

# Research Proposal

王梓奕 李一玄 刘俊豪 姚嘉浩 黄肖炜

- Research Proposal
  - 背景
  - 意义
  - 探索性分析
    - 结论
  - 思路
    - 神经网络
    - 抽样
      - 方案 1
        - 结果
      - 方案 2
        - 结果
      - 方案 3
        - 结果
  - TODOS
  - 未来

## 背景

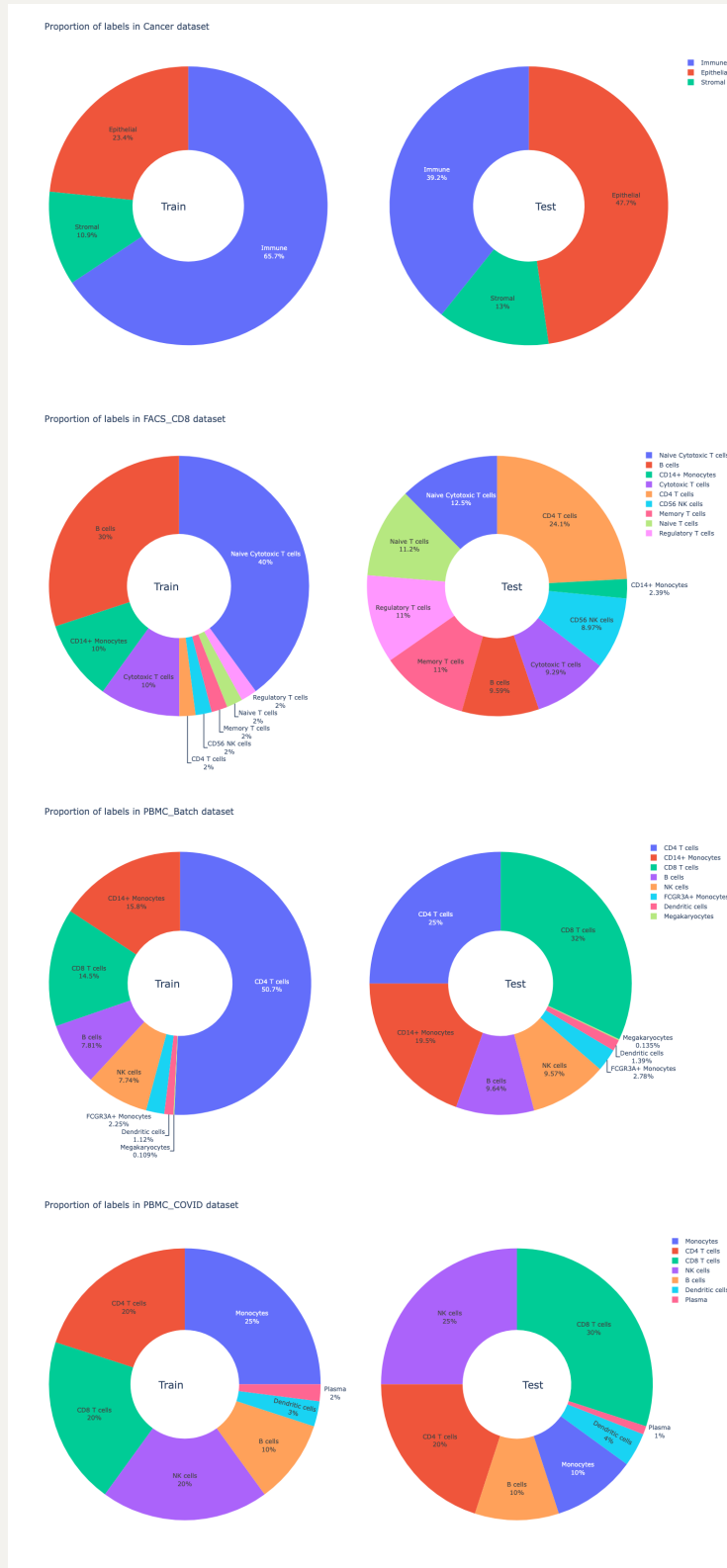
- 协变量偏移  
在假设标记函数不变的前提下，特征分布的变化导致的训练集和测试集存在本质上区别的分布偏移，称为协变量偏移。
- 标签偏移  
标签在训练集和测试集中出现的频率不同，例如，在训练集中，很少出现细胞类别为 q 的样本，但在测试集中，同时存在细胞类别为 p 和细胞类别 q 的样本，其中不变的是不同基因的表现。
- 概念偏移  
除上述两种偏移之外，还存在标签本身概念发生变化的偏移，例如，对三大类癌细胞的定义发生变化。

