

# Remote Sensing for Forest Recovery: Proposal

Benjamin Frizzell      Zanan Pech      Mavis Wong      Hui Tang

2025-05-05

## Table of contents

<b>1</b>	<b>Abstract</b>	<b>2</b>
<b>2</b>	<b>Introduction</b>	<b>2</b>
<b>3</b>	<b>Data Description</b>	<b>2</b>
3.1	Site Features . . . . .	3
3.2	Spectral Indices . . . . .	3
3.3	Out-of-range Values . . . . .	3
3.4	Previous Land Use and Species Composition . . . . .	4
3.5	Species Type . . . . .	4
3.6	Trends and Seasonality . . . . .	4
3.7	Collinearity . . . . .	4
<b>4</b>	<b>Data Engineering</b>	<b>7</b>
4.1	Train-Test Splitting . . . . .	7
4.2	Missing Data . . . . .	7
4.3	Feature Selection and Engineering . . . . .	8
4.4	Data Pivoting . . . . .	8
4.5	Conversion to Classifier Problem . . . . .	8
<b>5</b>	<b>Modelling Techniques</b>	<b>9</b>
5.1	Logistic Regression . . . . .	9
5.2	Aggregated Models: Random Forest and Gradient Boosting . . . . .	9
5.3	Sequential Deep Learning Models: LSTMs, GRUs, etc. . . . .	9
<b>6</b>	<b>Success Criteria</b>	<b>9</b>
<b>7</b>	<b>Timeline</b>	<b>9</b>

<b>8 Conclusion</b>	<b>10</b>
<b>References</b>	<b>10</b>

## 1 Abstract

Monitoring afforestation progress across hundreds of sites is a significant challenge. This project explores using site-level data and satellite-derived vegetation indices from Canada's 2 Billion Trees program to build machine learning models that predict tree survival over seven years. We will train models ranging from logistic regression to more advanced techniques like random forests, gradient boosting, and deep learning (RNNs, LSTMs) to capture temporal patterns. The goal is to evaluate predictive features and modeling strategies that enable scalable, cost-effective monitoring of afforestation efforts.

## 2 Introduction

Afforestation is crucial for combating climate change and supporting biodiversity. Trees help purify the atmosphere and provide food, nutrients, and habitats for countless species (Government of Canada 2023).

To advance these goals, the Canadian government launched the **2 Billion Trees** program, aiming to plant two billion trees nationwide by 2031 with support for provinces and territories (Natural Resources Canada 2021).

Monitoring survival at scale is a major challenge. While remote sensing enables broad environmental monitoring, detecting small or newly planted trees remains difficult due to limited canopy coverage and weak spectral signals (University of British Columbia Master of Data Science Program 2025).

This study aims to address two key questions: - Can remote sensing reduce the need for physical site visits? - Which modelling approach is most effective, and how soon after planting can survival be reliably predicted?

## 3 Data Description

The dataset used in this study includes field-measured survival rates of afforested sites collected by Forest Ontario (University of British Columbia Master of Data Science Program 2025), as well as satellite data products from the Harmonized Landsat Sentinel-2 (HLS) project (USGS 2024).

### 3.1 Site Features

Column Name	Description
ID	Site ID
PixelID	Pixel ID
Area_ha	Area of Site
Season	Planting Year
PlantDt	Planting Date
prevUse	Previous Land Use of Site
Planted	Number of Trees Planted
SpcsCmp	Species Composition of Site
Type	Species Type of Site
SrvvR_1, ..., SrvvR_7	Field Measured Survival Rate at Year 1-7 (Target)
AsssD_1, ..., AssD_7	Date of Field Survival Rate Measurement
NmbrPIO	Number of Trees Originally Planted
NmbrPIR	Number of Trees Replanted
NmbrPIT	Total Number of Trees Planted
ImgDate	Image Date of the Remote Sensing Data
Year	Image Year of Remote Sensing Data
DOY	Day of Year of the Remote Sensing Data

### 3.2 Spectral Indices

Type	Index	Description
Vegetation Index	NDVI	Normalized Difference Vegetation Index
	SAVI	Soil-Adjusted Vegetation Index
	MSAVI	Modified Soil-Adjusted Vegetation Index
	EVI	Enhanced Vegetation Index
	EVI2	Two-band Enhanced Vegetation Index
Water Index	NDWI	Normalized Difference Water Index
Fire Index	NBR	Normalized Burn Ratio
Tasseled Cap Transformed	TCB	Tasseled Cap Brightness
	TCG	Tasseled Cap Greenness
	TCW	Tasseled Cap Wetness

### 3.3 Out-of-range Values

During EDA, we noticed out-of-range values in the vegetation indices and survival rates. With the exception of TCB, TCW and TCG, the vegetation indices should range between -1 to 1

(USGS; Mondal 2011; Singh et al. 2015). Survival rates should not exceed 100%. Out-of-range records will be removed from the dataset.

