

Datawhale AI夏令营 第三期

一行代码上大分

AutoGluon: 小白上分神器

分享嘉宾: 骆秀韬

Datawhale /01 一行代码?

```
1 from autogluon.tabular import TabularDataset, TabularPredictor
2 import pandas as pd
3
4 train_data = TabularDataset("./data/train.csv")
5 test_data = TabularDataset("./data/test.csv")
6 submit = pd.DataFrame()
7 submit["uuid"] = test_data["uuid"]
8 label = "target"
9
10 predictor = TabularPredictor(
11     label=label,
12     problem_type="binary",
13     eval_metric="f1",
14     ).fit(
15     train_data.drop(columns=["uuid"]),
16     )
17
18 submit[f"{label}"] = predictor.predict(test_data.drop(columns=
19     ["uuid"]))
20 submit.to_csv("submit.csv", index=False)
```

您当前最佳成绩为0.84311提交时间为
(注：以当前所在团队提交的最佳成绩)

具体含义：

- 指定预测标签
- 设定问题类型
- 指定评估指标
- 盘它! fit()

AutoGluon 是何方神圣?

AutoML for Image, Text, Time Series, and Tabular Data

官方网站: <https://auto.gluon.ai/stable/index.html>

Github地址: <https://github.com/autogluon/autogluon>

10行代码战胜90%数据科学家?

20.8万 56 2021-05-09 13:30:16 未经作者授权, 禁止转载



3人正在看, 已装填 56 条弹幕 发个友善的弹幕见证当下
3570 2152 4094 1196

<https://www.bilibili.com/video/BV1rh411m7Hb>

动手学深度学习 bilibili

<https://www.bilibili.com/video/BV1F84y1F7Ps>

AutoGluon背后的技术

6.6万 47 2021-05-17 07:53:42 未经作者授权, 禁止转载



1人正在看, 已装填 47 条弹幕 发个友善的弹幕见证当下
1103 612 619 213
稿件投诉 1篇笔记

跟李沐学AI bilibili

自动调参
名过其实?

autogluon背后的技术	34.16%
autogluon高级使用	21.23%
使用autogluon来解决其他应用	24.35%
加州房价数据是如何来的	20.26%

