Introduction to Minimum Bayes Risk Decoding

CyberAgent AI Lab

Yuu Jinnai

CyberAgent



Yuu Jinnai 陣内 佑

CyberAgent AI Lab, Reinforcement Learning Team

Research Interest

- Sequential Decision Making
 - Planning and Search (Text Generation)
 - Reinforcement Learning (RLHF)









Yuu Jinnai 陣内 佑

CyberAgent AI Lab, Reinforcement Learning Team

Research Interest

- Sequential Decision Making
 - Planning and Search (Text Generation)
 - Reinforcement Learning (RLHF)









Yuu Jinnai 陣内 佑

CyberAgent AI Lab, Reinforcement Learning Team

Research Interest

- Sequential <u>Decision Making</u>
 - Planning and Search (Text Generation)
 - Reinforcement Learning (RLHF)







講談社

4

Decoding Matters More than You Think!



Q1. Question Time!

What is He? A. Dog B. Cat C. Owl D. Bear



If you know him pretend that you don't

Problem: Text Generation





Many NLP Tasks Involve Text Generation



Many NLP Tasks Involve Text Generation



Many NLP Tasks Involve Text Generation



Text Generation Problem

Given a context x and a model *P*model, generate a *desired* output



h

Text Generation Problem

Given a context x and a model Pmodel, generate a *desired* output \rightarrow This process is called decoding!



h

Text Generation Problem

Given a context x and a model Pmodel, generate a desired output \rightarrow This process is called decoding!

...but what is *desired* output?



$$P_{\text{model}}(\mathbf{h}|\mathbf{x}) \longrightarrow \mathbf{r}$$

h

If we had a PERFECT language model that exactly captures $P_{model}(\mathbf{h}|\mathbf{x})$, would the text generation problem be considered solved?

- A. Yes text generation is trivial with a perfect model.
- B. Mostly yes rare edge cases may exist.
- C. No there are many other aspects to consider.
- D. It can never be perfect so the question has no point.

If we had a PERFECT language model that exactly captures $P_{model}(\mathbf{h}|\mathbf{x})$, would the text generation problem be considered solved?

- A. Yes text generation is trivial with a perfect model.
- B. Mostly yes rare edge cases may exist.
- C. No there are many other aspects to consider.
- D. It can never be perfect so the question has no point.

If we had a PERFECT language model that exactly captures $P_{model}(\mathbf{h}|\mathbf{x})$, would the text generation problem be considered solved?

- A. Yes text generation is trivial with a perfect model.
- B. Mostly yes rare edge cases may exist.
- C. No there are many other aspects to consider.
- D. It can never be perfect so the question has no point.



Hint is...

Okay, it's time to think about Him!

What is He? A. Dog B. Cat C. Owl D. Bear



Maximum-a-posteriori (MAP) Decoding

What is He? A. Dog B. Cat C. Owl D. Bear



MAP decoding (estimate) selects the most probable option (i.e. highest probability)

Why MAP Wouldn't be Perfect?



Why MAP Wouldn't be Perfect?





Why MAP Wouldn't be Perfect?

If it's actually a cat/dog/owl, it may not be a big problem, but...



Why MAP Wouldn't be Perfect?

If it's actually a cat/dog/owl, it may not be a big problem, but... What if he's a baby bear with an angry mom behind him?



Why MAP Wouldn't be Perfect?

If it's actually a cat/dog/owl, it may not be a big problem, but... What if he's a baby bear with an angry mom behind him?



Question should be what you would do



Question should be what you would do

A. Pet it!

It is most likely a cat Ο



Question should be what you would do

A. Pet it!

It is most likely a cat Ο

B. Pull back the kid and take a closer look

Most likely a cat but it might be a Ο bear



Question should be what you would do

A. Pet it!

- It is most likely a cat Ο
- B. Pull back the kid and take a closer look
 - Most likely a cat but it might be a Ο bear

C. Ruuuuun!!!

Worst case, it might be a bear! Ο



Question should be what you would do

A. Pet it!

- It is **most likely** a cat Ο
- B. Pull back the kid and take a closer look
 - Most likely a cat but it might be a Ο

bear

C. Ruuuuun!!!

Worst case, it might be a bear! Ο



Question should be what you would do

- A. Pet it!
 - He is most likely a cat Ο

B. Pull back the kid and take a closer look

Most likely a cat but it might be Ο

some wild animal

C. Ruuuuun!!!

• There is a chance he is a bear

 \rightarrow Minimax

 \rightarrow Bayes risk

minimization

 \rightarrow Maximum-a-

posteriori (MAP)





If we had a PERFECT language model that exactly captures $P_{model}(\mathbf{h}|\mathbf{x})$, would the text generation problem be considered solved?

