



Sparse regression

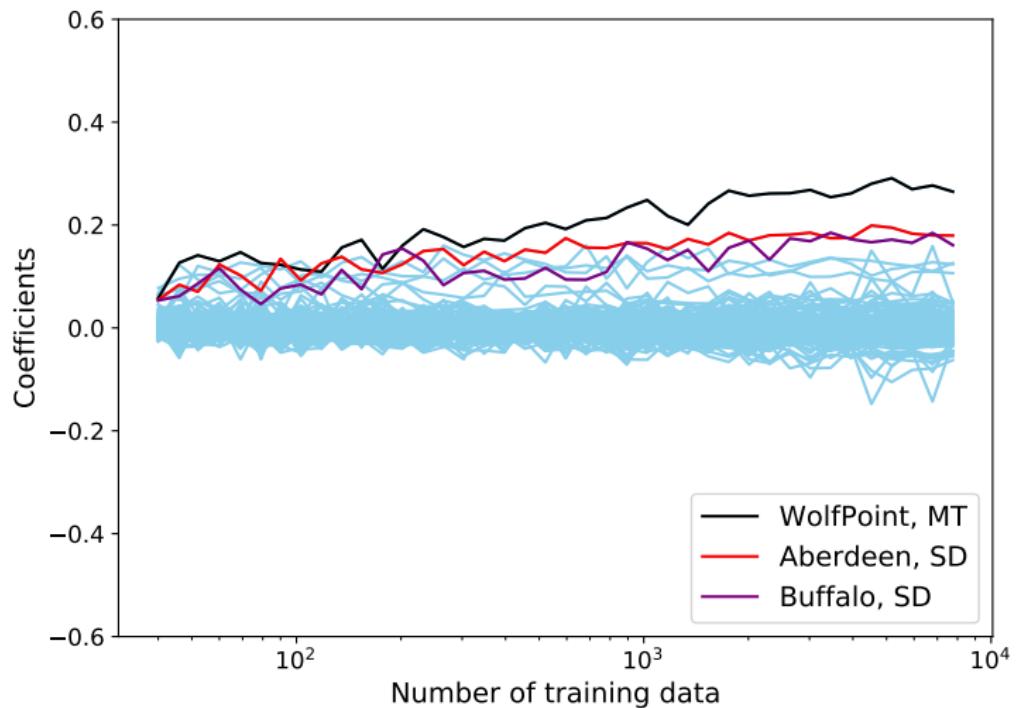
DS-GA 1013 / MATH-GA 2824 Mathematical Tools for Data Science

Carlos Fernandez-Granda

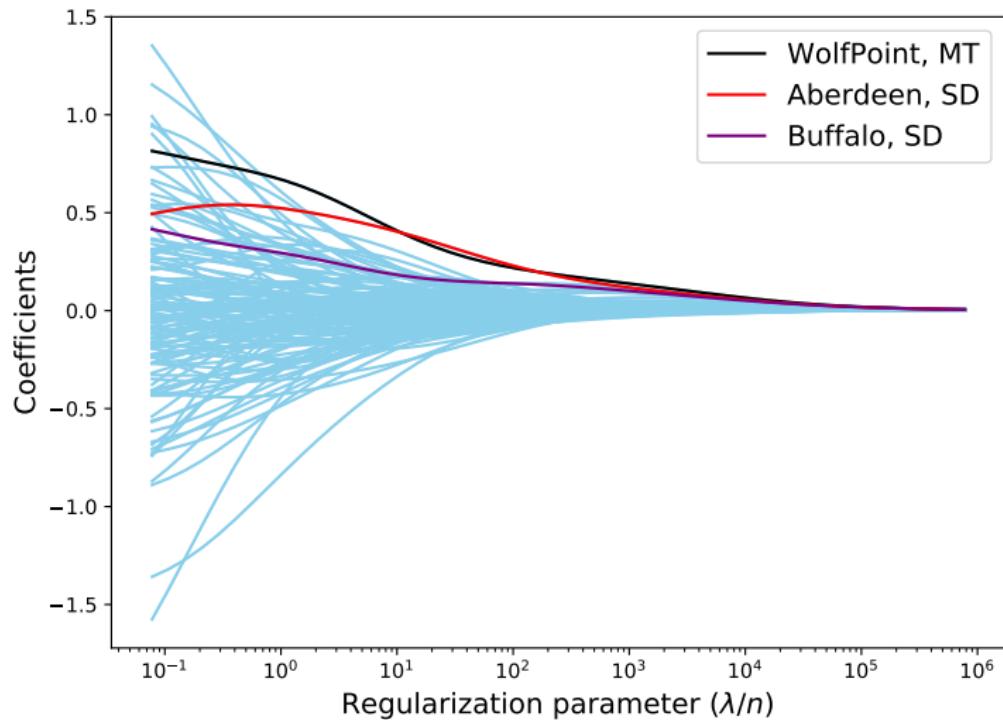
Temperature prediction via linear regression

