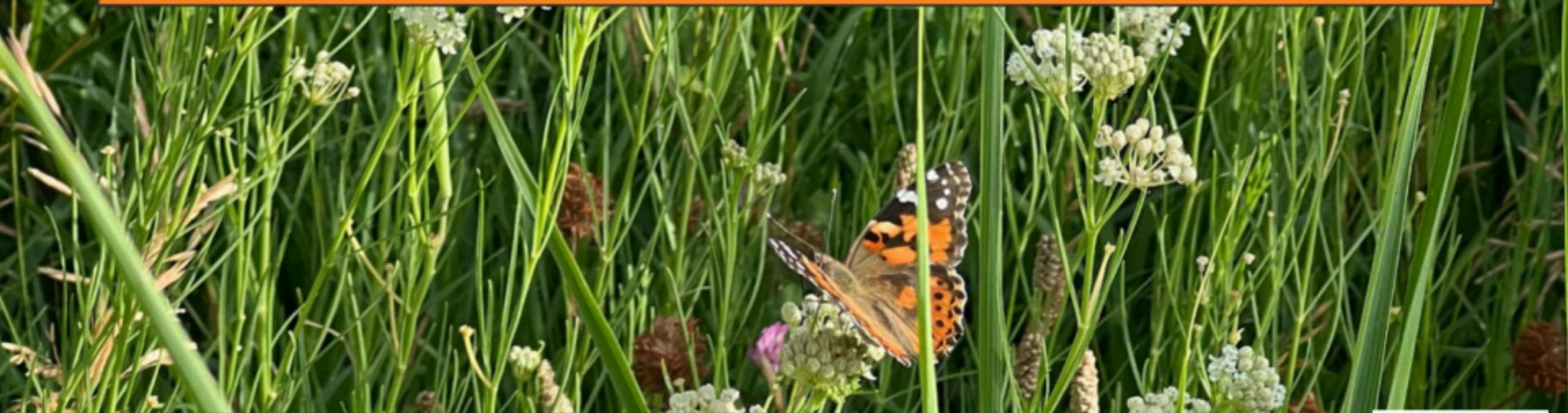


# Biodiversity through the lens of Keystone Species in Albuquerque



CNM Deep Dive Data Science Bootcamp  
Cohort 10 Capstone Project  
9/15/23



# Problem Definition

Citizen science may be an under recognised data source. With tree-based and linear regression modeling we are hoping to show how citizen science can be used to understand the relationship between soil type, land cover classification and the predictive factors for keystone species in the urban environment of Albuquerque, New Mexico.

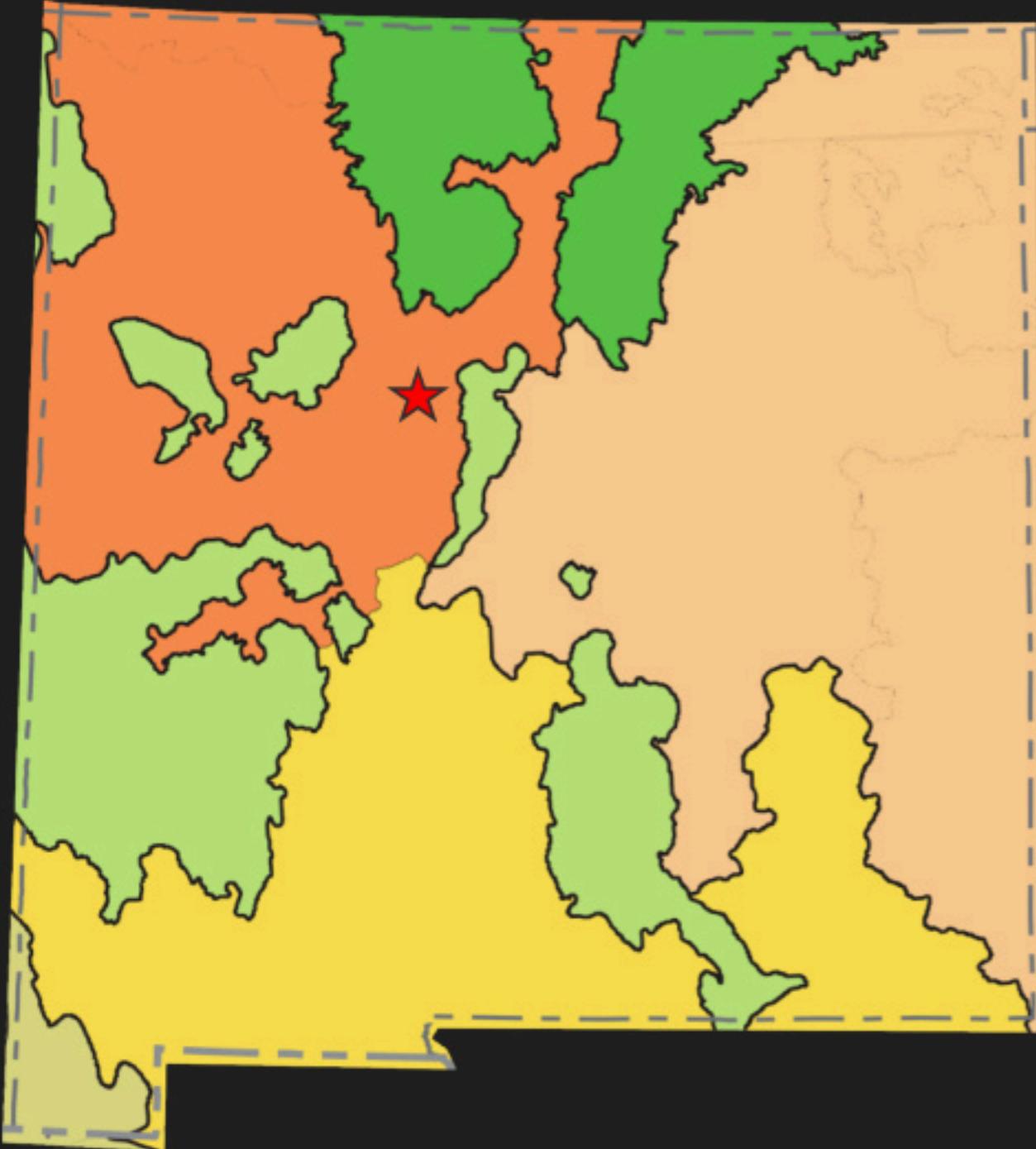


# Concept of Keystone Species



- Robert Paine experiment
  - Early 1960's
  - Conducted at Mukkaw Bay, WA
  - Similar observations made by other Ecologists
- Organisms that keep ecosystem in equilibrium

# Ecoregions of New Mexico



Arizona / New Mexico  
Plateau

Chihuahuan Desert

Southern Rockies

Great Plains

Arizona / New Mexico  
Mountains

Madrean Archipelago

# Keystone Species - Arizona / NM Plateau

- Predator
  - Coyote
- Prey
  - Cottontail Rabbit
- Ecosystem Engineer
  - Beaver
  - Prairie Dog
- Mutualist
  - Hummingbird
  - Bee
- Plants
  - Coyote Willow
  - Chamisa
  - Cottonwood



# Data Sources and Infrastructure



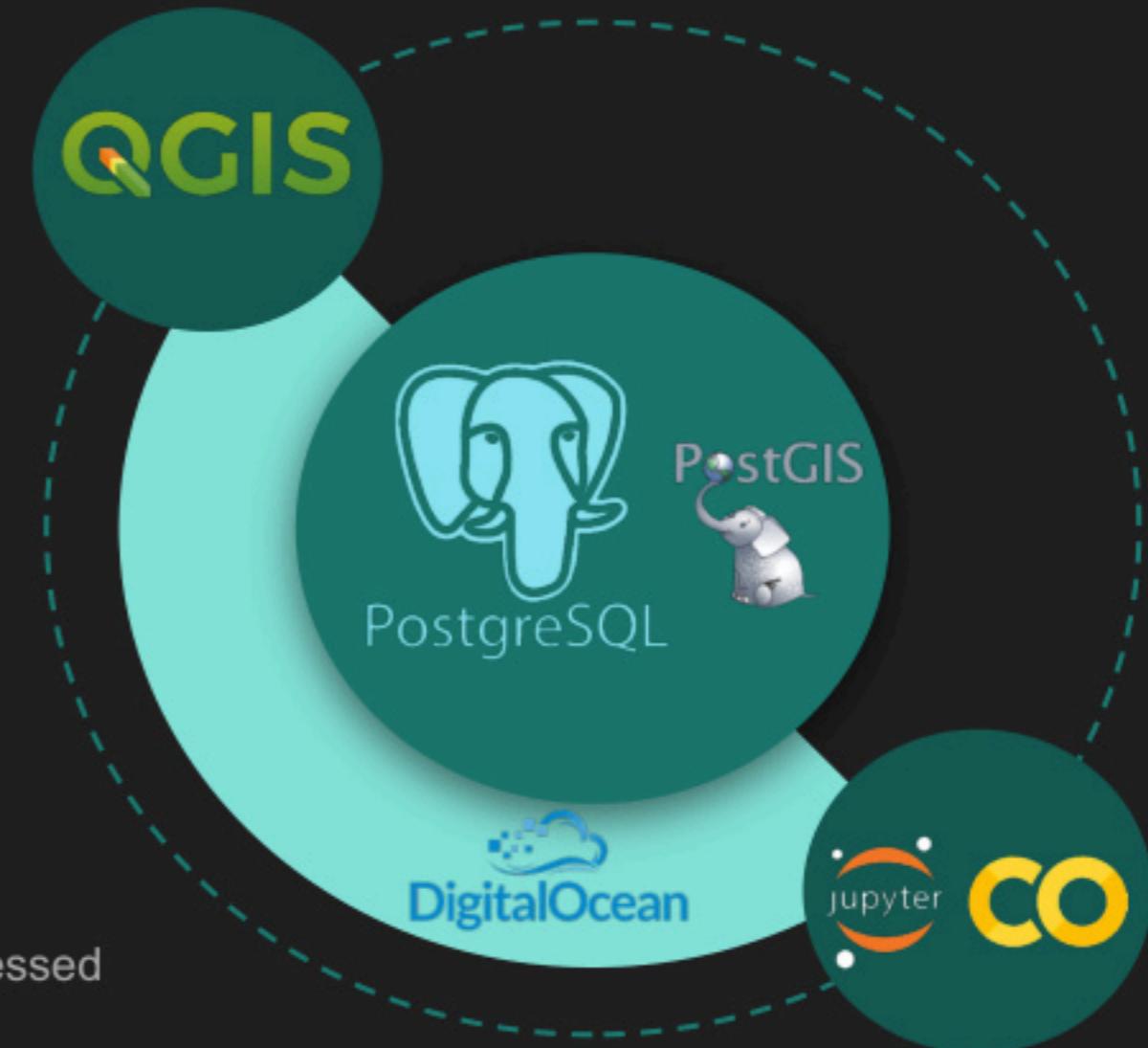
iNaturalist: an online social network of people sharing biodiversity information to help each other learn about nature.



USGS Multi-Resolution Land Characteristics Consortium (MRLC) land cover data from latest 2021 survey



USDA Web Soil Survey: tiled area of interest over City of Albuquerque processed in QGIS



# pgAdmin 4



pgAdmin 4 - Object Tools Edit Window Help

defaultdb/deadadmin@capstone\_10\* Untitled\*

Object Explorer

- Languages
- Publications
- Schemas (1)
- public
- Aggregates
- Collations
- Domains
- FTS Configurations
- FTS Dictionaries
- FTS Parsers
- FTS Templates
- Foreign Tables
- Functions
- Materialized Views
- Operators
- Procedures
- Sequences
- Tables (11)
- ABQ\_ciliate\_2018\_10
- lat
- Naturalet
- Keystone\_lat\_long
- landcover\_vectorized
- landuse
- lat\_long
- rectangle
- soil\_area
- spatial\_ref\_sys
- trigger\_menus
- Trigger Functions
- Types
- Views
- Subscriptions
- Login/Group Roles (16)
- arch
- deadadmin
- ds
- pg\_checkpoint
- pg\_database\_owner
- pg\_statistic\_snapshot

