# Cross-validation for comparing qSIP prediction models trained on same or other groups
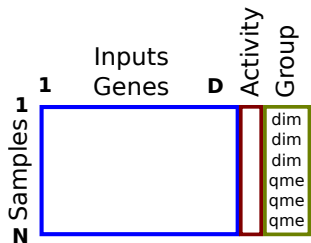
Toby Dylan Hocking
toby.hocking@nau.edu
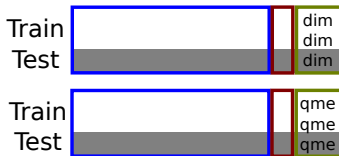toby.hocking@r-project.org

January 11, 2024

# Machine learning predictive analysis of qSIP data

- ▶ Inputs/features $\mathbf{x} \in \mathbb{R}^D$ is vector of TODO for D genes (Amplicon Sequence Variants / ASVs, range from 0 to 10).
- ▶ Output $y \in \mathbb{R}$ is relative activity/growth per day from qSIP (excess atom fraction/EAF normalized by maximum isotope enrichment and incubation length, ranging from 0 to 0.3315).
- ▶ Want to learn $f(\mathbf{x}) = y$ (predict growth from genes).
- ▶ Hypothesis: expect we can learn $f$ on mixed conifer (MC) controls in experiment=dim (room temp), and accurately predict experiment=qme at temp=15C (or vice versa). TODO Jeff what is qme/dim?
- ▶ Question: is this expectation consistent with the data?
- ▶ Answer by using 10-fold cross-validation: train on one experiment or other, quantify prediction error on held out test set.
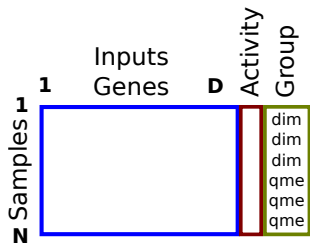
# Comparison 1: controls in different experiments

- ▶ Data table with $N = 7710$ rows/observations (TODO), across two experiments dim=3120, qme=4590.
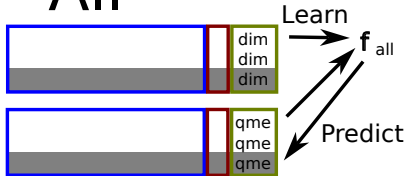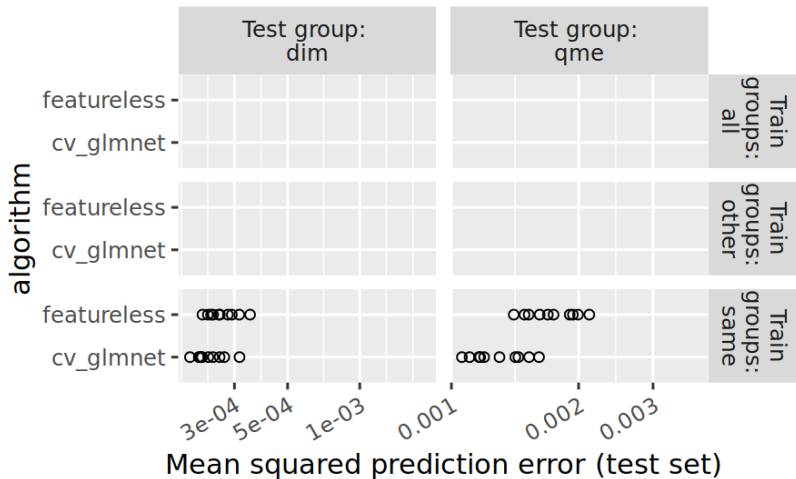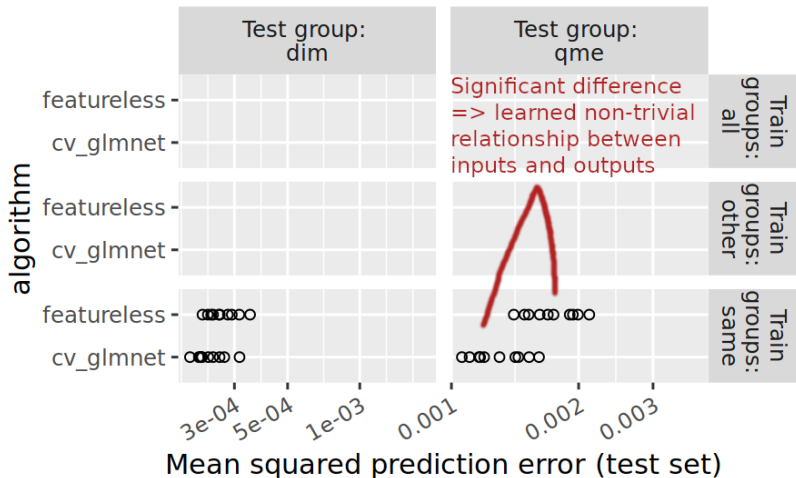
- ▶ $D = 8380$ gene features.

