

## Logistics:

- HW0 graded, for regrade request submit it through GradeScope within 7 days from release of grade.
- HW1 due Tuesday Jan 25th midnight

# Lecture 9:

## Simple variable selection:

## LASSO for sparse regression

---

- Yet another hyper-parameter/family of model classes, but with a special property
  - # of features in polynomial regression
  - Regularization coefficient  $\lambda$  for ridge regression
  - Regularization coefficient  $\lambda$  for LASSO



# Sparsity

---

$$\hat{w}_{LS} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

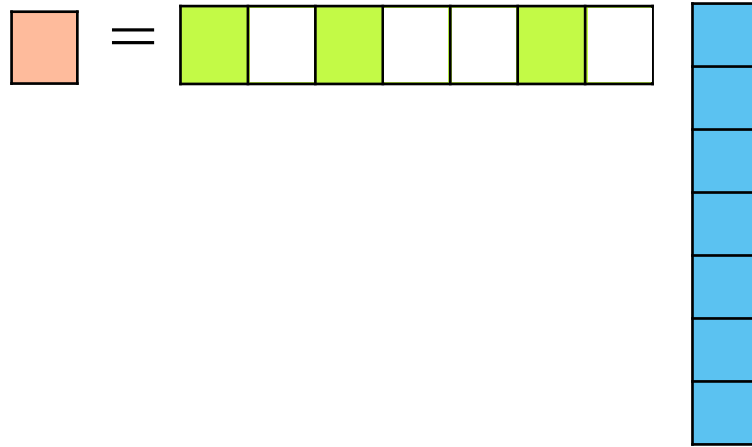
- Vector  $w$  is **sparse**, if many entries are zero
  - A vector  $w$  is said to be  $k$ -sparse if at most  $k$  entries are non-zero
  - We are interested in  $k$ -sparse  $w$  with  $k \ll d$
  - Why do we prefer sparse vector  $w$  in practice?

# Sparsity

$$\hat{w}_{LS} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

- Vector  $w$  is **sparse**, if many entries are zero
  - Efficiency**: If  $\text{size}(w) = 100$  Billion, each prediction  $w^T x$  is expensive:
    - If  $w$  is sparse, prediction computation only depends on number of non-zeros in  $w$

$$\hat{y}_i = \hat{w}_{LS}^T x_i$$



$$= \sum_{j=1}^d \hat{w}_{LS}[j] \times x_i[j] = \sum_{j: \hat{w}_{LS}[j] \neq 0} \hat{w}_{LS}[j] \times x_i[j]$$

Computational complexity decreases from  $2d$  to  $2k$  for  $k$ -sparse  $\hat{w}_{LS}$

# Sparsity

$$\hat{w}_{LS} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

- Vector  $w$  is **sparse**, if many entries are zero
  - Interpretability**: What are the relevant features to make a prediction?



Lot size
Single Family
Year built
Last sold price
Last sale price/sqft
Finished sqft
Unfinished sqft
Finished basement sqft
# floors
Flooring types
Parking type
Parking amount
Cooling
Heating
Exterior materials
Roof type
Structure style

Dishwasher
Garbage disposal
Microwave
Range / Oven
Refrigerator
Washer
Dryer
Laundry location
Heating type
Jetted Tub
Deck
Fenced Yard
Lawn
Garden
Sprinkler System

- How do we find “best” subset of features useful in predicting the price among all possible combinations?

# Finding best subset of features that explain the outcome/label: Exhaustive

- Try all subsets of size 1, 2, 3, ... and one that minimizes validation error
  - Problem?
  - Any Ideas?

