

本课件仅用于教学使用。未经许可,任何单位、组织和个人不得将课件用于该课程教学之外的用途(包括但不限于盈利等),也不得上传至可公开访问的网络环境

-

数据科学导论 Introduction to Data Science

第二章 数据分析

黄振亚, 陈恩红

Email: cheneh@ustc.edu.cn, huangzhy@ustc.edu.cn

课程主页:

http://staff.ustc.edu.cn/~huangzhy/Course/DS2025.html

An Introduction to Data Science

9/28/2025



数据分析基础

6

□数据采集

□数据预处理

□特征工程

Data Collection

Data Preprocessing

Feature Engineering

数据采集

2

□ 无时无刻产生数据,获得数据的方式多种多样





测量

网页



数据库



监控



传统媒体

9/28/2025



数据采集

Δ

- □数据检索与开源
- □批量数据获取
 - □网络爬虫
- □数据生成



数据采集:数据检索

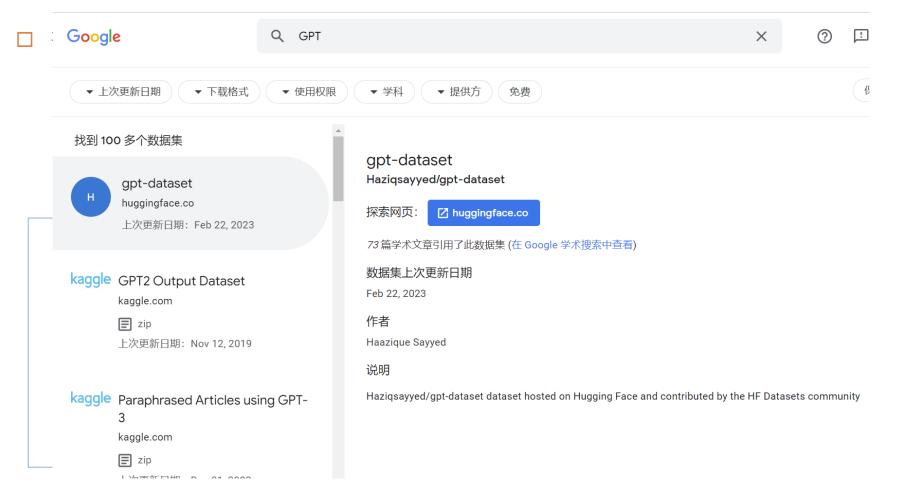
- □最简单、最灵活的数据获取方式就是依靠检索
- □搜索引擎



数据采集:数据检索

6

」最简单、最灵活的数据获取方式就是依靠检索



数据采集:公开数据

- □代表性公开数据集
 - □ ImageNet: http://www.image-net.org/
 - □ Amazon从2008年开始就为开发者提供几十TB的开发数据
 - http://aws.amazon.com/datasets
 - □ YouTube视频与社交网络: http://netsg.cs.sfu.ca/youtubedata/
 - □ 用户评分MovieLens: https://grouplens.org/datasets/movielens/
 - □ 网络数据-Large scale network: https://snap.stanford.edu/data/
 - □ 教育数据:
 - ASSISTmentsData: https://sites.google.com/site/assistmentsdata/home/
 - BASEGroup: https://github.com/bigdata-ustc/EduData
 - □ 阿里天池数据-数据平台: https://tianchi.aliyun.com/dataset/
 - □ 公开大数据竞赛的数据: KDDCup, NeurIPS Challenge



数据采集: 批量数据获取

- □ 大量数据的获取难以手动实现, 需借助**爬虫程序**
 - □ 也有可能通过交易(购买)"数据"而得
- □ 网络爬虫是一个自动在网上抓取数据的程序
 - □ 爬虫本质上就是**下载**特定网站网页的HTML/JSON/XML数据 ,并对数据进行**解析、提取与存储**
 - □ 通常先定义一组入口URL,根据页面中的其他URL,深度优 先或广度优先的遍历访问,逐一抓取数据





