
Using Machine Learning Methods to Classify Urban Heat Islands

Jennifer Xu and Susan Chen
Stanford University
jennxu23@stanford.edu; peeka5@stanford.edu

1 Introduction

2 1.1 Background

3 In past decades, global warming has had a significant effect on society. More extreme weather events,
4 such as the intense forest fires in California or the unusually cold weather in Texas, has illustrated the
5 urgency of stopping climate change. Increasing urbanization, which involves both the movement of
6 people into urban areas as well as construction of urban infrastructure, has contributed to climate
7 change. The term "urban heat island" refers to a metropolitan area that is significantly warmer than
8 its surrounding rural counterparts.

9 Specific characteristics of urban areas that contribute to the urban heat island effect. For example,
10 urban areas are more dense than rural areas. People, buildings, and vehicles are packed together and
11 constantly burning energy, generating heat. Man-made structures, such as buildings and parking
12 garages, reflect light more than natural structures eventually warming up the surrounding air. As
13 a result, urban areas also tend to be significantly warmer at night compared their rural neighbors
14 (National Geographic 2012). As cities develop and more people move into urban areas, the urban
15 heat island effect could have dire consequences. Additionally, a higher population density as well as
16 warmer temperatures lead to resource strain within urban heat islands.

17 For our project, we decided to study the influence of population, land surface temperature, land usage,
18 and reflectivity as possible contributors to the urban heat island effect. Using relevant data sets and
19 applicable machine learning algorithms, our goal for this project is to build a model that can classify
20 a given city as an urban heat island.

21 1.2 Related Work

22 Existing research has been done on specific urban regions. In a recent study, researchers used
23 multi-temporal Landsat satellite images to simulate land cover and temperature changes overtime in
24 Bangladesh (Kafy et. al 2021). From the LandSat data, Kafy used a Support Vector Machine (SVM)
25 model to simulate urban development and changes in land surface temperature. Another study used
26 machine learning methods to determine which factors contribute most to the development of urban
27 heat islands (Yoo, 2018). Using data from the cities of Marion County, Indiana, Yoo used Principal
28 Component Analysis (PCA) to determine that imperviousness, vegetation, and building footprint
29 were significant contributors to the Urban Heat Island Effect. A similar study was done on a larger
30 scale with urban areas in China, using a linear regression model to illustrate the trend between urban
31 development and increases in surface temperature (Zhou et. al, 2015)

32 We decided to take an approach similar to the above studies and build an urban heat island classifier
33 based on environmental data. Similar to Yoo's study, we also analyze the relative influences of
34 particular environmental factors.

35 2 Methodology

36 2.1 Data set and Feature Extraction

37 The data for this project is sourced from three different data sets: The Global Urban Heat Island
38 (UHI) data set v1 (2013) from the Socioeconomic Data and Applications Center (SEDAC), the
39 LandSat 8 data set from the United States Geological Survey (USGS), and the Moderate Resolution
40 Imag-154ing Spectroradiometer (MODIS) Land Cover Type (MCD12Q1) Version 6 data set from the
41 USGS.

42 Each entry in the Global UHI data set represented a specific urban area. In addition to identifying
43 information, the Global UHI dataset included population data from 1990, 1995, and 2000; area
44 of urban extent in square kilometers, daytime average temperature data, and nighttime average
45 temperature data. To develop the training labels, we used the absolute difference between the daytime
46 urban average temperature and the daytime buffer (average temperature of an urban area and a 10km
47 buffer) average temperature. If a given area had greater than a 1° difference in daytime average
48 temperature as compared to its buffer zone, it was classified as an urban heat island (NASA, 2013).

49 Using Google Earth Engine, we located 2013 LandSat images and composed an extraction script to
50 obtain bands 4, 5, and 6 data corresponding to each city. Bands 4, 5 and 6 represent the reflectance
51 from short-wave infrared (1.55 - 1.75 μm), near-infrared (0.77 - 0.90 μm), and visible (0.63 - 0.69
52 μm) light respectively. For the land usage data, we used the National Land Cover Database (NLCD).
53 A city normally has multiple land cover types. Each land cover type is represented as a number.
54 Dominant land cover types were selected for analysis. To analyze the impact of each dominant land
55 cover type, we used hot-encoding. For a given city and a given land cover type, a 1 indicated that the
56 land cover type was present in the city and a 0 indicated otherwise.

