

Final Project

RJ Cass

Table of contents I

- 1 Abstract
- 2 Introduction
- 3 Exploratory Data Analysis
- 4 Methodology
- 5 Results
- 6 Conclusion
- 7 Appendix

Abstract

Tulips form a significant portion of the exports from the Netherlands. We want to understand what changes can be made to ensure tulip growth is not impacted by the changing climate.

We used a logistic model and found:

- Effect of chilling times on germination is dependent on species (some improve with longer chilling, some get worse with chilling)
- Ideal chilling time for each population (for most, range of 8-10 weeks)
- Predicted impact of the chilling time decreasing from 10 weeks to 9 (some do marginally better, but some do much worse)

Context

Tulip Production:

- Tulip products form 25% of agricultural exports from the Netherlands
- Changing climate puts the tulip industry at risk
- Want to understand how to adapt to these changes and protect the industry

Dataset of sample tulip growth populations:

- Year they were grown
- Number of weeks the bulbs were chilled
- Whether or not the bulb germinated (binary: yes/no)
- Indices (can be removed from dataset)

Questions of Interest

We want to use the provided data to answer the following questions:

- ① What is the effect of chilling time for the different species of tulips? Is it the same across the species? Which species are the same/different?
- ② Is there an ideal chilling time for each species? If so, is it the same for all species?
- ③ Given climate change conditions, winters are expected to decrease from 10 to 9 weeks in the coming few years. What effect will this decrease in chilling time have on the probability of germination for each species? Is it the same for all species?

EDA - Population 12

None of population 12 germinated. We removed it from the dataset to not dilute the rest of the data

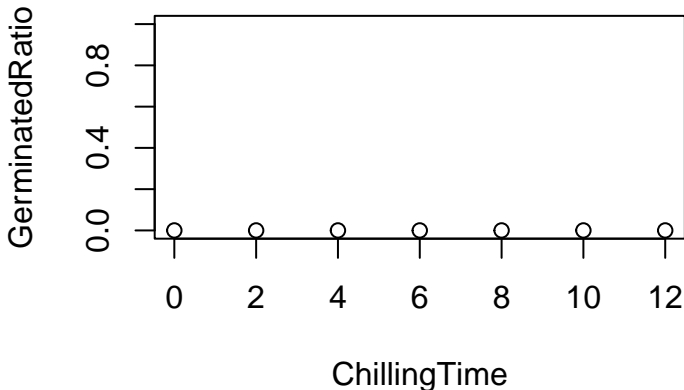


Figure 1: *Population 12 had a 0% germination rate across all chilling times*

EDA - Year

- Each population was tested in only 1 year (ie. no crossing with different years having an effect on one population)
- Physically, given testing conditions, we don't expect year to have an impact on germination
- In variable selection, Year was not important ($p > .05$)
- Removing year from models

EDA - Interactions - Population vs. Chilling

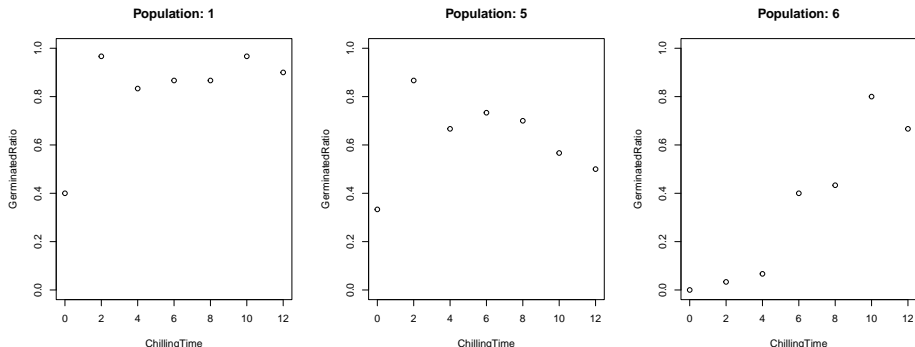


Figure 2: *Sample interaction plots of population and chilling. Some behave similarly, others are vastly different. Note the non-linearity of some populations.*

EDA - Summary

Removing from dataset:

- Population 12
- 'Year' variate

Need to account for the following:

- Interaction between Chilling Time and Population
 - If not included, resulting model will not capture the full impact of each variate
- Non-linearity of relationship between chilling time and germination rate
 - If not included, model will not represent the correct relationship of this variate
- Binary Response
 - If not included, model will not represent correct behavior of outcome variable

Proposed Models - 1

Logistic Model

$$Y_n = \ln\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 X_{pop} + \beta_c \text{poly}(\text{ChillingTime}) + \beta_i (X_{pop} * \text{poly}(\text{ChillingTime})) + \epsilon \quad (1)$$

Strengths:

- The concept of logistic (change in log-odds) is reasonably interpretable
- Accounts for interactions of population and year on ChillingTime
- Accounts for non-linearity of ChillingTime

