Sequence Discriminative Training

Hao Tang

Automatic Speech Recognition—ASR Lecture 14 09 March 2022

Hao Tang Sequence Discriminative Training









• It is not about whether one is better than the other.

- It is not about whether one is better than the other.
- Should we use the samples (and computation) to learn the decision boundary or the data distribution?

- It is not about whether one is better than the other.
- Should we use the samples (and computation) to learn the decision boundary or the data distribution?
- The discriminative approach might be a better solution when the boundary is simple to learn.

- It is not about whether one is better than the other.
- Should we use the samples (and computation) to learn the decision boundary or the data distribution?
- The discriminative approach might be a better solution when the boundary is simple to learn.
- If the goal is to do prediction, we should focus on learning the bounary.

• Map words to a sequence of phones $\label{eq:speech} \text{speech} \to \text{s } p \text{ iy ch}$

• Map words to a sequence of phones

 $\mathsf{speech} \to \mathsf{s} \mathsf{ p} \mathsf{ iy ch}$

• Chain phone HMMs



• Map words to a sequence of phones

 $\mathsf{speech} \to \mathsf{s} \mathsf{ p} \mathsf{ iy ch}$

• Chain phone HMMs



• Find parameters that maximize p(X)

• p(X) really should be p(X|W).

'▲ 문 ► ▲ 문 ►

- p(X) really should be p(X|W).
- p(X, Q) really should be p(X, Q|W).

• • = • • = •

э

- p(X) really should be p(X|W).
- p(X, Q) really should be p(X, Q|W).
- *Q* is a valid sequence for *W* if the phone sequence produced by *Q* is the pronunciation of *W*.

- p(X) really should be p(X|W).
- p(X, Q) really should be p(X, Q|W).
- *Q* is a valid sequence for *W* if the phone sequence produced by *Q* is the pronunciation of *W*.
- s1 s1 s2 s2 s2 s3 s3 p1 p2 p3 iy1 iy1 iy1 iy2 iy2 iy2 iy3 iy3 iy3 ch1 ch2 ch3 is a valid sequence for the word "speech" if there are 22 frames.

• p(X, Q|W) = 0 when Q is not a valid state sequence for W.

()

- p(X, Q|W) = 0 when Q is not a valid state sequence for W.
- $p(X, Q|W) = p(q_1)p(x_1|q_1)\prod_{t=2}^{T} p(q_t|q_{t-1})p(x_t|q_t)$ when Q is a valid state sequence for W.

- p(X, Q|W) = 0 when Q is not a valid state sequence for W.
- $p(X, Q|W) = p(q_1)p(x_1|q_1)\prod_{t=2}^{T} p(q_t|q_{t-1})p(x_t|q_t)$ when Q is a valid state sequence for W.
- Use B(W) to denote the set of valid state sequences for W.

- p(X, Q|W) = 0 when Q is not a valid state sequence for W.
- $p(X, Q|W) = p(q_1)p(x_1|q_1)\prod_{t=2}^{T} p(q_t|q_{t-1})p(x_t|q_t)$ when Q is a valid state sequence for W.
- Use B(W) to denote the set of valid state sequences for W.

•
$$p(X|W) = \sum_{Q \in B(W)} p(X, Q|W)$$

- p(X, Q|W) = 0 when Q is not a valid state sequence for W.
- $p(X, Q|W) = p(q_1)p(x_1|q_1)\prod_{t=2}^{T} p(q_t|q_{t-1})p(x_t|q_t)$ when Q is a valid state sequence for W.
- Use B(W) to denote the set of valid state sequences for W.

•
$$p(X|W) = \sum_{Q \in B(W)} p(X, Q|W)$$

B(W) is represented as an FST U ∘ H ∘ C ∘ L ∘ W, where U represent the frames and W represent the word sequence.

Recap of HMM Training

Hao Tang Sequence Discriminative Training

æ

・日・ ・ ヨ・・

• $\operatorname{argmax}_{\lambda} p(X|W)$ can be solved with EM or gradient descent.

A B M A B M

- $\operatorname{argmax}_{\lambda} p(X|W)$ can be solved with EM or gradient descent.
- $\operatorname{argmax}_{\lambda} p(X|W)$ is a generative approach.

