# Sparse regression via the lasso

**DS-GA 1013 / MATH-GA 2824 Mathematical Tools for Data Science**

Carlos Fernandez-Granda

# Prerequisites

Ordinary least squares (OLS)

Ridge regression

# Sparse regression

Linear regression is challenging when the number of features $p$ is large

Possible solution: Select subset of relevant features $\mathcal{I} \subset \{1, \ldots, p\}$, so that

$$y \approx \sum_{i \in \mathcal{I}} \beta[i] x[i]$$

Problem: How do we find this subset?

Equivalently, how do we find a sparse coefficient vector $\beta \in \mathbb{R}^p$ such that

$$y \approx \langle x, \beta \rangle$$

# Toy problem

Find $t$ such that

$$v_t := \begin{bmatrix} t \\ t - 1 \\ t - 1 \end{bmatrix}$$

is sparse

Equivalently, find $\arg \min_t \|v_t\|_0$

# $\ell_0$ "norm"
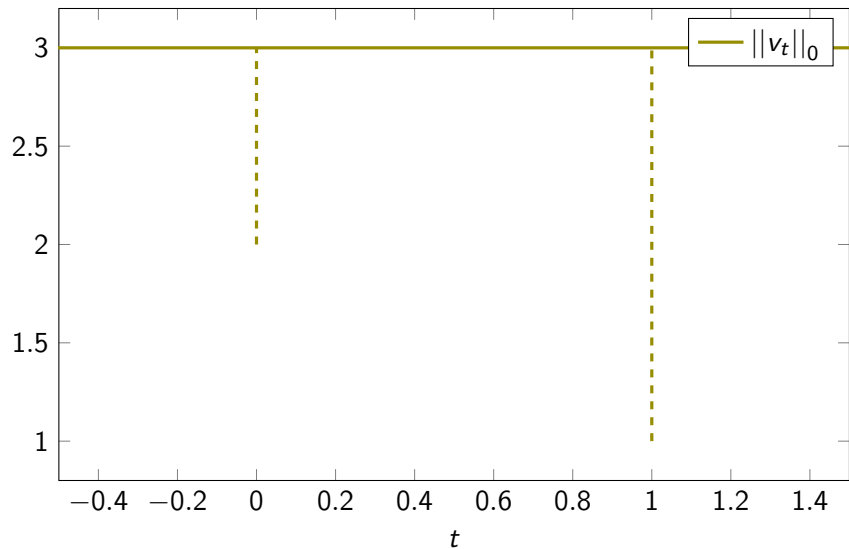
Number of nonzero entries in a vector

Is this a norm?

All norms are homogeneous, $||\alpha x|| = \alpha\, ||x||$ for any $x$ and any scalar $\alpha > 0$

