

Statistical Dialogue Systems NPFL099 Statistické Dialogové systémy

5. Dialogue State Tracking

Ondřej Dušek & Vojtěch Hudeček

http://ufal.cz/npfl099

31.10.2019



Dialogue State Tracking

- Dialogue management consist of:
 - **State update** ← here we need DST
 - Action selection (later)
- Dialogue State needed to remember what was said in the past
 - tracking the dialogue progress
 - summary of the whole dialogue history
 - basis for action selection decisions

U: I'm looking for a restaurant in the <u>city centre</u>. S: OK, what kind of food do you like? U: Chinese.

- **X** S: What part of town do you have in mind?
- X S: Sure, the Golden Dragon is a good Chinese restaurant. It is located in the west part of town.

S: Sure, the Golden Dragon is a good Chinese restaurant. It is located in the <u>city centre</u>.



Dialogue State Contents

- "All that is used when the system decides what to say next" (Henderson, 2015)
- User goal/preferences ~ NLU output
 - slots & values provided (search constraints)
 - information requested

Past system actions

- information provided
 - slots and values
 - list of venues offered
- slots confirmed
 S: OK, Chinese food. [...]
- slots requested

U: Give me the address of <u>the first one</u> you talked about. U: Is there <u>any other</u> place in this area?

— S: What time would you like to leave?

- Other semantic context
 - user/system utterance: bye, thank you, repeat, restart etc.



Problems with Dialogue State

- NLU is unreliable
 - takes unreliable ASR output
 - makes mistakes by itself some utterances are ambiguous
 - output might conflict with ontology
- Possible solutions:
 - detect contradictions, ask for confirmation
 - ignore low-confidence NLU input
 - what's "low"?
 - what if we ignore 10x the same thing?
- Better solution: make the state probabilistic **belief state**

ASR: 0.5 I'm looking for an expensive hotel
 0.5 I'm looking for inexpensive hotels

NLU: 0.3 inform(type=restaurant, stars=5)

only hotels have stars!

Belief State



- Assume we don't know the true current dialogue state s_t
 - states (what the user wants) influence **observations** o_t (what the system hears)
 - based on observations o_t & system actions a_t, we can estimate a probability distribution b(s) over all possible states – belief state
- More robust than using dialogue state directly
 - accumulates probability mass over multiple turns
 - low confidence if the user repeats it, we get it the 2nd time
 - accumulates probability over NLU n-best lists
- Plays well with probabilistic dialogue policies (POMDPs)
 - but not only them rule-based, too





(from Milica Gašić's slides)

(=belief state)

Basic Discriminative Belief Tracker



- simplify assume each slot is independent:
 - state $\mathbf{s} = [s^1, \dots s^N]$, belief $b(\mathbf{s}_t) = \prod_i b(s_t^i)$
- Always trust the NLU
 - this makes the model parameter-free
 - ...and basically rule-based
 - but very fast, with reasonable performance

user silent about slot *i*

 $p(s_t^i | a_{t-1}^i, s_{t-1}^i, o_t^i) = \begin{cases} p(o_t^i) \text{ if } s_t^i = o_t^i \wedge o_t^i \neq \textcircled{s} \\ p(o_t^i) \text{ if } s_t^i = s_{t-1}^i \wedge o_t^i = \textcircled{s} \\ 0 \text{ otherwise} \end{cases}$

update
$$b(s_t^i) = \sum_{\substack{s_{t-1}^i, o_t^i \\ \text{model}}} p(s_t^i | a_{t-1}^i, s_{t-1}^i, o_t^i) b(s_{t-1}^i)$$
 substitution
discriminative model $b(s_t^i) = \begin{cases} p(s_t^i = \textcircled{r}) p(o_t^i = \textcircled{r}) \\ p(o_t^i = s_t^i) + p(o_t^i = \textcircled{r}) p(s_t^i = s_{t-1}^i) \\ p(o_t^i = s_t^i) + p(o_t^i = \textcircled{r}) p(s_t^i = s_{t-1}^i) \end{cases}$ otherwise
NPFL099 L5 2019 $(\overset{(\check{z})}{\underset{\text{http://www.aclweb.org/anthology/W13-4070}{}} the belief state update rule is deterministic$

Basic Feed-forward Tracker

- a simple feed-forward network
 - input features (w.r.t. slot-value v & time t)
 - SLU score of *v*
 - n-best rank of v
 - user & system act type
 - ... domain-independent, low-level NLU outputs
 - 3 tanh layers
 - output softmax (= probability distribution over values
- static does not model dialogue as a sequence
 - uses a sliding window current time t + few steps back + \sum previous



 $\mathbb{P}(s=v) = \frac{e^{E(t,v)}}{Z} \\
\mathbb{P}(s \notin S_{t,s}) = \frac{e^{B}}{Z} \\
Z = e^{B} + \sum_{v' \in S_{t,s}} e^{E(t,v')}$

