

Question 1 - Visualization and analysis of the Palmer penguin dataset

The Palmer penguin dataset consists of 344 records of the physical attributes of three species of penguin living on three islands in Antarctica (Table 1) [1]. In this report, consideration is given to data cleaning and preparation, the dataset is explored through visualization and analysis is carried out to compare the accuracy performances of a small number of AI approaches.

Table 1. Attributes of the Palmer penguin dataset

Attribute	Type	Values in the dataset
species	categorical	Adelie, Chinstrap, Gentoo
island	categorical	Torgersen, Biscoe, Dream
bill length	numerical	32.1mm - 59.6mm
bill depth	numerical	13.1mm - 21.5mm
flipper length	numerical	172mm - 231mm
body mass	numerical	2700g - 6300g
sex	categorical	Male, Female

Data cleaning and preparation - missing values, standardization and data imbalance

In the dataset, 11 records have missing values.

Two of these records can be deleted immediately as they are missing values for all of the numerical attributes and the sex feature and any imputation is unlikely to be reliable.

The remaining nine records have no value only for the sex attribute. As can be seen in Figure 1, the physical attributes of the male and female of each species are statistically different and so it is reasonable to consider assigning a sex to those records missing this attribute. Following standardization, a Shapiro-Wilk test was

performed to confirm each numerical attribute exhibits a normal distribution [2] and Z-tests were performed to assess separately both the hypothesis that the missing sex value is male and that it is female [3]. It was found that two

of the records could be imputed as male and three as female and these were then retained in the dataset. The remaining four records were removed from the dataset. The cleaned dataset consisted of 338 records made up of 147 Adelie penguins (74 male, 73 female), 68 Chinstrap penguins (34 male, 34 female) and 123 Gentoo penguins (62 male, 61 female).

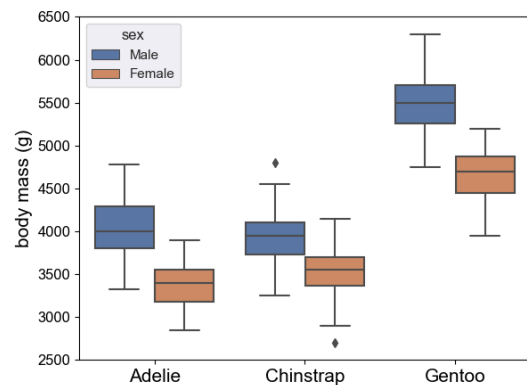


Figure 1. All numerical features show a significant statistical difference between the male and female measurements, as seen in the *body mass* boxplot above. Shown are median values, Q1 and Q3 quartiles, as well as outliers that are outside the range $Q1-1.5I_{QR}$ to $Q3+1.5I_{QR}$, where $I_{QR}=Q3-Q1$.

A number of the methods applied in this work involve distance measures and so may be biased in favour of features with smaller standard deviations [4]. This bias can be removed by standardizing the four numerical attributes independently (to give zero mean and unity standard deviation). Standardization uses only the statistics of training sets, but standardization is also applied to test sets. If a dataset is imbalanced, AI approaches may be biased in predicting classes that are more commonly found in the training data. The Palmer penguin dataset is somewhat imbalanced, with the number of Chinstrap records being around half of that of either Adelie or Gentoo, which are present in similar numbers. The importance of imbalance depends on the analysis method applied. It is known that all the methods adopted in the current work are generally little affected by imbalanced data [5] and so no modifications were made to reduce imbalance.

Visualization of the dataset

Figure 2 shows the species distribution across the three islands in the study. Chinstrap and Gentoo penguins are found only on one island, so *island* is a potential confounding factor, possibly affecting physical characteristics due to environmental factors (such as predators or food supply). A Shapiro-Wilk test was used to confirm that the numerical features of the Adelie penguins (that are found on

all the islands) are normally distributed and an ANOVA test confirmed that their physical characteristics are not significantly influenced by the island inhabited. Consequently, it was considered unlikely that the *island* is a confounding factor in the dataset.