### 3.4 Previous Land Use and Species Composition

Class imbalance was observed in `SpcsCmp`. With over 300 categories in `SpcsCmp`, it would not be practical to use it as a predictor for our model.

Similarly, severe class imbalance was observed in `prevUse`, as such, it will also be excluded from our model.

### 3.5 Species Type

From `SpcsCmp`, we are able to impute the missing values in `Type`. Sites are classified as “Conifer” (“Deciduous”) if  $\geq 80\%$  of the species are “Softwood” (“Hardwood”); otherwise, it would be classified as “Mixed” (Canada 2025).

As shown in Figure 1, survival rates and spectral signals vary by species type, suggesting ‘Type’ could be a viable feature for our model.

### 3.6 Trends and Seasonality

From Figure 2, we observed clear seasonality in the spectral indices, where spectral signals peaked during the summer and dropped during winter. This would need to be accounted for during model development.

Figure 3 reveals a positive relationship between vegetation indices and tree age. Minimal changes in spectral signals are observed between ages 1 to 4, indicating potential difficulties in predicting the survival rates for younger trees.

### 3.7 Collinearity

As shown in Figure 4, there is strong collinearity between the vegetation indices, except for TCB. The strongest correlation between survival rate and vegetation indices was observed in Year 7.

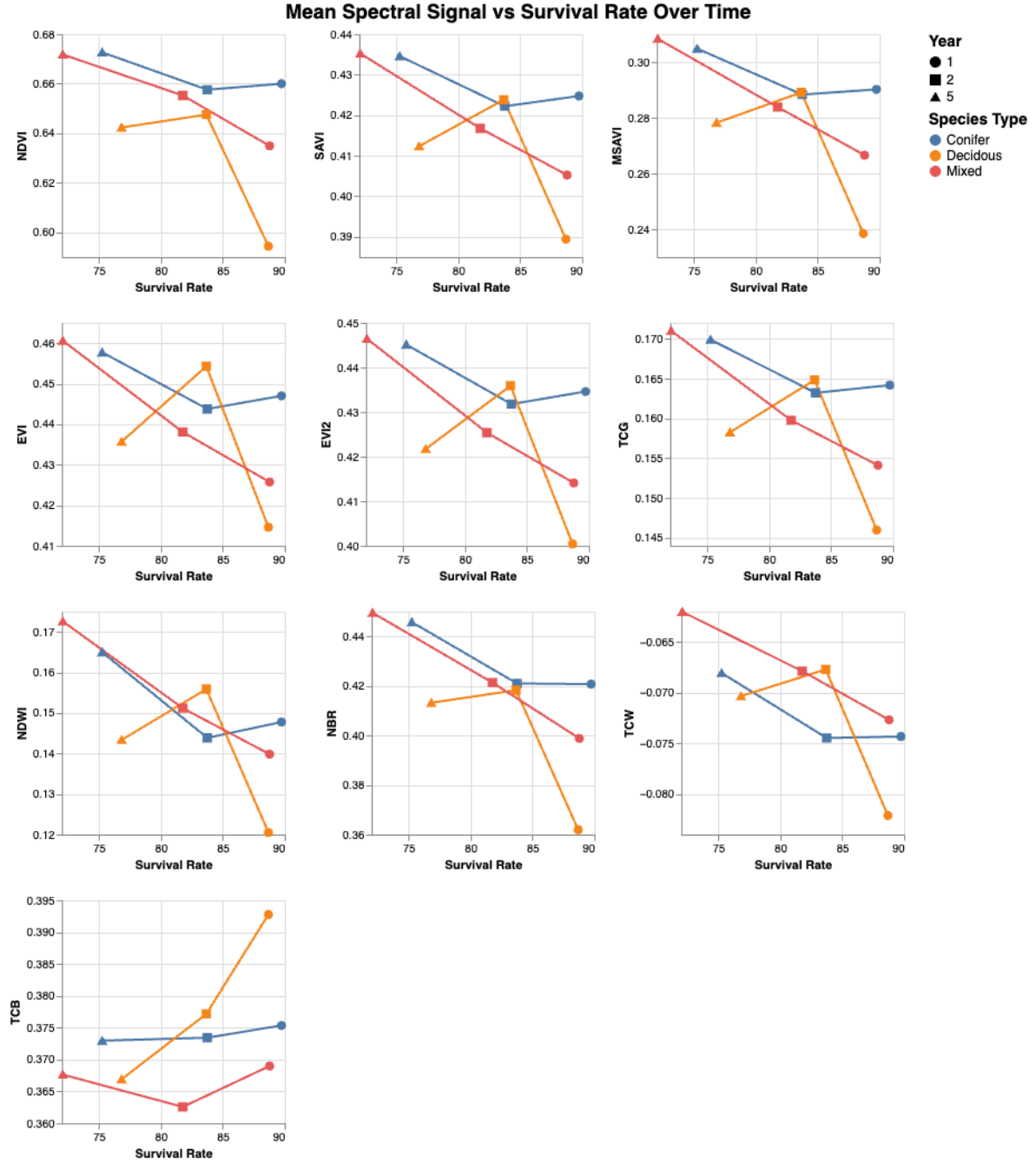


Figure 1: Plot showing mean survival rate and vegetation index signals for different species types across Years 1, 2 and 5. There is a significant difference in the relationship between survival rate and spectral signals for different species types. Conifers show a weaker signal response to changes in survival rate. Deciduous shows the strongest response during the first two years. Mixed type shows a linear relationship between survival rate and spectral signal.

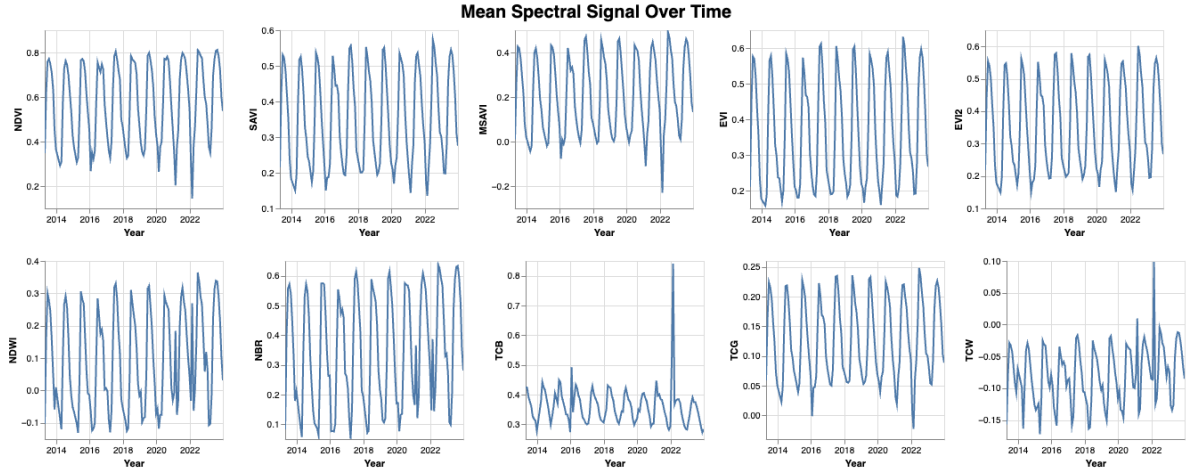


Figure 2: Plot showing seasonality in spectral indices.