数据采集: 网络爬虫

7

- □ 网络爬虫是什么?
 - □ 网络爬虫(又被称为网页蜘蛛,网络机器人,网页追逐者) ,是一种按照一定的规则,自动的抓取万维网信息的程序或 者脚本。
 - ■请求网站并提取数据的自动化程序
 - □爬虫的行为可以划分为:载入、解析、存储,

■最复杂的部分为载入



载入

发起请求(Request) 获得响应(Response)-类型多样

获得数据 存储数据-类型多样

存储

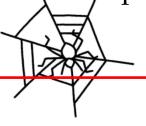
解析

解析内容 提取数据-类型多样



10

- □访问网页示例
 - □ 网站数据主要依托于网页(html, 超文本标记语言)展示
 - □ 用户Request服务器,服务器response信息(html等)



存储数据



浏览器解析 展示数据



发起请求(Request)



获得响应(Response) 网页数据



11

- □ 网页示例
 - □ 网站数据主要依托于网页(html, 超文本标记语言)展示
 - □ 用户Request服务器,服务器response信息(html等)
 - □ 课程主页: http://staff.ustc.edu.cn/~huangzhy/Course/DS2025.html
 - □查看网页源代码

Introduction to Data Science 数据科学导论

课程代码: CS4023.01

学院: 011计算机科学与技术学院

课程教师: 黄振亚, 陈恩红

上课时间: (1-18周)每周二下午第8、9节: 15:55-17:30, 教室: 高新校区 GT-B212



- □网页示例
 - □ 网站数据主要依托于网页(html, 超文本标记语言)展示
 - □ 用户Request服务器,服务器response信息(html等)
 - □ 课程主页<u>http://staff.ustc.edu.cn/~huangzhy/Course/DS2023.html</u>
 - □ 右键"检查"查看网页源代码



13

网页

```
<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Strict//EN" "http://www.w3.org/TR/xhtml1/DTD/xhtml1-strict.dtd">
</head>
<font size="8px"><strong> <center>Introduction to Data Science </center></strong> <font>
<font size="8px"><strong> <center>数据科学导论 </center>
<font size="6px"><strong> <center>课程代码: CS4023.01 </center></strong> <font>
<font size="6px"><strong> <center>学院: 011计算机科学与技术学院 </center></strong> <font>
<font size="6px">
   <center>课程教师:
   <a href="http://staff.ustc.edu.cn/~huangzhy/" target= "-parent" >黄振亚</a>,
   <a href="http://staff.ustc.edu.cn/~cheneh/" target= "-parent" >陈恩红</a>
<br/>font size="6px">
   <center>上课时间: (1-18周)每周二下午第8、9节: 15:55-17:30,教室: 高新校区 GT-B212 </center>
<hr size="1px" noshade>
助教(含作业提交):
<br/>\dir>\font size="5px">
   <a href="mailto:ds intro2024@163.com">data science 2025@163.com</a></a></a>
   <1i>QQ群: <a href="https://qm.qq.com/q/Fv9uQZVG8w">931835645</a> </font>
<hr size="1px" noshade>
Recommended Textbooks:
       <br><dir><font size="5px">
      <a href="http://book.douban.com/subject/5377669/" target= "-parent" > (Introduction to Data Mining) </a>
      <a href="http://book.douban.com/subject/10590856/" target= "-parent" > 《统计学习方法》</a></a>
       <1i><a href="https://book.douban.com/subject/26708119/" target= "-parent" > 《机器学习》</a>
       <1i><a href="https://item.jd.com/31931077735.html" target= "-parent" > 《大数据导论》</a></font>
```



- □ 载入:将目标网站数据下载到本地
 - □ Html, HyperText Markup Language
 - □爬虫程序向服务器发送网络请求 Request, 获取相应的网页
 - 网站常用网络协议: http, https
 - 数据常用请求方式: get, post
 - get: 参数常放置在URL中
 - http://www.abc.com?p=1&q=2&r=3,
 - 问号后为参数
 - post:参数常放置在一个表单中(报文头(header))
 - 在向目标URL发送请求时,将参数放置在一个网络请求的报文头中
 - 更安全



