# EVOS data project - Analysis

*Jessica Couture*

*10/22/2017*
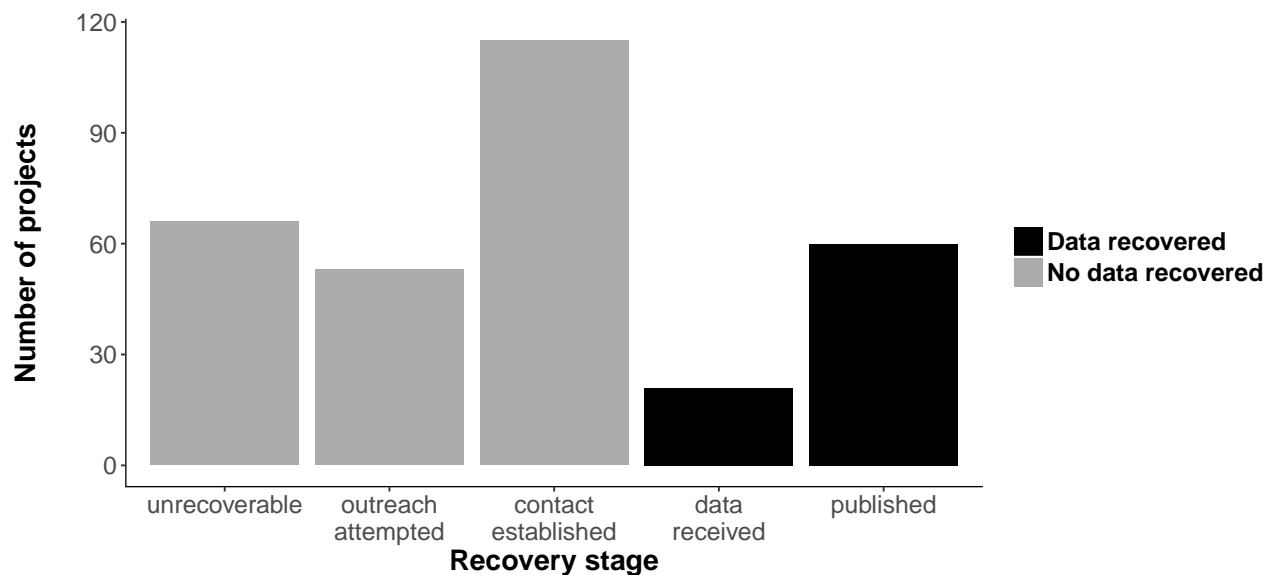
## Contents

## Background

This analysis accompanies publication the manuscript "How effective are funder imposed data sharing requirements?", assessing results of a two year data archiving effort by a group of researchers and students at the National Center for Ecological Analysis and Synthesis at UC Santa Barbara. The Exxon Valdez Oil Spill Trustee Council (EVOSTC) was formed following the Exxon Valdez oil spill in Alaska in 1989. Since then, the EVOSTC has funded hundreds of projects and in 2012 we initiated an effort to recover and archive the data collected through these EVOSTC funded projects. The recovery effort spanned two years.

For this paper we ask 3 main questions about the data collected from the Exxon Valdez Oil Spill Trustee Council funded projects:

1. Twenty-five years after the EVOS, for how many projects funded by EVOSTC can we collect data?
2. Are there differences in data reporting based on characteristics of the data?

- Research field
- Sector of researching body
- Year data projects ended

3. Which of these characteristics are most *important* in determining if a dataset will be successfully recovered and how do the important characteristics influence the output (success)?

## 1. Project Status Reporting

Twenty-five years after the EVOS, for how many projects funded by EVOSTC can we collect data?

```
# Number of projects requested
sum(overall)
```

```
## [1] 315
```

```
# percent successful
sum(overall[c("Published","SentData")])/sum(overall)
```

```
## [1] 0.2571429
```

---

## 2. Are there differences in data reporting based on data characteristics?

**Chi-square tests for each variable (characteristic)**

**Research field**

**Are there certain research fields that are more likely to make data available than others?**

Chi-squared test for equal proportions between research fields

```
##
##  Pearson's Chi-squared test
##
## data:  bioProps
## X-squared = 25.58, df = 9, p-value = 0.002392
```

**SIGNIFICANT: Reject the H0 that there are no differences in recovery in different research fields**

---

**Research sector**

**Test for equal propportions between PI's home institution type/sector:**

government and private split: govFed, govState, nonProf, forProf

```
##
##  Pearson's Chi-squared test
##
## data:  subSecProps
## X-squared = 2.6061, df = 5, p-value = 0.7604
```