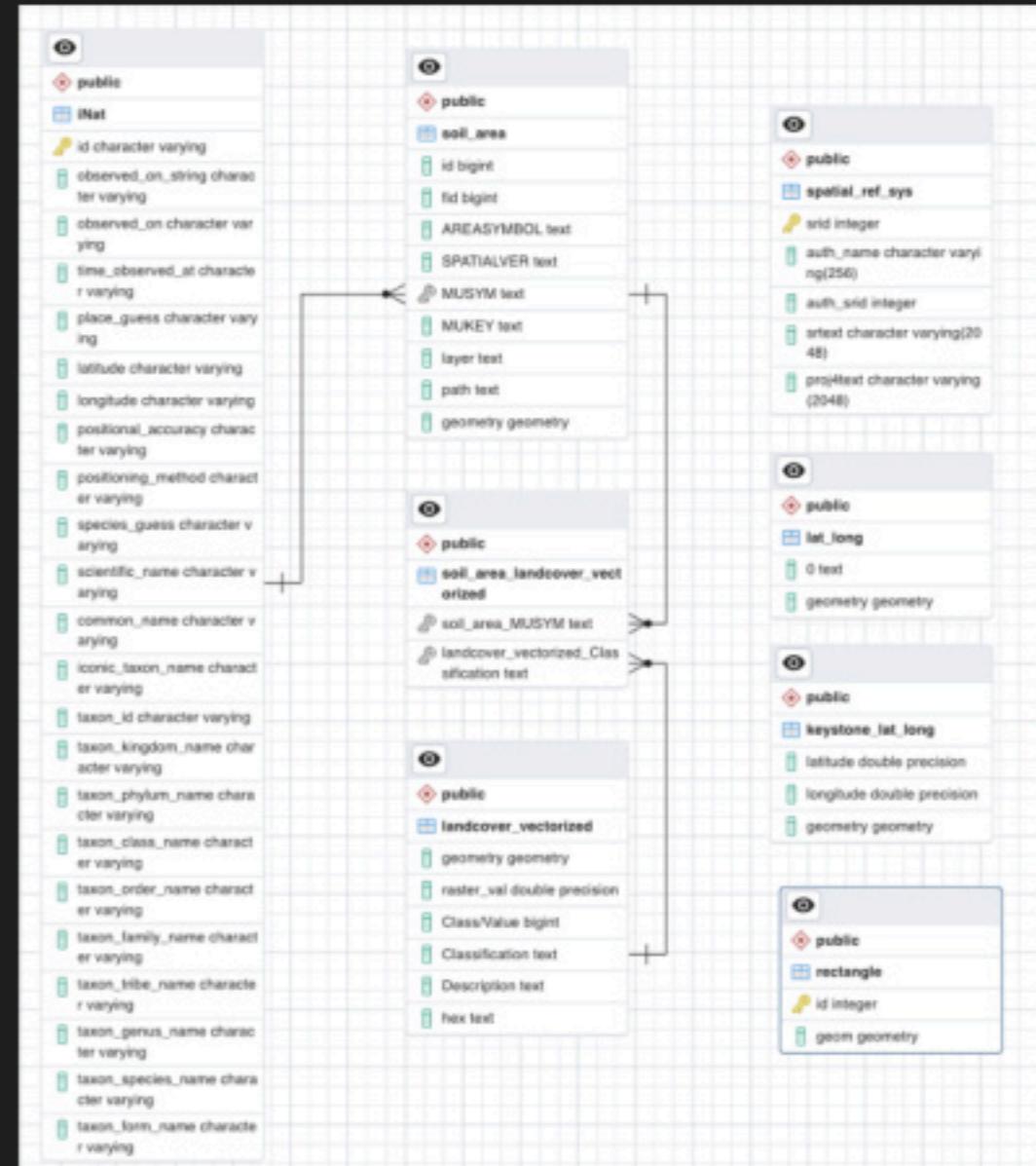
Query Query History

```
1 SELECT * FROM public."iNat" WHERE (scientific_name LIKE 'Populus%' AND scientific_name != 'Populus tremuloides')
2 OR scientific_name LIKE 'Salix exigua'
3 OR scientific_name LIKE 'Canis latrans'
4 OR scientific_name LIKE 'Castor canadensis'
5 OR scientific_name LIKE 'Cynocephalus'
6 OR common_name LIKE 'Hummingbirds'
7 OR scientific_name LIKE 'Hesperocamelus'
8 OR taxon_family_name LIKE 'Aptidae'
9
10 ORDER BY scientific_name ASC
```

Data Output Messages Notifications

id	name	common_name	isotax_base_name	taxon_id	taxon_kingdom_name	taxon_phylum_name	taxon_class_name	taxon_order_name	taxon_family_name
2505	or canadensis	American Beaver	Mammalia	43794	Animalia	Chordata	Mammalia	Rodentia	Castoridae
2506	or canadensis	American Beaver	Mammalia	43794	Animalia	Chordata	Mammalia	Rodentia	Castoridae
2507	or canadensis	American Beaver	Mammalia	43794	Animalia	Chordata	Mammalia	Rodentia	Castoridae
2508	or canadensis	American Beaver	Mammalia	43794	Animalia	Chordata	Mammalia	Rodentia	Castoridae
2509	or canadensis	American Beaver	Mammalia	43794	Animalia	Chordata	Mammalia	Rodentia	Castoridae
2510	or canadensis	American Beaver	Mammalia	43794	Animalia	Chordata	Mammalia	Rodentia	Castoridae
2511	or canadensis	American Beaver	Mammalia	43794	Animalia	Chordata	Mammalia	Rodentia	Castoridae
2512	or canadensis	American Beaver	Mammalia	43794	Animalia	Chordata	Mammalia	Rodentia	Castoridae
2513	or canadensis	American Beaver	Mammalia	43794	Animalia	Chordata	Mammalia	Rodentia	Castoridae
2514	or canadensis	American Beaver	Mammalia	43794	Animalia	Chordata	Mammalia	Rodentia	Castoridae
2515	or canadensis	American Beaver	Mammalia	43794	Animalia	Chordata	Mammalia	Rodentia	Castoridae
2516	or canadensis	American Beaver	Mammalia	43794	Animalia	Chordata	Mammalia	Rodentia	Castoridae
2517	or canadensis	American Beaver	Mammalia	43794	Animalia	Chordata	Mammalia	Rodentia	Castoridae
2518	or canadensis	American Beaver	Mammalia	43794	Animalia	Chordata	Mammalia	Rodentia	Castoridae
2519	or canadensis	American Beaver	Mammalia	43794	Animalia	Chordata	Mammalia	Rodentia	Castoridae
2520	or canadensis	American Beaver	Mammalia	43794	Animalia	Chordata	Mammalia	Rodentia	Castoridae

Query Tool with SQL Console



Entity Relationship Diagram (ERD)

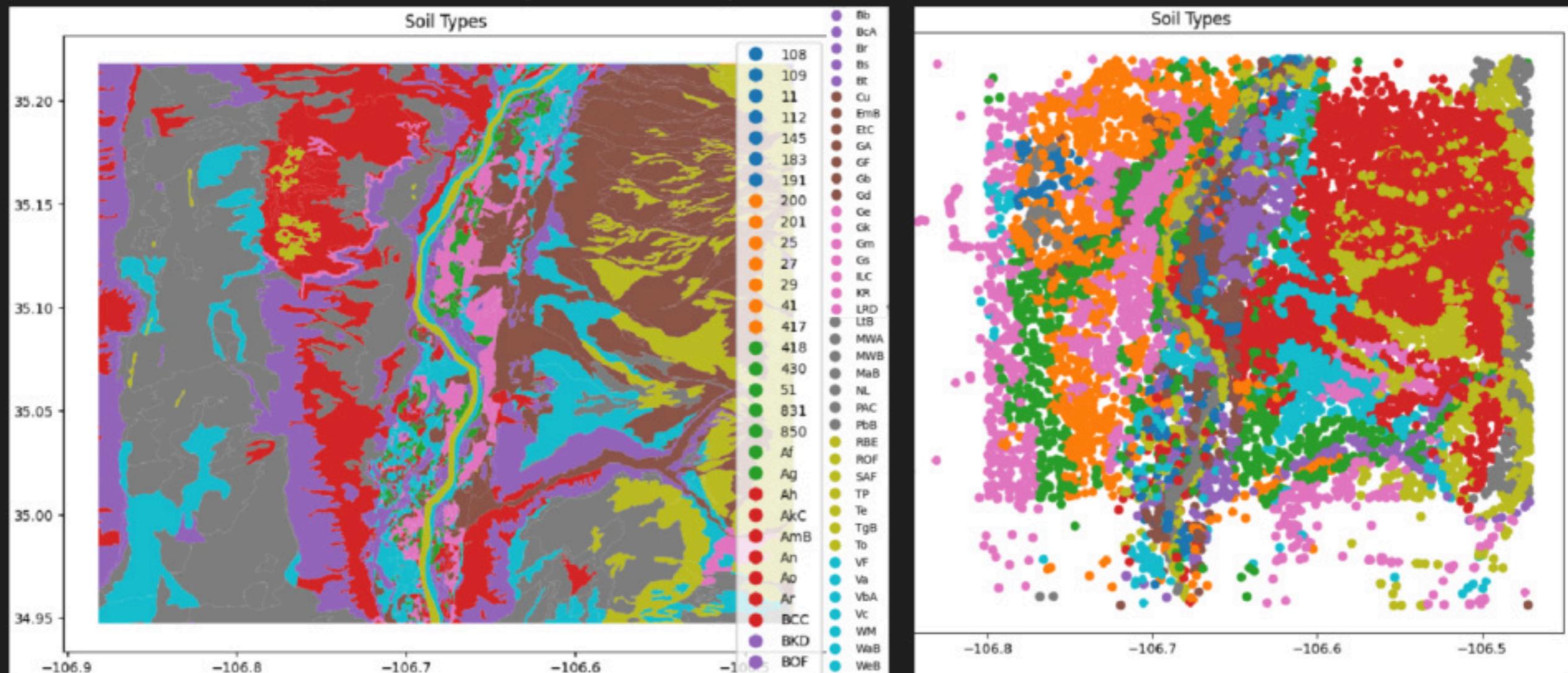
# Technology Overview



## Machine Learning Techniques Used

- Decision Trees
- Random Forest
- Logistic Regression
- K-Means Clustering
- XGBoost
  - Feature Engineering
- Hierarchical Clustering
- RMSE

# Feature Engineering: Soil Type



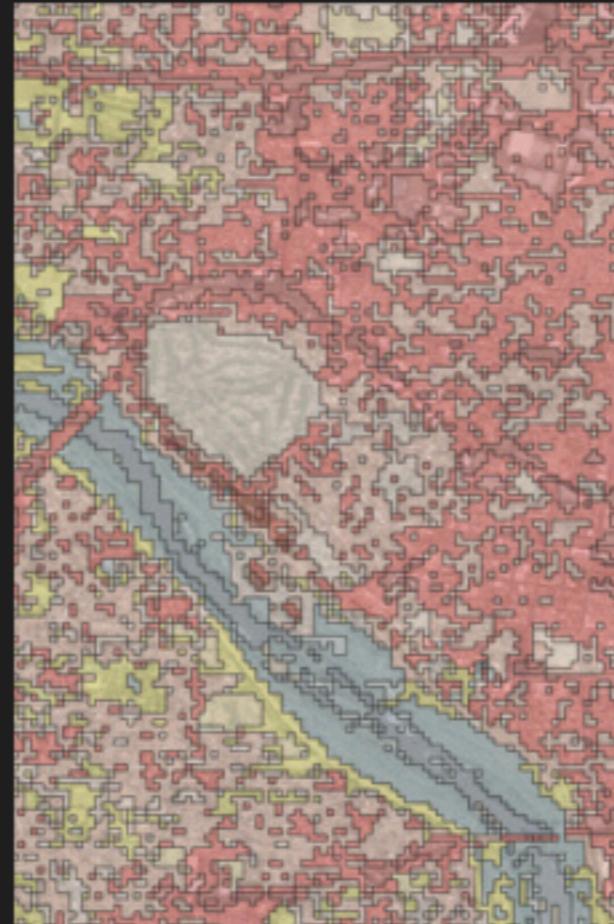
# Feature Engineering: Land Cover



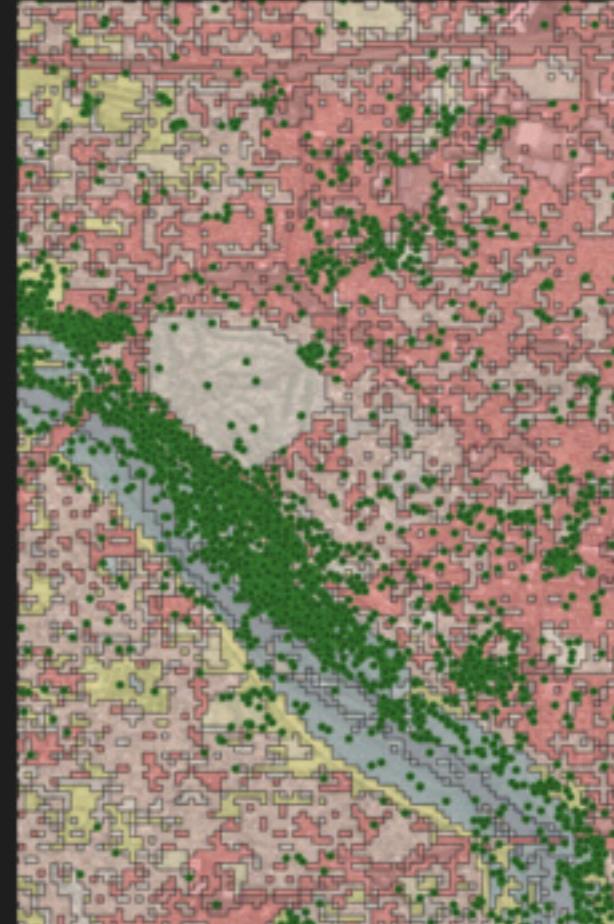
Aerial Image  
(Raster)



Land Cover  
(Raster)