$$\begin{aligned}||2x||_0 &= ||x||_0 \\ &\neq 2\,||x||_0\end{aligned}$$

# Toy problem: $||v_t||_0$?

$$v_t := \begin{bmatrix} t \\ t-1 \\ t-1 \end{bmatrix}$$

# Toy problem: $||v_t||_0$?

# Alternative strategy
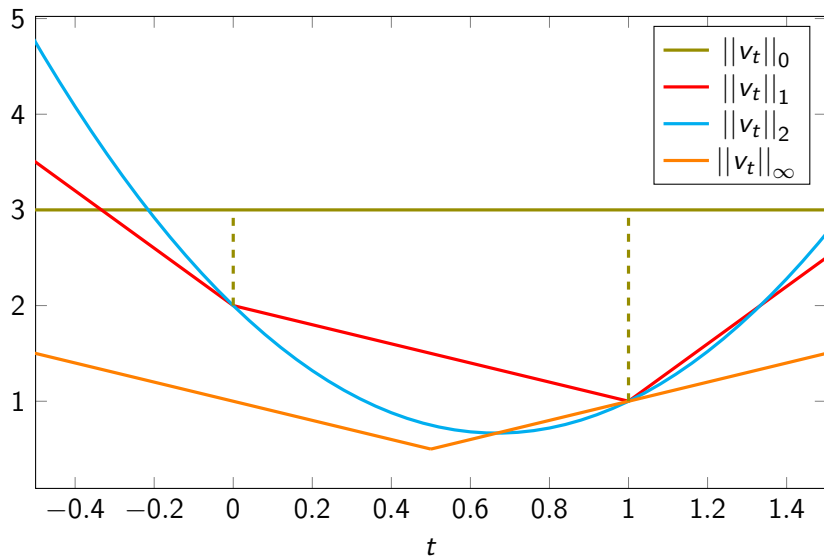
Minimize an actual norm

$$f(t) := \|v_t\| :$$

$$\|x\|_1 := \sum_{i=1}^{d} |x_i|$$

$$\|x\|_2 := \sqrt{\sum_{i=1}^{d} x_i^2}$$

$$\|x\|_\infty := \max_{1 \le i \le d} |x_i|$$

# Toy problem

# Sparse linear regression

Find a small subset of useful features

Model selection problem

Two objectives:

- Good fit to the data; $\left|\left|X^T\beta - y\right|\right|_2^2$ should be as small as possible

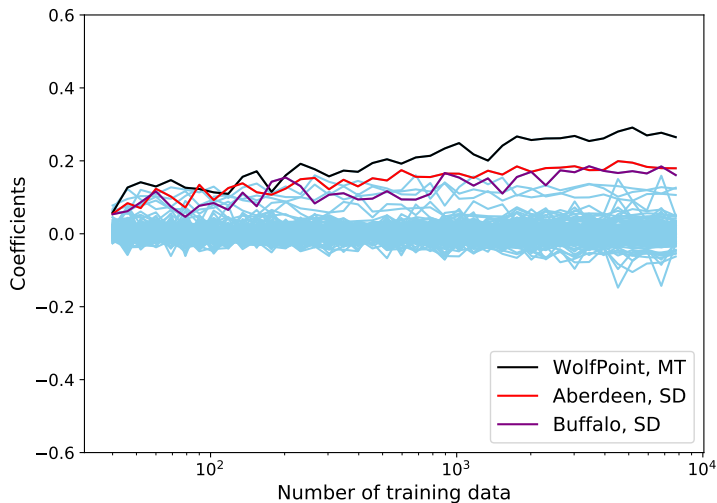- Using a small number of features; $\beta$ should be as sparse as possible

# The lasso

Uses $\ell_1$-norm regularization to promote sparse coefficients

$$\beta_{\mathsf{lasso}} := \arg\min_{\beta} \frac{1}{2} \left\| y - X^T \beta \right\|_2^2 + \lambda \left\| \beta \right\|_1$$

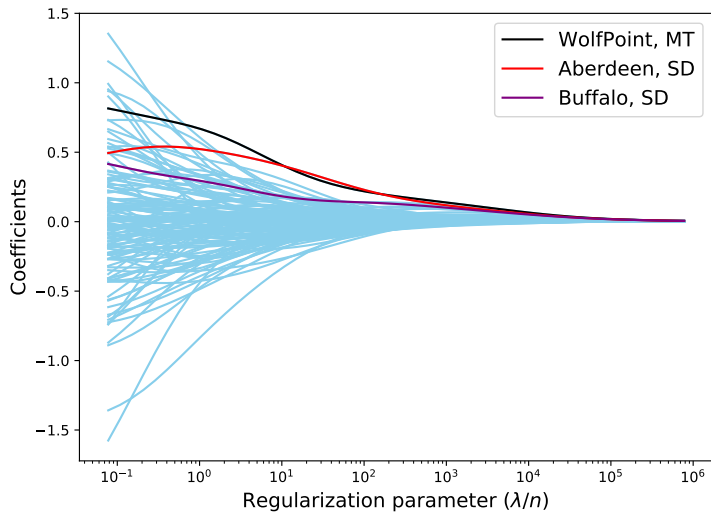# Temperature prediction via linear regression

