

Home Work 1

Principles of Statistical Data Analysis

Jiacheng Shen

Sajjad Ayashm

2024-10-22

Question 1: Data Exploration

In this section, we explore the data using summary statistics and visualizations to understand the distribution of ant abundance and soil moisture.

```
load("ants.Rdata")
summary(ants$abundance)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      0.00   6.75   35.50   86.31   99.25   767.00
```

```
summary(ants$moisture)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##     12.00   15.40   17.95   17.52   19.90   21.50
```

```
var(ants$abundance)
```

```
## [1] 19882.49
```

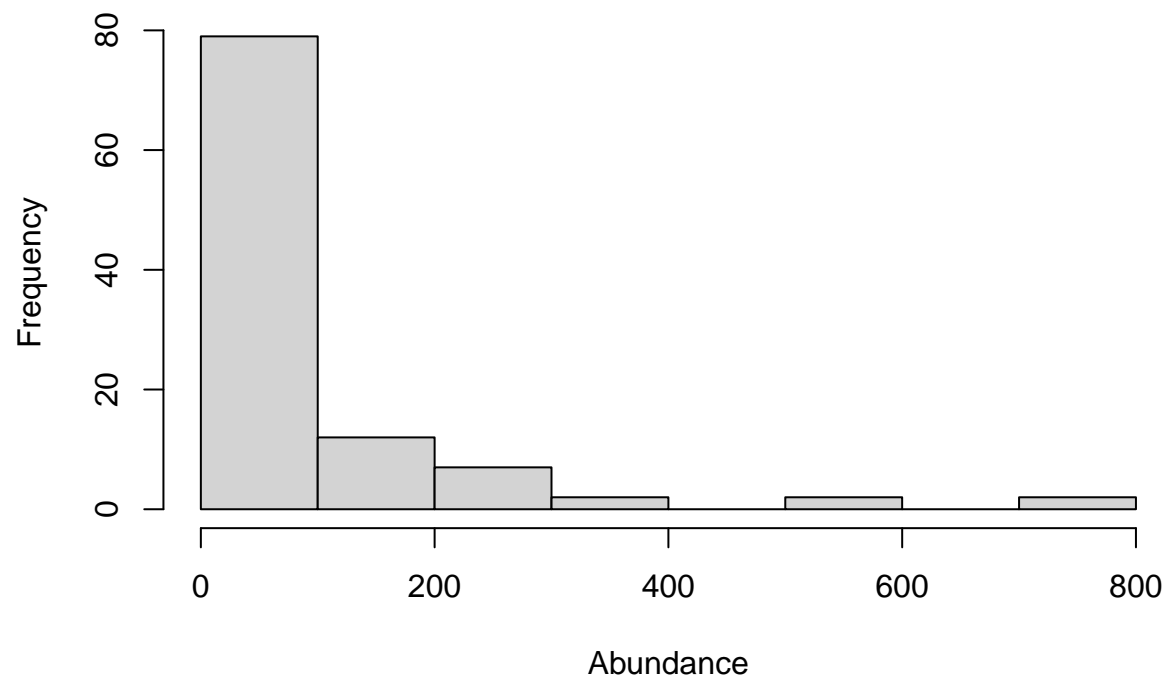
```
var(ants$moisture)
```

```
## [1] 7.083322
```