- A. Yes text generation is trivial with a perfect model.
- B. Mostly yes rare edge cases may exist.
- C. No there are many other aspects to consider.
- D. It can never be perfect so the question has no point.

If we had a PERFECT language model that exactly captures $P_{model}(\mathbf{h}|\mathbf{x})$, would the text generation problem be considered solved?

- A. Yes text generation is trivial with a perfect model.
- B. Mostly yes rare edge cases may exist.
- C. No there are many other aspects to consider.
 - **D.** It can never be perfect so the question has no point.

How you ACT given the probability estimate is up to you

If we had a PERFECT language model that exactly captures $P_{model}(\mathbf{h}|\mathbf{x})$, would the text generation problem be considered solved?

- A. Yes text generation is trivial with a perfect model.
 - B. Mostly yes rare edge cases may exist.
 - C. No there are many other aspects to consider.
 - D. It can never be perfect so the question has no point.

But can't we think of a "language model" that <u>takes into account of</u> <u>human preference?</u>

"Language Model" with Human Preference



But can't we think of a "language model" that <u>takes into account of</u> <u>human preference?</u> Yes. Reinforcement learning.


▲ CyberAgent Al Lab

He is Atchoum

What is He? A. Dog **B. Cat** & & & C. Owl



CyberAgent Al Lab

He is Atchoum

What is He? A. Dog **B. Cat** C. Owl D. Bear

Image from https://www.atchoumthecat.com/



He is Atchoum

What is He?

A. Dog B. Cat 🞉 🞉 C. Owl D. Bear

Image from https://www.atchoumthecat.com/







▲ CyberAgent Al Lab

You Have Multiple Options even with a Perfect Probability Model

Question should be what you would do

- A. Pet it!
 - He is most likely a cat Ο
- B. Pull back the kid and take a closer look
 - Most likely a cat but it might be Ο \rightarrow Bayes risk some wild animal minimization
- C. Ruuuuun!!!
 - There is a chance he is a bear

 \rightarrow Minimax

(MBR)

 \rightarrow Maximum-a-

posteriori (MAP)



CyberAgent Al Lab

Algorithm: Minimum Baye Risk Decoding



Procedure of Minimum Bayes Risk (MBR) Decoding (Kumar+ '04, Eikema+ '20)

1. Sample outputs randomly



Procedure of Minimum Bayes Risk (MBR) Decoding (Kumar+ '04, Eikema+ '20)

- 1. Sample outputs randomly
- Estimate the utility between the outputs using a function $u(\mathbf{h}, \mathbf{y})$ 2.



A cat

Prompt: "What's in the picture?"

kitty 🧺

.....

A cat

kitty 🧺

Procedure of Minimum Bayes Risk (MBR) Decoding (Kumar+ '04, Eikema+ '20)

- 1. Sample outputs randomly
- 2. Estimate the <u>utility</u> between the outputs using a function $u(\mathbf{h}, \mathbf{y})$

= - risk



Prompt: "What's in the picture?"

Procedure of Minimum Bayes Risk (MBR) Decoding (Kumar+ '04, Eikema+ '20)

- 1. Sample outputs randomly
- 2. Estimate the utility between the outputs using a function $u(\mathbf{h}, \mathbf{y})$
- 3. Select the output that maximizes the average utility to the others



Procedure of Minimum Bayes Risk (MBR) Decoding (Kumar+ '04, Eikema+ '20)

- 1. Sample outputs randomly
- 2. Estimate the utility between the outputs using a function $u(\mathbf{h}, \mathbf{y})$
- 3. Select the output that maximizes the average utility to the others



Interpretation of MBR Decoding

Assuming the generated samples are the possible "true answers", **minimize the average risk over them**



Utility Function Matters!