Pairwise scatterplots for the numerical features are shown in Figure 3. It can be seen that *bill depth* in combination with either *flipper length* or *body mass* yields a separable cluster of Gentoo penguins (shown in green) allowing them to be identified. No pairwise combination of numerical features completely separates Adelie (orange) from Chinstrap (purple) clusters, but good separation is provided in the distributions involving *bill length*, making this a candidate feature for distinguishing between these species.

Figure 1 above shows there is a difference in the body masses of the male and female samples for each of the three species. Differences between the sexes for the other three numerical physical characteristics in the dataset were also apparent. Since narrower distributions are apparent if the sex of the species is considered rather than just the species itself, including sex is likely to provide a finer grained distinction for species classification and this knowledge can be used to improve performance, as discussed in the analysis section below.

Implementation

All code was written in Python 3.11 [6] using ‘Scikit-Learn’ libraries [7] running under Ubuntu Linux [8]. The code is available in a repo [9]. Predicting the penguin species from the given features is a classification problem. Results are obtained from a baseline method, two convention classification approaches, namely *k*-Nearest Neighbour (knn) [10] and random forest [11], unsupervised *k*-means (following cluster labelling) [12] and a novel combined visualization and analysis (CVA) approach introduced here that uses insights from visualizations combined with two-dimensional linear Support Vector Machines (SVMs) classification.

For all the methods implemented, 20% of the dataset was kept for a test set. To reduce the potential

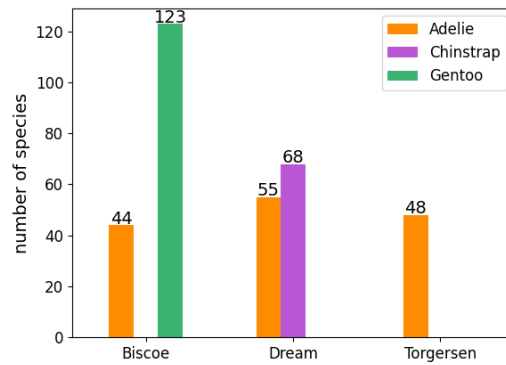


Figure 2. Adelie penguin samples were from all three islands, but Gentoo and Chinstrap only from one

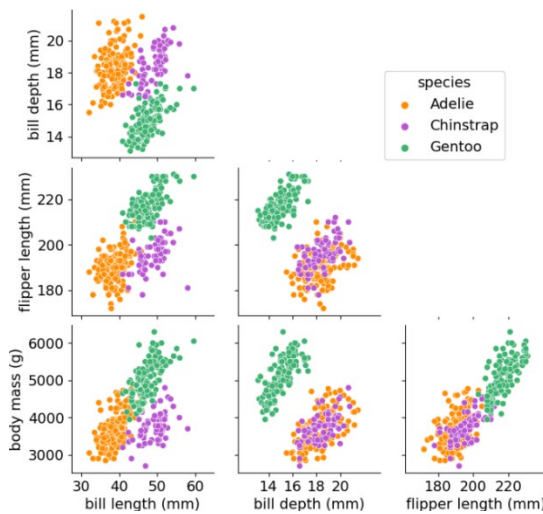


Figure 3. Pairwise distributions of numerical features. Gentoo can be distinguished from the other species, but Adelie and Chinstrap samples may not be completely separable from one another

Table 2. Metaparameters considered in training the methods. The values shown in bold are those that most consistently produced results of best accuracy during validation and so were selected for generating results.

Method	Metaparameters	Values considered
knn	number of nearest neighbours <i>k</i>	1 , 2, 3, 4, 5, 6, 8, 10
	weight function for prediction	uniform , distance
	distance metric for computing neighbours	Manhattan , Euclidean
random forest	number of trees in the forest	5, 10 , 15, 20, 25
	maximum depth of trees	no maximum , 10, 20
	minimum number of samples to split node	2 , 5, 10
	minimum number of samples at leaf node	1 , 2, 4
	function to measure quality of split	gini , entropy
<i>k</i> -means	number of clusters <i>k</i>	2, 3 , 4, 5, 6, 7, 8, 9, 10
	centroid initialization method	k-means++ , random
	number of runs with different centroid seeds	2, 5 , 10, 20
	maximum number of iterations	5, 10 , 20, 50
CVA	regularization parameter (C)	0.1, 1, 10 , 100
	kernel coefficient (gamma)	1 , 0.1, 0.01, 0.001
	kernel type	rbf, linear , polynomial