- ▶ 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 $n := 135$

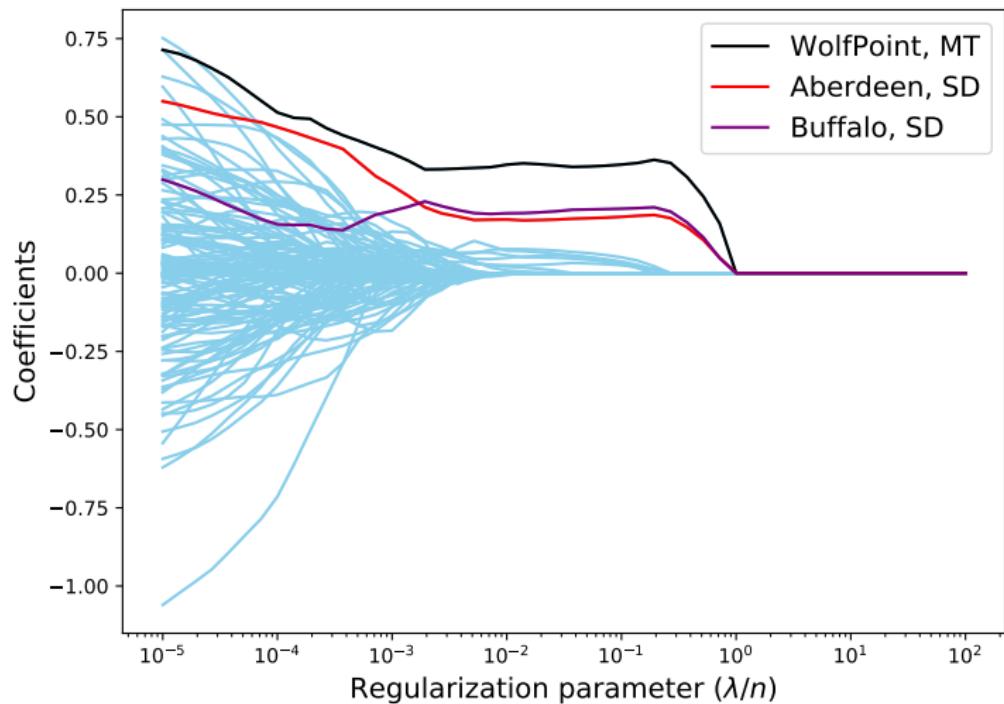


The lasso

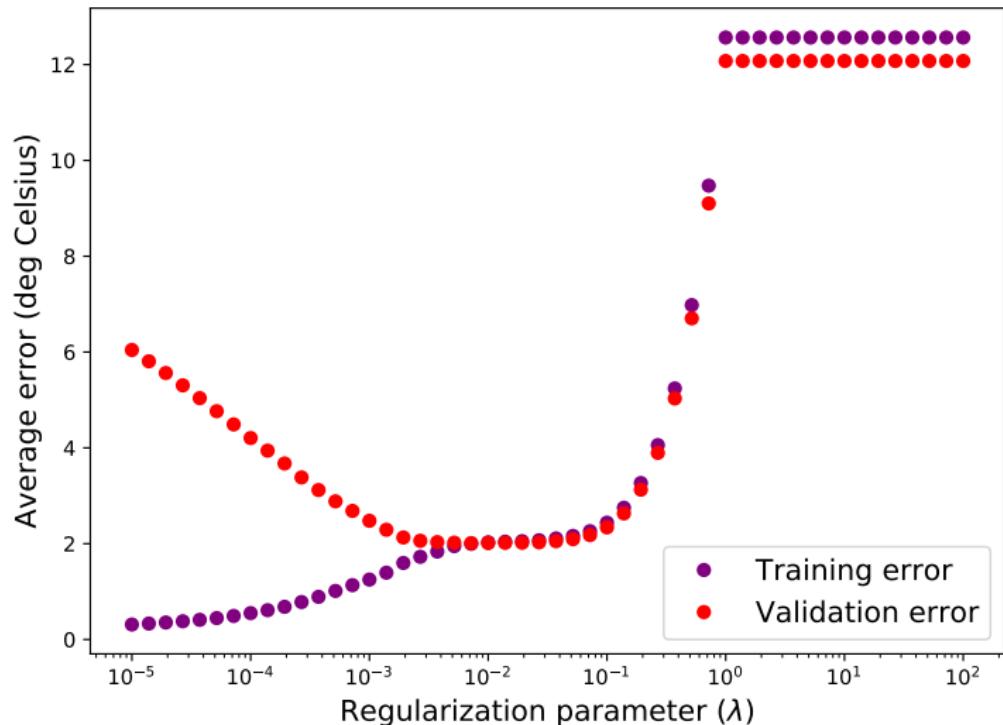
Uses ℓ_1 -norm regularization to promote sparse coefficients