Datawhale /02

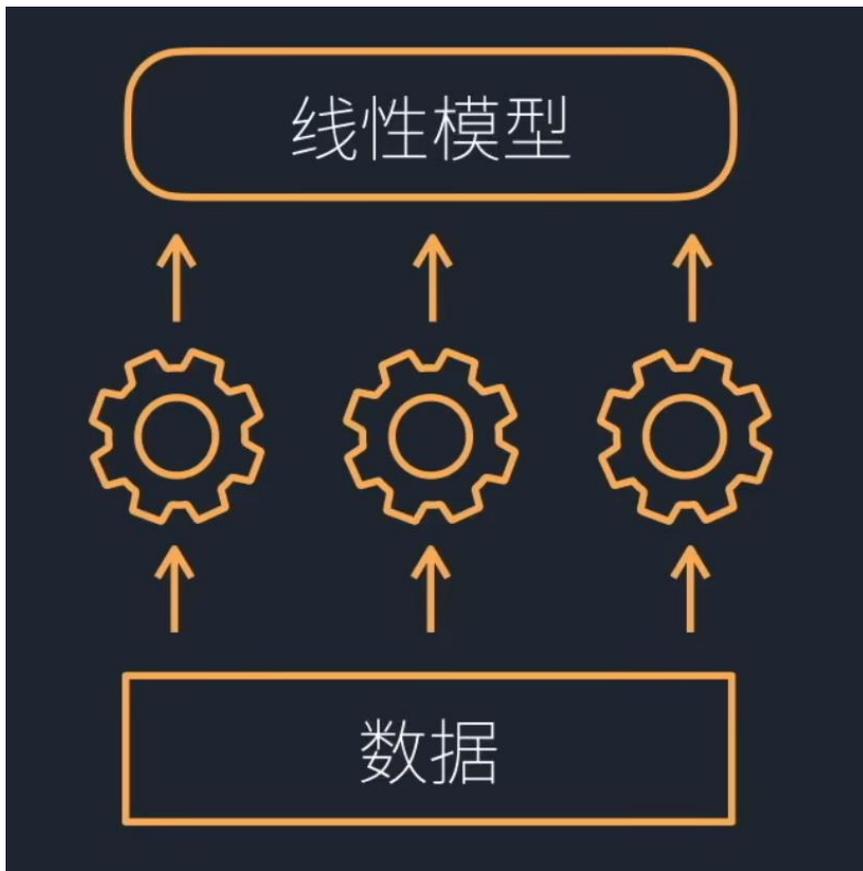
AutoML 技术杂谈

大部分automl框架是基于超参数搜索技术

Autogluon则依赖于融合多个无需超参数搜索的模型



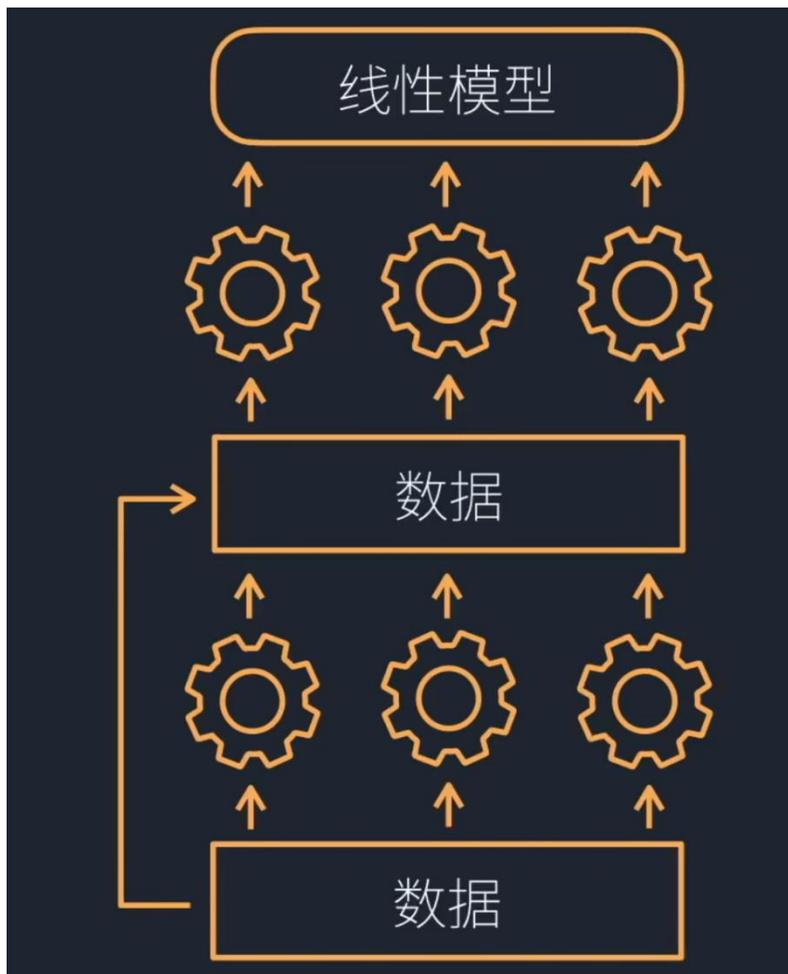
- **Stacking**
- **K-折交叉 Bagging**
- **多层 Stacking**



