A Site is Not a Centroid

Modeling Archaeological Landforms

Kernel Regression on Focal Mean Embeddings

Matthew D. Harris, **AECOM**







Code and non-sensitive data used for this study are openly available online at:









*Special thanks to **Zoltán Szabó** (Center of Applied Mathematics, École Polytechnique, Paris, France) for inspiring this approach and kindly enduring a long sporadic email dialog. Thanks to Chester Cunanan (AECOM) for poster design and support. All errors, misrepresentations, and omissions are solely my own.

PROBLEM: MODELING THE RICHNESS OF LANDSCAPES

Landscape level archaeological sensitivity models are typically conceptualized as projecting the landform patterns observed at known sites to unsurveyed areas. However, most methods fail to characterize the richness of these landforms and the environmental variation is typically lost. Most commonly variation is lost when the entire landscape within a site is reduced to a single point at the center or summarized as to a mean for each variable measured. Alternatively, the variation is retained as independent observations per site area thereby violating the assumption of identically and independently distributed (i.i.d) observations due to spatial correlation.

While either of these two approaches can produce successful models, they beg the question: Can we model a representation of a landforms Richness without losing variation? More technically, the question is can we approximate a function to map distributions to presence/absence without assuming the shape of the distribution or costly density estimation? The answer presented here is to incorporate the methods of Distribution Regression (Szabó et al. 2016) along with focal Kernel Logistic Regression and focal window prediction.

TWO STAGE SAMPLING

PROBLEM STATEMENT

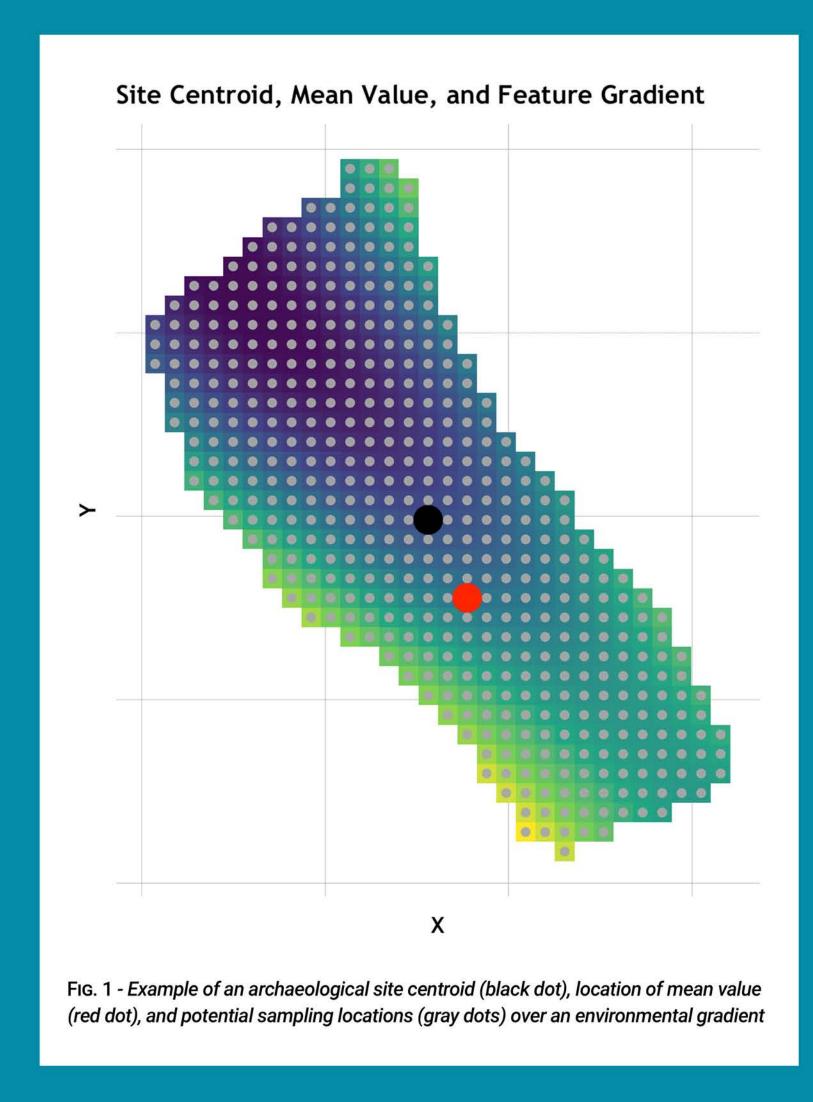
A two stage sampling strategy is the key to modeling the distribution of a feature (e.g. distance to water) as representing a single label; this is the Distribution Regression problem (Szabó et al. 2016). This approach assumes that all site locations are samples from some meta-distribution `M` and that each site is represented by a distribution of features and a label (e.g. site presence or absence). Since we do not measure the distribution directly, we take samples to approximate it.