# Finding best subset: Greedy

## Forward stepwise:

Starting from simple model and iteratively add features most useful to fit

### Forward Greedy

1:  $T \leftarrow \emptyset$

2: **For**  $j = 1, \dots, k$  **do**

3:  $j^* \leftarrow \arg \min_{\ell} \min_w \sum_{i=1}^n \left( y_i - \sum_{j \in T \cup \{\ell\}} w[j] \times x_i[j] \right)^2$

4:  $T \leftarrow T \cup \{j^*\}$

## Backward stepwise:

Start with full model and iteratively remove features least useful to fit

## Combining forward and backward steps:

In forward algorithm, insert steps to remove features no longer as important

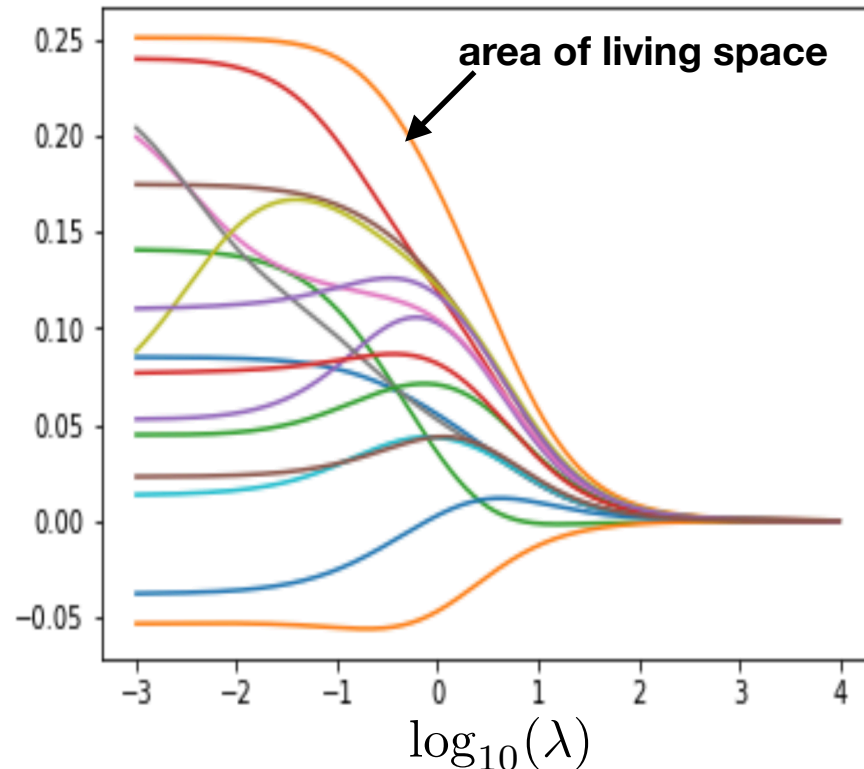
*Lots of other variants, too.*

# Finding best subset: Regularize

**Recall** that Ridge regression makes coefficients small

$$\hat{w}_{ridge} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 + \lambda ||w||_2^2$$

