# Changelog - July 5th, 2023

**Dane:**
- [Visualization](#)
    - Ergot
    - Station distribution
    - Station data (hly / dly)
    - Soil moisture
    - Soil
- Updated importErgot script (rejects if incomplete)
- Merged districts 4840 and 4841/persisted changes
- Soil data (condensed and added to readme)
- Aggregated soil, ergot and weather data
- Feature engineering
- [Started research](#)/modeling

**Daniel:**
- Peer review
- Ergot visualization
- Ergot statistics
- Research on models
- Continue implementing autoencoder and ML models
- Cross validation techniques

**Dharmit:**
- Ergot Visualisation
- Soil Data Aggregation
- Satellite Soil Moisture Data
- Started research on models

**Jay:**
- Process soil moisture data
- Aggregate soil moisture data
- Soil moisture Visualization
- Soil moisture statistics
- Did some research and build a basic MLP model
- Tried different variation to merge different datasets

**Joseff:**
- Copernicus satellite data retrieval / aggregation
- Copernicus skip existing data
- Data Synchronization between databases
- Peer review
- Git/Python mentoring
- Data aggregation strategy
- Repo Linter Actions
    - Code format consistency - Black
    - Code type consistency - Mypy
    - ~~Code PEP8 conformance - PyLint~~

# Aggregation strategy

**Primary aggregation attributes**
- Year
- District

**Ergot:**
- Future predictors
- Districts:
  - MB = crop_district + 4600
  - SK = crop_district + 4700 - 1
  - AB = (crop_district)(10) + 4800

**Soil:**
- Very granular and many different soils per polygon
- Using weighted averages based on the percentage they occupy within their polygon
  - Districts have multiple polygons
- Non numerical values become booleans, currently ignoring them