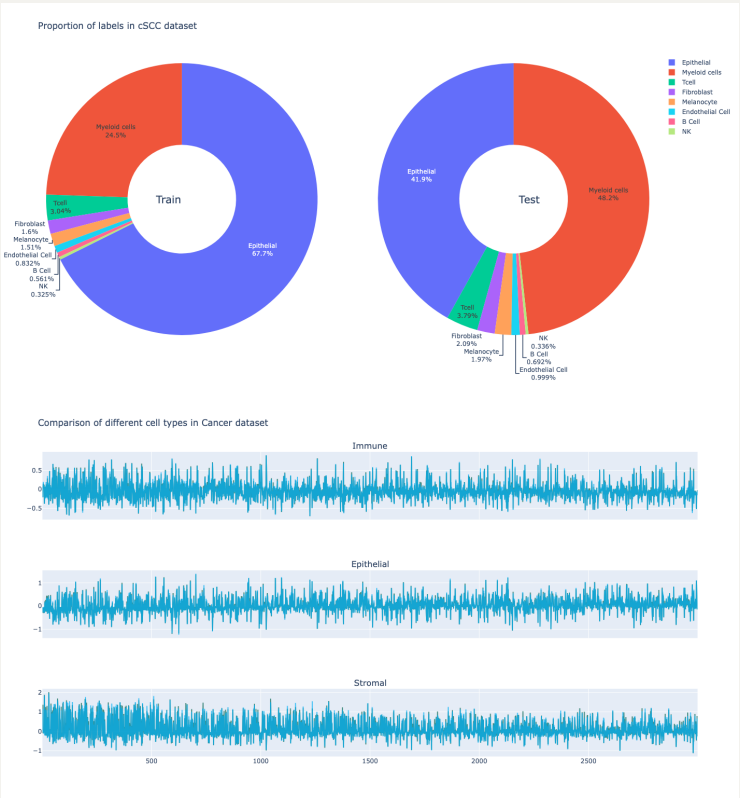
## 意义

- 提高模型在实际应用中的泛化能力和性能。
- 提高模型在不同标签分布下的预测准确性。
- 通过引入假设检验和超参数等方法，可以进一步解决不同类型的偏移问题，并提高模型的性能，改进预测方法。

## 探索性分析

- 统计五组训练集和测试集中各类别细胞的频率，绘制饼图，可看出各数据集均涉及协变量偏移





- 统计五组训练集和测试集中各类别细胞各基因的均值，绘制折线图，深色代表训练集与测试集的均值相近，浅蓝色与红色则代表该基因下均值存在差异，结果发现在五组数据集中差异均不显著，给出 **Cancer** 数据集的折线图作为参考

