

ADAPTIVE CONFORMAL SETS UNDER DISTRIBUTION SHIFT: USING ENSEMBLE DISAGREEMENT AS AN EPISTEMIC NORMALIZER FOR EEL-GRASS SEGMENTATION UNDER TEMPORAL DRIFT

A Thesis  
presented to  
the Faculty of California Polytechnic State University,  
San Luis Obispo

In Partial Fulfillment  
of the Requirements for the Degree  
Master of Science in Statistics

by  
Dillon Murphy  
March 2026

© 2026

Dillon Murphy

ALL RIGHTS RESERVED

## COMMITTEE MEMBERSHIP

TITLE: Adaptive Conformal Sets Under Distribution Shift: Using Ensemble Disagreement as an Epistemic Normalizer for Eelgrass Segmentation under Temporal Drift

AUTHOR: Dillon Murphy

DATE SUBMITTED: March 2026

COMMITTEE CHAIR: Kelly Bodwin, Ph.D.

Associate Professor of Statistics and Data Science

COMMITTEE MEMBER: Jonathon Ventura, Ph.D.

Associate Professor of Computer Science and Software  
Engineering

COMMITTEE MEMBER: Andrew Fricker, Ph.D.

Associate Professor of Geography

## ABSTRACT

Adaptive Conformal Sets Under Distribution Shift: Using Ensemble Disagreement as an Epistemic Normalizer for Eelgrass Segmentation under Temporal Drift

Dillon Murphy

Semantic segmentation of eelgrass (*Zostera marina*) from high-resolution drone-based imagery is essential for coastal habitat monitoring, restoration, and management, as these habitats see rapid changes due to climate change and human-induced influences. However, the reliability of generalizing these classification models not only relies on high-accuracy segmentation but also on rigorous uncertainty quantification that holds up when conditions change across years or locations. Conformal prediction (CP) converts a classifier's output into prediction sets with finite-sample marginal coverage; however, the standard score  $s = 1 - p_y$  can under-cover in hard or out-of-distribution (OOD) regions under temporal drift. Inspired by heteroskedastic conformal regression, this study proposes normalizing conformal scores by an epistemic proxy, ensemble disagreement. These scores re-scale nonconformity by a data-driven estimation of local difficulty learned on the calibration set, so that the global quantile is conservative where OOD risk is high. This study evaluates (i) a parametric form  $s' = (1 - p_y)/(1 + \lambda V)$ , where  $V$  is ensemble variance of probabilities, and (ii) a nonparametric normalization  $s' = s/\hat{\sigma}(V)$ , where  $\hat{\sigma}(V) = E[s|V]$  is estimated through quantile binning and interpolation. On drone imagery of the Morro Bay estuary, California (2018-2022), the models were trained on 2018-2021, calibrated on 2021, and evaluated on 2022 on OOD points. Normalized scores recover much of the lost coverage seen by vanilla split conformal prediction and approximate target, increased the percentage of singletons, and reduced the spatial coverage variability. The study provides a simple-to-implement drop-in framework to regain coverage in difficult regions and out-of-distribution data without retraining or labels at test time.

Keywords: conformal prediction, adaptive sets, uncertainty, eelgrass, remote sensing.

## ACKNOWLEDGMENTS

This page is not required, but if you have received funding for your research or assistance or guidance that you feel should be noted, it belongs on this page.

## TABLE OF CONTENTS

	Page
LIST OF TABLES . . . . .	vii
LIST OF FIGURES . . . . .	viii
CHAPTER	
1 Introduction . . . . .	1
2 Background and Motivation . . . . .	3
2.1 Uncertainty in Deep Learning . . . . .	3
2.2 Deep Ensembles . . . . .	4
2.3 Semantic Segmentation in Environmental Monitoring . . . . .	4
2.4 Conformal Prediction Foundations . . . . .	4
2.5 Adaptive or Shift-Aware CP . . . . .	4
3 Methodology . . . . .	5
3.1 Data . . . . .	5
3.2 Model Training + Pipeline . . . . .	5
3.3 Uncertainty Analysis . . . . .	5
3.4 Conformalizing . . . . .	5
3.4.1 Split CP for Classification . . . . .	5
3.4.2 Variance-Aware Score Normalization . . . . .	5
3.5 CP Evaluation . . . . .	5
4 Results . . . . .	6
4.1 Ensemble Results . . . . .	6
4.2 In-Distribution Evaluation . . . . .	6
4.3 Temporal OOD (2022) . . . . .	6
4.4 Class Conditional Coverage . . . . .	6
4.5 Spatial Robustness within 2022 . . . . .	6