- Dataset of hourly temperatures measured at weather stations all over the US

- Goal: Predict temperature in Jamestown (North Dakota) from other temperatures

- Response: Temperature in Jamestown

- Features: Temperatures in 133 other stations ($p = 133$) in 2015

- Test set: $10^3$ measurements
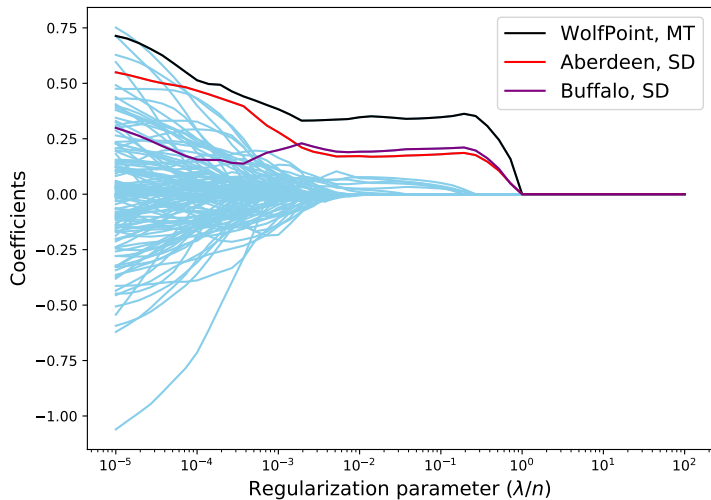
- Additional test set: All measurements from 2016
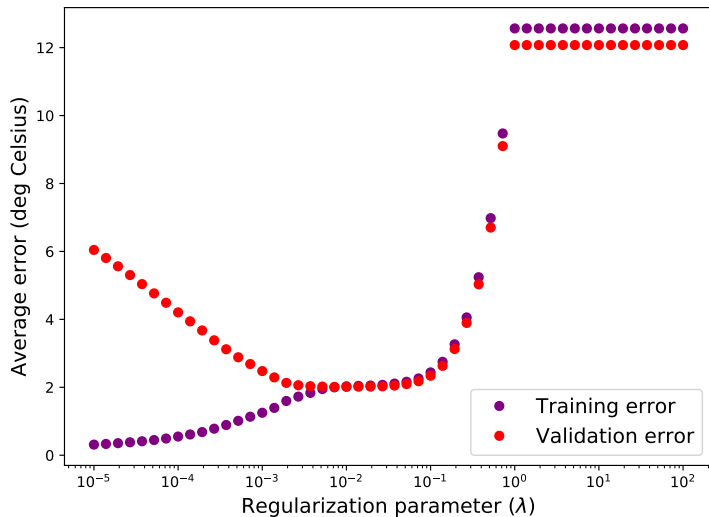
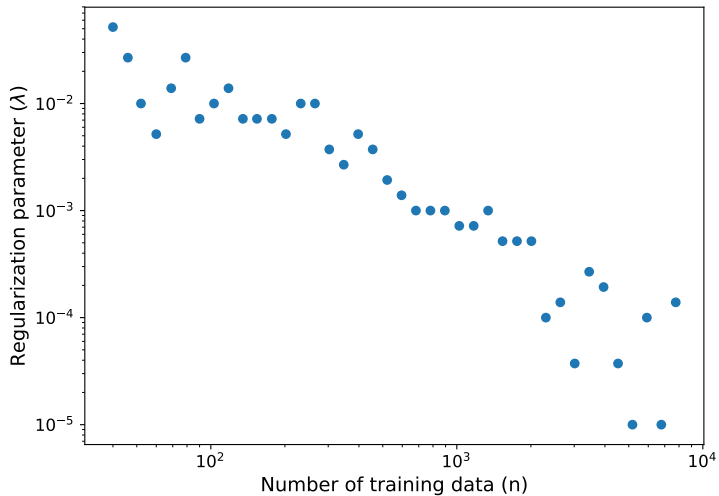# Ridge-regression coefficients