- ▶ We compare two learning algorithms

  cv_glmnet: L1 regularized linear model (LASSO), small subset of important genes selected and used for prediction (other un-important genes are not used for prediction).

  featureless ignore all genes/features, and always predict mean output in train set.

- ▶ If there is any non-trivial relationship/pattern learned between inputs and outputs, then **linear model should have smaller prediction error than featureless**.

- ▶ If patterns are similar in different groups/experiments (dim and qme), then **linear model should have similar prediction error, when trained on other groups/experiments**.
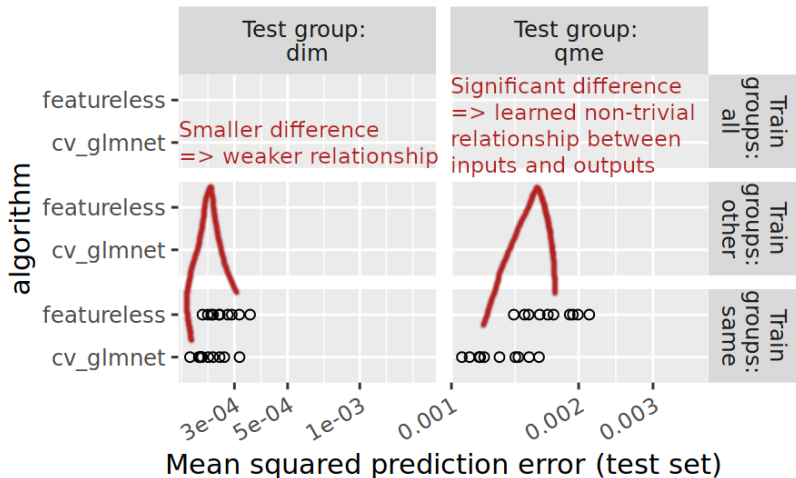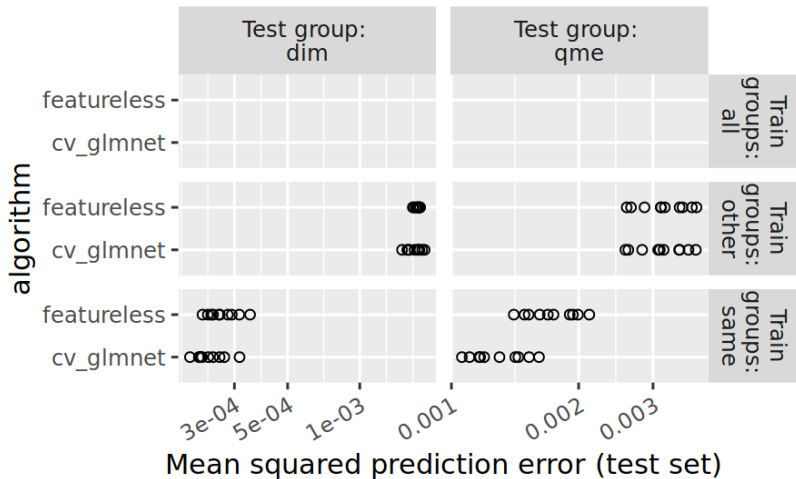
controls between experiments
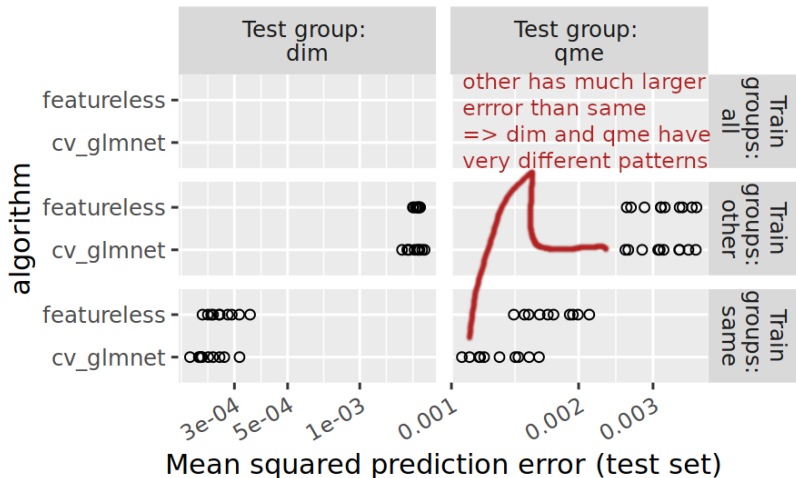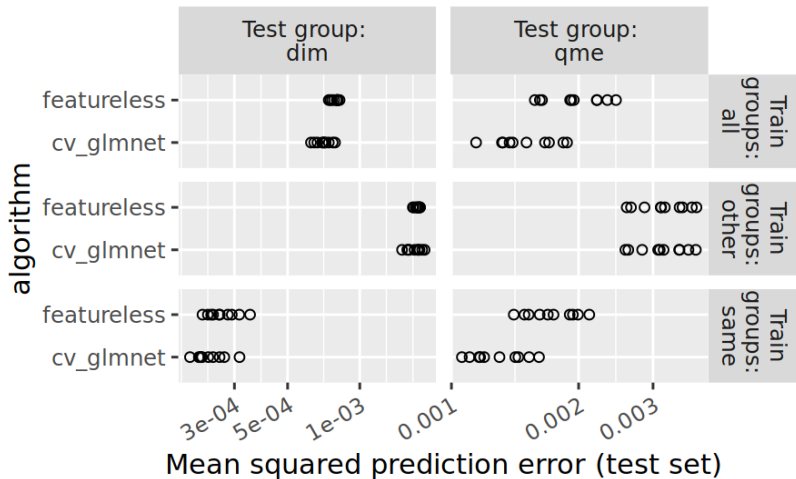
# controls between experiments