**Weather:** 365 days x features

**Soil Moisture:**
- Replaced null values
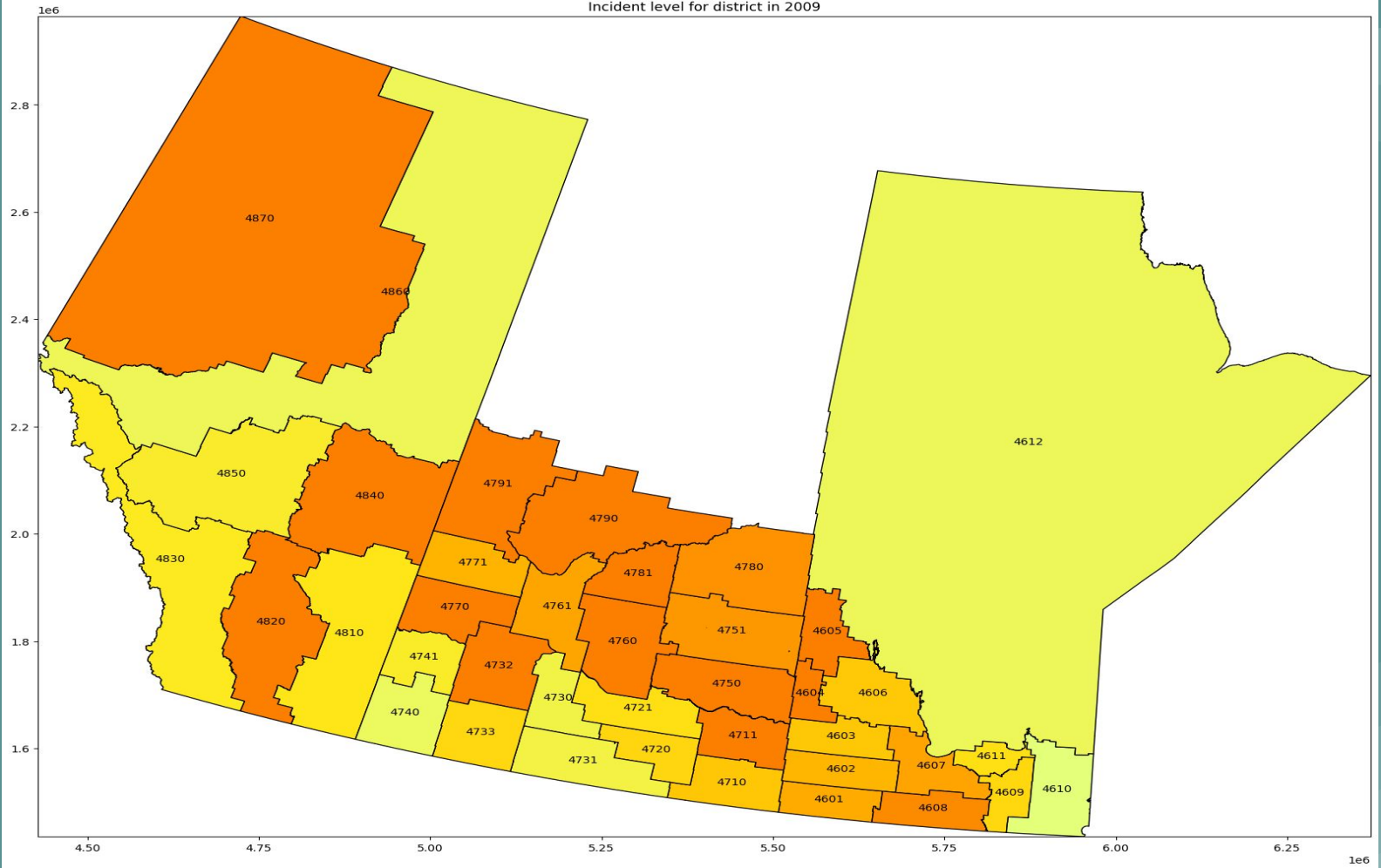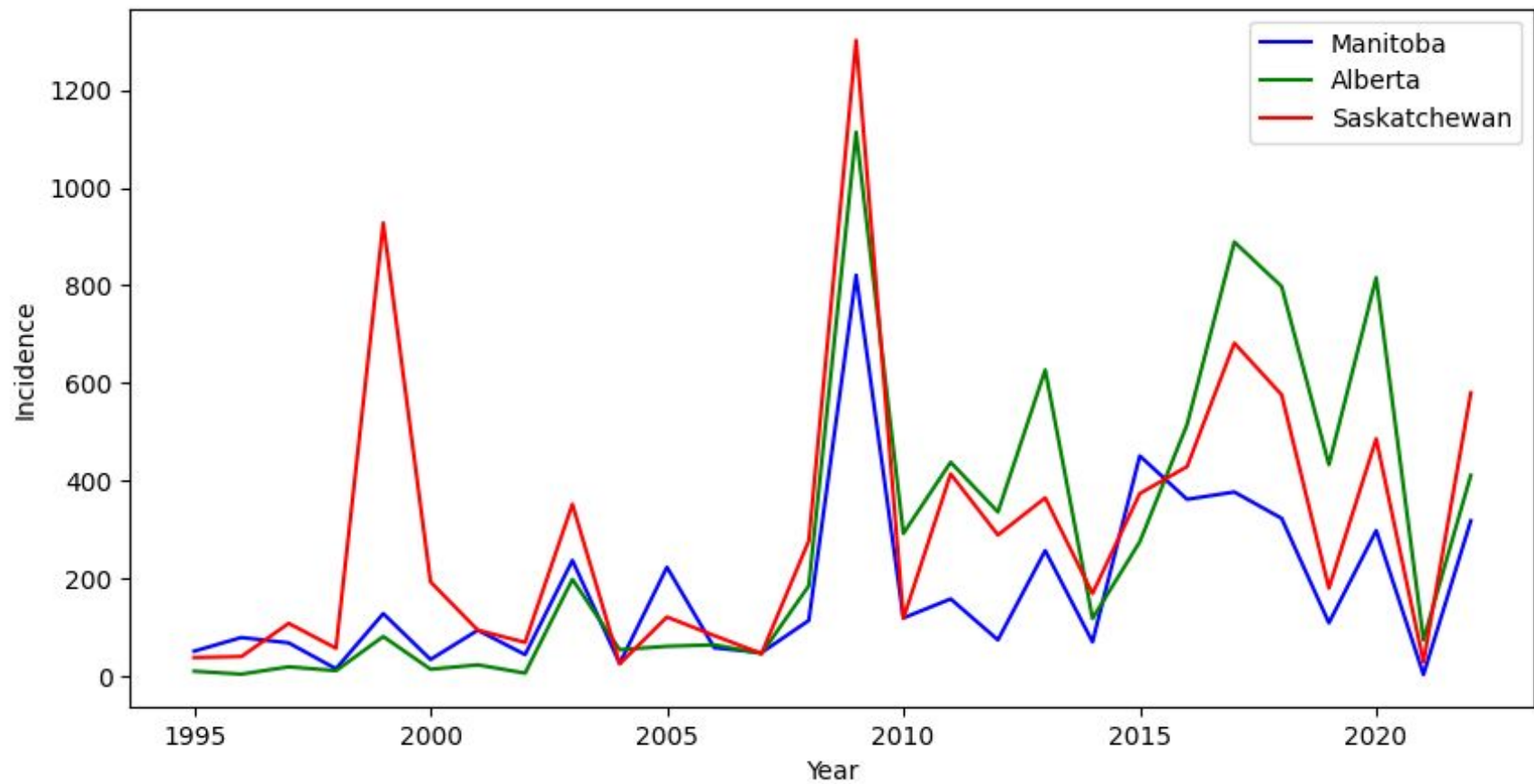- Using year and district to aggregate and generated separate column for min, max and mean

