

大数据分析 - 罪犯减刑时长预测

DJH WDZ WGY WMH

【摘要】海量公开的司法数据，使得 AI 辅助解决法律领域问题成为一个有价值的研究课题。通过 AI 方法对减刑结果进行预测，也就是利用深度学习等相关技术，建立减刑预测模型，预测罪犯将获得的减刑时长，可以作为对司法过程的非人工监督，以促进执法公正的建设。本次选择 CCF BDCI 数据竞赛：罪犯减刑时长预测作为大数据分析选题，依次通过赛题分析、数据分析与预处理、算法设计、结果分析和优化思考五个步骤，对该选题进行了深入研究，研究结果达到 TOP10 左右，为该选题提供了有效的解决方案。

1. 赛题分析

1.1 赛题背景

我国为完善法制社会建设，贯彻落实依法治国的政治方针，将司法相关的文书都要在网上公开，以便受人民监督，彰显执法透明。海量公开的司法数据，使得 AI 辅助解决法律领域问题成为一个有价值的研究课题。

减刑、假释、暂予监外执行是我国重要的刑罚执行制度，也是司法实践中容易滋生腐败、产生执法司法不公的重点环节，党中央高度重视，社会普遍关注。

减刑是最常遇到的一种情况。而通过 AI 方法对减刑结果进行预测，可以作为对司法过程的非人工监督，以促进执法公正的建设。



图 CCF BDCI 数据竞赛：罪犯减刑时长预测

1.2 赛题任务

可以把该问题简单抽象为：根据判决书中的事实描述部分，利用深度学习等相关技术，建立减刑预测模型，预测罪犯将获得的减刑时长。

进一步抽象为：

- 输入：判决书中的事实描述部分
- 算法：利用深度学习，建立减刑预测模型
- 输出：预测罪犯将获得的减刑时长

该任务的一个实例如下：

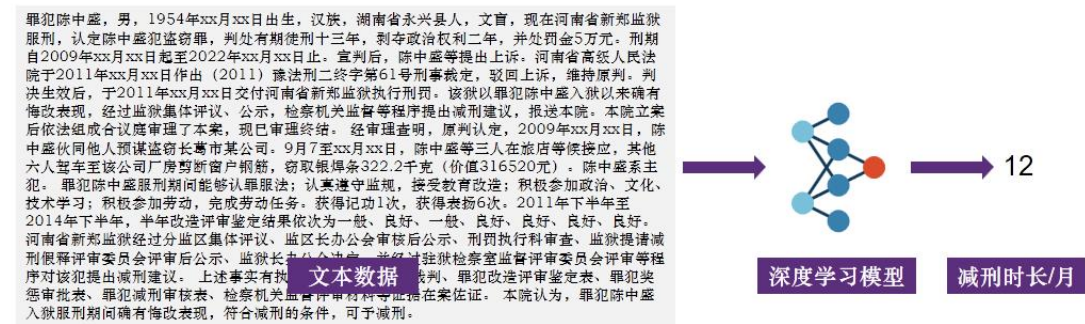


图 赛题任务实例

1.3 解决思路

为解决该问题，对上述流程做进一步具体化。也就是对一个 `string` 类型的数据。首先通过编码模型将其转化为计算机可以理解的数值型向量。然后通过预测模型计算一个减刑时长整数。具体如下：

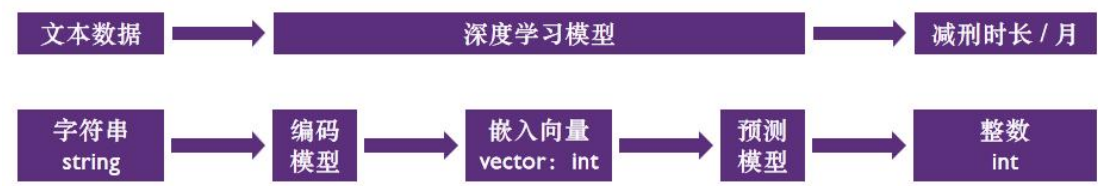


图 赛题任务解决思路图

其中，在具体的模型应用之前。需要首先进行评价指标选择，数据集划分，数据预处理，模型设计和模型训练工作。这些内容将依次在 1.4、1.5、2、3 节进行详细介绍。



图 问题解决流程图

1.4 评价指标选择

本赛题选择的评价指标包括相对准确率，绝对准确率和最终得分。设每一个样本的序号为 i ，样本由两部分构成，分别是用于案件描述的一个文本段和标准减刑月份这一数值。它们对应于后续模型的标准输入与参考输出。后续模型的输出为预测出的减刑月份这一数值。