- 在**同一个数据集**上训练**多个**不一样的模型
- 这些模型的输出然后进入到一个**线性的模型**
- 得到**最终的输出**



- 不同的初始权重与数据训练多个同类模型
- K-折交叉验证的相同做法
- 整合多个同类模型的结果
 - 回归：取平均值
 - 分类：投票，少数服从多数

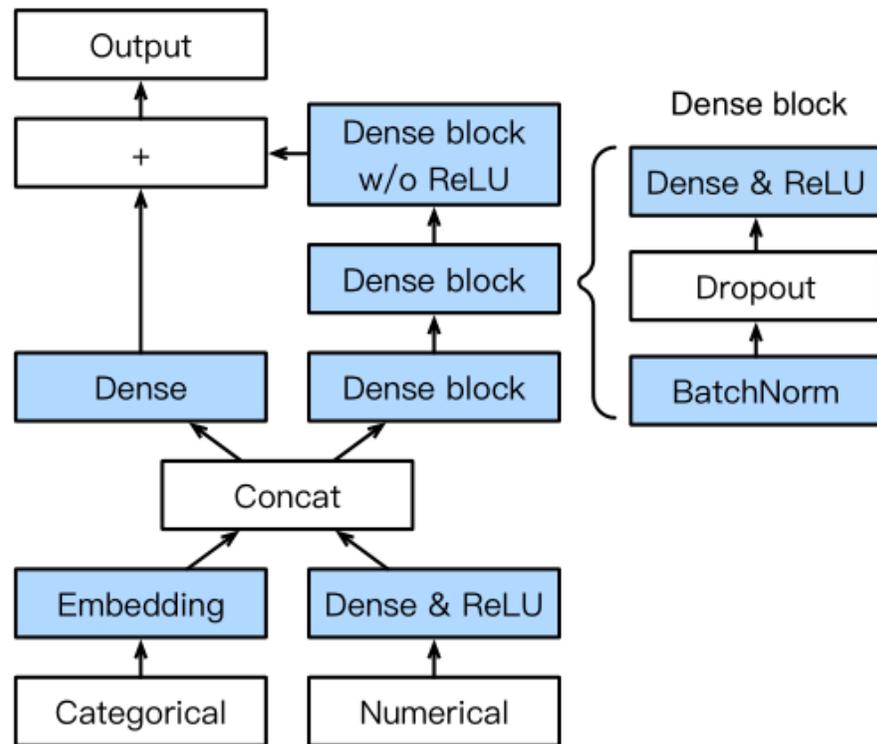


- 将多个模型的输出和数据合并起来再做一次 Stacking
- 防止过拟合，多层 Stacking 往往配合 K-则交叉 Bagging 使用
- 计算开销剧增！不环保 :(

Algorithm 1 AutoGluon-Tabular Training Strategy
(multi-layer stack ensembling + n -repeated k -fold bagging).

Require: data (X, Y) , family of models \mathcal{M} , # of layers L

- 1: Preprocess data to extract features
- 2: **for** $l = 1$ **to** L **do** {Stacking}
- 3: **for** $i = 1$ **to** n **do** { n -repeated}
- 4: Randomly split data into k chunks $\{X^j, Y^j\}_{j=1}^k$
- 5: **for** $j = 1$ **to** k **do** { k -fold bagging}
- 6: **for each** model type m in \mathcal{M} **do**
- 7: Train a type- m model on X^{-j}, Y^{-j}
- 8: Make predictions $\hat{Y}_{m,i}^j$ on OOF data X^j
- 9: **end for**
- 10: **end for**
- 11: **end for**
- 12: Average OOF predictions $\hat{Y}_m = \{\frac{1}{n} \sum_i \hat{Y}_{m,i}^j\}_{j=1}^k$
- 13: $X \leftarrow \text{concatenate}(X, \{\hat{Y}_m\}_{m \in \mathcal{M}})$
- 14: **end for**