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

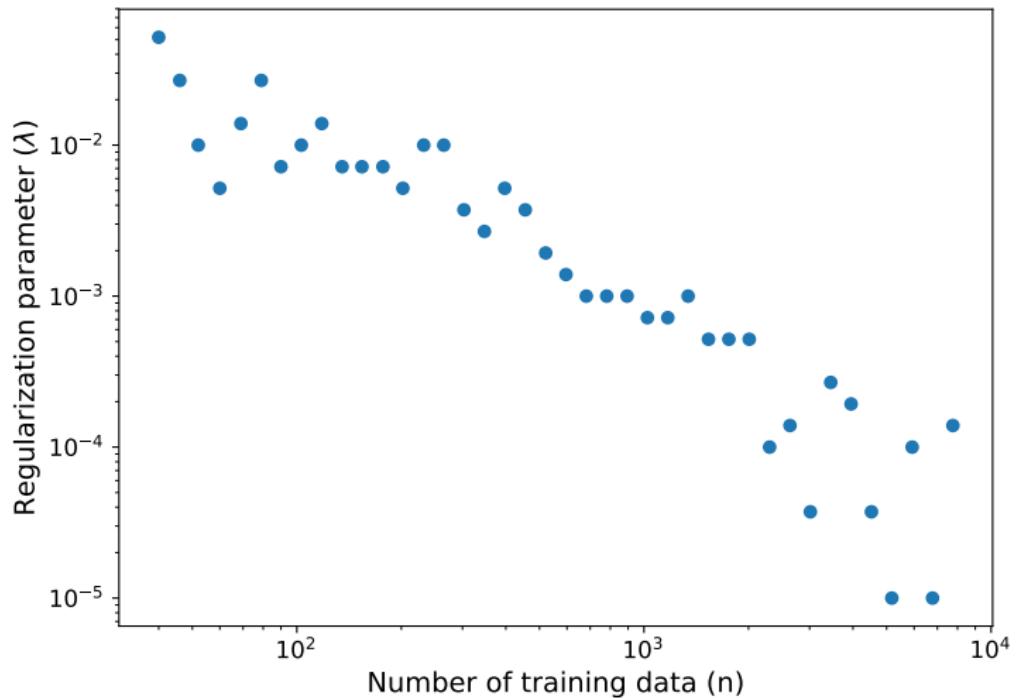
Lasso $n := 135$



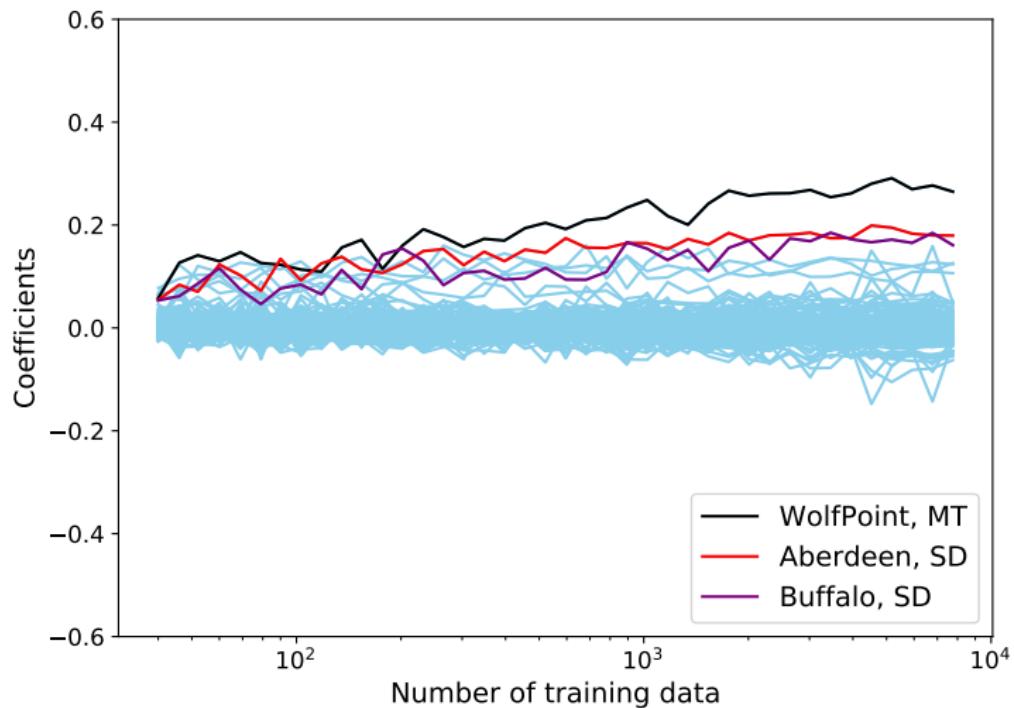
Lasso $n := 135$



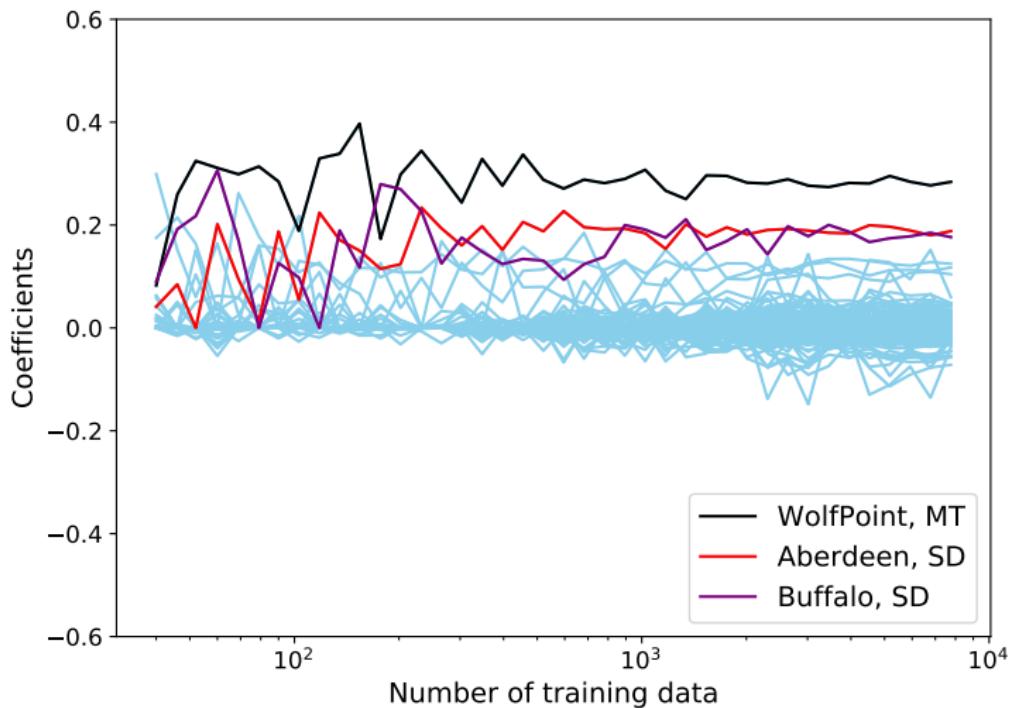
Lasso $n := 135$



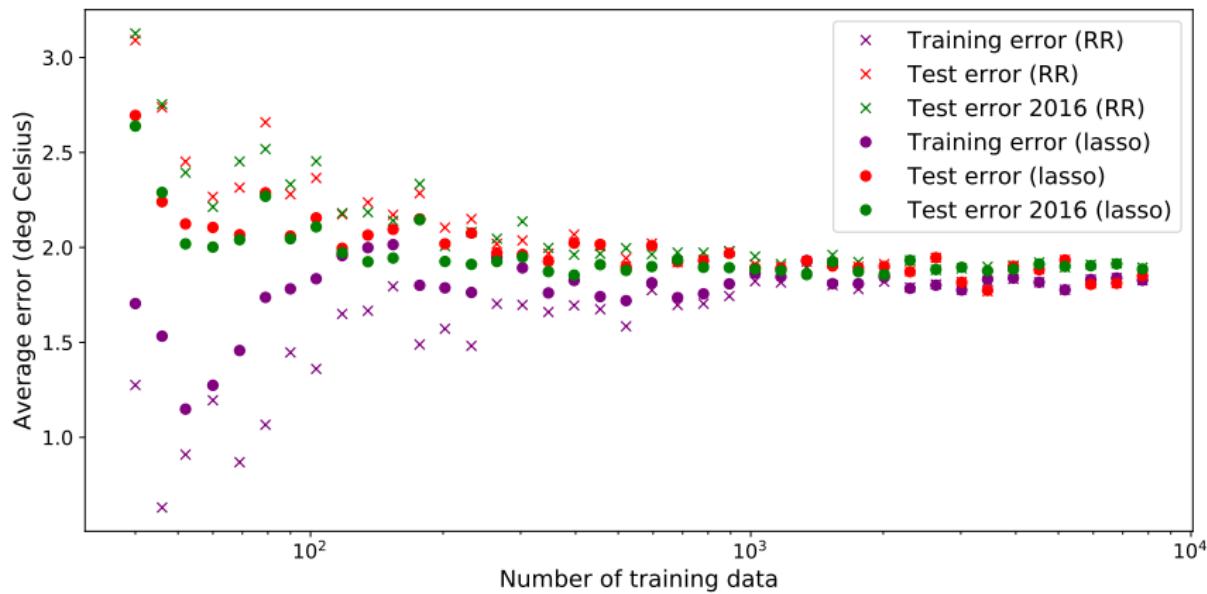
Ridge-regression coefficients



Lasso coefficients



Results



Sparse regression with two features

Feature vectors and noise are fixed n -dimensional vectors

$$y := a + z$$

We fit a model using an additional feature

$$X := \begin{bmatrix} a & b \end{bmatrix}^T$$

$$\beta_{\text{true}} := \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\|a\|_2 = \|b\|_2 = 1$$

Ridge regression

$$\min_{\beta} \left\| X^T \beta - y \right\|_2^2 + \lambda \|\beta\|_2^2$$

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

Solution

$$\beta_{RR} =$$

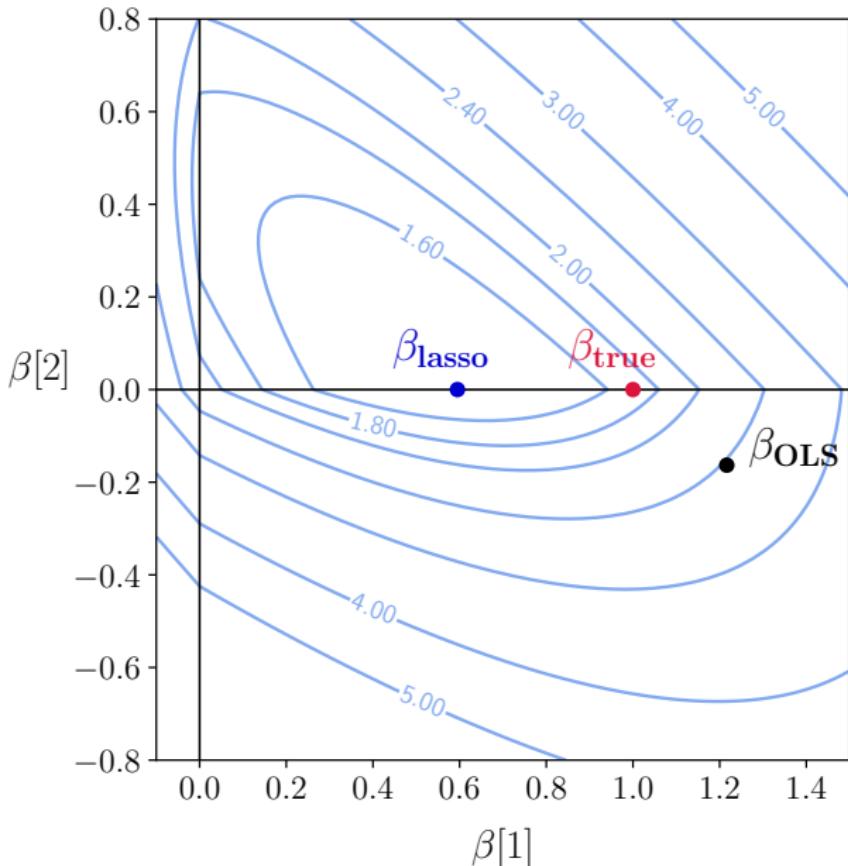
Ridge regression

$$\min_{\beta} \left\| X^T \beta - y \right\|_2^2 + \lambda \|\beta\|_2^2 \quad \begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