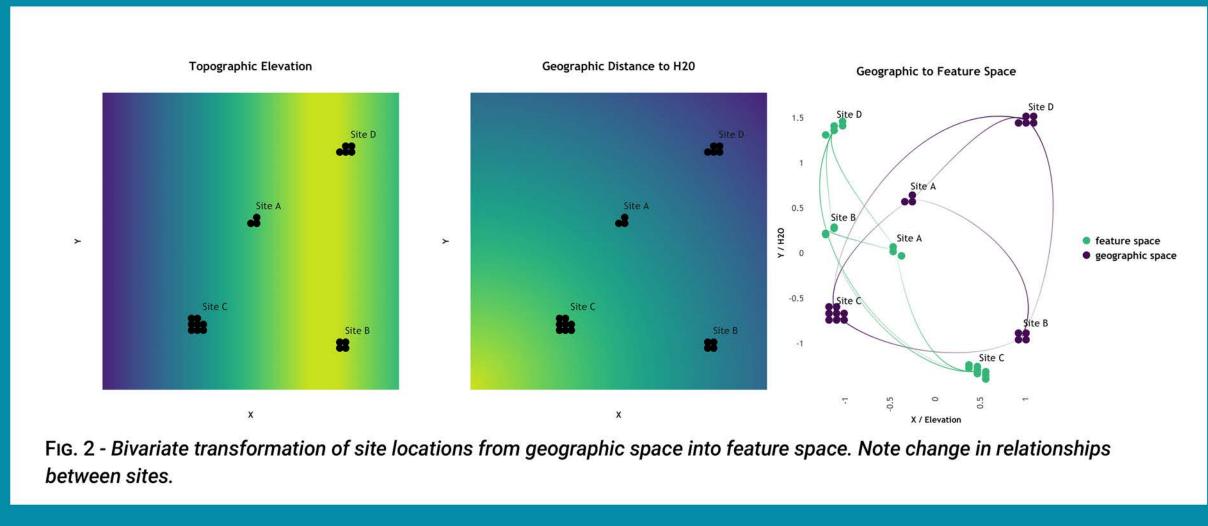
Thus, we take representative samples from the distribution on a site or background that are themselves samples from the meta-distribution of all possible sites and background (see Equation 1). Aggregating the intra-site observations to the inter-site level mitigates the issue of intra-site spatial correlation and moves the i.i.d assumption to the site level. This is the two-stage sampling approach. From this, the representative samples from within sites and backgroundcan be used to model the similarity between all sites and background.