1.4.1 相对准确率 Score

相对准确率 Score 根据预测出的减刑月份与案件标准减刑月份之间的差值距离作为评价指标。其计算方式如下：

1. 设预测出的减刑月份为 l_p ，标准答案为 l_a ，则：

$$v = |\log(l_p + 1) - \log(l_a + 1)|$$

2. 利用 v 换算每一个样本的得分 $Score_i$ ：

v	$v \leq 0.2$	$0.2 < v \leq 0.4$	$0.4 < v \leq 0.6$	$0.6 < v \leq 0.8$	$0.8 < v \leq 1.0$	else
$Score_i$	1	0.8	0.6	0.4	0.2	0

3. 最终相对准确率 Score 为：

$$Score = \frac{\sum_{i=1}^n Score_i}{n} \times 100\%$$

1.4.2 绝对准确率 ExtAcc

绝对准确率 ExtAcc 是指预测结果与真实值的绝对匹配准确率。其计算公式如下：

$$ExtAcc = \frac{\sum_{i=1}^n 1(if \text{预测值} = \text{真实值}) + 0(if \text{预测值} \neq \text{真实值})}{n} \times 100\%$$

1.4.3 最终得分 FinalScore

最终得分 FinalScore 是由 Score 和 ExtAcc 分别按 0.7 和 0.3 进行加权求和得到，作为一个预测结果的得分。其计算公式如下：

$$FinalScore = Score \times 0.7 + ExtAcc \times 0.3$$

1.5 数据集划分

1.5.1 数据来源

数据整理自全国减刑、假释、暂予监外执行信息网（<http://jxjs.court.gov.cn>）公开的减刑判决书文书（已对罪犯个人信息脱敏处理），均为真实案件减刑信息，数据集中包含罪犯个人信息、犯罪信息、在狱表现等字段。

表 数据字段解释

字段名	类型	取值范围	字段解释
id	string	-	案例唯一 id 标识
fact	string	-	案例描述（罪犯个人信息、犯罪经历、在狱表现等）
label	int	{0, 1, ...}	准许减刑的时长，单位为月

1.3.2 数据划分

赛题提供共计 75,000 条数据。按照以下形式加以划分和利用：

表 数据划分与利用

数据集	数据标识	数据量	数据利用
train	id_1 - id_50000	50000 条	数据本来无序，不需随机打乱 按照 8:2 划分训练集和验证集 用于模型训练
test	id_50001-id_75000	25000 条	用于模型测试 测试结果和其他打榜人员进行比较

