

选题范围

选题范围：同学们自行选题，在机器学习和人工智能范畴内的主题均可。可以选择扩展一个课程案例，或者探索新案例。

需要考虑：自身时间安排、题目算力要求、背景知识要求等等，选择一个自己能够完成的题目。相比问题的难度高低，我们更重视问题的完成度，报告的完整性、科学性。

报告内容

1 介绍选题背景、问题定义：

选题背景：介绍要解决一个什么问题，问题的背景是什么，问题的意义是什么等等。

正式定义：给出问题的正式定义，明确解决该问题涉及到解决哪类机器学习问题：分类、回归、检测、分割、生成等等。阐明问题涉及的输入、输出和相关变量，其类型与范围或取值等信息。

2. 介绍所用数据集：

- (1) 数据来源：说明所使用的数据集来自哪个公开数据集或自行采集。如果使用公开数据集，应给出数据集名称及数据量相关信息，如图像数量、特征维度等。如果自行采集，应简述采集方法和过程。
- (2) 数据预处理：说明对数据进行的预处理步骤，如噪音过滤、缺失值填补、特征缩放、数据增强等。这些步骤应该是可重复的。如果没有进行预处理，也应说明原因。
- (3) 数据划分：说明如何划分训练/验证/测试集。划分比例和数量应该具体给出。如果使用的公开数据集自带划分，也应说明。
- (4) 数据分析：进行一定的数据分析,如变量类型、特征值范围、类别分布等，必要时给出用以分析的数据可视化结果。这可以让读者对数据有一个大致了解，也有助于理解后续的方法和结果。

3. 明确衡量问题的指标：

明确解决问题需要关注哪些指标（metric），例如：

分类问题：准确率、精度、召回、F1 score 等。
回归问题：L1/L2 误差、均方根误差(RMSE)等。
图像检测：mAP、准确率、精度、召回、F1 score、mIoU 等。
图像分割：Pixel Accuracy、mIoU、F1 score 等。
图像生成：FID、Inception score 等。
机器翻译：BLEU score 等。
语音识别：词错率(WER)、字符错率(CER)等。
其他：根据具体问题而定。

参考：<https://scikit-learn.org/stable/modules/classes.html#module-sklearn.metrics>

4. 设计具体方案：

基于前面的问题定义、数据、评价指标等，设计解决问题的具体方案，包括但不限于：

- (1) 选择哪些算法或机器学习模型，解释为什么这些算法/模型适合解决前面的问题。如果涉及多个算法/模型，说明如何串接组织这些算法/模型。
- (2) 网络/模型结构：如果用神经网络模型，考虑通过或文字、或代码、或结构图，来解释神经网络的架构。
- (3) 训练参数：详细给出模型训练所用的优化算法、损失函数、学习率、epoch 轮数、batch size 等相关超参数设置及取值理由。

5. 代码说明：

要求见《代码要求》

6. 分析实验结果：

对实验结果进行详细分析，有多个角度可以考虑，比如：

- (1) 指标：分析选择使用哪些指标评价结果，这些评价指标的具体评价结果如何。
- (2) 参数：分析超参数设置对结果的影响，尝试找到最优的参数区间或值。
- (3) 数据：如果有不同的数据划分方式，分析其对结果的影响。如果有更多的数据,会如何提高性能。
- (4) 模型：分析不同模型结构或不同模型种类对结果的影响。
- (5) 算法：分析如果使用不同的机器学习算法，会对最终结果产生什么影响。

7. 参考文献

对于报告所参考的学术论文、在线技术资料等，应该在参考文献部分给出，并在正文中引用。