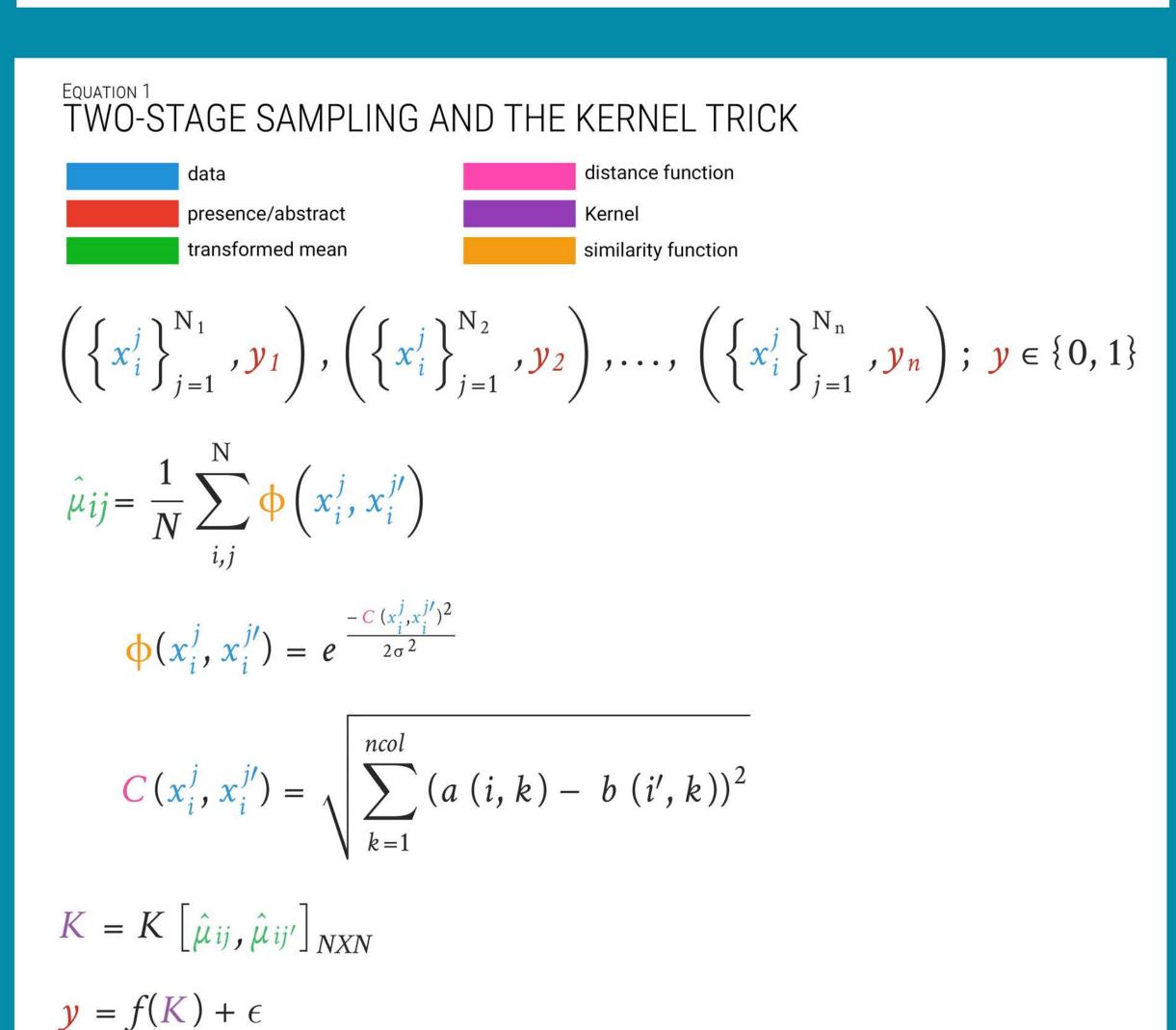
MODELING FEATURE SPACE

The principle concept to this approach is comparing landforms in their feature space as opposed to geographic space; this is called the "Kernel Trick" in Machine Learning. Once mathematically projected into feature space, higher dimensional nonlinear functions can be approximated without explicit (and expensive) mapping; that is the trick!

In high dimensional space, a linear model is able separate complex data. Once projected back to low dimensional space, the separation is non-linear. This property makes kernel methods convenient and computationally tractable.







MODEL VALIDATION AND COMPARISON

COMPARING KLRFOME TO OTHER MODELS

It is important to know how well this approach works when compared to other common or state-of-the-art models. Figure 5 is a comparison between the results of the KLR model against both Logistic Regression (LR) and a Support Vector Machine (SVM) ith Gaussian Kernel. These models are compared on two ifferent metrics (Area Under the Curve (AUC) and Youden's J), within five different physiographic zones, and over 100 different uns of the models on random subsets of the data. LR is used as a trustworthy and commonly known benchmark, while the SVM is used because it is the most well-known machine learning algorithm that also uses the kernel-trick.

In this comparison, KLR did equally as well based on the AUC and Youden's J metrics suggesting that this approach does not suffer in accuracy versus well-known models (Table 1). Interestingly, the KLR model does appreciably better than better than LR and SVM physiographic region 1, which is the most difficult to model.

		1	2	6	8	12
model	KLR	0.707	0.815	0.766	0.819	0.866
	LR	0.618	0.867	0.740	0.834	0.892
	SVM	0.609	0.850	0.745	0.811	0.868

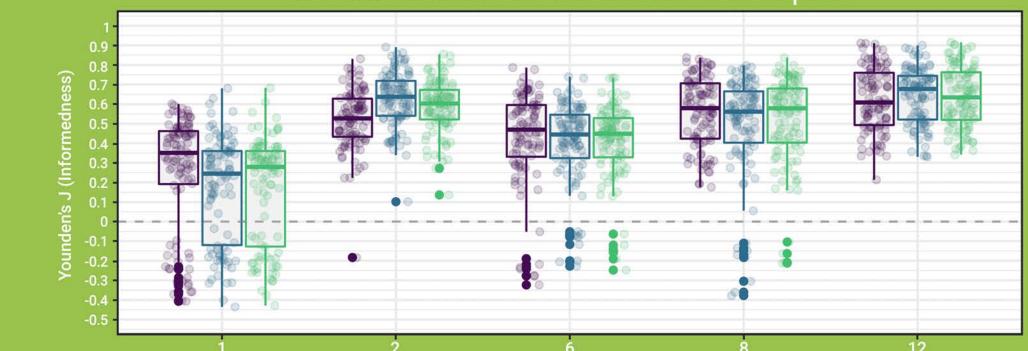
TESTING KLRFOME MODEL PERFORMANCE

Before the model is projected onto the study area, there are numerous internal tests of model performance and validity base on held-out test samples and the 100 random resamples shown n Figure 5. These tests give a very good approximation on how the model performs in a general sense. However, in use, most ensitivity models are thresholded into high, moderate, and low ensitivity and the true performance of a model depends not onl on these sensitivity classes, but also on a utility function defining cost" for errors. When defining appropriate threshold, it is very nformative to define the model across many possible thresholds and minimize or maximize the desired utility.