Removed:
Water table (groundwater)
Root restriction types
Drainage
Parent material textures
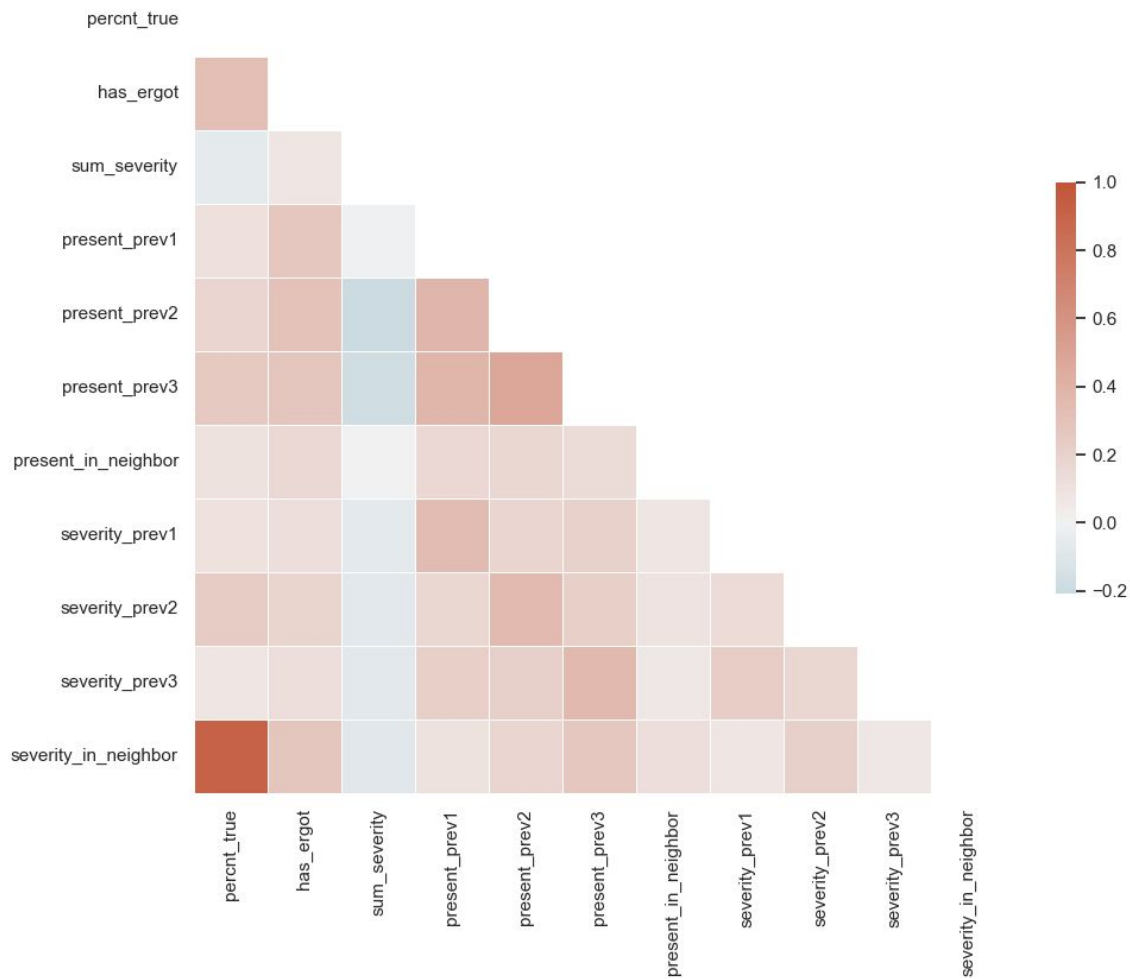Parent material chemicals
Mode of depositions
Sand texture

# Visualization - Ergot

1. What kind of statistics we have done?
   a. Correlation
   b. Incidence ratio

2. What visualization we have created?
   a. Correlation plots
   b. Pair plots
   c. Region plot (show the incidence level)
   d. Line plot (to show the number of ergot over the year)