2. 数据分析与预处理

2.1 数据分析

关键数据是一段文本，且为真实网站上的犯罪减刑文书，该数据没有缺失、重复、异常等情况，不需要整合、标准化、降维操作。

但是因为是一段长文本，所以会有一些不相关的文字信息，需要尽量提取关键信息，剔除无关信息，从而帮助模型更好的处理数据，提取特征，这可能需要对数据进行离散化操作。

文本按照一定的结构进行编写，包含罪犯个人信息、减刑原因、审理结果、在狱表现等部分。

id	fact	label
id	<p>【罪犯个人信息、犯罪经历】罪犯陈中盛，男，1954 年 xx 月 xx 日出生，汉族，湖南省永兴县人，文盲，现在河南省新郑监狱服刑，认定陈中盛犯盗窃罪，判处有期徒刑十三年，剥夺政治权利二年，并处罚金 5 万元。刑期自 2009 年 xx 月 xx 日起至 2022 年 xx 月 xx 日止。宣判后，陈中盛等提出上诉。河南省高级人民法院于 2011 年 xx 月 xx 日作出（2011）豫法刑二终字第 61 号刑事裁定，驳回上诉，维持原判。判决生效后，于 2011 年 xx 月 xx 日交付河南省新郑监狱执行刑罚。</p> <p>【减刑原因】该狱以罪犯陈中盛入狱以来确有悔改表现，经过监狱集体评议、公示，检察机关监督等程序提出减刑建议，报送本院。本院立案后依法组成合议庭审理了本案，现已审理终结。</p> <p>【审理结果】经审理查明，原判认定，2009 年 xx 月 xx 日，陈中盛伙同他人预谋盗窃长葛市某公司。9 月 7 至 xx 月 xx 日，陈中盛等三人在旅店等候接应，其他六人驾车至该公司厂房剪断窗户钢筋，窃取银焊条 322.2 千克（价值 316520 元）。陈中盛系主犯。</p> <p>【在狱表现】罪犯陈中盛服刑期间能够认罪服法；认真遵守监规，接受教育改造；积极参加政治、文化、技术学习；积极参加劳动，完成劳动任务。获得记功 1 次，获得表扬 6 次。2011 年下半年至 2014 年下半年，半年改造评审鉴定结果依次为一般、良好、一般、良好、良好、良好、良好。</p> <p>【减刑评审】河南省新郑监狱经过分监区集体评议、监区长办公会审核后公示、刑罚执行科审查、监狱提请减刑假释评审委员会评审后公示、监狱长办公会决定，并经过驻狱检察室监督评审委员会评审等程序对该犯提出减刑建议。上述事实有执行机关提供的生效裁判、罪犯改造评审鉴定表、罪犯奖惩审批表、罪犯减刑审核表、检察机关监督评审材料等证据在案佐证。本院认为，罪犯陈中盛入狱服刑期间确有悔改表现，符合减刑的条件，可予减刑。</p>	12

2.2 数据预处理

通过对数据文本的分析，发现可以分成两大部分处理，第一部分包括罪犯个人信息、减刑原因，第二部分则是剩下的法院的审查结果。之所以分成这两部分处理，是因为通过查询资料发现，比较多的关键信息都出现在第二部分，第一部分关键信息较少，所以分成两部分分别进行处理。