57 The extracted data from the LandSat and NLCD data sets were appended onto the the Global UHI
58 data set. Due to the geographical and temporal restrictions of the three original data sets, we limited
59 the scope of our final data set to be urban areas in the United States as measured in 2013. Our final
60 data set included the original features of the Global UHI data set (excluding identifying information
61 and the daytime average difference); band 4,5, and 6 reflectance data; and presence of land cover
62 type 82, 90, 41, 52, 81, 71, 42, 11, 95, 43, 21, 23, 22, 24, and 31. The final data set has a total of 28
63 features.

64 Finally, we randomly split our data set into a training set, a testing set, and a validation set based on a
65 70:15:15 ratio. The training set contains 3500 samples, the testing set contains 750 samples, and the
66 validation set contains 749 samples. Each sample represents one metropolitan area.

67 2.2 Models

68 Because our ultimate goal is to build a classifier for urban heat islands, and because we have
69 categorical features, we decided to use three basic classification algorithms: logistic regression,
70 Gaussian discriminant analysis, and a support vector machine (SVM) model.

71 **Logistic Regression** Our logistic regression model directly predicts whether or not a city is an
72 urban heat island based on a probabilistic analysis. The basis of a logistic regression model is the
73 Sigmoid function, which maps our real valued predictions to a value between 0 and 1. We can write
74 the probability of a given data point as

$$p(y|x; \theta) = h_{\theta}(x)^y (1 - h_{\theta}(x))^{1-y} \quad (1)$$

75 where x is a given training example, y is the associated label, and $h_{\theta}(x)$ is the sigmoid function.

76 From the above equation, we can derive the weights of each feature by gradient descent, where the
77 log-likelihood function is

$$\ell(\theta) = \sum_{i=1}^n y^{(i)} \log h(x^{(i)}) + (1 - y^{(i)}) \log(1 - h(x^{(i)})) \quad (2)$$

78 In the above equation, $x^{(i)}$ is a training example, $y^{(i)}$ is its associated label, and $h(x^{(i)})$ is the sigmoid
79 function. We used SciKit-Learn's Logistic Regression class to generate our logistic regression model.

80 **Gaussian Discriminant Analysis** We also developed a Gaussian discriminant analysis model
 81 (GDA). Since the features we collected are environmental data, we can possibly fit a Gaussian
 82 Distribution to our model. GDA is predicated on the assumption that the probability of the training
 83 set given the labels ($p(x|y)$) can be modeled as a multivariate Gaussian distribution. From the law of
 84 total probability and the underlying Gaussian assumption, we can derive and maximize the following
 85 likelihood function:

$$\ell(\phi, \mu_0, \mu_1, \Sigma) = \log \prod_{i=1}^n p(x^{(i)}|y^{(i)}; \phi, \mu_0, \mu_1, \Sigma)p(y^{(i)}; \phi, \mu_0, \mu_1, \Sigma) \quad (3)$$

86 where ϕ is the probability that a training example is an urban heat island ($p(y = 1)$), μ_0 and μ_1
 87 are the expected values of the negative (non-urban heat islands) and positive (urban heat islands)
 88 classes respectively, and Σ is the covariance matrix for both classes. We used Scikit-Learn’s Linear
 89 Discriminant class to generate our Gaussian Discriminant model. We also used Scikit-Learn’s
 90 Gaussian Discriminant class to explore the possibility of a quadratic decision boundary.

91 **Support Vector Machine** We also developed an SVM model based on the data set. The incentive
 92 to use SVM is to find a separating hyperplane that maximizes the total distance from the data points.
 93 Our project can be treated as a linear classifier for a binary classification problem. Therefore, we use
 94 weights w and intercept b , and write our classifier as

$$h_{w,b}(x) = g(w^T x + b) \quad (4)$$

95 The functional margin of (w, b) with respect to the training example can be written as:

$$\gamma^{(i)} = y^{(i)}(w^T x^{(i)} + b) \quad (5)$$

96 The SVM model minimized the sum of functional margins for the training set. We tried implementing
 97 both a polynomial kernel with degree 3 and radial basis function kernel to explore different model
 98 complexities.

99 **Feature Analysis** As a part of our analysis, we wanted to see if feature selection would lead to
 100 an increase in model performance. We used a Random Forest Classifier to select features. Random
 101 decision forests are an ensemble learning method for classification. The model operates by construct-
 102 ing a multitude of decision trees at training time and outputting the class that is the mode of the
 103 classes (classification) (Ho, 1995). Since we have multiple features belong to different data types,
 104 random forests serve as an effective tool to rank the importance of features and narrow down to the
 105 features that matters. Our Random Forest model ranks features based on Mean Decrease in Impurity
 106 (MDI), which measures the average decrease in impurity based on the Gini impurity across all the
 107 trees (Louppe et. al, 2013).

