

2020

Machine Learning

Fourth Assignment

姓名	班级	学号
王成航	统计 71	2176122248
张申铎	统计 71	2176112379
王泽昊	统计 71	2176112782

ŀì?l 集ŀř?ŋ??ă

目录

1 引言

集成学习(Ensemble Learning) 用多种不同的模型来预测目标变量, 然后将这些模型的预测结果进行组合, 以提高预测的准确性。集成学习在机器学习中有着广泛的应用, 如分类、回归、排序等。集成学习的主要思想是“三个臭皮匠, 顶个诸葛亮”。集成学习可以分为 Bagging、Boosting 和 Stacking 三种主要方法。Bagging 通过并行训练多个模型, 然后对它们的预测结果进行平均或投票。Boosting 通过串行训练多个模型, 每个模型都试图纠正前一个模型的错误。Stacking 通过训练多个模型, 然后将它们的预测结果作为输入, 训练一个元模型来做出最终预测。

1979年, Dasarthy 提出了集成学习(Ensemble system) 的概念, 并用它来预测股票价格。1988年, Kearns 和 Valiant 提出了弱学习(Weak Learning) 的概念, 并证明了强学习(Strong Learning) 可以通过组合弱学习来实现。1990年, Schapire 提出了 boosting 的概念, 并证明了它可以用于分类和回归。1996年, Breiman 提出了 Bagging 的概念, 并证明了它可以用于分类和回归。Bagging 和 Boosting 是集成学习的两种主要方法。Bagging 通过并行训练多个模型, 然后对它们的预测结果进行平均或投票。Boosting 通过串行训练多个模型, 每个模型都试图纠正前一个模型的错误。

集成学习(Ensemble Learning) 是一种机器学习方法, 它通过组合多个模型的预测结果来提高预测的准确性。集成学习可以分为 Bagging、Boosting 和 Stacking 三种主要方法。Bagging 通过并行训练多个模型, 然后对它们的预测结果进行平均或投票。Boosting 通过串行训练多个模型, 每个模型都试图纠正前一个模型的错误。Stacking 通过训练多个模型, 然后将它们的预测结果作为输入, 训练一个元模型来做出最终预测。

2 集成学习(Ensemble Learning)

集成学习(Ensemble Learning) 是一种机器学习方法, 它通过组合多个模型的预测结果来提高预测的准确性。集成学习可以分为 Bagging、Boosting 和 Stacking 三种主要方法。Bagging 通过并行训练多个模型, 然后对它们的预测结果进行平均或投票。Boosting 通过串行训练多个模型, 每个模型都试图纠正前一个模型的错误。Stacking 通过训练多个模型, 然后将它们的预测结果作为输入, 训练一个元模型来做出最终预测。

2.1 集成学习(Ensemble Learning) 的概述

集成学习(Ensemble Learning) 是一种机器学习方法, 它通过组合多个模型的预测结果来提高预测的准确性。

Bagging 概述

Bagging 是一种集成学习方法, 它通过并行训练多个模型, 然后对它们的预测结果进行平均或投票。Bagging 的主要思想是“三个臭皮匠, 顶个诸葛亮”。Bagging 可以分为 Bagging for Classification 和 Bagging for Regression 两种主要方法。Bagging for Classification 通过并行训练多个分类模型, 然后对它们的预测结果进行投票。Bagging for Regression 通过并行训练多个回归模型, 然后对它们的预测结果进行平均。

wwz

2020

Machine Learning

Fourth Assignment

姓名	班级	学号
王成航	统计 71	2176122248
张申铎	统计 71	2176112379
王泽昊	统计 71	2176112782

图 1: 三种?U??ZÍ?η??ă?Ů??şT.