3. From the visualization, what we have learned about the dataset?
   - 2009 was when ergot happened significantly.
   - We have 139 outliers out of 1092 data points with q1, q2 (median), q3 are 0.02, 0.13, 0.55, respectively
   - Percent of having ergot when prev year had ergot:  0.6813186813186813
   - Percent of having ergot when prev 2 year had ergot:  0.673992673992674
   - Percent of having ergot when prev 3 year had ergot:  0.6446886446886447
   - Percent of having ergot when neighbor is having ergot:  0.7976190476190477
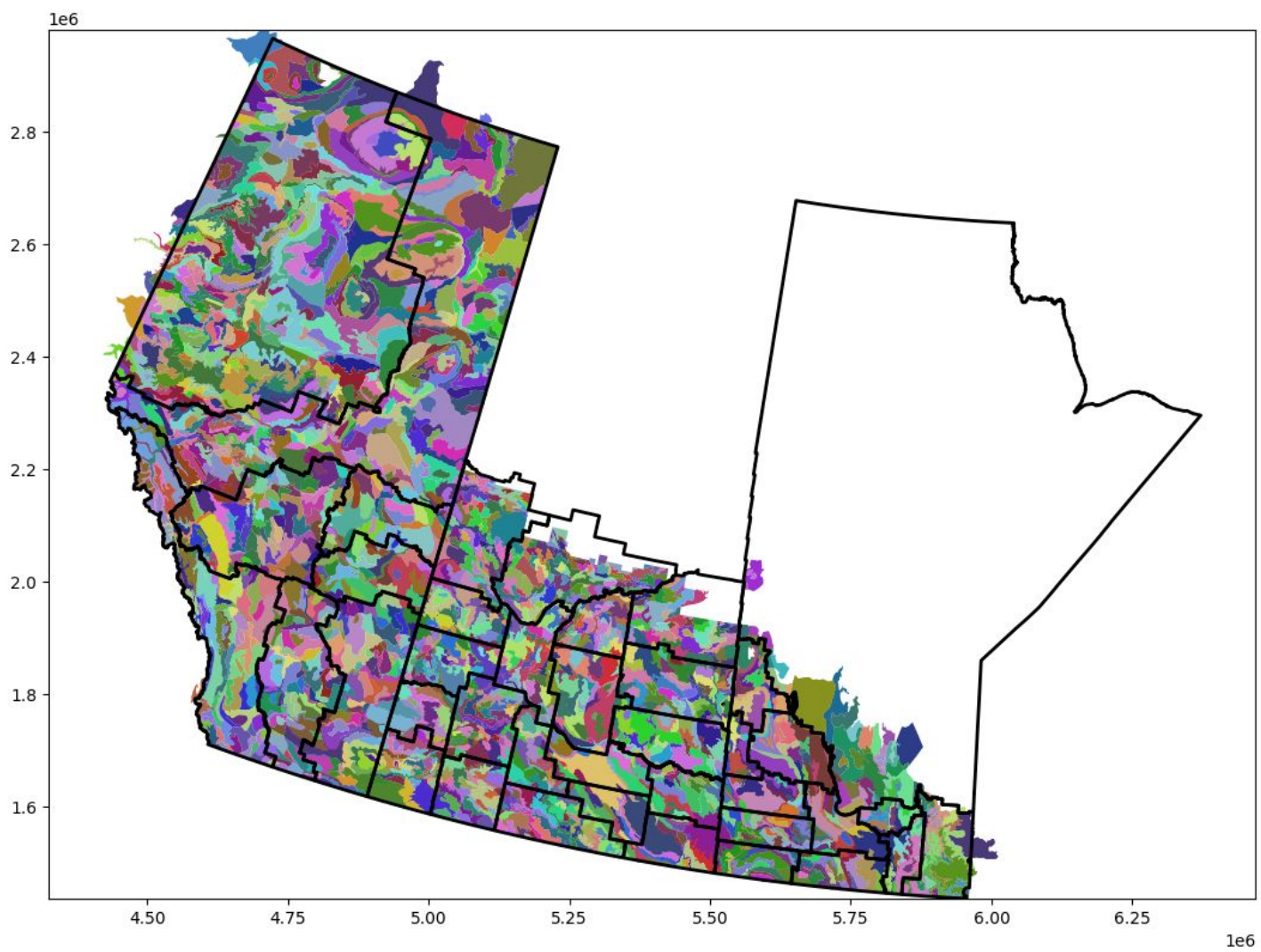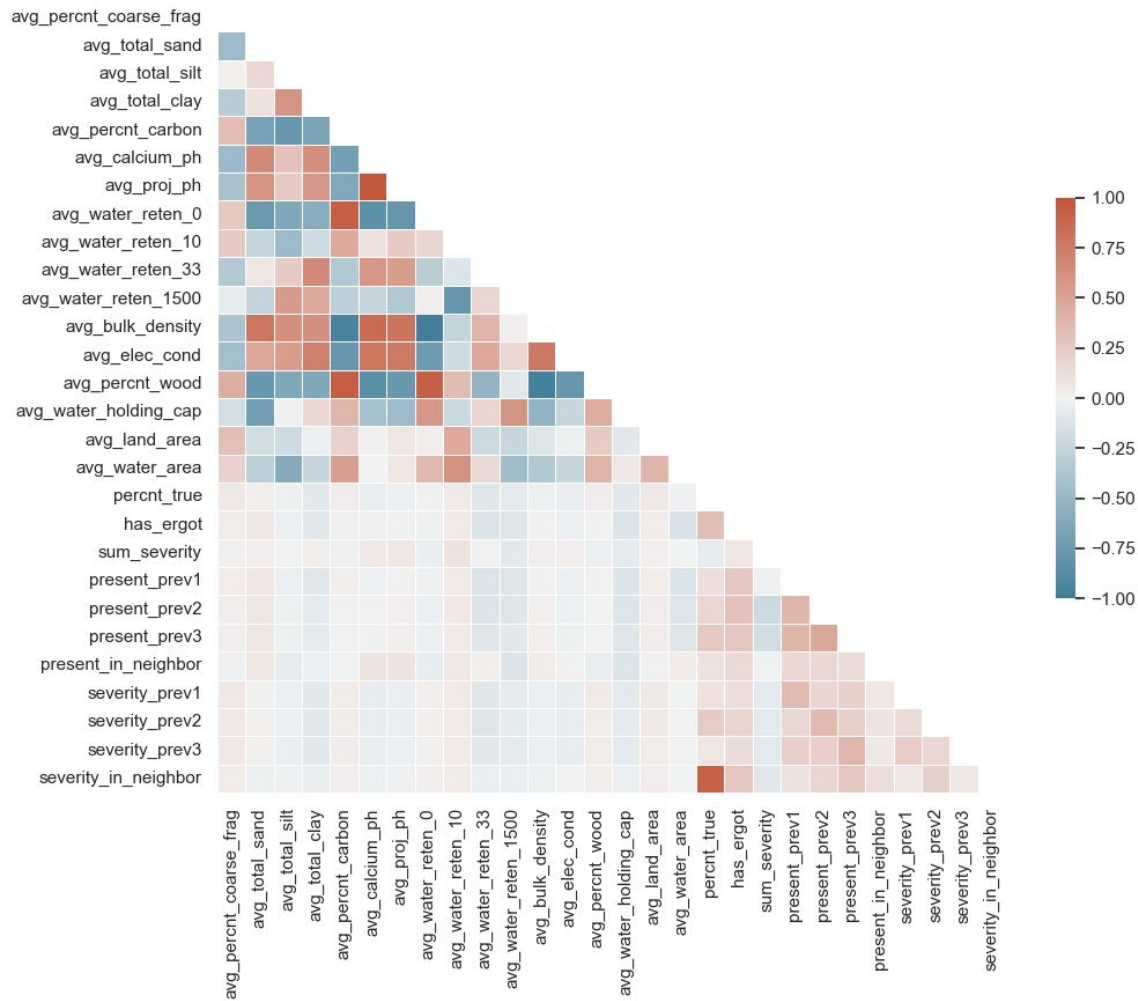
Incident level for district in 2009
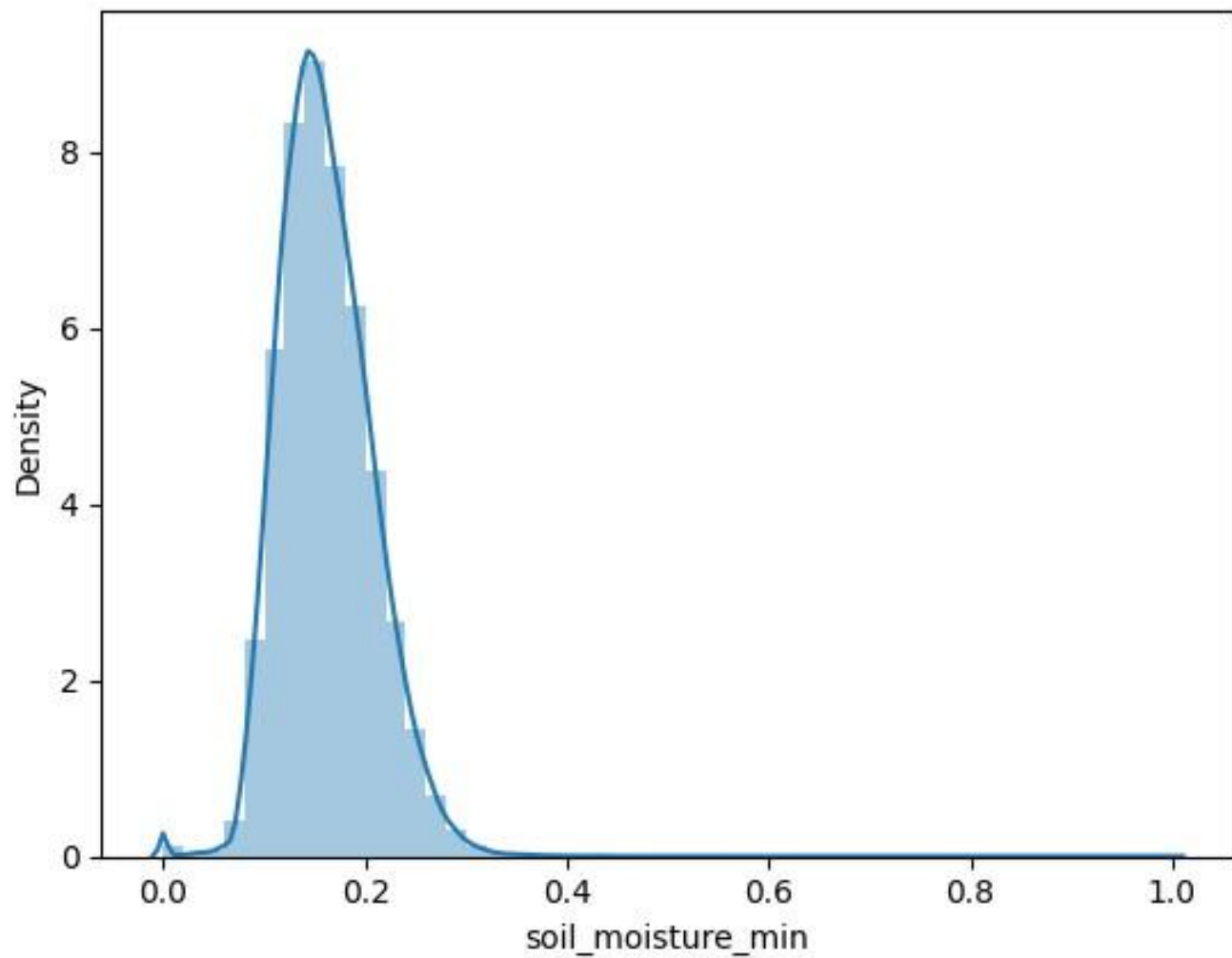
# Visualization - Soil Data

