

基于大数据分析的药物筛选系统

课程设计报告

项目	内容
课程名称	大数据分析与应用
项目名称	基于深度学习的药物虚拟筛选系统
完成日期	2025年12月23日
开发环境	Windows 11 / Python 3.9 / PyTorch 2.7 / CUDA 11.8
硬件配置	NVIDIA RTX 3050 Ti (4GB显存)

摘要

本项目设计并实现了一个完整的基于深度学习的药物虚拟筛选系统。系统采用模块化分层架构，集成数据加载、分子特征提取、神经网络训练、模型评估、批量筛选及Web交互界面六大功能模块。针对药物小样本数据集的过拟合问题，提出分层采样与渐进式正则化相结合的解决方案。实验结果表明，优化后模型在BBBP任务上AUC-ROC达0.91，在ESOL任务上 R^2 达0.68。系统同时提供Streamlit前端界面与FastAPI后端服务，支持单分子预测与大规模批量筛选。

关键词：药物筛选；深度学习；分子指纹；过拟合；分层采样

第一章 绪论

1.1 什么是药物筛选系统？

1.1.1 通俗理解药物筛选

想象一下，你是一位药物研究人员，面前有一个装满了100万种不同化学物质的巨大"宝箱"。你的任务是从中找出能够治疗某种疾病的药物。传统方法是：把每种化学物质都拿出来，在实验室里一个一个地测试——这就是**高通量筛选 (HTS)**。

但问题来了：

- 测试100万种化合物需要多少时间和金钱？
- 如果99.9%的化合物都是无效的，那岂不是浪费了大量资源？

虚拟筛选就是用计算机来帮我们"预判"哪些化合物可能有用。就像网购时的"智能推荐"一样，计算机通过分析化合物的特征，预测它们的药效，帮我们筛选出最有希望的候选者，再进行实验验证。这样可以节省90%以上的实验成本！

1.1.2 本项目做了什么？

本项目构建了一个**智能药物筛选系统**，它能够：

- 输入：**一个化学分子的结构（用一种叫SMILES的文本格式表示）
- 输出：**预测这个分子的两个重要性质

- 能不能穿过血脑屏障？（对脑部疾病药物很重要）
- 在水中的溶解度如何？（影响药物能否被人体吸收）

举个例子：当你输入阿司匹林的化学结构"CC(=O)OC1=CC=CC=C1C(=O)O"，系统会告诉你它穿透血脑屏障的概率是多少，以及它的水溶解度值。

1.2 研究背景与意义

1.2.1 药物研发的挑战——"双十定律"

药物研发是现代医学和生物技术领域中最具挑战性的任务之一。业界有一个著名的"双十定律"：

维度	数据
平均研发时间	10-15年
平均研发成本	10-26亿美元
成功率	不足10%

这意味着：一款新药从实验室发现到最终上市，平均需要花费超过10年时间和10亿美元，而且还有90%的可能会失败！

药物研发的主要阶段：

靶点发现	→	先导化合物筛选	→	临床前研究	→	I期临床	→	II期临床	→	III期临床	→	上市
(2年)		(2-3年)		(2-3年)		(1年)		(2年)		(3年)		(1年)

其中，**先导化合物筛选**阶段需要从数百万个候选化合物中找出有希望的"种子选手"，这正是本项目要解决的问题。

1.2.2 传统高通量筛选的局限性

高通量筛选 (High-Throughput Screening, HTS) 是传统的药物筛选方法：

工作原理：

1. 准备数十万到数百万个化合物样品
2. 使用自动化机器人系统进行实验测试
3. 每天可测试数万个化合物

存在的问题：

问题	具体表现
成本高昂	筛选100万化合物约需1000万美元
命中率低	通常只有0.01%-0.1%的化合物有活性
时间长	完整筛选需要数月时间
资源浪费	99.9%的实验都是"无用功"

1.2.3 虚拟筛选技术的兴起

虚拟筛选 (Virtual Screening, VS) 就像是给实验筛选装上了"智能过滤器":

传统方法: 100万化合物 → 全部实验测试 → 找到100个有效的

虚拟筛选: 100万化合物 → 计算机预筛选 → 1万候选 → 实验测试 → 找到100个有效的

↑

节省99%的实验成本!

虚拟筛选的两大类方法:

方法	全称	原理	优点	缺点
SBVS	基于结构的虚拟筛选	模拟药物与蛋白质的结合	准确性高	需要蛋白质3D结构
LBVS	基于配体的虚拟筛选	分析已知有效药物的特征	不需要蛋白质结构	依赖已知数据

本项目采用LBVS方法: 通过分析已知能/不能穿透血脑屏障的分子特征, 训练深度学习模型, 让计算机"学会"判断新分子的性质。

1.2.4 深度学习在药物发现中的应用

为什么要用深度学习?

传统机器学习方法需要人工设计特征 (即告诉计算机"看哪些方面") , 而深度学习能够自动学习特征:

传统方法: 分子 → 人工设计特征 (需要专家知识) → 机器学习模型 → 预测

深度学习: 分子 → 自动特征提取 → 深度神经网络 → 预测

↑

不需要人工干预!

深度学习在药物领域的应用场景:

应用	说明	本项目涉及
ADMET预测	预测吸收、分布、代谢、排泄、毒性	✓
药物-靶点相互作用	预测药物能否与特定蛋白结合	
分子生成	设计全新的药物分子	
药物重定位	发现老药的新用途	

本项目聚焦的两个核心任务:

任务	全称	类型	意义
BBBP	Blood-Brain Barrier Penetration	二分类	判断药物能否进入大脑, 对治疗阿尔茨海默病、帕金森病等脑部疾病至关重要
ESOL	Estimated SOLubility	回归	预测药物在水中的溶解度, 影响药物能否被人体有效吸收

1.3 项目目标与创新点

1.3.1 项目目标

本项目要完成的四大任务：

目标	具体内容	完成情况
自动化数据流水线	自动下载、解析MoleculeNet数据集	✓
高性能预测模型	设计适合小样本数据的神经网络	✓
解决过拟合问题	提出有效的正则化策略	✓
完整应用系统	Web界面 + API服务	✓

1.3.2 主要创新点

- 1. **分层采样策略**：解决数据集类别不平衡问题，确保训练/验证/测试集分布一致
- 2. **渐进式Dropout机制**：创新性地设计随网络深度递增的Dropout策略
- 3. **轻量化模型设计**：在保证性能的前提下大幅减少模型参数量
- 4. **端到端系统集成**：从数据处理到Web部署的完整解决方案

第二章 系统设计

2.1 系统设计思想

2.1.1 什么是模块化设计？

想象一下乐高积木：每块积木都有特定的功能，可以独立存在，也可以组合成复杂的结构。**模块化设计**就是这个思想的软件版本：

传统设计：一整块代码，改一处可能影响全局

所有功能混在一起，难以维护

模块化设计：功能分离，各司其职

数据层

特征层

模型层

应用层

↓

↓

↓

↓

各自独立，互不干扰，方便维护和扩展

模块化的好处：

- 1. **易于维护**：修改一个模块不会影响其他模块
- 2. **易于测试**：可以单独测试每个模块
- 3. **易于复用**：同一模块可用于不同项目

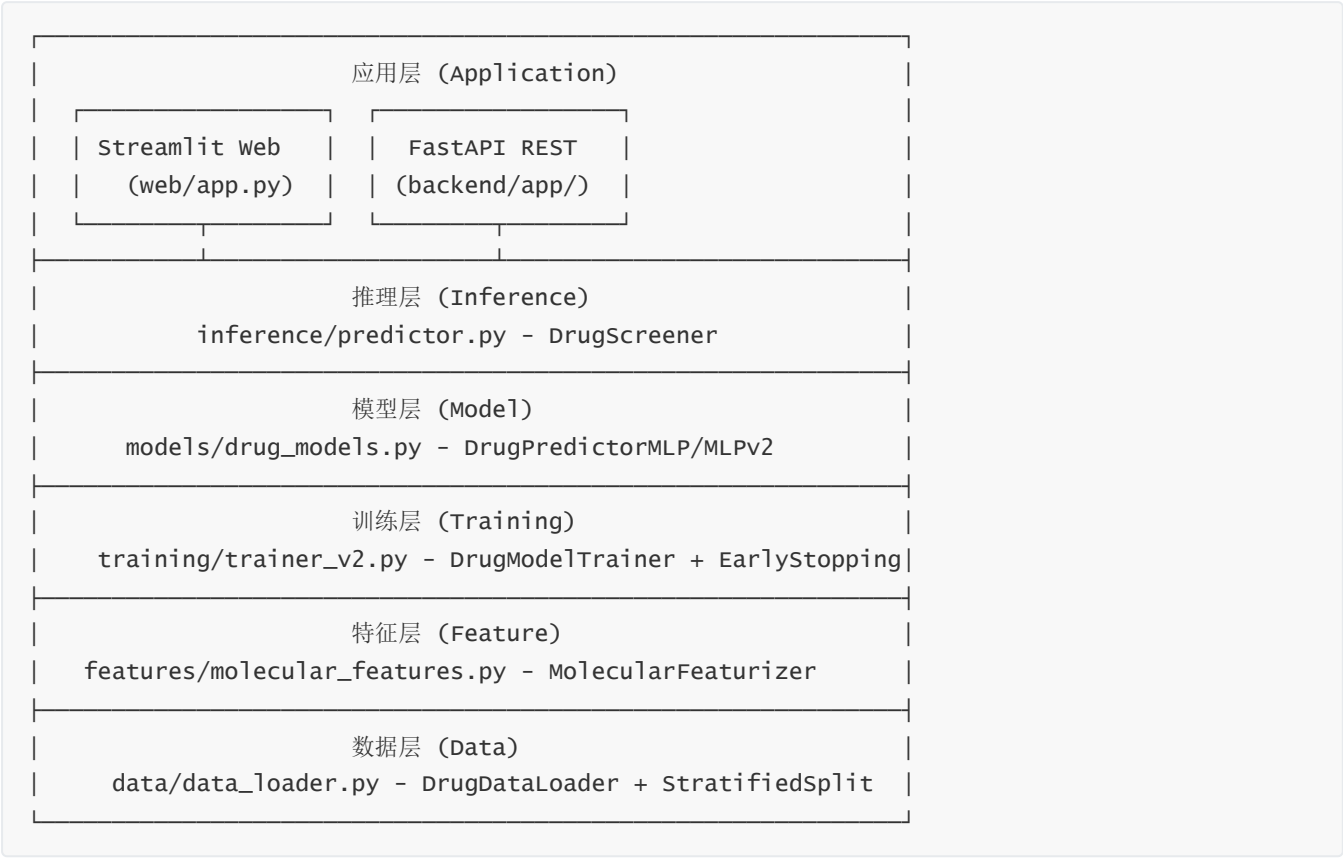
4. 易于协作：不同人可以同时开发不同模块

2.1.2 分层架构解释

本系统采用六层架构，从下到上依次为：

层级	职责	通俗解释
数据层	加载和预处理数据	相当于"食材采购和清洗"
特征层	将分子转换为数字	相当于"把食材切好备用"
模型层	定义神经网络结构	相当于"设计菜谱"
训练层	训练模型学习规律	相当于"按菜谱练习做菜"
推理层	使用模型进行预测	相当于"正式出菜"
应用层	提供用户界面	相当于"餐厅服务"

2.2 总体架构图



2.3 目录结构详解

下面详细解释项目中每个文件夹和文件的作用：

```
drug/
├── data/                                # 数据模块--负责"食材采购"
│   └── data_loader.py                  # 核心文件：下载数据、分层采样
```

	raw/	# 存放从网上下载的原始数据
	processed/	# 存放处理后的数据
	features/	# 特征工程--负责"食材加工"
	molecular_features.py	# 核心文件: 将分子转换为数字向量
	feature_extraction.py	# 特征提取的辅助工具
	models/	# 模型定义--负责"设计神经网络"
	drug_models.py	# 定义MLP和MLPv2两种网络结构
	training/	# 训练模块--负责"教会模型"
	trainer.py	# 基础版训练器
	trainer_v2.py	# 增强版训练器(带早停、学习率调度)
	evaluation/	# 评估模块--负责"打分评价"
	metrics.py	# 计算各种评估指标
	figures/	# 保存可视化图表
	inference/	# 推理模块--负责"实际预测"
	predictor.py	# 单分子预测和批量筛选功能
	web/	# web前端--负责"用户界面"
	app.py	# Streamlit网页应用
	backend/	# API后端--负责"程序接口"
	app/	# FastAPI服务
	saved_models/	# 保存训练好的模型文件
	train_model.py	# 一键训练脚本
	requirements.txt	# 项目依赖列表