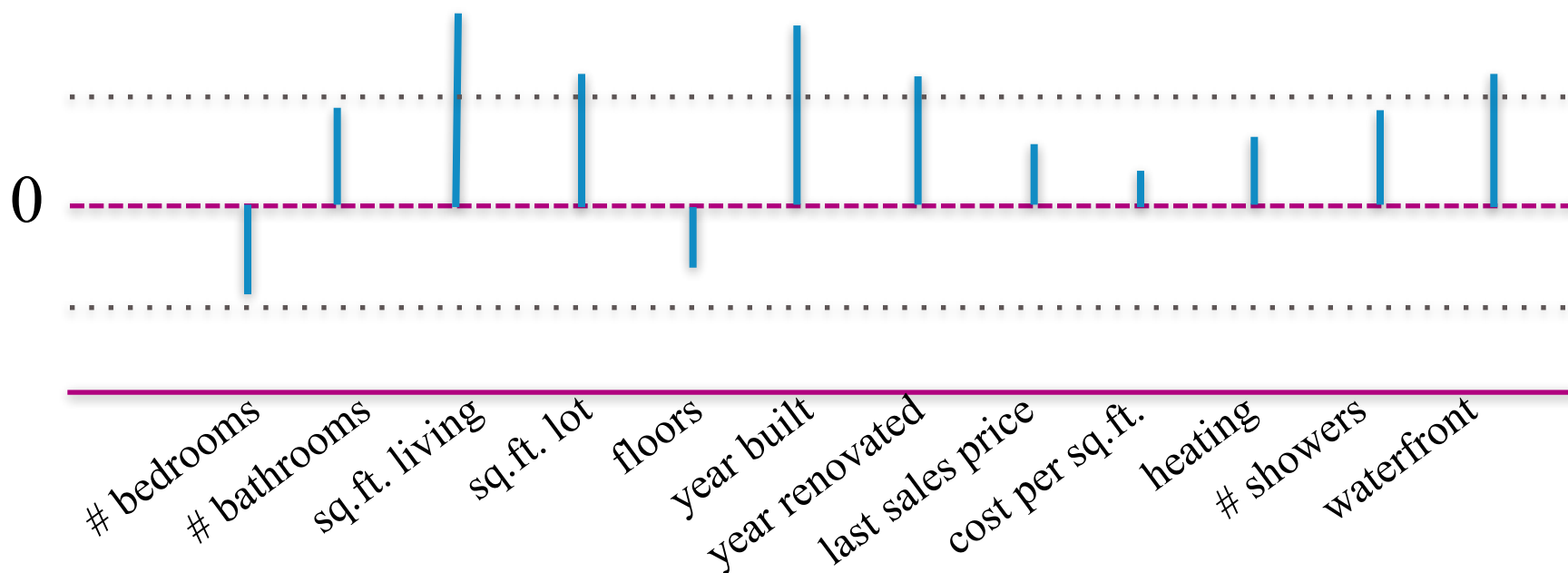
$w_i$ 's



# Thresholded Ridge Regression

$$\hat{w}_{ridge} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 + \lambda ||w||_2^2$$

- Why don't we just set **small** ridge coefficients to 0?
  - Any issues?

