# How cryptic is cryptic diversity? Machine learning approaches to classifying morphological variation.

Peter D Smits[1], Kenneth D Angielczyk[2], and James F Parham[3]

[1]Committee on Evolutionary Biology, University of Chicago

[2]Department of Geology, Field Museum of Natural History

[3]Department of Geological Sciences, California State University – Fullerton

August 1, 2013

**Corresponding author:** Peter D Smits, Committee on Evolutionary Biology, University of Chicago, 1025 E. 57th Street, Culver Hall 402, Chicago, IL, 60637, USA; E-mail: psmits@uchicago.edu

## Abstract

Cryptic diversity is the phenomenon where some taxa are believed to be identifiable only based on molecular data. This is concerning because the majority of extant taxa and virtually all extinct taxa are delimited entirely via morphology. Here we address questions about whether it is possible to determine, based on morphology, if one classification hypotheses can be considered better than others. Using a combination of unsupervised and supervised machine learning methods we demonstrate a suite of approaches for better understanding differences in morphology between classes, the

1

odds of classifying one class relative to another, and what aspects of morphology best describe the differences between classes. These approaches are applied to the classification of the emydid turtle, *Emys marmorata*. This species has conflicting hypotheses of the number of meaningful subclades based on either morphological or molecular information. We compared multiple explicit classification hypotheses by characterizing variation in plastral shape and how it may be identifiably different between different classes. By splitting a large dataset of specimens into both training and testing datasets, we were also able to determine which of the classification hypotheses best corresponded to the observed plastral variation. The results from our analysis shows that the best classification of plastral variation in *Emys marmorata* is in accordance with the molecularly based hypothesis. This demonstrates that, by using alternative methods for characterizing variability, it is possible to estimate the classification scheme which most agrees with observed morphological variation. Additionally, we demonstrate that it is possible that not all examples of cryptic variation may be truly cryptic, just a product of sample size or methodology because of the extremely fine scale variation between the different classes.

(Keywords: Testudines, Emydidae, morphology, geometric morphometrics, random forests)

Cryptic diversity is the phenomenon that not all taxa can be recognized based on morphology and can only be delimited using molecular information (Clare 2011; Funk et al. 2012; Pfenninger and Schwenk 2007; Stuart et al. 2006). Conceringly, most extant taxa, and nearly all extinct taxa, are delimited based solely via morphological analysis. This phenomenon is of great concern when studying variation and diversity dynamics over long periods of time, where apparent morphological stasis may not actually reflect true diversity CITATIONS. In the case of endangered or conserved taxa, morphometric approaches for classifying and identifying taxa and populations of importance would greatly improve the ability to maintain these high risk groups. Additionally, this would allow for better classifying

extinct taxa.

Much work has been devoted to species delimitation via sequence difference (Fujita et al. 2012; Yang and Rannala 2010) while comparatively little has been devoted for case of purely morphological data. The majority of this effort has focused on identifying differences between already identified taxa (Zelditch et al. 2004) and automated taxon identification (MacLeod 2007).

Here, we address the question of how can alternative approaches and methodology improve morphology based classification. From this approach, we ask if it is possible to determine which amongst a set of classification hypotheses is best.

## *Background and system*

Differences in morphological variation between different classes has previously been analyzed using methods like linear discriminate analysis and canonical variates analysis (Zelditch et al. 2004). These methods are comparatively simple and straight forward ways of understanding the differences in morphology between classes. Also, they are very visual methods which aides in interpretation and presentation of information. Neural network models have also been introduced and applied frequently in the context of automated taxon identification along with more general applications (MacLeod 2007). Here, we used multiple alternative machine learning methods, both unsupervised and supervised, in order to compare different classification hypotheses. These methods provide different and unique advantages for understanding how to classify taxa, with what accuracy, and what these classifications are based on. Additionally, we investigate variation in continuous traits, and do not search for discrete differences between each class, instead focusing on suites of traits together.

The two major classes of machine learning methods, unsupervised and supervised, are essentially extensions of known statistical methods. Unsupervised learning methods are analogous to clustering and density estimation methods, while supervised learning methods

are analogous classification and regression models. In both cases, many of these methods are not fit via maximum likelihood and are supplemented by randomization, sorting, and partitioning algorithms along with the maximization or minimization of summary statistics in order to best estimate a general model for all data, both sampled and unsampled. The application of the alternative approaches used in this study illustrates only a sampling of the various previously derived methods for clustering observations and fitting classification models. Additionally, instead of pure classification accuracy, here we use a statistic of classification strength that reflects the rate at which taxa are both accurately and inaccurately classified (see Methods).

In this study, we investiage the subspecific classification of the western pond turtle, *Emys marmorata*. *E. marmorata* is distributed from northern Washington State, USA to Baja California, Mexico. Traditionally, *E. marmorata* was classified into three groups: the northern *E. marmorata marmorata*, the southern *E. marmorata palida*, and a central Californian intergrade zone (Seeliger 1945) UNPUBLISHED MASTERS THESIS. *E. marmorata marmorata* is differentiated from *E. marmorata palida* by the presence of a pair of triangular inguinal plates and darker neck markings. It should be noted that the triangular inguinal plates can sometimes be present in *E. marmorata palida* though they are considerably smaller.

Previous work on morphological differentiation between subspecies of *E. marmorata* focused on just the known subspecies of Seeliger (1945). SEE MASTERS THESIS AND INCLUDE INFORMATION ON THE LINEAR MORPHOMETRIC CLASSIFICATION WORK DONE PRIOR ON E MARMORATA.

More recently, *E. marmorata* was divided into four clades based on mitochondrial DNA: a northern clade, a southern clade, and eastern and western central Californian clades (Spinks and Shaffer 2005; Spinks et al. 2010). While nuclear DNA supports two major clades, one northern and one southern, Spinks et al. (2010) argue that the four clade classification is of greater conservation utility. While the mitochondrially based classification is considered

4

robust, there is no known morphological differentiation between these clades.

In this study, we attempt to estimate the best classification scheme of *E. marmorata* based on variation in plastral shape. Because of unclear geographic boundaries between subgroups of *E. marmorata*, we compare two hypotheses of morphologically based classification and two hypotheses of molecularly based classification. We hypothesize that if morphological variation corresponds to class assignment, then it should be possible to determine the best classification hypothesis of *E. marmorata* from amongst multiple candidate hypotheses. However, if morphological variation variation does not correspond to any classification hypothesis, then supervised learning model generalization performance will be poor and reflect how variation may not follow along with any of the candidate classification hypotheses.

# Materials and Methods

## *Specimens*

We collected landmark-based morphometric data from 524 adult *E. marmorata* museum specimens. Specimen classification was based on known specimen geographic information which was recorded from museum collection information. When precise latitude and longitude information was not available it was estimated from whatever locality information was present. Because the specimens used to define the subclades in Spinks and Shaffer (2005) and Spinks et al. (2010) were not available for study, all specimen classifications were based solely on this geographic information and not from explicit assignment in previous studies. Because the exact barriers between different biogeographic regions are unknown and unclear, two assignments for both the morphologically and molecularly based hypotheses were used. Each morphologically based hypothesis had three classes, while each molecular-based had four classes. In total, each specimen was given four different classifications.