Visualize the distribution of ant abundance

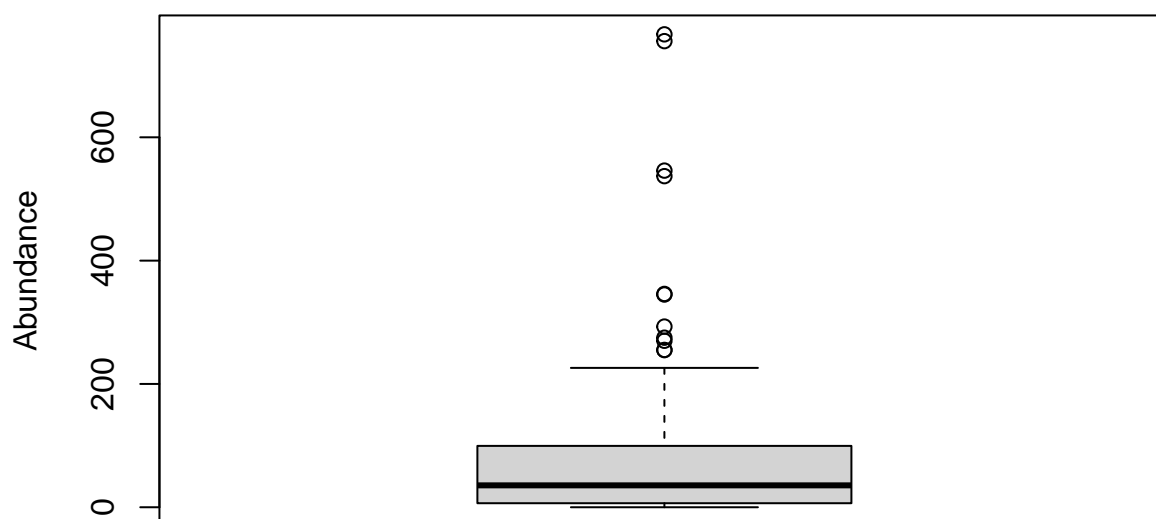
```
hist(ants$abundance, main = "Histogram of Ant Abundance", xlab = "Abundance", ylab = "Frequency")
```

Histogram of Ant Abundance



```
boxplot(ants$abundance, main = "Boxplot of Ant Abundance", ylab = "Abundance")
```

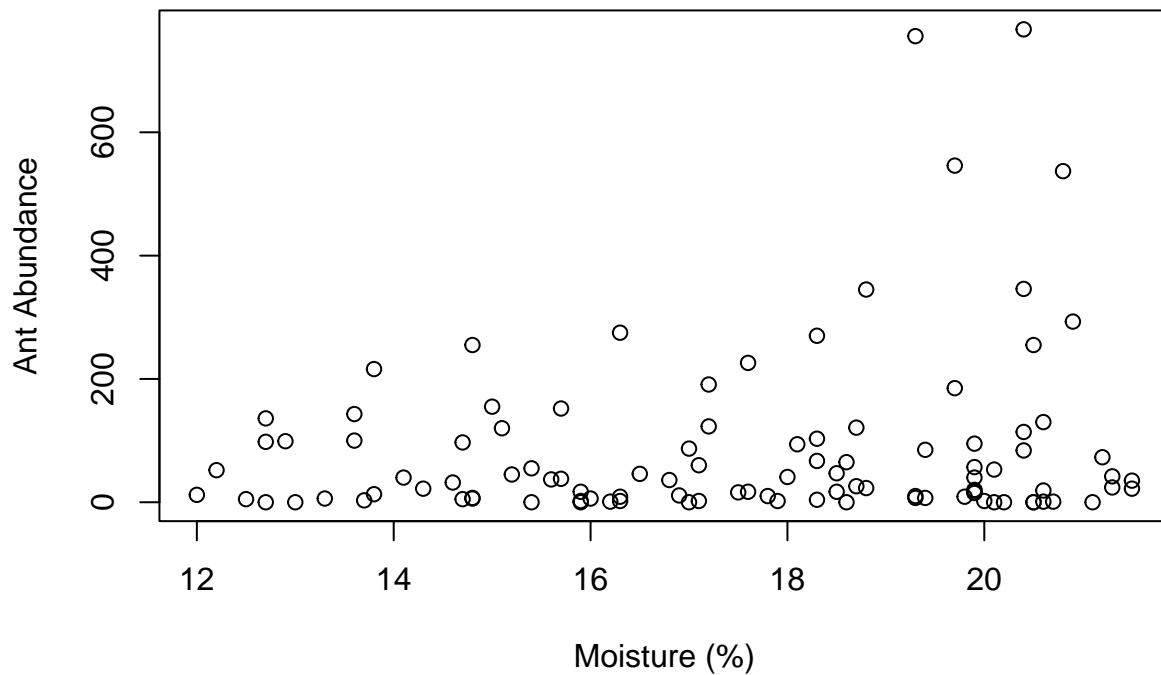
Boxplot of Ant Abundance



Scatter plot to explore the relationship between moisture and ant abundance

```
plot(ants$moisture, ants$abundance,  
     main = "Scatter Plot of Moisture vs Ant Abundance",  
     xlab = "Moisture (%)", ylab = "Ant Abundance")
```