Median AUC Value for 100 Resamples



Median Youden's J value for 100 Resamples

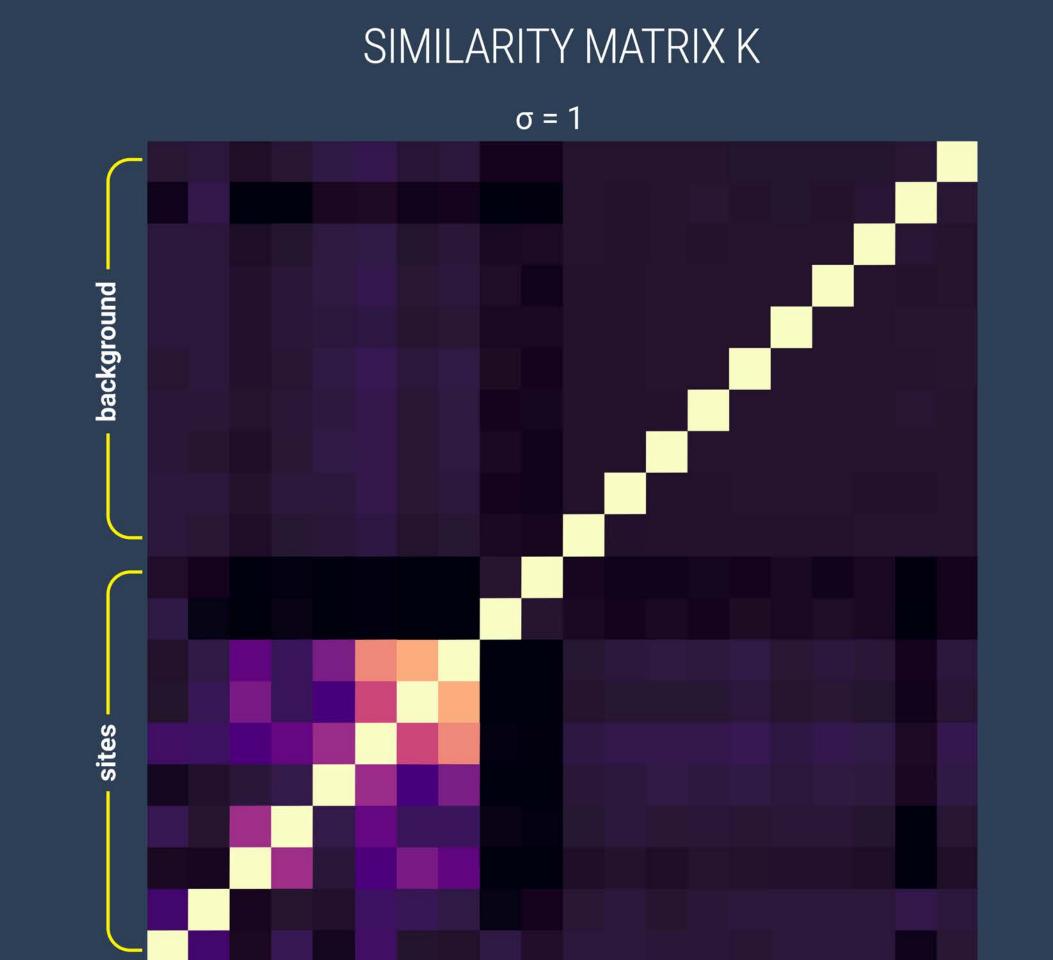


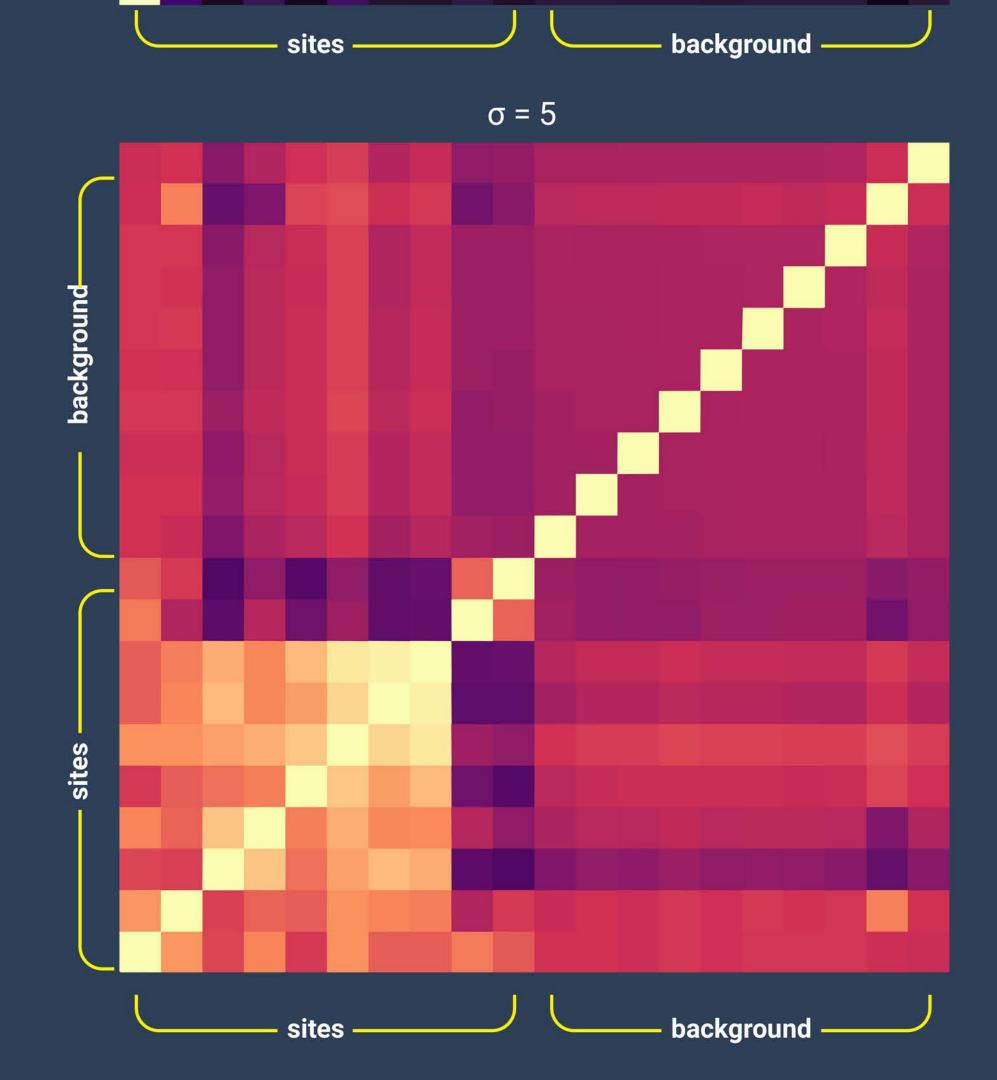
model type 📥 KLR 📥 LR 📥 SVM J = sensitivity + specificity - 1