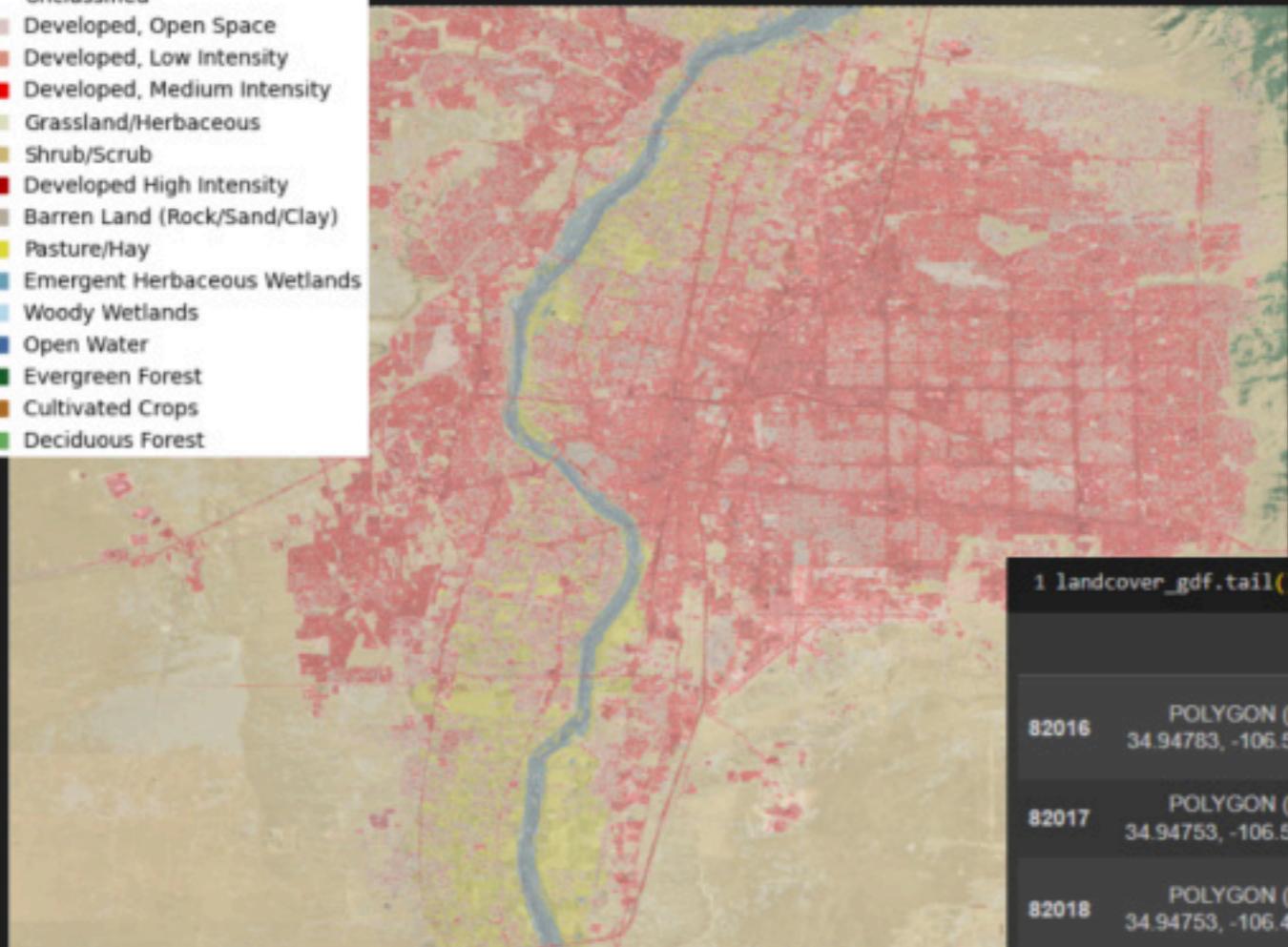
Land Cover  
(Vectorized)



iNaturalist  
Observations

# Feature Engineering: Land Cover

Unclassified
Developed, Open Space
Developed, Low Intensity
Developed, Medium Intensity
Grassland/Herbaceous
Shrub/Scrub
Developed High Intensity
Barren Land (Rock/Sand/Clay)
Pasture/Hay
Emergent Herbaceous Wetlands
Woody Wetlands
Open Water
Evergreen Forest
Cultivated Crops
Deciduous Forest

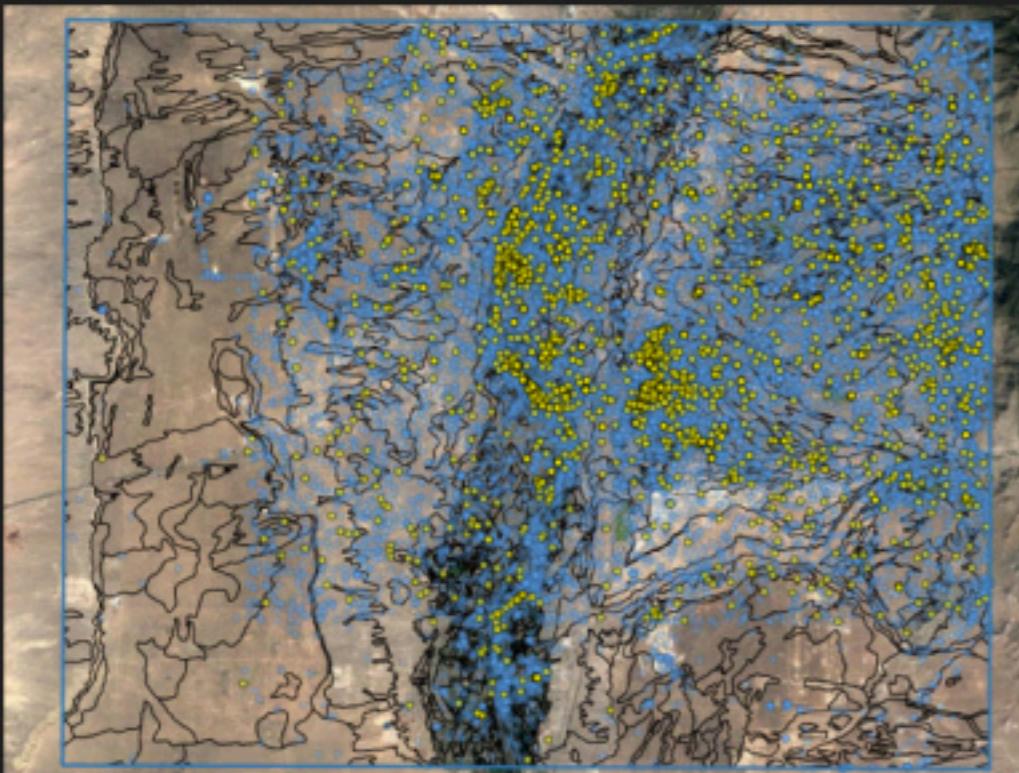


**City of Albuquerque  
2021 Land Cover**

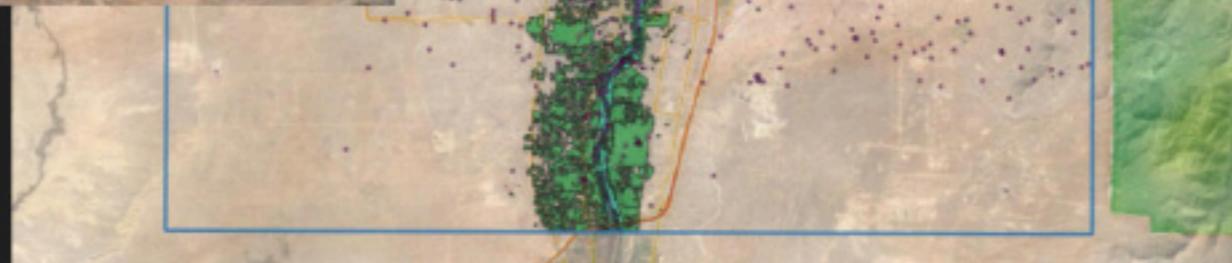
## Land Cover Classification to GeoDataFrame

		geometry	raster_val	Class/Value	Classification	Description	hex
82016		POLYGON ((-106.50144 34.94783, -106.50114 34.9...	71.0	71	Grassland/Herbaceous	areas dominated by graminoid or herbaceous veg...	#dtdfc2
82017		POLYGON ((-106.50024 34.94753, -106.50024 34.9...	52.0	52	Shrub/Scrub	areas dominated by shrubs; less than 5 meters ...	#ccb879
82018		POLYGON ((-106.49934 34.94753, -106.49934 34.9...	71.0	71	Grassland/Herbaceous	areas dominated by graminoid or herbaceous veg...	#dtdfc2
82019		POLYGON ((-106.48465 34.95143, -106.48375 34.9...	71.0	71	Grassland/Herbaceous	areas dominated by graminoid or herbaceous veg...	#dtdfc2
82020		POLYGON ((-106.48855 35.06504, -106.48795 35.0...	52.0	52	Shrub/Scrub	areas dominated by shrubs; less than 5 meters ...	#ccb879

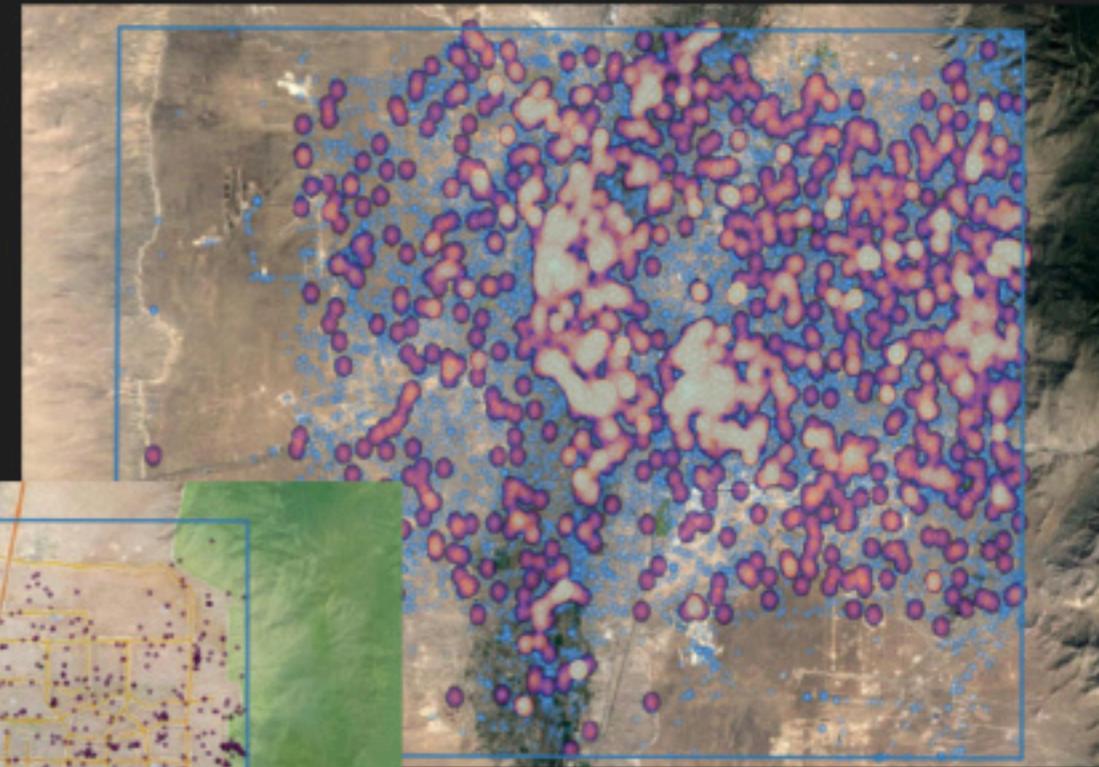
# Geospatial EDA



Soil Polygons with  
iNaturalist Observations



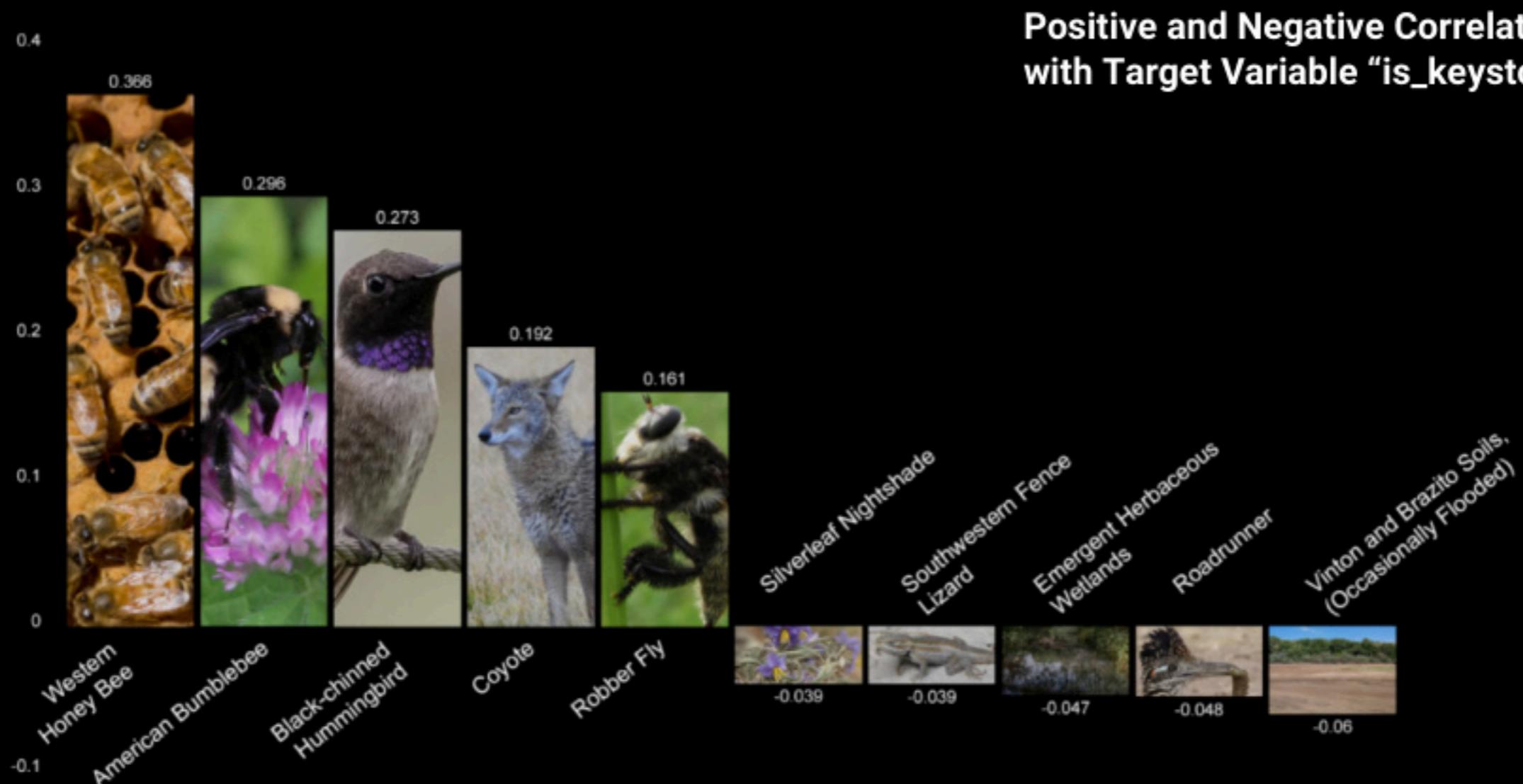
Keystone Species vs. Riparian Land Cover