- 15
- □ 载入:将目标网站数据下载到本地
 - □ 数据常用请求方式: get, post
 - ■get: 参数常放置在URL中
 - http://www.abc.com?p=1&q=2&r=3,问号后为参数
 - 例如,https://www.baidu.com/s?ie=[中国科学技术大学]



×	Headers	Payload	Preview	Response	Initiator
			=%E4%B8	%AD%E5%9	B%BD%E/%
			7%91%E5%	%AD%A6%E	6%8A%80%
			6%9C%AF	%E5%A4%A [*]	7%E5%AD%
			&pbs=%E4	1%B8%AD%	E5%9B%BD
			7%A7%919	%E5%AD%A	6%E6%8A%
			0%E6%9C%	%AF%E5%A4	1%A7%E5%
			D%A6&csc	or=0&sugm	ode=2&hot
			unch=0&c	b=jQuery11	020160099
				1758608069	
			08069384		_
Requ	est Method		GET		
Statu	s Code		200 OK		
Remo	te Address		127.0.0.1:1	0080	
Refer	rer Policy		origin-whe	n-cross-orio	jin
				0 /	28/2025



16

- □ 载入:将目标网站数据下载到本地
 - □ 数据常用请求方式: get, post
 - post: 参数常放置在一个表单中
 - 在向目标URL发送请求时,将参数放置在一个网络请求的报文头中
 - 相比于Get,多了Form Data部分(请求体)
 - 更安全: 登录操作常用(不会放在URL后面)

Bai do 百度 · 用户名密码登录	
手机号/用户名/邮箱	
密码	
登录	
忘记密码?	
扫码登录 🛔 💣 🗞 立即注册	

×	Headers	Preview	Response	>>	
▼ G	eneral				
	Request UR	L: https:/	/passport.	baidu.	
	n				
	Request Method: POST				
	Status Code: 200 OK				
	Remote Ad	dress: 127.	0.0.1:9999)	
	Referrer Po	licy: unsaf	e-url		

请求体

Form Data	view source	view URL-encoded			
staticpage: https://www.baidu.com/cache/user/html/v3Jump.html					
charset: UTI	charset: UTF-8				
token: ddfa	token: ddfac7e17ce70dc6187ce33dffee73ed				
tpl: mn			_		
username:					
password:					
JyR3lzjQ6W	KlglOsWOjf125		_		
3izl5ofHOt	3izl5ofHOtMTeUneEHAbKW(



17

- □ 载入:将目标网站数据下载到本地
 - □ 数据常用<mark>请求方式: get</mark>,post
 - □ 获得服务器的响应: Response, 即获取网页源代码

响应头

▼ Response Headers View source

Bdpagetype: 2

Bdqid: 0xdabb3fd4000073b3

×	Headers	Preview	Response	Initiator	Timing	
▼ Ge	▼ General					
Request URL:			http://st	aff.ustc.edu	.cn/~huangzhy/	
Request Method:			GET	GET		
Status Code:			200 C	200 OK (from disk cache)		
Remote Address:			127.0.0.	127.0.0.1:10001		
Refe	Referrer Policy:			strict-origin-when-cross-origin		