Fig. 5 - Comparison of AUC and Youdens'J performance metrics for KLR, Logistic Regression (LR), and Support Vector Machine (SVM) models in different geographic settings. J = sensitivity + specificity - 1

Table 2 illustrates the true/false positive and true/false negative ed quantities, as well as relevant metrics for each potential probability threshold. If maximizing for the Youden's J statistic, the optimal threshold between site-present and site-absent classes is around 0.4. At this threshold, the KLR model given this area and parameters correctly classified 91% of the site-present cells into an area covering 53% of the study area. In practice, a finer sequence of thresholds would be calculated, and a weighted utility function would be the objective.

2 KERNEL LOGISTIC REGRESSION (KLR) ON FEATURE SPACE MEAN EMBEDDING





SIMILARITY KERNEL

Unlike most statistical learning approaches, this approach does not model the individual observations of presence and absence. Instead, via the two-stage sampling approach and kernel trick, it models the similarity between all sites and the background environment. This similarity is based on many samples of the characteristics that define the landscape of each site and the general environment in which they are found.

By modeling similarity, we are explicitly trying to project the same rich landscapes on which we have found sites to new areas of similar richness? Depending on the research question, specific definitions of "similarity" or "distance" can be used to model the relationship between sites. This study uses the Exponentiated Quadratic (a.k.a Gaussian kernel) to define similarity. With the use of a similarity matrix as the objective, as opposed to simple tabular data, a Kernel Logistic Regression (KLR) approach (Zhu and Hastie 2005) is used.

LEARNING PARAMETERS WITH NEWTON'S METHOD (IRLS)

In order to project the similarity of known sites into areas that have not been surveyed we must learn a set of functions to fit our training data in the similarity kernel. Statistical and Machine Learning are terms for the process of learning functions from existing data in order to apply them to new data. In this research, we use the Iterative Reweighted Least Squares (IRLS) (aka Newton's Method) to fit the model parameters.

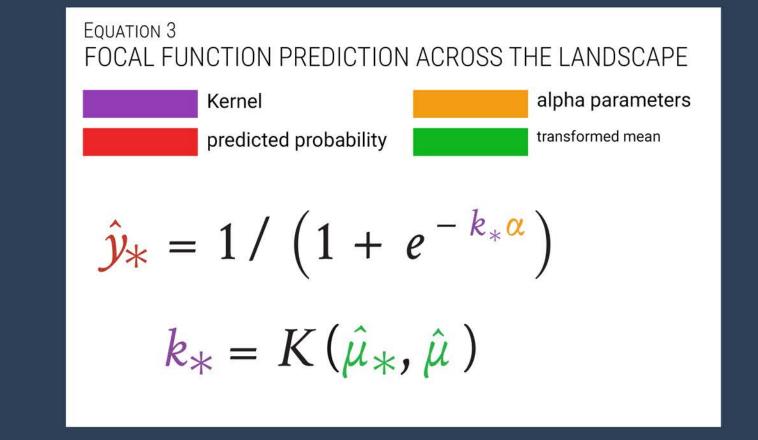
ITERATIVELY REWEIGHTED LEAST SQUARES (IRLS) $p = 1 / \left(1 + e^{-K\alpha} \right)$ alpha parameters $\alpha = [1/N]$ probability function estimated response $z = K\alpha + (y - p)/W$ $W = diag \left(\frac{p}{p} (1 - \frac{p}{p}) \right)_{[N]}$ $\alpha^{new} = (K + \lambda I)^{-1} z$ $\hat{f}(K) = \arg\min \left[\sum_{\alpha} \left(\frac{\alpha^{new}}{\alpha^{new}} - \alpha \right) + \lambda \left\| f \right\|_{H_K}^2 \right]$ $\hat{\mathbf{y}} = 1 / \left(1 + e^{-K\alpha} \right)^{new}$

Generally, the IRLS approach iteratively minimizes the difference between old and new parameters regularized by the lambda hyperaparemeter. The result of IRLS is a set of parameters that can be used to compare the similarity of any new location to those that we already know and produce a probability score of whether is should be in the site-present or background class.