## Geometric morphometrics

Following Angielczyk et al. (2011), 19 landmarks were digitized using TpsDig 2.04 (Rohlf 2005). These landmarks were chosen to maximize the description of general plastral variation. 17 of these landmarks are at the endpoints or intersection of the keratinous plastral scutes that cover the platron. 12 of these landmarks were chosen to be symmetrical across the axis of symmetry and, in order to prevent degrees of freedom and other concerns (Klingenberg et al. 2007), prior to analysis these landmarks were reflected across the axis of symmetry (i.e. midline) and the average position of each symmetrical pair was used. In cases where damage or incompleteness prevented symmetric landmarks from being determined, only the single member of the pair was used. Analysis was conducted on the resulting "half" plastra. Plastral landmark configurations were superimposed using generalized Procrustes analysis (Dryden and Mardia 1998) after which, the principal components (PC) of shape were calculated. This was done using the `shapes` package for R (Dryden 2013; R Core Team 2013).

## Machine learning analyses

*Unsupervised learning.*— In order to preserve the relationship between all landmark configurations in shape space, the dissimilarity between observations was measured using Kendall's Riemanninan shape distance or $\rho$ (Dryden and Mardia 1998; Kendall 1984). This metric was chosen because shape space, or the set of all possible shape configurations following Procrustes superimposition, is a Riemannian manifold and thus non-Euclidean (Dryden and Mardia 1998). $\rho$ varies between 0 and $\pi/2$ when there is no reflection invariance, which should not be a concern in the case of the half plastral landmark configurations used in the study.

The $\rho$ dissimilarity matrix was divisively clustered using partitioning around mediods clustering (PAM), a method similar to $k$-means clustering except that instead of minimizing the sum of squared Euclidean distances between observations and centroids, the sum of squared

dissimilarities between observations and mediods is minimized (Kaufman and Rousseeuw 1990). Because the optimal number of clusters of shape configurations in the study was unknown, being possibly three, four, or some other value, clustering solutions were estimated with the number of clusters varied between one and 40. Clustering solutions were compared using the gap statistic, which is a measure of goodness of clustering (Tibshirani et al. 2001). The gap statistic is defined

$$Gap_n(k) = E_n^*[\log(W_k)] - \log(W_k)$$

where $W_k$ is

$$W_k = \sum_{r=1}^{k} \frac{1}{2n_r} \left( \sum_{i,i' \in C_r} d_{ii'} \right)$$

. $d_{ii'}$ is the dispersion of the clustering solution or the sum of the pairwise dissimilarities between observations in each cluster and there respective mediods $(C)$ for all clusters $r$. This value is averaged and compared to the expected dispersion $(E_n^*)$ of a sample $n$ from a reference distribution. In this case, the reference distribution was estimated from 500 resamples of the dataset taking into account the original structure of the data. This analysis was conducted using the `cluster` package for R (Maechler et al. 2013) using all 524 observations.

*Supervised learning.—* The total dataset of 524 observations was split into training and testing datasets. The training dataset represented 75% of the total dataset, split proportionally by class, and was used for model fitting. The testing dataset represented the remaining 25% of the total dataset and was used after model fitting to estimate the effectiveness of each classification hypothesis and generalizability of the supervised learning models (i.e. performance in the wild).

Two different supervised learning methods were used to model the relationship between plastral shape and class: multinomial logistic regression and random forest. These methods were chosen because of various properties of these methods which allow for useful interpretations about the quality and structure of the classification.

7

Multinomial logistic regression is an extension of logistic regression, where instead of a binary response there are three or more response classes (Venables and Ripley 2002). Effectively, this type of model can be viewed as multiple, simultaneous logistic regression models for each class and the final classification of the observation being the most probable of all the constituent model results. Similar to the odds ratios calculated from the coefficients of a logistic regression, the relative risk of a classification with reference to a baseline class can be determined from the coefficients of the model. Multinomial logistic regression models were fit using the `nnet` package for R (Venables and Ripley 2002)

Random forest models are an extension of classification and regression trees (CART) (Breiman 2001; Breiman et al. 1984). Because this study relies on classification models, CARTs are explained with reference to classification but the approach is equally valid for regression. The goal of CARTs are to use a series of different features to estimate the final class. In top-down induction of decision trees for each member of a given set of predictor variables, attribute value test are used to estimate the differences between classes. This process is then repeated on each subset, called recursive partitioning. The recursion continues until the resulting observations all share the same class or no more meaningful partitions are possible. The resulting model is a tree structure by which observations are classified at each intersection via the estimated cutoff points from the attribute tests made during model fitting.

In a random forest model, many CARTs are built from a random subsample of both the features and the observations. This process is then repeated many times and the parameters of the final model was chosen as the mode of estimates from the distribution of CARTs (Breiman 2001). In addition to fitting a classification model, this procedure allows for the features to be ranked in order of importance. In the context of this study, this means that the PCs most important for describing the difference between classes can be estimated, and thus illustrate the most important variation amongst classes as opposed to just the greatest

amount of variation in the entire dataset. This is a generally important property that is useful for many other studies which want to describe and model the differences between classes and the relative importance different features. Random forest models were fit using the `randomForest` package for R (Liaw and Wiener 2002).

Supervised learning models have tuning parameters which help to increase the genearlizability of the model and prevent them from being overfit. For the supervised learning models fit in this study, tuning parameters were estimated via 10 rounds of 10-fold cross-validation (CV) across a grid search of all tuning parameter combinations. Optimal tuning parameter values were selected based on area under the receiver operating characteristic (ROC) curve. The area under the multiclass ROC curves was estimated using the all-against one strategy derived by Hand and Till (2001). This tuning process was implemented following the default grid search implemented in the `caret` package for R (Kuhn 2013).

ROC is a confusion matrix (Table 1) statistic that is a descriptor the relationship between the false positive rate ($FPR$, Eq. 1) of a classification model and the true positive rate ($TPR$, Eq. 2) of a classification model (Hastie et al. 2009). The area under the ROC curve (AUC) is a summary statistic of the quality of the classification and varies between 0.5 and 1, with values of 0.5 indicating a model that classifies no better than random and a value of 1 indicating perfect classification (Hastie et al. 2009). AUC can be used as a model selection criterion for classification models and is especially useful in cases where some if not all of the models in question were not fit via maximum likelihood where a criterion such as AICc (see below) or similar can be used (Hastie et al. 2009). It is important to note that, unlike AICc, AUC is not calculated with reference to the complexity of the model.

$$FPR = \frac{FP}{FP + TN} \tag{1}$$

$$TPR = \frac{TP}{TP + FN} \tag{2}$$

For the multinomial logistic regression models, 10 different models were fit each having sequentially more PCs as predictors in order to have models representing different levels of overall amount of shape variation and to estimate how much was necessary and sufficient to best estimate class. The maximum number of PCs allowed as predictors was 10 because of both the large number of parameters estimated per model and the necessary sample size needed to estimate that many parameters accurately. The final model was that with the lowest AICc (Akaike 1974; Burnham and Anderson 2002; Hurvich and Tsai 1989). AICc is a model selection criterion where the model with lowest AICc has the fairest variance–bias tradeoff (Burnham and Anderson 2002). Model selection was performed in this manner because the optimal number of PCs to use as predictors was not know *a priori*, and while including all of the PCs of shape would mean that all shape variability would be used to estimate class, this may cause the model to be overfit and not provide an accurate estimate of unsampled plastral variation. In addition to the AICc of each model the ΔAICc and Akaike weights are also reported. ΔAICc values are the different in AICc between the AICc best model and that model while Akaike weights are a transformation of the AICc of a model with relation to all other models being compared and measures the relative amount of information explained by that model compared to all other models (Burnham and Anderson 2002).