Scatter Plot of Moisture vs Ant Abundance



Question 2: Likelihood and Log-Likelihood Functions

Likelihood Function

To define the likelihood function in R, we use a function that computes the product of probabilities for all observations.

```
# Define the likelihood function
likelihood_function <- function(params, abundance, moisture) {
  beta_0 <- params[1] # The intercept parameter
  beta_1 <- params[2] # The coefficient for moisture
  phi <- params[3]     # The overdispersion parameter

  # Calculate the mean (mu) for each observation
  mu <- exp(beta_0 + beta_1 * moisture)

  # Compute the negative binomial likelihood for each observation
  likelihoods <- dnbinom(abundance, size = 1/phi, mu = mu, log = FALSE)

  # Calculate the total likelihood by multiplying individual likelihoods
  total_likelihood <- prod(likelihoods)

  return(total_likelihood)
}
```

```
initial_params <- c(1, 0.1, 0.5) # Example starting values for beta_0, beta_1, and phi
likelihood_value <- likelihood_function(initial_params, ants$abundance, ants$moisture)
likelihood_value
```

```
## [1] 0
```

Note: If we use the set of parameters described in the example above, the likelihood value results 0. One possible reason is that the number is too small to exhibit. That's why we derive the following log likelihood function.

Log-Likelihood Function

The log-likelihood function computes the sum of the log-probabilities instead of multiplying the probabilities directly, providing more numerical stability.

```
# Define the log-likelihood function
log_likelihood <- function(params, abundance, moisture) {
  beta_0 <- params[1] # The intercept parameter
  beta_1 <- params[2] # The coefficient for moisture
  phi <- params[3]    # The overdispersion parameter

  # Calculate the mean (mu) for each observation
  mu <- exp(beta_0 + beta_1 * moisture)

  # Compute the log-likelihood for each observation
  log_lik <- sum(
    lgamma(abundance + 1 / phi) - lgamma(1 / phi) - lgamma(abundance + 1) +
    (1 / phi) * log(1 / (1 + mu * phi)) +
    abundance * log(mu * phi / (1 + mu * phi))
  )

  return(log_lik) # Return the log-likelihood
}

initial_params <- c(1, 0.1, 0.5) # Example starting values for beta_0, beta_1, and phi
log_likelihood_value <- log_likelihood(initial_params, ants$abundance, ants$moisture)
log_likelihood_value
```

```
## [1] -1148.303
```

Using the same set of parameters, the log likelihood values result -1148.3025463“, which makes more sense.

Question 3: Estimating Equation for β_1

We are tasked with deriving the estimating equation for β_1 , which involves maximizing the log-likelihood function. The maximum likelihood estimate (MLE) is found by solving this equation using numerical optimization.

Step 1: Negative Log-Likelihood Function

Since the `optim()` function in R minimizes functions, we use the negative log-likelihood function for the optimization process. The negative log-likelihood is simply the negative of the log-likelihood function.

```
# Define the negative log-likelihood function for optimization
neg_log_likelihood <- function(params) {
  -log_likelihood(params, ants$abundance, ants$moisture)
}
```

Step 2: Perform Optimization to Estimate Parameters

We now use the `optim()` function to find the maximum likelihood estimates (MLE) for β_0 , β_1 , and ϕ . The initial guesses for the parameters are provided, and the optimization is performed using the BFGS method.

```
# Initial parameter guesses: beta_0 = 1, beta_1 = 0.1, phi = 0.5
initial_params <- c(1, 0.1, 0.5)

# Perform optimization to estimate beta_0, beta_1, and phi
optim_result <- optim(initial_params, neg_log_likelihood, method = "BFGS")
```

```
## Warning in log(1/(1 + mu * phi)): NaNs produced
```

```
# Extract the estimated parameters
estimated_params <- optim_result$par

# Assign the estimates to beta_0_hat, beta_1_hat, and phi_hat
beta_0_hat <- estimated_params[1]
beta_1_hat <- estimated_params[2]
phi_hat <- estimated_params[3]

# Display the estimated parameters
beta_0_hat
```

```
## [1] 2.504213
```

```
beta_1_hat
```

```
## [1] 0.1091764
```

```
phi_hat
```

```
## [1] 2.288232
```

Conclusion

The estimated parameters β_0 , β_1 , and ϕ are calculated using numerical optimization. The MLE of β_1 is particularly important for the analysis, as it describes the relationship between moisture and ant abundance.

These values will be used for further analysis in the next questions.