Boosting?Ů?şT

Boosting?Ů?şT?Ÿ? 一种能?đş?řE???ŋ?ă?ŽÍ??ŃŮ????ŋ?ă?ŽÍ, ?Ŏ 而?ŘŘ?ŃĜ?ŘĎ 种?ŋ?ă?Ů?şT 的?Ů?şT. 理??上, boosting ?Ř????Ÿ? 著?ĜŘ?řR???ŋ?ă?ŽÍ 的?AŘ?·?, ?Ž?Ž????ŋ?ă?ŽÍ 的?TŁ?đIJ?R??Ÿ? 稍???Ÿ?Ŏ 随?IJ?猜?tŃ, ??Ť??C?řR?Eş 策?ăS——?Ťř?Ń??Lă?İČ?Í??đŃ. ??ŘŎBoosting?IJÍ?Ř?Ń?Ů??řE?Žř?đŽ 的?İČ 重赋?Ă??Ž?Ů??IJş?ŋ?Č 中错 ???IJĂ?đŽ 的?Ťř?Ń? 集, 通 ?Ĝ?Ş?ŘL?Lă?İČ?đŽ?Ťř?ŁŤ?Í?L?LE???L?LŮ?Lă?İČ??C 生?IJĂ?L?Ď?tŃ. Boosting 的??R?Í??đŃ 都?Ÿ??ŃŤ??Ř?Ń, ??R?Í??đŃ?Eş??Ž 下一 ??Í??đŃ?A?Ěş?şÍ 的 ???A, ?IJĂ?ŘŎ?IJÍ 不?AŘ?ŘS?? 何?Í??đŃ 的?LŃ?ŘŘ 下聚?ŘL?Ş?Ĝ??Ş?đIJ.

Stacking?Ů?şT

Stacking?Ů?şT?Ÿ? 一种集?LŘ?ŋ?ă?LĂ?IJ?, 用 ?Ŏ?IJĂ?řR?ŃŮ — ??LŮ?đŽ??şŽ?ŃŮ?ŽÍ 的?şŽ?ŃŮ??? · ? 率?Ů?şT. ??Č 通 ?Ĝ?ŎÍ????şŽ?ŃŮ?ŽÍ 相????Ŏ?LĂ?ŘŘ?Ž 的?ŋ?ă 集的?AŘ?· ??İ??RS?Ń??Ě? 作用. ?Ž??ŎÍ??? 的 ?Ĝ 程?ŃĚ?Ń??Ž?IJÍ??Ń??C 中?řE?? —??C 的?Ŏş?ĝŃ?şŽ?ŃŮ?ŽÍ??? 部?LE?ŋ?ă 集的猜?tŃ?Ž?Ń?şŽ?ŃŮ, ???RL?řİ?Ť????ŋ?ă 集的?L? 余部?LE?Ž?Ń 猜?tŃ, ??? 且 ?Ş?Ĝ??ŋč?? 的 ?Ş?đIJ. ?;Ş 与?đŽ??şŽ?ŃŮ?ŽÍ 一起 ? 用?Ů?, ?ăE?Řă?şŽ?ŃŮ?R???? 看作?Ÿ? — ??đ?RL?Ń?A 的?đŃ?İČ 版?IJ?, ?L? 用??Ť?đ?RL?Ń?A?Žř?đŃ?İČ 的策 ??İ??Ď?ŘL?ŘĎ??şŽ?ŃŮ?ŽÍ. ?;Ş 与?ŃŤ??şŽ?ŃŮ?ŽÍ 一起 ? 用?Ů?, ?ăE?Řă?şŽ?ŃŮ?Ÿ? 一种用 ?Ŏ?ř????L??ŘŎ?ă?ŋč?L?şŽ?ŃŮ?ŽÍ 的错 ?? 的?Ů?şT, ???şŽ?ŃŮ?ŽÍ · ş?R?IJÍ??Ž?ŋ?ă 集上 ?Ž?Ń?E?ŋ?Č???????E???Ž??Ÿ.