对于第一部分，因为关键信息比较少，直接构造规则，利用正则表达式，抽取罪名、原判刑期等信息。

对于第二部分，关键信息比较多，处理方式是利用 jieba 分词，对分词结果利用停用词表删除地名等无关的信息，形成最后的分词结果。

这是处理后的部分结果示例，可以看到文本大大减少，去除了很多不重要的信息，帮助模型能好的提取特征，从而提高性能。

```
text,label
强奸罪 有期徒刑 四年 执行 机关 康平 监狱 服刑 期间 认罪 悔罪 认真 遵守 法律法规 监规 接受 教育 改造 积极 参加 思想 文化 职业 技术 教育 积极 参加 劳动 努力完成 劳
诈骗罪 有期徒刑 十五年 执行 机关 黎塘 监狱 服刑 期间 认罪 悔罪 服从 管理 认真学习 积极 劳动 确有 悔改 表现 建议 予以 减刑 有期徒刑 八个 检察机关 报请 减刑 无异议
贩卖毒品罪 有期徒刑 三年 服刑 期间 认罪服法 服从 管教 认真 遵守 法律法规 监规 接受 教育 改造 积极 参加 思想 文化 职业 技术 教育 积极 参加 劳动 努力完成 劳动改造
受贿罪 有期徒刑 三年 执行 机关 广西 女子监狱 服刑 期间 认罪 悔罪 服从 管理 认真学习 积极 劳动 确有 悔改 表现 建议 减刑 有期徒刑 二个月 符合 减刑 法定条件 请求 法
抢劫 盗窃罪 有期徒刑 九年 入监 认罪服法 遵守 监规 劳动纪律 积极 参加 劳动 认真学习 政治 文化 技术 积极 缴纳 罚金 1500 元 本院认为 该犯 符合 减刑 法律.6
故意伤害罪 有期徒刑 六年 执行 机关 提出 服刑 期间 认罪 悔罪 服从 管理 认真学习 积极 劳动 确有 悔改 表现 符合 减刑 条件 检察机关 报请 减刑 符合 法律 条件 报请 程序
贩卖毒品罪 有期徒刑 九年 服刑 期间 认罪服法 服从 管教 认真 遵守 法律法规 监规 接受 教育 改造 积极 参加 思想 文化 职业 技术 教育 积极 参加 劳动 努力完成 劳动改造
盗窃罪 有期徒刑 六年 执行 期间 确有 悔改 表现 执行 机关 呈报 减刑 材料 符合 法律 相关 本院认为 符合 法定 减刑 条件 执行 机关 呈报 符合 法定程序 改造 表现 原 犯罪
抢劫罪 有期徒刑 十年 执行 机关 监狱 服刑 期间 认罪服法 接受 改造 确有 悔改 表现 建议 犯 减刑 有期徒刑 一年 四个 服刑 期间 遵守 监规 认真学习 劳动 任务 确有 悔改 表
非法制造枪支罪 故意伤害罪 聚众斗殴罪 有期徒刑 五年 执行 机关 弄窝 监狱 服刑 改造 期间 认罪 悔罪 积极 改造 确有 悔改 表现 建议 减刑 人民检察院 执行 机关报 请 减刑
抢劫罪 有期徒刑 十一年 执行 机关 监狱 减刑 建议书 中 提出 入监 认罪 悔罪 积极 接受 教育 改造 积极 参加 劳动 累计 受 监狱 表扬 奖励 四次 嘉奖 奖励 提供 相关 证据 材
强奸罪 有期徒刑 七年 执行 机关 曲沃 监狱 服刑 改造 中 认罪 悔罪 认真 遵守 法律法规 监规 积极 参加 思想 文化 职业 技术 教育 劳动 中 服从 管理 严格遵守 操作规程 保
组织 领导黑社会性质组织罪 故意伤害罪 寻衅滋事罪 开设赌场罪 有期徒刑 十七年 服刑 改造 期间 认罪 悔罪 服从 民警 管理 教育 认真 遵守 监规 纪律 接受 教育 改造 积极
强奸罪 有期徒刑 四年 执行 机关 监狱 减刑 建议书 服刑 期间 认罪服法 遵守 监规 积极 劳动 参加 三课 学习 9 1 获得 计分 考核 积分 转换 表扬 三个 服刑 期间 确有 悔改 表
强奸罪 有期徒刑 十三年 执行 机关 监狱 刑罚 执行 期间 确有 悔改 表现 建议 减刑 有期徒刑 五个 剥夺 政治权利 改为 二年 人民检察院 执行 机关报 请 减刑 建议 无异议 服
贩卖毒品罪 有期徒刑 九年 执行 机关 提出 服刑 期间 认罪 悔罪 服从 管理 认真学习 积极 劳动 确有 悔改 表现 符合 减刑 条件 检察机关 提请 减刑 符合 法律 条件 提请 程序
盗窃罪 有期徒刑 五年 执行 机关 服刑 期间 认罪服法 积极 参加 思想 文化 技术 教育 良好 成绩 积极 参加 劳动 任务 确有 悔改 表现 建议 减刑 入监 认罪服法 认真 遵守 监
故意伤害罪 无期徒刑 执行 机关 监狱 提出 服刑 期间 确能 认罪 悔罪 严格遵守 法律 法规 监规 接受 教育 改造 积极 参加 思想 文化 职业 技术 教育 学习成绩 合格 积极 参加
强奸罪 有期徒刑 四年 执行 机关 提出 服刑 期间 认罪 悔罪 服从 管理 认真学习 积极 劳动 确有 悔改 表现 符合 减刑 条件 检察机关 报请 减刑 符合 法律 条件 报请 程序 符
抢劫罪 有期徒刑 二年 服刑 期间 认罪 悔罪 服从 管教 认真 遵守 法律法规 监规 接受 教育 改造 积极 参加 思想 文化 职业 技术 教育 积极 参加 劳动 努力完成 劳动改造 任务
抢劫罪 有期徒刑 十一年 执行 机关 刑罚 执行 期间 确有 悔改 表现 建议 犯 减刑 刑罚 执行 期间 认罪 悔罪 认真 遵守 法律法规 监规 接受 教育 改造 积极 参加 思想 文化 职
诈骗罪 有期徒刑 四年 执行 机关 提出 服刑 期间 认罪服法 服从 管理 教育 遵守 监规 纪律 积极 参加 三课 文化 学习 上课 专心 听讲 作业 劳动 中能 吃苦耐劳 踏实 肯干 超
故意杀人罪 死刑 执行 机关 提出 服刑 期间 确有 悔改 表现 以此为由 犯 提出 减刑 八个 认罪 悔罪 接受 改造 参加 劳动 学习 认真 卫生 整洁 上次 减刑 执行 机关 表扬 5 次
```