Weaknesses:

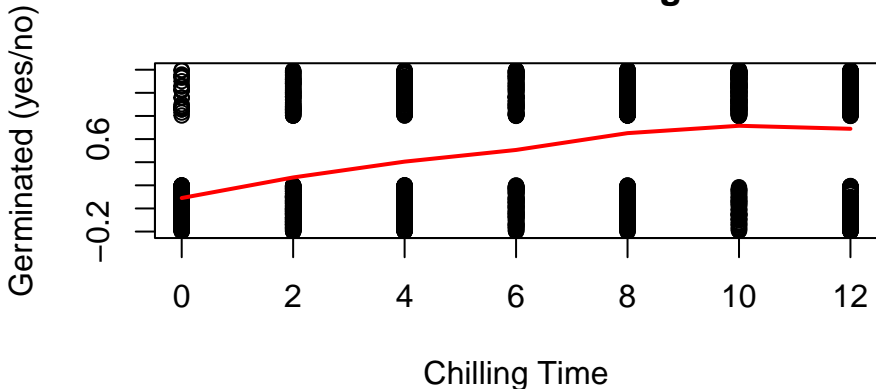
- Using 2nd degree polynomial loses interpretability

Proposed Models - 1 - Cont'd

Assumptions - Independence, Monotonicity

- Independence: Assumed due to the design of the experiment
- Monotonicity

Germination vs. Chilling Time



Proposed Models - 2

Random Forest

Strengths:

- Relatively explainable (lots of trees, each tree gets a vote, average the votes, compare to cutoff)
- Inherently explores non-linearity/interactions. Particularly good at the 'step' functions
- No inherent assumptions (besides the data being 'good')

Weaknesses:

- Less direct at answering questions (would require much more computation)

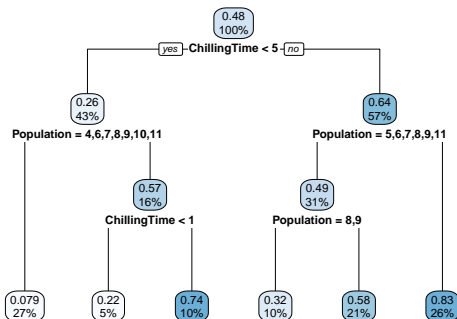


Figure 3: *Sample tree with $cp = .01$*

Model Evaluation/Selection

Model	In-Sample Accuracy	Out-of-Sample Accuracy
Logistic	0.7927	0.7662
Forest	0.8475	0.7749

- Random Forest does marginally better in accuracy
- Logistic model is much more interpretable
- Logistic model more clearly answers research questions

Will use the Logistic model to answer research questions

- Coefficients provided in Appendix
- Highlight a coefficient:
 - ChillingTime - Poly1: 30.68
 - ChillingTime - Poly2: -25.36

Effect of Chilling Time

3 types:

- Step-functions:
above a certain value appears to be steady rate
- Increase up to a value, then decrease
- Primarily decreasing

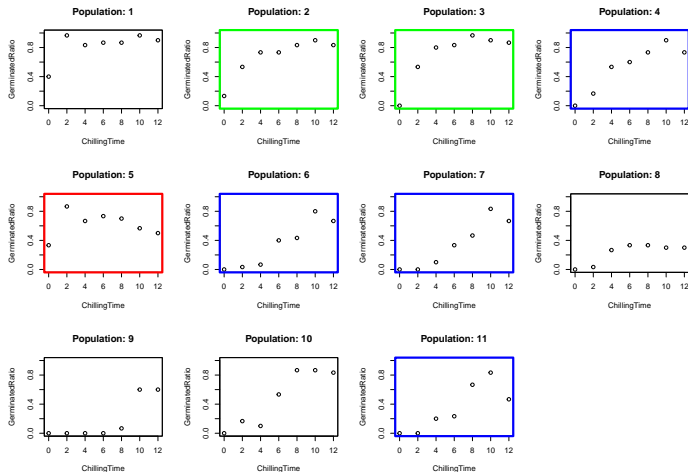
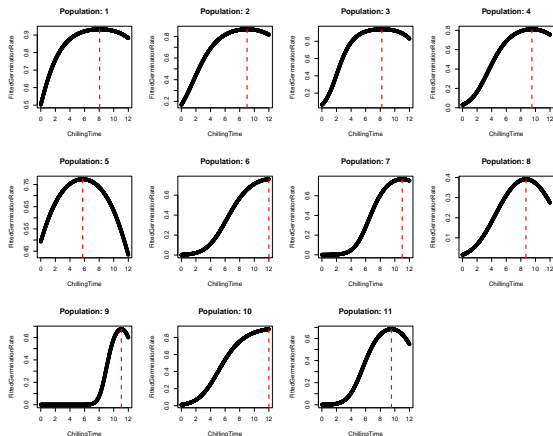