**NOT SIGNIFICANT: there are no differences in recovery in different sectors**

---

**Year project ended**

**Test for equal proportions between years:**

Chi-squared test for equal proportions between age of data

```
##
##  Pearson's Chi-squared test
##
## data:  tempProps
## X-squared = 30.577, df = 21, p-value = 0.08099
```

**NOT SIGNIFICANT: there are no differences in recovery based on when data were collected**

---

## 3. Which characteristics are most important in determining if a dataset will be successfully recovered?

We use the "party" package in R to run a random forests analysis to determine which variables are most important. I use the same model as the glm, then create a classification tree below to show how the important variables influence the outcome. This package is be better than randomForest when independent variables are different types (Strobl et al. 2009)

**Random forests**

```
rslt2$dataType<-as.factor(rslt2$dataType)
partyForBio<-cforest(statSucc~agSubGrp+end+dataType,data=rslt2,controls = cforest_unbiased(mtry = 2, nt:
varimp(partyForBio)
```

```
##     agSubGrp          end      dataType
## -0.002606187   0.005361820   0.014385360
```

**Based on these results the most important variable in determining the outcome is research field**

---

## How do the important characteristics influence the output (success)?
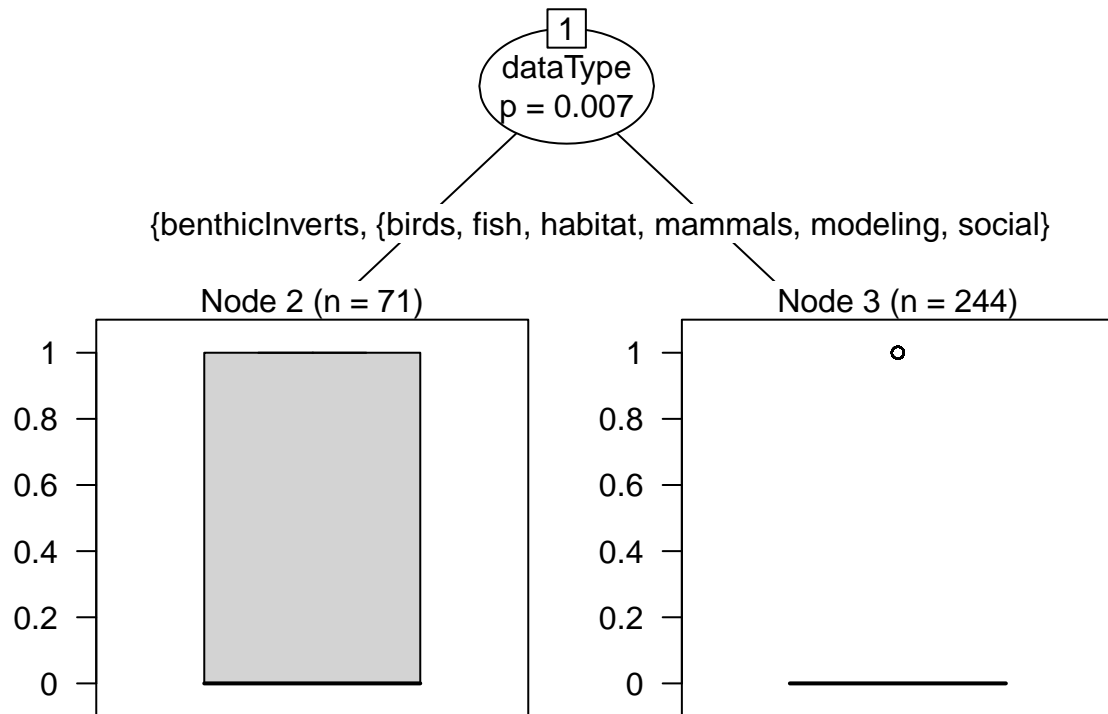
**Classification tree**

Here we run one iteration of the forest analysis above to display which variables whithin each classification determine positive or negative results.

```
partreeBio<-ctree(statSucc~agSubGrp+end+dataType,data=rslt2)

partreeBio

##
##   Conditional inference tree with 2 terminal nodes
##
## Response:  statSucc
## Inputs:  agSubGrp, end, dataType
## Number of observations:  315
##
## 1) dataType == {benthicInverts, oil, physical, plankton}; criterion = 0.993, statistic = 25.499
##   2)*  weights = 71
## 1) dataType == {birds, fish, habitat, mammals, modeling, social}
##   3)*  weights = 244
```

```
plot(partreeBio)
```



*Birds, fish, habitat, mammal* and *modeling data* result in negative results. *Benthic invertebrates, plankton, oil,* and *physical data* result in positive results.