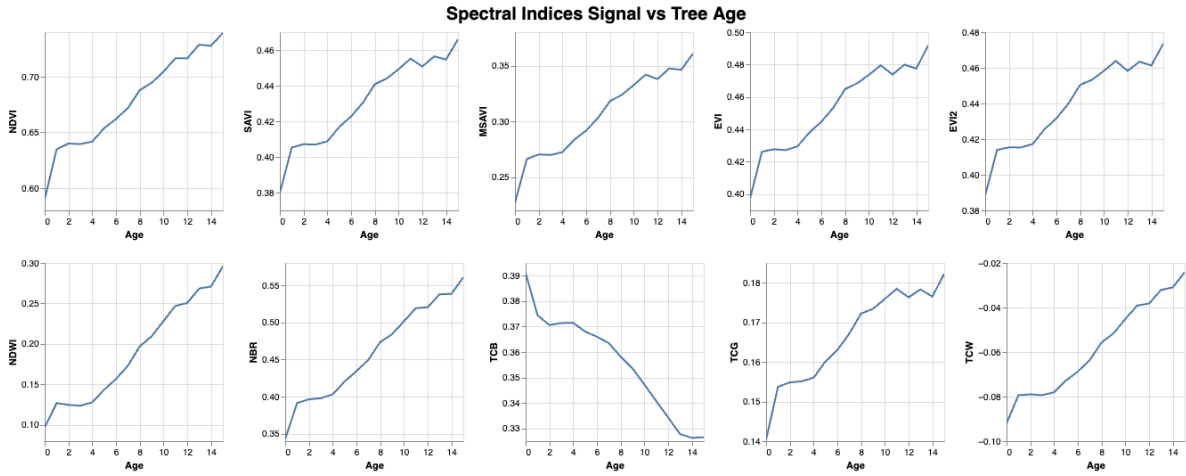


Figure 3: Plot showing mean spectral signal by tree age. With the exception of TCB, the spectral signal increases with age. A negative relationship was observed for TCB due to lower surface brightness from canopy cover.

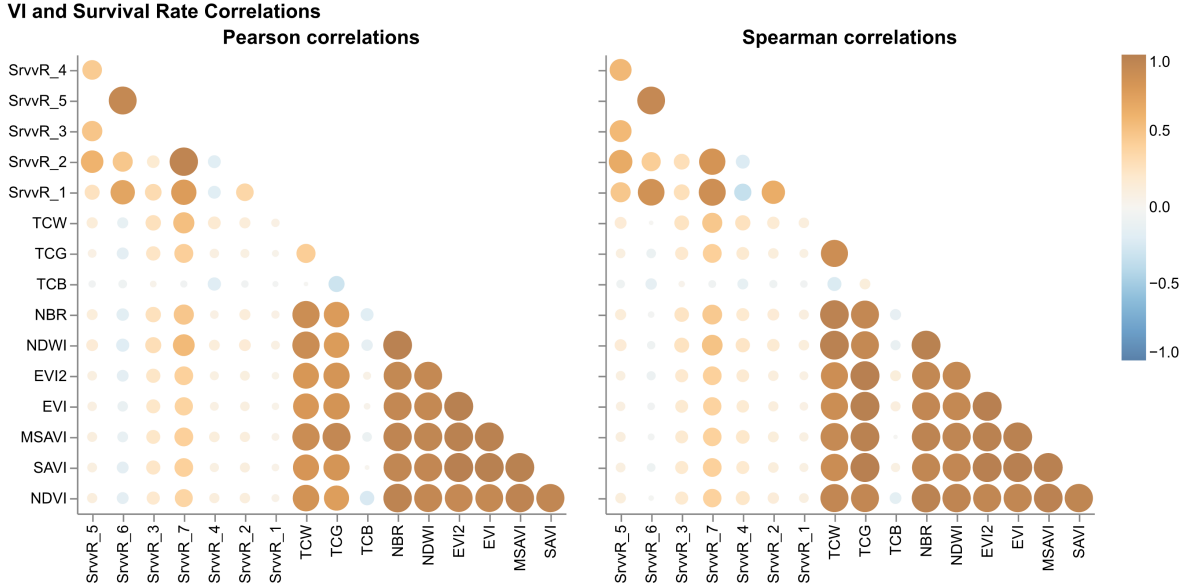


Figure 4: Correlation plot showing strong collinearity between vegetation indices.

## 4 Data Engineering

### 4.1 Train-Test Splitting

Splitting the dataset into training and testing subsets is necessary to prevent data leakage, but requires nuance due to its hierarchical structure. We perform the train-test split by **unique sites** to ensure pixels and time step records for a particular site appear in only one of the two subsets.