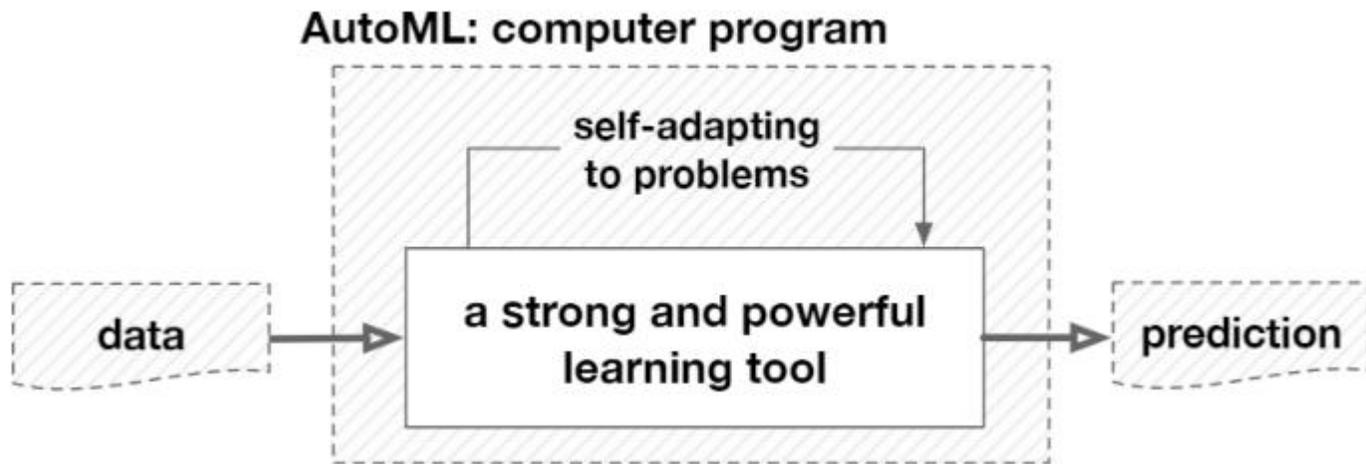


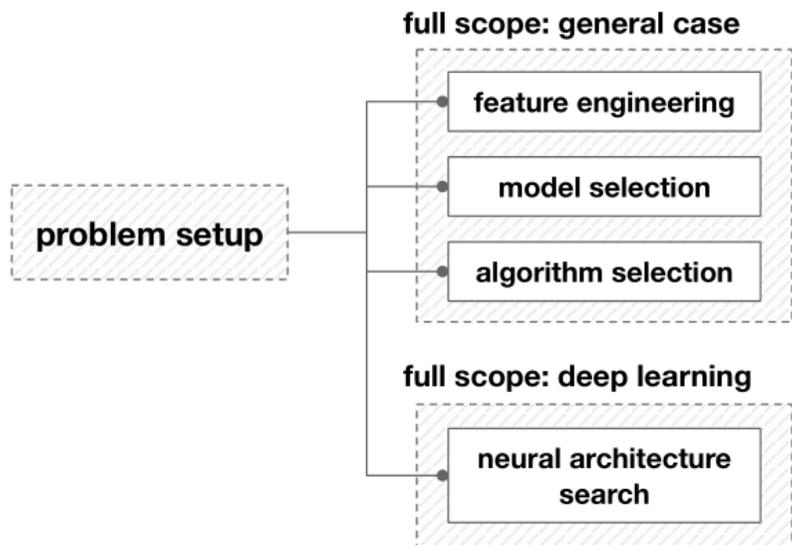
- 机器学习的大量工作仍然需要人工干预：特征提取、模型选择、参数调节
- 希望自动化进行机器学习中的特征、模型、优化、评价，甚至部署等环节

$$\begin{aligned} & \max_{\text{configurations}} \text{ performance of learning tools,} \\ & \text{s.t.} \begin{cases} \text{limited (or no) human assistance} \\ \text{limited computational budget} \end{cases} \end{aligned}$$