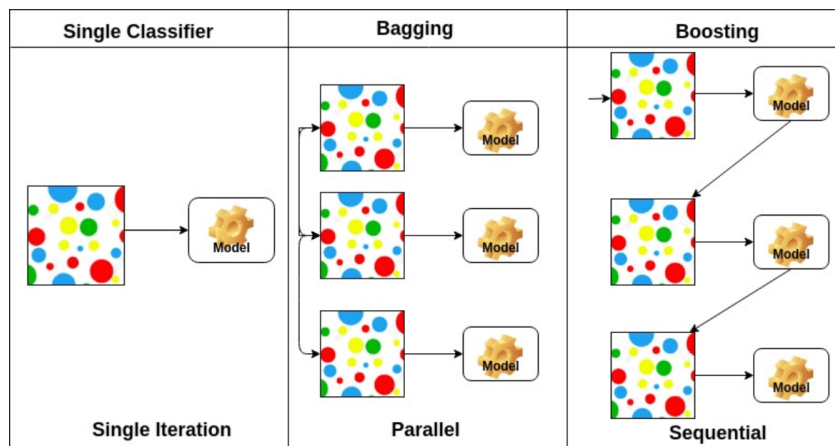


图 2: şş?ŎStacking?Ů?şT 的?LE??.

???ş?ŋ?ă?ŽÍ 的?LE?Ÿč??ş??Ă, stacking?Ů?şT?Ř????LE? 两 ?Ď?Ž?IJÍ??ş?Ń 集?LŘ?Ů??şT 中, ?ş?ŋ?ă?ŽÍ??ş?Ń 的 ?ğ 生, ?Ń??C 随?IJ??č??đŮ. ?IJÍ????R 集?LŘ?Ů??şT 中, ?ş?ŋ?ă?ŽÍ?Ÿ????R 生?LŘ 的, ?Ń??CAdaBoost?Ž

???ş?ŋ?ă?ŽÍ 的 ???đŃ?ş?R????LE? 两 ?Ď?Ž?ŘŃ?đĎ 集?LŘ?Ů??şT?IJÍ??R????ŋ?č 中 ? 用相?ŘŃ????đŃ 的?ş?ŋ?ă?ŽÍ. ??C?đĎ 集?LŘ?Ů??şT?IJÍ??R????ŋ?č 中 ? 用不?ŘŃ????đŃ 的?ş?ŋ?ă?ŽÍ.

3 弱分类器集成

AdaBoost 是一种弱分类器集成算法。AdaBoost 通过不断训练弱分类器，并将它们的预测结果加权组合，最终得到一个强分类器。AdaBoost 背后的思想是：如果一个弱分类器在训练集上的表现不佳，那么通过调整其权重，并重新训练，可以使其表现更好。AdaBoost 通过迭代的方式，不断调整弱分类器的权重，直到达到预定的精度要求。

1. 弱分类器在训练集上的表现不佳，通过调整权重，使其表现更好。
2. 在弱分类器集成中，每个弱分类器的权重都根据其表现进行调整。表现越好的弱分类器，其权重越大；表现越差的弱分类器，其权重越小。