### 4.2 Missing Data

We find excessive missingness in the columns `PlantDt`, `Type`, `NmbrPIO`, `NmbrPIR`, and `NmbrPIT`. `PlantDt`, `NmbrPIO`, `NmbrPIR`, and `NmbrPIT` relate to sites where replanting has occurred, and can be removed as they are outside of the project scope, and `Type` can be fully imputed via string processing on the `SpcsCmp` column. There is a direct correspondence between missingness in survival rate and assessment time, allowing for easier tracking of temporal dependence across survival records.

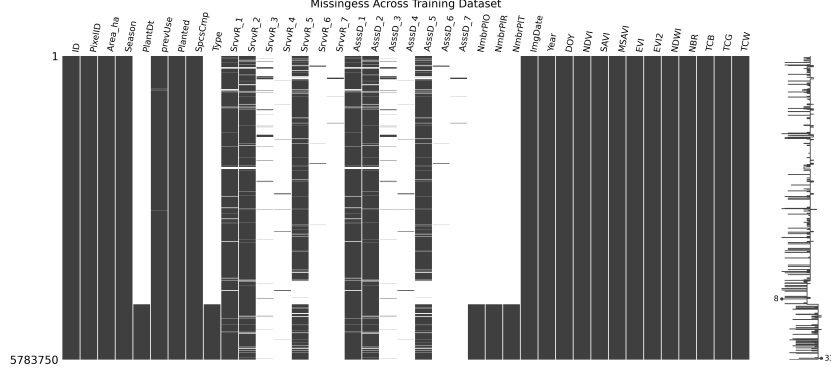


Figure 5: A plot visualizing missing record patterns across the training dataset. Grey-coloured records indicate rows that have recorded values.

### 4.3 Feature Selection and Engineering

As mentioned in Section 3, strong collinearity between vegetation indices indicates a need for feature selection. We will begin with **Recursive Feature Elimination**, due to its compatibility with nonlinear ML models (Pedregosa et al. 2011). We also propose the use of **Bayesian Model Averaging**, as it is suitable for an analysis of multiple competing models (Hoeting et al. 1998). Domain knowledge of the vegetation indices as given by Zeng et al. (2022) will also be leveraged in this process. Very little feature engineering will be performed, but we aim experiment with tree density (number of trees per unit area) as a predictor.

### 4.4 Data Pivoting

Most machine learning models require the input data to contain just one target column. We will pivot the seven target columns into one, keeping track of temporality using column names and assesment dates. We will then remove rows with missing survival rates, and with mismatching assesment and imaging dates.

### 4.5 Conversion to Classifier Problem

Since the survival rates are given as percentage proportions, they will be converted to binary classes to simplify the analysis and emphasize high-risk sites. Since most survival rates are within 70% - 100%, usefulness alongside class imbalance must be considered when deciding on a suitable threshold.



## 5 Modelling Techniques

### 5.1 Logistic Regression

Logistic regression will serve as a baseline to contrast performance against other, more complex models. It may also provide insight into feature importance through interpretable coefficients.

### 5.2 Aggregated Models: Random Forest and Gradient Boosting

Studies by Bergmüller and Vanderwel (2022) suggest that aggregated tree models may effectively utilize vegetation indices to predict tree mortality. Random Forests and GBMs complement each other well, as the former produces many ‘overfit’ trees in tandem to produce more accurate responses, whereas the latter iteratively produces ‘underfit’ trees which correct on mistakes made by the previous (Pedregosa et al. 2011).

### 5.3 Sequential Deep Learning Models: LSTMs, GRUs, etc.

If time permits, we may consider implementing sequential RNN-based deep learning models such as LSTMs and/or GRUs. These models have the benefit of capturing non-linear structure and may exploit changes in vegetation indices across time (Paszke et al. 2019).

## 6 Success Criteria

The success of this project will be evaluated based on the usefulness of the selected predictors and the performance of key evaluation metrics, including accuracy, log loss, F1 score, and ROC and PR curves. If our predictors are shown to contribute meaningfully to modeling tree survival rates, we aim to achieve at least 60% accuracy in correct predictions. Accuracy will serve as the primary metric for communicating results to general audiences—particularly government officials without data science expertise—as it provides a straightforward summary of model performance. However, given the imbalanced nature of our dataset, metrics such as log loss, F1 score, and the ROC and PR curves will be especially important for technical stakeholders and researchers within the federal government who are equipped to interpret more nuanced performance indicators.