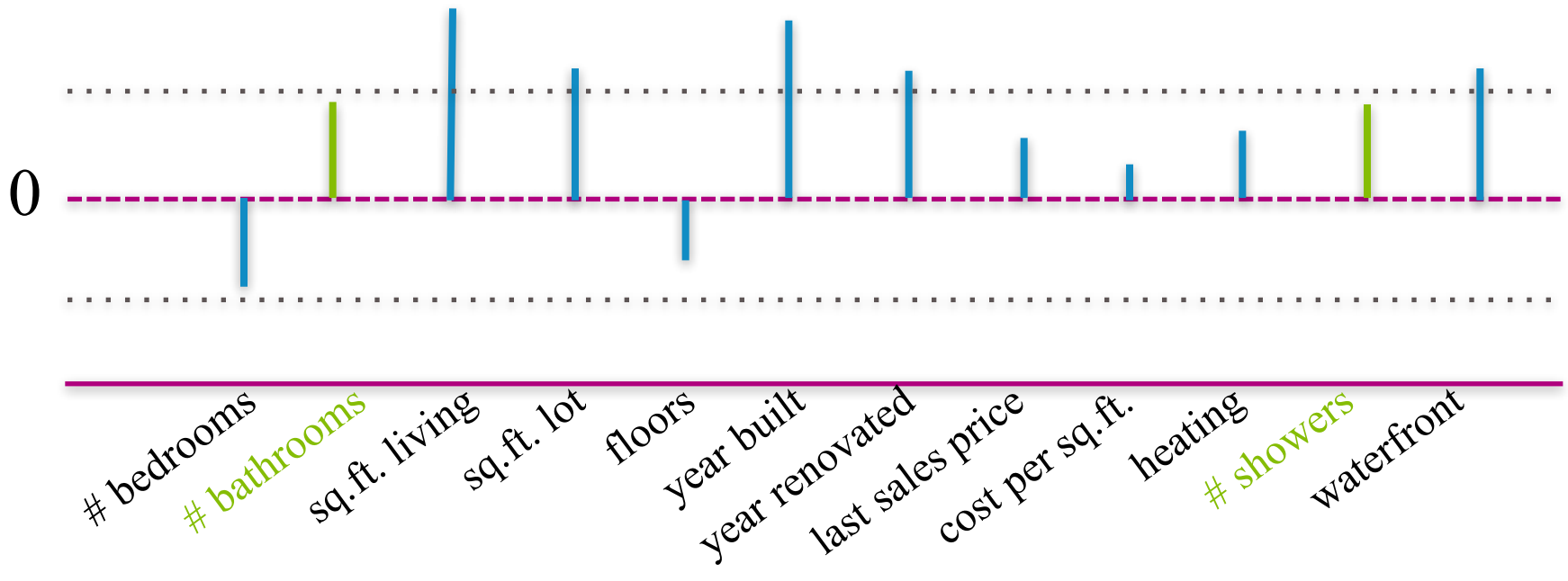




# Thresholded Ridge Regression

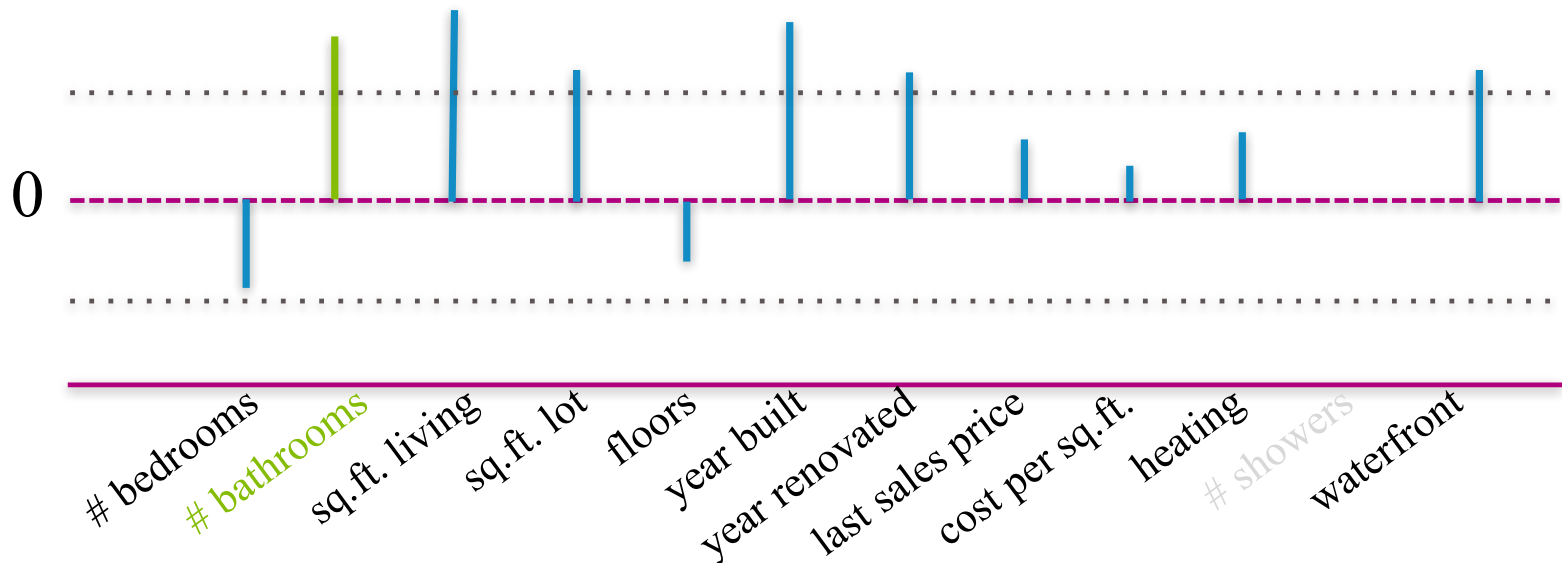
$$\hat{w}_{ridge} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 + \lambda ||w||_2^2$$