3. 算法设计

3.1 模型设计

首先进行问题定义，可以认为减刑时长预测问题既可定义为回归问题又可定义为多分类问题，但经过一些简单的实验，可以发现回归问题的效果稍差，因此最终将该问题作为多分类问题。

模型选择上，基于离散特征和连续特征进行两种建模，对于离散特征，基于上述数据预处理获得的离散值加规则模版匹配，然后将其输入到机器学习模型中。对于连续特征，对文本进行一些简单的预处理，如去除标点符号和罪犯个人信息等无关特征，然后直接输入深度学习模型进行训练。

3.2 模型介绍

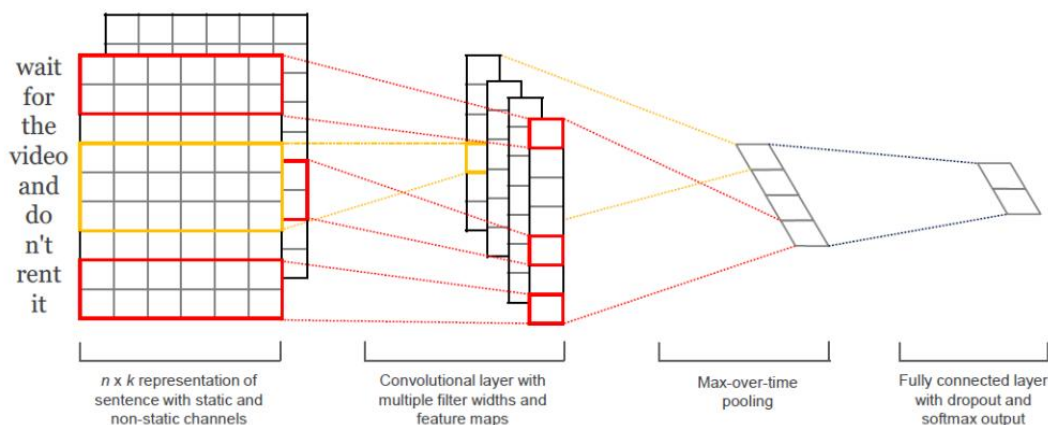
3.2.1 多分类器

模型简称	模型名称	模型原理	模型优势	模型劣势
Random Forest	随机森林	在训练时构建大量决策树,在决策时输出大多数树选择的分类	纠正了决策树过度拟合的问题	准确性低于梯度提升树,受数据特征影响大
XGBoost	极端梯度提升	并行树提升,将叶节点按比例收缩、引入额外的随机化参数	模型复杂度可控,防止过拟合	调参复杂,且对于高维数据的特征捕获能力弱
SVM	支持向量机	使用广义最优超平面对向量集合进行分割	多用于解决高维特征的分类问题	计算量大,对脏数据敏感
KNN	K 近邻	使用与该样本最接近的K个邻居来决定样本的取值	思想简单,对异常值不敏感	计算量大,分类速度慢
MLP	多层感知机	多层全连接神经网络	课模拟任何复杂分类函数	参数量大,容易过拟合