结论

五组数据集主要问题在于较为严重的协变量偏移，因此针对该现象训练模型并拟合

思路

How to resist distribution shift

神经网络

- 调用 **PyTorch** 的 API 建立模型

Architecture of Neural Network

Input and Embedding layers

`Linear ( in_features=n_features, out_features=n_features, bias=True)`

`Linear ( in_features=n_features, out_features=300, bias=True)`

MLP

`Linear ( in_features=300, out_features= inter_features, bias=True)`

`BatchNorm1d ( batch_size= inter_features, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)`  
`ELU ( alpha=1.0)`  
`Dropout ( p=0.5, inplace=False)`

```

Linear (in_features=_inter_features,
         out_features=_inter_features, bias=True)

BatchNorm1d (batch_size=_inter_features, eps=1e-05,
              momentum=0.1, affine=True, track_running_stats=True)
ELU (alpha=1.0)
Dropout (p=0.5, inplace=False)

Linear (in_features=_inter_features, out_features=n_labels,
         bias=True)

Softmax (dim=1)

```

## Optimizer

```

Adam (lr=LR, betas=(BETA1, BETA2), eps=EPS)

```

## LossFuntion

```

CrossEntropyLoss

```

## Learning-rate Scheduler

```

ReduceLROnPlateau ("min", patience=PATIENCE,
                    threshold=THRESHOLD)

```

- $\text{\_inter\_features} = \frac{2}{3} (\text{\_in\_features} + \text{\_out\_features})$
- `nEpoch = 10`

## 抽样

```

For BGD @property@abstractmethod

```

## 方案 1

使用一种概率分布 (默认为标准高斯分布) 为训练集中的所有样本生成一个概率 (权重), 基于这些权重抽取样本

```

1  def sample(self, distribution: scipy.stats.rv_continuous =
    scipy.stats.uniform, *args
2      ) -> tuple[torch.Tensor, torch.Tensor]:
3      """Use the generated weights to sample the data
4
5      Parameters
6      -----
7      distribution : scipy.stats.rv_continuous, optional
8          The distribution to generate each probability from,
9      by default scipy.stats.uniform
10
11     Returns
12     -----
13     tuple[torch.Tensor, torch.Tensor]
14         X_Batch, y_Batch
15     """
16     self.weights = prosGenerator(distribution=distribution,
    size=self.size, *args)
    self.choices = np.random.choice(

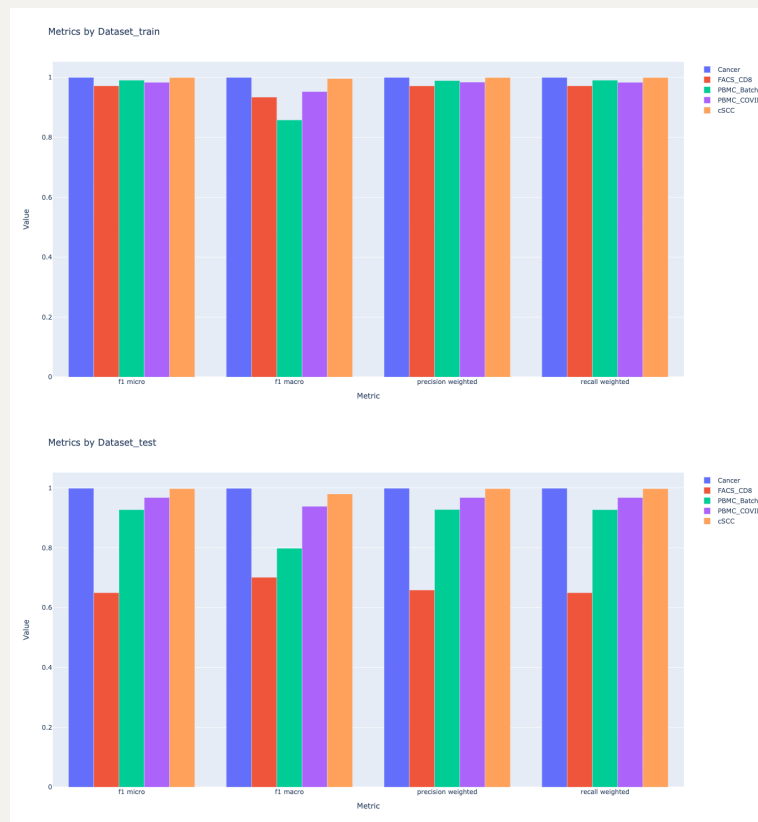
```

```

17         range(self.size), self.batch_size, False,
        self.weights
18     )
19     return featureLabelSplit(self.data[self.choices])