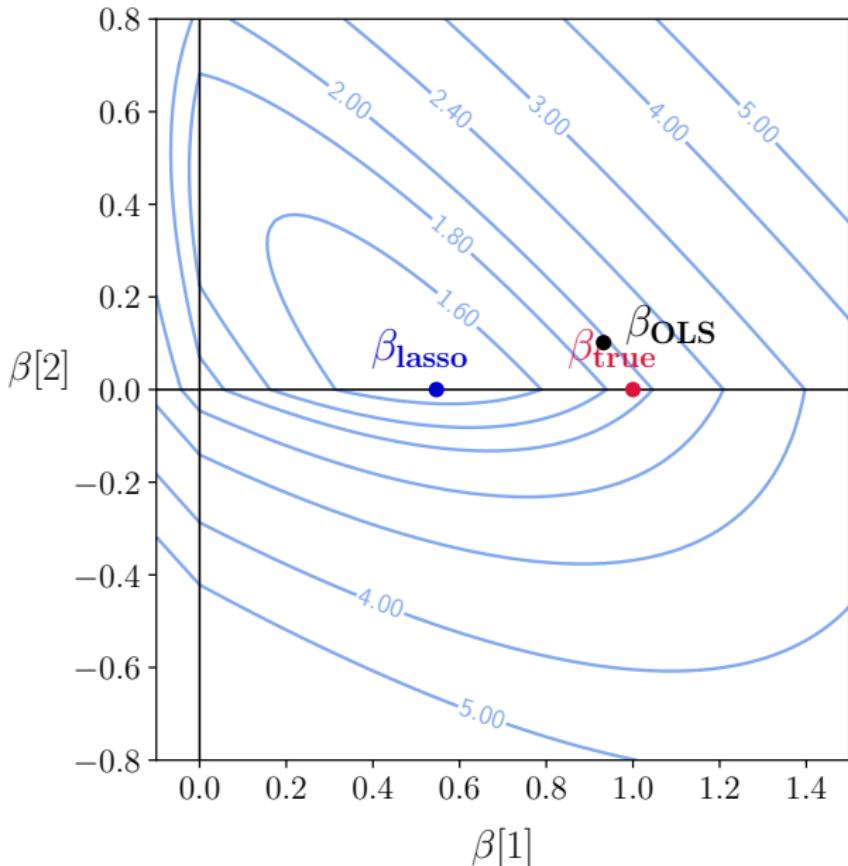
Solution

$$\begin{aligned}\beta_{RR} &= (X X^T + \lambda I)^{-1} X y \\ &= \begin{bmatrix} 1 + \lambda & \rho \\ \rho & 1 + \lambda \end{bmatrix}^{-1} \begin{bmatrix} 1 + a^T z \\ \rho + b^T z \end{bmatrix} \\ &= \frac{1}{(1 + \lambda)^2 - \rho^2} \begin{bmatrix} 1 + \lambda & -\rho \\ -\rho & 1 + \lambda \end{bmatrix} \begin{bmatrix} 1 + a^T z \\ \rho + b^T z \end{bmatrix} \\ &= \frac{1}{(1 + \lambda)^2 - \rho^2} \begin{bmatrix} (1 + \lambda)(1 + a^T z) - \rho(\rho + b^T z) \\ b^T z(1 + \lambda) + \rho(\lambda - a^T z) \end{bmatrix}\end{aligned}$$

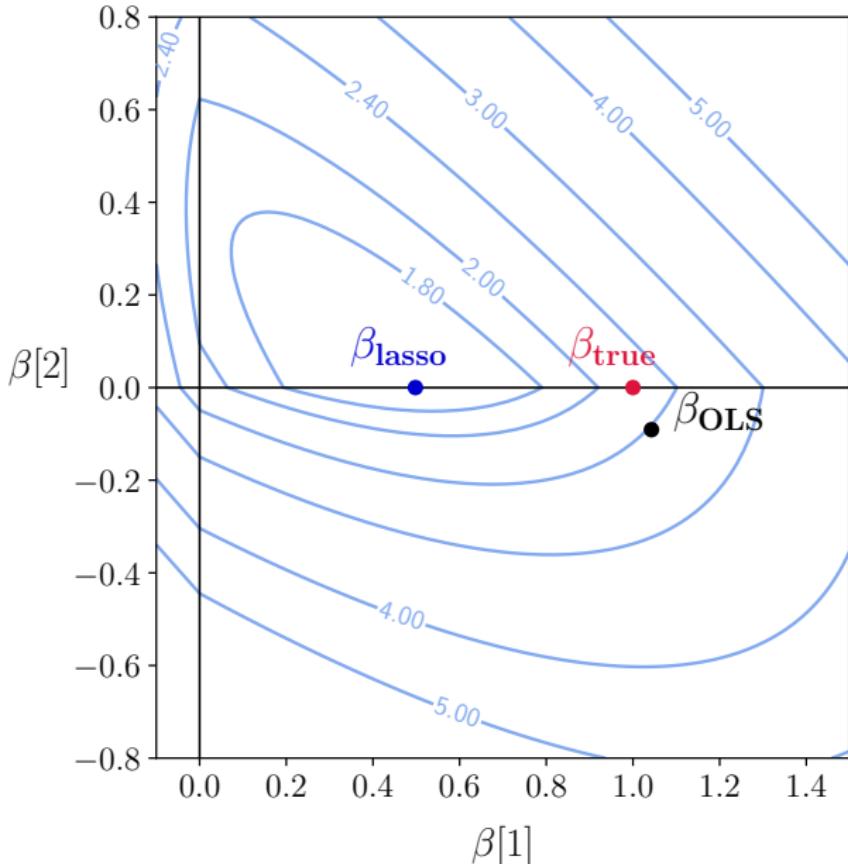
$$\frac{1}{2} \left\| X^T \beta - y \right\|_2^2 + \lambda \|\beta\|_1$$



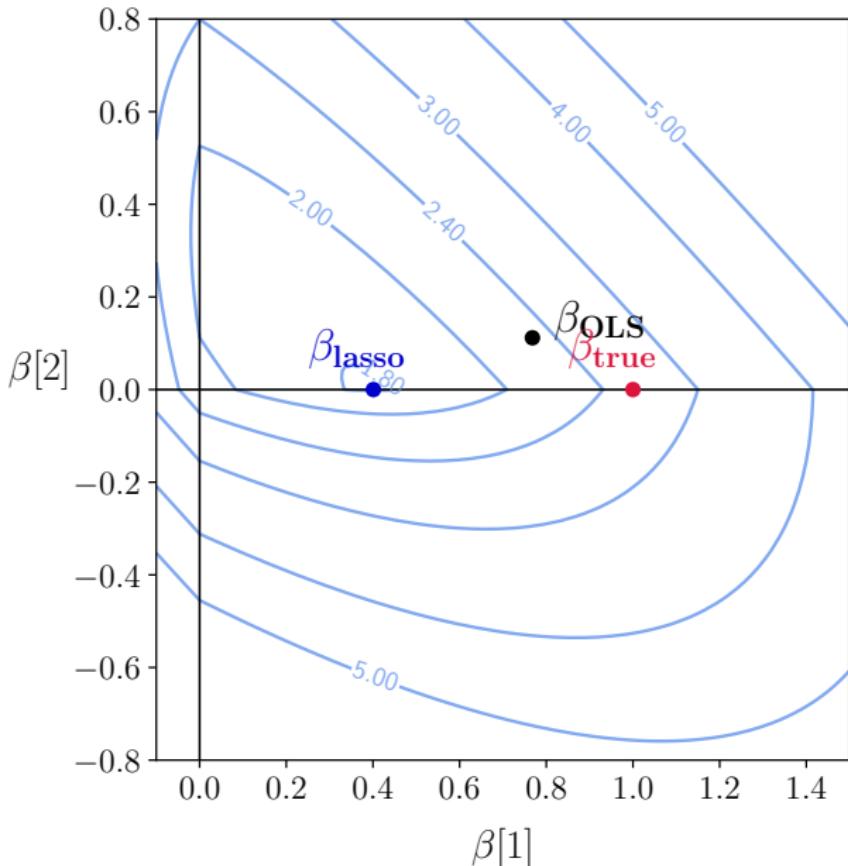
$$\frac{1}{2} \left\| X^T \beta - y \right\|_2^2 + \lambda \|\beta\|_1$$