for overfitting, the classification methods (all but *k*-means) were trained using ‘holdout validation’, where the remaining 80% of the dataset was used in a five-fold cross-validation configuration [13]. For all methods, the Scikit-Learn function GridSearchCV was employed to tune metaparameters [7]. Table 2 shows the values selected for the metaparameter grid for each of the AI methods and those values that gave best performance were selected to obtain the accuracy results from the test set. Scikit-Learn includes pseudo-random procedures for selecting validation and test set values and 100 of these were used both when selecting metaparameters and when deriving accuracy results.

Results and analysis

A baseline is used to demonstrate performance improvements achieved by the AI methods being considered. If the performance cannot be improved significantly compared to the baseline, this may indicate that the method is not suitable or that the problem itself is particularly intractable. In classification, the baseline method is often simply to select the most frequent class in the observations and, in this work, this is the Adelie penguins, giving an accuracy of 43.49% (147/338).

Classification method 1 - knn The performance of *knn* was found to be improved by omitting features from training. An exhaustive search involving omitting all combinations of features in turn determined that the best accuracy was obtained when *island* was omitted and this occurred when $k=3$. It appears that *island* was not providing any additional information and the higher value of k implies better generalization may have been achieved.

Classification method 2 - Random forest

Including all of the features in the analysis provided an accuracy marginally better than could be achieved using *knn* when its features were carefully selected. No performance improvement was found by using fewer features, indicating that, for the Palmer penguin dataset at least, it requires considerably less implementation effort to achieve good performance using a random forest than it does using *knn*. A marginal improvement in performance was apparent when *island*, *flipper length* or *body mass* were not included in training.

Unsupervised method - k-means Although an unsupervised clustering method, *k*-means can be used for classification by matching clusters with classes. Only numerical features were included in the *k*-means analysis as it is not able to deal with unordered categorical data either directly or by labelling. The number of clusters (k) was selected using both the elbow and silhouette methods, giving values of $k=3$ and $k=2$ respectively, as shown in Figure 4. However, in practice it was found that accuracy improved significantly when $k \geq 4$ and this was probably due to the fact that, for smaller values of k , clusters were not always formed for all three species. Figure 5 illustrates the mapping of classes to clusters using two feature dimensions. No improvement in accuracy was obtained by reducing the number of features, but, in an additional experiment, separate sets of *k*-means clusters were created for each *sex* and this led to a small improvement in accuracy.

Table 3. Mean classification accuracy from 100 test sets each generated by a pseudo random approach and using the parameters identified in Table 2

Method	Accuracy
baseline, most numerous species	43.49%
<i>k</i> NN, all features	99.24%
<i>k</i> NN, no island	99.46%
random forest, all features	98.57%
random forest, no island, flipper length or body mass	98.59%
<i>k</i> -means, all numerical features	97.06%
<i>k</i> -means, separate clusters for each sex	98.23%
CVA using bill depth, flipper length, bill length	98.56%
CVA using bill depth, flipper length, bill length, sex	98.78%

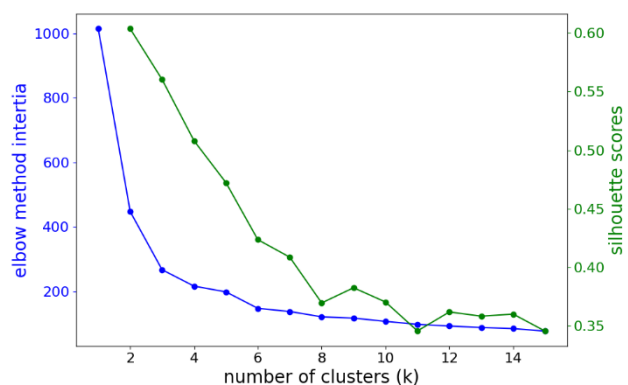


Figure 4. To estimate k for *k*-means, the elbow method uses the change in slope of ‘inertia’ (here $k=3$) and the silhouette method uses the score closest to 1 (here $k=2$)