108 3 Results

109 For our preliminary experiments, we built a Logistic Regression Model, a GDA model, and a SVM
 110 model on the entire data set. To explore hyper parameters and evaluate model fit, we utilized SciKit-
 111 Learn’s GridSearch CV algorithm to exhaustively test and determine optimal hyper parameters for
 112 each model. We evaluated the effectiveness of our model based on the accuracy and the ROC-AUC
 113 score (generated from SciKit-Learn’s metrics package). Error was calculated by taking the norm of
 114 the differences between the predicted and true labels.

115 Next, we explored the possibility of feature selection to both improve our models and to determine
 116 which features have the greatest impacts on the models. We used a Random Forest model to select the
 117 most impactful features. Then, we re-built each model and calculated optimal hyper parameters based
 118 on the reduced-dimensionality data set. Finally, we used SHAP to analyze the impact of individual
 119 features on individual models.

120 **Model Results** Before hyperparameter tuning, the Logistic Regression model had an accuracy of
 121 0.536 with an error of 18.65. Using GridSearchCV, combinations of different degrees of the inverse
 122 of regularization ($\frac{1}{\lambda}$ and the penalty term $l1$ and $l2$) were tested over 5 folds. The optimal hyper
 123 parameter combination was determined to be the $l2$ penalty term with an inverse of regularization of
 124 10^{-2} . Optimizing for the best accuracy, the optimal model had an increased accuracy of 0.644 and a
 125 reduced error of 15.71.

126 Similarly, the Linear Discriminant Analysis model had an accuracy of 0.524 with an error of 18.894.
 127 In order to improve the accuracy of the Linear Discriminant analysis, we tried two different solvers:
 128 Singular value Decomposition (SVD) and Least Squares (LSQR). Using GridSearchCV over 5 folds,
 129 the optimal solver was determined to be SVD. With SVD, the accuracy of the Linear Discriminant
 130 Analysis was determined to be 0.620 with a reduced error of 16.248.

131 Finally, the SVM model has an accuracy of 0.536 and an error of 18.65. For our initial SVM model,
 132 we chose a radial basis function kernel. After using GridSearchCV and exhausting all combinations
 133 of gamma and regularization strength vaues ($\frac{1}{\lambda}$), the accuracy of SVM model increased to 0.647 and
 134 error reduced to 16.03. The optimum hyper parameters for accuracy were an inverse regularization
 135 value of 10^{-4} and a gamma value of 10^{-3} .

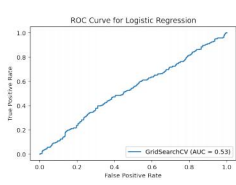


Figure 1: ROC Curve for Logistic Regression

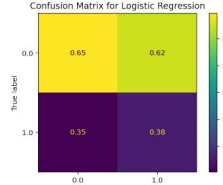


Figure 2: Normalized Confusion Matrix, Logistic Regression

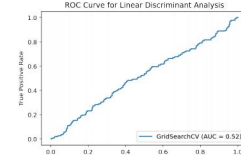


Figure 3: ROC Curve for LDA

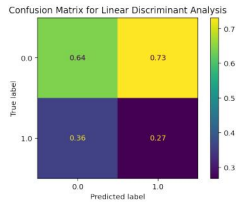


Figure 4: Normalized Confusion Matrix, LDA

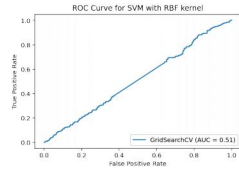


Figure 5: ROC Curve for SVM

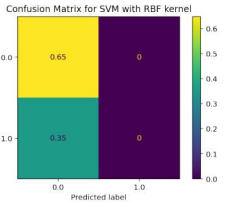


Figure 6: Normalized Confusion Matrix, SVM

136 Based on the ROC-AUC cures, all three models perform about the same as a random classifier. Closer
 137 analysis via the confusion matrices show that all three models correctly classify negative examples
 138 (non urban heat islands) a majority of the time, but poorly classifies positive examples (urban heat
 139 islands).