Heat Map Showing  
Distribution of Keystone  
Species

**QGIS**

## Combined Features Analysis (Excluding Decision Tree)

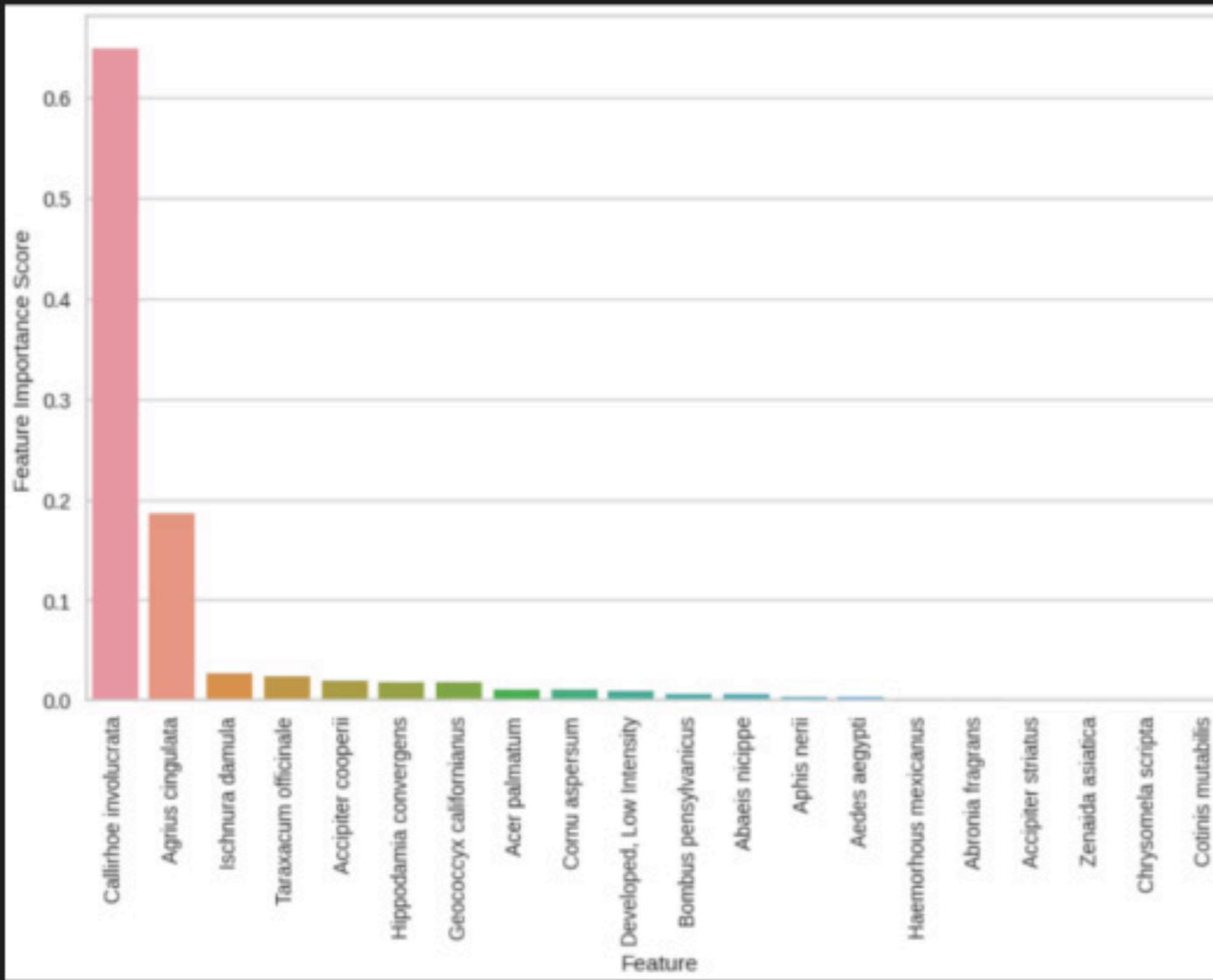


This updated bar chart combines the top 5 positively and negatively correlated features with the target variable 'is\_keystone', excluding the decision tree important features. The x-axis labels are directly below the bars for easier identification.

# Honey Bees!

RMSE: 0.09681156163532739

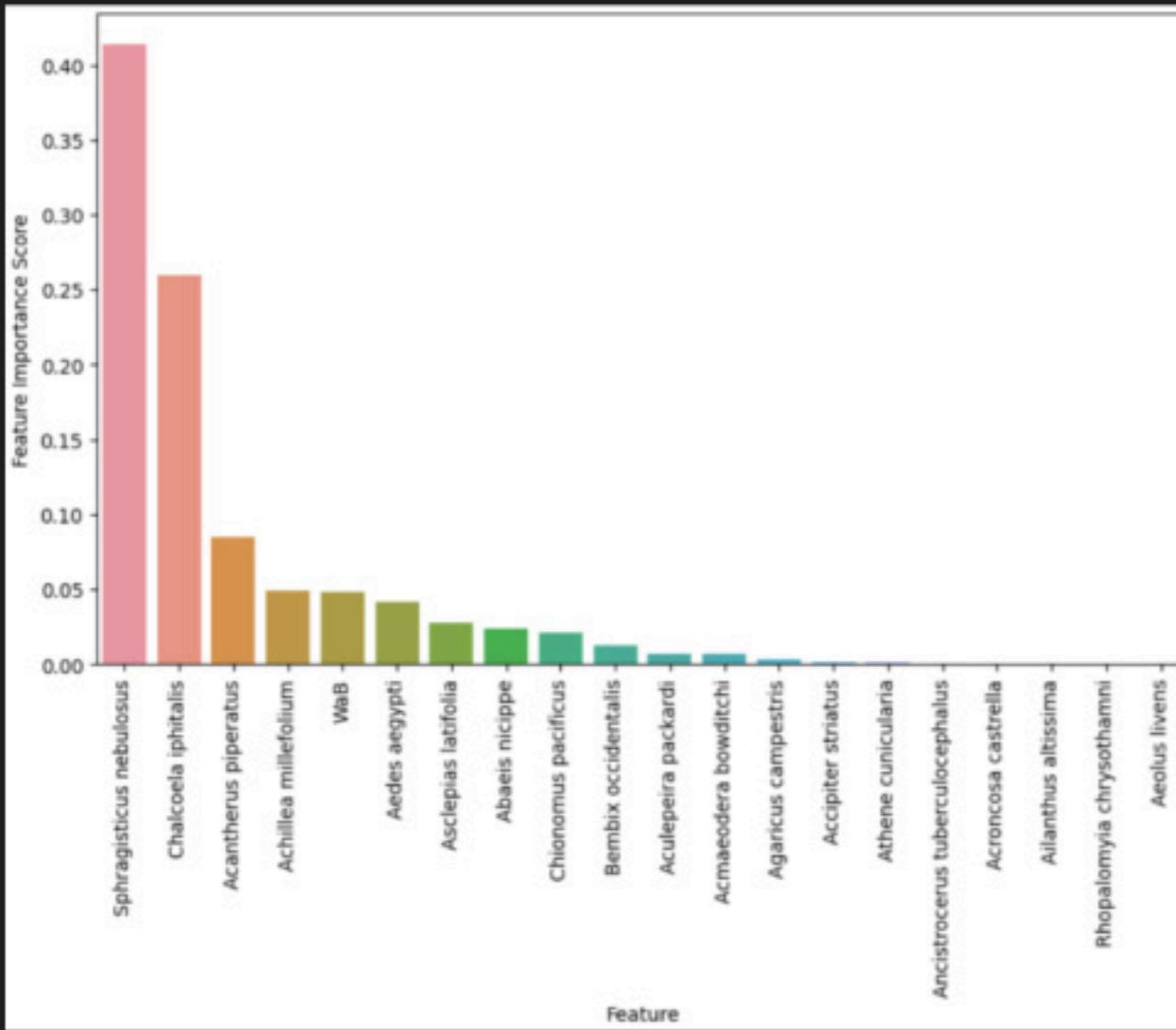
- The most important feature for the presence of honey bees in XGBoost is unsurprisingly a flower



Purple Poppy-Mallow, courtesy of  
<https://en.wikipedia.org>

# Prairie Dogs!

RMSE: 0.1611310832860667



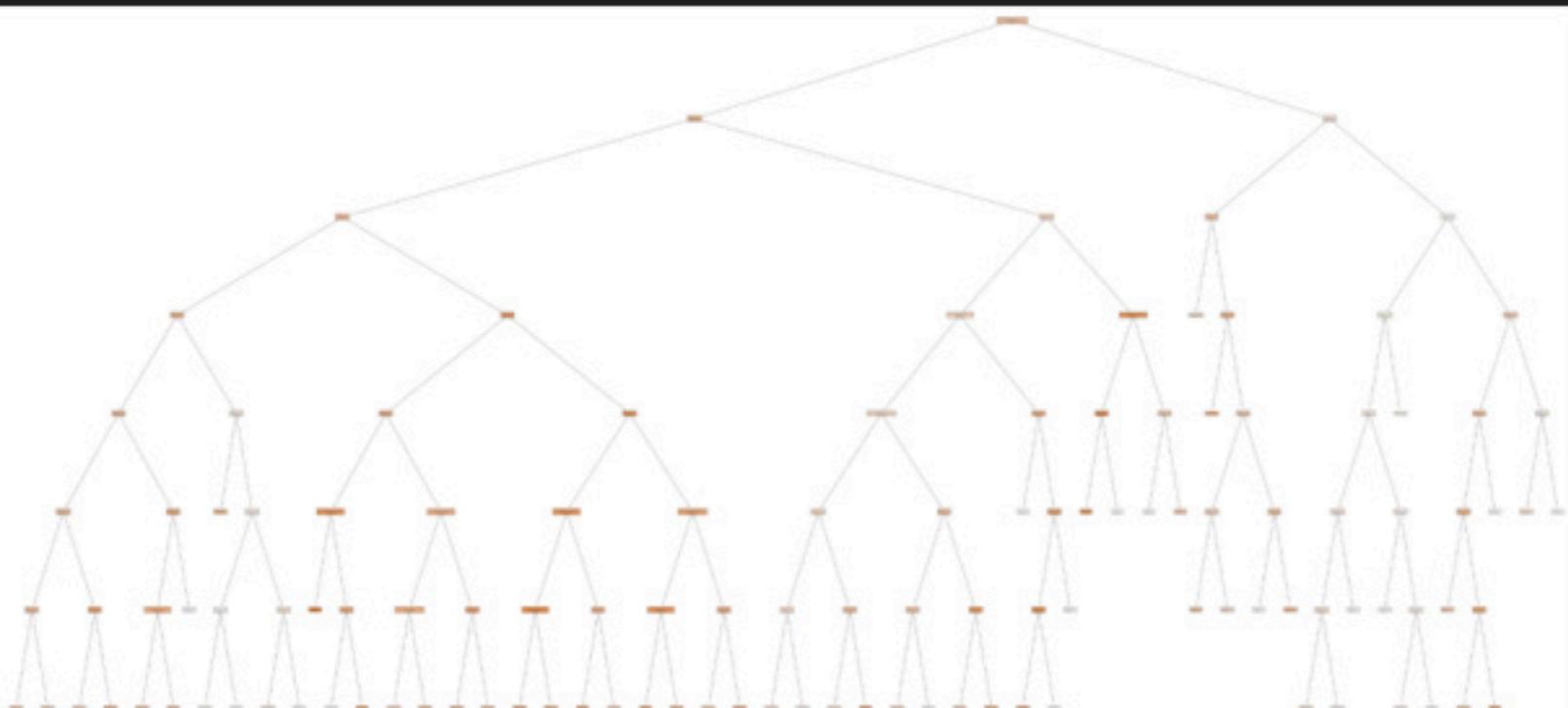
Seed bug courtesy of  
<https://www.britishbugs.org.uk>