3.1 AdaBoost 算法流程

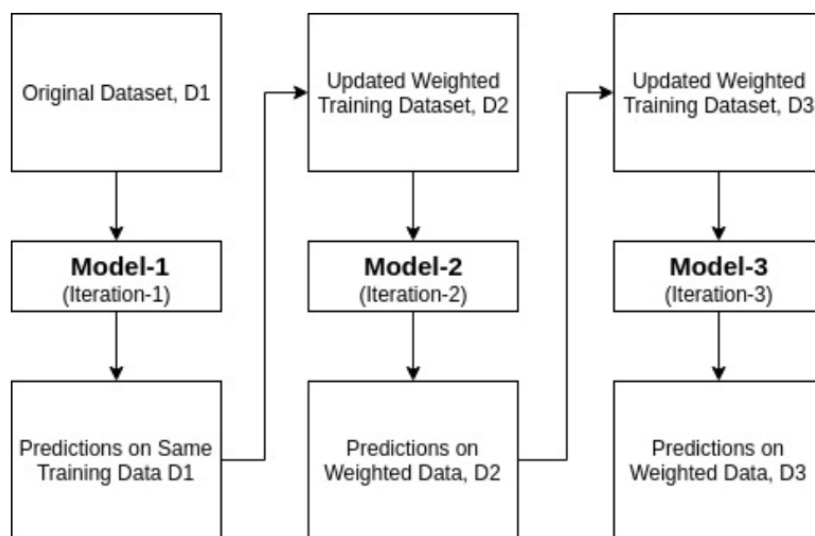


图 3: AdaBoost 算法流程。

1. AdaBoost 算法通过不断训练弱分类器，并将它们的预测结果加权组合，最终得到一个强分类器。
2. 在弱分类器集成中，每个弱分类器的权重都根据其表现进行调整。表现越好的弱分类器，其权重越大；表现越差的弱分类器，其权重越小。
3. 弱分类器在训练集上的表现不佳，通过调整权重，使其表现更好。
4. 在弱分类器集成中，每个弱分类器的权重都根据其表现进行调整。表现越好的弱分类器，其权重越大；表现越差的弱分类器，其权重越小。
5. 弱分类器在训练集上的表现不佳，通过调整权重，使其表现更好。
6. 在弱分类器集成中，每个弱分类器的权重都根据其表现进行调整。表现越好的弱分类器，其权重越大；表现越差的弱分类器，其权重越小。

3.2 数据集

我们使用 E sklearn 中 ensemble 里的 AdaBoostClassifier。我们数据集选用 EUCI 集中的 IRIS 数据集。IRIS 数据集 — 一个包含 150 个样本的数据集，其中包含 3 种 Iris (Iris Setosa, Iris Versicolour, Iris Virginica) 的 4 个特征 (sepal length, sepal width, petal length, petal width)。我们使用这 4 个特征来区分这 3 种 Iris。我们使用 sklearn 中的 load_iris 函数来加载这个数据集。我们使用 sklearn 中的 train_test_split 函数来将数据集分为训练集和测试集。我们使用 sklearn 中的 AdaBoostClassifier 来训练模型。我们使用 sklearn 中的 score 函数来评估模型的性能。

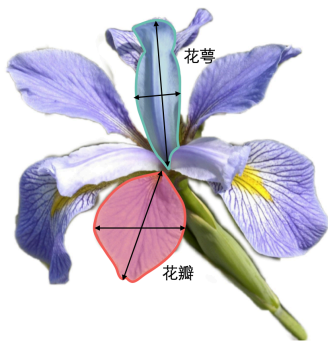


图 4: Iris 数据集的样本与特征

4 模型

4.1 数据集

4.1.1 数据集

我们使用 sklearn 中的 load_iris 函数来加载这个数据集。我们使用 sklearn 中的 train_test_split 函数来将数据集分为训练集和测试集。我们使用 sklearn 中的 AdaBoostClassifier 来训练模型。我们使用 sklearn 中的 score 函数来评估模型的性能。

我们使用 sklearn 中的 load_iris 函数来加载这个数据集。我们使用 sklearn 中的 train_test_split 函数来将数据集分为训练集和测试集。我们使用 sklearn 中的 AdaBoostClassifier 来训练模型。我们使用 sklearn 中的 score 函数来评估模型的性能。

我们使用 sklearn 中的 load_iris 函数来加载这个数据集。我们使用 sklearn 中的 train_test_split 函数来将数据集分为训练集和测试集。我们使用 sklearn 中的 AdaBoostClassifier 来训练模型。我们使用 sklearn 中的 score 函数来评估模型的性能。