140 **Feature Analysis** We explored feature reduction as a way to possibly improve our models and
 141 hyper parameters. Using a Random Forest model, we found that land cover data doesn't significantly
 142 impact the model. Then, we retrained our models on a reduced data set including only population,
 143 temperature, and reflectivity data. The reduced data set did not improve the fit or accuracy of the
 144 models. However, analyzing the features that were most influential on the models reveal interesting
 145 trends.

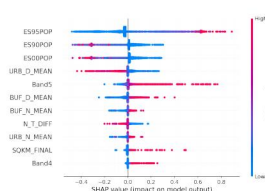


Figure 7: SHAP, LDA

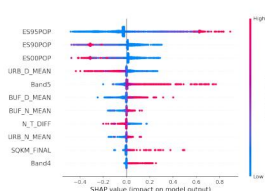


Figure 8: SHAP, SVM

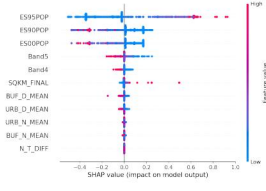


Figure 9: SHAP, Logistic Regression

146 We used the SHAP python library for model explainability. SHAP (SHapley Additive exPlanation)
 147 leverages the idea of Shapley values for model feature influence scoring. A Shapley value represents
 148 the average marginal contribution of a given feature across of all features. This exhaustive approach
 149 would guarantee SHAP's consistency and local accuracy.

150 4 Discussion

151 Our initial models performed poorly, with an accuracy within the .50 - .60 range. With hyperparameter
152 tuning, our models' accuracies improved into the .60-.70 range with lower error. There is not a
153 significant difference in accuracy between the models. As shown by the ROC-AUC curves, all three
154 models perform about as accurately as a random classifier. Trying more complex models, feature
155 reduction, and hyperparameter tuning did not improve our models' performances. The confusion
156 matrices indicate that all three models perform decently at classifying negative examples (non UHIs),
157 but do not accurately classify positive examples (UHIs).

158 However, SHAP analyses reveals useful information about the impact of particular features. We see
159 from the random forest analyses that land cover data is not particularly impactful. The SHAP analyses
160 of all three models show population data as highly impactful. Specifically, a larger population in 1995
161 seems to be correlated with a greater chance of being classified as an urban heat island [ES95POP].
162 The opposite relationship is true for the population in 1990 and 2000 [ES90POP and ES00POP].
163 Population data seems to be particularly impactful in the Linear Regression model, as it accounts
164 for most of the variance in the model. The impact and positive correlation is expected, as a higher
165 population within an urban area would mean more movement. As a result, more energy and heat is
166 generated.

167 Regarding temperature, we see that, across all models, daytime temperature metrics [URB_D_MEAN
168 and BUF_D_MEAN] seem to have a greater impact on models than nighttime daytime temperature
169 metrics [URB_N_MEAN, BUF_N_MEAN and N_T_DIFF]. This may be due to the fact that there's
170 a greater difference in activity, and therefore generated energy and heat, between urban areas and
171 their surroundings during the day as compared to the night.

172 Band 5 data, which corresponds with the intensity of reflectivity of Short-wave Infrared (1.55 - 1.75
173 μm) light, was the most impactful reflective feature wavelength as compared to Near-Infrared (Band 4)
174 or visible light (Band 6). Interestingly, high Band 5 values were positively correlated with a positive
175 urban heat island classification in the SVM and LDA models, while it was negatively correlated in
176 the Logistic Regression model. Previous studies indicated that reflectivity is positively associated
177 with the urban heat islands. Urban heat islands often have artificial materials with high reflectance
178 such as buildings and pavements, that can lead to an increase in temperature.

179 5 Conclusion

180 **Main takeaway** In summary, we attempted to build an urban heat island classifier using three
181 models: Logistic Regression, Gaussian Discriminant Analysis, and Support Vector Machine. We
182 build our models based on remote sensing data. While none of our models were particularly good
183 classifiers of urban heat islands, feature analysis revealed interesting and novel conclusions.

184 **Future Directions** As shown by our previous work, we tested for different models with and without
185 feature selection. Based on our results, more features should be integrated into this model. Examples
186 of such parameters include a different population dataset (i.e population density, a more recent
187 population data set, etc.), reflectivity of light from other wavelength in the spectrum, economic
188 parameters, vegetation indices, or calculated measures of urbanization. This would increase model
189 complexity and potentially lead to more accurate models and a more informative feature reduction.