# XGBoost Derived Environmental Feature Importance (Top 3)

- Coyote - RMSE 0.0645
  - (0.1968) Gila Loam, moderately alkali
  - (0.0677) Rock outcrop orthids complex, 40 to 80 percent slopes
  - (0.0406) Emergent herbaceous wetlands
- Beaver - RMSE 0.0440
  - (0.1451) Pajarito loamy fine sand, 1 to 9 percent slopes
  - (0.0998) Grassland/herbaceous
  - (0.0955) Shrub/scrub
- Cottontail Rabbit - RMSE 0.0898
  - (0.1324) Cut and fill land
  - (0.1011) Shrub/scrub
  - (0.0740) Brazito complex soil
- Gunnison's Prairie Dog - RMSE 0.2475
  - (0.7294) Emergent herbaceous wetlands
  - (0.2455) Pasture/hay
  - (0.0071) Shrub/scrub
- Black-chinned Hummingbird - RMSE 0.0903
  - (0.1944) Emergent herbaceous wetlands
  - (0.0730) Woody wetlands
  - (0.0614) Grassland/herbaceous
- Western Honey Bee - RMSE 0.1202
  - (0.3506) Cut and fill land
  - (0.1218) Embudo-Tijeras complex, 0 to 9 percent slopes
  - (0.0542) Tijeras gravelly fine sandy loam, 1 to 5 percent slopes
- American Bumblebee - RMSE 0.0983
  - (0.1008) Gila loam, 0 to 1 percent slopes
  - (0.0657) Tijeras gravelly fine sandy loam, 1 to 5 percent slopes
  - (0.0651) Gila clay loam
- Coyote Willow - RMSE 0.0266
  - (0.2508) Emergent herbaceous wetlands
  - (0.1345) Torrifluvents, frequently flooded
  - (0.0815) Pasture/hay
- Chamisa - RMSE 0.0506
  - (0.6287) Cultivated crops
  - (0.0830) Developed - high intensity
  - (0.0291) Bluepoint-Kokan association, hilly
- Rio Grande Cottonwood - RMSE 0.0401
  - (0.1288) Tome very fine sandy loam
  - (0.1249) Developed - low intensity
  - (0.0894) Pasture/hay

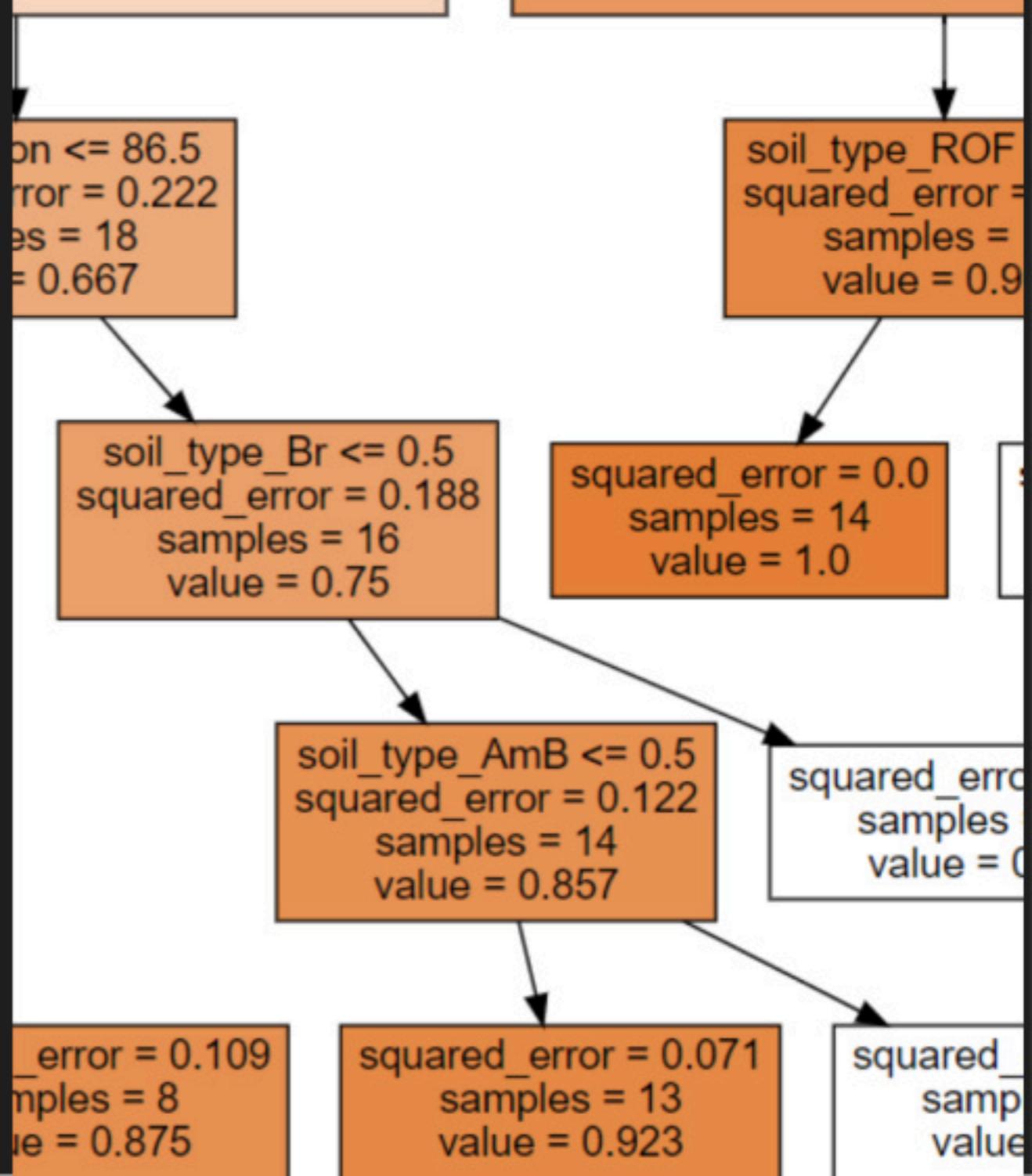
# Regression Comparisons

- iNaturalist - Soil - Land Use
  - Filtered for target keystone genera - target
- XGBoost
  - Calculated RMSE and top 20 most important features
- Decision Tree
  - Filtered out all features except for top 20



## Regression Comparisons Cont.

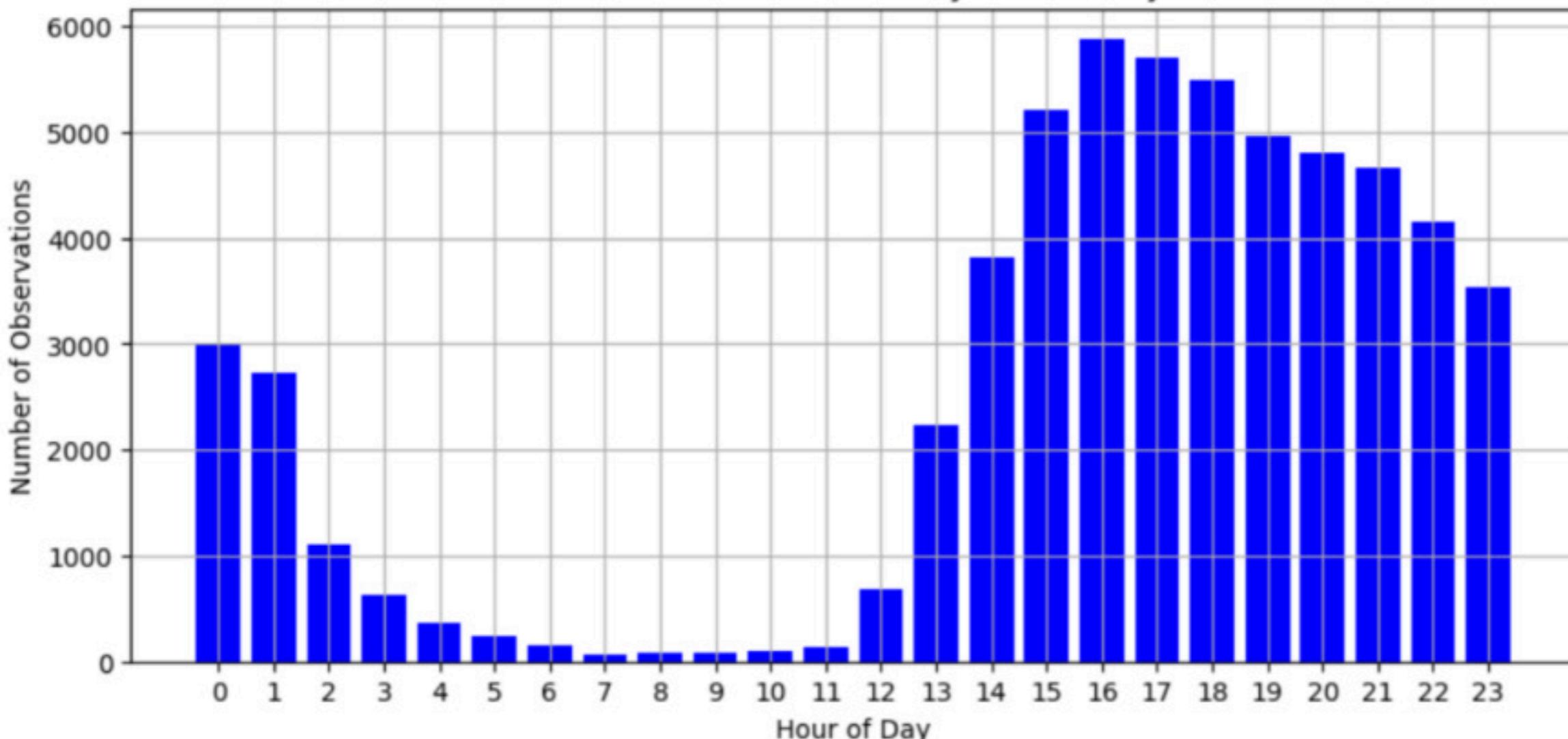
- Decision Tree Cont.
  - RMSE - Bees, Hummingbirds
  - Max Depth
- Random Forest
  - RMSE - Bees



# Biases/Data Struggles (UTC)

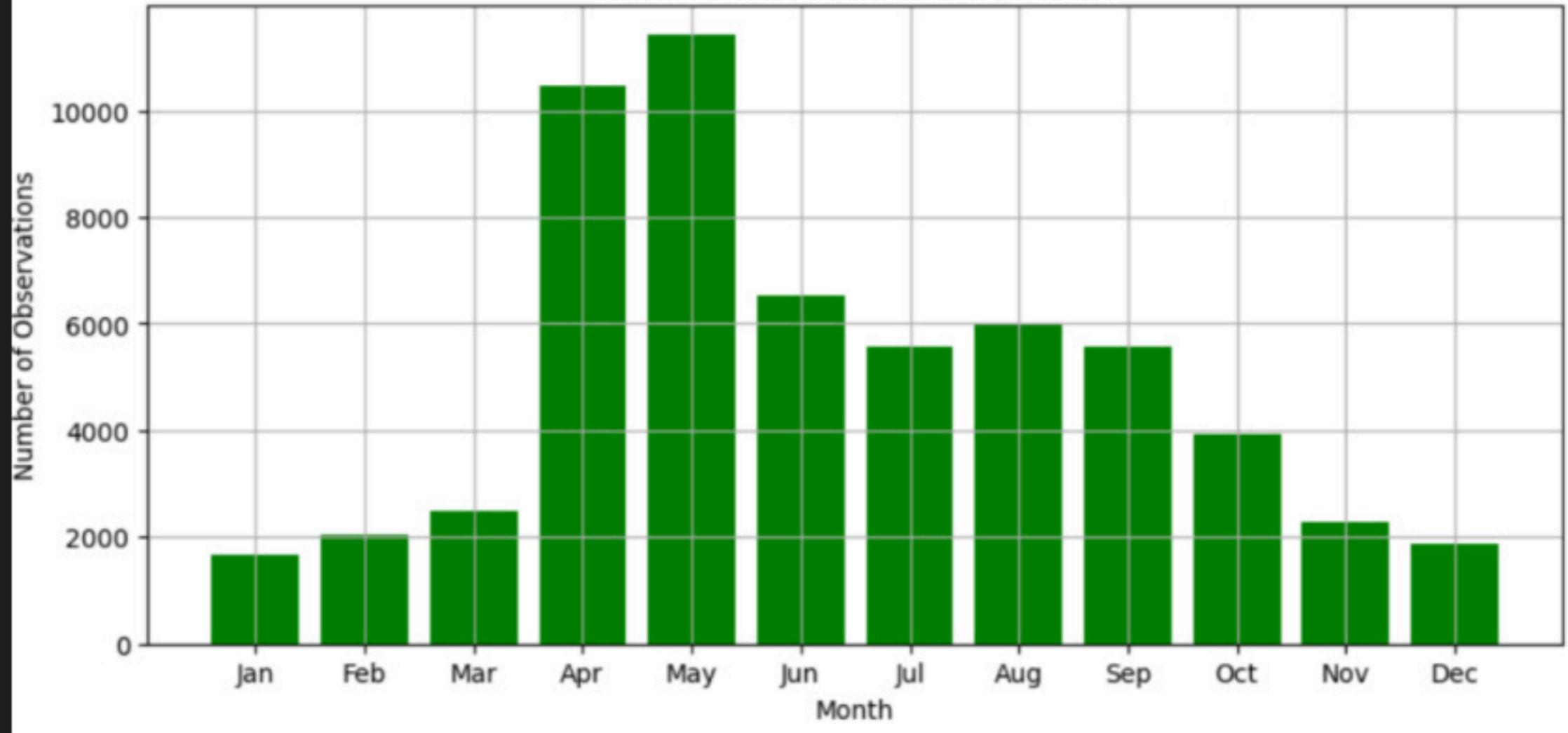
Mountain Standard Time (North America) is 7 hours behind **UTC** universal time

Observation Distribution by Time of Day



# Biases/Data Struggles

Observation Distribution by Month



## Future Study

- Due to time/data/budget constraints this project is limited in scope and precision, but with limited resources we have shown the incredible power of the tools of data science.
- Different data could easily be plugged into the same model to show factors in and prevalence of mental health problems, homelessness or disease to name a few.
- A few other geospatial data leads that would bring a richer dimension to the project would be climate data, terrain analysis, solar exposure, microclimates associated with structures, hydrography, etc.