A B > A B >

- $\operatorname{argmax}_{\lambda} p(X|W)$ can be solved with EM or gradient descent.
- $\operatorname{argmax}_{\lambda} p(X|W)$ is a generative approach.
- The discriminative approach solves $\operatorname{argmax}_{\lambda} p(W|X)$.

(E)

Maximum Mutual Information (MMI) (Bahl et al., 1986)

$$p(W|X) = \frac{p(X|W)p(W)}{p(X)} = \frac{p(X|W)p(W)}{\sum_{W'} p(X|W')p(W')}$$

Maximum Mutual Information (MMI) (Bahl et al., 1986)

$$p(W|X) = \frac{p(X|W)p(W)}{p(X)} = \frac{p(X|W)p(W)}{\sum_{W'} p(X|W')p(W')}$$

- How do we compute the numerator p(X|W)p(W)?
- How do we compute the denominator $\sum_{W'} p(X|W')p(W')$?
- Why is this called maximum mutual information (MMI)?

p(X|W)p(W)

Hao Tang Sequence Discriminative Training

p(X|W)p(W)

• Map words to a sequence of phones

 $\mathsf{speech} \to \mathsf{s} \mathsf{ p} \mathsf{ iy ch}$

• Chain phone HMMs



- Compute p(X|W)
- Compute p(W) with a language model

• It's computationally expensive to compute the denominator exactly.

$$\sum_{W'} p(X|W')p(W')$$

• Instead we can approximate it with a set of high-probability word sequences *D*.

$$\sum_{W'\in D} p(X|W')p(W')$$

• The set of high-probability sequences *D* (stored in a compact format) is called a **lattice**.





Lattice



→ time

▲御▶ ▲ 臣▶ ▲ 臣▶

















・ロト ・四ト ・ヨト ・ヨト





・ロト ・四ト ・ヨト ・ヨト



→ time

・ロト ・四ト ・ヨト ・ヨト





・ロト ・四ト ・ヨト ・ヨト

Forward-Backward on Graphs



문 문 문

Forward-Backward on Graphs



$$\alpha(\mathbf{v}) = \sum_{\mathbf{e} \in \mathsf{in}(\mathbf{v})} p(X_{\mathbf{e}} | W_{\mathbf{e}}) \alpha(\mathsf{tail}(\mathbf{e}))$$

문 문 문

• Running the forward algorithm on a lattice only gives

$$\alpha(\mathsf{final}) = \sum_{w' \in D} p(X|W')$$

• Running the forward algorithm on a lattice composed with an LM gives

$$\alpha(\mathsf{final}) = \sum_{w' \in D} p(X|W')p(W')$$

$$p(W|X) = \frac{p(X|W)p(W)}{\sum_{W'} p(X|W')p(W')}$$

- Generate lattice (through beam search)
- Run the forward algorithm
- Compute the gradient
- Do gradient update

 $\frac{\partial L}{\partial p(X_e|W_e)}$



æ

(*) * 문 * * 문 *

$$\frac{\partial L}{\partial p(X_e|W_e)} = \sum_{v} \frac{\partial L}{\partial \alpha(v)} \frac{\partial \alpha(v)}{\partial p(X_e|W_e)}$$

æ

▷ ▲ 글 ▶

$$\frac{\partial L}{\partial p(X_e|W_e)} = \sum_{v} \frac{\partial L}{\partial \alpha(v)} \frac{\partial \alpha(v)}{\partial p(X_e|W_e)}$$
$$= \sum_{v} \frac{\partial L}{\partial \alpha(v)} \alpha(\mathsf{tail}(e)) \mathbb{1}_{v=\mathsf{head}(e)}$$

æ

▷ ▲ 글 ▶

$$\begin{aligned} \frac{\partial L}{\partial p(X_e|W_e)} &= \sum_{v} \frac{\partial L}{\partial \alpha(v)} \frac{\partial \alpha(v)}{\partial p(X_e|W_e)} \\ &= \sum_{v} \frac{\partial L}{\partial \alpha(v)} \alpha(\mathsf{tail}(e)) \mathbb{1}_{v=\mathsf{head}(e)} \\ &= \frac{\partial L}{\partial \alpha(\mathsf{head}(e))} \alpha(\mathsf{tail}(e)) \end{aligned}$$

