

千帆ModelBuilder

开启大模型卓越之门：模型优化的关键钥匙

千帆大模型平台产品经理-徐星雨

千帆大模型平台算法-江锦

目录

01	大模型的能力及问题	03
02	优化大模型的关键步骤	08
03	大模型优化的流程-以作文批改场景为例	17
04	Q&A	20

一、大模型的能力及问题

大模型的发展



大模型技术快速发展，百度推出大模型开发平台ModelBuilder，支撑多来源、多模态大模型的开发



大模型的能力

大模型具备强大的生成能力，可应用于各行各业的应用场景中

知识覆盖广

大模型通过在海量数据上训练，积累了广泛的知识，可以应用于多种场景中

生成质量高

大模型能够根据积累的知识、历史问答信息等生成高质量的内容

自动化

大模型能够自动化的完成大量重复性的内容生成工作，提高效率

个性化

大模型可以根据指令要求，生成个性化的内容，提升体验

渗透各 行各业



教育



电商



医疗



文娱



办公



金融

大模型的局限性



然而，大模型可能在一些特定的任务类型或领域上输出质量不高，需要进一步提升



特定任务适应性不足

大模型在某些应用场景中表现差，不能根据指定要求输出

例如：

合同信息提取的场景中，大模型不能按照要求抽取关键信息



缺乏领域知识

在一些垂直领域下，大模型未学习过相关知识

例如：

询问大模型医疗相关知识，大模型回复不清楚相关内容



回复不按照指定要求

大模型不遵循指令要求，如输出格式存在问题

例如：

Prompt:

请列出2个关于环境保护的具体措施。

Response:

1. 我喜欢绿色的植物。
2. 天气不错，适合户外活动。
3. 科学技术可以帮助我们。

借助千帆ModelBuilder精调提升模型输出效果



千帆大模型平台提供从数据管理、精调工具、评估工具等提供模型优化的工具

精调工具链	数据管理	训练模式	训练能力	模型评估	平台预置
	数据洞察	Post-pretrain	表单式开发	自动评估	预置训练集
	数据清洗	SFT	自定义开发 (敬请期待)	人工评估	预置评估集
	数据增强	偏好对齐-DPO	通用、垂直混合语料	单个评估	精调样板间
	智能标注	偏好对齐-KTO	增量训练	对比评估	
	数据回流	偏好对齐-RLHF	闲时资源训练 (限免)	基线评估 (敬请期待)	

精调模型	ERNIE 模型				开源模型	
	ERNIE 4.0 Turbo	ERNIE 3.5	ERNIE Speed Pro	ERNIE Lite Pro	ERNIE Speed	Baichuan
			ERNIE Lite	ERNIE Tiny	ERNIE Character	LLaVA
						...

二、优化大模型的关键步骤

精调的价值



Prompt调优后大模型输出效果依旧不满足要求，可进一步准备问答数据发起精调

Prompt调优

- ❖ **适用场景：**大模型输出效果差，并且不需要添加新知识时。
- ❖ **方法：**
 - 精准描述需求，增加规则
 - 添加Fewshot示例
 -

调优后的Prompt可
用于进一步精调

依旧效
果差？



模型精调

- ❖ **适用：**经过调整Prompt后模型输出效果差，或者需要引入新知识。
- ❖ **可解决的问题：**
 - 纠正大模型输出的格式、口吻
 - 场景比较复杂，Prompt难约束
 - 需要大模型处理一些边界Case
 - 处理垂直领域/行业知识
 - 解决小参数模型的输出差的问题，比肩大参数模型，降低调用成本

模型训练矩阵：模型生态丰富，支持文生文及多模态



支持25+大语言模型，并支撑图像生成和图像理解业务场景



ERNIE包括不同长度及参数规模的大模型，如何选择



- 效果：大模型在任务中回复的准确性和性能表现
- 成本：训练以及调用大模型所需的资源
- 时延：模型从接收输入到回复过程中产生的时延

定位	模型名称	上下文长度	效果/成本/时延
旗舰大模型	ERNIE 4.0 Turbo ERNIE 3.5	8K 8K	效果最好、成本最高、时延最长
主力大模型	ERNIE Speed ERNIE Lite	8K、128K 8K、128K	效果、成本、时延均衡 成本相对较低时延较短
轻量级大模型	ERNIE Tiny	8K、128K	成本最低时延最短
垂直场景模型	ERNIE Character	8K	角色扮演场景表现好 成本、时延均衡