# Acknowledgments

We want to thank CNM staff, Job Training Albuquerque, Dekker/Perich/Sabatini, Tricore and Hewlett Packard for their support in this half year endeavour.

Special thanks to Robert Citek! We will miss the mind expanding Fridays!

We've used a constellation of tools, libraries and data sources:



**DEEP DIVE  
Coding**  
powered by CNM In-enuity



## Q&A

Nicky Hughes

nbhughes27@gmail.com

LinkedIn:



Drew Seavey

drew.seavey@gmail.com

LinkedIn:



Elizabeth Fitter

elizagethgfitter@gmail.com

LinkedIn:



Project Repository on GitHub:



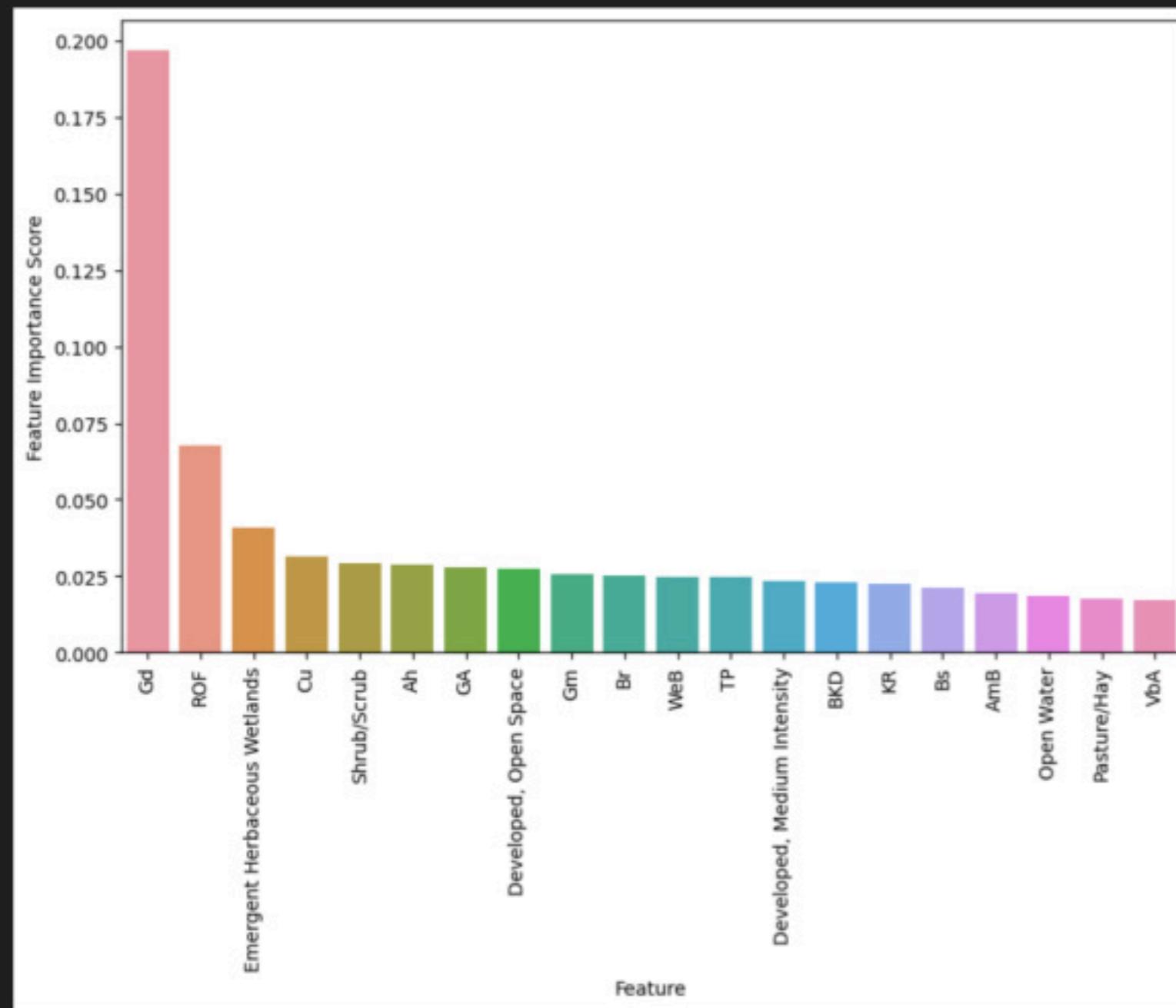
**DEEP DIVE  
Coding**

powered by CNM Ingenuity

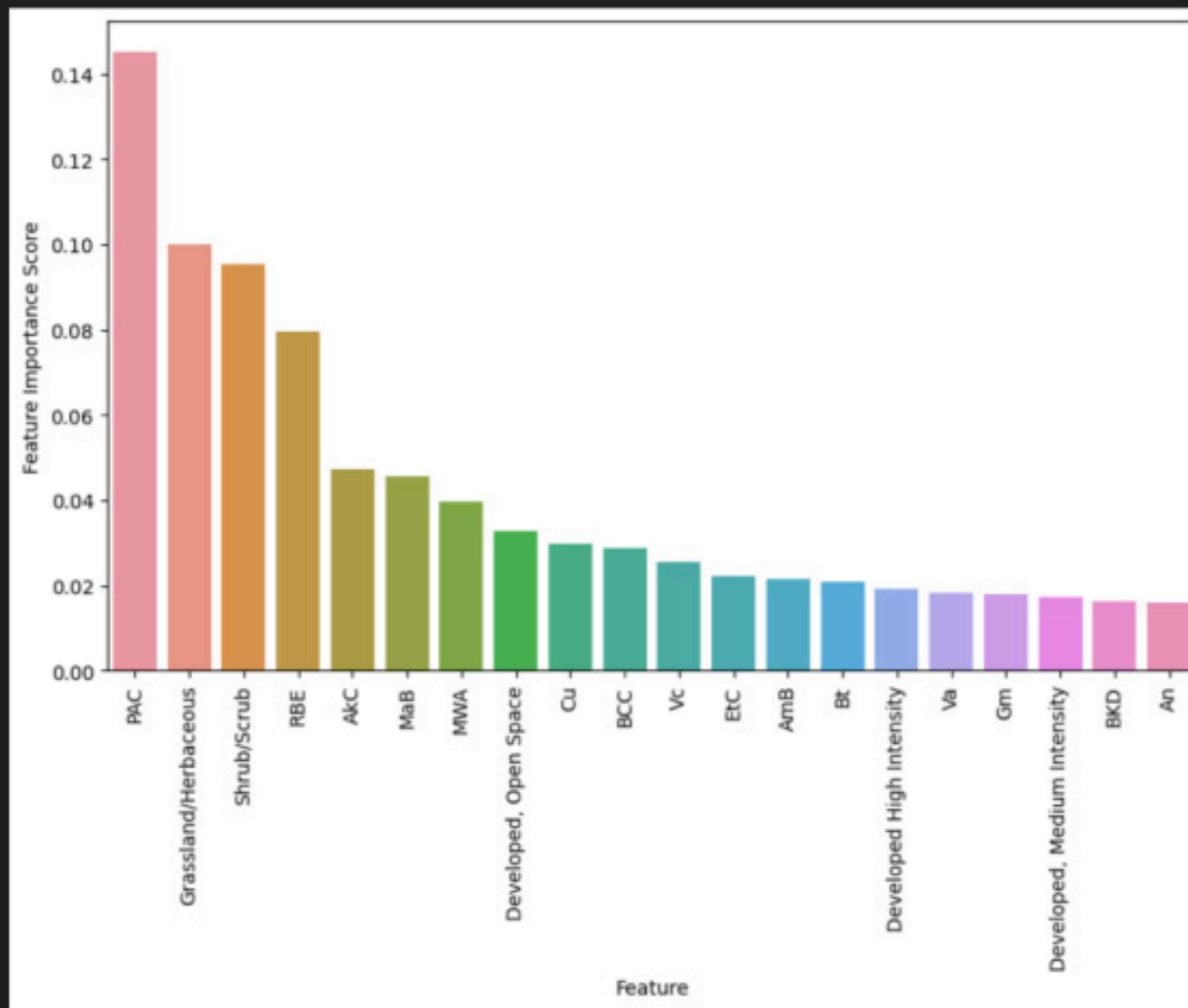
# Soil Types

Symbol	Description	Slope %	Composition
AmB	Alameda sandy loam	0 to 5	sandy loam, gravelly sandy loam, very cobbly loam, bedrock
BOF	Borolls-Rock outcrop association	very steep	Very stony loam, very stony fine sandy loam, bedrock
Br	Brazito fine sandy loam		Fine sandy loam, course sand
Cu	Cut and fill land		
Gm	Glendale clay loam	0 to 1	Clay loam, silt loam, fine sand, clay
LtB	Latene sandy loam	0 to 5	Sandy loam, gravelly sandy loam
ROF	Rock outcrop-Orthids complex	40 to 80	Bedrock
VFL	Flagstone		Very flaggy

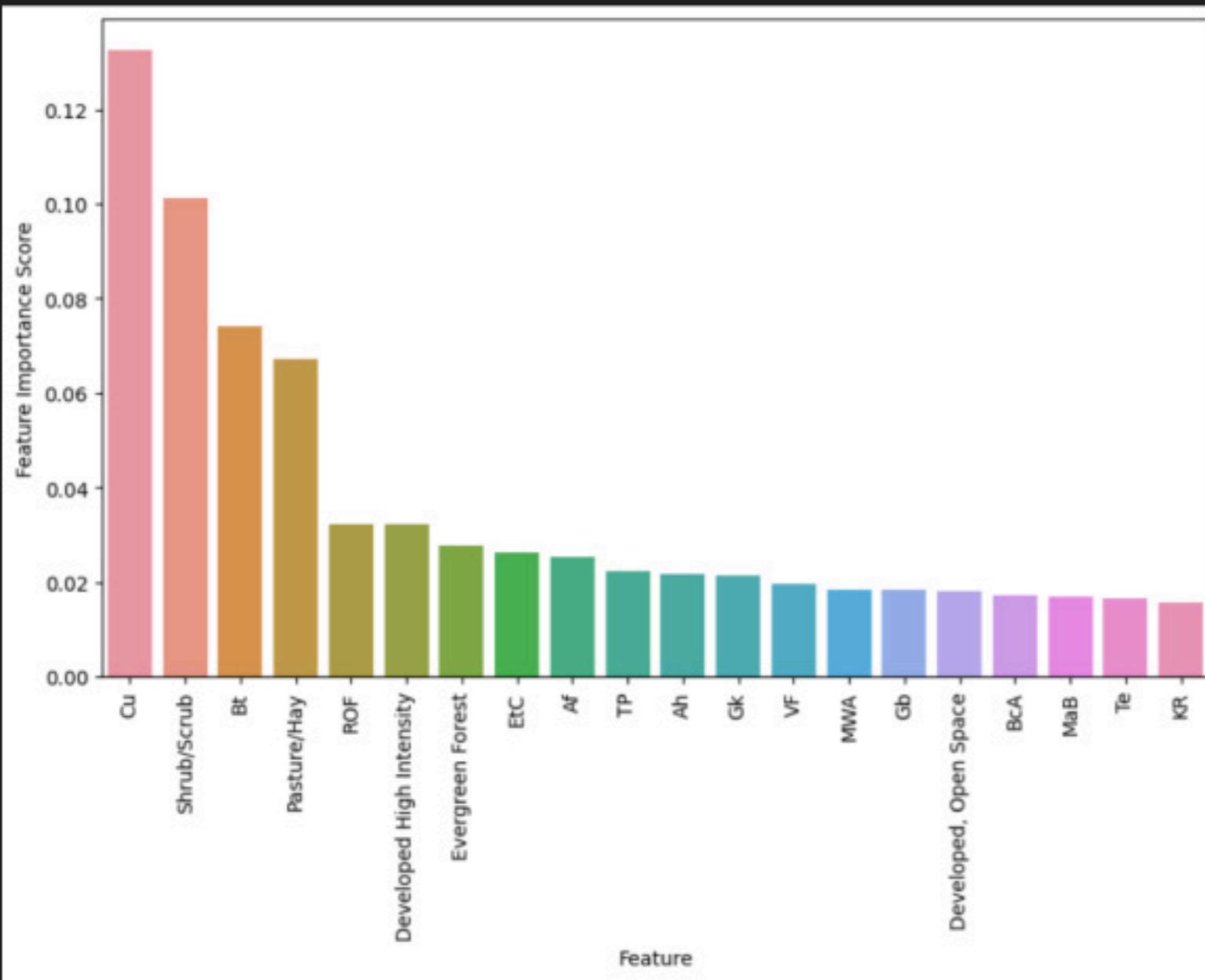
# Environmental Feature Importance - Coyote