4.6 Sensitivity . . . . .	6
5 Discussion . . . . .	7
6 Limitations . . . . .	8
7 Conclusion . . . . .	9
7.1 Computational Details . . . . .	9
REFERENCES . . . . .	10
APPENDICES	
A Super Cool Fancy Function . . . . .	11

## LIST OF TABLES

Table	Page
-------	------

## LIST OF FIGURES

Figure

Page

## Chapter 1

### INTRODUCTION

Eelgrass (*Zostera marina*) is a seagrass species of temperate waters that provides an anchor to coastal ecosystems by providing sediment stabilization, carbon sequestration, eliminating contamination, and providing habitats for protected species. Warming oceans, storm regimes, and local disturbance have quickly reshaped the intertidal morphology and habitat extent. As these threats increase, there is a need to monitor changes to determine appropriate environmental management. In recent years, deep networks have made this intertidal mapping practical in complex imagery, achieving strong pixel-wise accuracy in areas such as Morro Bay, California (Tallam et al., 2023). However, the reliability of general utilization of these classification models not only relies on high-accuracy segmentation, but also on rigorous uncertainty quantification under temporal and spatial distribution shift. Imagery collected over different years, at different tides or seasons, can erode the model's calibration, causing it to be over-confident. This study builds off these models to create a post hoc uncertainty framework.

Uncertainty quantification is a requirement for flagging ambiguous areas. To answer this, split conformal prediction provides an appealing framework. It offers a wrapping of any probabilistic predictor to produce prediction sets with finite-sample marginal coverage under exchangeability, simply by using a calibration set and a quantile nonconformity score (Angelopoulos & Bates, 2022). In classification, a natural choice of nonconformity score is  $s = 1 - p_y$ , yielding a prediction set of all labels whose scores do not exceed a global threshold learned on calibration. With stable data, this can be simple and effective; however, in reality, the distribution of scores can differ by difficulty. A threshold calibrated over one year can underperform when a distributional shift occurs, such as temporal drift.

To remedy this, it is necessary to make the conformal prediction adaptive to drift or local difficulty. First, work in adaptive or locally weighted CP in regression has shown that normalizing scores by a measure of local difficulty can make a single global quantile more appropriate across heterogeneous inputs, which improves approximate conditional behavior (Dewolf, De Baets, & Waegeman, 2025). Second, research on uncertainty in deep learning has demonstrated that diversified deep ensembles provide practical, shift-aware signals of difficulty that tend to grow on out-of-distribution inputs (Lakshminarayanan, Pritzel, & Blundell, 2017). Using these ideas, this study implements a simple strategy for rescaling the classification score by a calibration-learned difficulty proxy, so that the global threshold becomes conservative under drift, without the need for recalibration or labels.

This thesis develops this design for eelgrass segmentation of multi-year drone imagery. The proposed normalizers are deliberately simple to adopt into existing pipelines. This includes: 1. a linear normalization  $s' = (1 - p_y)/(1 + \lambda V)$  where  $V$  summarizes ensemble disagreement and  $\lambda$  controls degree of conservativeness and 2. a learned normalizer  $s' = s/\hat{\sigma}(V)$  where  $\hat{\sigma}(V) = E[s|V]$  is estimated on the calibration split via quantile binning and interpolation. Both approaches retain standard CP workflow and require no labels at test time or retraining.

## Chapter 2

### BACKGROUND AND MOTIVATION

Effective eelgrass monitoring requires models that both provide accurate predictions, as well as how certain they are for those predictions. This section reviews the uncertainty concepts used in deep learning.