Using better utility function results in better generation

Method	Automatic Evaluation					Model	Human Eval		
		BLEU	sB leu	Chrf	YISI	BL.1	BL.2	logP	$MQM\downarrow$
Human Transl.	Ref-D	31.5	31.6	60.9	84.7	37.1	75.6	-38.0	0.388^\dagger
Beam 4		34.3	34.2	62.5	85.3	26.8	71.6	-11.5	2.030
	SBLEU	34.7	<u>34.8</u>	62.5	85.4	23.4	70.5	-11.2	1.855
	Chrf	34.2	34.3	<u>64.1</u>	85.7	25.8	71.4	-13.2	2.139
MBR	Yisi	34.2	34.2	62.8	<u>86.0</u>	26.4	71.6	-11.4	2.445
	BLEURT VO.1	29.2	29.4	60.0	84.3	<u>50.0</u>	77.1	-18.7	1.571^\dagger
	Bleurt v0.2	25.4	26.0	57.7	83.1	43.9	<u>79.0</u>	-24.4	1.661^{\dagger}

Sampling Algorithm Matters!

de-en	en-de	ru-en	en-ru
85.82	<u>86.32</u>	82.11	<u>86.13</u>
88.47	89.32	84.16	89.44
88.51	89.47	84.36	90.17
88.01	89.12	83.76	89.96
88.02	89.04	83.98	89.57
	de-en <u>85.82</u> 88.47 88.51 88.01 88.02	de-enen-de85.8286.3288.4789.3288.5189.4788.0189.1288.0289.04	de-enen-deru-en85.8286.3282.1188.4789.3284.1688.5189.4784.3688.0189.1283.7688.0289.0483.98

5)	Pseudo-Reference	de-en	en-de	ru-en	en-ru
0.0	Ancestral	85.82	87.51	82.02	88.41
II	Beam	<u>85.62</u>	87.40	<u>81.64</u>	<u>87.78</u>
n (e	Epsilon ($\epsilon = 0.02$)	85.89	87.74	82.01	88.46
ilo	Epsilon ($\epsilon = 0.02$)*	85.87	87.74	81.98	88.46
Eps	Nucleus ($p = 0.6$)	85.69	87.57	81.76	88.26
	Nucleus ($p = 0.9$)	86.04	87.82	82.18	88.61
	Beam Search	84.38	86.13	80.76	85.69
	Beam Search (ensemble)	84.30	86.06	80.91	85.74

Applications of MBR Decoding

APPLICATIONS OF MBR DECODING

SUMMARIZATION A artcificles the artcific the at is on the ground A artcific the at is on the ground A art the ground A art the ground A art the ground A art the ground A at the ground A art the		after noon.	It is a lovely after noon.	Die Katze schläft auf dem Sofa.
SUMMARIZATION A artclicles The cat is on the The cat is on the ground		A tabby cat is indoors.	It is lovely afternoon.	
MBR DECIDE Char United Char Char Char Char Char Char Char Char	hoose	PTIONING	Acat sitting on the floor The ground	SUMMARIZATION Aartciicles The cat is on the

MBR Decoding for Machine Translation

Many submissions to WMT'24 use MBR Decoding

	en→xx				
Models	$METRICX \downarrow$	XCOMET ↑	CometKiwi↑		
Baselines					
NLLB-54B	7.61 7	66.907	57.01 7		
GPT-40	1.50 6	83.74 6	77.04 5		
CLAUDE-SONNET-3.5	1.40 5	84.85 5	78.09 4		
DeepL	—	—	—		
TOWER					
TOWER-V2 7B	1.48 5	83.77 5	77.02 5		
Tower-v2 70B	1.32 4	84.87 4	78.29 4		
TOWER + QAD					
TOWER-V2 70B+MBR	0.92 2	88.78 2	81.393		
TOWER-V2 70B+TRR	1.03 3	87.95 3	82.13 2		
TOWER-V2 70B 2-step	0.891	89.25 1	82.54		



Rei et al., WMT 2024

MBR Decoding for Machine Translation

MBR Decoding is better than beam search



candidate list size (log scale)

Freitag et al., TACL 2022

MBR with Chain-of-Thought a.k.a. Self-Consistency



MBR for Distillation from Teacher LLM



MBR for Self-Distillation



Section 4: MBR Distillation

Why does MBR Decoding Work?