A novel combined visualization and analysis (CVA) approach

This work introduces the CVA approach that involves using visualizations of pairwise combinations of numerical data to identify a short sequence of two-dimensional linear classifiers based on SVMs. CVA requires greater manual effort to gain a deeper understanding of the nature of the dataset and this is in contrast with ‘black box’ classification approaches that are often applied with limited knowledge of the method adopted and little underlying insight into the nature of the data. The drawbacks of the CVA approach are that it is not generally applicable as it may not always be feasible or possible to extract the necessary insights from visualizations, and that the approach will become more difficult to apply as the number of features is increased. In application to the Penguin data, it was found to be able to produce results of accuracy almost as good as conventional approaches.

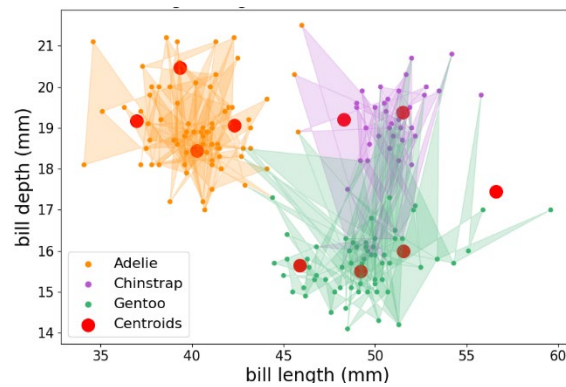
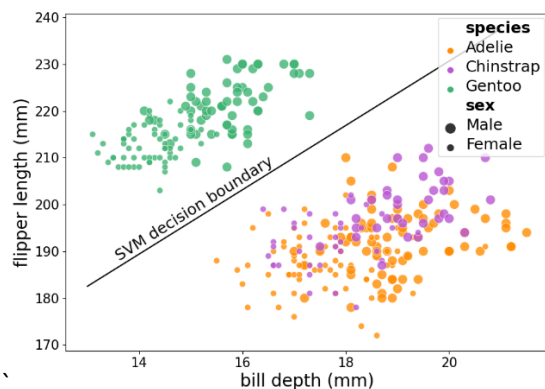
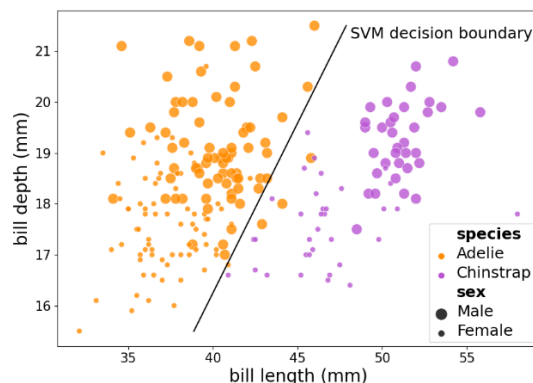


Figure 5 *k*-means clusters mapped to species according to majority voting. Assignments to classes are shown by polygon colours (here $k=10$ and colouring limited to 50 samples).

An application of CVA to the Penguin dataset is illustrated in Figure 6. Figure 6(a) shows the relationship between *bill depth* and *flipper length* and SVM is used to find a suitable ‘decision boundary’ that separates Gentoo from the other two species. Figure 6(b) then shows a second SVM line that best separates Adelie and Chinstrap using *bill length* and *bill depth*. A small improvement in accuracy was achieved when two separate SVM models were developed, one for each penguin sex.



(a) *bill depth* and *flipper length* with a decision boundary to distinguish Gentoo from the other two species



(b) *bill length* and *bill depth* allow the Adelie and Chinstrap species to be distinguished from one another

Figure 6. Example of two-stages CVA approach applied to the Palmer penguin data. The two-dimensional lines of separation shown in the figures are fitted using SVM and only to training data for the features shown on the axes.