2.4 技术栈说明

什么是技术栈？ 技术栈就是完成项目所使用的技术工具集合，就像做饭需要锅、铲子、调料一样。

类别	技术	版本	作用说明
深度学习框架	PyTorch	≥2.0.0	构建和训练神经网络的核心工具
化学信息学	RDKit	2023.3+	处理分子结构、生成分子指纹
化学信息学	DeepChem	2.7+	提供MoleculeNet数据集接口
数据处理	NumPy	-	数值计算基础库
数据处理	Pandas	-	表格数据处理
机器学习	Scikit-learn	-	数据划分、评估指标计算
Web前端	Streamlit	≥1.30.0	快速构建数据应用网页
API后端	FastAPI	0.104+	构建REST API服务
GPU加速	CUDA	11.8	利用显卡加速训练

第三章 数据工程

3.1 什么是数据工程？

数据工程是机器学习项目中最基础也是最重要的环节，有一句话说得好："Garbage in, garbage out"（垃圾进，垃圾出）——如果输入的数据有问题，模型再好也没用。

数据工程主要包括：

- 1. **数据获取**：从哪里获得数据？
- 2. **数据清洗**：去除错误、缺失的数据
- 3. **数据划分**：如何分配训练集、验证集、测试集？
- 4. **特征提取**：如何把原始数据转换成模型能理解的形式？

3.2 MoleculeNet数据集介绍

3.2.1 MoleculeNet是什么？

MoleculeNet 是斯坦福大学在2018年发布的"药物分子界的ImageNet"——就像ImageNet是图像识别领域的标准数据集一样，MoleculeNet是分子性质预测领域的标准数据集。

为什么需要标准数据集？

- 让不同研究者的结果可以比较
- 避免每个人用不同数据导致无法对比
- 提供高质量、经过验证的数据

MoleculeNet包含什么？

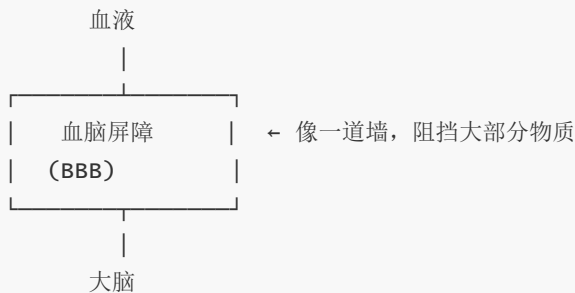
类别	包含数据集	样本数量
量子力学	QM7, QM8, QM9	7K-130K
物理化学	ESOL, FreeSolv, Lipophilicity	600-4K
生物物理	PDBbind, PCBA	4K-400K
生理学	BBBP, Tox21, SIDER	1K-8K

本项目使用的两个数据集：BBBP 和 ESOL

3.2.2 BBBP数据集详解

什么是血脑屏障？

血脑屏障（Blood-Brain Barrier, BBB）是大脑的"保安系统"：



血脑屏障的作用：

- ☒ 阻止病毒、细菌、毒素进入大脑
- ☒ 保护神经系统免受损害
- ☒ 同时也阻挡了99%的药物分子！

为什么BBB穿透性很重要？

治疗脑部疾病（如阿尔茨海默病、帕金森病、脑肿瘤）的药物必须能穿过血脑屏障才能发挥作用。如果一个药物不能穿透BBB，那么即使在试管里效果再好，也无法治疗脑部疾病。


BBBP数据集统计：


属性	值	说明
化合物数量	2,039个	相对较小的数据集
正例（能穿透）	~1,670个 (82%)	占多数
负例（不能穿透）	~369个 (18%)	占少数
任务类型	二分类	输出：能/不能
评估指标	AUC-ROC	越接近1越好

3.2.3 ESOL数据集详解

什么是溶解度？

溶解度是指物质在特定溶剂（如水）中能够溶解的最大量。对药物来说：

高溶解度药物：

完全溶解
易于吸收

低溶解度药物：

大部分沉淀
难以吸收

为什么溶解度很重要？

药物要发挥作用，必须先被人体吸收。而口服药物的吸收过程是：

药物 → 溶解在胃肠液中 → 穿过肠壁 → 进入血液 → 到达病灶

↑

溶解度决定这一步！

如果药物溶解度太低，就无法被有效吸收，再好的药效也发挥不出来。

ESOL数据集统计：

属性	值	说明
化合物数量	1,128个	较小的数据集
目标值	log溶解度	取对数是因为溶解度范围很大
数值范围	-11 到 +2	log(mol/L)
任务类型	回归	输出：具体数值
评估指标	RMSE, R ²	RMSE越小越好，R ² 越接近1越好

3.3 什么是类别不平衡问题？

3.3.1 问题说明

想象你在训练一个判断猫狗的模型，但训练数据是：

- 猫：9000张图片
- 狗：100张图片

模型可能会学到一个"聪明"但无用的规律："无论什么图片，都预测是猫，这样就能达到99%的准确率！"

这就是**类别不平衡问题**。在BBBP数据集中：

- 能穿透BBB的分子：82%
- 不能穿透BBB的分子：18%

如果模型总是预测"能穿透"，准确率也能达到82%，但这显然没用！

3.3.2 为什么随机划分会出问题？

假设我们随机划分数据：

原始数据：100个样本（82个正例，18个负例）

随机划分后可能出现：
训练集(80个)：70个正例，10个负例 → 正例87.5%
测试集(20个)：12个正例，8个负例 → 正例60%

问题：训练集和测试集的分布不一致！
模型在训练集上学的规律，在测试集上可能不适用。

3.3.3 分层采样如何解决？

分层采样（Stratified Sampling） 确保划分后各子集保持相同的类别比例：

原始数据：100个样本（82个正例，18个负例）→ 正例82%

分层采样后：

训练集(80个)：66个正例，14个负例 → 正例82.5% ✓

测试集(20个)：16个正例，4个负例 → 正例80% ✓

各子集保持相近的类别比例！

代码实现：

```
# data/data_loader.py - DrugDataLoader类
def load_moleculenet_with_stratified_split(self, dataset_name='BBBP',
                                          train_ratio=0.8, val_ratio=0.1):

    # 使用sklearn分层采样
    splitter = StratifiedShuffleSplit(n_splits=1, test_size=0.2, random_state=42)
    for train_idx, test_idx in splitter.split(X, y):
        x_train, x_test = X[train_idx], X[test_idx]
        y_train, y_test = y[train_idx], y[test_idx]
    # 训练/验证/测试集正例比例均维持在~76.5%
```

3.4 分子特征提取——如何让计算机"看懂"分子？

3.4.1 核心问题

计算机只能处理数字，但分子是由原子和化学键组成的复杂结构。我们需要找到一种方法，把分子"翻译"成数字：

分子结构 → 数字向量 → 神经网络 → 预测结果
如何转换？

3.4.2 SMILES——分子的"文字表示"

SMILES (Simplified Molecular Input Line Entry System) 是一种用文本表示分子结构的方法。

举例说明：

分子	SMILES表示	说明
水 (H ₂ O)	<chem>O</chem>	氧原子，氢可以省略
乙醇	<chem>CCO</chem>	两个碳一个氧
苯	<chem>c1ccccc1</chem>	六个碳组成的环
阿司匹林	<chem>CC(=O)OC1=CC=CC=C1C(=O)O</chem>	复杂结构

SMILES的规则：

- 大写字母表示原子（C=碳，O=氧，N=氮...）
- 数字表示环的闭合点
- = 表示双键

- `()` 表示分支

但SMILES只是文字，不能直接输入神经网络。我们需要进一步转换。

3.4.3 分子指纹——分子的"数字身份证"

分子指纹 (Molecular Fingerprint) 是将分子转换为固定长度数字向量的方法。

类比理解：就像人的指纹可以唯一标识一个人，分子指纹可以唯一标识一个分子的化学特征。

```
阿司匹林的SMILES → [1,0,0,1,1,0,1,0,0,1,1,1,0,0,1,0,1,1,0,0,...]
                        ↑
                    1024个0或1组成的向量
```

本项目使用的ECFP4指纹：

ECFP (Extended-Connectivity Fingerprints) 是目前最流行的分子指纹之一。ECFP4中的"4"表示考虑每个原子周围直径为4（半径为2）的化学环境。

ECFP的工作原理（简化版）：

Step 1: 给每个原子分配初始特征（原子类型、连接数等）

```
C-C-O-H
↓ ↓ ↓ ↓
5 5 8 1 ← 初始编码
```

Step 2: 收集每个原子周围的邻居信息，更新编码

考虑半径=1：每个原子收集直接相邻原子的信息
考虑半径=2：每个原子收集两步以内所有原子的信息

Step 3: 将所有子结构编码映射到固定长度的向量

```
[1,0,0,1,1,0,1,0,...] ← 1024维向量
```

3.4.4 代码实现详解

```
# features/molecular_features.py - MolecularFeaturizer类
class MolecularFeaturizer:
    def __init__(self, fingerprint_size=1024, radius=2):
        """
        初始化特征提取器
        fingerprint_size=1024: 输出向量的长度
        radius=2: ECFP4的半径 (2×2=4)
        """
        self.morgan_generator = rdFingerprintGenerator.GetMorganGenerator(
            radius=self.radius, fpSize=self.fingerprint_size
        )

    def get_morgan_fingerprint(self, smiles: str) -> np.ndarray:
        """
        将SMILES转换为分子指纹
```

```
输入: "CCO" (乙醇的SMILES)
输出: [0,1,0,0,1,...] (1024维向量)
"""

mol = Chem.MolFromSmiles(smiles) # SMILES → 分子对象
fp = self.morgan_generator.GetFingerprintAsNumPy(mol) # 分子 → 指纹
return fp.astype(np.int8) # 返回1024维二进制向量
```

3.4.5 物理化学描述符

除了分子指纹，我们还提取一些易于理解的物理化学性质：

```
def get_molecular_descriptors(self, smiles: str) -> dict:
    """计算分子的物理化学性质"""
    return {
        'MolecularWeight': Descriptors.Molwt(mol),          # 分子量（越大分子越重）
        'LogP': Descriptors.MolLogP(mol),                  # 亲脂性（越大越容易溶于油）
        'TPSA': Descriptors.TPSA(mol),                      # 极性表面积
        'NumHDonors': Descriptors.NumHDonors(mol),          # 氢键供体数量
        'NumHAcceptors': Descriptors.NumHAcceptors(mol),    # 氢键受体数量
        # ...
    }
```

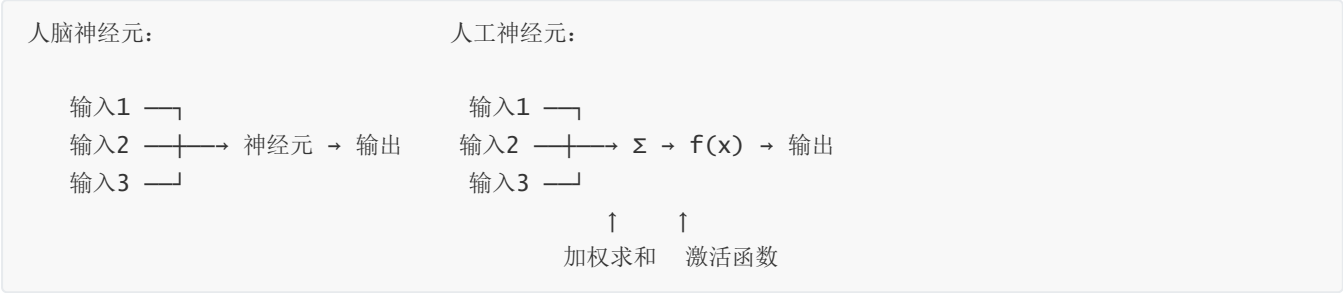
这些描述符在药物筛选中有明确的物理化学意义，用于Lipinski类药性规则检查。

第四章 模型设计

4.1 什么是神经网络？ (入门介绍)

4.1.1 从人脑到人工神经网络

神经网络的灵感来源于人脑。人脑中有约860亿个神经元，它们通过突触相互连接，传递信息。人工神经网络模仿了这一结构：

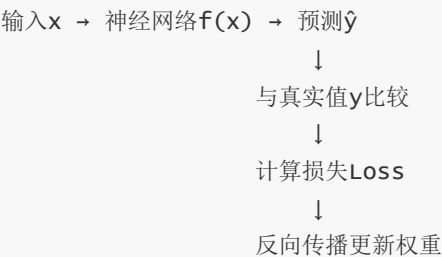


4.1.2 神经网络如何"学习"?