如何选择精调方法？



根据实际的业务场景分布准备数据，选择合适的训练方法

	数据类型	用途	适合的业务场景举例
Post-pretrain	❖ 纯文本数据	❖ 需要大量数据，用于行业或者领域大模型的训练	❖ 训练金融、医疗垂直大模型-知识问答
SFT	❖ 单轮/多轮对话-非排序类数据	❖ 调整模型输出的格式、口吻，处理，解决边界Case等	❖ 角色扮演、对话摘要总结、判题打分等
偏好对齐	❖ 单轮/多轮对话-存在偏好排序的数据	❖ 对模型的输出有排序/倾向时，可将接受/拒绝的信息训入模型	❖ 智能客服-分析用户反馈和偏好

如何配置训练数据？



在业务数据中混入平台预置的通用和垂直语料，提升业务场景效果、避免通用知识遗忘

“

SFT训练中建议使用1000条语料起；Post-pretrain中建议10亿tokens语料起。
具体需要根据场景的复杂度和分布情况调整

”



如何配置合适的参数？



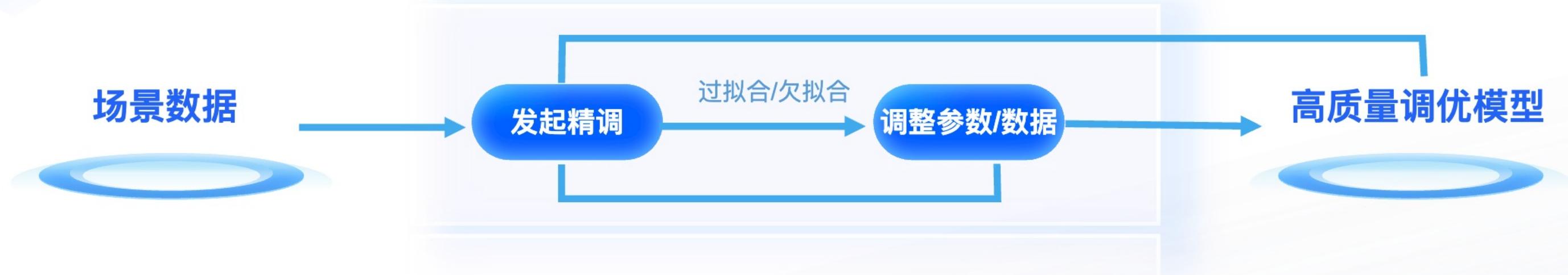
配置合适的训练参数是决定模型效果的关键因素，需要根据数据和训练的效果调整参数

参数名	含义/作用	过大的影响	过小的影响	推荐
迭代轮次-Epoch	一个Epoch是使用全部的训练数据集训练一个完整的迭代。Epoch增加可以加深模型对语言和语义关系的理解	Epoch太大可能会导致过拟合	Epoch太低可能会导致模型欠拟合，LLM不能从训练数据中学习到正确的Weights和biases	建议设置在1-5之间，小数据集可增大Epoch以促进模型收敛
学习率-Learning Rate	Learning Rate是控制LLM训练过程中根据Loss Function学习速度快慢的参数	训练过程中配置大的学习率可能会导致训练不稳定或者过拟合	使得学习过于缓慢	3e-4，模型收敛慢可增大值
批处理大小-Batch Size	Batch Size决定了每次训练时数据量的大小	大的Batch Size可以加速训练的过程，但可能导致过拟合	使得训练过于缓慢	建议10-20之间
序列长度-Seqlength	单条数据的最大长度，包括输入和输出。如果数据集中的文本普遍较短，建议选择较短的序列长度以提高计算效率	训练速度慢	比序列长度长的数据会被丢弃	建议根据最长训练语料的长度设置
Checkpoint保存个数	训练过程最终要保存的Checkpoint个数。训练完成后可以保存多个Checkpoint，选择表现好的Checkpoint发布模型	保存Checkpoint会增加训练时长	-	建议1-Epoch之间。具体可根据Step的数量调整
Checkpoint保存间隔数	训练过程中保存Checkpoint的间隔Step数	间隔太长则可能在故障时恢复较慢	间隔太短可能导致频繁的Checkpoint操作增加训练时长	建议50-100之间
验证步数-Validation Step	计算验证集指标的间隔步数，为0时不开启验证，没有相关指标	算验证集指标间隔步数大，指标打印少	训练时长增加	建议10-20之间
早停-Early Stop	监控精调任务的指标变化情况，指标连续不变则提前终止训练	-	-	-

