

By Tom Cairns, Gabriel Deacon, Ben
Kosiborod, Amy Tang

Oil has major impacts on wildlife

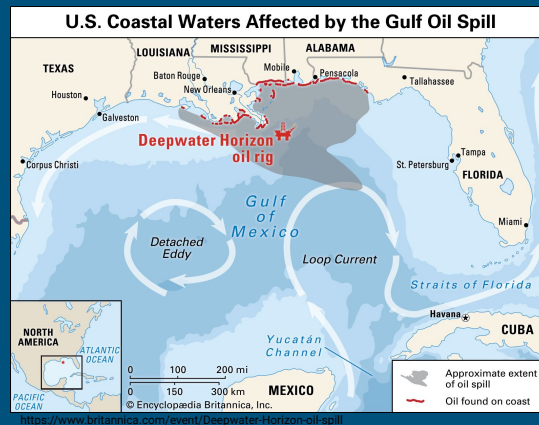
- Reduces water repellent properties of feathers
- Anemia
- Reduced immune function



<https://www.motherjones.com/environment/2019/04/deepwater-horizon-bp-oil-spill/>

Identify markers in birds most impacted by oil spills

Create a model to predict if a bird has been impacted



- Report generated to study Deepwater Horizon oil spill in 2010

- Deepwater Horizon (MC 252) Oil Spill Natural Resource Damage Assessment and Restoration
- **Report includes:**
 - PDF - written report including explanation of the variables
 - CSV - containing the data

[illegible]

Variables include: Band-number, species, classification (whether the bird was impacted or a reference), various physiological markers (PAH, Heinz Body Counts, Reticulocytes)

Shape: 1257 rows (samples) X 67 columns (variables)*
* prior to cleaning the data

Central Hypothesis:

Oil-spill impacted birds have different levels of physiological markers compared to non-impacted (reference) birds.

Questions and Hypotheses

Questions:

1. Does an increase of PAHs (ng/mL) in the blood indicate that a bird was impacted by an oil spill?
2. Are the reticulocyte counts between impacted and non-impacted birds different?
3. Are the metabolisms of impacted and non-impacted birds different, and can this be seen a physiological marker?

Hypotheses:

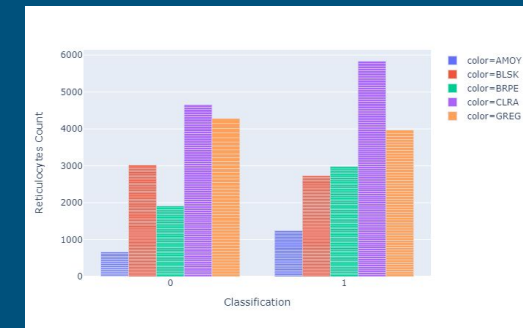
1. An increase in PAHs will indicate that a bird was impacted by oil pollution.
2. There is a decrease in reticulocyte counts for impacted birds compared to non-impacted birds.
3. There is a decrease in the metabolic marker, uric acid, in impacted birds compared to non-impacted birds.

Supervised Machine Learning Algorithms

- Tested 5 different algorithms:
 - K-Nearest Neighbors-based classification (KNeighborsClassifier)
 - Linear Support Vector Classification (LinearSVC)
 - Gaussian Naive Bayes (GaussianNB)
 - Decision Tree classifier (DecisionTreeClassifier)
 - Logistic Regression classifier (LogisticRegression)
- **Used Decision Tree classifier**

Results overview

- **Data Cleaning**
 - Dropped rows that had no useful data
 - Replaced empty/no record values with NaN
 - Cleaned up continuous data
 - Cleaned up categorical data
- **Data Exploration**
 - Explored relationship of data related to hypotheses



Results overview cont'd

- Feature Selection & Model Construction
 - Selected features using Decision Tree Classifier-based feature selection
 - Tested proposed hypotheses using aggregation by species and independent samples t-test
- Model Evaluation & Optimization
 - Scaled features
 - Cross validation approach *and* percentage-split approach
 - Evaluated accuracy on selected & scaled features on 5 different classifiers
 - Hyperparameter tuned all 5 models
- Model Testing
 - Tested our best scoring model (Decision Tree Classifier) on our test set
 - Used confusion matrix and classification report to report metrics of our selected model

Hypothesis Discussion

- Null hypothesis: means are not significantly different between reference and impacted birds (by species) for each feature we tested
 - PAH (ng/mL)
 - Reticulocytes Count
 - Uric Acid (mg/dL)
- Alternative hypothesis: means are significantly different between reference and impacted birds (by species) for each feature we tested
- We failed to reject all 3 of our null hypotheses, as for each hypothesis, not all the species showed a significant difference between means of the reference and impacted groups

ML Discussion

- Of the 5 classifiers that were evaluated, we ended up using the Decision Tree Classifier due to its high accuracy across both test/train split and cross validation metrics, as well as the fact that we believed it had the least likelihood of underfitting.
- There was no benefit from performing hyperparameter tuning on the Decision Tree Classifier, and even though other models did benefit from hyperparameter tuning, they still performed worse than the Decision Tree model.

ML Discussion

- The results from the classification report constructed from running the model against the testing data indicated that our model correctly labels reference birds 90% of the time and impacted birds 94% of the time.
- These scores could be improved with more data--this could either be in the form of more records (rows) in the dataset, or it could be more complete data in columns that we elected not to use as features for our model due to issues with the existing data in those columns.

Confusion Matrix

	Predicted reference	Predicted impacted
Actual reference	60	5
Actual impacted	7	74

Classification Report

	precision	recall	f1-score	support
0	0.90	0.92	0.91	65
1	0.94	0.91	0.92	81
accuracy			0.92	146
macro avg	0.92	0.92	0.92	146
weighted avg	0.92	0.92	0.92	146



Thank you