Random forest models are not fit using maximum likelihood so AICc based model selection was not possible. Instead, a recursive feature selection algorithm was used to choose the optimal number of PCs to include based on the AUC of the model. Following the backwards selection algorithm implemented in `caret` (Kuhn 2013), the maximum number of features were included in the initial model, there importance ranked, and the AUC of the model calculated. The lowest ranked feature was then removed, and the AUC of the model recalculated. This was repeated until only one feature, remained. Similar to the multinomial logistic regression models described above, the maximum number of PCs that could have been included in the model was 10. After each PC was removed , 10-fold CV was used to

estimate the optimal values of the tuning parameters as well as quantify the uncertainty of each model. Random forest model parameters were estimated from 1000 subtrees. Because PCs were kept in order of importance and not in relation to the amount of variance each PC described, this means that the exact PCs included in each model do not correspond to the PCs in each of the 10 multinomial logistic regressions models.

The final selected models were then used to estimate the class assignments of the training dataset. Model generality for both methods for all four classification schemes was measured using the AUC of the assignments. A distribution of AUC values was estimated for each classification scheme via 1000 nonparametric bootstrap resamples of the training dataset. The difference in distributions was assessed using pairwise Mann-Whitney U tests.

# Results

## Geometric morphometrics

The results of the PCA of the total dataset of *E. marmorata* pastral landmarks configurations demonstrates no clear or obvious groupings (Fig. 1). The first three PCs, which represent 55.0075790871272% of the total variation, are a cloud of points with no structure. Additionally, individual landmark variation is mostly circular around each landmark with some more elliptical variation observed along some midline landmarks and the most lateral landmark (Fig. 1). However, it is important to note that Procustes based superimposition attempts to evenly distribute variance around the mean landmark (Zelditch et al. 2004) and this observation should be considered cursory at best.

## Machine learning analyses

11

*Unsupervised learning.—*

Comparison of gap statistic values for the range of PAM solutions indicates that the optimal number of clusters is one (Fig. 3). The next best clustering solution had only two clusters, however there is no geographic structure to this classification scheme, with members of these clusters being seemingly randomly distributed (Fig. 2). Importantly, these clusters do not conform to the nothern and souther groups from the nuclear DNA hypothesis (Spinks et al. 2010).

Unfortunately, it was not possible to obtain enough detailed information on the sex of each *E. marmorata* specimen, thus it is difficult at best to determine if this clustering solution corresponds to sexual dimorphism between observations. Male emydid turtles are known to have a plastral concavity which may influence landmark position along the midline. However, the plastral concavity of *E. mamorata* males is considered less pronounced than in other emydid turtles. While we cannot completely rule out sexual dimorphism as the root cause of this observation, because of the very small degree of known sexual dimorphism in *E. marmorata* we are less concerned with possible effects of sexual dimorphism for the later supervised learning methods.

The gap statistic values for both three and four clusters are much lower than for one and two and are statistically identical. Interestingly, other solutions with a much greater number of clusters have higher gap statistic values though these are also not significantly different. Increasing the number of clusters does appear to improve the gap statistic enough compared to the best clustering solution to merit detailed discussion.

*Supervised learning.—*

The AICc best multinomial logistic regression models, for all four classification schemes, all include the first 10 possible PCs as features (Tables 2,3,4, and 5). The second best models for all classification schemes had the first 9 PCs as features. The ΔAICc values between

the optimal and second best model range from 2.0639 for the first morphological based classification hypothesis to 19.8349 for the second molecular based classification hypothesis (Tables 2,3,4, and 5). The first 10 PCs describe 88.6043874205075% of total variation in plastral shape.

While the $\Delta$AICc value between the optimal and second best model for the first morphological based classification hypothesis was within the range to be considered equally optimal (Burnham and Anderson 2002), for this analysis we chose to use only the AICc best model. While AICc values can not be compared between models with different responses (Burnham and Anderson 2002), we interpret the fact that the optimal model for all classification schemes is the global model as a reason to use only the AICc best model for all cases. Additionally, by using a single model for all classification schemes, this limits the number of comparisons between the bootstrap resampled distributions of the AUC values for the testing dataset (see below).

The selected number of features in the final random forest model for each classification scheme varied much more than for the multinomial logistic regression models, ranging from 6 for the first morphological based classification hypothesis to 10 for the second morphological based classification hypothesis (Fig. 4).

In the case of all models, there is a substantial increase in model performance as measured by AICc for the multinomial logistic models (Tables 2,3,4, and 5) or in AUC for the random forest models and illustrated for the multinomial logistic regression models as the number of features increases (Fig. 4).

The results from the generalization of the selected supervised learning models, measured by the distributions of the bootstrapped AUC values of the testing dataset, show that one of the molecular classification hypotheses was the best overall classification scheme (Fig. 5). For both methods, the distribution of bootstrapped AUC for the molecular hypothesis was significantly greater than all of the other classification schemes (Tables 6 and 7). Remarkably, the best

classification hypothesis was identical based on both the multinomial logistic regression and random forest models.

When the classification results of the training set for the best classification scheme based on the generalization results are compared with the references classes, the higher AUC value of the best multinomial logistic regression model compared to the best random forest model can be observed as the classifications are much closer to the reference classes (Fig. 6). The best random forest model misclassified many of the observations as the northern clade instead of the correct class. This pattern of misclassification is observable but not as exaggerated in the classifications of the multinomial logistic regression model.

This pattern of misclassification may have been caused by the subtle differences in mean shape between each of the different classes (Fig. 7). The mean shape of the northern clade is the most similar to the mean shape of the entire dataset (Fig. 7a), which may indicate that specimens that are closer to the mean shape will be systematically misclassified as the northern clade.

The results of fitting the final random forest model also include the variable importance for best separating the different classes. The selected random forest model for the best classification scheme had 7 PCs as features. The PCs included as features in the final random forest model, in descending order of importance, were PCs 3, 8, 6, 1, 4, 2 and 10. Of these 7 features, the first three are illustrated here (Fig. 8) in descending order of importance.

The first two most important features describe different aspects of variation (Fig. 9). The third and most important PC describes variation in the relative position of landmarks on anterior and posterior portions of the plastron and represents 9.6697% of total variation. The eighth and second most important PC mostly describes variation in landmarks along the midline of the plastron and represents 4.1061% of total variation. The major variations along these axes correspond well to the differences between the mean shape of each class (Fig. 7) where major class differences seem based on the relative ballooning or shrinking of the

14

anterior and posterior portions of the plastron together.

The relative risk values for classification from the multinomial logistic regression model, based on the three most important PCs, demonstrate that while an individual PC might not be sufficient to differentiate any individual class from the northern group, given multiple features the odds of determining the correct classification increase (Fig. 10). Changes along the second most important axes contribute very little to increasing the odds of classification. This is observable from the class histograms of PC 8 (Fig. 9). Changes along the first and third most important axes contribute more obviously to increasing the odds of correctly identifying the class of an observation, a result that is observable in both the relative risk (Fig. 10) and the different class histograms of the PCs (Fig. 9).

## DISCUSSION

The results of this study support the mitochrondial based classification hypothesis of *E. marmorata* (Spinks and Shaffer 2005; Spinks et al. 2010). This is contrary to the original classification of *E. marmorata* (Seeliger 1945) MASTERS THESIS and lends credence to the idea that at least some aspect of cryptic diversity is a product of sample size, methodology, or both.

The lack of coherent geographical subclass assignment from PAM clustering (Fig. 3) as well as the large number of features necessary before no increase in AUC for all models (Fig. 4) indicates that the morphological variation between classes is extremely fine grained. This was also exemplified by the small differences between mean class shapes of the final chosen classification scheme (Fig. 9).