æ

イロト イヨト イヨト イヨト

$$\frac{\partial L}{\partial p(X_e|W_e)} = \sum_{v} \frac{\partial L}{\partial \alpha(v)} \frac{\partial \alpha(v)}{\partial p(X_e|W_e)}$$
$$= \sum_{v} \frac{\partial L}{\partial \alpha(v)} \alpha(\operatorname{tail}(e)) \mathbb{1}_{v=\operatorname{head}(e)}$$
$$= \frac{\partial L}{\partial \alpha(\operatorname{head}(e))} \alpha(\operatorname{tail}(e))$$

$$\frac{\partial L}{\partial \alpha(u)} = \sum_{v} \frac{\partial L}{\partial \alpha(v)} \frac{\partial \alpha(v)}{\partial \alpha(u)}$$

æ

▷ ▲ 글 ▶

$$\frac{\partial L}{\partial p(X_e|W_e)} = \sum_{v} \frac{\partial L}{\partial \alpha(v)} \frac{\partial \alpha(v)}{\partial p(X_e|W_e)}$$
$$= \sum_{v} \frac{\partial L}{\partial \alpha(v)} \alpha(\operatorname{tail}(e)) \mathbb{1}_{v=\operatorname{head}(e)}$$
$$= \frac{\partial L}{\partial \alpha(\operatorname{head}(e))} \alpha(\operatorname{tail}(e))$$
$$\frac{\partial L}{\partial \alpha(u)} = \sum_{v} \frac{\partial L}{\partial \alpha(v)} \frac{\partial \alpha(v)}{\partial \alpha(u)}$$
$$= \sum_{v} \frac{\partial L}{\partial \alpha(v)} \sum_{e} P(X_e|W_e) \mathbb{1}_{u=\operatorname{tail}(e), v=\operatorname{head}(e)}$$

・ロト ・回ト ・ヨト ・ヨト

$$\frac{\partial L}{\partial p(X_e|W_e)} = \sum_{v} \frac{\partial L}{\partial \alpha(v)} \frac{\partial \alpha(v)}{\partial p(X_e|W_e)}$$
$$= \sum_{v} \frac{\partial L}{\partial \alpha(v)} \alpha(\operatorname{tail}(e)) \mathbb{1}_{v=\operatorname{head}(e)}$$
$$= \frac{\partial L}{\partial \alpha(\operatorname{head}(e))} \alpha(\operatorname{tail}(e))$$
$$\frac{\partial L}{\partial \alpha(u)} = \sum_{v} \frac{\partial L}{\partial \alpha(v)} \frac{\partial \alpha(v)}{\partial \alpha(u)}$$
$$= \sum_{v} \frac{\partial L}{\partial \alpha(v)} \sum_{e} P(X_e|W_e) \mathbb{1}_{u=\operatorname{tail}(e),v=\operatorname{head}(e)}$$
$$= \sum_{e \in \operatorname{out}(u)} \frac{\partial L}{\partial \alpha(\operatorname{head}(e))} P(X_e|W_e)$$

$$\underset{\lambda}{\operatorname{argmax}} \sum_{X,W} p(X,W) \log \frac{p(X,W)}{P(X)P(W)}$$

$$\operatorname{argmax}_{\lambda} \sum_{X,W} p(X, W) \log \frac{p(X, W)}{P(X)P(W)}$$

=
$$\operatorname{argmax}_{\lambda} \sum_{X,W} p(X, W) \log p(W|X) - \sum_{W,X} p(X, W) \log p(W)$$

$$\arg\max_{\lambda} \sum_{X,W} p(X, W) \log \frac{p(X, W)}{P(X)P(W)}$$

=
$$\arg\max_{\lambda} \sum_{X,W} p(X, W) \log p(W|X) - \sum_{W,X} p(X, W) \log p(W)$$

=
$$\arg\max_{\lambda} \sum_{X,W} p(X, W) \log p(W|X) - \sum_{W} p(W) \log p(W)$$

$$\begin{aligned} \arg\max_{\lambda} \sum_{X,W} p(X,W) \log \frac{p(X,W)}{P(X)P(W)} \\ &= \arg\max_{\lambda} \sum_{X,W} p(X,W) \log p(W|X) - \sum_{W,X} p(X,W) \log p(W) \\ &= \arg\max_{\lambda} \sum_{X,W} p(X,W) \log p(W|X) - \sum_{W} p(W) \log p(W) \\ &\approx \arg\max_{\lambda} \frac{1}{N} \sum_{n=1}^{N} \log p(W_n|X_n) \end{aligned}$$

- It is actually possible (just computationally expensive) to compute the denominator $\sum_{W'} p(X|W')p(W')$ exactly with the help of GPU.
- The trick is to realize that the forward algorithm is a matrix multiplication.

