who hear (.) .*
- Words are not explicitly encoded for number, subcat, or category

Training

- At any given point, the training set contained 10000 sentences, which were presented to the network 5 times
- The composition of sentences varied over time:
 - Phase 1: Only simple sentence (no relative clauses)
 - + 34,605 words forming 10000 sentences
 - Phase 2: 25% complex and 75% simple
 - + Sentence length from 3-13 words, mean: 3.92
 - Phase 3: 50/50, mean sentence length 4.38
 - Phase 4: 75% complex, 25% simple, max: 16, mean: 6
- WHY?: Pilot simulations showed the network was unable to learn the task when given the full range of complex data from the beginning.
- Focussing on simpler data first, the network learned quickly, and was then able to learn the more complex patterns.
- Earlier simple learning, usefully constrained later learning

Performance

- Weights are frozen and test on a novel set of data (as in phase 4).
- Since the solution is non-deterministic, the networks outputs were compared the context dependent likelihood vector of all words following the current input (as done in the previous simulation)
 - Error was 0.177, mean cosine: 0.852
 - High level of performance in prediction
- Performance on Specific Inputs
- Simple agreement:

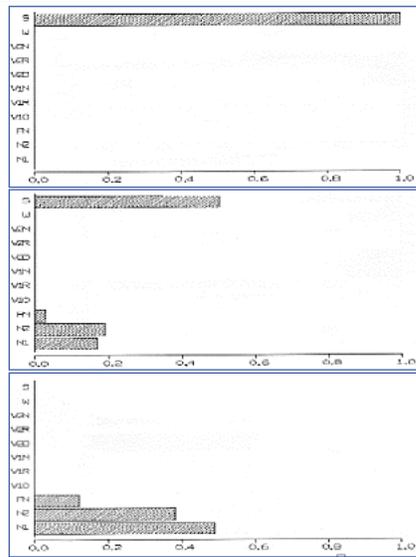


Subcategorization

- Intransitive: “Boy lives ...”
 - Must be a sentence, period expected

- Optional: “Boy sees ...”
 - Can be followed by either a period,
 - Or some NP

- Transitive: “Boy chases ...”
 - Requires some object



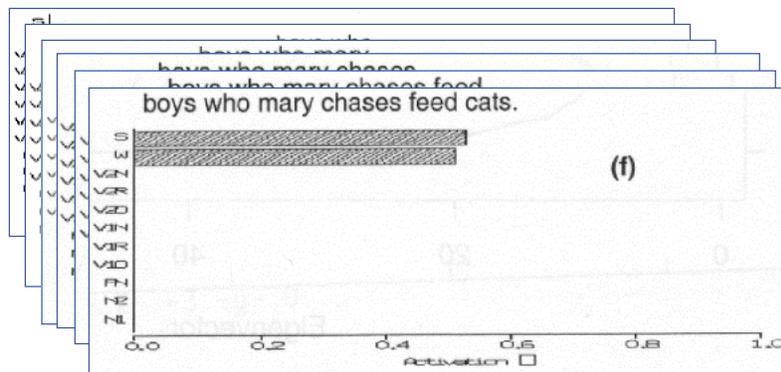
© Matthew W. Crocker

Connectionist and Statistical Language Processing

9

Processing complex sentences

- “Boys who many chases feed cats”
 - Long distance
 - + Agreement: Boys ... feed
 - + Subcategorization: chases is transitive but in a relative clause
 - + Sentence end: all outstanding “expectations” must be resolved



© Matthew W. Crocker

Connectionist and Statistical Language Processing

10

Prediction reconsidered

- SRNs are trained on the *prediction* task:
 - “Self-supervised learning”: no other teacher required
- Prediction forces the network to discover regularities in the temporal order of the input
- Validity of the the prediction tasks:
 - It is clearly not the “goal” of linguistic competence
 - But there is evidence that people can/do make predictions
 - Violated expectation results in distinct patterns of brain activity (ERPs)
- If children do make predictions, which are then falsified, this might constitute an indirect form of negative evidence, required for language learning.