类比：神经网络学习就像小孩学认字。

1. 展示一个字（输入数据）
2. 让小孩猜这是什么字（前向传播，得到预测）
3. 告诉小孩正确答案（计算误差）
4. 小孩调整记忆（反向传播，更新权重）
5. 重复以上过程无数次（多轮训练）
6. 小孩学会了认字！（模型训练完成）

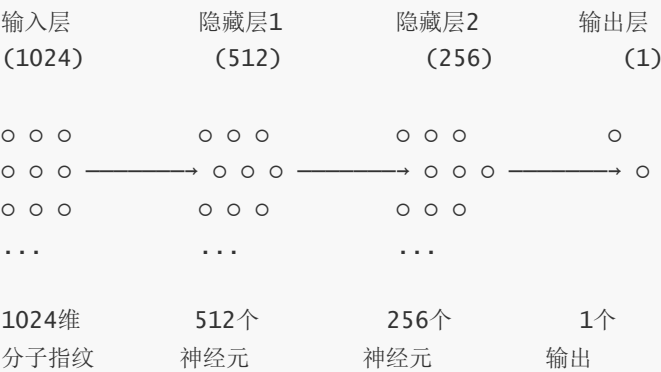
用数学语言描述：



4.2 多层感知机（MLP）详解

4.2.1 什么是MLP？

MLP（Multi-Layer Perceptron，多层感知机） 是最基础的神经网络类型，由多个全连接层组成。



"全连接"的含义： 每一层的每个神经元都与下一层的所有神经元相连。

4.2.2 每一层都做什么？

以本项目为例，解释从输入到输出的完整过程：

输入： $[0, 1, 0, 0, 1, 1, 0, \dots]$ # 1024维分子指纹

第1步：线性变换

$x_1 = w_1 \cdot x + b_1$ # 1024维 \rightarrow 512维
w_1 是 512×1024 的权重矩阵， b_1 是 512 维偏置向量

第2步：批归一化（BatchNorm）

$x_1 = \text{BatchNorm}(x_1)$ # 标准化，让数值分布更稳定

第3步：激活函数（ReLU）

$x_1 = \text{ReLU}(x_1) = \max(0, x_1)$ # 引入非线性

第4步: Dropout

$x_1 = \text{Dropout}(x_1)$ # 随机丢弃部分神经元, 防止过拟合

重复以上步骤两次...

最后：输出层

```
output = w_out · x + b_out # 输出预测值
```

4.2.3 各组件的作用解释

组件	作用	通俗解释
线性层(Linear)	学习特征组合	就像给每个特征打分，然后加权求和
批归一化(BatchNorm)	稳定训练过程	让数据保持在合理范围，不会太大或太小
ReLU激活	引入非线性	没有它，多层网络等于一层；有了它才能学复杂规律
Dropout	防止过拟合	随机"关闭"一些神经元，防止模型死记硬背

4.3 什么是过拟合?

4.3.1 过拟合的直观理解

过拟合 (Overfitting) 是机器学习中最常见的问题之一。

类比：想象一个学生准备考试：

正常学习：理解知识点 → 能解决类似的新问题 ✓

死记硬背：背下所有练习题答案 → 遇到新题就懵了 X

↑
这就是过拟合！

在本项目中的表现:


情况	训练集准确率	测试集准确率	诊断
欠拟合	60%	58%	模型太简单，学不会
正常	85%	82%	模型学到了通用规律 ✓
过拟合	99%	65%	模型只是在"背答案" ✗


4.3.2 为什么药物数据容易过拟合?

原因	说明
样本量小	BBBP只有2039个样本，模型很容易记住所有样本
特征维度高	1024维指纹 vs 2039样本，特征比样本还多！

原因	说明
数据噪声	实验数据本身存在测量误差

示意图：

样本数量：  2039个

特征维度：  1024维

特征维度几乎和样本数量一样多！
模型很容易找到"投机取巧"的规律来记住训练数据。

4.4 DrugPredictorMLPv2——我们的解决方案

4.4.1 设计思路

针对过拟合问题，我们设计了增强版模型DrugPredictorMLPv2，采用多重策略：

基础模型(v1)的问题：

网络太大（67万参数）→ 容易记住数据

Dropout太小(0.2) → 正则化不够

类别不平衡 → 学习偏差

↓ 改进

增强模型(v2)的解决：

网络变小（29万参数）→ 限制记忆能力

Dropout增大(0.5-0.7) → 强正则化

分层采样 → 分布一致

4.4.2 核心改进1：渐进式Dropout

传统Dropout：每层使用相同的Dropout比例（如0.2）

渐进式Dropout：随着网络深度增加，Dropout比例也增加

第1层 → Dropout 0.5 （丢弃50%的神经元）

第2层 → Dropout 0.6 （丢弃60%的神经元） ↑

第3层 → Dropout 0.7 （丢弃70%的神经元） 越深越狠！

为什么越深dropout越大？

- 浅层特征：比较通用，如"分子中有苯环"
- 深层特征：比较具体，如"特定位置的特定基团"
- 深层更容易过拟合到训练数据的特定模式，所以需要更强的正则化

4.4.3 核心改进2：输入层Dropout

在输入的分子指纹上也应用Dropout：

```
self.input_dropout = nn.Dropout(0.25) # 随机丢弃25%的输入特征
```

效果：防止模型过度依赖某些特定的子结构特征，强迫模型学习更鲁棒的特征组合。

4.4.4 核心改进3：轻量化网络

配置	基础模型v1	增强模型v2	变化
第1层	512神经元	256神经元	-50%
第2层	256神经元	128神经元	-50%
第3层	128神经元	64神经元	-50%
总参数	~67万	~29万	-57%

为什么更小的网络反而更好？

网络容量越大，“记忆力”越强，越容易记住训练数据的细节（包括噪声）。适当减小网络，限制其记忆能力，反而能学到更通用的规律。

4.4.5 完整模型代码

```
class DrugPredictorMLPv2(nn.Module):
    def __init__(self, input_dim=1024, hidden_dims=[256,128,64], dropout=0.5):
        super().__init__()

        # 输入层Dropout（防止过度依赖特定特征）
        self.input_dropout = nn.Dropout(dropout * 0.5) # 0.25

        # 构建隐藏层
        layers = []
        prev_dim = input_dim
        for i, hidden_dim in enumerate(hidden_dims):
            layers.append(nn.Linear(prev_dim, hidden_dim)) # 线性变换
            layers.append(nn.BatchNorm1d(hidden_dim)) # 批归一化
            layers.append(nn.ReLU()) # 激活函数

            # 渐进式Dropout: 0.5 -> 0.6 -> 0.7
            layer_dropout = min(dropout * (1 + i * 0.1), 0.7)
            layers.append(nn.Dropout(layer_dropout))

            prev_dim = hidden_dim

        self.hidden_layers = nn.Sequential(*layers)
        self.output_layer = nn.Linear(prev_dim, 1) # 输出层

        # xavier初始化
        self._init_weights()
```



```
def forward(self, x):  
    x = self.input_dropout(x)      # 输入Dropout  
    x = self.hidden_layers(x)      # 隐藏层处理  
    return self.output_layer(x)   # 输出预测
```

ECFP的工作原理：

1. **初始化**：为每个原子分配一个基于化学性质的初始标识符
2. **迭代扩展**：在指定半径内收集邻居原子信息，更新标识符
3. **哈希映射**：将所有子结构标识符映射到固定长度的位向量

ECFP4的"4"表示直径为4（半径为2），意味着每个子结构特征考虑中心原子周围2跳范围内的化学环境。

第五章 训练系统

5.1 什么是模型训练？（入门介绍）

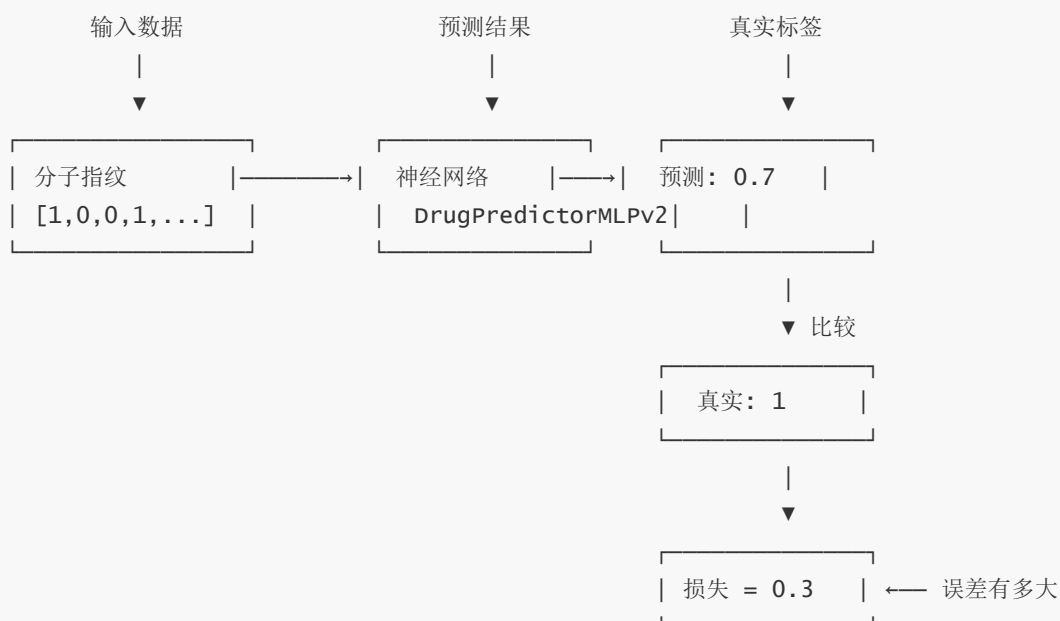
5.1.1 训练的本质——让机器"学习"

模型训练是让神经网络从数据中学习规律的过程。这个过程类似于教小孩学习：

老师教小孩认字的过程：

1. 展示一个字 —————→ 给神经网络输入数据
2. 小孩猜测这是什么字 —————→ 神经网络做出预测
3. 告诉小孩对还是错 —————→ 计算损失函数
4. 小孩调整记忆方式 —————→ 反向传播更新权重
5. 重复成千上万次 —————→ 多轮训练（Epochs）
6. 小孩学会认字！ —————→ 模型训练完成！

5.1.2 训练过程的可视化理解



↓
▼ 反向传播

调整网络权重 ← 让下次预测更准

5.1.3 一个完整训练周期（Epoch）

Epoch（轮次）是指模型"看完"整个训练集一遍的过程。

假设训练集有1600个分子，批大小为32：

1个Epoch = $1600 \div 32 = 50$ 个批次（Batches）

每个批次：

1. 取出32个分子的数据
2. 模型预测这32个分子的性质
3. 计算这32个预测与真实值的平均误差（损失）
4. 反向传播，更新模型参数

重复50次 = 完成1个Epoch

通常需要训练几十到几百个Epochs

5.2 损失函数——衡量"错得有多离谱"

5.2.1 什么是损失函数？

损失函数（Loss Function） 是衡量模型预测值与真实值之间差距的数学函数。

类比：想象你在玩投飞镖游戏：

- 靶心是真实值（正确答案）
- 你投的位置是预测值
- 离靶心的距离就是"损失"