响应状态: 200, 404等

HTTP状态码

响应体,即网页源代码

```
!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Strict//EN" "http://www.w3.org/TR/xhtml1/DTD/xhtml1-strict.dtd">
  font size="8px"><strong> <center>Introduction to Data Science </center></strong> <font>
  font size="<mark>8px"></mark><strong> <center>数据科学导论 </center></strong> <font
 〈font size="<mark>6px"</mark>><strong> <center>课程代码: CS4023.01 </center></strong> <font>
〈font size="<mark>6px</mark>"><strong> <center>学院: 011计算机科学与技术学院 </center></strong> <font>
      <a href="http://staff.ustc.edu.cn/huangzhy/" target= "-parent" >黄振亚</a>,
<a href="http://staff.ustc.edu.cn/cheneh/" target= "-parent" >陈思红</a>
       <center>上课时间: (1-18周)每周二下午第8、9节: 15:55-17:30, 教室: 高新校区 GT-B212 </center>
  hr size="lpx" noshade>
助教(含作业提交):
      <1i>程程, 于峻浩&nbsp; 〈a href="mailto:ds intro2024@163.com">data_science_2025@163.com</a>
〈li>Q即群: 〈a href="https://qm.qq.com/q/Fv9uQZVG8w">931835645</a> 〈li></font>
  hr size="lpx" noshade>
Recommended Textbooks:
             《ii〉<a href="<u>http://hook.douban.com/subject/5377669/</u>" target="-parent"> (Introduction to Data Mining) </a>/
(ii〉<a href="<u>http://hook.douban.com/subject/10590856/</u>" target="-parent"> (終計学习方法) </a>/
(ii〉<a href="<u>http://hook.douban.com/subject/10590856/</u>" target="-parent"> (桃学习方法) </a>/</a>/
(ii〉<a href="<u>https://hook.douban.com/subject/26708119/</u>" target="-parent"> (大数据号论) </a>/</a>/</a>/</a>/
(ii〉<a href="<u>https://item.jd.com/31931077735.html"</u> target="-parent"> (大数据号论) </a>/</a>/</a>/</a>/</a>/
```

- □ 实际操作: 抓取一个静态网页步骤
 - □ 首先确定URL,例如: http://www.baidu.com
 - □ 其次确定请求的方式以及相关参数:
 - 直接用浏览器实现: chrome, firefox浏览器抓包工具,详见
 - http://jingyan.baidu.com/article/3c343ff703fee20d377963e7.html
 - 或者抓包工具: charles等,详见
 - http://blog.csdn.net/jiangwei0910410003/article/details/41620363/
 - □ 最后在代码中按照特定的请求方式(get, post)向URL发送 参数,即可收到网页的结果





- 」但部分页面的数据是动态加载的
 - □ Ajax异步请求
 - 网页中的部分数据需要浏览器渲染 (JavaScript调用接口获取数据)
 - 用户的某些点击、下拉的<mark>操作</mark>触发才能获得
 - □解决方案:
 - ■借助抓包工具,分析Ajax某次操作所触发的请求,通过代码实现相应的请求
 - 有技术难度,但抓取速度快。
 - 利用智能化的工具: selenium webdriver
 - 用**程序控制驱动浏览器**,模拟浏览器
 - 可以模拟实现人的所有操作
 - 操作简单,但是速度慢
 - 因为爬虫需要启动浏览器,浏览器需要渲染页面,所以速度比较慢
 - 其他: Splash, Pyv8等

- 20
- □ 反爬虫: 随着网络爬虫对目标网站访问频率的加大,
 - 网站禁止爬虫程序继续访问
- □ 常见反爬手段:
 - □ 出现用户登录界面,需要验证码
 - □禁止某个固定帐号或ip一段时间内访问网站
 - □ 更有甚者,直接返回错误的无用数据

□ 应对措施:

- □ 优化爬虫程序,尽量减少访问次数,尽量不抓取重复内容
- □ 使用多个cookie(网站用来识别用户的手段,每个用户登录会 生成一个cookie)
- □ 使用多个ip (可以用代理实现)





网络爬虫:解析(对比JSON与XML)

```
"name": "中国",
"province": [{
  "name": "黑龙江",
  "cities": {
    "city": ["哈尔滨", "大庆"]
},
  "name": "广东",
  "cities": {
    "city": ["广州", "深圳", "珠海"]
},
}]
对象,成员:键值对
```