#### 2.1 Uncertainty in Deep Learning

Uncertainty in predictive modeling is commonly broken up into two components: *epistemic* and *aleatoric*. These two components collectively comprise the total uncertainty in a model’s predictions. Aleatoric uncertainty represents the noise inherent in observations, which is often impossible to remove, even with the addition of more data (e.g. sensor noise or true label ambiguity). Even a perfect model cannot remove it. Epistemic uncertainty represents the uncertainty of the model itself, which can be reduced with additional or higher-quality data (Kendall & Gal, 2017).

For modern vision models, raw softmax probabilities are often miscalibrated, meaning the outputted confidence does not match empirical accuracy. This can be especially prevalent under shift, where models can be confident yet wrong. Post-hoc calibration, such as temperature scaling, rescales logits via  $\text{softmax}(z/T)$  and improves in-distribution calibration without changing accuracy (Guo, Pleiss, Sun, & Weinberger, 2017). Under shift, however, calibration is not enough, as the structure of the errors can change. Thus, it requires an uncertainty representation that integrates decision rules that signal the difficulty of the prediction.

**2.2 Deep Ensembles**

**2.3 Semantic Segmentation in Environmental Monitoring**

**2.4 Conformal Prediction Foundations**

**2.5 Adaptive or Shift-Aware CP**

## Chapter 3

### METHODOLOGY

#### **3.1 Data**

#### **3.2 Model Training + Pipeline**

#### **3.3 Uncertainty Analysis**

#### **3.4 Conformalizing**

##### **3.4.1 Split CP for Classification**

##### **3.4.2 Variance-Aware Score Normalization**

###### **a) Parametric linear shrink**

###### **b) Nonparametric normalization**

#### **3.5 CP Evaluation**

Set Composition, Spatial Equality, and Conditional Coverage

## Chapter 4

# RESULTS

### **4.1 Ensemble Results**

### **4.2 In-Distribution Evaluation**

### **4.3 Temporal OOD (2022)**

global coverage and set composition

### **4.4 Class Conditional Coverage**

### **4.5 Spatial Robustness within 2022**

### **4.6 Sensitivity**

## Chapter 5

### DISCUSSION

## Chapter 6

### LIMITATIONS

## Chapter 7

### CONCLUSION

#### **7.1 Computational Details**

## REFERENCES

- Angelopoulos, A. N., & Bates, S. (2022). *A gentle introduction to conformal prediction and distribution-free uncertainty quantification*. Retrieved from <https://arxiv.org/abs/2107.07511>
- Dewolf, N., De Baets, B., & Waegeman, W. (2025). Conditional validity of heteroskedastic conformal regression. *Information and Inference: A Journal of the IMA*, 14(2), iaaf013. <https://doi.org/10.1093/imaiiai/iaaf013>
- Guo, C., Pleiss, G., Sun, Y., & Weinberger, K. Q. (2017). On calibration of modern neural networks. In D. Precup & Y. W. Teh (Eds.), *Proceedings of the 34th international conference on machine learning* (pp. 1321–1330). PMLR. Retrieved from <https://proceedings.mlr.press/v70/guo17a.html>
- Kendall, A., & Gal, Y. (2017). What uncertainties do we need in bayesian deep learning for computer vision? In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, & R. Garnett (Eds.), *Advances in neural information processing systems* (Vol. 30). Curran Associates, Inc. Retrieved from [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/2650d6089a6d640c5e85b2b88265dc2b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/2650d6089a6d640c5e85b2b88265dc2b-Paper.pdf)
- Lakshminarayanan, B., Pritzel, A., & Blundell, C. (2017). Simple and scalable predictive uncertainty estimation using deep ensembles. *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 6405–6416. Red Hook, NY, USA: Curran Associates Inc.
- Tallam, K., Nguyen, N., Ventura, J., Fricker, A., Calhoun, S., O’Leary, J., ... Walter, R. K. (2023). Application of deep learning for classification of intertidal eelgrass from drone-acquired imagery. *Remote Sensing*, 15(9). <https://doi.org/10.3390/rs15092321>

## APPENDICES

### Appendix A

#### SUPER COOL FANCY FUNCTION