1. What kind of statistics we have done?
   a. Correlation to ergot

2. What visualization we have created?
   a. Correlation plot
   b. Pair plots
   c. Region plot of soils across Canada (as many as 23 different soil types represented by a color in some cases)

3. From the visualization, what we have learned about the dataset?
   a. Complex relationships and complex data
      i. Exploring data (processing/features) and more complicated models will be important
      ii. Important to consider alternative aggregation strategies (possibly through combining factors together and feature engineering)
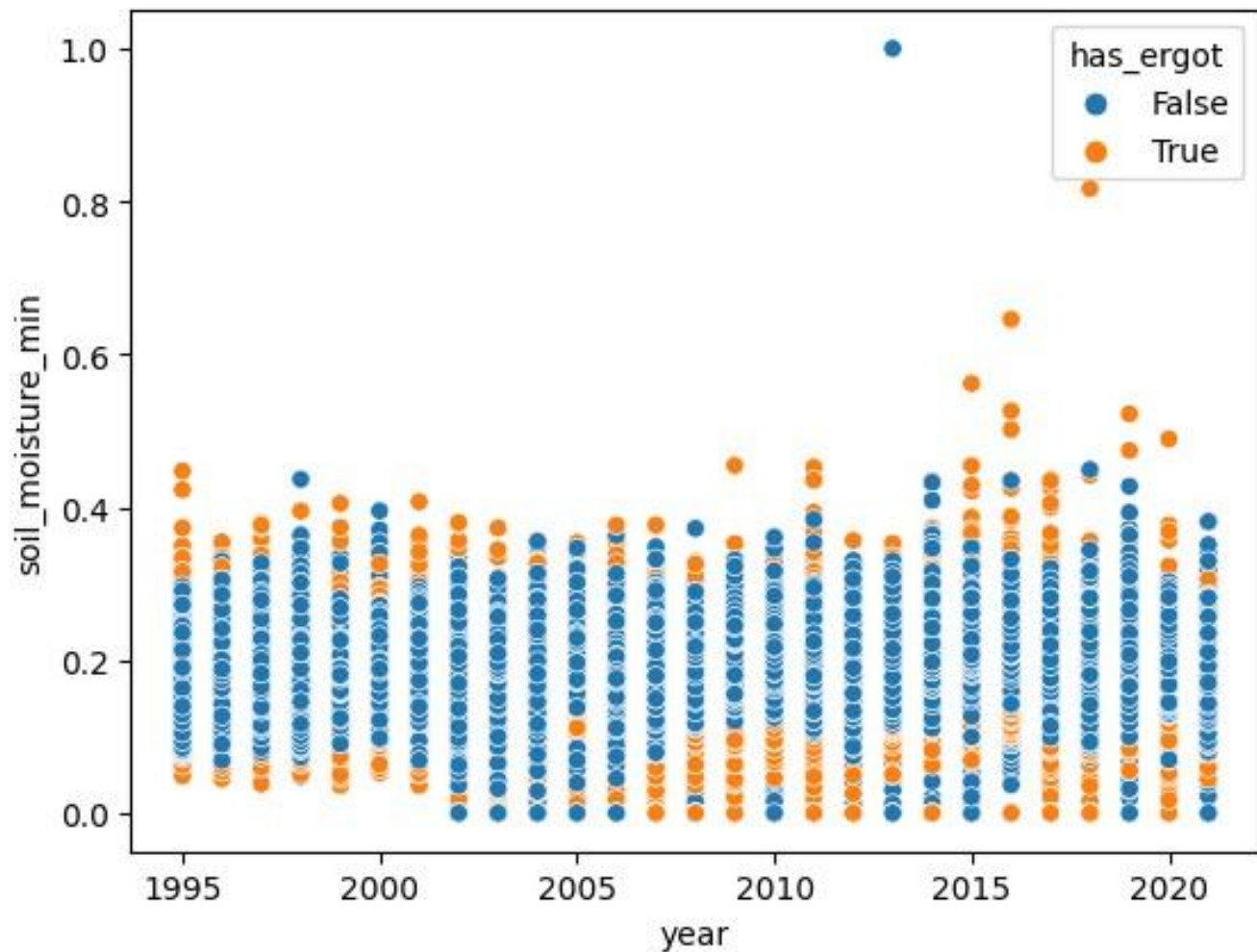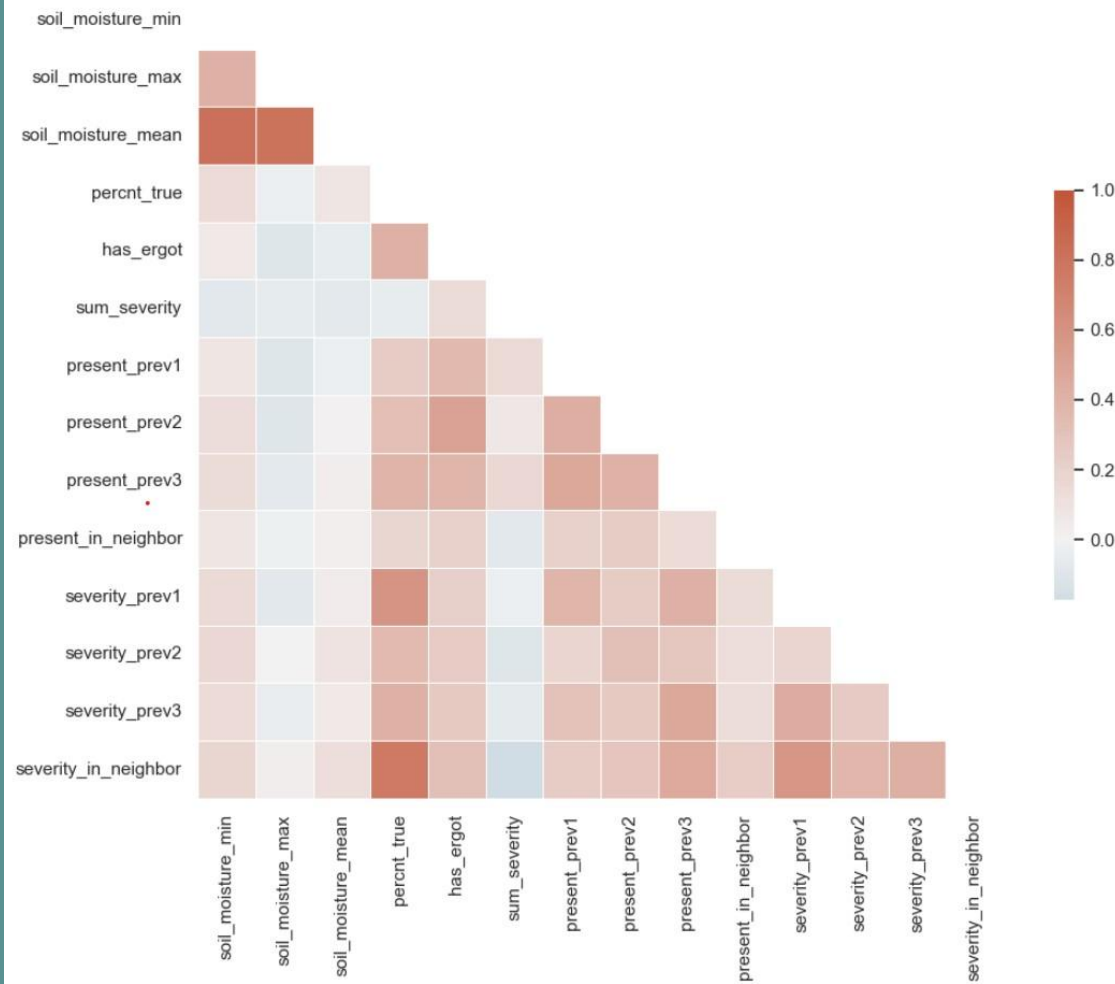
# Visualization - Soil Moisture