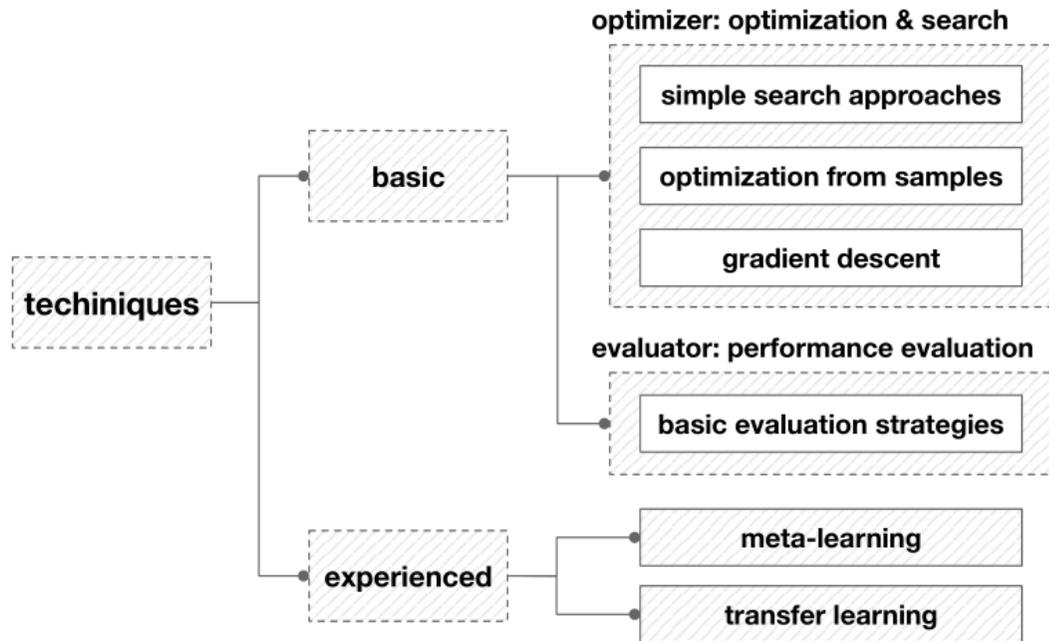
三个核心:

- 更好的训练效果
- 更少的人工参与
- 更低的计算资源





(a) "What to automate": by problem setup.



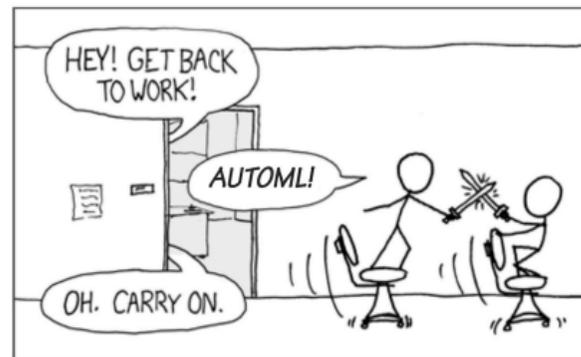
(b) "How to automate": by techniques.