controls between experiments
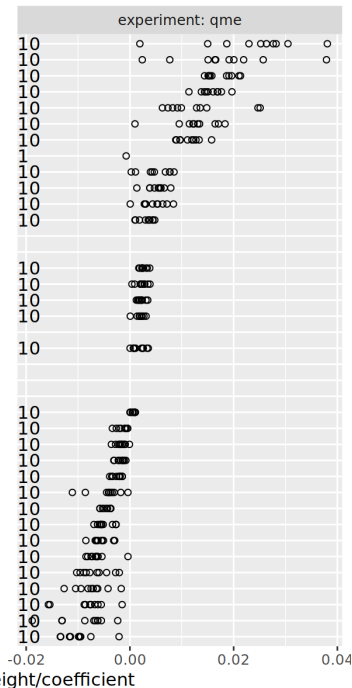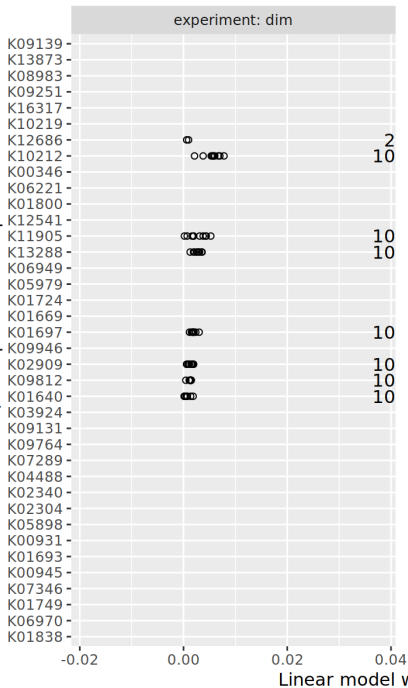
controls between experiments

controls between experiments

# Interpretation of linear model prediction error and weights/coefficientsc

▶ Hypothesis was: expect we can learn $f$ on mixed conifer (MC) controls in experiment=dim (room temp), and accurately predict experiment=qme at temp=15C (or vice versa).

▶ Prediction error cross-validation analysis is not consistent with that hypothesis.

▶ So there should be a different prediction function in each experiment, what is the difference?

▶ The L1 regularized linear model (LASSO) can be interpreted in terms of which genes are important/used for prediction (non-zero weights/coefficients) and others are ignored (weights=0, not used for prediction).