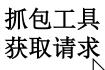
```
<?xml version="1.0" encoding="utf-8"?>
<country>
 <name>中国</name>
  ovince>
   <name>黑龙江</name>
   <cities>
     <city>哈尔滨</city>
     <city>大庆</city>
   </cities>
 </province>
 ovince>
    <name>广东</name>
   <cities>
     <city>广州</city>
     <city>深圳</city>
     <city>珠海</city>
   </cities>
 </country>
```



网络爬虫: 抓取微博评论

22





▼ General

Request URL: https://m.weibo.cn/api/comments/show?i

Request Method: GET Status Code: ● 200 0K

Remote Address: 123.125.106.67:443

Referrer Policy: no-referrer-when-downgrade

► Response Headers (14)

▼ Request Headers view source

Accept: application/json, text/plain, */*

Accept-Encoding: gzip, deflate, br

Accept-Language: zh-CN, zh; q=0.8, en; q=0.6

Connection: keep-alive

Cookie: _T_WM=d9a7dba4dd130f79eaecac13c8906050; ALbktAKLUXNkW1un7fu00CXjkppVYn1wGjJ3knF4g..; SUBP=0p5NHD95Q0So5Re0.cS020Ws4Dqcjn-fHBxHzLxK-LB.eLBK5L505136002; M_WEIBOCN_PARAMS=featurecode%3D200003236170084375%26uicode%3D20000061%26fid%3D414183617

Host: m.weibo.cn

Referer: https://m.weibo.cn/status/4141836170084375 User-Agent: Mozilla/5.0 (Macintosh; Intel Mac OS X

12.113 Safari/537.36

X-Requested-With: XMLHttpRequest



网络爬虫: 抓取微博评论

23

获得评论的json格式

```
- 京ICP备15025187号-1 邮箱: service@json.cn
       "mod_type":"mod/pagelist",
       "previous cursor":"",
       "next_cursor":"",
       "card_group":□[
          ⊟{
              "id":4142016554789113,
              "created at":"08-18 08:46",
              "source":"柔光自拍vivo X7",
              "user": ⊕Object{...},
              "text":"回复<a
href="/n/%E9%82%93%E8%B6%85">@邓超</a>:不管是谁,
请大家记住陈赫的话,他们很好,感情都很好。恳请各家粉丝不要
戏太多就好<i class="face face_1 icon_1">[微笑]</i>
没准你们那么嫌弃骂的那么难听,人家正主还是感情好的时不时去
吃火锅呢, 你们不累吗? 别用自己对他的爱去给他造成困扰",
              "reply_id":4142015488402958,
              "reply_text": "<a
href='/u/'5187664653>@邓超</a>: 我也不知道<i
class="face face_1 icon_20">[doge]</i>",
              "like_counts":10811,
              "liked":false,
              "mod type": "mod/single/infohoy"
```

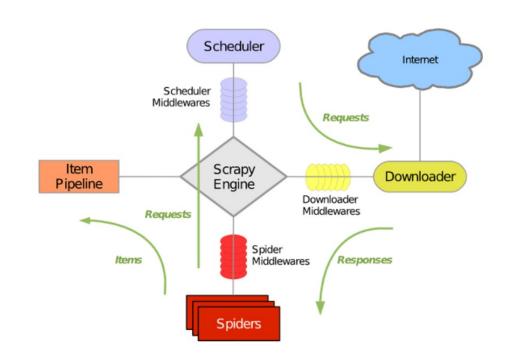
解析出需要的字段。

用户id	时间	内容
陈赫	08-18	天霸
邓超	08-18	我们都很好, 谢谢大家
邓超	08-18	我也不知道
贼亮zl	08-17	迪丽热巴
•••••	•••••	•••••



网络爬虫: 现有技术

- **24**
- □基于Java的工具
 - □ HttpClient
 - □ Jsoup
- □ 基于Python的工具
 - □ Scrapy
 - Beautiful-soup