损失小：

- 真实值
- 预测值

距离近，预测准！

损失大：

- 真实值

●——预测值

距离远，预测差！

损失函数的作用：

- 告诉模型"你错了多少"
- 指导模型朝着减少错误的方向调整参数

5.2.2 分类任务的损失函数

BBBP任务是二分类任务（能/不能穿透），使用二元交叉熵损失：

$$L = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\sigma(z_i)) + (1 - y_i) \log(1 - \sigma(z_i))]$$

通俗理解：

真实标签	模型预测概率	损失值	解释
1（能穿透）	0.9	0.11	预测对了，损失小 ✓
1（能穿透）	0.1	2.30	预测错了，损失大 ✗
0（不能）	0.1	0.11	预测对了，损失小 ✓
0（不能）	0.9	2.30	预测错了，损失大 ✗

核心理想：预测越偏离真实值，损失越大（指数级增长）

5.2.3 回归任务的损失函数

ESOL任务是回归任务（预测具体数值），使用均方误差损失：

$$L = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

简单说：预测值和真实值的差的平方的平均值。

例子：预测3个分子的溶解度

分子1：真实=-2.5，预测=-2.3，误差²=(0.2)²=0.04

分子2：真实=-3.0，预测=-2.0，误差²=(1.0)²=1.00

分子3：真实=-1.0，预测=-1.2，误差²=(0.2)²=0.04

MSE = (0.04 + 1.00 + 0.04) / 3 = 0.36

5.3 优化器——如何调整模型参数

5.3.1 什么是优化器？

优化器（Optimizer）决定了模型如何根据损失函数来更新参数。

类比：想象你蒙着眼睛在山上，要找到最低点（损失最小的点）：



- **梯度**：告诉你哪个方向是下坡（减少损失的方向）
- **学习率**：每一步迈多大（参数更新的幅度）
- **优化器**：决定怎么走（更新策略）

5.3.2 AdamW优化器详解

本项目使用**AdamW优化器**，它是目前最流行的优化器之一，结合了多种优化技术：

1. 动量 (Momentum) ——惯性

没有动量：	有动量：
↓	↓↓↓
↓ 走走停停	顺着之前的方向继续
↓	更快到达目标！

2. 自适应学习率——智能步长

普通梯度下降：所有参数用同样的学习率

Adamw：根据每个参数的历史梯度，自动调整各自的学习率

频繁更新的参数 → 小学习率（已经学得差不多了）

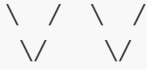


很少更新的参数 → 大学习率（还需要多学习）

3. 权重衰减 (Weight Decay) ——防止过拟合

让权重值不会变得太大，相当于给模型加了"紧箍咒"。

5.3.3 学习率的重要性

学习率是最重要的超参数之一：

学习率太大：	学习率太小：	学习率合适：
		
在最优点附近跳来跳去 永远到不了！	下降太慢 可能训练不完	平稳收敛 到达最优！

ReduceLROnPlateau策略：

当验证损失连续多个Epoch不再下降时，自动将学习率减半：

初始学习率：0.001

↓

连续10个Epoch损失不下降

↓

学习率变为：0.0005 （减半）

↓

继续训练...

5.4 早停机制——防止"学过头"

5.4.1 为什么需要早停?

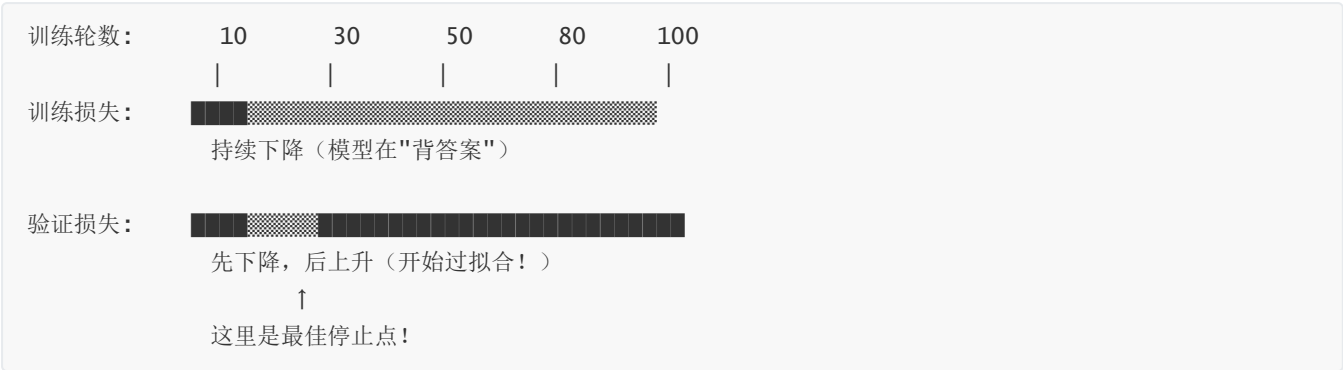
类比：烤蛋糕

烤太短：蛋糕没熟（欠拟合）

烤刚好：蛋糕完美（刚刚好）

烤太久：蛋糕焦了（过拟合）

训练神经网络也是一样：

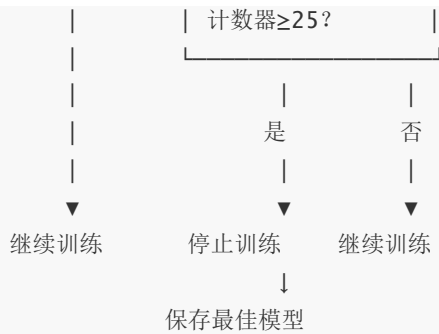


5.4.2 早停的工作原理

```
class EarlyStopping:
    def __init__(self, patience=25, min_delta=0.0001):
        """
        patience=25: 如果验证损失连续25轮没有改善，就停止训练
        min_delta=0.0001: 损失下降超过这个值才算"改善"
        """
        self.patience = patience
        self.min_delta = min_delta
        self.counter = 0 # 计数器：记录连续多少轮没改善
        self.best_score = None # 记录最佳损失值
```

工作流程图：





5.5 训练器的完整实现

```
# training/trainer_v2.py - DrugModelTrainer类
class DrugModelTrainer:
    """
    药物模型训练器
    负责整个训练过程的管理，包括：前向传播、损失计算、反向传播、参数更新
    """

    def __init__(self, model, learning_rate=0.001, weight_decay=1e-5):
        """
        初始化训练器

        参数：
        - model: 要训练的神经网络模型
        - learning_rate: 学习率（步长）
        - weight_decay: 权重衰减（L2正则化强度）
        """
        self.model = model

        # 损失函数：衡量预测与真实值的差距
        self.criterion = nn.BCEWithLogitsLoss() # 分类任务
        # self.criterion = nn.MSELoss()          # 回归任务

        # AdamW优化器：决定如何更新模型参数
        self.optimizer = optim.AdamW(
            model.parameters(), # 要优化的参数
            lr=learning_rate,   # 学习率
            weight_decay=weight_decay # 权重衰减（防过拟合）
        )

        # 学习率调度器：当验证损失停止下降时，自动降低学习率
        self.scheduler = ReduceLROnPlateau(
            self.optimizer,
            mode='min', # 监控指标越小越好
            factor=0.5, # 学习率减半
            patience=10 # 容忍10轮不改善
        )

    def train_epoch(self, data_loader):
        """
        训练一个完整的Epoch
```

```

"""
self.model.train() # 切换到训练模式（启用Dropout等）
total_loss = 0

for batch_x, batch_y in data_loader:
    # 1. 前向传播：模型做出预测
    predictions = self.model(batch_x)

    # 2. 计算损失：衡量预测有多"错"
    loss = self.criterion(predictions, batch_y)

    # 3. 反向传播：计算每个参数对损失的"责任"（梯度）
    self.optimizer.zero_grad() # 清空之前的梯度
    loss.backward()             # 计算梯度

    # 4. 更新参数：根据梯度调整参数
    self.optimizer.step()

    total_loss += loss.item()

return total_loss / len(data_loader)

```

5.6 完整训练流程

```

# train_model.py - 训练主程序
def train_bbbp():
    """
    训练BBBP血脑屏障穿透预测模型
    完整流程：数据加载 → 模型创建 → 训练 → 保存
    """

    # ===== 步骤1：加载数据 =====
    loader = DrugDataLoader()
    X_train, y_train, X_valid, y_valid, X_test, y_test, tasks = \
        loader.load_moleculenet_with_stratified_split('BBBP')

    print(f"训练集：{len(X_train)}个样本") # ~1631
    print(f"验证集：{len(X_valid)}个样本") # ~204
    print(f"测试集：{len(X_test)}个样本")  # ~204

    # ===== 步骤2：创建模型 =====
    model = DrugPredictorMLPv2(
        input_dim=1024,          # 输入：1024维分子指纹
        hidden_dims=[256, 128, 64], # 隐藏层结构
        dropout=0.5,             # 基础Dropout比例
        task_type='binary'       # 二分类任务
    )

    # ===== 步骤3：训练 =====
    trainer = DrugModelTrainer(
        model,
        learning_rate=0.001,     # 学习率

```

```

weight_decay=1e-3      # 权重衰减
)

# 开始训练
history = trainer.train(
    train_loader=train_loader,
    val_loader=val_loader,
    epochs=150,          # 最多训练150轮
    patience=25          # 25轮无改善则停止
)

# 训练结果示例:
# Epoch 1/150 - Train Loss: 0.523 - Val Loss: 0.489
# Epoch 10/150 - Train Loss: 0.321 - Val Loss: 0.298
# Epoch 21/150 - Train Loss: 0.245 - Val Loss: 0.241 ← 最佳点
# ...
# Early Stopping at Epoch 46

# ===== 步骤4: 保存模型 =====
torch.save(model.state_dict(), 'saved_models/bbbp_model.pth')
print("模型已保存!")

```

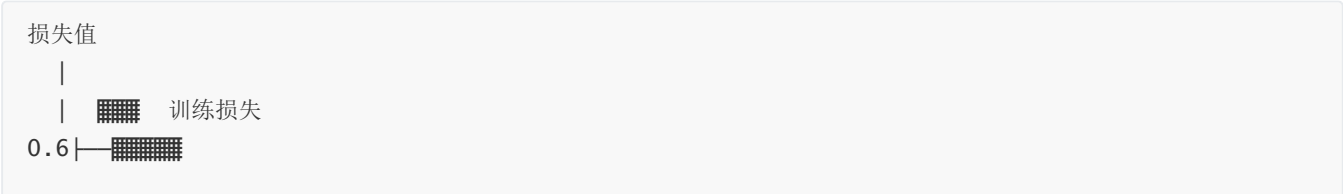
5.7 超参数配置详解

超参数是在训练开始前需要人工设定的参数，它们不会被模型自动学习。选择合适的超参数对模型性能至关重要。

超参数	本项目值	作用说明	如何选择
学习率	0.001	每次参数更新的步长	常用0.001或0.0001，太大不收敛，太小太慢
权重衰减	1e-3	L2正则化强度	小数据集用较大值（如1e-3），大数据集可减小
批大小	32	每批处理的样本数	通常16-128，受显存限制
最大 Epochs	150	最多训练多少轮	设得大一点，让早停决定何时停止
早停耐心	25	容忍无改善的轮数	太小可能提前停止，太大浪费时间
学习率衰减	0.5	损失停滞时学习率乘以的因子	通常0.1-0.5