## 7 Timeline

Date & Time	Deliverable	Description
<b>May 2, 2025</b>	Proposal Presentations	Group oral presentations
<b>May 6, 2025 12:00</b>	Proposal Report – Draft to Mentor	Ungraded draft
<b>May 9, 2025 17:00</b>	Final Proposal Report	Final version to partner & mentor
<b>June 9, 2025 16:00</b>	Data Product – Runnable Draft	Draft pipeline & code to mentor
<b>June 12–13, 2025 9:00–16:30</b>	Final Presentation	Group presentation of modelling approach & key results
<b>June 25, 2025 12:00</b>	Final Data Product & Report	Final pipeline, data product, and technical report

## 8 Conclusion

This project addresses the challenge of monitoring tree survival across hundreds of afforestation sites in Canada’s 2 Billion Trees program. By leveraging site-level data and satellite-derived spectral indices, we evaluate machine learning and deep learning models to predict survival over time. Our approach spans from interpretable models like logistic regression to advanced methods such as random forests, gradient boosting, and recurrent neural networks. Identifying key predictors is essential for building effective, scalable tools to support national afforestation efforts.

## References

- Bergmüller, Kai O, and Mark C Vanderwel. 2022. “Predicting Tree Mortality Using Spectral Indices Derived from Multispectral UAV Imagery.” *Remote Sensing* 14 (9): 2195.
- Canada, Natural Resources. 2025. “Forestry Glossary | Natural Resources Canada.” *Nrcan.gc.ca*. <https://cfs.nrcan.gc.ca/terms/read/782>.
- Government of Canada. 2023. “2 Billion Trees Commitment.” <https://www.canada.ca/en/campaign/2-billion-trees.html>.
- Hoeting, Jennifer A, David Madigan, Adrian E Raftery, and Chris T Volinsky. 1998. “Bayesian Model Averaging.” In *Proceedings of the AAAI Workshop on Integrating Multiple Learned Models*, 335:77–83. Citeseer.
- Mondal, Pinki. 2011. “Quantifying Surface Gradients with a 2-Band Enhanced Vegetation Index (EVI2).” *Ecological Indicators* 11 (3): 918–24.
- Natural Resources Canada. 2021. “2 Billion Trees Program.” <https://www.canada.ca/en/campaign/2-billion-trees/2-billion-trees-program.html>.

- Paszke, Adam, Sam Gross, Francesco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, et al. 2019. “PyTorch: An Imperative Style, High-Performance Deep Learning Library.” In *Advances in Neural Information Processing Systems*. Vol. 32. Curran Associates, Inc.
- Pedregosa, Fabian, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, et al. 2011. “Scikit-Learn: Machine Learning in Python.” *Journal of Machine Learning Research* 12: 2825–30.
- Singh, Kanwar Vivek, Raj Setia, Shashikanta Sahoo, Avinash Prasad, and Brijendra Pateriya. 2015. “Evaluation of NDWI and MNDWI for Assessment of Waterlogging by Integrating Digital Elevation Model and Groundwater Level.” *Geocarto International* 30 (6): 650–61.
- University of British Columbia Master of Data Science Program. 2025. “Remote Sensing for Forest Recovery.” [https://pages.github.ubc.ca/mds-2024-25/DSCI\\_591\\_capstone-proj\\_students/proposals/Remote\\_Sensing\\_for\\_Forest\\_Recovery.html](https://pages.github.ubc.ca/mds-2024-25/DSCI_591_capstone-proj_students/proposals/Remote_Sensing_for_Forest_Recovery.html).
- USGS. “Landsat Surface Reflectance-Derived Spectral Indices.” [http://w3techs.com/technologies/overview/content\\_language/all](http://w3techs.com/technologies/overview/content_language/all).
- . 2024. “HLS Overview.” *Usgs.gov*. <https://lpdaac.usgs.gov/data/get-started-data/collection-overview/missions/harmonized-landsat-sentinel-2-hls-overview/>.
- Zeng, Yelu, Dalei Hao, Alfredo Huete, Benjamin Dechant, Joe Berry, Jing M Chen, Joanna Joiner, et al. 2022. “Optical Vegetation Indices for Monitoring Terrestrial Ecosystems Globally.” *Nature Reviews Earth & Environment* 3 (7): 477–93.