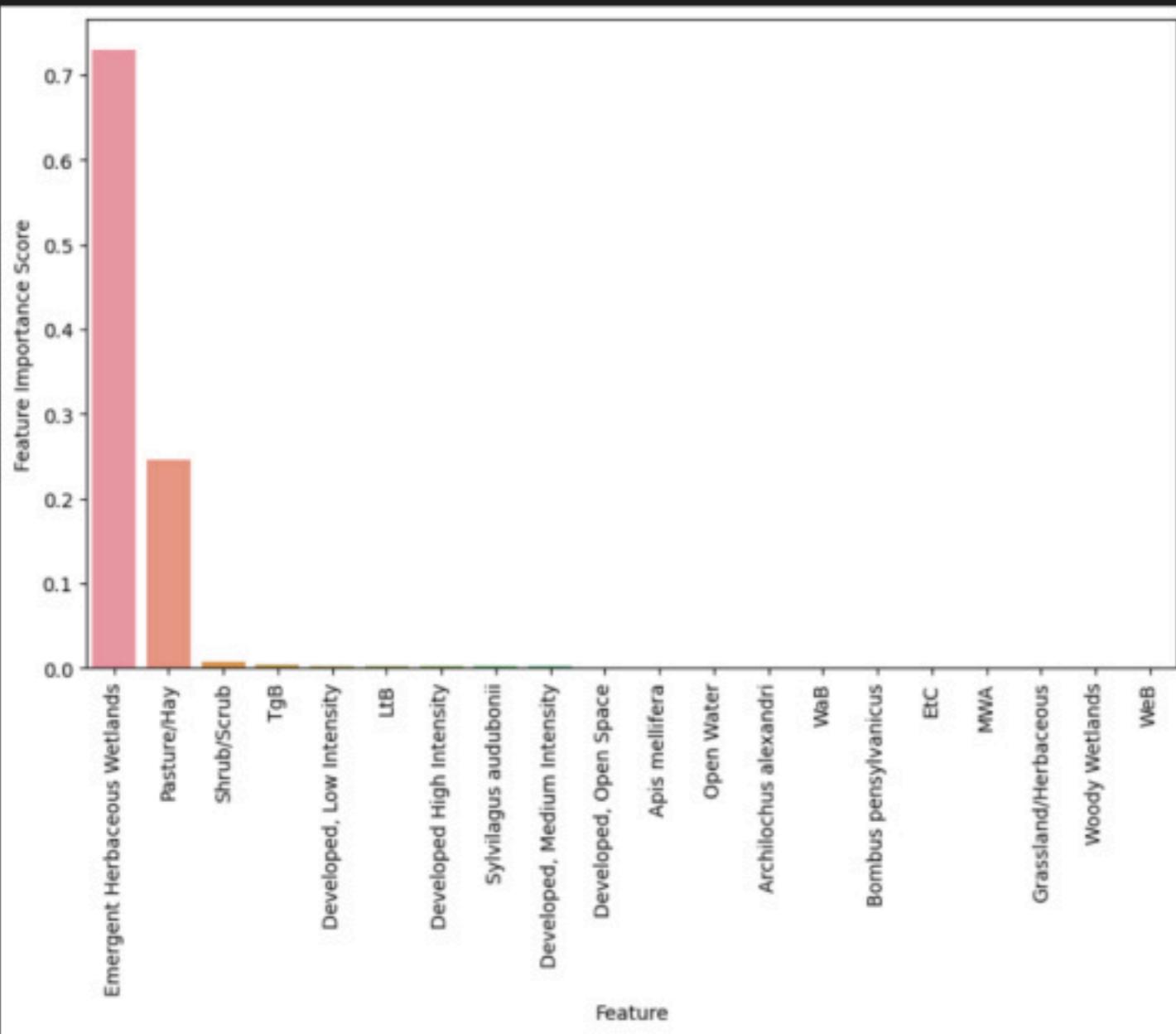
# Environmental Feature Importance - Beaver



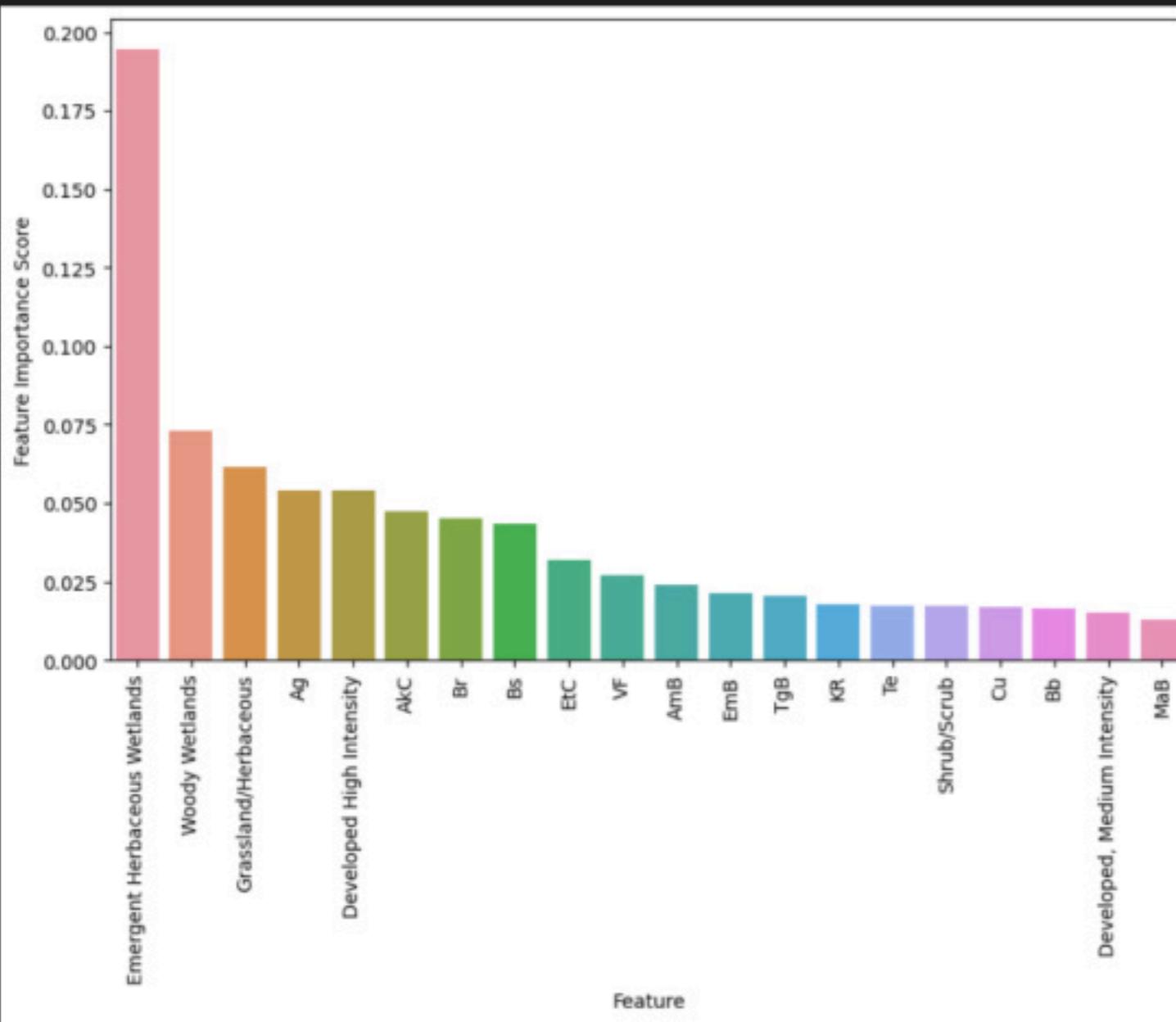
# Environmental Feature Importance - Cottontail



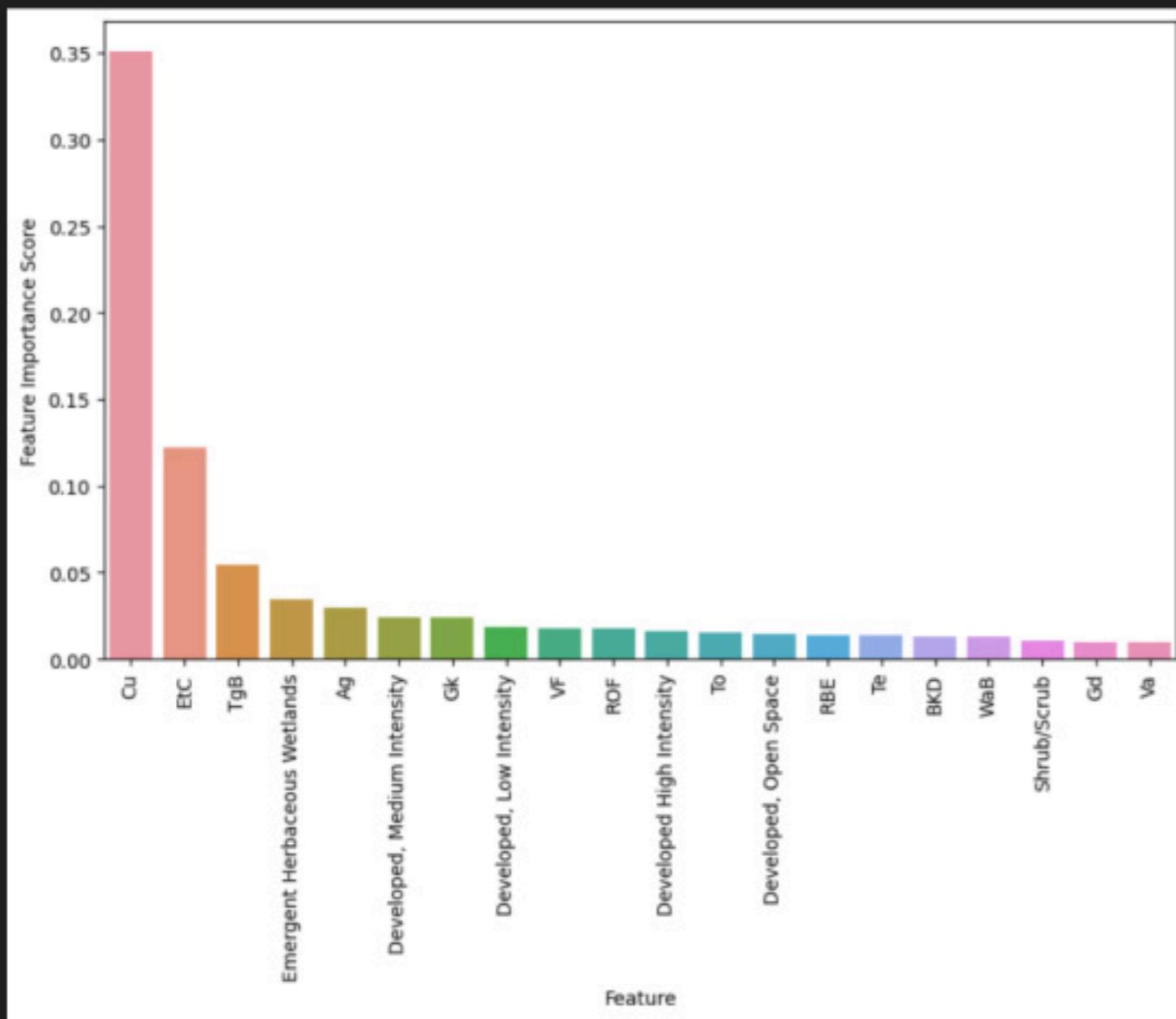
# Environmental Feature Importance - Gunnison's Prairie Dog



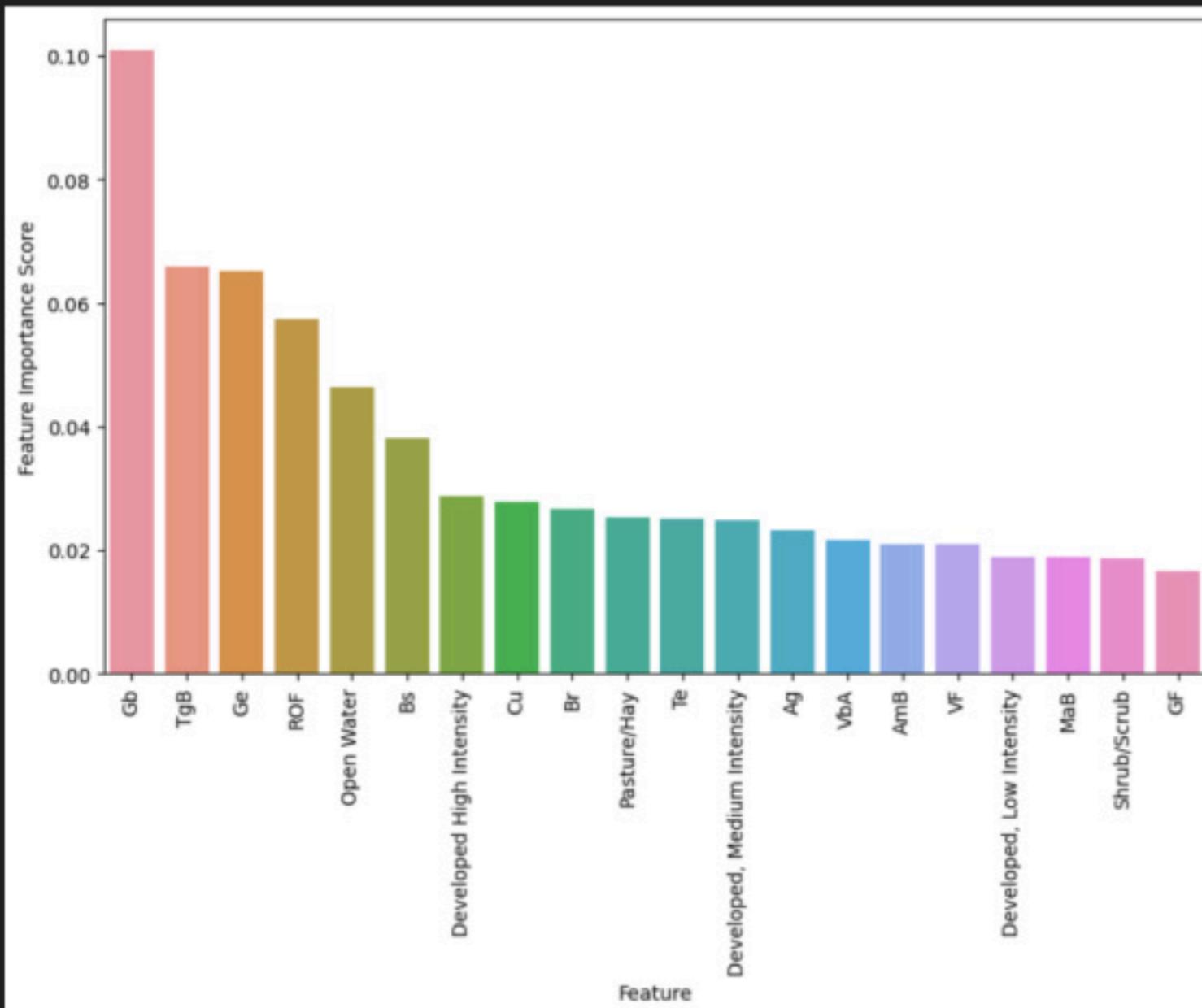
# Environmental Feature Importance - Black-chinned Hummingbird