- It is a discriminative approach.
- It considers a language model.
- It provides the same solution as minimizing the zero-one loss.

Properties of MMI

- It is a discriminative approach.
- It considers a language model.
- It provides the same solution as minimizing the zero-one loss.

 $\mathbb{E}_{W' \sim p(W'|X)}[\mathbb{1}_{W \neq W'}]$

- It is a discriminative approach.
- It considers a language model.
- It provides the same solution as minimizing the zero-one loss.

$$\mathbb{E}_{W' \sim \rho(W'|X)}[\mathbb{1}_{W \neq W'}] = 1 - \mathbb{E}_{W' \sim \rho(W'|X)}[\mathbb{1}_{W = W'}]$$

- It is a discriminative approach.
- It considers a language model.
- It provides the same solution as minimizing the zero-one loss.

$$\mathbb{E}_{W'\sim p(W'|X)}[\mathbbm{1}_{W
eq W'}] = 1 - \mathbb{E}_{W'\sim p(W'|X)}[\mathbbm{1}_{W=W'}] \ = 1 - p(W|X)$$

- It is a discriminative approach.
- It considers a language model.
- It provides the same solution as minimizing the zero-one loss.

$$\mathbb{E}_{W'\sim p(W'|X)}[\mathbbm{1}_{W
eq W'}] = 1 - \mathbb{E}_{W'\sim p(W'|X)}[\mathbbm{1}_{W=W'}] \ = 1 - p(W|X)$$

$$\operatorname*{argmin}_{\lambda} \mathbb{E}_{W' \sim p(W'|X)}[\mathbb{1}_{W \neq W'}] = \operatorname*{argmax}_{\lambda} p(W|X)$$

$\operatorname*{argmax}_{\lambda} \mathbb{E}_{W' \sim \mathcal{P}(W'|X)}[\operatorname{cost}(W,W')]$

æ

御 とくきとくきとう

$$\operatorname*{argmax}_{\lambda} \mathbb{E}_{W' \sim p(W'|X)}[\operatorname{cost}(W, W')]$$

- Allows partial credits
- Allows a user-defined cost function

< E

$\mathbb{E}_{W' \sim p(W'|X)}[\text{cost}(W, W')]$

Hao Tang Sequence Discriminative Training

æ

(* * 문 * * 문 *)

$\mathbb{E}_{\mathcal{W}' \sim \rho(\mathcal{W}'|X)}[\mathsf{cost}(\mathcal{W}, \mathcal{W}')] = \sum_{\mathcal{W}'} \rho(\mathcal{W}'|X)\mathsf{cost}(\mathcal{W}, \mathcal{W}')$

御 とくきとくきとう

э

$$\mathbb{E}_{W' \sim p(W'|X)}[\operatorname{cost}(W, W')] = \sum_{W'} p(W'|X) \operatorname{cost}(W, W')$$
$$= \frac{\sum_{W'} p(X|W') p(W') \operatorname{cost}(W, W')}{\sum_{W''} p(X|W'') p(W'')}$$

æ

イロト イ団ト イヨト イヨト

$$\mathbb{E}_{W' \sim p(W'|X)}[\operatorname{cost}(W, W')] = \sum_{W'} p(W'|X) \operatorname{cost}(W, W')$$
$$= \frac{\sum_{W'} p(X|W') p(W') \operatorname{cost}(W, W')}{\sum_{W''} p(X|W'') p(W'')}$$

- Both numerators and denominators require a lattice.
- The cost function needs to decompose according to a lattice, i.e., each edge having a cost.
- WER(W, W') does not decompose according to a lattice.



► time

→ < Ξ → <</p>

æ

Hao Tang Sequence Discriminative Training



æ

→ < Ξ → <</p>



æ

< ∃ →



- If the cost is at the phone level, the objective is called Minimum Phone Error (MPE) (Povey and Woodland, 2002).
- If the cost is at the word level, the objective is called Minimum Word Error (MWE) (Povey and Woodland, 2002).

- Discriminative vs Generative Training
- Maximum Mutual Information
- Forward-Backward on Graphs
- Minimum Bayes Risk