The approaches presented here for supervised learning analysis of the landmark variation represent a compromise between explicitly modeling all shape variation and preventing models from being overfit and ungeneralizable. While all aspects of shape may be evolving

simultaneously, and not along individual PCs, including all shape variation in each model might increase model complexity beyond a reasonable level for the sample size and possibly the necessary complexity to accurately model the response. Additionally, because only individual PCs are used as features in the models, this does not accurately represent shape evolution and how exactly different classes might be evolving in relation to each other. However, this compromise is not without its advantages. Because both AICc and AUC values improved rapidly with increased model complexity (Fig. 4), this helped demonstrate how fine scale the actual variation between classes was. Additionally, the relative risk values from the mulitinomial logistic regression models demonstrate that a single PC is probably not sufficient for estimating the class of an observation, but that given a set of PCs this classification would be more accurate (Fig. 10). The results of the relative risk values also shows that while a variable can be considered important for the random forest models (Fig. 8, 9) they each might not be considered useful on their own (Fig. 10).

Ultimately, it would be useful to not require such explicit classification hypotheses, especially when concerned about possible cryptic variation in extinct taxa. The only unsupervised method employed in this study, PAM, is rather simple and not model based. A more useful approach would be to employ various model based clustering approaches (Fraley and Raftery 2002; Zhong and Ghosh 2003). In this manner, a series of candidate models can be compared via model comparison methods, such as AIC or Bayes factors (Fraley and Raftery 2002), in order to asses the best clustering solution.

In this study we have demonstrated that, using alternative methodology to that which is most frequently applied, it is possible to determine which classification scheme best matches variation in a taxon amongst a set of alternative hypotheses. The observed plastral variation of *E. marmorata* is most consistent with the mitochondrial based hypothesis of Spinks and Shaffer (2005) and Spinks et al. (2010) and not with the original morphology based hypothesis of Seeliger (1945). We have also demonstrated the utility of various machine

learning approaches to understanding the structure of variation in morphometric data. Specifically, methods for better understanding misclassification and identifying which is the most important for delimiting different classes. These methods represent new applications which may be important for future studies on class-based morphological comparison and variation, both in the context of cryptic diversity and with known classifications.

## Acknowledgements

*

References

Akaike, H. 1974. A new look at the statistical model identification. IEEE Transactions on Automatic Control 19:716–723.

Angielczyk, K. D., C. R. Feldman, and G. R. Miller. 2011. Adaptive evolution of plastron shape in emydine turtles. Evolution 65:377–394.

Breiman, L. 2001. Random Forests. Machine Learning 45:5–32.

Breiman, L., J. Friedman, C. J. Stone, and R. A. Olshen. 1984. Classification and regression trees. Wadsworth International Group, Belmont.

Burnham, K. P. and D. R. Anderson. 2002. Model selection and multi-model inference: a practical information-theoretic approach. 2nd ed. Springer, New York.

Clare, E. L. 2011. Cryptic species? Patterns of maternal and paternal gene flow in eight neotropical bats. PloS one 6:e21460.

Dryden, I. L. 2013. shapes: Statistical shape analysis. R package version 1.1-8.

Dryden, I. L. and K. Y. Mardia. 1998. Statistical shape analysis. Wiley, New York.

Fraley, C. and A. E. Raftery. 2002. Model-Based Clustering, Discriminant Analysis, and Density Estimation. Journal of the American Statistical Association 97:611–631.

Fujita, M. K., A. D. Leaché, F. T. Burbrink, J. a. McGuire, and C. Moritz. 2012. Coalescent-based species delimitation in an integrative taxonomy. Trends in ecology & evolution 27:480–8.

Funk, W. C., M. Caminer, and S. R. Ron. 2012. High levels of cryptic species diversity uncovered in Amazonian frogs. Proceedings of the Royal Society B: Biological Sciences 279:1806–14.

Hand, D. J. and R. J. Till. 2001. A Simple Generalisation of the Area Under the ROC Curve for Multiple Class Classification Problems. Machine Learning 45:171–186.

Hastie, T., R. Tibshirani, and J. Friedman. 2009. The elements of statistical learning: data mining, inference, and prediction. 2nd ed. Springer, New York.

Hurvich, C. M. and C.-L. Tsai. 1989. Regression and time series model selection in small samples. Biometrika 76:297–307.

Kaufman, L. and P. J. Rousseeuw. 1990. Finding groups in data : an introduction to cluster analysis. Wiley, New York.

Kendall, D. G. 1984. Shape Manifolds, Procrustean Metrics, and Complex Projective Spaces. Bulletin of the London Mathematical Society 16:81–121.

Klingenberg, C. P., M. Barluenga, and A. Meyer. 2007. Shape analysis of symetric structures: quantifying variation among individuals and asymmetry. Evolution 56:1909–1920.

Kuhn, M. 2013. caret: Classification and Regression Training. R package version 5.15-61.

Liaw, A. and M. Wiener. 2002. Classification and regression by randomforest. R News 2:18–22.

MacLeod, N. 2007. Automated taxon identification in systematics: theory, approaches and applications. CRC Press, Boca Raton.

Maechler, M., P. Rousseeuw, A. Struyf, M. Hubert, and K. Hornik. 2013. cluster: Cluster Analysis Basics and Extensions. R package version 1.14.4.

Pfenninger, M. and K. Schwenk. 2007. Cryptic animal species are homogeneously distributed among taxa and biogeographical regions. BMC evolutionary biology 7:121.

R Core Team. 2013. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing Vienna, Austria.

Rohlf, F. J. 2005. TpsDig 2.04.

Seeliger, L. M. 1945. Variation in the Pacific Mud Turtle. Copeia 1945:150–159.

Spinks, P. Q. and H. B. Shaffer. 2005. Range-wide molecular analysis of the western pond turtle (Emys marmorata): cryptic variation, isolation by distance, and their conservation implications. Molecular ecology 14:2047–64.

Spinks, P. Q., R. C. Thomson, and H. B. Shaffer. 2010. Nuclear gene phylogeography reveals the historical legacy of an ancient inland sea on lineages of the western pond turtle, Emys marmorata in California. Molecular ecology 19:542–56.

Stuart, B. L., R. F. Inger, and H. K. Voris. 2006. High level of cryptic species diversity revealed by sympatric lineages of Southeast Asian forest frogs. Biology letters 2:470–4.

Tibshirani, R., G. Walther, and T. Hastie. 2001. Estimating the number of clusters in a data set via the gap statistic. Journal of the Royal Statistical Society: Series B (Statistical Methodology) 63:411–423.

Venables, W. and B. D. Ripley. 2002. Modern applied statistics with S. 4th ed. Springer, New York.

Yang, Z. and B. Rannala. 2010. Bayesian species delimitation using multilocus sequence data. Proceedings of the National Academy of Sciences 107:9264–9.

Zelditch, M. L., D. L. Swiderski, and H. D. Sheets. 2004. Geometric morphometrics for biologists: a primer. Elsevier Academic Press, Amsterdam.