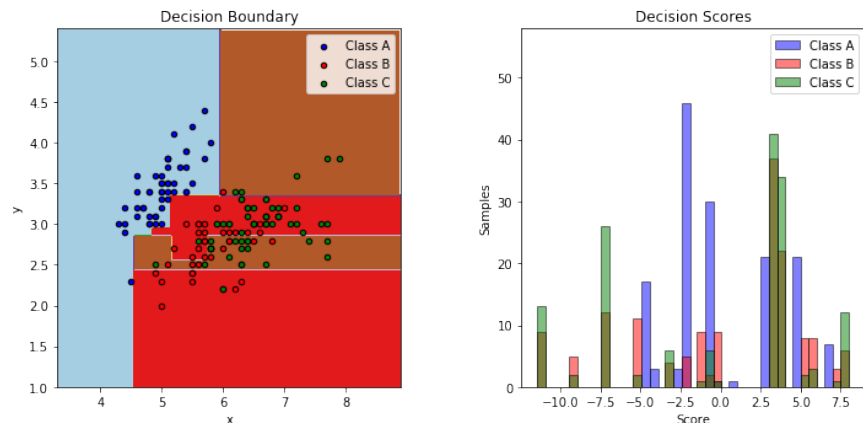


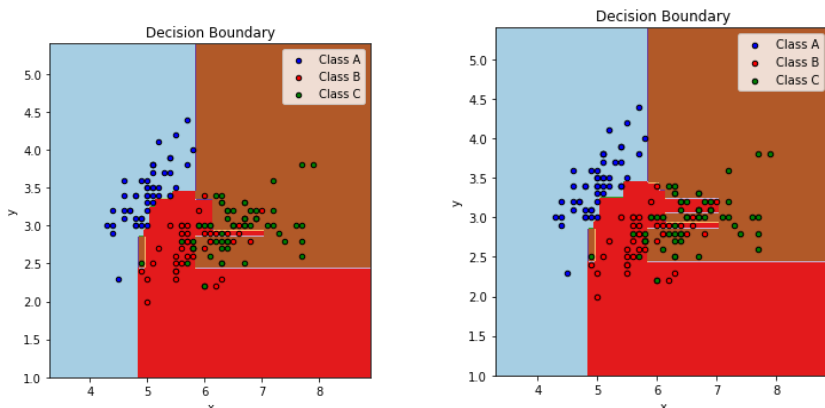
图 5: \mathbb{O} 采用 E_8 策 \mathfrak{a}_8 作 $\eta\mathfrak{a}_{IJ}$ 的 $\mathbb{L}_E\mathbb{Z}_I\mathbb{J}_I\mathfrak{a} \cdot \mathbb{J}$ 的两 \mathbb{A}_I 上 $\eta\mathbb{C}\mathbb{G}\mathfrak{r}_i$ 的 \mathbb{L}_E 界的 \mathbb{R} 视 $\mathbb{N}\mathbb{U}$, \mathbb{R} 看 $\mathbb{L}\mathfrak{r}\mathbb{E}\mathfrak{r}_i\mathbb{N}\mathbb{A}\mathbb{E}\mathbb{O}\mathfrak{s}\eta\mathfrak{a}_{IJ}$ 的 \mathbb{L}_E , $\mathbb{R}\mathfrak{r}_i\mathbb{Y}$ — $\mathbb{E}\mathfrak{g}\mathfrak{r}_i\eta\mathfrak{a}_{IJ}$.

\mathbb{R}^n 中 n 个线性无关的向量 $\alpha_1, \alpha_2, \dots, \alpha_n$ 的集合 $\{\alpha_1, \alpha_2, \dots, \alpha_n\}$ 称为 \mathbb{R}^n 的一个基。基的个数称为 \mathbb{R}^n 的维数，记为 $\dim \mathbb{R}^n$ 。显然 $\dim \mathbb{R}^n = n$ 。

4.2 适?;S 的?T??Z

4.2.1 量?RŘ?NĚ

一 一IAÄÄNT 的T??Z 的Äi??f??Y??Ö???d?Lä???η??ä?IJ? 的Tf 量, 通 ?G?Zf?dg 的η??ä 规í??Ö???G; Tf?Až?Žt??; 的Nš?ŘL. 那?L?LS???IJ?d?Lä?η??ä?IJ? 的规í? 的Ů??Äž???s???IJE?? 何的TŁ?dIJ 上的?ŘŘ?NĚ. 甚至 ?Y?IJE?LÄ 下降. 通 ?G?Ä?NŮ 版η??ä?IJ? 的?LE???? 界?R? 视?NŮ, ?LS?? 不 ?观?š?Lf?Ž?ä · 的??d?Lä???η??ä?IJ? 的Tf 量???s???IJE???LS?? 的?LE???? 界?g 生 ?? 何??d 质?Äg 的;??SŇ. 甚至都?s???IJE??d?Lä?LE???? 界的Zš?LŸ 程??, ?R??Y??? 的位 ??IJE?E 一??Ž 调Tf. ?LÄ???LS???s???IJ??R????d??ηd?Ů?í??dN 的限?L??s???IJ???N?Ěíí???Ř?Ö???η??ä?IJ? 的?LE????.