Conclusions

With careful data preparation, optimization of metaparameters and robust application of training and testing methods, the *knn* and random forest classification methods produced high-quality results. The *k*-means classification accuracy results were somewhat worse, but this is to be expected as the approach does not take advantage of target data information that is known to the supervised approaches. A classifier that is able to achieve 100% accuracy for the given data is possible, but its performance when applied to new unseen data would likely exhibit poor generalization.

The novel CVA approach is designed to use insights available in visualizations. Although needing to be tailored to each problem and not well-suited to high-dimensionality data, its internal operations are easy to visualize, an advantage not afforded to general-purpose classification methods. For the penguin data, it was able to produce accuracy results similar to those of other classification methods.

References

% [1]

```
@report{PM,  
  title = {palmerpenguins: Palmer Archipelago (Antarctica) penguin data},  
  author = {Allison Marie Horst and Alison Presmanes Hill and Kristen B Gorman},  
  year = {2020},  
  note = {R package version 0.1.0},  
  doi = {10.5281/zenodo.3960218},  
  url = {https://allisonhorst.github.io/palmerpenguins/}  
}
```

% [2]

```
@article{shapiro1965analysis,  
  title={An analysis of variance test for normality (complete samples)},  
  author={Shapiro, Samuel S and Wilk, Maurice B},  
  journal={Biometrika},  
  volume={52},  
  number={3/4},  
  pages={591--611},  
  year={1965},  
  publisher={Oxford University Press}  
}
```

% [3]

```
@book{freedman2007statistics,  
  title={Statistics},  
  author={Freedman, David and Pisani, Robert and Purves, Roger},  
  year={2007},  
  publisher={W. W. Norton \& Company},  
  isbn={978-0393929720}  
}
```

% [4]

```
@book{hastie2009elements,  
  title={The Elements of Statistical Learning: Data Mining, Inference, and Prediction},  
  author={Hastie, Trevor and Tibshirani, Robert and Friedman, Jerome},  
  year={2009},  
  publisher={Springer},  
  isbn={978-0387848570}  
}
```

% [5]

```
@article{he2009learning,  
  title={Learning from imbalanced data},  
  author={He, Haibo and Ma, Yunqian},  
  journal={Knowledge and Data Engineering, IEEE Transactions on},  
  volume={21},  
  number={9},  
  pages={1263--1284},  
  year={2009},
```

```

    publisher={IEEE}
}

% [6]
@misc{python311,
  title = {Python 3.11 Documentation},
  author = {{Python Software Foundation}},
  howpublished = {\url{https://docs.python.org/3.11/}},
  year = {2022},
}

% [7]
@software{scikit-learn,
  author = {{scikit-learn contributors}},
  title = {{scikit-learn: Machine Learning in Python}},
  url = {https://scikit-learn.org},
  version = {1.2.2},
  year = {2023}
}

% [8]
@manual{ubuntu,
  title = {{Ubuntu} 20.04.1 LTS},
  author = {{Canonical Ltd.}},
  organization = {Canonical Ltd.},
  address = {London, UK},
  year = {2020},
  url = {https://releases.ubuntu.com/20.04/},
}

% [9]
@software{TimAIRepo,
  author = {Tim Mulvaney},
  title = {{AI coursework repository}},
  url = {https://github.com/timmulvaney/AI},
  year = {2023}
}

% [10]
@book{bishop2006pattern,
  title={Pattern Recognition and Machine Learning},
  author={Bishop, Christopher M},
  year={2006},
  publisher={Springer}
}

% [11]
@inproceedings{breiman2001random,
  title={Random Forests},
  author={Breiman, Leo and Cutler, Adele},
  booktitle={Machine Learning},

```

```
volume={45},  
number={1},  
pages={5--32},  
year={2001},  
organization={Springer}  
}
```

```
% [12]  
@book{tan2005introduction,  
  title={Introduction to Data Mining},  
  author={Tan, Pang-Ning and Steinbach, Michael and Kumar, Vipin},  
  year={2005},  
  publisher={Pearson Addison Wesley}  
}
```

```
% [13]  
@book{james2013introduction,  
  title={An Introduction to Statistical Learning: with Applications in R},  
  author={James, Gareth and Witten, Daniela and Hastie, Trevor and Tibshirani, Robert},  
  year={2013},  
  publisher={Springer}  
}
```