1.  What kind of statistics we have done?
    a.  Correlation to ergot
    b.  Data distributions - min, mean, max

2.  What visualization we have created?
    a.  Correlation plot
    b.  Distribution plot (histogram + kde)
    c.  Scatter plot

3.  From the visualization, what we have learned about the dataset?
    a.  The values of soil moisture is skewed
    b.  Complex relationships
        i.  Exploring data (processing/features) and more complicated models will be important
        ii.  There is some kind of reaction between soil moisture and ergot
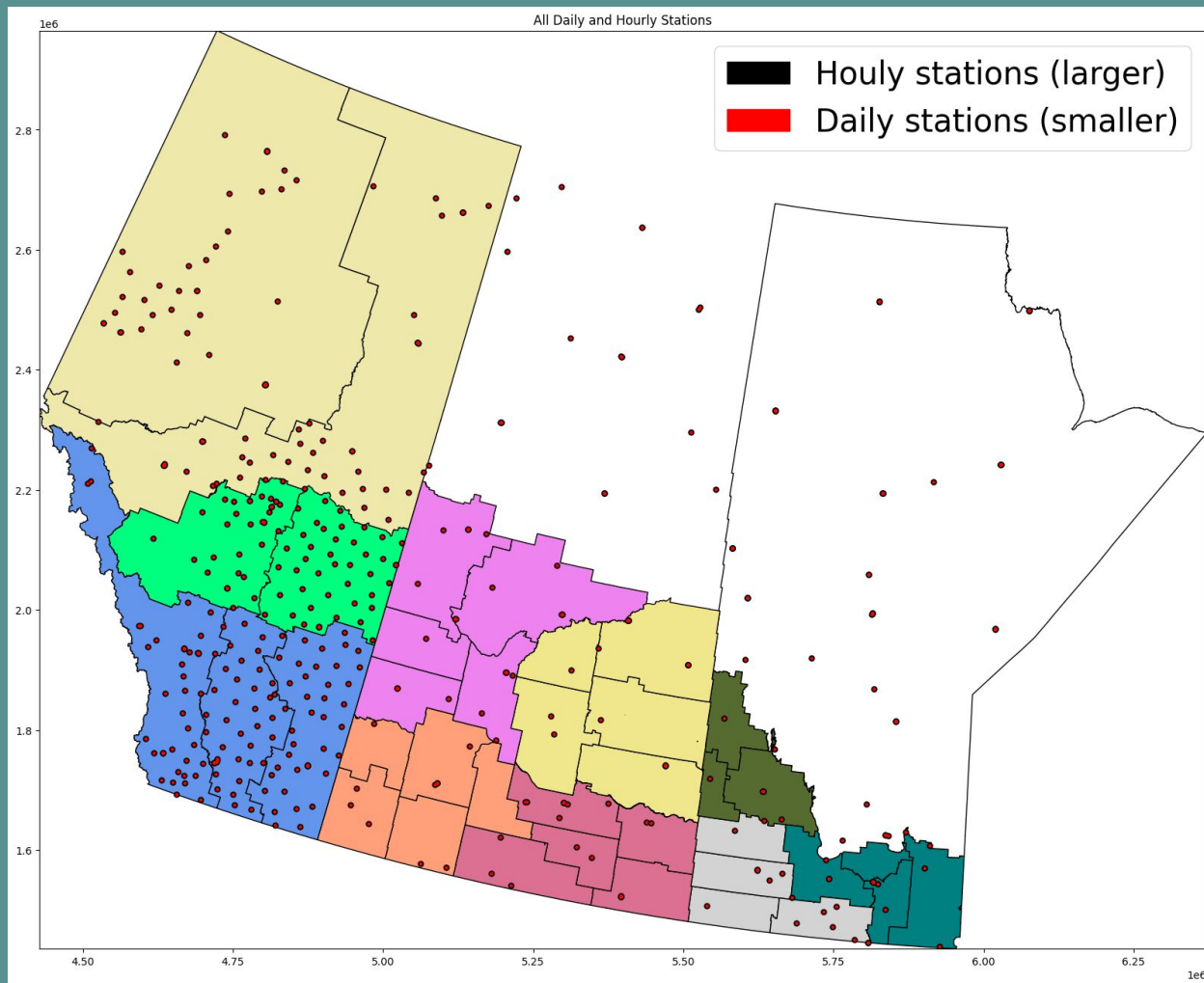
# Visualization - Weather Station