Zhong, S. and J. Ghosh. 2003. A unified framework for model-based clustering. The Journal of Machine Learning Research 4:1001–1037.

|  |  | Predicted class | |
| --- | --- | --- | --- |
|  |  | 1 | 0 |
| Actual class | 1 | TRUE POSITIVE | FALSE NEGATIVE |
|  | 0 | FALSE POSITIVE | TRUE NEGATIVE |

Table 1: Example confusion matrix. The columns correspond to the predicted class of an observation, while the rows correspond to the actual class of that observation. Depending on the type match between the prediction and reality, four different outcomes are possible: true positive (TP), false negative (FN), false positive (FP), and true negative (TN). These four quantities are used for calculating all confusion matrix statistics. Each of these values is an integer and the sum of the number of occurrences of that event during classification.
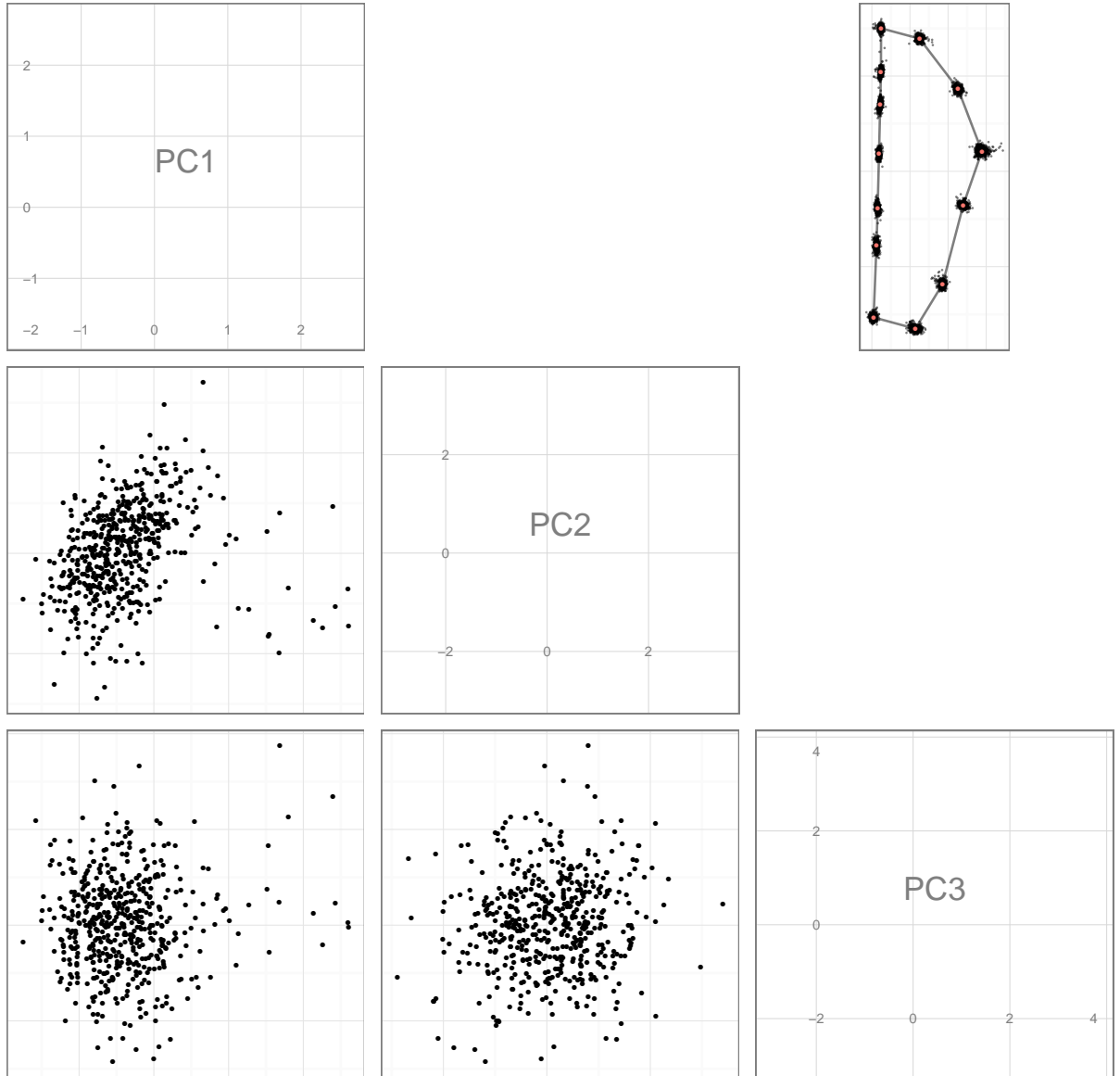
Figure 1: Results from PCA of the Procrustes superimposed "half" plastral landmarks. Depicted here are the for three PCs (lower triangle) and the mean shape with observed variance around each point (upper right). The first three PCs account for total 0 of the variance in plastral shape.
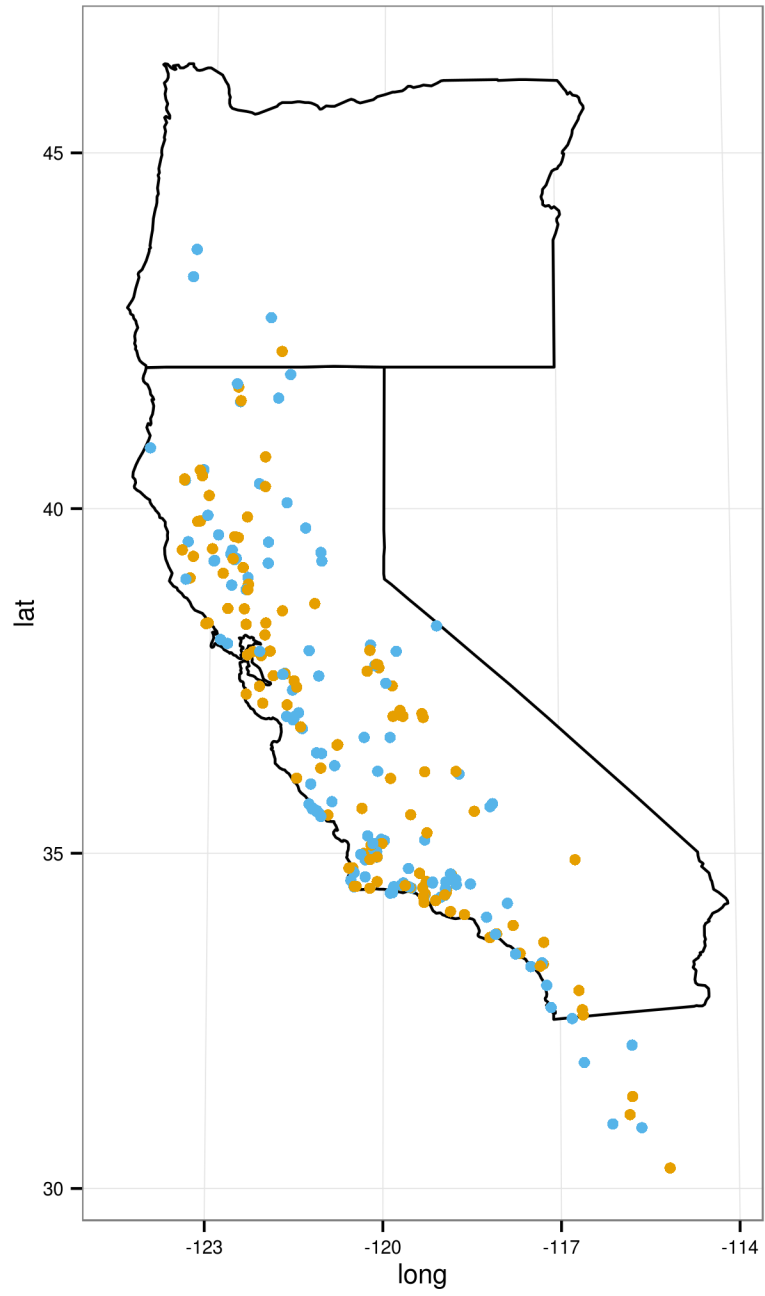