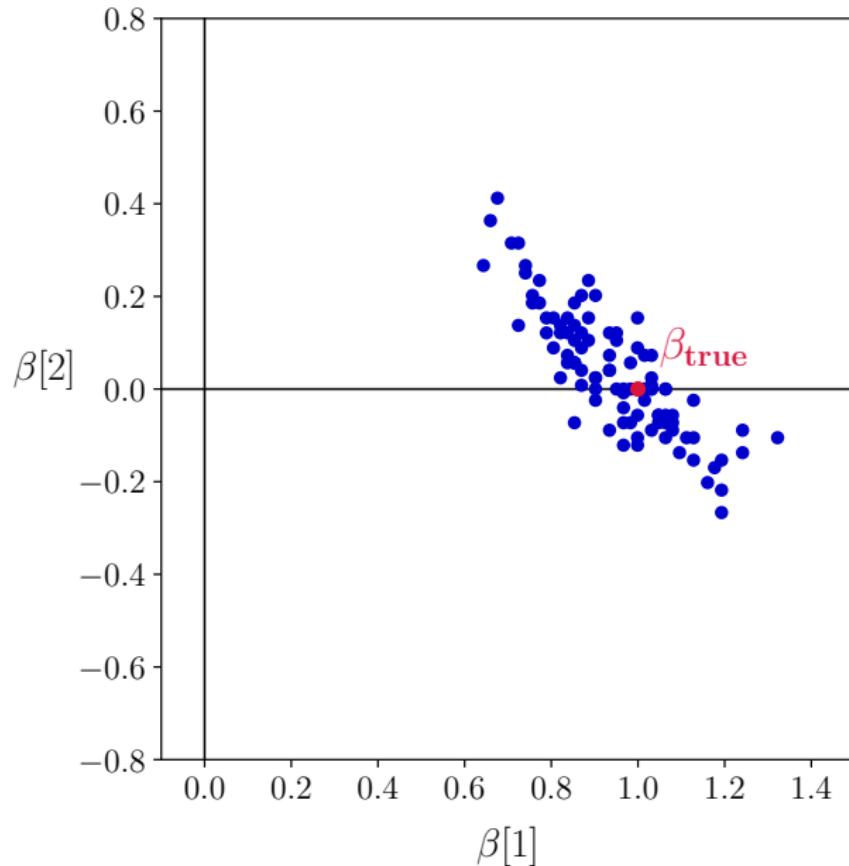
$$\frac{1}{2} \left\| X^T \beta - y \right\|_2^2 + \lambda \|\beta\|_1$$



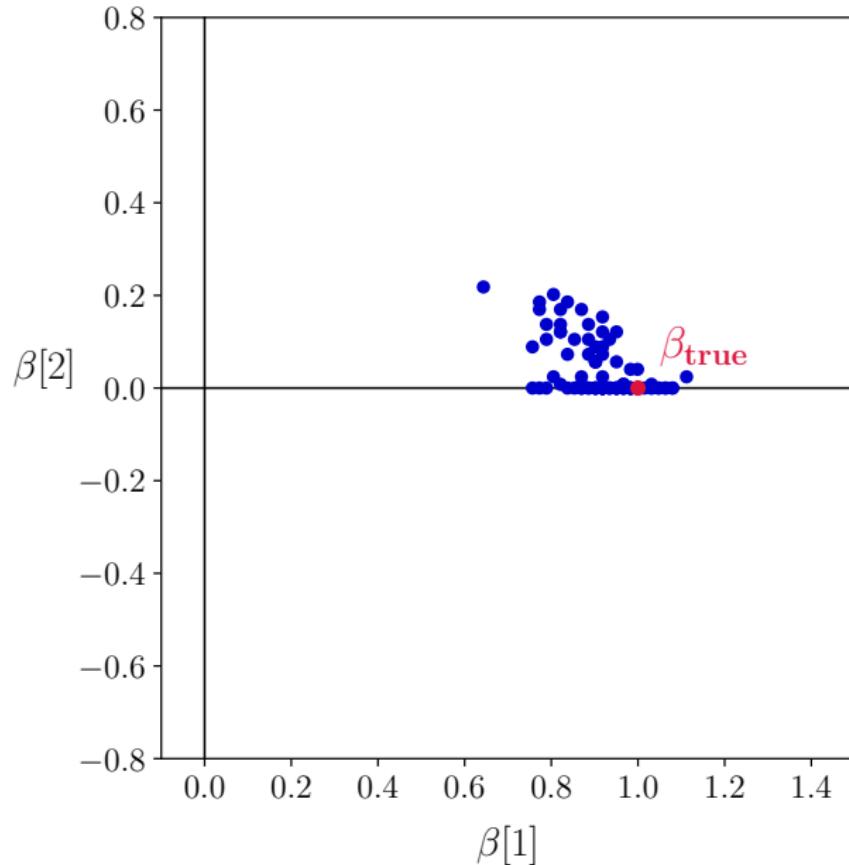
$$\frac{1}{2} \left\| X^T \beta - y \right\|_2^2 + \lambda \|\beta\|_1$$