1. What kind of statistics we have done?
   a. Station elevation
   b. Which stations are still active?
   c. Which stations are hourly and which are daily?
   d. Amount of data collected

2. What visualization we have created?
   a. Station summaries for each district
   b. Region plots for stations

3. From the visualization, what we have learned about the dataset?
   a. Usually multiple stations located at the same coordinates
   b. All locations collect both hourly and daily data

All Daily and Hourly Stations

# Visualization - Weather Station Data

1. What kind of statistics we have done?
   a. Correlation to ergot
   b. Data distributions - min, mean, max

2. What visualization we have created?
   a. Correlation plot
   b. Histograms
   c. Pair plots
   d. Box plots (subsetting caused issues)

3. From the visualization, what we have learned about the dataset?
   a. Most exciting correlation was between max_temperature and ergot
   b. Outliers
   c. Complex relationships
      i. Exploring data (processing/features) and more complicated models will be important

# Engineered Features

| | | |
|---|---|---|
| percnt_true | FLOAT | |
| has_ergot | BOOL | |
| median_severity | FLOAT | |
| sum_severity | FLOAT | |

Note that checks for previous years are not accumulative

| | |
|---|---|
| present_in_neighbor | BOOL |
| sum_severity_in_neighbor | FLOAT |

| | |
|---|---|
| median_prev1 | FLOAT |
| median_prev2 | FLOAT |
| median_prev3 | FLOAT |

| | |
|---|---|
| present_prev1 | BOOL |
| present_prev2 | BOOL |
| present_prev3 | BOOL |

| | |
|---|---|
| percnt_true_in_q1 | BOOL |
| percnt_true_in_q2 | BOOL |
| percnt_true_in_q3 | BOOL |
| percnt_true_in_q4 | BOOL |

| | |
|---|---|
| percnt_true_prev1 | FLOAT |
| percnt_true_prev2 | FLOAT |
| percnt_true_prev3 | FLOAT |

| | |
|---|---|
| sum_severity_in_q1 | BOOL |
| sum_severity_in_q2 | BOOL |
| sum_severity_in_q3 | BOOL |
| sum_severity_in_q4 | BOOL |

| | |
|---|---|
| sum_severity_prev1 | FLOAT |
| sum_severity_prev2 | FLOAT |
| sum_severity_prev3 | FLOAT |

# Regression models - ML

- Types of models used:
    - Logistic Regression
    - Random Forest
    - Decision Tree
    - Gradient Boosting
    - Support Vector Classifier (took long)
    - Linear Support Vector Classifier (took long)

- How it performed so far:
    - Current features used are from the soil moisture table (including min, max, mean of the moisture)
    - Most of them has good r-square score (>80%) and f1 score (>90%) using different cross validation techniques (kfold and kfold stratified)
    - Outperform so far: Random Forest

# Models - MLP

- Preparing data for the model means converting categorical columns into numeric column with one-hot encoding like approach. And dealing with missing data and merging different data set.
- We are building Multi-layer perceptron with 1 hidden layer using sigmoid function for the output layer and relu function for the hidden layer.
- As this a binary classification problem we are using binary cross entropy as a loss function and Adam as a optimizer
- Precision and accuracy are too good to be true, which is likely >95% of the time. The suspect is that the number of true samples dominate the number of false samples.
- Metrics are used to measure the models: accuracy, precision, f1 score

# Cross-validation techniques

- Normal train test split
- KFold
- KFold Stratified
- Leave one out

# Regularization strategy

Normal distribution parameters
- Normalization scaling

Skewed distribution parameters
- Log scaling
- Then normalization

Biased data
- Random sampling with
    - Over sampling (of less dominant attribute)
    - Under sampling (of more dominant attribute)

Other options?

# Goals for next 2 weeks

**Dane:**
- Models
- Dimensionality Reduction
- Update documentation i.e missing tables
- Interacting with the system?
    - Front-end?
    - Improved pipeline?

**Daniel:**
- Merge dataset - looking into what kind of datasets we should experiment on
- Continue modelling
- Write script to run all the models and comparing model on different of kind metrics (worth looking the custom metric?)

**Dharmit:**
- Research into data and models
- Data Visualisation on correlated data
- Look into aggregation methods
- Models

**Jay:**
- Figure out the best way to merge dataset with leaking information to model
- Understand data to decide which attribute to keep
- Create/improve more model (complex) if MLP starts giving promising results
- More research on data and models

**Joseff:**
- Data visualization not already done (e.g. yearly corr plots)
- Validate aggregation methods
- Models