PROJECTION OF SIMILARITY ACROSS THE LANDSCAPE

Prediction involves the calculation of a new similarity kernel between the new test data and all of the training data. The new prediction kernel is then multiplied by the learned parameters to derive a probability of site-presence. In this research, prediction is carried out as a roving focal window function across the study area.

This means instead of predicting based on the information contained in single raster cell, the prediction is based on a N by N neighborhood around each cell; in this case 3 by 3 cells. In this sense, the prediction is that of a landscapes similarity and the view of this landscape is controlled by the size of the N by N window. The larger the focal window, the greater the perspective of the landscape and the smoother



3) FOCAL PREDICTION

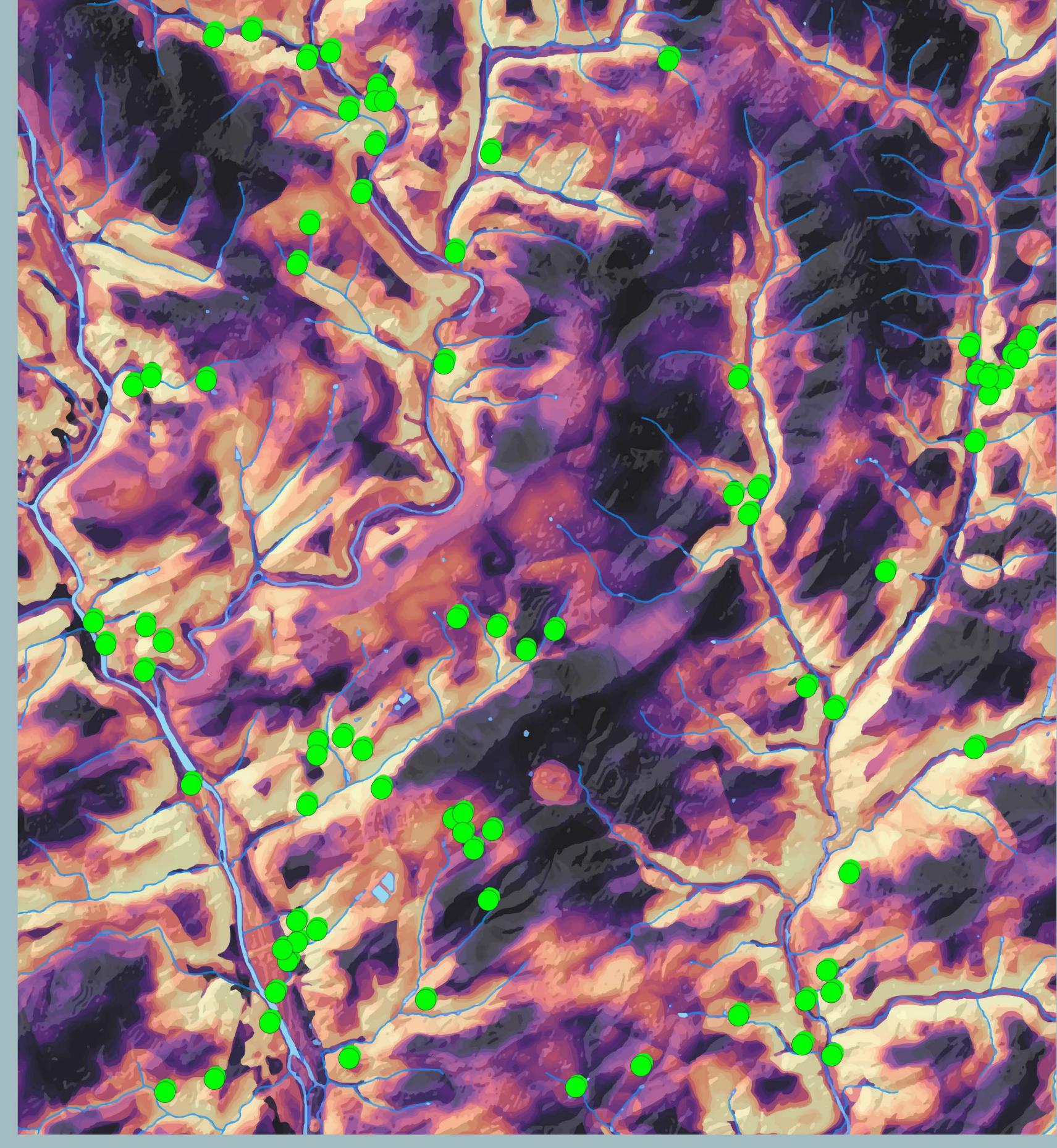
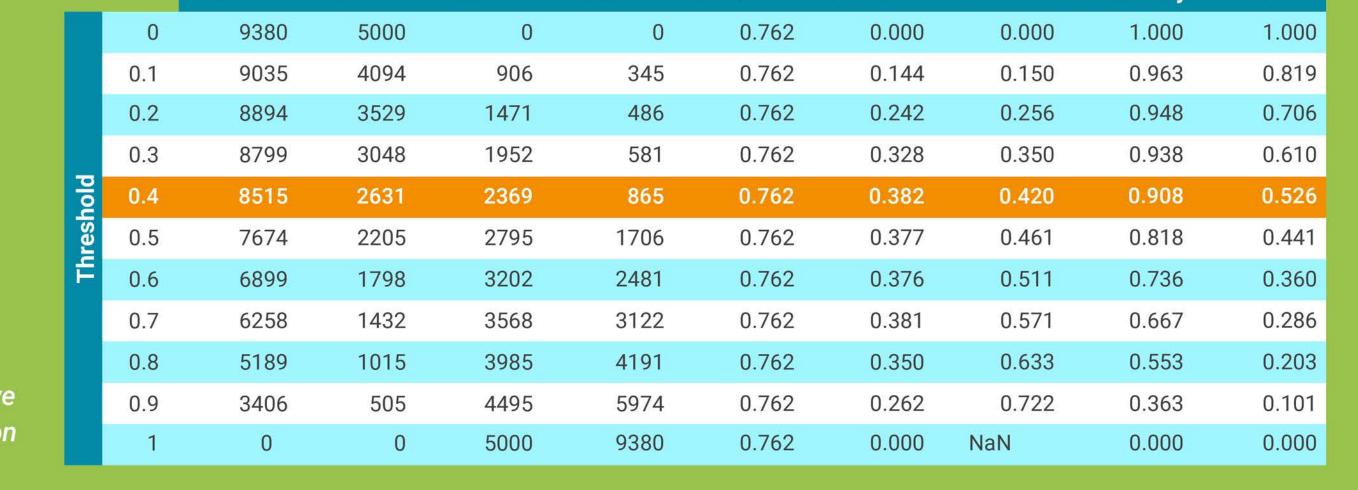