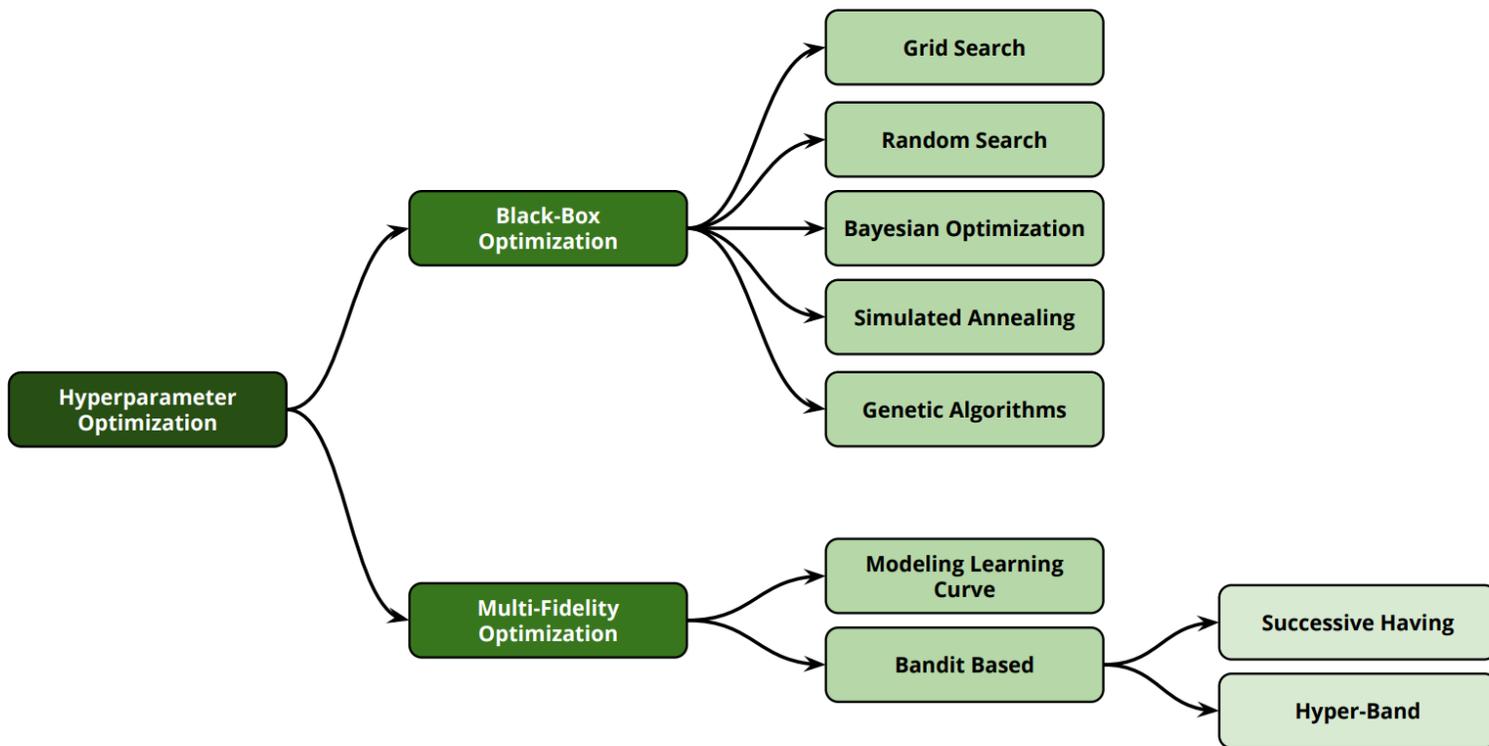
Automated Machine Learning (AutoML)