如何根据评估指标判断训练的效果？



根据训练报告中的指标评估训练效果，并根据指标进一步调整数据或者参数训练



以Loss指标为例，根据训练过程指标分析训练情况

平滑度

平滑下降说明模型的效果是稳定提升的，出现抖动说明可能数据存在异常

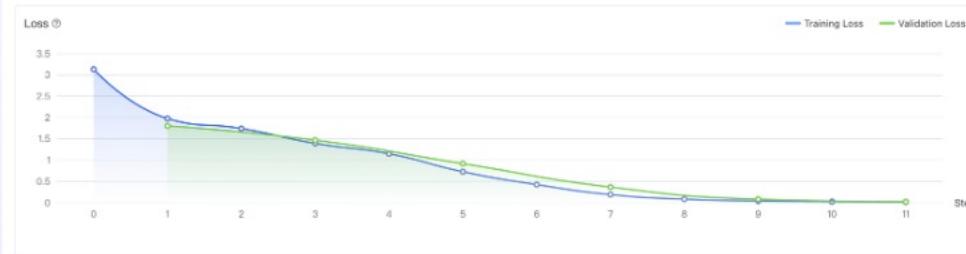
收敛性

曲线收敛到一个点，训练增加时指标不会增加/减小

泛化性

不仅在训练集上有好的效果，在验证集上表现同样优异

大模型训练效果分析（分析TraningLoss和ValidationLoss）



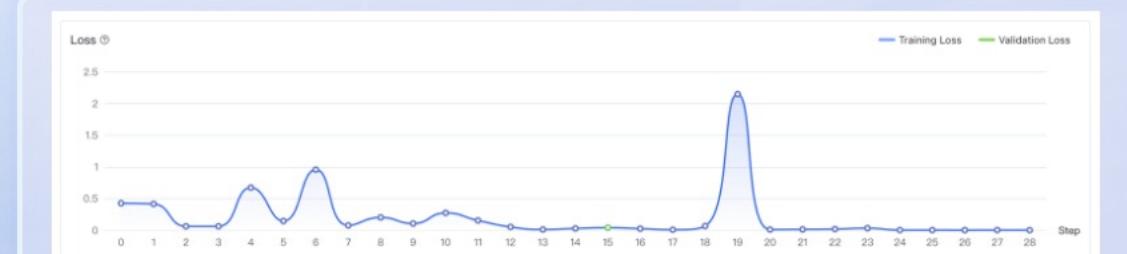
Case1:两条曲线同时收敛，可进一步评估模型效果



Case2:两条曲线同时在下降，可进一步调整参数/数据训练



Case3:训练集指标下降，验证集指标上升，过拟合



Case4:曲线出现抖动，说明存在异常数据

精调模型效果衰减怎么办？



模型上线后，如果发现在实际业务场景中的效果发生变化，可进一步增量训练模型



线上业务分布变更

如果随着业务的变化，实际分布发生了变化，可以进一步增加场景数据增量训练



新增业务场景

如果需要新增部分业务场景的回复，可以准备场景的数据集，进一步增量训练



出现Badcase

可以收集线上的Badcase或者需要进一步加强的场景数据，进一步训练

增量训练

基于已经精调的模型继续训练，强化在某些场景的效果

增量训练:



* 选择基准模型:

SFT > cha1120 > V1



bleu-4=98.60 edit-distance=0.22 embedding-distance=0.01 rouge-1=98.61 rouge-2=98.15 rouge-l=99.31

数据来源:

平台数据集 对象存储BOS 分布式存储AFS

数据格式:

Prompt+Response Role (user+assistant)

* 选择数据集:

预置数据集>小说... × 预置数据集>用户... ×

预置数据集>劳动... ×

创建数据集

▲ 数据配比设置

数据集	字符数	样本数	采样率	数据占比
用户情感分析评估数据集 V1	283873	500	0.10	25.38%
小说人物角色扮演V2评估数据集 V1	36031	30	0.10	3.22%

防止业务效果损失，可以在增量训练时，增加基准模型训练时的部分数据

三、大模型优化的流程-以作文批改场景为例

模型训练：精调样板间，指引精调场景落地



以行业领域落地场景为例，从数据准备、模型选择、参数调整等方面提供精调样板，协助您体验调优过程和效果

千帆ModelBuilder

概览 模型广场 体验中心

Prompt工程 模型服务 应用接入 在线推理

批量推理 调用统计 模型调优 我的模型

模型精调 精调样板 Post-pretrain

SFT 偏好对齐 模型评估 模型压缩

数据管理

粘贴件

精调样板

请输入样板名称

小说人物角色扮演V2
社交文娱 KTO
通过调用如ERNIE 4.0旗舰版模型产生的问答对，经过筛选得到有效的训练数据。在降低数据标注成本的同时，可通过精调得到成本更低、性能更优、且特...

劳动合同关键信息提取
法律 SFT
在企业和法律事务中，大模型可以扮演合同审查专家的角色，自动识别劳动合同是否包含所有必备条款并提取合同中的关键信息。通过模型精调，可以...

客服对话多标签生成
电商营销 SFT
在客服对话场景中，可以通过大模型分析用户与客服之间的对话信息，准确识别用户的意图和对应原因，生成对应标签为后续回复和营销策略服务。

文本创作字数控制
交互助手 SFT
文本创作字数控制场景中，大模型可以扮演高效的写作助手的角色。其具备出色的文本理解、生成及编辑能力，通过精确设定的评分标准和模型精调，大模...

数学判题
在线教育 SFT
在线教育中，判定学生数学答案会考察解题步骤、方法和答案的准确性。大模型同样可以扮演数学题目评判专家的角色。通过精确设定评分标准和模型的训...

购物平台客服对话摘要
电商营销 SFT
购物平台中，顾客通常会咨询各种购物相关的问题。我们可以借助大模型的生成能力抽取顾客咨询问题的类别、顾客对客户回答的满意度等等，进而做进一...

英语口语练习
在线教育 SFT

作文自动点评或批改
在线教育 SFT

精调样板-高考作文批改



以作文自动批改场景为例，从准备高考作文及点评数据、生成数据，选择模型及训练方法，调整参数、评估模型效果等方面讲解精调的全流程



四、Q&A

课后作业



使用上节课准备的数据集（未准备可使用预置数据集或者其他数据），并选择ERNIE 系列模型发起SFT精调。

作业发布要求：

- 1、说明SFT精调的思路（如调参）
- 2、展示精调的报告详情（含Loss曲线等）

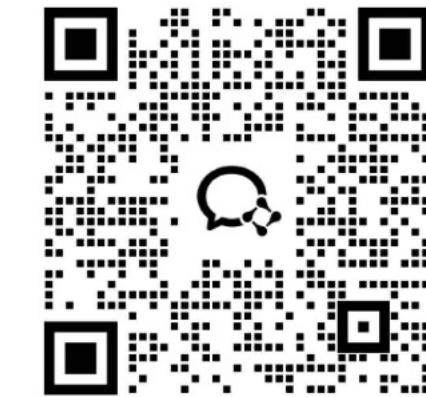
Tips:

- 1、在资源配比中选择闲时调度，可免费精调

*将上述作业在「百度智能云千帆社区」进行发布，发布时选择“千帆大模型训练营”话题



扫码进入
百度智能云千帆社区



扫码进入课程群

谢谢

百度智能云千帆ModelBuilder