

Methodology: Cherry Blossom Bloom-Date Prediction Linear Model and LASSO

Peak Bloom Prediction Workflow

February 21, 2026

1 Objective

This document summarizes the modeling methodology used to predict annual cherry blossom peak bloom timing (day-of-year) for each location, with a focus on the standard **linear model (LM)** and a secondary **LASSO** model.

2 Data Sources and Initialization

The workflow loads and harmonizes:

- Bloom history data by location and year (peak bloom date and day-of-year).
- Daily climate station data (temperature, precipitation, station metadata).

Key preprocessing decisions include:

- Restrict records to years ≥ 1973 .
- Map NOAA stations to bloom locations.
- Keep selected climate variables (temperature and precipitation inputs for the standard model pipeline).
- Fill only short internal temperature gaps (less than 3 consecutive days) by linear interpolation.

3 Feature Engineering for Standard Model

From daily climate records and bloom events, yearly pre-bloom predictors are constructed.

3.1 Altitude/Location Adjustment

For each location, station-to-bloom-site altitude difference is computed using:

- Station altitude from climate station elevation.
- Bloom-site altitude from bloom history metadata.

A lapse-rate temperature correction is applied:

$$\Delta T = \gamma \cdot \frac{h_{\text{station}} - h_{\text{bloom}}}{1000}, \quad \gamma = 6.5 \text{ }^{\circ}\text{C/km.} \quad (1)$$

Adjusted daily temperatures are:

$$T_{\max}^{\text{adj}} = T_{\max} + \Delta T, \quad T_{\min}^{\text{adj}} = T_{\min} + \Delta T. \quad (2)$$

3.2 Pre-bloom Annual Aggregates

For each location-year, using all days from Jan 1 through bloom date:

- mean_tmax_adj_prebloom
- mean_tmin_adj_prebloom
- total_prcp_prebloom
- Bloom-site altitude feature bloom_alt_m

The target is bloom day-of-year, $y = \text{bloom_doy}$.

4 Train / Validation / Test Design

The standard model uses location-based domain split:

- **Training locations:** Kyoto, Washington DC, Liestal.
- **Holdout locations:** all remaining locations.

For holdout locations, years are ordered chronologically and split into:

- first half → validation,
- second half → test.

This tests geographic transferability while preserving temporal order within each holdout location.

5 Primary Model: Linear Regression (LM)

The baseline model is:

$$\text{bloom_doy} = \beta_0 + \beta_1 \text{mean_tmin_adj_prebloom} + \beta_2 \text{mean_tmax_adj_prebloom} + \beta_3 \text{total_prcp_prebloom} + \beta_4 \text{bloom_alt_m} \quad (3)$$

Evaluation metrics on validation/test:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i|, \quad \text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2}. \quad (4)$$

6 Secondary Model: LASSO

A LASSO regression is fit using the same predictors and target, with λ selected by cross-validation on training data only.

Optimization objective:

$$\min_{\beta_0, \beta} \left\{ \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - x_i^\top \beta)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}. \quad (5)$$

The selected model uses λ_{\min} from CV and is scored on validation/test using the same MAE and RMSE metrics as LM.

7 2026 Forecasting and Uncertainty

7.1 Feature Construction for 2026

For each location, the most recent available feature row is used as the predictor baseline and assigned year = 2026.

7.2 LM 90% Confidence Bounds

For the linear model, 90% confidence intervals are produced directly from the regression prediction interval output for the mean response:

- predicted DOY,
- lower/upper 90% DOY bounds,
- converted calendar-date bounds,
- uncertainty width as \pm days.

7.3 LASSO 90% Confidence Bounds

Because standard closed-form confidence bounds are not directly provided for penalized fits, uncertainty is estimated via bootstrap on training rows:

1. Resample training rows with replacement.
2. Refit LASSO at fixed λ_{\min} .
3. Predict 2026 DOY for each location.
4. Repeat (e.g., 500 replicates) and take empirical 5th/95th percentiles for a 90% interval.

The reported outputs include lower/upper bounds and \pm day summaries.

8 Model Comparison Outputs

The process stores:

- Validation/test metrics by model (LM vs LASSO).
- Prediction-level residual tables by split and model.
- 2026 prediction tables for LM and LASSO.
- A side-by-side 2026 comparison including DOY difference (LASSO – LM).

9 Outputs for Final Report

This section can be used directly in the final competition report to summarize model performance and 2026 predictions.

9.1 2026 Prediction Results (LM vs LASSO)

The following values are taken from:

- `data/model_outputs/predictions_2026_linear.csv`
- `data/model_outputs/predictions_2026_lasso.csv`
- `data/model_outputs/predictions_2026_comparison.csv`

Location	LM Date	LM \pm Days	LASSO Date	LASSO \pm Days
Kyoto	2026-03-29	4.10	2026-03-28	2.52
Washington DC	2026-04-01	1.69	2026-04-01	1.31
Liestal	2026-04-03	1.91	2026-04-03	1.66
Vancouver	2026-04-10	2.52	2026-04-09	1.93
New York City	2026-04-10	3.46	2026-04-08	1.98

Table 1: Predicted 2026 peak bloom dates with uncertainty summarized as \pm days.

9.2 Model Difference (LASSO - LM)

Location	DOY Difference (LASSO - LM)
Kyoto	-1.27
Washington DC	+0.08
Liestal	+0.43
Vancouver	-1.26
New York City	-1.51

Table 2: Difference in predicted bloom day-of-year between LASSO and linear models for 2026.

10 Reproducibility Notes

- A fixed random seed is used for splitting and bootstrap reproducibility.
- All intermediate and final objects are saved as RData artifacts.
- Modeling scripts are modularized into data prep, feature construction, and model fitting stages.