Figure 2: Clustering solution for PAM with two mediods for the entire set of observed *E. marmorata*. Clustering was based entirely on the $\rho$ dissimilarity matrix of "half" plastral landmark configurations following Procrustes superimposition.
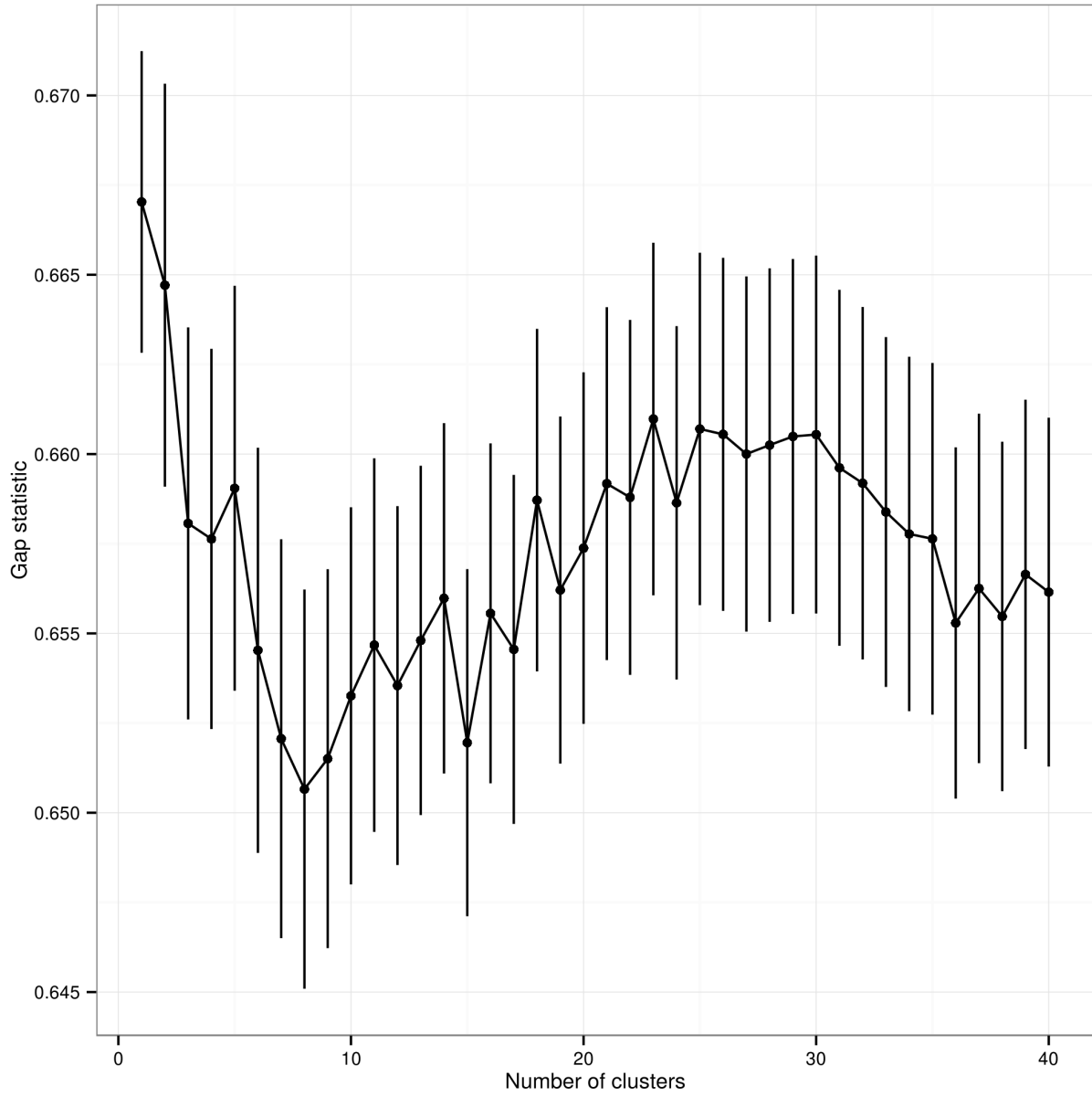
Figure 3: Gap statistic values for PAM clustering results for the $\rho$ dissimliarity matrix of plastron shape. Error bars are standard errors estimated via 500 bootstrap samples.
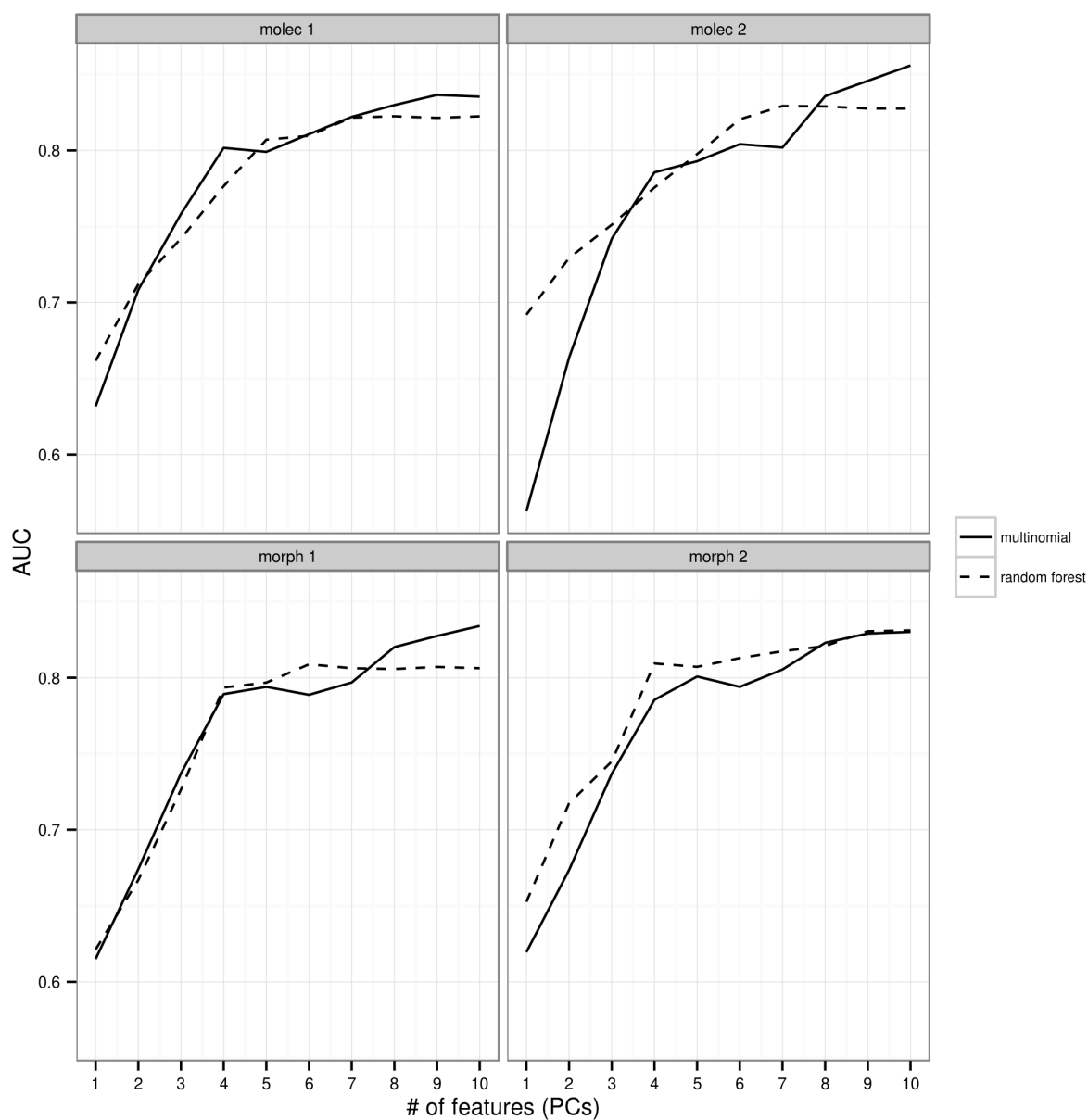
Figure 4: Effect of increasing the number of PCs as features, or predictors, of classification of plastra for all four classification schemes. As the total number of features increase, AUC increases until eventually leveling off. Both multinomial logisitic regression and random forest models are illustrated here, though AUC based model selection was only performed for random forest models.
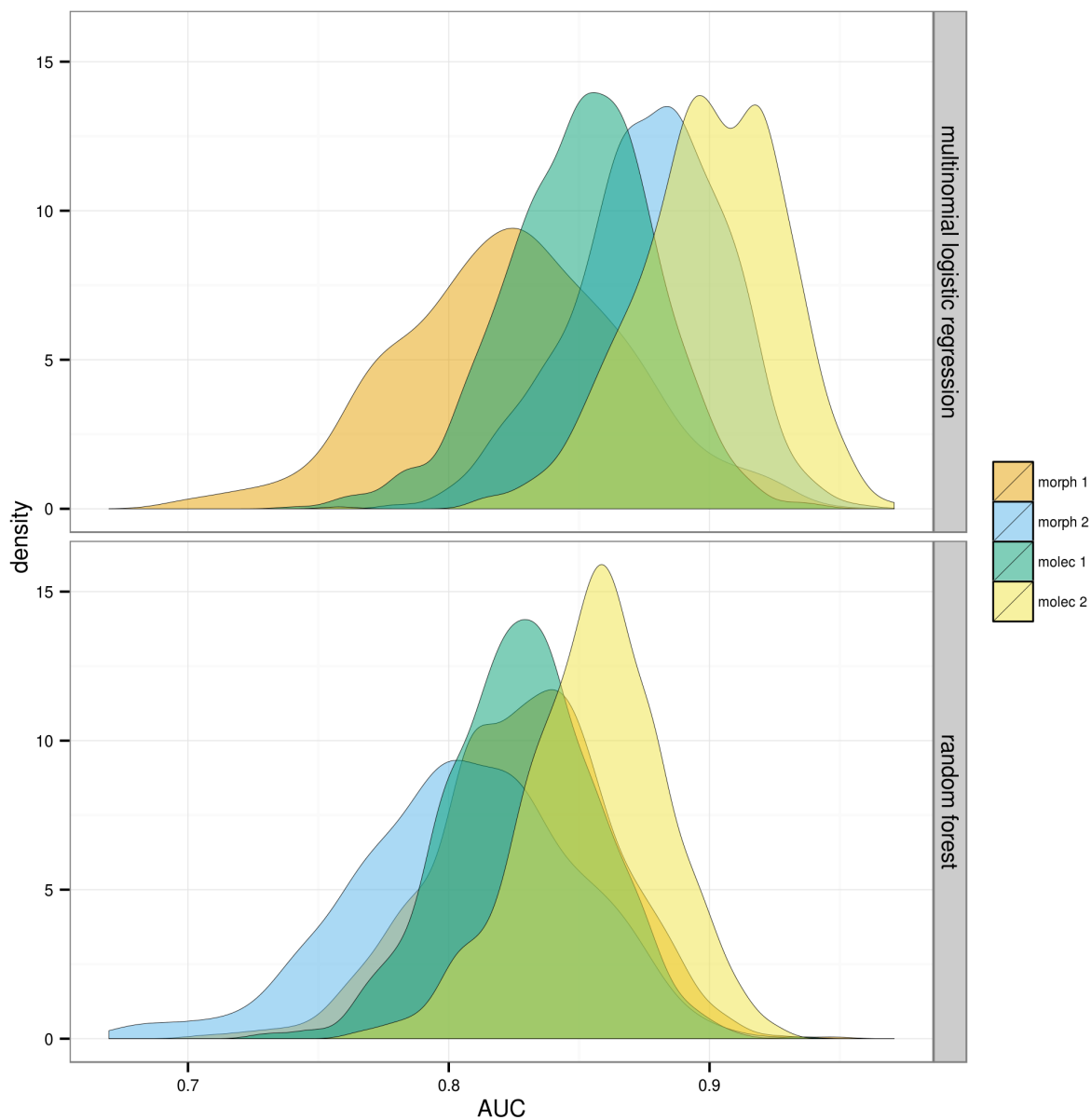
Figure 5: Density estimates of AUC values of predictions of the testing dataset of plastra from 1000 bootstrap resamples. The top facet corresponds to values using the optimal multinomial logistic regression model, as chosen by minimum AICc value. The bottom facet corresponds to the values using the optimal random forest model, as chosen by maximum AUC value.