190 Additionally, more research should be done to integrate existing features in a way that reduces
191 colinearity and better reflects their impact on given models. For example, different reflectivity bands
192 can be mathematically combined and used to indicate different land cover types.

193 Furthermore, our model indicates that hot-encoded dominant land cover do not greatly impact
194 classification. Our models could possibly be improved by using a larger data set in which cities are
195 classified as urban and rural, rather than being separated by specific land cover types. This could
196 generate clearer data clustering.

197 Lastly, remote sensing data is often generated and transformed by complicated algorithms, whose
198 complexities could be difficult to capture with a linear model. This could've contributed to the high
199 bias in our models. We could try different, possibly more complex, methods of analysis such as deep
200 learning or unsupervised training.

201 **6 Code and Contributions**

202 Both Susan and Jennifer worked equally on the project. Susan extracted and cleaned the reflectance
203 and land use data. Jennifer cleaned the Global UHI dataset. Both members worked together to
204 determine and fit the best classifiers, as well as do the feature analysis. Jennifer focused on the
205 logistic regression and GDA model. Susan worked on the SVM and Random Forest model. Both
206 members did hyperparameter tuning and feature explanation together.

207 The code can be found at the following link: [https://colab.research.google.com/drive/
208 1u7sGnJNoVKsp_oAr3edPJsUn8eprpSMU?usp=sharing](https://colab.research.google.com/drive/1u7sGnJNoVKsp_oAr3edPJsUn8eprpSMU?usp=sharing)

209 **References**

210 [1] Ho, Tin Kam (1995). Random Decision Forests. Proceedings of the 3rd International Conference on
211 Document Analysis and Recognition, Montreal, QC, 14–16 August 1995. pp. 278–282.

212 [2] Kafy, A. A., Shuvo, R. M., Naim, M. N. H., Sikdar, M. S., Chowdhury, R. R., Islam, M. A., ... Kona, M. A.
213 (2021). *Remote sensing approach to simulate the land use/land cover and seasonal land surface temperature
214 change using machine learning algorithms in a fastest-growing megacity of Bangladesh*. Remote Sensing
215 Applications: Society and Environment, 21, 100463. [3] Louppe, G., Wehenkel L., Suter A., Geurts P.
216 (2013). *Understanding variable importances in forests of randomized trees*. NIPS'13: Proceedings of the 26th
217 International Conference on Neural Information Processing Systems - Volume 1. December 2013. pp. 431-439

218 [4] National Aeronautics and Space Association. (2000). Population Density [Data File]. Retrieved from
219 https://neo.sci.gsfc.nasa.gov/view.php?datasetId=SEDAC_POP.

220 [5] National Aeronautics and Space Association. (2013). Global Urban Heat Island (UHI) Data Set,
221 v1 [Data File]. Retrieved from [https://sedac.ciesin.columbia.edu/data/set/sdei-global-uhi-
222 2013/data-download](https://sedac.ciesin.columbia.edu/data/set/sdei-global-uhi-2013/data-download).

223 [6] National Geographic Society. (2012, October 09). *Urban heat island*. Retrieved from [https://www.
224 nationalgeographic.org/encyclopedia/urban-heat-island/](https://www.nationalgeographic.org/encyclopedia/urban-heat-island/).

225 [7] United States Geological Survey. (2013). Landsat 8 Surface Reflectance Tier 1 [Data File]. Retrieved
226 from [https://developers.google.com/earth-engine/datasets/catalog/LANDSAT_LC08_C01_T1_
227 SR#bands](https://developers.google.com/earth-engine/datasets/catalog/LANDSAT_LC08_C01_T1_SR#bands).

228 [8] United States Geological Survey. (2019). Terra and Aqua combined Moderate Resolution Imag-
229 ing Spectroradiometer (MODIS) Land Cover Type (MCD12Q1) Version 6 [Data File]. Retrieved from
230 <https://lpdaac.usgs.gov/products/mcd12q1v006/>.

231 [9] Yoo, S. (2018). *Investigating important urban characteristics in the formation of urban heat islands: a
232 machine learning approach*. Journal of Big Data, 5(1), 1-24.

233 [10] Zhou, Decheng; Zhang, Liangxia; Hao, Lu; Sun, Ge; Liu, Yongqiang; Zhu, Chao. (2016). *Spatiotemporal
234 trends of urban heat island effect along the urban development intensity gradient in China*. Science of The Total
235 Environment, Vol. 544: 10 pages.: 617-626. doi:10.1016/j.scitotenv.2015.11.168