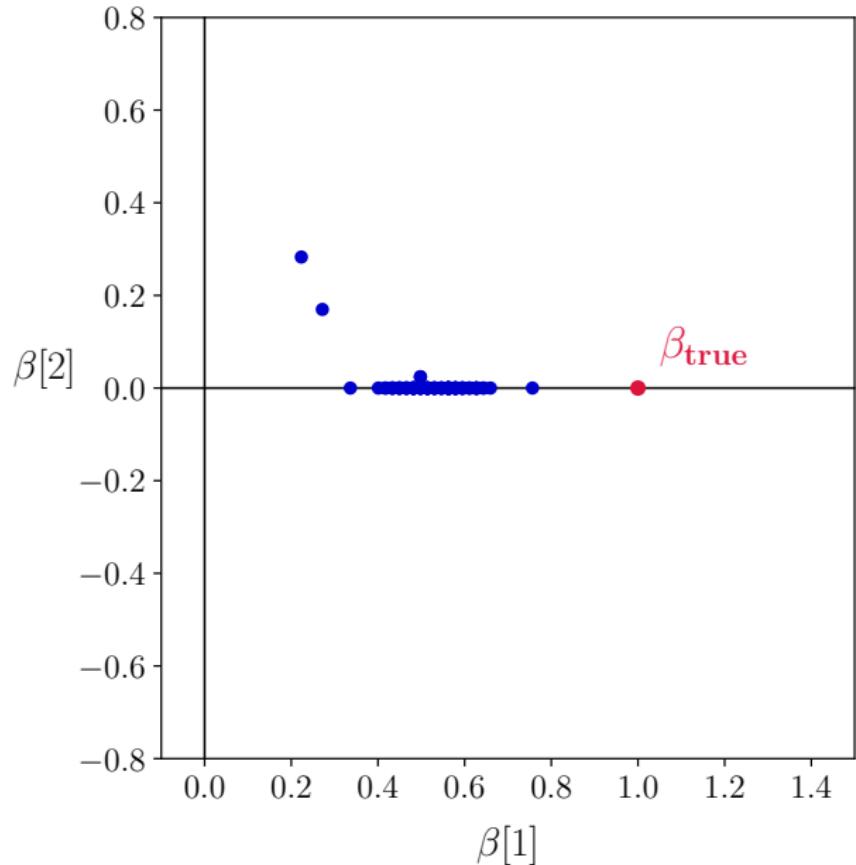
$$\lambda = 0.02$$



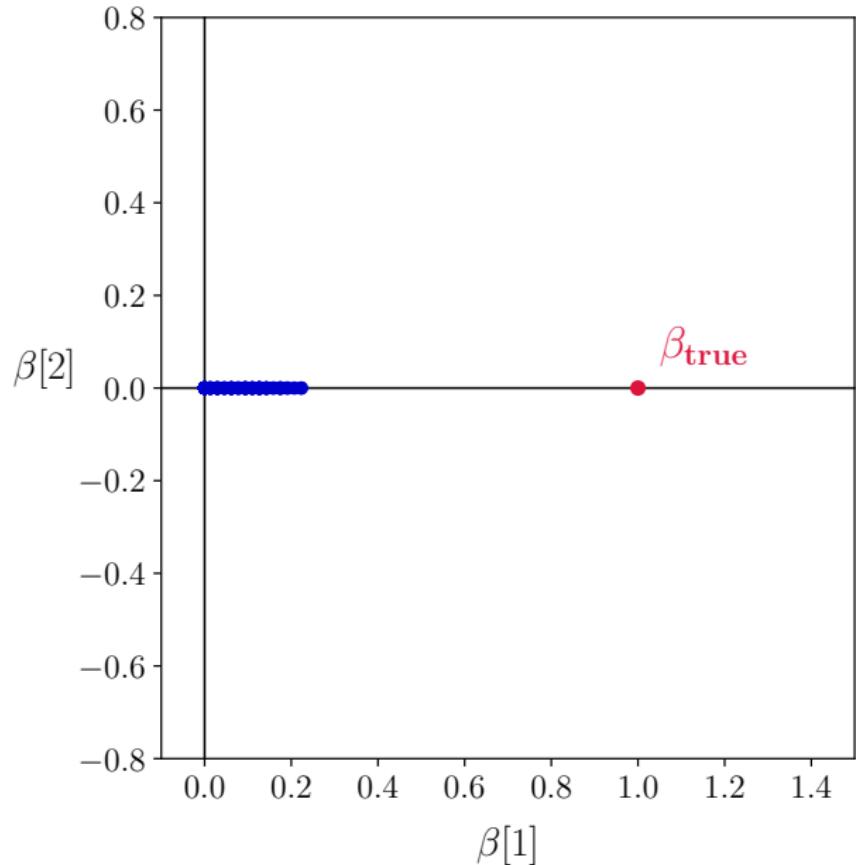
$$\lambda = 0.2$$



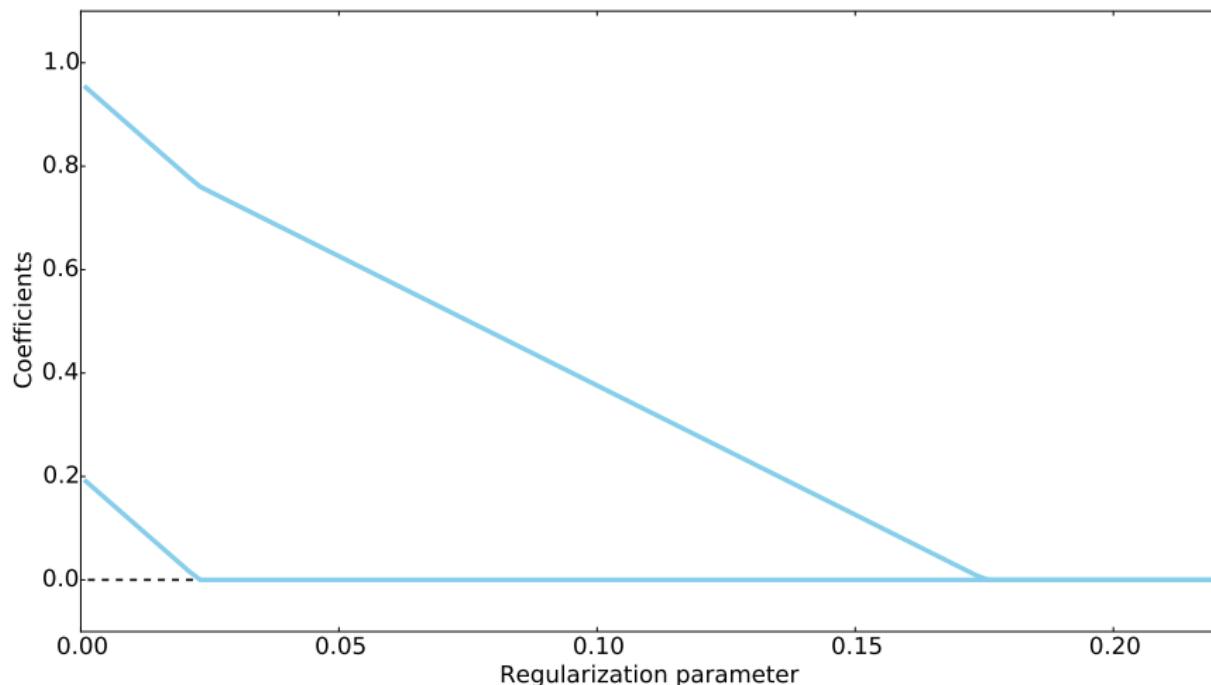
$$\lambda = 2$$



$$\lambda = 4$$



Coefficients



Lasso

$$\min_{\beta} \frac{1}{2} \left\| X^T \beta - y \right\|_2^2 + \lambda \|\beta\|_1$$

Does this function have a unique minimum? Why?

How can we prove that the minimum is sparse?

Lasso

$$\min_{\beta} \frac{1}{2} \|X^T \beta - y\|_2^2 + \lambda \|\beta\|_1$$

Does this function have a unique minimum? Why?

Yes if X is full rank and $n \geq p$ because it is strictly convex

How can we prove that the minimum is sparse?