Basic RNN Tracker



- plain sigmoid RNN with a memory vector
 - not quite LSTM/GRU, but close
 - memory updated separately, used in belief update
- does not need NLU
 - turn features = lexicalized + delexicalized n-grams from ASR n-best list, weighted by confidence
- delexicalization very harsh: <slot> <value>
 - you don't even know which slot it is
 - this apparently somewhat helps the system generalize across domains
- dynamic explicitly models dialogue as sequence
 - using the network recurrence



(Vodolán et al., 2017) http://arxiv.org/abs/1702.06336

Neural/Rule Hybrid

- explicit update over belief
 - per-slot model (separate for each slot)
 - simple update rule *R*
 - for a value: add $a \cdot \text{current NLU}$ confidence, normalize
 - differentiable, can be trained end-to-end
 - trained models *F*, *G* provide *a*
 - F is generic LSTM, G is value specific feed-forward
- needs an NLU, but postprocesses it
 - input & output of tracker NLU step
 = prob. dist. of informs over slot values in current turn
 - generic & specific part again



Incremental Recurrent Tracker



- simple: LSTM over words + classification on hidden states
 - runs over the whole dialogue history (user utterances + system actions)
 - classification can occur after each word, right as it comes in from ASR
- also dynamic/sequential
- also doesn't use any NLU
 - infrequent values are delexicalized (otherwise it can't learn them)
- slightly worse performance possible causes:
 - only uses ASR 1-best
 - long recurrences (no hierarchy)

(Žilka & Jurčíček, 2015) https://dl.acm.org/citation.cfm?id=2955040 http://arxiv.org/abs/1507.03471

NPFL099 L5 2019



(Mrkšić et al., 2017) https://www.aclweb.org/anthology/P17-1163











- Better use of data: Getting rid of delexicalization
 - pre-trained word vectors important!
 - shared parameters
 - RNN/CNN feature extractors
- Discriminative learning slot values

Belief state update $\mathbb{P}(s, v \mid h^{1:t}, sys^{1:t-1}) = \lambda \mathbb{P}(s, v \mid h^t, sys^{t-1}) + (1 - \lambda) \mathbb{P}(s, v \mid h^{1:t-1}, sys^{1:t-2})$

| DST Model | DSTC2 | | WOZ 2.0 | |
|--|-------|----------|---------|----------|
| | Goals | Requests | Goals | Requests |
| Delexicalisation-Based Model | 69.1 | 95.7 | 70.8 | 87.1 |
| Delexicalisation-Based Model + Semantic Dictionary | 72.9* | 95.7 | 83.7* | 87.6 |
| NEURAL BELIEF TRACKER: NBT-DNN | 72.6* | 96.4 | 84.4* | 91.2* |
| NEURAL BELIEF TRACKER: NBT-CNN | 73.4* | 96.5 | 84.2* | 91.6* |

Candidate Ranking

(Rastogi et al., 2017) https://arxiv.org/abs/1712.10224



- Previous systems consider all values for each slot
 - this is a problem for open-ended slots (e.g. restaurant name)
 - enumerating over all takes ages, some are previously unseen
- Alternative: always consider just K candidates
 - use last *K* candidates from system actions and NLU output
 - NB: only way history is incorporated here!
 - select from them using a per-slot softmax





Candidate Ranking – representation UFA^{L}

- Using BiGRU over lexicalized & delexicalized utterance
- Features:

NPFL099 L5 2019

- **utterance** last GRU state + indicators for non-slot DAs (user & prev. system)
- slot indicators for DAs with this slot (user & prev. system) + last turn scores for null & dontcare
- candidate GRU states over matched value words inform(slot=value)
 + indicators for DAs with this slot & value (user & prev. system)



bye(), affirm()

Multi-value Candidate Ranking



- What if multiple values are true?
 - previous approach picks one (softmax)
 - use set of binary classifiers (log loss) instead
- + more flexible regarding candidates
 - can be past k from NLU, but also just current ASR n-grams
 - this model keeps track of context by itself

(Goel et al., 2018) http://arxiv.org/abs/1811.12891





Hybrid Classify/Rank

- Ranking is faster & more flexible
- Classification over all values is more accurate
 - at least for most slots, where # of values is limited
- Solution: combine classification & ranking
 - choose best model for each slot based on dev data performance

metric: joint goal accuracy

- Ranking approach multi-value from previous slide
- Classification approach straightforward: hierarchical LSTM + per-slot feed-forward + softmax

- exact match on dialogue state (most probable value only) Accuracy Method Majority Baseline 1.5%MultiWOZ-2.0 Benchmark 25.83% **Ranking only** 31.11% (29.73%) **Classification only** 40.74% (38.42%) Hybrid 44.24% (42.33%) ensemble (majority vote of 3 models) single model separate for each slot Values for Slot_E Values for Slot

Dialogue LSTM (E

entence LSTM (E

shared

across slots

(Goel et al., 2019) http://arxiv.org/abs/1907.00883

NPFL099 L5 2019

Using BERT

(Chao & Lane, 2019) http://arxiv.org/abs/1907.03040



- Very basic:
 - run BERT over previous system & current user utterance
 - from 1st token's representation, get a decision: *none/dontcare*/span
 - per-slot (BERT is shared, but the final decision is slot-specific)
 - span = need to find a concrete value as a span somewhere in the text
 - predict start & end token of the span using 2 softmaxes
 - carry-over across multiple turns is rule-based:
 - if *none* is predicted, keep previous value, otherwise change it