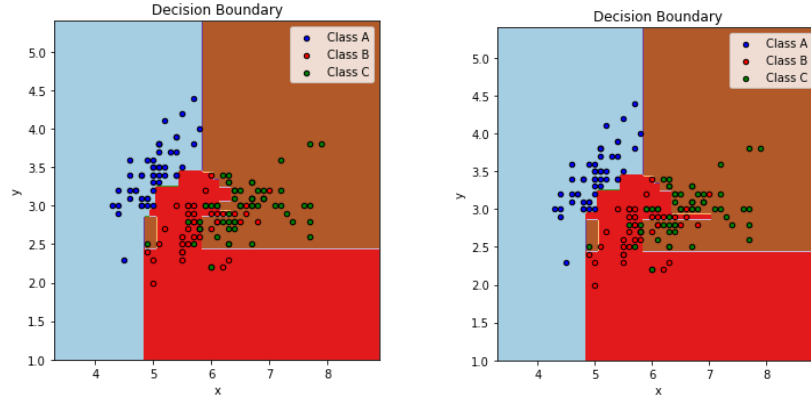


图 6: 图 6 展示了在 2D 空间 (x 轴范围 4 到 8, y 轴范围 1.0 到 5.0) 中, 使用 SVM 模型对三个类别 (Class A, Class B, Class C) 的数据进行决策边界划分的结果。图中显示了三个类别的数据点分布以及由 SVM 模型生成的非线性决策边界。决策边界将空间划分为三个区域, 分别对应 Class A (蓝色)、Class B (红色) 和 Class C (绿色)。

4.2.2 使用 AdaBoost 提升 SVM 模型性能

在 2D 空间上, 使用 AdaBoost 提升 SVM 模型的性能。AdaBoost 是一种集成学习方法, 通过多次迭代, 每次迭代使用一个弱分类器 (如 SVM) 对数据进行分类, 然后将多个弱分类器的结果进行加权组合, 得到最终的强分类器。在 AdaBoost 提升后的 SVM 模型中, 决策边界变得更加复杂, 能够更好地拟合训练数据。图 7 展示了 AdaBoost 提升后的 SVM 模型对三个类别的数据进行分类的结果。图 7 左侧的散点图显示了决策边界, 右侧的柱状图显示了决策得分 (Decision Scores)。

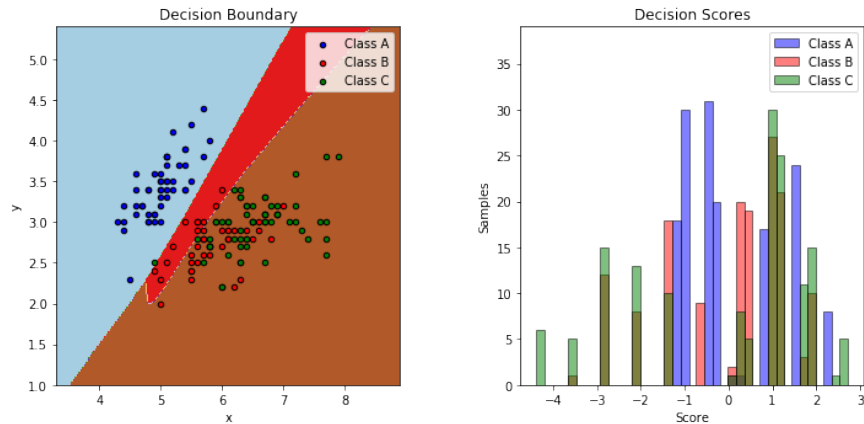


图 7: 图 7 展示了使用 AdaBoost 提升后的 SVM 模型对三个类别的数据进行分类的结果。左侧的散点图显示了决策边界, 右侧的柱状图显示了决策得分 (Decision Scores)。