Procedure of Minimum Bayes Risk (MBR) Decoding (Kumar+ '04, Eikema+ '20)

- 1. Sample outputs randomly
- 2. Estimate the utility between the outputs using a function $u(\mathbf{h}, \mathbf{y})$
- 3. Select the output that maximizes the average utility to the others



Why does MBR Decoding Work?

MBR Decoding only need finite samples (e.g., 100) to surpass the performance of beam search (state-of-theart) whereas the number of possible sequences is infinite.



MBR decoding Beam

Search



Freitag et al., TACL 2022

Hypothesis: MBR decoding is a Convolutional Filter



https://suzyahyah.github.io/bayesian%20inference/machine%20translation/2022/02/15/mbr-decoding.html

Hypothesis: MBR decoding is a Convolutional Filter



https://suzyahyah.github.io/bayesian%20inference/machine%20translation/2022/02/15/mbr-decoding.html

So far I could not find evidence

- MBR decoding only using k-neighbors did not improve over MBR
- Neighbors of the MBR output do not always have high probability

Minimum Bayes Risk Decoding Minimizes Bayes Risk (Ichihara et al., ACL 2025)

Which objective functions are easier to optimize, MAP or MBR?

- With large enough number of samples, MBR is likely to be better

under assumptions



Minimum Bayes Risk Decoding Minimizes Bayes Risk (Ichihara et al., ACL 2025)

MBR decoding converges to the optimal solution with high probability at a rate of $O(1/\sqrt{n})$ where n is the number of samples under assumptions



candidate list size (log scale)

Freitag et al., TACL 2022

MBR Decoding as a Medoid Identification Problem (Jinnai&Ariu, Findings 2024)



MBR Decoding as a Medoid Identification Problem (Jinnai&Ariu, Findings 2024)



MBR Decoding as a Noisy Signal Decoding

Random sampling has no bias but high variance

Noise (variance) can be ignored by sample-and-aggregate strategy



Why does MBR Decoding Work?

MBR Decoding only need finite samples (e.g., 100) to surpass the performance of beam search (state-of-theart) whereas the number of possible sequences is infinite.

Still an open question!



Open Problems



MBR x Computational Linguistics

- Information density, surprisal, exposure bias...

- These theories are often mentioned to explain why MAP decoding fails
- Does it explain why MBR decoding is good?

e.g. The uniform information density hypothesis claims that human prefers to distribute information uniformly

How big is $[_{NP}$ the family $_{i} [_{RC} (\frac{\text{that}}{\text{that}}) \text{ you cook}$

for __i]]? Native English speakers don't want to omit this "that"

MBR for En-Ja, Ja-En Translation

- pfnet/plamo-2-translate is available for research!
 - But not for me... :(
- pfnet/plamo-2-translate-eval is a pairwise evaluation model which may be used in a different way than MBR decoding

Faster Inference: MBR Decoding is SLOW

MBR Decoding is 10-100 times more slower than beam search

Possible solutions

- Faster computation of utility function (e.g., Cheng&Vlachos, EMNLP 2023)
- Efficient use of autoregressive model
 - e.g. Speculative decoding (e.g., Sun et al., NeurIPS 2025)
- Efficient sampling
 - Non-iid sampling algorithm?

Summary

Questions: jinnai_yu@cyberagent.co.jp

Decoding Matters More than You Think!





Where Should I Start?

Starter kit for MBR decoding

- Sampling-Based Approximations to Minimum Bayes Risk Decoding for Neural Machine Translation (Eikema & Aziz, EMNLP 2022)
- High Quality Rather than High Model Probability: Minimum Bayes Risk Decoding with Neural Metrics (Freitag et al., TACL 2022)
- Minimum Bayes-Risk Decoding for Statistical Machine Translation (Kumar & Byrne, NAACL 2004)

Implementations (Library)

- <u>https://github.com/naist-nlp/mbrs</u>
- https://github.com/ZurichNLP/mbr