Figure 6: Comparison between reference classification of testing data set and the estimated classifications based on the selected multinomial logistic regression and random forest models, from left to right respectively. Classification corresponds to the four classes as suggested by the hypothesis of Spinks and Shaffer (2005) and Spinks et al. (2010).

(a) Northern

(b) Eastern

(c) Western

(d) Southern

Figure 7: Thin-plate splines for each of the four classes from the best classification hypothesis based on the generalization results (Fig. 5). The four different classes are labeled according to the biogeographic groups as depicted in figure 6. The deformations are depicted with 2x magnification from base.

Figure 8: Pairs plot of the first three most important variables of the optimal random forest model of turtle plastral shape. The variables descend in importance from the upper left to the lower right. The observations are colored as in figure 6. The bottom row are histograms of classification occurrences along the PCs.

Figure 9: Landmark variation along the two most important features (PCs) based on the final random forest model. The first row corresponds to the third PC and the second corresponds to eighth PC. Landmark configurations are minimum observed on that PC, mean shape, and maximum observed on that PC.

Figure 10: Forest plot of the relative risk, with 95% confidence intervals, of classifying a give specimen based on the first three most important variables according to the random forest model. Relative risk values are calculated from the coefficients of the multinomial logistic regression model. All risks are relative to the northern group from Spinks and Shaffer (2005); Spinks et al. (2010). Variable importance is from left to right.

| (Intercept) | PC1 | PC2 | PC3 | PC4 | PC5 | PC6 | PC7 | PC8 | PC9 | PC10 | df | logLik | AICc | delta | weight |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| + | + | + | + | + | + | + | + | + | + | + | 22.00 | -245.34 | 537.41 | 0.00 | 0.74 |
| + | + | + | + | + | + | + | + | + | + |  | 20.00 | -248.61 | 539.48 | 2.06 | 0.26 |
| + | + | + | + | + | + | + | + | + |  |  | 18.00 | -255.73 | 549.28 | 11.87 | 0.00 |
| + | + | + | + | + | + | + | + |  |  |  | 16.00 | -267.59 | 568.62 | 31.21 | 0.00 |
| + | + | + | + | + | + |  |  |  |  |  | 12.00 | -283.44 | 591.70 | 54.29 | 0.00 |
| + | + | + | + | + |  |  |  |  |  |  | 10.00 | -286.02 | 592.61 | 55.20 | 0.00 |
| + | + | + | + | + | + | + |  |  |  |  | 14.00 | -282.24 | 593.59 | 56.17 | 0.00 |
| + | + | + | + |  |  |  |  |  |  |  | 8.00 | -307.43 | 631.23 | 93.81 | 0.00 |
| + | + | + |  |  |  |  |  |  |  |  | 6.00 | -340.94 | 694.09 | 156.68 | 0.00 |
| + | + |  |  |  |  |  |  |  |  |  | 4.00 | -345.95 | 700.01 | 162.60 | 0.00 |