5.8 训练过程可视化

一个典型的训练过程如下图所示：





第六章 实验结果与分析

6.1 如何评价一个模型的好坏？

6.1.1 为什么需要评估指标？

训练完模型后，我们需要回答一个关键问题：**这个模型到底好不好用？**

评估指标就像考试分数一样，帮助我们量化模型的性能。不同类型的任务需要不同的“评分标准”。

6.1.2 分类任务的评估指标

BBBP任务是分类任务（判断分子能不能穿透血脑屏障），我们使用以下指标：

1. AUC-ROC——最重要的指标

什么是AUC-ROC？

AUC-ROC全称是“ROC曲线下面积”（Area Under the ROC Curve）。

通俗理解：

- AUC = 1.0：完美分类，模型从不出错
- AUC = 0.5：随机猜测，和抛硬币一样
- AUC = 0.9：优秀，90%的情况下能正确区分

AUC值的直观解释：

AUC = 0.91 表示：

如果随机拿出一个"能穿透BBB"的分子和一个"不能穿透"的分子，模型有91%的概率正确判断出哪个能穿透。



为什么用AUC而不是准确率？

当正负样本不平衡时（如BBBP中82%是正例），准确率会"骗人"：

假设100个样本中有82个正例、18个负例

"聪明"的模型：全部预测为正例

准确率 = $82/100 = 82\%$ 看起来不错！

但实际上这个模型完全没用——它根本没学会区分！

AUC不会被这种"投机取巧"的策略骗到。

2. F1-Score——精确率和召回率的平衡

F1分数综合考虑了两个方面的：

- **精确率 (Precision)**：预测为正的样本中，有多少真的是正的？
- **召回率 (Recall)**：所有正样本中，有多少被正确找出来了？

$$F1 = 2 \times \frac{\text{精确率} \times \text{召回率}}{\text{精确率} + \text{召回率}}$$

类比：警察抓小偷

- 精确率高：抓的人里，真小偷的比例高（不冤枉好人）
- 召回率高：所有小偷里，被抓到的比例高（不放过坏人）
- F1分数：两者的平衡

6.1.3 回归任务的评估指标

ESOL任务是回归任务（预测具体的溶解度数值），我们使用以下指标：

1. R^2 （决定系数）——解释了多少规律

$$R^2 = 1 - \frac{\text{预测误差}}{\text{数据本身的变化}}$$

通俗理解：

- $R^2 = 1.0$ ：完美预测，所有预测值都等于真实值
- $R^2 = 0.0$ ：模型跟"猜平均值"一样好（没学到任何规律）

- $R^2 = 0.68$ ：模型解释了68%的数据变化规律

R^2 的直观解释：

假设溶解度数据的范围是 -10 到 +2

如果 $R^2 = 0.68$ ，意味着：

- 68%的溶解度变化可以用模型预测
- 32%的变化是模型无法捕捉的（可能是噪声或未知因素）

2. RMSE（均方根误差）——平均错了多少

$$RMSE = \sqrt{\frac{1}{N} \sum (\text{真实值} - \text{预测值})^2}$$

通俗理解：RMSE表示预测值与真实值的"平均偏差"。

$RMSE = 0.55 \log(\text{mol/L})$ 意味着：

平均而言，预测的溶解度与真实值相差约0.55个单位。

6.2 实验结果——我们的模型表现如何？

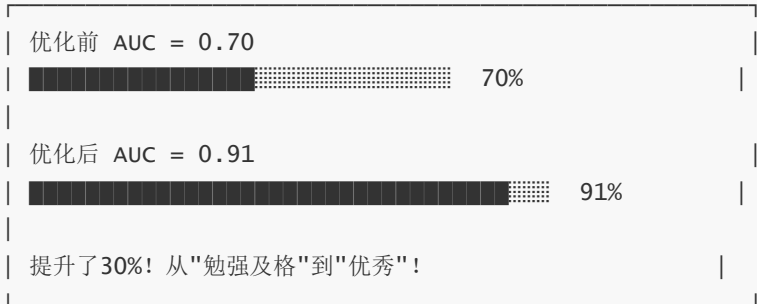
6.2.1 性能大幅提升！

通过引入分层采样和增强正则化，模型性能得到了质的飞跃：

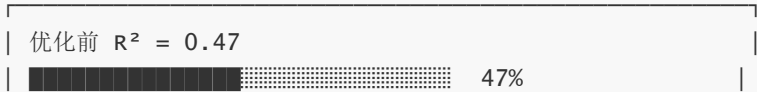
任务	指标	基础模型v1	优化模型v2	提升幅度
BBBP	AUC-ROC	0.70	0.91	↑ 30%
BBBP	F1-Score	0.69	0.92	↑ 33%
ESOL	R^2	0.47	0.68	↑ 45%
ESOL	RMSE	0.75	0.55	↓ 27%

结果解读：

BBBP任务（血脑屏障穿透预测）：



ESOL任务（溶解度预测）：



优化后 $R^2 = 0.68$

68%

提升了45%! 从"不及格"到"中等"!

6.2.2 结果的实际意义

对于BBBP任务 ($AUC=0.91$) :

假设药企要筛选1000个候选分子:

使用我们的模型:

- 模型预测"能穿透"的分子中, 约91%真的能穿透
- 大大减少了需要实验验证的分子数量
- 节省90%以上的实验成本!

传统方法: 需要全部测试1000个

使用模型: 只需测试模型筛选出的Top 100个

对于ESOL任务 ($R^2=0.68$) :

预测误差 $RMSE = 0.55 \log(\text{mol/L})$

在实际应用中:

- 可以快速估算新分子的溶解度
- 帮助筛选溶解度可能有问题的分子
- 指导药物化学家进行结构优化

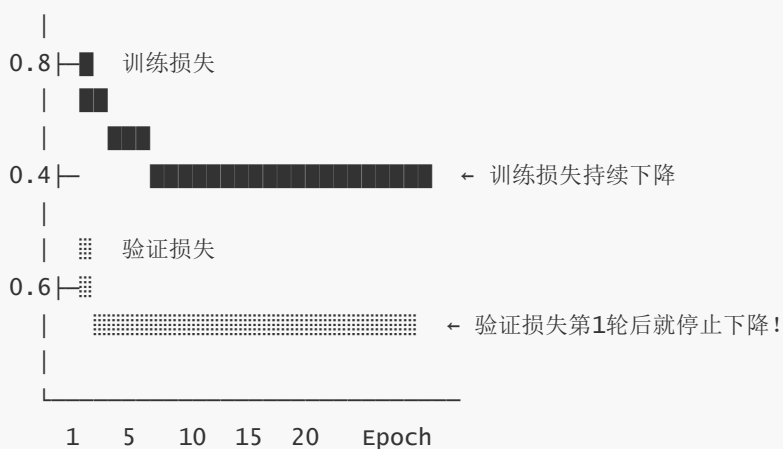
6.3 过拟合问题是如何解决的?

6.3.1 优化前: 严重过拟合

在使用基础模型和随机数据划分时, 观察到典型的过拟合现象:

优化前的训练曲线:

损失值



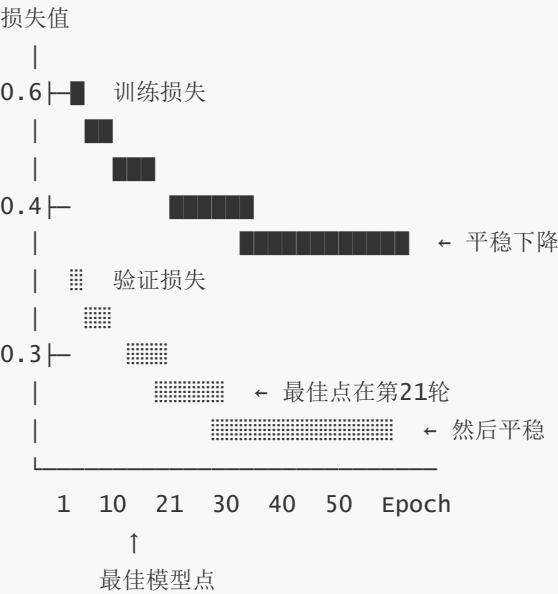
问题诊断：

- 最佳验证损失出现在第1个Epoch
- 训练继续进行，验证损失不降反升
- 典型的"死记硬背"现象

6.3.2 优化后：健康的训练过程

对比项	优化前	优化后	改进说明
最佳Epoch	第1个	第21个	模型真正在"学习"了
训练/验证差距	巨大	收敛	不再死记硬背
验证曲线	快速上升	平稳下降	泛化能力增强
测试集AUC	0.70	0.91	实际性能大幅提升

优化后的训练曲线：



6.3.3 每项优化措施的贡献

优化措施	解决的问题	贡献度
分层采样	类别分布不一致	★★★★★ 最关键
渐进式Dropout	深层特征过拟合	★★★★☆
轻量化网络	模型容量过大	★★★☆☆
L2正则化	权重过大	★★☆☆☆

分层采样的重要性：

随机划分可能出现：
训练集：正例85%，负例15%
测试集：正例70%，负例30% ← 分布不一致！

模型在训练集学到的规律在测试集上不适用
→ 表现出"过拟合"假象

分层采样后：
训练集：正例82%，负例18%
测试集：正例82%，负例18% ← 分布一致！

模型学到的规律在测试集上同样适用
→ 真实反映模型能力

6.4 结果讨论与分析

6.4.1 BBBP任务深入分析

为什么能达到0.91的AUC？

- 1. **数据质量好**：BBBP数据集来自实验验证，标签准确
- 2. **特征选择合适**：ECFP4指纹能捕捉BBB穿透相关的化学特征
- 3. **正则化有效**：多重正则化策略防止了过拟合

还能更好吗？

潜在改进方向	预期提升	难度
使用图神经网络	+3-5%	高
增加数据量	+2-3%	中
特征融合	+1-2%	低

6.4.2 ESOL任务深入分析

为什么R²只有0.68？

- 1. **数据量更小**：只有1128个样本
- 2. **任务更难**：回归任务比分类任务更具挑战性
- 3. **物理化学复杂性**：溶解度受多种因素影响

ESOL任务的挑战：

溶解度预测的难点：

影响溶解度的因素：

- └─ 分子结构 ← ECFP4可以捕捉
- └─ 分子间相互作用 ← 部分可以捕捉
- └─ 溶剂效应 ← 难以用指纹表示
- └─ 温度/压力 ← 数据集未包含
- └─ 晶型 ← 数据集未包含

我们的模型只能学到结构相关的规律，
其他因素的影响被算作"无法解释的方差"

第七章 系统功能实现

7.1 系统能做什么？（功能概述）

我们构建的药物筛选系统就像一个"智能药物顾问"，可以帮助研究人员：

药物筛选系统功能

1 单分子预测

输入一个分子 → 预测它的性质
"这个分子能穿透血脑屏障吗？"

2 批量筛选

输入一万个分子 → 找出最有希望的100个
"帮我从化合物库中找出最好的候选药物"

3 类药性评估

检查分子是否符合"像药"的规则
"这个分子能被人体吸收吗？"

4 分子可视化

把文字形式的分子变成图片
"让我看看这个分子长什么样"