Results

- Learning was only possible when the network was forced to begin with simpler input
 - This effectively restricted the range of data to which the networks were exposed during initial learning
 - Contrasts with other results showing the entire dataset is necessary to avoid getting stuck in local minima (e.g. XOR)
- This behaviour partially resembles that of children:
 - Children do not begin by mastering language in all its complexity
 - They begin with simplest structures, incrementally building their “grammar”
- But the simulation achieves this by manipulation the environment:
 - This does not seem an accurate model of the situation in which children learn language
 - While adults do modify their speech, it is not clear they make such grammatical modifications
 - Children hear all exemplars of language from the beginning

General results

- Limitations of the simulations/results:
 - Memory capacity remains un-probed
 - Generalisation is not really tested
 - ✦ Can the network inferentially extend what is know about the types of NPs learned to NPs with different structures
 - Truly a “toy” in terms of real linguistic complexity and subtlety
 - ✦ E.g. lexical ambiguity, verb-argument structures, structural complexity and constraints
- Successes
 - Representations are distributed, which means less rigid resource bounds
 - Context sensitivity, but can respond to contexts which are more “abstractly” defined
 - ✦ Thus can exhibit more general, abstract behaviour
 - ✦ Symbolic models are primarily context insensitive
- Connectionist models begin with local, context sensitive observations
- Symbolic models begin with generalisation and abstractions

A Second Simulation

- While it's not the case that the environment changes, it true that the child changes during the language acquisition period
- Solution: keep the environment constant, but allow the network to undergo change during learning
- Incremental memory:
 - Evidence of a gradual increase in memory and attention span in children
 - In the SRN, memory is supplied by the “context” units
 - Memory can be explicitly limited by depriving the network, periodically, access to this feedback
- In a second simulation, training began with limited memory span which was gradually increased:
 - Train began from the outset with the full “adult” language (which was previously unlearnable)

Training with Incremental Memory

- Phase 1:
 - Training on corpus generated from the entire grammar
 - Recurrent feedback was eliminated after every 3 or 4 words, by setting all context units to 0.5
 - Longer training phase (12 epochs, rather than 5)
- Phase 2:
 - New corpus (to avoid memorization)
 - Memory window increased to 4-5 words
 - 5 epochs
- Phase 3: 5-6 word window
- Phase 4: 6-7 word window
- Phase 5: no explicit memory limitation implemented

- Performance: as good as on the previous simulation

Analysing the solution

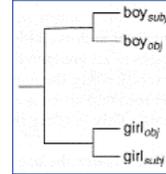
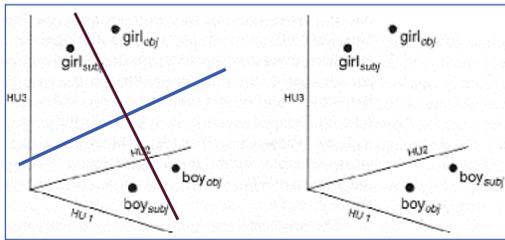
- Hidden units permit the network to derive a *functionally-based* representation, in contrast to a *form-based* representation of inputs

- Various dimensions of the internal representation were used for:
 - Individual words, category, number, grammatical role, level of embedding, and verb argument type
 - The high-dimensionality of the hidden unit vectors (70 in this simulation) makes direct inspection difficult

- Solution: Principle Component Analysis can be used to identify which dimensions of the internal state represent these different factors
 - This allows us to visualise the movement of the network through a state space for a particular factor, by discovering which units are relevant

Principle Component Analysis

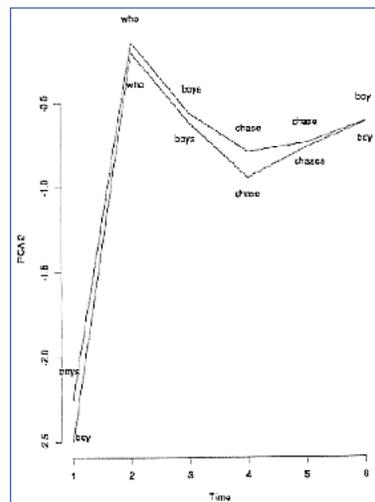
- Suppose we're interested in analysing a network with 3 hidden units and 4 patterns of activation, corresponding to: boy_{subj} , $girl_{subj}$, boy_{obj} , $girl_{obj}$
- Cluster analysis might reveal the following structure:
 - But nothing of the subj/obj representation is revealed
- If we look at the entire space, however, we can get more information about the representations:



- Since visualising more than 3 dimensions is difficult, PCA permits us to identify which "units" account for most of the variation.
 - Reveals partially "localist" representations in the "distributed" hidden units

Examples of Principle Components: 1

- Agreement
 - *Boy who boys chase chases boy*
 - *Boys who boys chase chase boy*
- The 2nd PCA encodes agreement in the main clause



Examples of Principle Components: 2

Transitivity

- Boy chases boy*
- Boy sees boy*
- Boy walks*

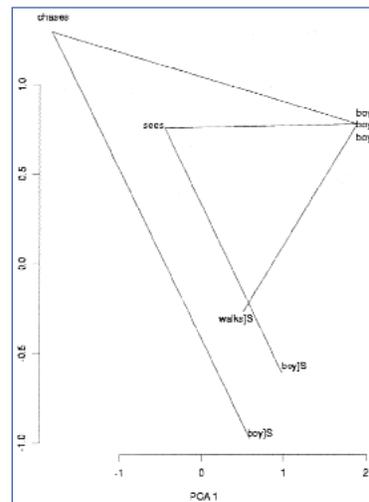
Two principle components: 1 & 3

PCA 1:

- Nouns on the right
- Verbs left

PCA 2:

- Intrans: low
- Optional trans: mid
- Transitive: high



© Matthew W. Crocker

Connectionist and Statistical Language Processing

19

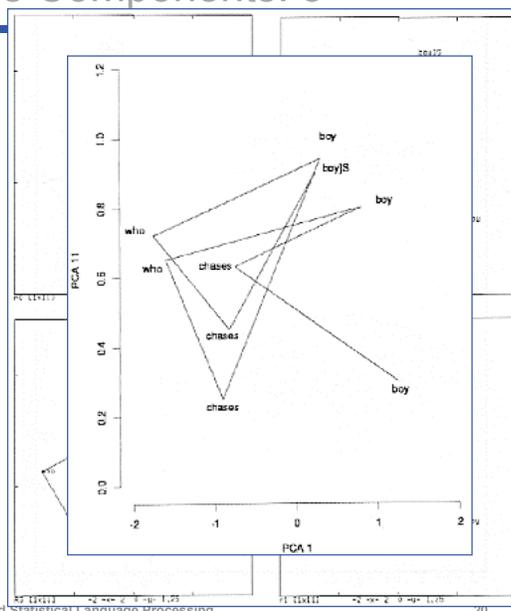
Examples of Principle Components: 3

Right embedding:

- Boy chases boy*
- Boy who chases boy chases boy*
- Boy chases boy who chases boy*
- Boy chases boy who chases boy who chases boy*

PCA 11 and 1:

- "Embedded clause are shifted to the left"
- "RCs appear nearer the noun they modify"



© Matthew W. Crocker

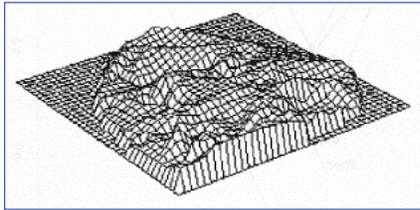
Connectionist and Statistical Language Processing

20

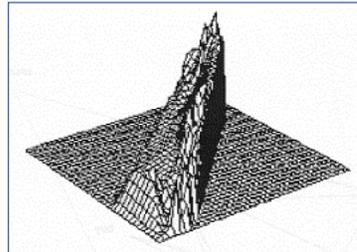
PCA analysis of “Starting Small”

- We can use “Principle Component Analysis” to examine particularly important dimensions of the networks solutions more globally:
 - Sample of the points visited in the hidden unit space as the network processes 1000 random sentences
- The results of PCA after training:

Training on the full data set



Incremental training



The right plot reveals are more clearly “organised” use of the state space

Comments

- To solve the task, the network must learn the sources of variance (number, category, verb-type, and embedding)
- If the network is presented with the complete corpus from the start:
 - The complex interaction of these factors, long-distance dependencies, makes discovering the sources of variance difficult
 - The resulting solution is imperfect, and internal representation don't reflect the true sources of variance
- When incremental learning takes place (in either form):
 - The network begins with exposure to only some of the data
 - ✦ Limited environment: simple sentences only
 - ✦ Limited mechanisms: simple sentences + noise (hence longer training)
 - Only the first 3 sources of variance, and no long-distance dependencies
- Subsequent learning is constrained (or guided) by the early learning of, and commitment to, these basic grammatical factors
 - Thus initial memory limitations permit the network to focus on learning the subset of facts which lay the foundation for future success