Table 2: Model selection table for the multinomial logistic regression models of the first morphologically based classification hypothesis. This classification hypothesis corresponds to "morph 1" also depicted in figures 4 and 5. This hypothesis is based on Seeliger (1945). The column "delta" corresponds to the $\delta$AICc values of each model, while "weights" correspond to the Akaike weight of that model relative to all others.

| (Intercept) | PC1 | PC2 | PC3 | PC4 | PC5 | PC6 | PC7 | PC8 | PC9 | PC10 | df | logLik | AICc | delta | weight |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| + | + | + | + | + | + | + | + | + | + | + | 22.00 | -242.26 | 531.25 | 0.00 | 0.86 |
| + | + | + | + | + | + | + | + | + | + | | 20.00 | -246.34 | 534.94 | 3.70 | 0.14 |
| + | + | + | + | + | + | + | + | + | | | 18.00 | -252.76 | 543.34 | 12.10 | 0.00 |
| + | + | + | + | + | + | + | + | | | | 16.00 | -263.19 | 559.83 | 28.59 | 0.00 |
| + | + | + | + | + | + | | | | | | 12.00 | -275.48 | 575.79 | 44.54 | 0.00 |
| + | + | + | + | + | + | + | | | | | 14.00 | -273.87 | 576.85 | 45.60 | 0.00 |
| + | + | + | + | + | | | | | | | 10.00 | -280.27 | 581.11 | 49.86 | 0.00 |
| + | + | + | + | | | | | | | | 8.00 | -303.54 | 623.46 | 92.22 | 0.00 |
| + | + | + | | | | | | | | | 6.00 | -342.40 | 697.01 | 165.76 | 0.00 |
| + | + | | | | | | | | | | 4.00 | -348.78 | 705.67 | 174.42 | 0.00 |

Table 3: Model selection table for the multinomial logistic regression models of the first morphologically based classification hypothesis. This classification hypothesis corresponds to "morph 2" also depicted in figures 4 and 5. This hypothesis is based on Seeliger (1945). The column "delta" corresponds to the $\delta$AICc values of each model, while "weights" correspond to the Akaike weight of that model relative to all others.

| (Intercept) | PC1 | PC2 | PC3 | PC4 | PC5 | PC6 | PC7 | PC8 | PC9 | PC10 | df | logLik | AICc | delta | weight |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| + | + | + | + | + | + | + | + | + | + | + | 33.00 | -297.91 | 668.06 | 0.00 | 0.97 |
| + | + | + | + | + | + | + | + | + | + |  | 30.00 | -305.12 | 675.37 | 7.30 | 0.03 |
| + | + | + | + | + | + | + | + | + |  |  | 27.00 | -314.81 | 687.76 | 19.70 | 0.00 |
| + | + | + | + | + | + | + | + |  |  |  | 24.00 | -329.53 | 710.30 | 42.24 | 0.00 |
| + | + | + | + | + | + | + |  |  |  |  | 21.00 | -342.93 | 730.35 | 62.28 | 0.00 |
| + | + | + | + | + | + |  |  |  |  |  | 18.00 | -353.86 | 745.55 | 77.49 | 0.00 |
| + | + | + | + | + |  |  |  |  |  |  | 15.00 | -357.66 | 746.59 | 78.53 | 0.00 |
| + | + | + | + |  |  |  |  |  |  |  | 12.00 | -385.36 | 795.54 | 127.48 | 0.00 |
| + | + | + |  |  |  |  |  |  |  |  | 9.00 | -424.50 | 867.48 | 199.41 | 0.00 |
| + | + |  |  |  |  |  |  |  |  |  | 6.00 | -442.85 | 897.91 | 229.85 | 0.00 |

Table 4: Model selection table for the multinomial logistic regression models of the first morphologically based classification hypothesis. This classification hypothesis corresponds to "molec 1" also depicted in figures 4 and 5. This hypothesis is based on Seeliger (1945). The column "delta" corresponds to the $\delta$AICc values of each model, while "weights" correspond to the Akaike weight of that model relative to all others.

| (Intercept) | PC1 | PC2 | PC3 | PC4 | PC5 | PC6 | PC7 | PC8 | PC9 | PC10 | df | logLik | AICc | delta | weight |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| + | + | + | + | + | + | + | + | + | + | + | 33.00 | -253.75 | 579.72 | 0.00 | 1.00 |
| + | + | + | + | + | + | + | + | + | + |   | 30.00 | -267.22 | 599.55 | 19.83 | 0.00 |
| + | + | + | + | + | + | + | + | + |   |   | 27.00 | -275.24 | 608.60 | 28.88 | 0.00 |
| + | + | + | + | + | + | + | + |   |   |   | 24.00 | -302.99 | 657.22 | 77.51 | 0.00 |
| + | + | + | + | + | + | + |   |   |   |   | 21.00 | -307.70 | 659.88 | 80.16 | 0.00 |
| + | + | + | + | + |   |   |   |   |   |   | 15.00 | -327.52 | 686.30 | 106.59 | 0.00 |
| + | + | + | + | + | + |   |   |   |   |   | 18.00 | -324.35 | 686.51 | 106.79 | 0.00 |
| + | + | + | + |   |   |   |   |   |   |   | 12.00 | -350.14 | 725.11 | 145.39 | 0.00 |
| + | + | + |   |   |   |   |   |   |   |   | 9.00 | -390.78 | 800.03 | 220.32 | 0.00 |
| + | + |   |   |   |   |   |   |   |   |   | 6.00 | -405.14 | 822.49 | 242.77 | 0.00 |

Table 5: Model selection table for the multinomial logistic regression models of the first morphologically based classification hypothesis. This classification hypothesis corresponds to "molec 2" also depicted in figures 4 and 5. This hypothesis is based on Seeliger (1945). The column "delta" corresponds to the $\delta$AICc values of each model, while "weights" correspond to the Akaike weight of that model relative to all others.

|         | morph 1 | morph 2 | molec 1 | molec 2 |
|---------|---------|---------|---------|---------|
| morph 1 |         |         |         |         |
| morph 2 | 0.00    |         |         |         |
| molec 1 | 0.00    | 0.00    |         |         |
| molec 2 | 0.00    | 0.00    | 0.00    |         |

Table 6: Results from pairwise Mann-Whitney U test between the AUC distributions of the generalizations of the multinomial logistic regression models. Labels correspond to those in Figure 5. Values of 0 correspond to p-values lower than 0.01. P-values were corrected for multiple comparison using the Holm method CITATION.

|         | morph 1 | morph 2 | molec 1 | molec 2 |
|---------|---------|---------|---------|---------|
| morph 1 |         |         |         |         |
| morph 2 | 0.00    |         |         |         |
| molec 1 | 0.51    | 0.00    |         |         |
| molec 2 | 0.00    | 0.00    | 0.00    |         |

Table 7: Results from pairwise Mann-Whitney U test between the AUC distributions of the generalizations of the random forest models. Labels correspond to those in Figure 5. Values of 0 correspond to p-values lower than 0.01. P-values were corrected for multiple comparison using the Holm method CITATION.