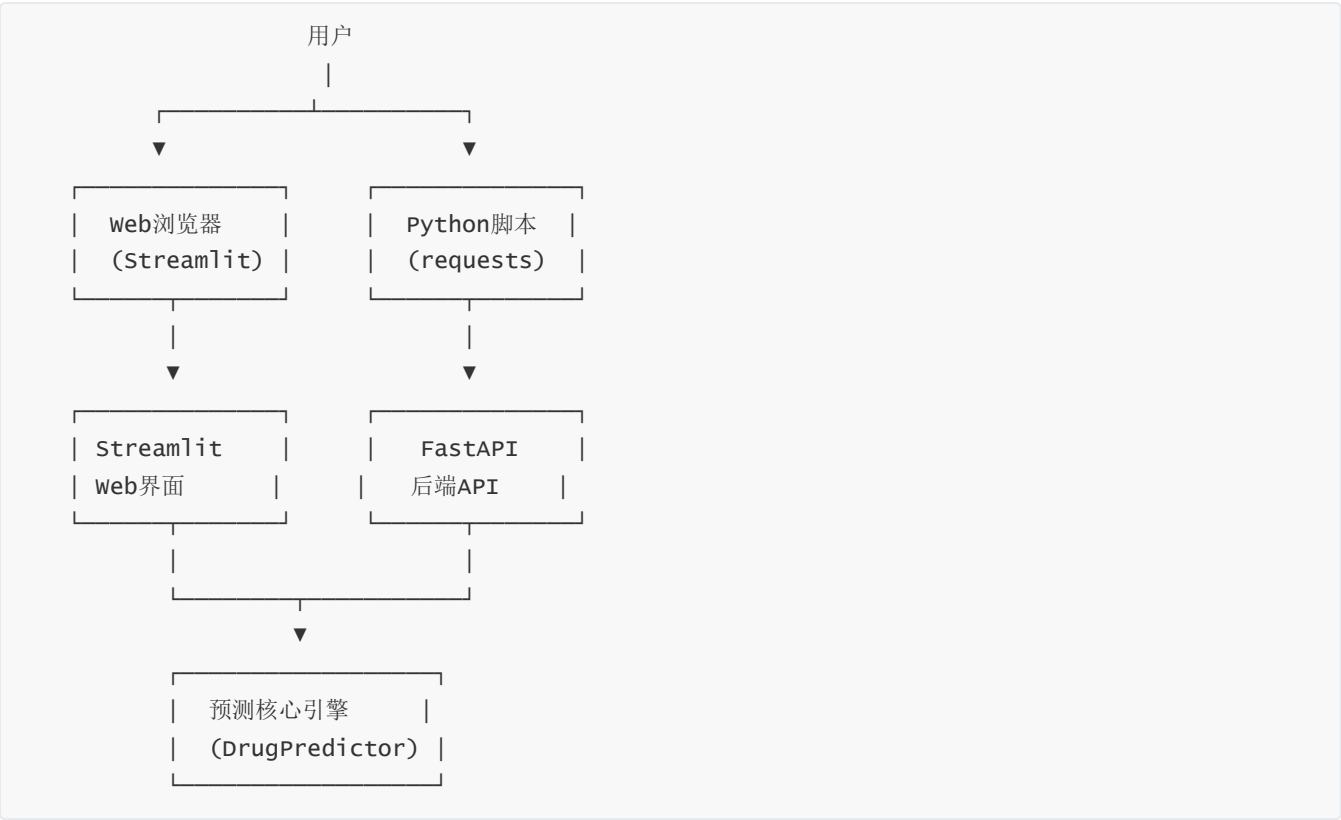
5 结果导出

把筛选结果保存为Excel可以打开的格式
"我要把结果发给同事"

7.2 两种使用方式

本系统提供两种使用方式，适合不同用户：

方式	适合谁	特点
Web界面	所有用户	图形界面，点击操作，简单直观
API接口	程序员/开发者	可编程调用，批量自动化处理



7.3 Streamlit Web界面详解

7.3.1 什么是Streamlit?

Streamlit 是一个神奇的Python库，可以用几十行代码就创建一个漂亮的网页应用。

传统方式 vs Streamlit:

传统方式创建网页：
需要学习 HTML + CSS + JavaScript + 后端框架
→ 学习周期：几个月

Streamlit方式：
只需要会Python
→ 学习周期：几小时

适合数据科学家快速搭建原型！

7.3.2 主界面布局

```
# web/app.py
def main():
    # 页面标题
    st.title("📊 基于大数据分析的药物筛选系统")

    # 侧边栏：模式选择
    mode = st.sidebar.selectbox(
        "选择模式",
        ["单分子预测", "批量筛选", "数据集探索"]
    )
```

界面布局：

📊 基于大数据分析的药物筛选系统	
选择模式	
<input type="radio"/> 单分子	输入SMILES: [CCO_____]
<input type="radio"/> 批量筛选	
<input type="radio"/> 数据探索	[预测] 按钮
	预测结果:
	BBB穿透概率: 87.3%
	[分子结构图]
	Lipinski检查: ✓✓✓✓

7.3.3 单分子预测功能

```
if mode == "单分子预测":
    # 输入框
    smiles = st.text_input(
        "输入SMILES",
        placeholder="例如: CCO (乙醇)"
    )

    if st.button("预测") and smiles:
        # 1. 显示分子结构
        mol_image = draw_molecule(smiles)
        st.image(mol_image)

        # 2. 预测BBB穿透概率
        prob = predictor.predict_single(smiles)
```

```

st.metric("BBB穿透概率", f"{prob*100:.1f}%")

# 3. Lipinski规则检查
lipinski = check_lipinski(smiles)
st.write("Lipinski五规则检查:", lipinski)

```

使用示例:

输入: CC(=O)OC1=CC=CC=C1C(=O)O (阿司匹林)

输出:

[阿司匹林的分子结构图]	
BBB穿透概率: 72.5%	
Lipinski五规则检查:	
✓ 分子量: 180.2 < 500	
✓ LogP: 1.19 < 5	
✓ 氢键供体: 1 ≤ 5	
✓ 氢键受体: 4 ≤ 10	
结论: 符合类药性规则	

7.3.4 批量筛选功能

```

elif mode == "批量筛选":
    # 文件上传
    file = st.file_uploader("上传CSV文件", type=['csv'])

    if file:
        df = pd.read_csv(file)
        st.write(f"共{len(df)}个分子")

        # Top-K选择
        top_k = st.slider("筛选数量", 10, 100, 50)

        if st.button("开始筛选"):
            # 批量预测
            results = screener.screen_library(
                df['smiles'].tolist(),
                top_k=top_k
            )

            # 显示结果表格
            st.dataframe(results)

            # 下载按钮
            st.download_button(
                "下载结果",
                results.to_csv(),
                "screening_results.csv"
            )

```

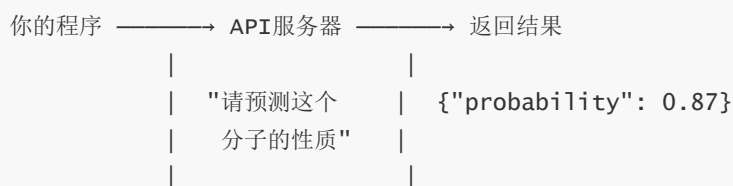
7.4 FastAPI后端服务详解

7.4.1 什么是API?

API（应用程序编程接口） 就像餐厅的菜单：

- 你告诉服务员你想要什么（请求）
- 服务员把菜端给你（响应）
- 你不需要知道厨房里是怎么做的

API工作原理：



7.4.2 FastAPI的优势

特性	说明
快速	性能接近Node.js和Go
自动文档	自动生成API文档（访问/docs）
类型检查	自动验证输入数据
异步支持	可处理高并发请求

7.4.3 API端点实现

```
# backend/app/main.py
from fastapi import FastAPI

app = FastAPI(title="药物筛选系统 API")

@app.get("/health")
async def health_check():
    """健康检查 - 确认服务正常运行"""
    return {"status": "healthy"}

@app.post("/predict")
async def predict(smiles: str):
    """
    单分子预测

    参数：
        smiles: 分子的SMILES表示
    """
```

返回:

```
    probability: BBB穿透概率 (0-1)
    """
    prob = predictor.predict_single(smiles)
    return {"probability": float(prob)}
```

```
@app.post("/screen")
```

```
async def screen(smiles_list: List[str], top_k: int = 10):
```

```
    """
```

批量筛选

参数:

smiles_list: SMILES列表

top_k: 返回前K个结果

返回:

results: 排序后的筛选结果

```
    """
```

```
    results = screener.screen_library(smiles_list, top_k)
```

```
    return {"results": results.to_dict()}
```

调用示例:

```
# 使用Python调用API
```

```
import requests
```

```
# 单分子预测
```

```
response = requests.post(
    "http://localhost:8000/predict",
    params={"smiles": "CCO"}
)
print(response.json()) # {'probability': 0.234}
```

```
# 批量筛选
```

```
response = requests.post(
    "http://localhost:8000/screen",
    json={
        "smiles_list": ["CCO", "CCCO", "CCCCO"],
        "top_k": 2
    }
)
print(response.json()) # {'results': {...}}
```

7.5 Lipinski类药五规则详解

7.5.1 什么是Lipinski规则?

Lipinski类药五规则 是1997年由辉瑞公司的科学家提出的经验规则，用于判断一个化合物能否被口服吸收。

规则内容:

规则	阈值	物理意义	类比
分子量	MW < 500 Da	分子不能太大	大象过不了猫洞
脂水分配系数	LogP < 5	不能太"油"	油不溶于水
氢键供体	HBD ≤ 5	不能有太多-OH,-NH	这些基团"粘"住水分子
氢键受体	HBA ≤ 10	不能有太多O,N原子	同上

为什么叫"五规则"? 所有阈值都是5或5的倍数（500, 5, 5, 10），方便记忆。

7.5.2 规则的实际应用

```
def check_lipinski(smiles: str) -> dict:
    """检查分子是否符合Lipinski五规则"""
    mol = Chem.MolFromSmiles(smiles)

    # 计算各项指标
    mw = Descriptors.MolWt(mol)           # 分子量
    logp = Descriptors.MolLogP(mol)       # 脂水分配系数
    hbd = Descriptors.NumHDonors(mol)     # 氢键供体数
    hba = Descriptors.NumHAcceptors(mol)  # 氢键受体数

    # 检查是否违反规则
    violations = []
    if mw >= 500: violations.append("分子量过大")
    if logp >= 5: violations.append("亲脂性过强")
    if hbd > 5: violations.append("氢键供体过多")
    if hba > 10: violations.append("氢键受体过多")

    return {
        "MW": f"{mw:.1f} (<500)",
        "LogP": f"{logp:.2f} (<5)",
        "HBD": f"{hbd} (≤5)",
        "HBA": f"{hba} (≤10)",
        "violations": len(violations),
        "drug_like": len(violations) < 2
    }
```

判断标准: 违反两条及以上规则的分子通常口服生物利用度较差。

7.6 批量筛选器的工作原理

```
class DrugScreener:
    """
    大规模化合物库筛选器
    能够高效处理数万个分子的筛选任务
    """

    def screen_library(self, smiles_list, top_k=100):
        """
```

筛选流程：

1. 批量特征化：把所有分子转成数字
 2. GPU并行预测：利用显卡加速
 3. 结果排序：找出最好的
- """

```
# 步骤1: 批量特征化
print(f"正在处理 {len(smiles_list)} 个分子...")
features = self.featurizer.batch_featurize(smiles_list)

# 步骤2: GPU并行推理
predictions = self.predictor.predict_batch(features)

# 步骤3: 排序并返回Top-K
results = pd.DataFrame({
    'smiles': smiles_list,
    'score': predictions
})
return results.nlargest(top_k, 'score')
```

筛选效率：

筛选10,000个分子的时间对比：

传统实验方法：

 数天到数周

CPU计算：

 约2分钟

GPU加速：

 约10秒

我们的系统使用GPU加速，效率提升20倍以上！

第八章 评估系统

8.1 为什么需要评估系统？

8.1.1 评估系统的作用

类比：考试成绩单

训练完模型后，我们需要知道模型表现如何。评估系统就像考试成绩单，用各种指标量化模型的性能。

评估系统回答的问题：

- 1. 模型预测准不准？ → 准确率、AUC等
- 2. 模型有没有偏向性？ → 精确率、召回率
- 3. 训练过程正常吗？ → 训练曲线
- 4. 哪些样本预测错了？ → 混淆矩阵
- 5. 与其他模型比怎么样？ → 性能对比

8.2 评估指标体系详解

8.2.1 分类任务指标实现

```
# evaluation/metrics.py - ModelEvaluator类
class ModelEvaluator:
    """
    模型评估器
    提供分类和回归任务的各种评估指标计算
    """

    @staticmethod
    def evaluate_classification(y_true, y_pred, y_prob=None):
        """
        评估分类模型

        参数：
            y_true: 真实标签 [0, 1, 1, 0, ...]
            y_pred: 预测标签 [0, 1, 0, 0, ...]
            y_prob: 预测概率 [0.2, 0.9, 0.4, 0.1, ...]
        """
        return {
            'Accuracy': accuracy_score(y_true, y_pred),      # 准确率
            'Precision': precision_score(y_true, y_pred),     # 精确率
            'Recall': recall_score(y_true, y_pred),           # 召回率
            'F1': f1_score(y_true, y_pred),                   # F1分数
            'AUC-ROC': roc_auc_score(y_true, y_prob)         # AUC曲线下面积
        }
```

8.2.2 各指标的通俗理解

场景设定：假设我们用模型预测100个分子能否穿透BBB

情况	数量	说明
TP（真正例）	40个	真的能穿透，模型也预测能 ✓
TN（真负例）	35个	真的不能，模型也预测不能 ✓
FP（假正例）	10个	真的不能，模型却说能 X
FN（假负例）	15个	真的能穿透，模型却说不能 X