The importance of starting small

- Networks rely on the representativeness of the training set:
 - Small samples may not provide sufficient evidence for generalisation
 - ✦ Possibly poor estimates of the populations statistics
 - ✦ Some generalisations may be possible from a small sample, but are later ruled out
 - Early in training the sample is necessarily small
- The representation of experience:
 - Exemplar-based learning models store all prior experience, and such early data can then be re-accessed to subsequently help form new hypotheses
 - SRNs do not do this: each input has its relatively minor effect on changing the weights (towards a solution), and then disappears. Persistence is only in the change made to the network.
- Constraints on new hypotheses, and continuity of search:
 - Changes in a symbolic systems may lead to suddenly different solutions
 - ✦ This is often ok, if it can be checked against the prior experience
 - Gradient descent learning makes it difficult for a network to make dramatic changes in its solution: search is continuous, along the error surface
 - Once committed to an erroneous generalisation, the network might not escape from a local minima

Starting small (continued)

- Network are most sensitive during the early period of learning:
 - Non-linearity (the logistic activation function) means that weight modifications are less likely as learning progresses
 - ✦ Input is “squashed” to a value between 0 and 1
 - ✦ Non-linearity means that the function is most sensitive for inputs around 0 (output is 0.5)
 - ✦ Nodes are typically initialised randomly about 0, so netinput is also near 0
 - ✦ Thus the network is highly sensitive
 - Sigmoid function become “saturated” for large +/- inputs
 - ✦ As learning proceeds units accrue activation
 - ✦ Weight change is a function of the error and slope of the activation function
 - ✦ This will become smaller as units activations become saturation, regardless of how large the error is
 - Thus escaping from local minima becomes increasingly difficult
- Thus most learning occurs when information is least reliable

Conclusions

- Learning language is difficult because:
 - Learning linguistic primitives is obscured by the full complexity of grammatical structure
 - Learning complex structure is difficult because the network lacks knowledge of the basic primitive representations
- Incremental learning shows how a system can learn a complex system by having better initial data:
 - Initially impoverished memory provides a natural filter for complex structures early in learning so the network can learn the basic forms of linguistic regularities
 - As the memory is expanded, the network can use what it knows to handle increasingly complex inputs
 - Noise, present in the early data, tends to keep the network in a state of flux, helping it to avoid committing to false generalisations

Summary of SRNs ...

- Finding structure in time/sequences:
 - Learns dependencies spanning more than a single transition
 - Learns dependencies of variable length
 - Learns to make partial predictions from structure input
 - + Prediction of **consonants**, or particular lexical **classes**
- Learning from various input encodings:
 - Localist encoding: XOR and 1 bit per word
 - Distributed:
 - + Structured: letter sequences where consonants have a distinguished feature
 - + Random: words mapped to random 5 bit sequence
- Learns both general categories (types) and specific behaviours (tokens) based purely on distributional evidence
- What are the limitations of SRNs
 - Do they simply learn co-occurrences and contingent probabilities?
 - Can they learn more complex aspects of linguistic structure?

Summary

- Implicit representation of time, reflected in the dynamic behaviour of the network: not explicitly encoded.
- The importance of starting small:
 - Learning the more complex language was only possible by first learning simpler aspects of the grammar
- Outstanding problems:
 - Is grammatical structure really being learned?
 - Full linguistic complexity
 - + Ambiguity: lexical, syntactic, semantic
 - + Structural: subadjacency, islands, extraction, ...
 - + Scale: large lexicons, large structures
- Statistical/Probabilistic Models
 - Connectionist models have a highly probabilistic nature:
 - + Learn regularities in a way which is sensitive to and reflect frequency
 - We can model language by directly applying probabilistic theory
 - We can combine symbolic and probabilistic approaches to achieve hybrid symbolic/sub-symbolic systems.