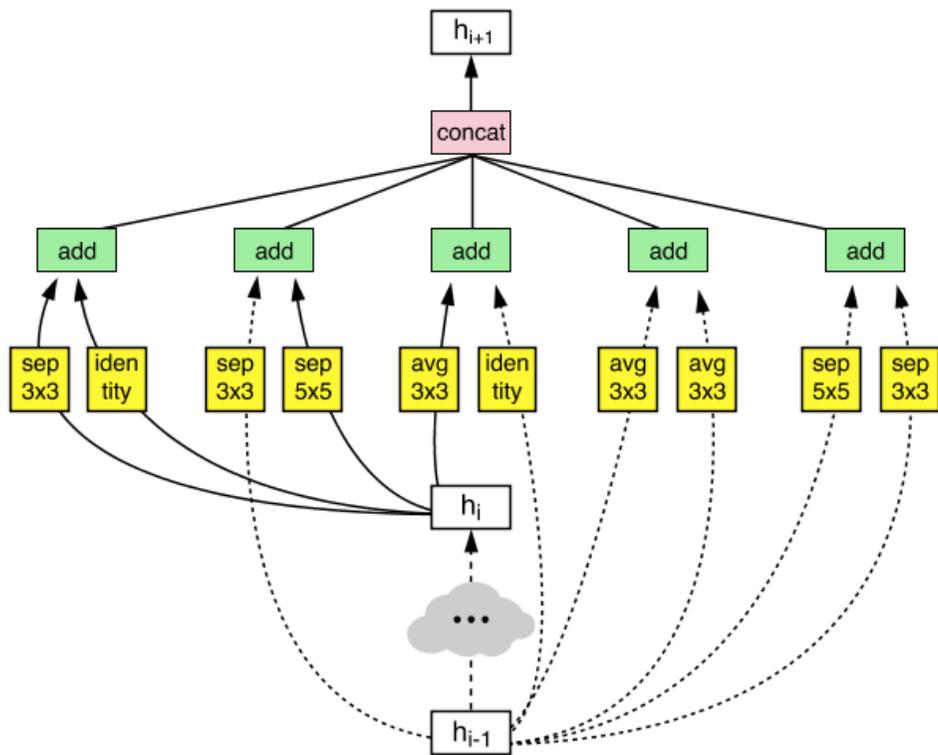


- Automate every step in applying ML to solve real-world problems: data cleaning, feature extraction, model selection...
- **Hyperparameter optimization (HPO):** find a good set of hyperparameters through search algorithms
- **Neural architecture search (NAS):** construct a good neural network model

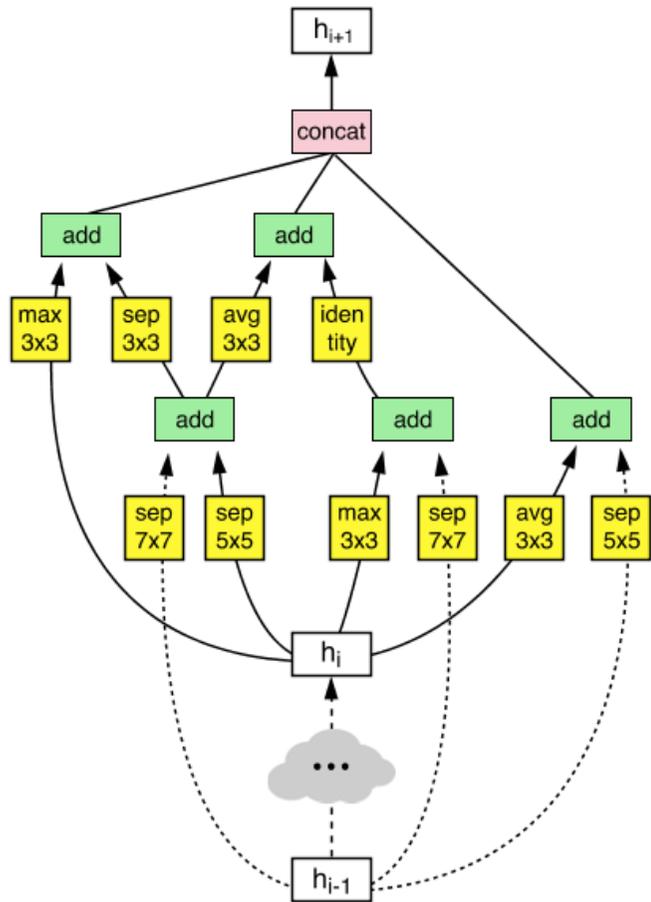
*THE DATA SCIENTIST'S #1 EXCUSE FOR LEGITIMATELY SLACKING OFF:
"THE AUTOML TOOL IS OPTIMIZING MY MODELS!"*







Normal Cell



Reduction Cell

结果对比

跟李沐学AI bilibili

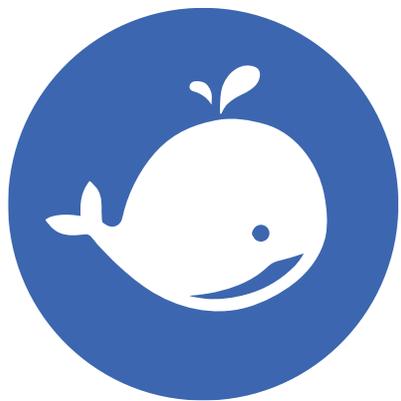
50个数据集

限制4小时运行时间

	冠军
AutoGluon	30
TPOT	5
GCP	7
Auto-sklearn	4
H2O	2
Auto-WEKA	1



for the learner,
和学习者一起成长



Thank
you



很高兴跟大家一起在夏令营中度过这个愉快的暑假，我们一起继续向人工智能训练大师之路进发！