Figure 4: *Plots of chilling time on germination by population*

Ideal Chilling Time

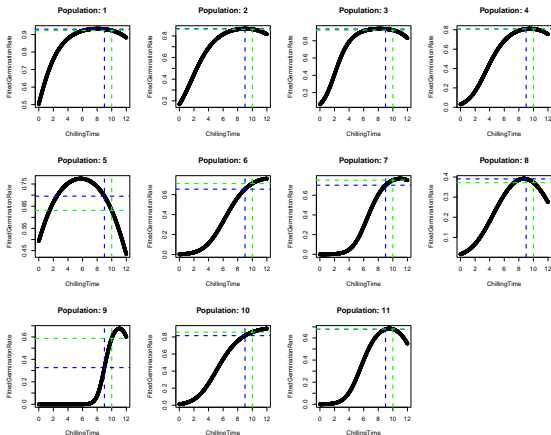
Used model to identify ideal values:



Population	IdealChillWeeks
1	8.072
2	9.009
3	8.216
4	9.514
5	5.742
6	12.000
7	11.003
8	8.697
9	11.039
10	12.000
11	9.526

Effect of Decrease in Chilling Time

Used model to calculate difference in Germination rate at 9 vs. 10 weeks:



Population	GerminationDiff
1	0.77%
2	0.5%
3	1.36%
4	-0.03%
5	6.41%
6	-5.69%
7	-5.07%
8	1.85%
9	-26.03%
10	-4.09%
11	-0.11%

Summary

Used the provided tulips data to:

- Identify that impact of chilling time on germination rate depends on species (most increase with chilling time, some decrease)
- Identify the ideal chilling time for each species (for most, range of 8-10)
- Predict the impact of chilling time decreasing from 10 to 9 weeks on germination rate (mixture of increases and decreases, largest is a decrease)

Next Steps

To improve our understanding of tulip chilling time on germination rate we suggest:

- Increased resolution of chilling times (ie. get samples of every week, maybe even per day)
- With risk of rising sea levels, test different humidity levels, soil saturation, etc. (effects of more moisture)

Table of Coefficients I

	Estimate	Lower2.5	Upper97.5
(Intercept)	1.812	1.386	2.296
Population2	-1.020	-1.613	-0.456
Population3	-0.764	-1.394	-0.152
Population4	-1.906	-2.529	-1.327
Population5	-1.197	-1.766	-0.665
Population6	-3.271	-4.248	-2.481
Population7	-4.131	-5.663	-2.979
Population8	-3.392	-4.122	-2.756
Population9	-15.604	-26.123	-7.778
Population10	-2.231	-2.970	-1.563
Population11	-3.756	-4.860	-2.867
poly(ChillingTime, degree = 2)1	30.680	13.963	49.718
poly(ChillingTime, degree = 2)2	-25.362	-44.436	-6.778
Population2:poly(ChillingTime, degree = 2)1	16.467	-8.364	40.798

Table of Coefficients II

Population3:poly(ChillingTime, degree = 2)1	34.300	7.630	61.369
Population4:poly(ChillingTime, degree = 2)1	39.333	12.039	67.691
Population5:poly(ChillingTime, degree = 2)1	-34.167	-57.191	-12.660
Population6:poly(ChillingTime, degree = 2)1	76.815	37.558	126.867
Population7:poly(ChillingTime, degree = 2)1	125.457	65.052	207.282
Population8:poly(ChillingTime, degree = 2)1	17.467	-13.289	52.606
Population9:poly(ChillingTime, degree = 2)1	613.209	236.097	1107.435
Population10:poly(ChillingTime, degree = 2)1	69.645	37.643	105.913
Population11:poly(ChillingTime, degree = 2)1	92.153	45.449	151.018

Table of Coefficients III

Population2:poly(ChillingTime, degree = 2)2	-1.597	-26.898	23.674
Population3:poly(ChillingTime, degree = 2)2	-25.097	-54.341	3.003
Population4:poly(ChillingTime, degree = 2)2	-8.875	-36.372	17.887
Population5:poly(ChillingTime, degree = 2)2	1.278	-22.021	24.747
Population6:poly(ChillingTime, degree = 2)2	-3.846	-41.257	29.982
Population7:poly(ChillingTime, degree = 2)2	-28.420	-78.899	14.819
Population8:poly(ChillingTime, degree = 2)2	-5.347	-35.504	22.795
Population9:poly(ChillingTime, degree = 2)2	-	-	-39.108
Population10:poly(ChillingTime, degree = 2)2	194.673	384.036	
Population10:poly(ChillingTime, degree = 2)2	2.897	-29.602	33.780

Table of Coefficients IV

Population11:poly(ChillingTime, degree = 2)	-34.579	-74.750	1.471
---	---------	---------	-------