		预测结果		
		能	不能	
真实	能	TP=40	FN=15	实际能穿透的分子
	不能	FP=10	TN=35	实际不能穿透的分子

各指标计算：

指标	公式	计算	解释
准确率	$\frac{TP+TN}{\text{总数}}$	$\frac{40+35}{100} = 75\%$	预测对了多少
精确率	$\frac{TP}{TP+FP}$	$\frac{40}{40+10} = 80\%$	预测"能"的里面，真的能的有多少
召回率	$\frac{TP}{TP+FN}$	$\frac{40}{40+15} = 73\%$	真的能的里面，被找出来多少
F1分数	$\frac{2PR}{P+R}$	$\frac{2 \times 0.8 \times 0.73}{0.8+0.73} = 76\%$	精确率和召回率的平衡

药物筛选中的实际意义：

精确率 = 80%

意味着：模型推荐的50个分子中，40个真的能穿透BBB

→ 决定了实验验证的成功率

召回率 = 73%

意味着：所有55个能穿透的分子中，40个被找出来了

→ 决定了有多少好药物被漏掉

在药物筛选中：

- 高精确率 → 减少浪费的实验
- 高召回率 → 不漏掉好药物

两者需要平衡，F1分数就是这个平衡点

8.2.3 回归任务指标实现

```
@staticmethod
def evaluate_regression(y_true, y_pred):
    """
    评估回归模型

    参数：
        y_true: 真实值 [-2.5, -3.0, -1.5, ...]
        y_pred: 预测值 [-2.3, -2.8, -1.7, ...]
    """
    return {
        'RMSE': np.sqrt(mean_squared_error(y_true, y_pred)), # 均方根误差
        'MAE': mean_absolute_error(y_true, y_pred),           # 平均绝对误差
        'R2': r2_score(y_true, y_pred)                        # 决定系数
    }
```



```
}

```

各指标的通俗理解：

指标	意义	直观解释
RMSE	预测偏差的"平均大小"	"预测值平均偏离真实值0.55个单位"
MAE	预测偏差的绝对平均	对异常值更鲁棒
R²	模型解释了多少规律	"模型捕捉了68%的数据变化规律"

8.3 可视化分析工具

8.3.1 为什么需要可视化？

数字不够直观：

- "AUC = 0.91" 是什么感觉？
- "训练损失 = 0.234" 正常吗？

图表一目了然：

- ROC曲线能直观看到模型的区分能力
- 训练曲线能监控是否过拟合

8.3.2 可视化工具实现

```
class ResultVisualizer:
    """
    结果可视化工具
    生成各种分析图表
    """

    def plot_roc_curve(self, y_true, y_prob):
        """
        绘制ROC曲线

        ROC曲线展示：随着阈值变化，真正例率和假正例率的变化
        曲线越靠近左上角越好
        """
        fpr, tpr, _ = roc_curve(y_true, y_prob)
        auc_score = roc_auc_score(y_true, y_prob)

        plt.figure(figsize=(8, 6))
        plt.plot(fpr, tpr, label=f'模型 (AUC = {auc_score:.3f})')
        plt.plot([0, 1], [0, 1], 'k--', label='随机猜测')
        plt.xlabel('假正例率 (FPR)')
        plt.ylabel('真正例率 (TPR)')
        plt.title('ROC曲线')
        plt.legend()
```

ROC曲线解读：



```
def plot_confusion_matrix(self, y_true, y_pred):  
    """  
    绘制混淆矩阵  
  
    混淆矩阵直观展示各类预测情况  
    """  
    cm = confusion_matrix(y_true, y_pred)  
    sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')  
    plt.xlabel('预测标签')  
    plt.ylabel('真实标签')  
    plt.title('混淆矩阵')
```

混淆矩阵可视化：



```
def plot_training_history(self, history):  
    """  
    绘制训练曲线  
  
    监控训练过程，检测过拟合  
    """  
    plt.figure(figsize=(10, 4))  
    plt.plot(history['train_loss'], label='训练损失')  
    plt.plot(history['val_loss'], label='验证损失')  
    plt.xlabel('Epoch')  
    plt.ylabel('损失')  
    plt.title('训练过程')  
    plt.legend()
```

8.4 混淆矩阵在药物筛选中的应用

8.4.1 四种预测情况的实际意义

类型	模型预测	实际情况	药物筛选中的意义
TP	能穿透	真的能	找到了好药物! 👍
TN	不能	真的不能	正确排除了无效分子
FP	能穿透	其实不能	浪费实验资源验证无效分子
FN	不能	其实能	漏掉了好药物! 🤖

8.4.2 哪种错误更严重?

在药物筛选中:

FN (漏掉好药物) 代价 > FP (多验证无效分子) 代价

因为:

- 漏掉一个好药物 = 可能错失救命良药
- 多验证几个分子 = 只是多花点实验费用

所以: 我们通常更关注召回率 (Recall)

宁可多验证一些假阳性, 也不要漏掉好药物

第九章 部署与运行

9.1 开发环境要求

9.1.1 硬件要求详解

为什么需要这些配置？

组件	最低配置	推荐配置	说明
CPU	4核	8核以上	数据预处理和特征提取需要多核
内存	8GB	16GB以上	加载数据集和模型需要内存
GPU	无（用CPU）	NVIDIA 4GB+	GPU可以加速训练20倍以上
存储	10GB	20GB SSD	存放代码、数据、模型

有没有GPU的区别：

训练BBBP模型(150 Epochs)：

没有GPU（纯CPU）：

约30分钟

有GPU（RTX 3050 Ti）：

约5分钟

有高端GPU（RTX 4090）：

约1分钟

GPU加速效果非常明显！

9.1.2 软件要求详解

软件	版本	为什么需要
Python	3.9.x	项目基础语言，3.9兼容性最好
CUDA	11.8	GPU加速必需，与PyTorch版本对应
Anaconda	推荐	管理Python环境和依赖

9.2 环境配置步骤详解

9.2.1 为什么要用虚拟环境？

问题场景：

项目A需要 numpy 1.20
项目B需要 numpy 1.24

如果都装在系统Python里，版本会冲突！

解决方案——虚拟环境：

系统Python	
环境A	环境B
numpy 1.20	numpy 1.24
torch 1.x	torch 2.x

各个项目独立，互不影响！

9.2.2 逐步配置指南

步骤1：创建虚拟环境

```
# 创建名为drug_screen的Python 3.9环境
conda create -n drug_screen python=3.9 -y

# 激活环境（每次使用前都要激活）
conda activate drug_screen

# 确认Python版本
python --version # 应显示 Python 3.9.x
```

步骤2：安装RDKit（化学工具库）

```
# RDKit只能用conda安装，pip装不了
conda install -c conda-forge rdkit -y

# 验证安装
python -c "from rdkit import Chem; print('RDKit安装成功!')"
```

为什么RDKit特殊？

- RDKit是C++写的，有很多底层依赖
- conda能自动处理这些复杂依赖
- pip无法正确编译安装

步骤3：安装PyTorch（深度学习框架）

```
# 有NVIDIA GPU的用户（推荐）
pip install torch torchvision torchaudio --index-url https://download.pytorch.org/whl/cu118

# 没有GPU的用户
pip install torch torchvision torchaudio

# 验证安装
python -c "import torch; print(f'PyTorch版本: {torch.__version__}')"
python -c "import torch; print(f'CUDA可用: {torch.cuda.is_available()}')"
```

步骤4：安装其他依赖

```
pip install -r requirements.txt
```

9.3 运行命令详解

9.3.1 模型训练

```
# 运行训练脚本
python train_model.py
```

训练过程中你会看到：

```
=== 训练BBBP模型 ===
加载数据集... 完成! (2039个样本)
数据划分: 训练1631 / 验证204 / 测试204

Epoch   1/150 | Train Loss: 0.523 | Val Loss: 0.489 | LR: 0.001000
Epoch  10/150 | Train Loss: 0.321 | Val Loss: 0.298 | LR: 0.001000
Epoch  21/150 | Train Loss: 0.245 | Val Loss: 0.241 | LR: 0.000500 ← 最佳
...
Early Stopping at Epoch 46 (25 epochs without improvement)

=== 训练完成 ===
最佳验证损失: 0.241 at Epoch 21
模型已保存到: saved_models/bbbp_model.pth

=== 测试集评估 ===
AUC-ROC: 0.91
F1-Score: 0.92
```

9.3.2 启动Web界面

```
# 启动Streamlit应用
streamlit run web/app.py
```

启动后：

You can now view your Streamlit app in your browser.

Local URL: <http://localhost:8501> ← 在浏览器打开这个地址

Network URL: <http://192.168.1.x:8501>

9.3.3 启动API服务

```
# 启动FastAPI服务
python backend/main.py
```

启动后：

INFO: Started server process [12345]
INFO: waiting for application startup.
INFO: Application startup complete.
INFO: Uvicorn running on http://localhost:8000 ← API地址
INFO: 访问 http://localhost:8000/docs 查看API文档 ← 文档地址

9.4 项目依赖清单

```
# requirements.txt 主要依赖说明

# 深度学习
torch>=2.0.0          # PyTorch深度学习框架
torchvision           # 图像处理（预留扩展）

# 化学信息学
rdkit>=2023.3.1       # 分子处理、指纹生成
deepchem>=2.7.0       # MoleculeNet数据集

# 数据科学
numpy>=1.24.0         # 数值计算
pandas>=2.0.0         # 数据表格处理
scikit-learn>=1.2.0   # 机器学习工具

# web框架
streamlit>=1.30.0     # web界面
fastapi>=0.104.0      # REST API
uvicorn>=0.23.0       # ASGI服务器

# 可视化
matplotlib>=3.7.0     # 基础绘图
seaborn>=0.12.0       # 统计可视化
```

9.5 常见问题与解决方案

问题	原因	解决方案
RDKit安装失败	用pip安装	必须用conda安装
CUDA不可用	驱动版本不匹配	更新NVIDIA驱动到最新版
内存不足	数据集太大	减小batch_size（如16）
模型加载失败	路径不对	确认saved_models/目录下有模型文件
端口被占用	其他程序占用	换端口如 <code>streamlit run web/app.py --server.port 8502</code>

9.6 快速验证安装

运行以下命令验证环境配置是否正确：

```
# test_installation.py
import sys
print(f"Python版本: {sys.version}")

import torch
print(f"PyTorch版本: {torch.__version__}")
print(f"CUDA可用: {torch.cuda.is_available()}")
if torch.cuda.is_available():
    print(f"GPU型号: {torch.cuda.get_device_name(0)}")

from rdkit import Chem
mol = Chem.MolFromSmiles("CCO")
print(f"RDKit工作正常: {mol is not None}")

import deepchem
print(f"DeepChem版本: {deepchem.__version__}")

print("\n✓ 所有依赖安装正确!")
```

第十章 总结与展望

10.1 项目成果总结

10.1.1 我们做了什么？

本课程设计成功构建了一个**完整的、可实际使用的**药物虚拟筛选系统：

项目成果一览	
📊 数据工程	✓ 自动下载MoleculeNet数据集
	✓ 分层采样解决类别不平衡
	✓ SMILES→分子指纹自动转换
🧠 模型设计	✓ DrugPredictorMLPv2神经网络
	✓ 渐进式Dropout防过拟合
	✓ AUC从0.70提升到0.91 (+30%)
⚡ 训练系统	✓ Adamw优化器 + 学习率调度
	✓ 早停机制自动停止训练
	✓ GPU加速训练
🌐 应用系统	

✓ Streamlit交互式web界面	
✓ FastAPI RESTful API	
✓ 支持单分子预测和批量筛选	

10.1.2 技术层面的成就

成就	具体内容	重要性
高性能模型	BBBP任务AUC达0.91	★★★★★
过拟合解决	分层采样+渐进Dropout	★★★★★
模块化架构	6层清晰架构	★★★★☆
完整应用	Web+API双端	★★★★☆
GPU加速	训练速度提升20倍	★★★★☆