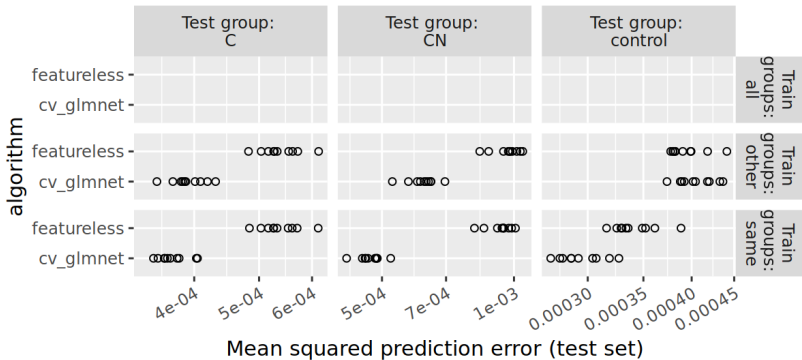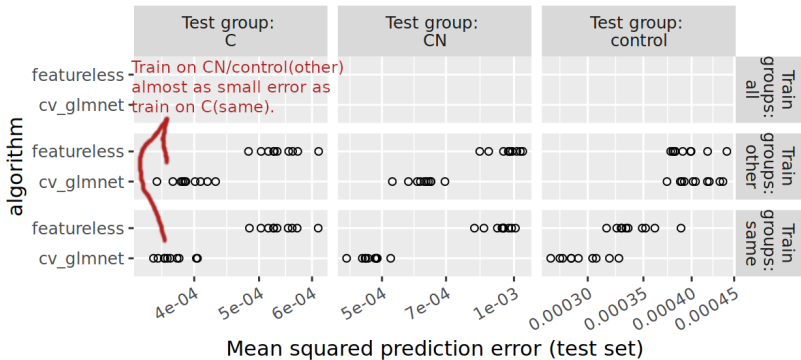▶ Compute and plot weights which are non-zero/important in all 10 train/test splits of cross-validation.

# Comparison 2: control versus carbon additions

- $N = 60877$ samples total, in 3 groups/treatments: control=17225, C=23214, CN=20438.
- Same $D = 8380$ gene features.
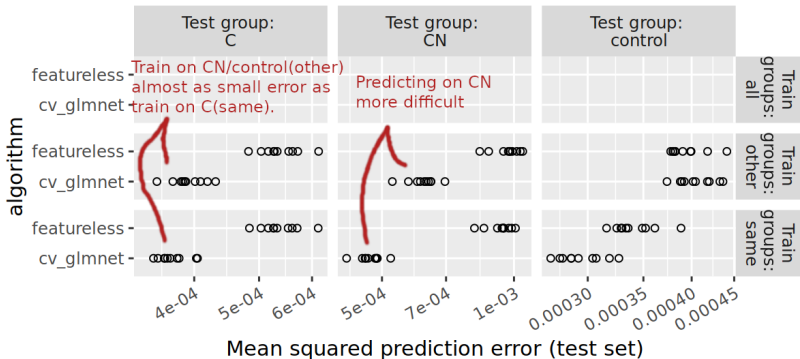- Can we train on one group/treatment, and predict accurately on another?
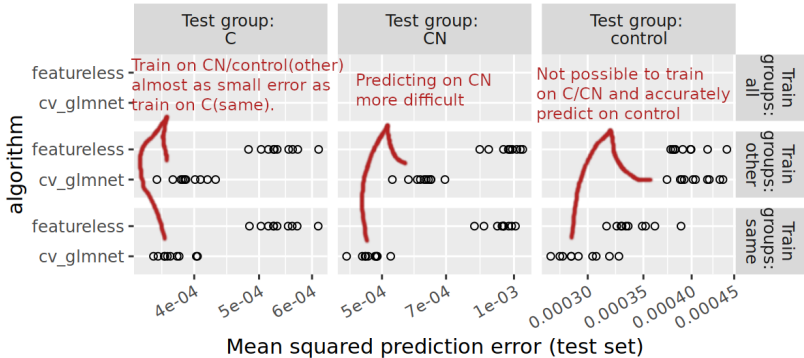
control vs carbon additions

# control vs carbon additions

# control vs carbon additions

Mean squared prediction error (test set)

Train on CN/control(other) almost as small error as train on C(same).

Predicting on CN more difficult

control vs carbon additions

# Discussion and conclusions

- TODO
- Free/open-source software available: mlr3resampling R package.