现有的爬虫框架很成熟,能够合理的控制爬取的过程,并有效的 处理爬取过程中出现的各种异常,推荐使用Scrapy



数据采集

- □注意网站规定
- □注意法律规定
 - □ 2021年6月1日, 《中华人民共和国数据安全法》
- □注意数据使用规范
- etc

课达编程 (CODIA)·隐私政策

引言

课达编程重视用户的隐私。您在使用我们的服务时,我们可能会收集和使用您的相关信息。我们希望通过本《隐私政策》向您说明,在使用我们的服务时,我们如何收集、使用、储存和分享这些信息,以及我们为您提供的访问、更新、控制和保护这些信息的方式。本《隐私政策》与您所使用的课达编程服务息息相关,希望您仔细阅读,在需要时,按照本《隐私政策》的指引,作出您认为适当的选择。本《隐私政策》中涉及的相关技术词汇,我们尽量以简明扼要的表述,并提供进一步说明的链接,以便您的理解。

您使用或继续使用我们的服务,即意味着同意我们按照本《隐私政策》收集、使用、储存和分享 您的相关信息。

如对本《隐私政策》或相关事宜有任何问题,请通过 codia-support@bdaa.pro 与我们联系。



数据采集:课后练习(不考核)

26

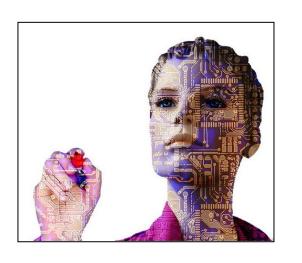
□ 写一个爬虫,爬取一个感兴趣的网站



27

➤ AI: 1993 – 现在

- □ 信息检索、互联网搜索
- □ 机器学习、自然会语言处理
- □ 数据挖掘、大数据分析
- □ "深蓝"战胜人类国际象棋冠军卡斯帕罗夫













OpenAI于2022年11月30 日发布的文本模态聊天机 器人,被社会广泛关注。









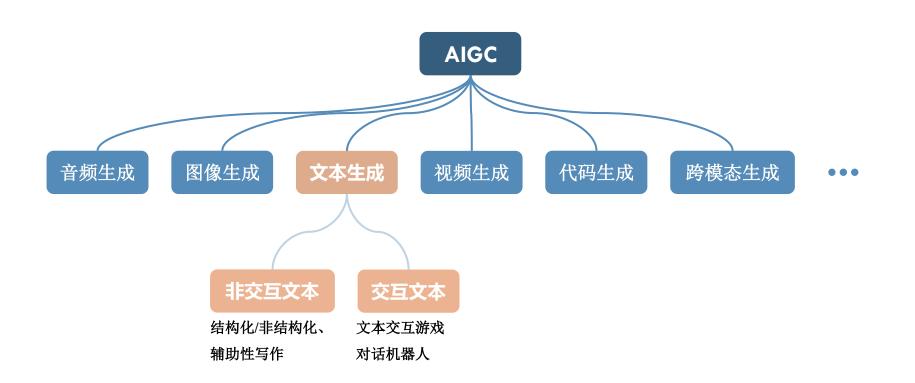




国内大模型从大模型多重能力构建,多模态语义理解等 角度快速发展,在多个领域超过ChatGPT。



AI Generated Content (AIGC, 人工智能自动生成内容) 技术,是一种新型内容创作方式,可以定制化生成需要的数据。





30

ChatGPT 4 v

You

请以鲁迅的写作风格,写一段对春天的描述

文本数据生成

鲁迅先生的文风犀利深刻,同时带有一丝讽刺,他对于社会现象的描写往往深入骨髓。试以他的风 格,来描述春天:

春天,这个听起来应当充满生机的季节,总是在无声无息中悄然到来。你看,那嫩绿的芽尖透过泥 土、似乎是在告诉人们、即使经历了漫长冷冽的冬日、生命依旧顽强地挣扎着要站直。田野里、那 些归来的燕子在空中划过弧线,似乎是在嘲笑地上那些还未觉醒的生灵,它们或许并不理解,即使 春光明媚,人间的烦恼和苦楚并未因此减少半分。