```

## 结果



- 在 FACS\_CD8 与 PBMC\_Batch 上表现很差

## 方案 2

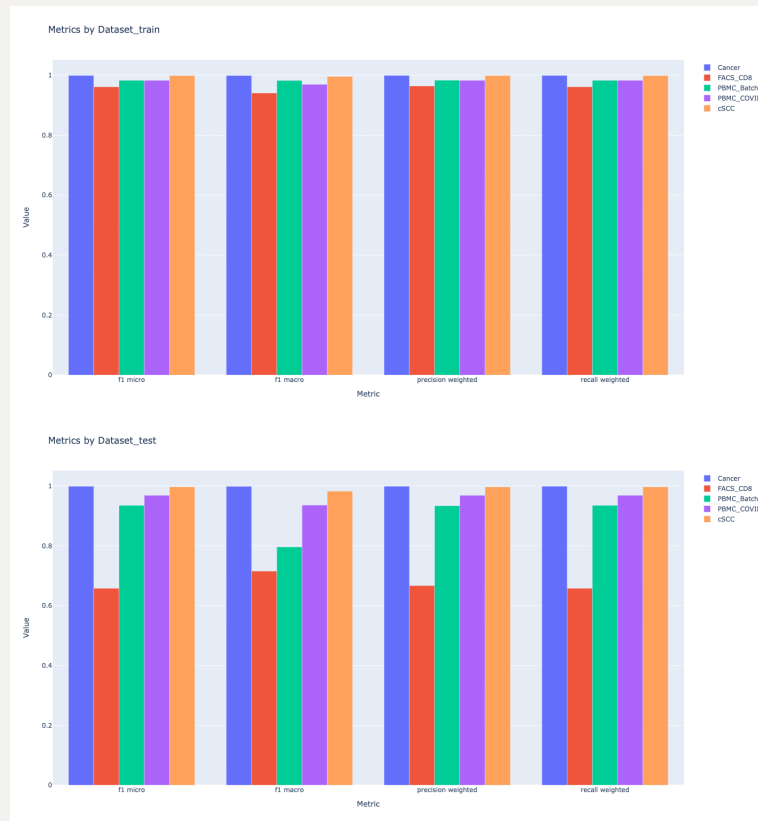
将要抽取的样本数随机分为 n 组，其中 n 为数据集中不同标签个数，每组中的样本数对应于各个标签中采用 Bootstrap 采样的样本数

```

1     def sample(self) -> tuple[torch.Tensor, torch.Tensor]:
2         """Utilize the generated counts to sample the data by
        Bootstrap
3
4         Returns
5         -----
6         tuple[torch.Tensor, torch.Tensor]
7             X_Batch, y_Batch
8         """
9         nums = self.getNum # A property returns the nums to be
        sampled for each label
10        for i in range(len(self.changeIndexes) - 1):
11            self.choices += list(
12                np.random.choice(
13                    range(self.changeIndexes[i],
        self.changeIndexes[i + 1]),
14                    nums[i],
15                    True,
16                )
17        )

```

## 结果



- 相比方案 1 有提升但在 FACS\_CD8 与 PBMC\_Batch 上表现依然很差

## 方案 3

利用训练集中各标签下数据的均值及标准差伪造样本并植入到训练集中

```
1 def sample(self) -> NotImplemented
2     return NotImplemented
```

## 结果

## TODOS

1. 实现方案 3
2. 改善研究背景与研究意义
3. 编写函数删除所有 .html 文件

## 未来

1. 混合各个抽样方案
2. 对表现较好的数据集部署更简单的模型以节省成本