

# Student Cup 2022

## 3rd place solution

Team : UZIA

# チーム紹介

チーム名 : UZIA (ユージア)

ユーザー名 : ktask(ケイタスク)、jt(ジェイティー)、gregley(グレグリー)

所属 : 会津大学大学院コンピュータ理工学研究科、会津大学コンピュータ理工学部

研究 : ヘルスケア分野の統計解析

# 解法

- ① KaggleのNLPコンペの解法を参考にBERT系のモデルをいろいろ試す
- ② 大量アンサンブル

- 前処理

htmlタグの除去

trainデータの重複データの除去

encoding errorの除去 (フォーラムの706番目のデータについてを解決) → KaggleのFeedback Prize 4th place solutionを参考

- モデル

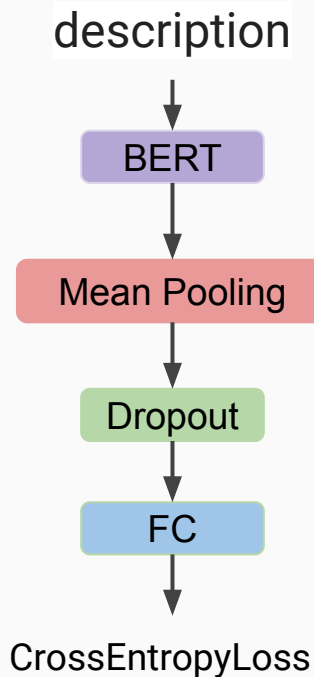
deberta-v3-base

deberta-v3-large

deberta-large

roberta-large

# モデルの構築



検証方法 : Stratified KFold (jobflag) fold=5

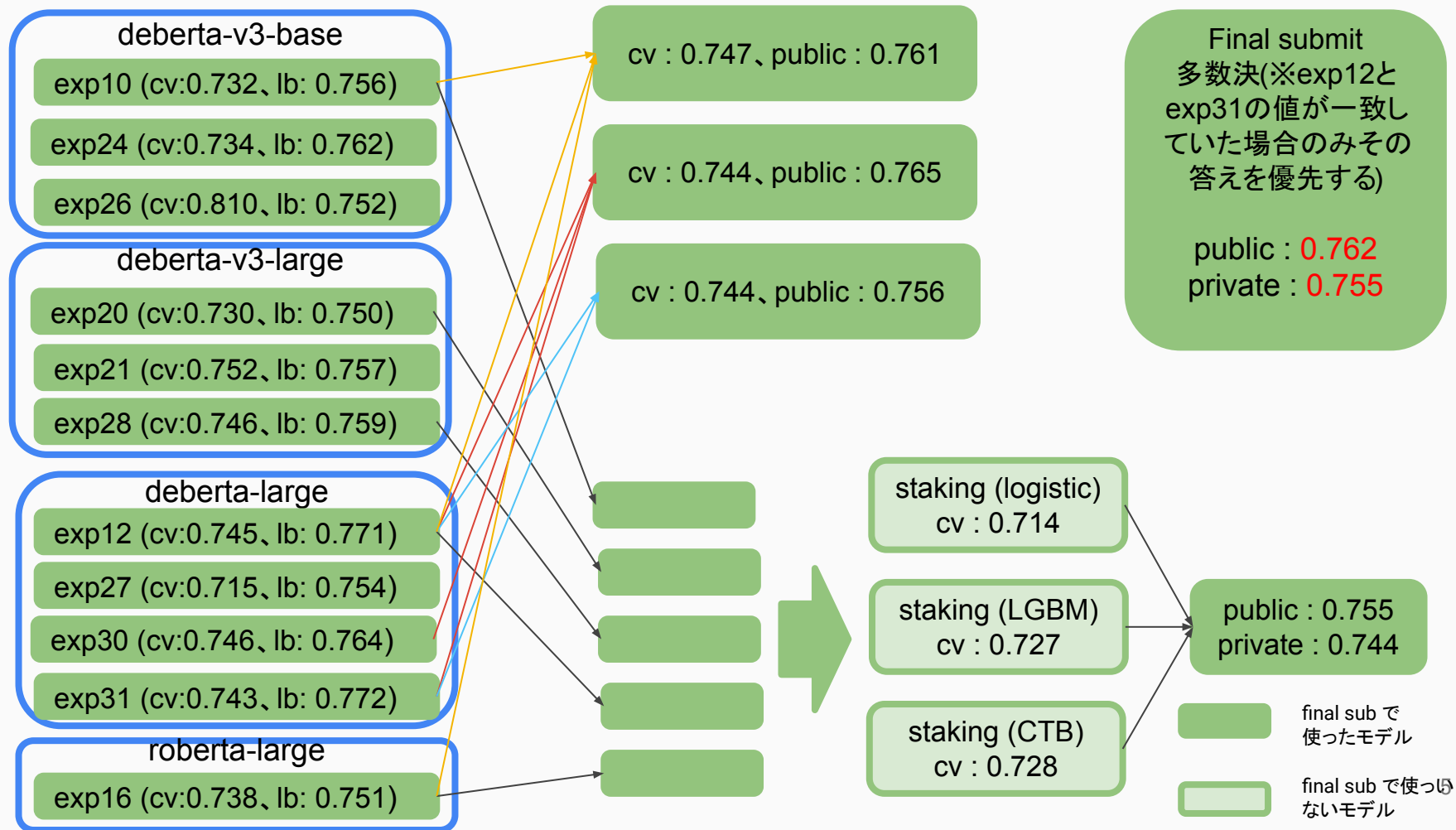
損失関数 : CrossEntropyLoss

optimizer : AdamW

epoch : 10

max\_len : 1024

# 全体像



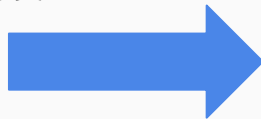
# 重複データを削除することで精度は変わるのか？

今回のデータセットにはtrainデータ内に12件、trainデータとtestデータに40件の重複データが含まれていた

全データ

	cv	public	private
roberta-large	0.738	0.751	0.739
deberta-v3-base	0.732	0.756	0.731
deberta-v3-large	0.752	0.757	0.743

deberta-v3-baseのみ  
精度が上がった



重複データを削除  
してもcvとlbの乖離  
は改善されなかつ  
た

重複削除

	cv	public	private
roberta-large	0.722	0.728	0.713
deberta-v3-base	0.734	0.762	0.747
deberta-v3-large	0.746	0.759	0.712

# その他取り組んだこと

- testデータによるデータの水増し

予測した各jobflagが0.9を超えるものを追加して学習を行う

→あまり精度の向上はしなかったがlbが高かったものは最後のアンサンブルに使用

- 逆翻訳によるデータの水増し

trainデータに対して英語→フランス語→英語の逆翻訳を行ったが精度が悪化した

運営の皆様、参加者の皆様  
ありがとうございました！