在这看似温暖的阳光下,农夫的背影显得更加佝偻,他们在田间辛勤地耕作,希望收获的不仅仅是 庄稼,更是对未来的一线希望。而城市里,那些街道和小巷,则被春雨洗刷得略显清新,却也掩盖 不了那永恒的尘埃和疲惫。春风并未带来太多温柔,反倒是像是在嘲讽那些还奢望着彻底变革的愚 人们。

于是,春天并不只是季节的轮回,它更像是一面镜子,折射出每一个人内心的希望与绝望。在这复 苏的季节里,人们的心却不见得能与大自然一样,得以更新。

< 2/2 > (1) (1) (5) (7)













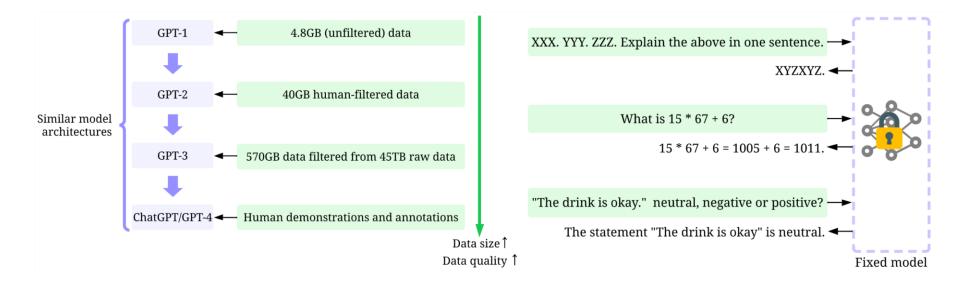


【提示词】逼真的特写视频展现了两艘海盗船在一杯咖啡中激烈交战,它们在咖啡液中航行。



数据科学基础

- □人工智能逐渐从以模型为中心过渡到以数据为中心
 - □ **GPT成功的数据基石**: **GPT**进化中,模型结构保持相似,训练数据的规模、质量得到极大提升
 - □模型对数据的需求急速上升





大模型的扩展定律

Scaling law

□ 损失值随着模型规模、训练数据量、训练计算量的增加 ,呈现出幂律下降的趋势

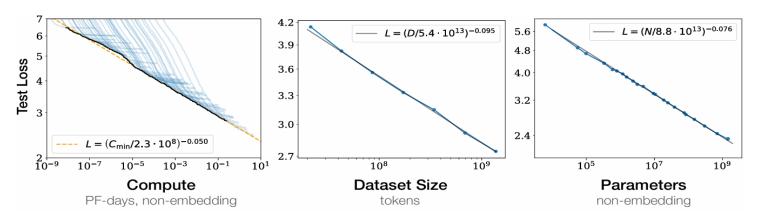


Figure 1 Language modeling performance improves smoothly as we increase the model size, datasetset size, and amount of compute used for training. For optimal performance all three factors must be scaled up in tandem. Empirical performance has a power-law relationship with each individual factor when not bottlenecked by the other two.

$$L(N,D)=E+rac{A}{N^lpha}+rac{B}{D^eta}$$



大模型的扩展定律

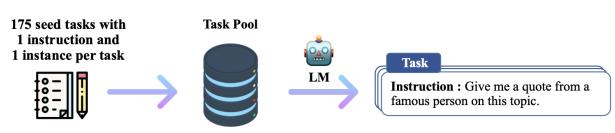
□通义干问:多样高质量训练数据

□数据过滤: 使用大模型多维度数据分析和过滤

□数据合成:使用大模型合成和过滤高质量训练数据

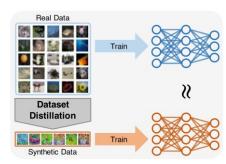
□数据混合:平衡低质量社交、电商数据 vs 高质量的技术、学术类数据

增加数据量



数据生成

改善数据质量



数据蒸馏



数据选择



数据分析基础

37

- □数据采集
- □ 数据预处理
- 」特征工程

Data Propression

Data Preprocessing

Feature engineering