在 2D 空间上, 使用 AdaBoost 提升 SVM 模型的性能。AdaBoost 是一种集成学习方法, 通过多次迭代, 每次迭代使用一个弱分类器 (如 SVM) 对数据进行分类, 然后将多个弱分类器的结果进行加权组合, 得到最终的强分类器。在 AdaBoost 提升后的 SVM 模型中, 决策边界变得更加复杂, 能够更好地拟合训练数据。图 7 展示了 AdaBoost 提升后的 SVM 模型对三个类别的数据进行分类的结果。图 7 左侧的散点图显示了决策边界, 右侧的柱状图显示了决策得分 (Decision Scores)。

并且从直观上看, 提升后的 SVM 模型在 2D 空间上的决策边界更加复杂, 能够更好地拟合训练数据。提升后的 SVM 模型在 2D 空间上的决策边界更加复杂, 能够更好地拟合训练数据。

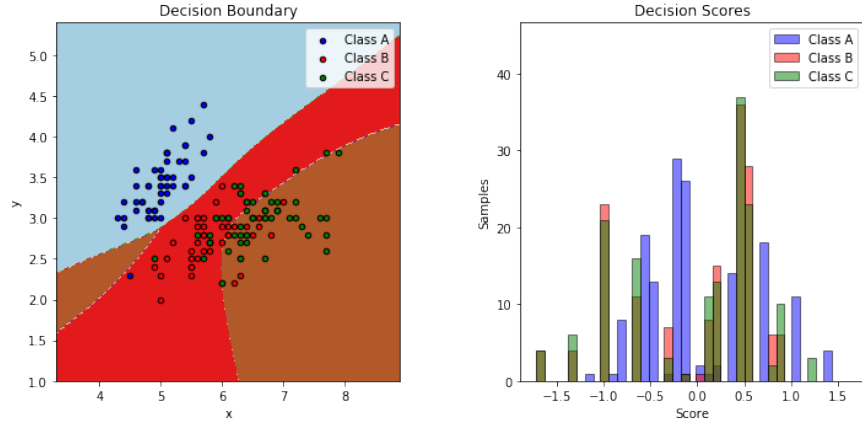


图 8: SVM 决策边界的可视化, 显示了 100 个样本在 2D 空间上的分布。决策边界将三个类别 (Class A, Class B, Class C) 分离开。右侧的直方图显示了每个类别的决策得分分布, 证实了 SVM 模型的有效性。

界。相對於決策邊界而言, SVM 作爲一個非線性的分類器, 其決策邊界通常比直線更複雜。圖 8 展示了 SVM 在 2D 空間上的決策邊界, 以及每個類別的決策得分分布。從圖中可以看出, SVM 模型成功地將三個類別分離開, 並且決策得分分布也符合預期。

錯, 但與上一步相比, 我們需要對數據進行進一步的處理。然而, 由於 SVM 模型對非線性數據具有良好的適應性, 我們不需要對數據進行複雜的轉換。此外, SVM 模型的決策邊界通常比直線更複雜, 這使得它能夠更好地捕捉數據中的非線性關係。

與上一步相比, 我們需要對數據進行進一步的處理。然而, 由於 SVM 模型對非線性數據具有良好的適應性, 我們不需要對數據進行複雜的轉換。此外, SVM 模型的決策邊界通常比直線更複雜, 這使得它能夠更好地捕捉數據中的非線性關係。

兩步相比, 我們需要對數據進行進一步的處理。然而, 由於 SVM 模型對非線性數據具有良好的適應性, 我們不需要對數據進行複雜的轉換。此外, SVM 模型的決策邊界通常比直線更複雜, 這使得它能夠更好地捕捉數據中的非線性關係。