Fig. 4 - Projection of KLRfome model onto geographic space using roving focal window with 3 by 3-cell

QUANTIFY UNCERTAINTY

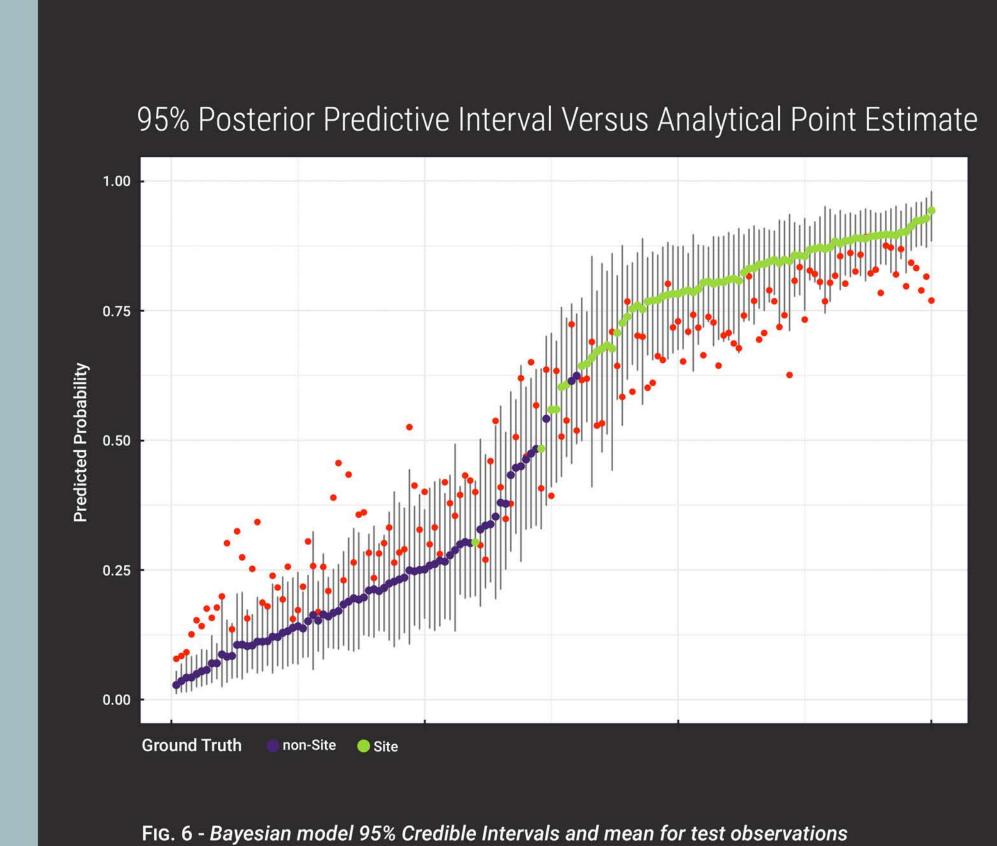


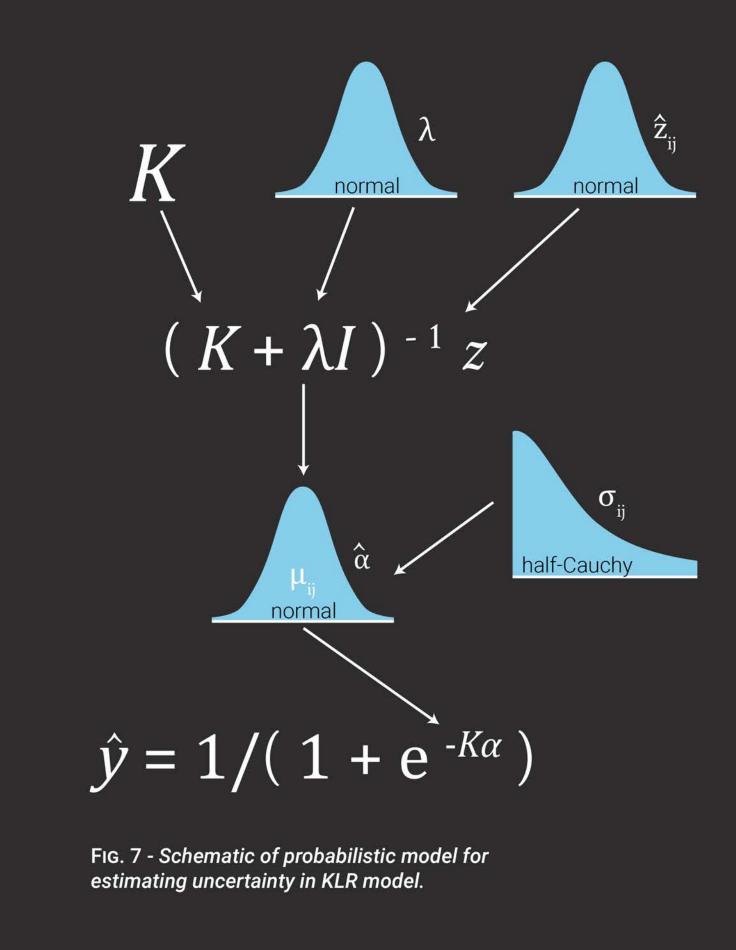
5 BAYESIAN INTERPRETATION OF KLR

BAYESIAN INTERPRETATION TO

Beyond creating and characterizing a model of explicit similarity modeling, an objective of this study is to recast the model into a Bayesian framework to better propagate uncertainty into the posterior predictive distribution. In this approach, uncertainty in model parameters, location measurements, and sample sizes can be approximated with distributions on order to better understand the implications of what we do not know.

Presented here is the initial framework of a Bayesian KLR and an illustration of the predictive distribution. Coded in the Stan probabilistic Programming Language, this model places a range of many possible values over the lambda, alpha, and sigma parameters from the KLR model. The Bayesian model is still a work in progress.





Flaxman, Seth, Yu-Xiang Wang, and Data Mining, pages 289-298. ACM, 2015.

Law, BHo Chung Leon, Dougal J. Sutherland, Dino Sejdinovic, and Seth Flaxman. Bayesian Approaches to Distribution Regression, 2017. https://arxiv.org/abs/1705.04293 [accessed 03/21/2018]

Sciaini, Marco, Matthias Fritsch, Craig E. Simpkins (2018). {NLMR}: R package version 0.2.0. URL

Magazine, 30(4):98-111, 2013.

Szabó, Zoltán, Arthur Gretton, Barnabás Póczos, and Bharath Sriperumbudur. Two-stage sampled learning theory on distributions. In International Conference on Artificial Intelligence and Statistics (AISTATS), pages 948-957, San Diego, California, USA, 9-12 May 2015.

Szabó, Zoltán, Bharath Sriperumbudur, Barnabás Póczos, and Arthur Gretton. Learning theory for distribution regression. Journal of Machine Learning Research, 17(152):1-40, 2016.

Zhu, Ji and Trevor Hastie. Kernel logistic regression and the import-vector machine. Journal of Computational and Graphical Statistics, 14(1):185-205, 2005.

REFERENCES

Alexander J Smola. Who Supported Obama in 2012?: Ecological inference through distribution regression. In Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and

Simulating neutral landscape models. https://CRAN.R-project.org/package= Song, Li, Kenji Fukumizu and Arthur Gretton. Kernel Embeddings of Conditional Distributions: A Unified Kernel Framework for Nonparametric Inference in Graphical Models, in *IEEE Signal Processing*