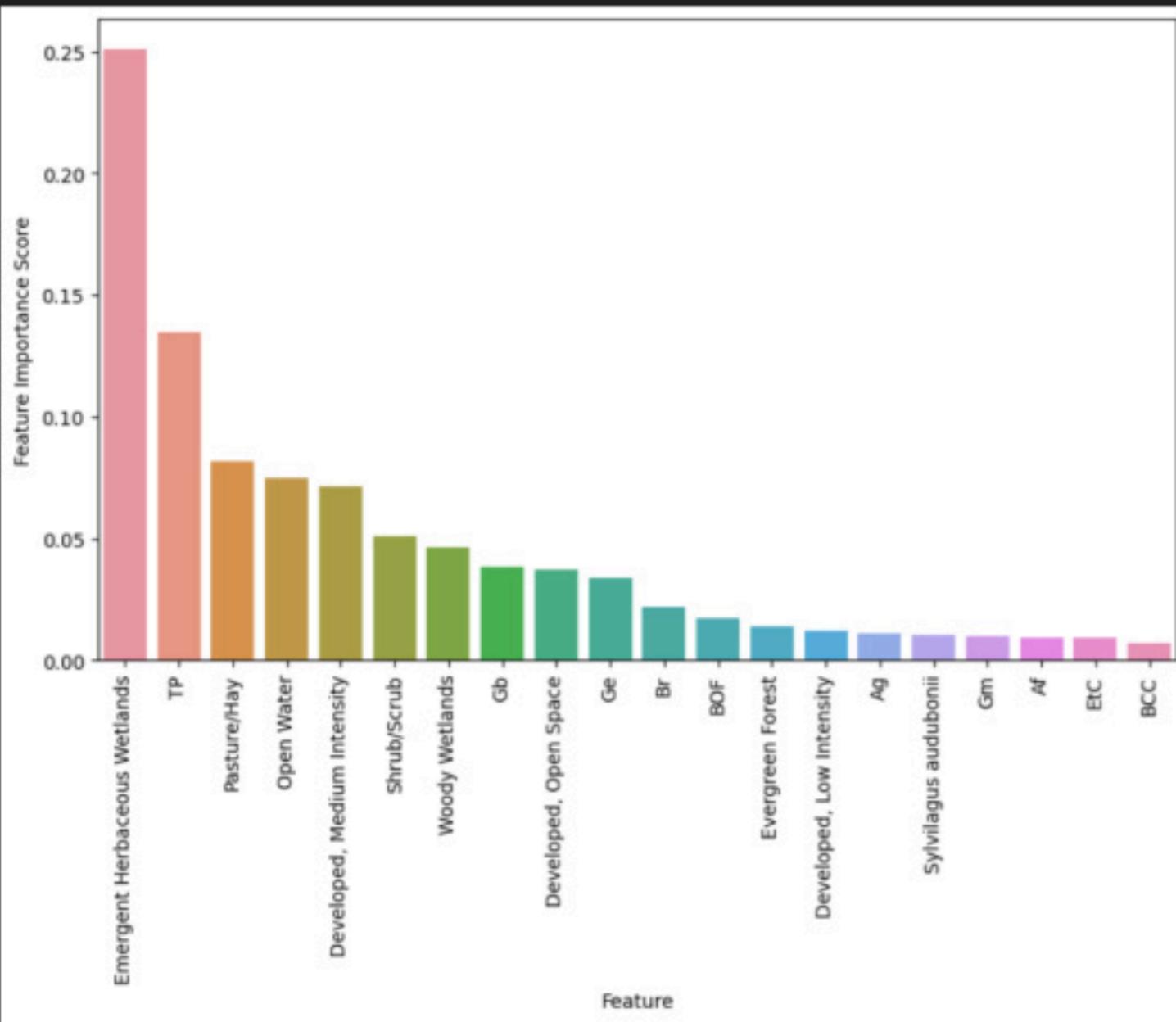
# Environmental Feature Importance - Western Honey Bee



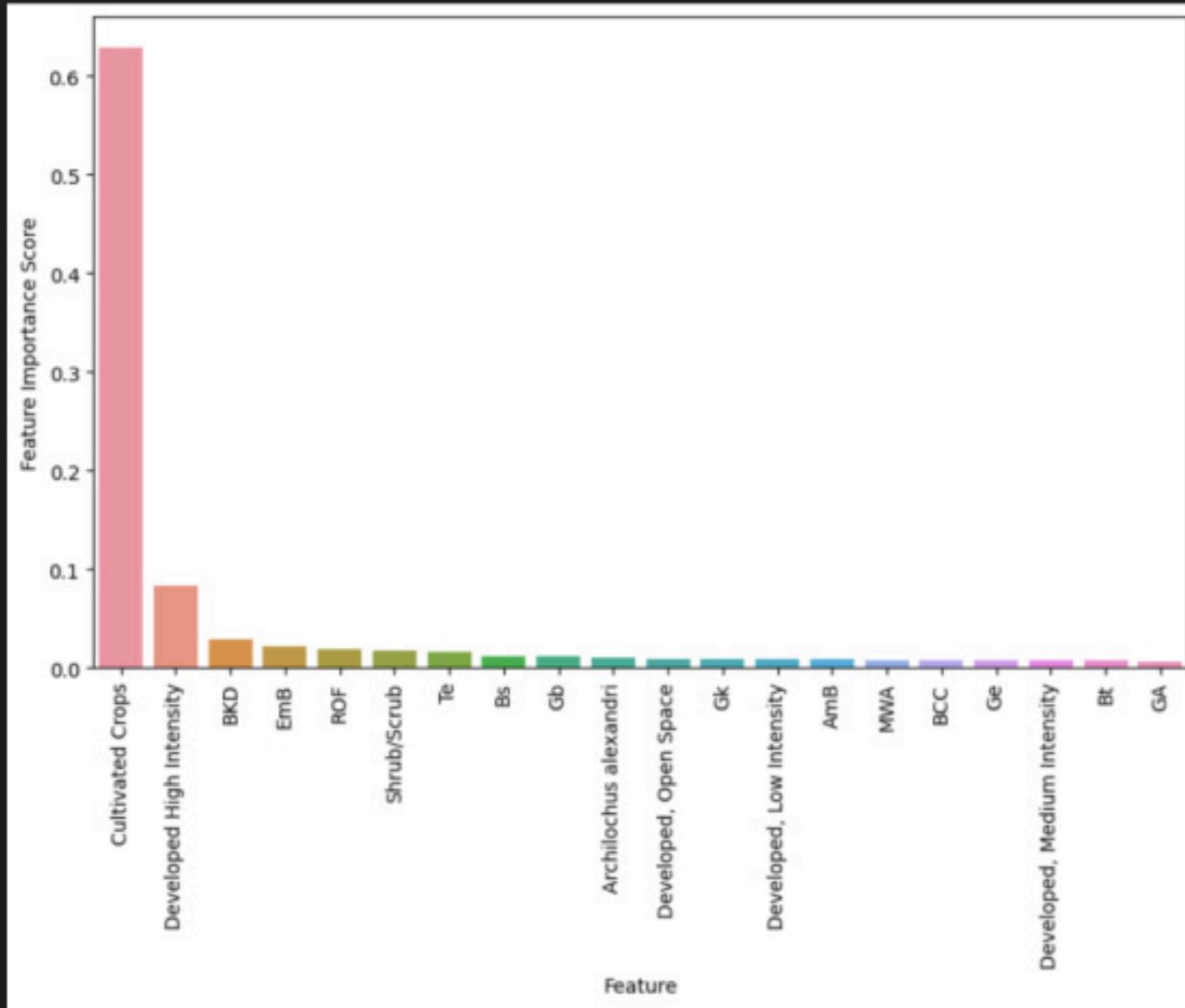
# Environmental Feature Importance - American Bumblebee



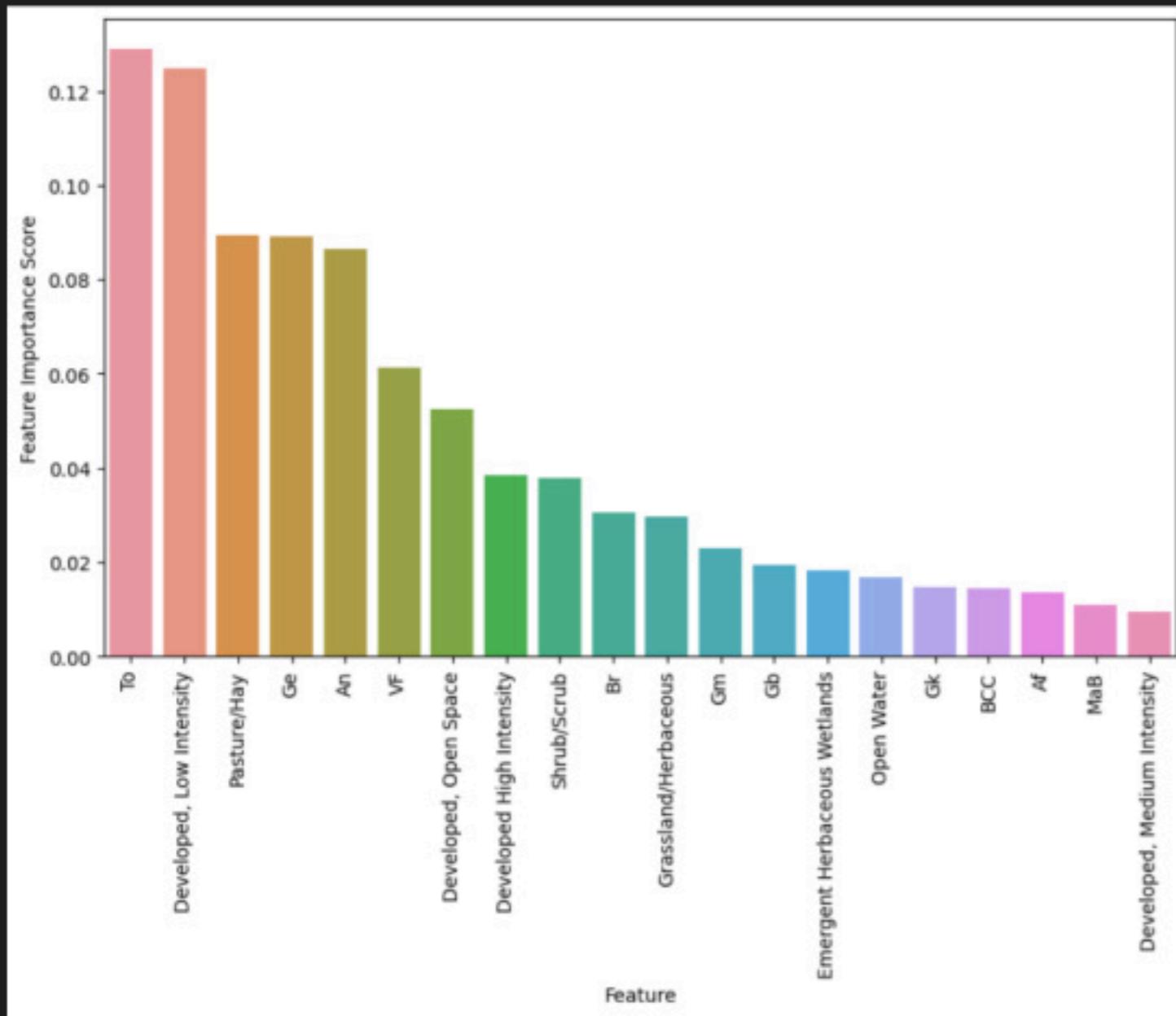
# Environmental Feature Importance - Coyote Willow



# Environmental Feature Importance - Chamisa



# Environmental Feature Importance - Rio Grande Cottonwood



# Regression Comparisons Cont.

- Logistic Regression
  - Correlation Matrix
  - Distribution of predictors
  - RMSE - bees