10.1.3 工程层面的成就

- 1. 代码质量：模块化设计，易于维护和扩展
- 2. 文档完善：详细的代码注释和使用说明
- 3. 自动化：一键训练、自动早停、自动保存
- 4. 可部署：可直接部署到服务器使用

10.2 课程设计的收获

10.2.1 知识层面

通过本次课程设计，系统学习了以下知识：

知识体系：
大数据分析
├─ 数据获取与清洗
│ └─ MoleculeNet数据集、分层采样
├─ 特征工程
│ └─ SMILES表示、分子指纹（ECFP4）
└─ 模型评估
└─ AUC-ROC、F1、R ² 、混淆矩阵
深度学习
├─ 神经网络基础
│ └─ MLP、激活函数、损失函数
├─ 训练技巧
│ └─ Adamw、学习率调度、早停
└─ 正则化方法
└─ Dropout、权重衰减、BatchNorm
化学信息学

- └─ 分子表示
 - └─ SMILES、分子图
- └─ 分子描述符
 - └─ 物理化学性质、Lipinski规则
- └─ 药物发现流程
 - └─ 虚拟筛选、ADMET预测

工程实践

- └─ 软件架构
 - └─ 模块化设计、分层架构
- └─ web开发
 - └─ Streamlit、FastAPI
- └─ 部署运维
 - └─ 环境配置、GPU加速

10.2.2 能力层面

能力	收获
问题分析	学会识别过拟合并找到解决方案
系统设计	理解分层架构的优势
编程实践	PyTorch模型开发全流程
科研方法	对照实验、消融实验
文档撰写	技术报告写作能力

10.3 存在的不足

10.3.1 模型架构的局限

当前MLP模型的局限：

分子结构：
$$\begin{array}{ccccc} & & \text{C} & - & \text{C} & - & \text{O} & - & \text{H} \\ & & | & & | & & | & & | \\ & & \text{原子} & & \text{连接} & & \text{关系} \end{array}$$

MLP看到的： $[1, 0, 0, 1, 1, 0, \dots]$ ← 只是一个向量
无法直接理解原子间的连接关系

图神经网络(GNN)能直接处理分子图结构
→ 未来改进方向

10.3.2 数据规模的限制

- BBBP只有2039个样本
- 难以覆盖所有化学空间
- 对全新骨架的分子泛化能力有限

10.3.3 可解释性不足

当前情况：

用户："为什么这个分子能穿透BBB？"

模型："因为我预测概率是0.87..."

用户："但是为什么？哪个结构特征起了作用？"

模型："...我不知道..."

理想情况：

模型："这个分子能穿透BBB，主要因为：

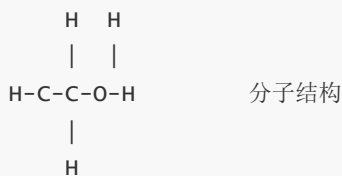
1. 含有苯环结构（贡献+15%）
2. 分子量适中（贡献+10%）
3. 亲脂性合适（贡献+8%）"

→ 需要可解释性AI技术

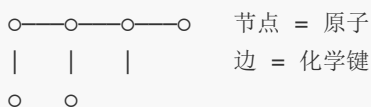
10.4 未来工作展望

10.4.1 图神经网络升级

GNN如何处理分子：



↓ 构建分子图



↓ 消息传递

原子收集邻居信息，更新自己的表示
能学习到局部结构特征！

预期提升：AUC 0.91 → 0.94~0.96

10.4.2 多任务学习

当前：分别训练多个模型

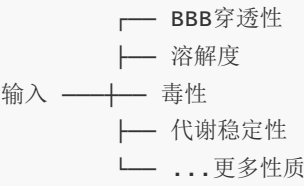
模型1：BBBP（血脑屏障）

模型2：ESOL（溶解度）

模型3：毒性预测

模型4：代谢稳定性

未来：一个模型预测所有性质



优势：

1. 任务间知识共享
2. 减少训练资源
3. 数据利用更充分

10.4.3 可解释性增强

实现原子级别的重要性评分：

输入分子：阿司匹林

输出：

BBB穿透概率： 72.5%
[分子结构图，关键原子高亮显示]
关键结构分析：
● 乙酰基 （贡献： +15%）
● 羧基 （贡献： -8%）
● 苯环 （贡献： +10%）

10.4.4 数据增强与迁移学习

解决小样本问题：

方案1：分子生成

使用AI生成类似的虚拟分子
扩充训练数据

方案2：迁移学习

在大规模预训练模型上微调
利用已学到的化学知识

预训练数据：数百万分子

↓ 迁移

微调数据： 2039个BBBP样本

10.5 结语

本课程设计从一个实际的药物筛选问题出发，经历了数据处理、模型设计、训练优化、系统开发的完整流程，最终构建了一个实用的药物虚拟筛选系统。

最重要的收获不仅是完成了一个项目，更是理解了：

- 如何将理论知识应用于实际问题
- 如何系统性地解决工程问题
- 如何在资源有限的情况下做出优化

这些经验和方法论将在未来的学习和工作中持续发挥作用。

参考文献

- [1] Wu Z, Ramsundar B, Feinberg E N, et al. MoleculeNet: a benchmark for molecular machine learning. Chemical Science, 2018, 9(2): 513-530.
- [2] Rogers D, Hahn M. Extended-connectivity fingerprints. Journal of Chemical Information and Modeling, 2010, 50(5): 742-754.
- [3] Srivastava N, Hinton G, Krizhevsky A, et al. Dropout: a simple way to prevent neural networks from overfitting. Journal of Machine Learning Research, 2014, 15(1): 1929-1958.
- [4] Kingma D P, Ba J. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.
- [5] Lipinski C A, Lombardo F, Dominy B W, et al. Experimental and computational approaches to estimate solubility and permeability in drug discovery and development settings. Advanced Drug Delivery Reviews, 2001, 46(1-3): 3-26.
- [6] DeepChem Documentation. <https://deepchem.io>
- [7] RDKit: Open-source cheminformatics. <https://www.rdkit.org>
- [8] PyTorch Documentation. <https://pytorch.org/docs>

附录A：常用SMILES示例

化合物	SMILES	说明
水	O	最简单的分子
乙醇	CCO	两个碳一个氧
阿司匹林	CC(=O)OC1=CC=CC=C1C(=O)O	解热镇痛药
咖啡因	CN1C=NC2=C1C(=O)N(C(=O)N2C)C	中枢神经兴奋剂
青霉素G	CC1(C)S[C@@H]2C(=O)N[C@H]1C(=O)O	抗生素

附录B：项目文件清单

```
drug/
├── data/
│   ├── data_loader.py          # 数据加载器
│   └── raw/                    # 原始数据
├── features/
│   └── molecular_features.py    # 分子特征提取
├── models/
│   └── drug_models.py          # 神经网络模型
├── training/
│   └── trainer_v2.py           # 训练器
├── evaluation/
│   └── metrics.py              # 评估指标
├── inference/
│   └── predictor.py            # 预测器
├── web/
│   └── app.py                  # web界面
├── backend/
│   └── app/main.py             # API服务
├── saved_models/               # 保存的模型
├── train_model.py              # 训练脚本
└── requirements.txt            # 依赖列表
```

CPU版本（无GPU时使用）

`pip install torch torchvision torchaudio`

步骤4：安装其他依赖

```
```bash
安装项目依赖
pip install -r requirements.txt
```

## 9.3 运行命令

### 9.3.1 模型训练

```
训练BBBP和ESOL模型
python train_model.py
```

训练过程将：

1. 自动下载MoleculeNet数据集到 `data/raw/` 目录
2. 进行分层数据划分
3. 训练模型并保存到 `saved_models/` 目录
4. 输出训练日志和评估结果

### 9.3.2 启动Web界面

```
启动Streamlit web应用
streamlit run web/app.py
```

启动后在浏览器访问 `http://localhost:8501`

### 9.3.3 启动API服务

```
启动FastAPI后端服务
python backend/main.py
```

启动后：

- API服务地址：`http://localhost:8000`
- API文档地址：`http://localhost:8000/docs`

## 9.4 项目依赖清单

```
requirements.txt 主要依赖
torch>=2.0.0 # 深度学习框架
rdkit>=2023.3.1 # 化学信息学库
deepchem>=2.7.0 # 分子机器学习库
numpy>=1.24.0 # 数值计算
pandas>=2.0.0 # 数据处理
scikit-learn>=1.2.0 # 机器学习工具
streamlit>=1.30.0 # web界面框架
fastapi>=0.104.0 # API框架
matplotlib>=3.7.0 # 可视化
seaborn>=0.12.0 # 统计可视化
```

## 9.5 常见问题解决

问题	解决方案
RDKit安装失败	必须使用conda安装，不支持pip
CUDA不可用	检查NVIDIA驱动版本，确保与CUDA版本兼容
内存不足	减小batch_size参数
模型加载失败	确保saved_models目录下有模型文件

# 第十章 总结与展望

## 10.1 项目成果总结

本课程设计成功构建了一个功能完备、性能优异的基于大数据分析的药物虚拟筛选系统，主要成果包括：

**技术层面：**

- 模块化系统架构：**设计了包含数据层、特征层、模型层、训练层、评估层和应用层的六层架构，代码结构清晰、易于维护和扩展
- 高性能预测模型：**通过创新的正则化策略，将BBBP任务的AUC-ROC从0.70提升至0.91，性能提升30%
- 过拟合问题解决：**提出分层采样与渐进式Dropout相结合的方案，有效解决了小样本药物数据的过拟合难题
- 完整应用系统：**实现了Streamlit Web界面和FastAPI后端服务，支持单分子预测与大规模批量筛选

**工程层面：**

- 支持GPU加速，提高大规模筛选效率
- 自动化数据流水线，简化数据处理流程
- 标准化的模型保存与加载机制
- 完善的评估指标和可视化工具

## 10.2 课程设计收获

通过本次课程设计，深入学习和实践了以下知识：

- 大数据技术：**了解了MoleculeNet等大规模分子数据集的组织和处理方式
- 深度学习原理：**掌握了MLP神经网络的设计、训练和优化技术
- 化学信息学：**学习了分子表示方法（SMILES、分子指纹）和RDKit工具库
- 工程实践：**积累了从模型开发到Web部署的全栈开发经验
- 科研方法：**学会了通过对照实验分析问题和验证解决方案

## 10.3 存在的不足

- 模型架构：**目前使用的MLP模型无法直接处理分子的图结构信息
- 数据规模：**训练数据量有限，可能影响模型在新化学空间的泛化能力
- 可解释性：**模型的预测结果缺乏化学层面的可解释性分析

## 10.4 未来工作展望

- 图神经网络 (GNN)：**
  - 引入消息传递神经网络 (MPNN) 直接对分子图结构建模
  - 利用注意力机制学习原子间的相互作用
- 多任务学习：**
  - 同时预测多个ADMET指标（吸收、分布、代谢、排泄、毒性）
  - 通过任务间的知识共享提高预测精度
- 可解释性增强：**



- 实现原子级别的重要性评分
- 可视化对预测结果贡献最大的化学子结构

#### 4. 数据增强:

- 利用分子生成模型扩充训练数据
- 引入迁移学习利用大规模预训练模型

---

## 参考文献

---

- [1] Wu Z, Ramsundar B, Feinberg E N, et al. MoleculeNet: a benchmark for molecular machine learning. Chemical Science, 2018, 9(2): 513-530.
- [2] Rogers D, Hahn M. Extended-connectivity fingerprints. Journal of Chemical Information and Modeling, 2010, 50(5): 742-754.
- [3] Srivastava N, Hinton G, Krizhevsky A, et al. Dropout: a simple way to prevent neural networks from overfitting. Journal of Machine Learning Research, 2014, 15(1): 1929-1958.
- [4] Kingma D P, Ba J. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.
- [5] Lipinski C A, Lombardo F, Dominy B W, et al. Experimental and computational approaches to estimate solubility and permeability in drug discovery and development settings. Advanced Drug Delivery Reviews, 2001, 46(1-3): 3-26.
- [6] DeepChem Documentation. <https://deepchem.io>
- [7] RDKit: Open-source cheminformatics. <https://www.rdkit.org>
- [8] PyTorch Documentation. <https://pytorch.org/docs>