3.2.2 CNN——TextCNN

与 CNN 相比,TextCNN 在网络结构的改变不大,甚至可以说更加简单,引入词向量使得该方法在文本分类任务中拥有比较好的效果。

在本次实验中,预训练词向量选择利用搜狐新闻训练的中文词向量,嵌入长度为 300。网络结构是将输入数据分别经过三种大小的卷积核(2, 3, 4),每个卷积核数量为 32,后接最大池化向量。实验中学习率为 0.005, batch_size 为 256, 优化器为 Adam。



3.2.3 RNN——BERT、RoBERTa、Lawformer

TextCNN 等简单网络的准确率结果不太令人满意，调整网络结构也需要很大的设计和训练成本。因此考虑两个可能的优化方向：构建更复杂、更深的网络结构，增加参数量，提高网络的拟合能力；使用更前沿的模型，如加入注意力机制，使用 Transformer 结构等。综合这两点，决定使用 BERT 等预训练模型。它们的优势在于，可提供不同大小的预训练模型和预训练分词器，只需要在其基础上进行 fine-tune，就免去了调节网络结构，以及大量从头训练的时间。

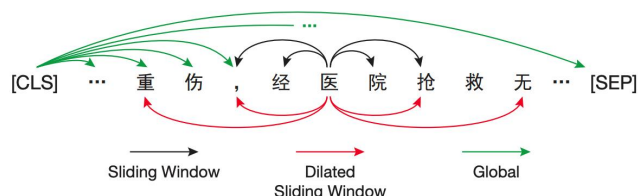
● BERT、RoBERTa

BERT 模型是谷歌提出的基于双向 Transformer 构建的预训练模型，有双向编码的能力。并且它天然支持文本分类任务，只需将预处理后的文本构造为 BERT 的输入形式，最后再添加一个额外的输出层对应到多个分类即可。

使用 BERT 模型的中文版本，共 12 层 Transformer 编码器，768 个隐单元，12 个注意力 head，约 110M 的参数。学习率为 $5e-5$ ，序列长度为 512，batch_size 为 32。由于进行的是 fine-tune，因此训练五轮左右即可达到很好的效果，继续训练可能发生拟合。当然，由于模型很大，训练过程十分缓慢。

RoBERTa 是 BERT 模型的扩展，是一个优化后更加鲁棒的 BERT 版本。它用更大的 batch 和更长的序列进行训练，通过大批量的数据提高任务的准确性。

● Lawformer



Lawformer 是基于 Longformer 训练的预训练模型，模型在大规模的中文法律长文本上预训练得到，基于 Longformer，可以处理上千个 token。模型特点在于没有采用标准的 self-attention，而是结合局部滑窗和全局 attention 机制来捕获长程依赖。

在实际实现上，基于 RoBERTa-wwm-ext 的 checkpoint 训练，以 MLM 为任务。学习率为 $5e-5$ ，序列长度 4096，batch_size 为 32。为充分利用长序列，将一些短文书拼接。优化器为 Adam，模型采用 8 块 32G V100 训练。

理论上 Lawformer 对于中文法律文本，尤其是长文本会有较好的效果。

3.3 模型训练

具体的训练过程分为数据读取、数据划分、数据预处理、模型训练和调参、模型测试几个步骤。将训练集和验证集按 8:2 的比例进行划分。训练时还使用了网格化搜索和五折交叉验证。最终使用在验证集上表现最好的模型进行测试。

4. 结果分析

4.1 测试结果

测试结果主要包括两部分，第一部分为传统机器学习模型的测试结果，第二部分为深度学习模型的测试结果。

传统机器学习的测试结果如下，包括 XGBoost、随机森林、支持向量机、KNN 与朴素 MLP 算法，其中算法性能较为相近，并且各有优劣，综合效果为 XGBoost 最佳。

模型	最终得分 FinalScore	相对准确率 Score	绝对准确率 ExtAcc
XGBoost	0.7587	0.8671	0.5058
Random Forest	0.7460	0.8812	0.4305
SVM with RBF kernel	0.7525	0.8492	0.5269
KNN	0.7372	0.8406	0.4959
MLP	0.7492	0.8413	0.5343