Show that $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$ is a subgradient of the cost function at $\begin{bmatrix} \alpha \\ 0 \end{bmatrix}$ for some $\alpha > 0$

Subgradients

The **subgradient** of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at $x \in \mathbb{R}^n$ is a vector $g \in \mathbb{R}^n$ such that

$$f(y) \geq f(x) + g^T (y - x), \quad \text{for all } y \in \mathbb{R}^n$$

Subgradients

The **subgradient** of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at $x \in \mathbb{R}^n$ is a vector $g \in \mathbb{R}^n$ such that

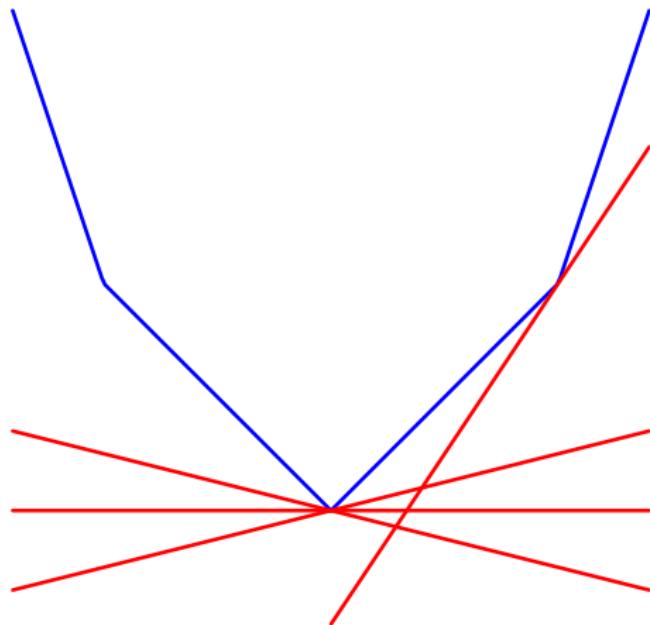
$$f(y) \geq f(x) + g^T (y - x), \quad \text{for all } y \in \mathbb{R}^n$$

The hyperplane

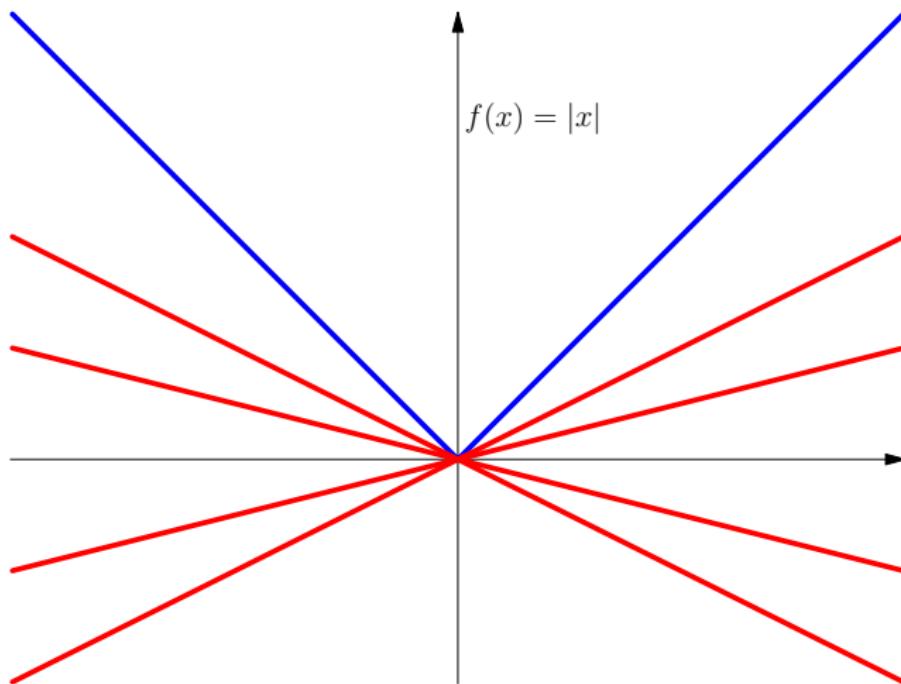
$$\mathcal{H}_g := \left\{ y \mid y[n+1] = f(x) + g^T \begin{pmatrix} [y[1]] \\ \dots \\ [y[n]] \end{pmatrix} - x \right\}$$

is a supporting hyperplane of the epigraph of f at $\begin{bmatrix} x \\ f(x) \end{bmatrix}$