Slot-Utterance Matching Belief Tracker

- different take on BERT trackers
 - inspired by reading comprehension
 - considers "domain slot" a question & tries to find the value in the input utterance
- tracker over BERT
 - attention + turn-based RNN
 - attention in current utterance
 - RNN (LSTM/GRU) for carry-over of past values
 - layer normalization to match BERT outputs
 - BERT includes layer normalization by default
 - trained to match the correct values in the utterance
 - loss: distance of true value BERT encoding from the tracker output (Euclidean/Cosine)

(Lee et al., 2019) http://arxiv.org/abs/1907.07421





slot BERT encoder

utterance BERT encoder (prev. system + user)

Even More Reading Comprehension



(Gao et al., 2019) https://www.aclweb.org/anthology/W19-5932/

- Also uses BERT, but not necessarily
 - works slightly worse with random-initialized word embeddings
- sequence of 3 decisions
 - do we carry over last turn's prediction? (Yes/No)
 - if no: what kind of answer are we looking for? (*yes/no/dontcare*/span of text)
 - if span: predict span's start and end



Dialogue State as SQL

User goal is a query → why not SQL query?

(Yu et al., 2019)

http://arxiv.org/abs/1909.05378 http://arxiv.org/abs/1906.02285

- Text-to-SQL models used for tracking
 - with contextual enhancements, input:
 - all user inputs so far
 - previous system response
 - database schema
- Seq2seq-based model example:
 - hierarchical LSTM for encoding user & system
 - database column embeddings
 averaged embeddings over table + column name
 - decoder:
 - decide between SQL keyword vs. column
 - then select which keyword / column via softmax
- So far, experimental performance is low



WELCOM

| D_1 : Database about student dormitories conta | aining 5 tables |
|---|---|
| Q_1 : What are the names of all the dorms? | INFORM_SQL |
| S_1 : SELECT dorm name FROM dorm | |
| A_1 : (Result table with many entries) | |
| R_1 : This is the list of the names of all the dorms. | CONFIRM_SQL |
| Q_2 : Which of those dorms have a TV lounge? | INFORM_SQL |
| S2 : SELECT T1.dorm_name FROM dorm AS T1 JOIN T2 ON T1.dormid = T2.dormid JOIN dorm_ame T2.amenid = T3.amenid WHERE T3.amenity_na Lounge' | has_amenity AS enity AS T3 ON me = `TV |
| A_2 : (Result table with many entries) | |
| R_2 : This shows the names of dorms with TV lounges. | CONFIRM_SQL |
| Q ₃ : What dorms have no study rooms as amenities? | AMBIGUOUS |
| R ₃ : Do you mean among those with TV Lounges? | CLARIFY |
| $Q_4: Yes.$ | AFFIRM |
| S4: SELECT T1.dorm_name FROM dorm AS T1 JOIN 1 AS T2 ON T1.dormid = T2.dormid JOIN dorm_ ON T2.amenid = T3.amenid WHERE T3.amenity Lounge' EXCEPT SELECT T1.dorm_name FROM d JOIN has_amenity AS T2 ON T1.dormid = T2. dorm_amenity AS T3 ON T2.amenid = T3.amen T3.amenity_name = `Study Room' | has_amenity amenity AS T3 _name = `TV orm AS T1 dormid JOIN id WHERE |
| A ₄ : Fawlty Towers | |
| R ₄ : Fawlty Towers is the name of the dorm that has a TV lounge but not a study room as an amenity. | CONFIRM_SQL |
| Q_8 : Thanks! | THANK_YOU |
| | |

 R_8 : You are welcome.

Summary



- State tracking is needed to maintain user goal over multiple turns
- Best to make the state probabilistic **belief state**
- Architectures many options
 - good NLU + rules works well!
 - neural, hybrid
 - static (sliding-window) vs. dynamic (recurrent, modelling dialogue as sequence)
 - with/without NLU
 - classifiers vs. candidate rankers vs. reading comprehension
 - classifiers are more accurate than rankers but slower, limited to seen values
 - reading comprehension is a very new approach, works nicely but probably slow
 - BERT & co. as usual good but slow
 - incremental not used too much so far
- Alternative/experimental: SQL instead of slots/values





Contact us:

odusek@ufal.mff.cuni.cz hudecek@ufal.mff.cuni.cz

Get these slides here:

http://ufal.cz/npfl099

References/Inspiration/Further:

- Filip Jurčíček's slides (Charles University): <u>https://ufal.mff.cuni.cz/~jurcicek/NPFL099-SDS-2014LS/</u>
- Milica Gašić's slides (Cambridge University): <u>http://mi.eng.cam.ac.uk/~mg436/teaching.html</u>
- Henderson (2015): Machine Learning for Dialog State Tracking: A Review <u>https://ai.google/research/pubs/pub44018</u>