

Lec 12 (Reinforcement) Learning w Human Feedback

Getting human feedback (for generation task)

So far: maximum likelihood training

- Problems:
- Some mistakes are worse than others
 $a \leftrightarrow \text{the}$ maybe less bad] but to likelihood, these are equal
 $\text{Tokyo} \leftrightarrow \text{Pittsburgh}$ bad
 - Gold standard can be bad
 e.g. Google translate posted online
 - Overexposure bias
 previously appeared word more likely to appear again
 model repeats itself
 model not seen bad text

Measure output quality

- Objective assessment
 - annotated data
- Subjective human eval
 - generate hypotheses, human rate them
 - trait-based assessments
 - fluency
 - adequacy (is it semantically correct)
 - factuality (in general or w.r.t. input)
 - coherence
 - preference rating - which is better?
 - good for relative comp but not globally how good
 - hybrid: preference rating + score rating
 - use ELO ranking
 - pairwise contest to rank many models
 - from MT: multi-dim quality metric
 - annotate spans in seq, categorise them, and rate their severity (minor/major)
 - + consistent, fine-grained
 - time consuming
- Machine predict human preference "reward model", "automatic eval"
 - reference human feedback for some outputs
 - Embedding based (BERTScore)
 - unsupervised, calc embedding similarity btwn output & ref
 - embed each token, pairwise cos sim btwn two seqs, mk sim matrix, take max along a dim, ...
 - Regression-based (COMET)
 - supervised training of emb-based regressor
 - QA-based
 - ask LM how good it is wrt reference (GEMBA)
 - e.g. GPT
 - ask LM what the mistakes are (AutoMQM)
 - + more consistent results
- Downstream task
 - Intrinsic: quality of downstream output itself
 - Extrinsic: by utility

Human feedback postproc

- Normalise per annotator
- Average
- Drop examples human disagree on

Meta evaluation

- Are metrics associated?
 - .. human/auto eval correlated?
- Can use large shared decks to do this

Error and Risk

Generate output \hat{Y}

Calculate the badness error $(Y, \hat{Y}) = 1 - \text{eval}(Y, \hat{Y})$

Try to minimise error, usually can't do gradient-based opt

$$\text{risk}(X, Y, \theta) = \sum_{\tilde{Y}} P(\tilde{Y}|X; \theta) \text{error}(Y, \tilde{Y})$$

$\tilde{Y} \leftarrow$ but too many things in possible outputs. intractable

Try minimise risk instead

Sampling (5-50 sentences) (to make tractable)

$$\text{risk}(X, Y, \theta) = \sum_{\tilde{Y} \in S} \frac{P(\tilde{Y}|X; \theta)}{Z} \text{error}(Y, \tilde{Y})$$

→ Sampling, n-best search, etc.

Reinforcement Learning, Policy Grad

Given: env X
 action A
 delayed reward R

Generation: X - conversation history
 A - next token
 R - human pressing buttons, call center outcome

Cool: X - code, env
 A - token
 R - unit tests, speed, readability, etc.

Why reinforcement in NLP

- When having delayed reward info
- Have a latent variable
- Have sequence-level evaluation

Supervised MLE loss $\ell_{\text{super}}(X, Y) = -\log P(Y|X)$

for imitation learning

Self training $\ell_{\text{self}}(X) = -\log P(\hat{Y}|X)$

optimises toward both good and bad output but actually works in some situation

- Co-training
- ...

▷ Policy Grad / REINFORCE

scale loss by reward $\ell_R(X, \hat{Y}) = -R(Y, \hat{Y}) \log P(Y|X)$

when reward is 1

Problem: how do we know which action led to reward?

- best: immediate reward
- worse: eventually reward

→ Just not deal with it usually ppl just do this

→ Decay reward for future event (viz. scale by time)

Reinforcement Learning Stabilisation

- Unstable reasons
- Sample individual output
 - Very large output space

▷ Pre-train with MLE, then switch to RL

just always do this

▷ Regularise to existing LM

- not let the model diverge too far away

→ KL regularisation

Regularisation = reward stuff - β KL with other model

→ Proximal policy opt

use clipping to not reward large jumps

push model back if model goes too bad

▷ Adding a baseline

account for hard vs easy examples by having expectation on reward

opt with reward - baseline

baseline can be arbitrary

- use model to predict reward
 - sentence level
 - decoder state level
- calc mean reward of batch

▷ Contrasting Pairwise Examples

- learn directly from pairwise human judgement

→ DPO - upweight better output, downweight the other

▷ Increase batch size

- RL samples high variance so higher batch size helps