基于深度学习的模型主要有以下几种，包括文本 CNN、BERT、RoBERTa 与 Lawformer。其中 TextCNN 与第二项 BERT 使用了特殊的分词处理模块，如表所示，综合效果为 BERT 最优，而且加了预处理分词模块后，效果提升了很多，其中 Lawformer 与 BERT 的表现差异不太大，这可能是由于 Lawformer 除了专注于法律文本训练，还特别对长文本的输出有优化，而本题的文本输出大部分都属于短文本，不能很好的发挥作用，而 RoBERTa 的表现相对较差，这一点略微令我们奇怪，不过表现的轻微波动也属正常。

模型	最终得分 FinalScore	相对准确率 Score	绝对准确率 ExtAcc
TextCNN	0.4631	0.5450	0.2718
BERT	0.7531	0.8854	0.4446
RoBERTa	0.7367	0.8794	0.4040
Lawformer	0.7524	0.8906	0.4300
Bert+预处理分词	0.7909	0.9098	0.5134

以上的模型硬件环境为 NVIDIA A100 GPU，训练时每个模型训练 30 个 epoch，取验证集表现最好的一个，BERT 等模型初始学习率为 5e-5。每个 epoch 的训练时长大约为 500 秒。

4.2 竞赛排名

竞赛排名见下表，结果为预估值，排名 10 名左右。需要说明的是，由于时间原因，2022年12月9日后无法提交结果，结果是在自行分割的测试集预估的，并不是实际的测试集，因此结果可能有偏差。

排名	最终得分 FinalScore	相对准确率 Score	绝对准确率 ExtAcc	提交次数
1	0.8433	0.9487	0.5975	32
2	0.8344	0.9420	0.5834	35
3	0.8343	0.9495	0.5655	23
4	0.8168	0.9349	0.5412	5

5	0.8163	0.9341	0.5413	42
6	0.8161	0.9381	0.5314	57
7	0.8139	0.9317	0.5392	23
8	0.8131	0.9378	0.5221	8
9	0.8109	0.9331	0.5258	29
10	0.8040	0.9289	0.5124	13
Our	0.7909	0.9098	5134	3

5. 优化思考

5.1 存在问题

最直观的问题是：最终得分仍有提升空间。一方面，和最优 TOP1 差 5% 百分点，现在是 79.09%。另一方面，TOP1 达不到 100%，现在是 84.33%。

5.2 优化方向

针对上述问题：未来，拟在以下方面考虑做进一步处理与优化。

1. 数据预处理：更细致的数据预处理，当前的数据预处理相对粗糙。希望未来将不同信息的表现形式，如上文提到的离散化信息和连续化信息做对齐统一。
2. 知识增强：引入法律领域知识，根据专业知识增加数据描述词语，扩大信息量。如通过引入法律词典库优化文本特征抽取效果，引入法律知识图谱和图谱推理技术对减刑时长进行推理。
3. 模型优化：对模型方法进一步改进。如对模型输入做调整，通过特征工程等方式优化模型效果。
4. 模型融合：使用投票模型支持多个模型同时预测，聚合模型能力。

6. 附录

6.1 成员分工

序号	姓名	专业	算法/模型分工	项目分工
Leader	DJH	计算机	BERT、多分类器	研究框架、赛题分析、优化思考、汇总
2	WDZ	计算机	TextCNN	数据分析与预处理
3	WGY	计算机	BERT、RoBERTa	算法设计
4	WMH	GIX 全球创新学院	LawFormer	结果分析

6.2 代码仓库

全部代码在 github 上进行了开源，具体访问地址如下：

- 远程仓库的名称 greengrasscugb/BDA_PredictionOfCommutationTimeForOffenders
- 远程仓库的地址
https://github.com/greengrasscugb/BDA_PredictionOfCommutationTimeForOffenders.git
- 不同成员负责的代码在不同分支中，分支命名选择了成员的姓名首字母缩写，如DJH