Subgradients



Subgradients of absolute value



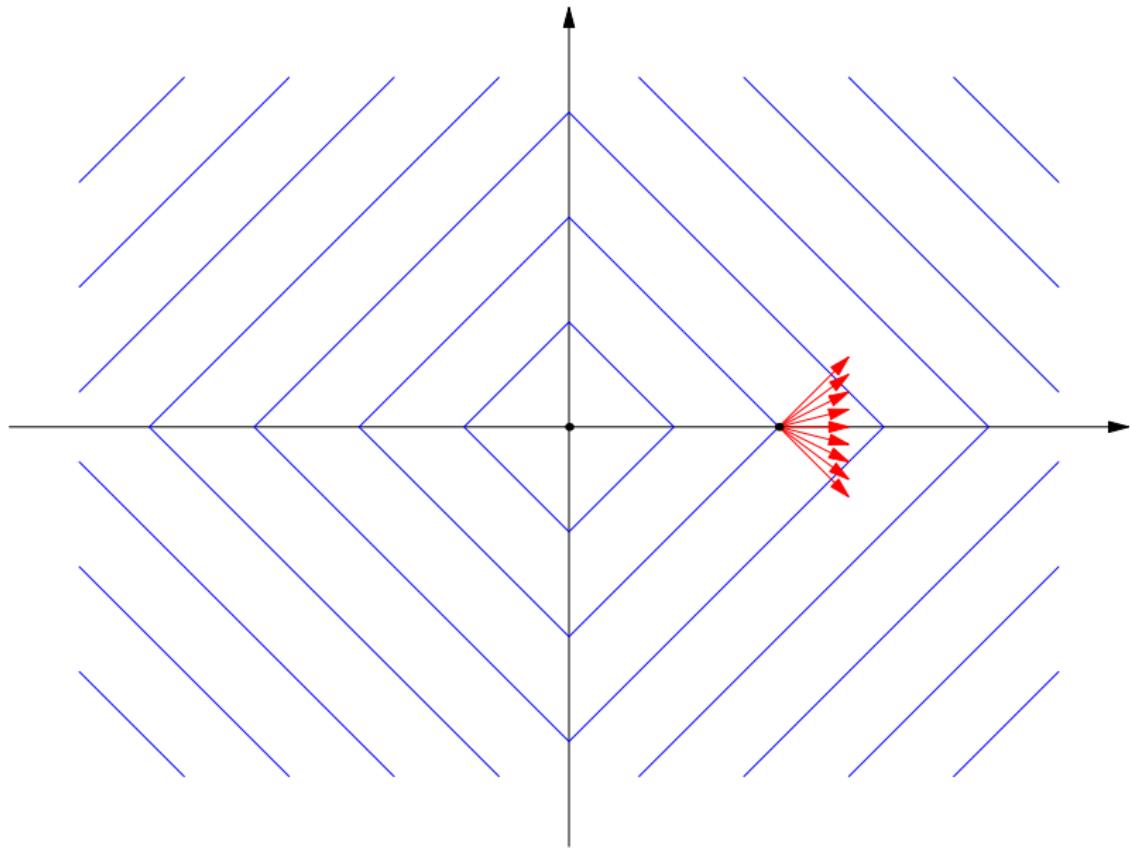
Subgradients of ℓ_1 norm

g is a subgradient of the ℓ_1 norm at $x \in \mathbb{R}^n$ if and only if

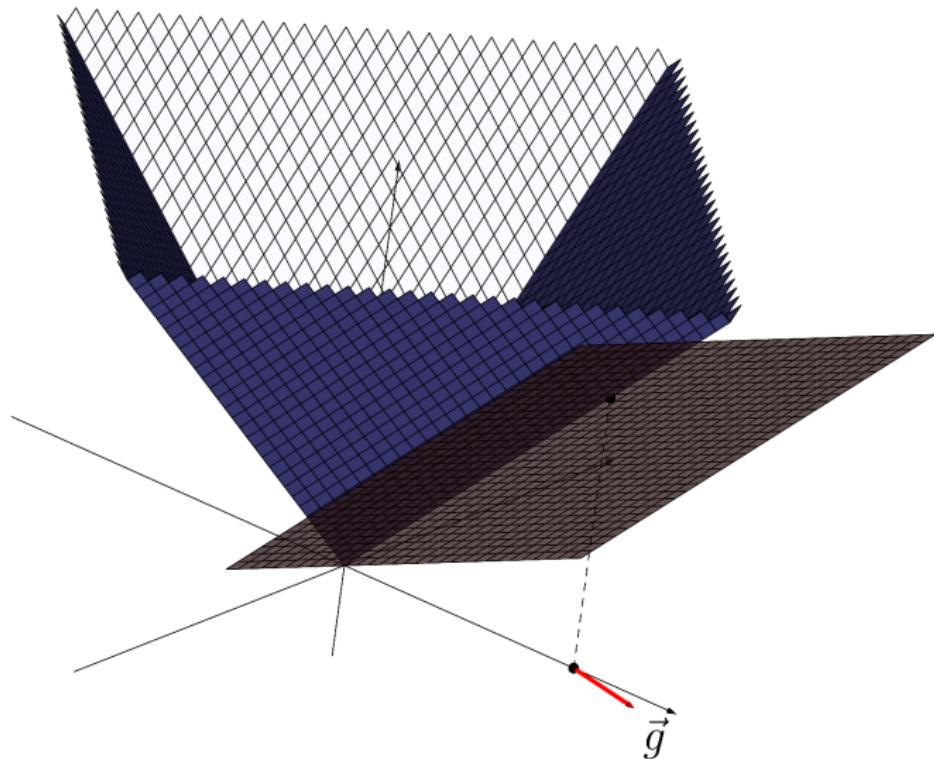
$$g[i] = \text{sign}(x[i]) \quad \text{if } x[i] \neq 0$$

$$|g[i]| \leq 1 \quad \text{if } x[i] = 0$$

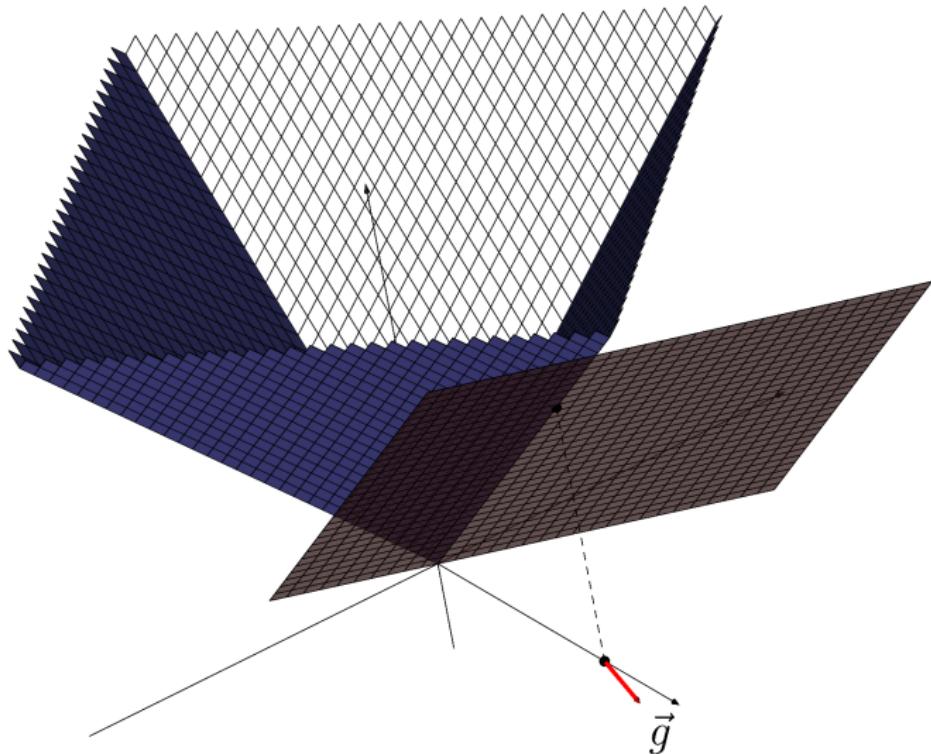
Subgradients of ℓ_1 norm



Subgradients of ℓ_1 norm



Subgradients of ℓ_1 norm



Subgradients of lasso cost function

Gradient of $\frac{1}{2} \|X^T \beta - y\|_2^2$ at β_{lasso} :

Subgradient of ℓ_1 norm at $\begin{bmatrix} \alpha \\ 0 \end{bmatrix}$ for some $\alpha > 0$:

Subgradient of lasso cost function at $\begin{bmatrix} \alpha \\ 0 \end{bmatrix}$ for some $\alpha > 0$:

Subgradients of lasso cost function

Gradient of $\frac{1}{2} \|X^T \beta - y\|_2^2$ at β_{lasso} :

$$X(X^T \beta_{\text{lasso}} - y)$$

Subgradient of ℓ_1 norm at $\begin{bmatrix} \alpha \\ 0 \end{bmatrix}$ for some $\alpha > 0$:

$$g_{\ell_1} := \begin{bmatrix} 1 \\ \gamma \end{bmatrix} \quad |\gamma| \leq 1$$

Subgradient of lasso cost function at $\begin{bmatrix} \alpha \\ 0 \end{bmatrix}$ for some $\alpha > 0$:

$$g_{\text{lasso}} := X \left(X^T \begin{bmatrix} \alpha \\ 0 \end{bmatrix} - y \right) + \lambda \begin{bmatrix} 1 \\ \gamma \end{bmatrix} \quad |\gamma| \leq 1$$

Subgradients of lasso cost function

At $\begin{bmatrix} \alpha \\ 0 \end{bmatrix}$ for some $\alpha > 0$

$$g_{\text{lasso}} := X \left(X^T \begin{bmatrix} \alpha \\ 0 \end{bmatrix} - y \right) + \lambda \begin{bmatrix} 1 \\ \gamma \end{bmatrix}$$

Subgradients of lasso cost function

At $\begin{bmatrix} \alpha \\ 0 \end{bmatrix}$ for some $\alpha > 0$

$$\begin{aligned} g_{\text{lasso}} &:= X \left(X^T \begin{bmatrix} \alpha \\ 0 \end{bmatrix} - y \right) + \lambda \begin{bmatrix} 1 \\ \gamma \end{bmatrix} \\ &= X(\alpha a - a - z) + \lambda \begin{bmatrix} 1 \\ \gamma \end{bmatrix} \\ &= \begin{bmatrix} a^T ((\alpha - 1)a - z) + \lambda \\ b^T ((\alpha - 1)a - z) + \lambda \gamma \end{bmatrix} \\ &= \begin{bmatrix} \alpha - 1 - a^T z + \lambda \\ \rho(\alpha - 1) - b^T z + \lambda \gamma \end{bmatrix} \end{aligned}$$

Is zero a subgradient for some α ?

Is zero a subgradient for some α ?

Setting $g_{\text{lasso}} = 0$

$$\begin{aligned}\alpha &= 1 - \lambda + a^T z \\ \gamma &= \frac{\rho + b^T z - \rho \alpha}{\lambda} \\ &= \frac{b^T z - \rho a^T z}{\lambda} + \rho\end{aligned}$$

We need $\alpha \geq 0$

$$\lambda \leq 1 + a^T z$$

We need $|\gamma| \leq 1$

$$\frac{|b^T z - \rho a^T z|}{1 - |\rho|} \leq \lambda$$