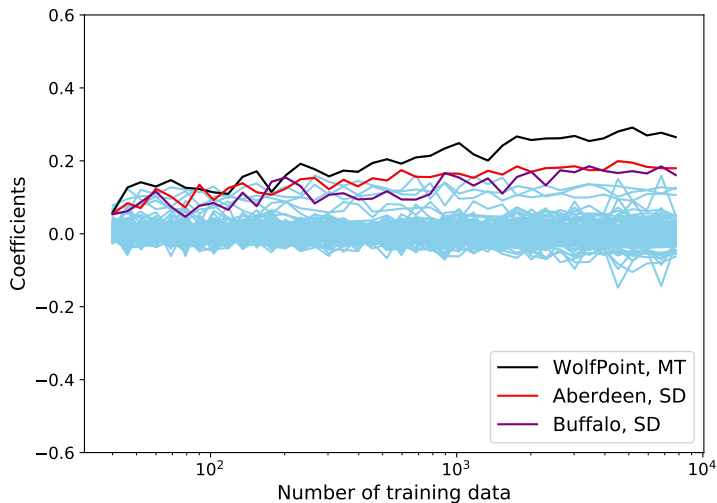
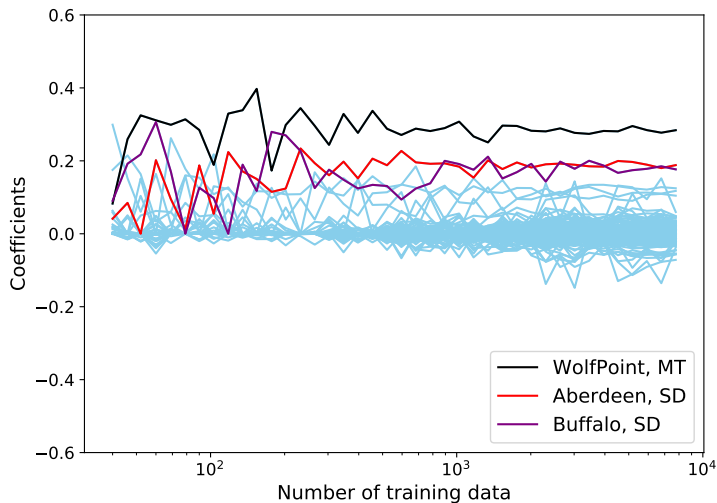# Ridge regression $n := 135$
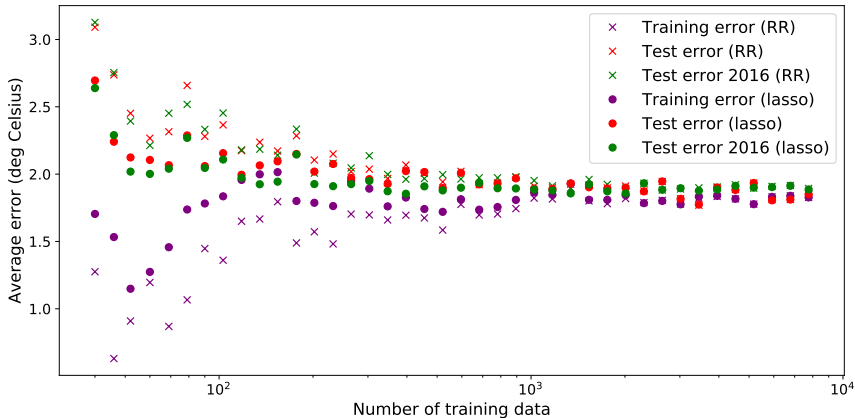
Lasso $n := 135$

Lasso $n := 135$

# Ridge-regression coefficients

# Lasso coefficients

# Results

# What have we learned

- Sparse regression is the problem of fitting a linear model that includes only a subset of relevant features

- Regularization based on $\ell_1$ norm makes it possible to do this automatically!