- Consider two **related** features (bathrooms, showers)
- Consider  $w[\text{bath}] = 1$  and  $w[\text{shower}] = 1$ , and  
 $w[\text{bath}] = 2$  and  $w[\text{shower}] = 0$ ,  
which one does ridge regression choose?  
(assuming #bathroom=#showers in every house)



# Thresholded Ridge Regression

- Consider two **related** features (bathrooms, showers)
- Issue with thresholded ridge regression is that ridge regression prefers balanced weights between similar features
- What if we **didn't** include showers? Weight on bathrooms increases, and it should have been selected.
- We want a feature selection scheme that selects one of (#bathroom) or (#showers) automatically, using the fact that if you delete #showers #bathroom is an important feature



- There is a better regularizer for sparse regression, that can perform the feature selection automatically.

# Ridge vs. Lasso Regression

- Recall Ridge Regression objective:

$$\hat{w}_{ridge} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 + \lambda \|w\|_2^2$$

- sensitivity of a model  $w$  is measured in squared  $\ell_2$  norm  $\|w\|_2^2$
- A principled method to get sparse model is **Lasso** with regularized objective:

$$\hat{w}_{lasso} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 + \lambda \|w\|_1$$

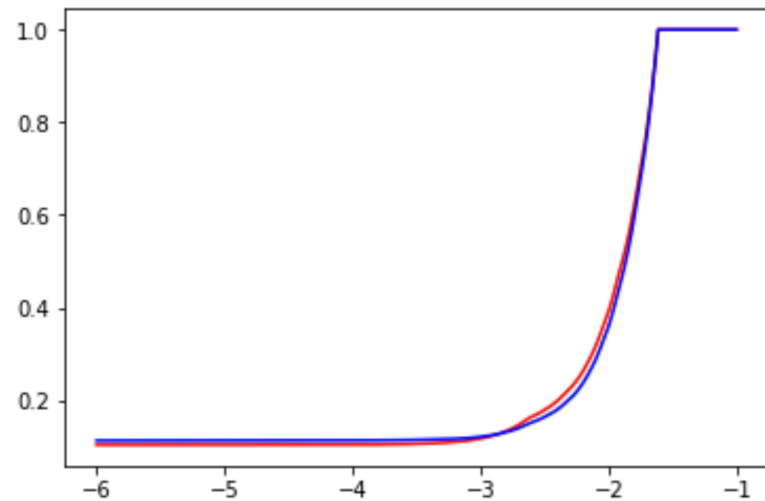
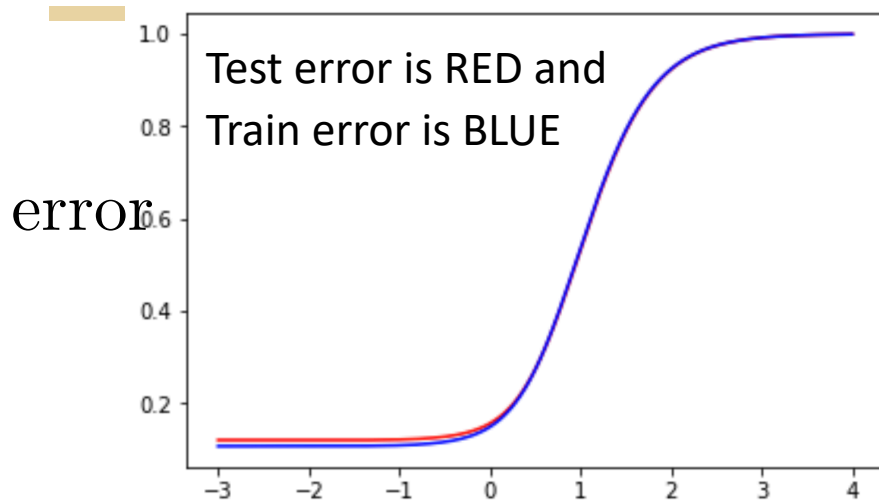
- sensitivity of a model  $w$  is measured in  $\ell_1$  norm:

$$\|w\|_1 = \sum_{j=1}^d |w[j]|$$

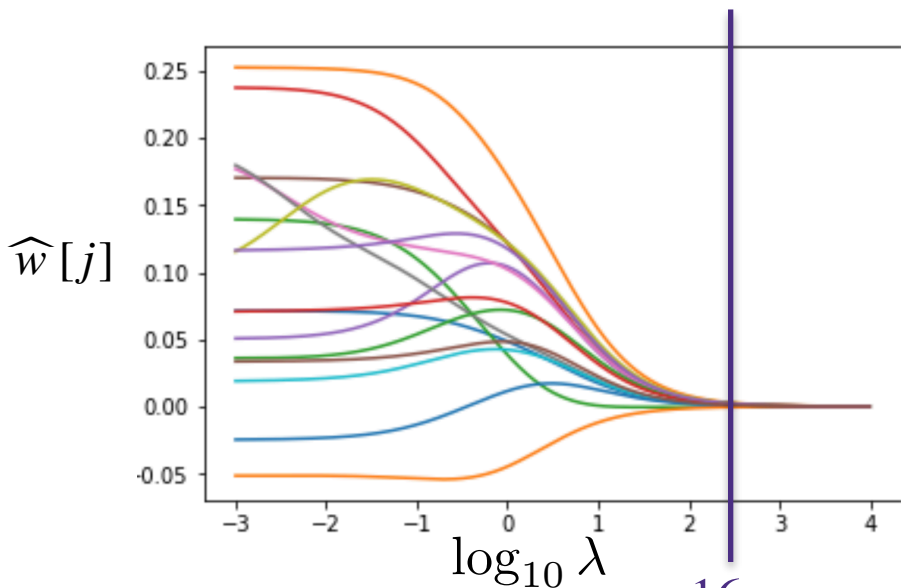
$\ell_p$ -norm of a vector  $w \in \mathbb{R}^d$  is

$$\|w\|_p \triangleq \left( \sum_{j=1}^d |w[j]|^p \right)^{1/p}$$

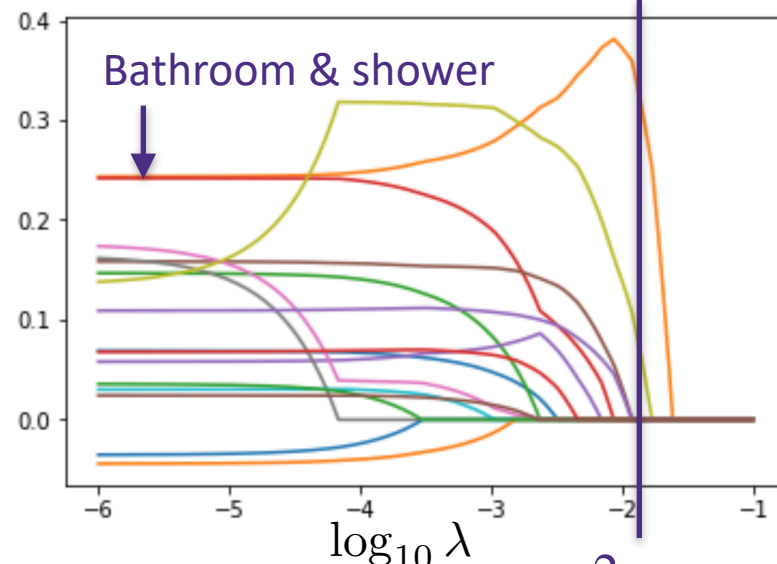
# Example: house price with 16 features



- Regularization path for Lasso shows that weights drop to exactly zero as  $\lambda$  increases



Ridge regression



Lasso regression

# Lasso regression naturally gives sparse features

- **feature selection** with Lasso regression
  1. **Model selection**: choose  $\lambda$  based on cross validation error
  2. **Feature selection**: keep only those features with non-zero (or not-too-small) parameters in  $w$  at optimal  $\lambda$
  3. **retrain** with the sparse model and  $\lambda = 0$

why do we need to retrain?

# Example: piecewise-linear fit

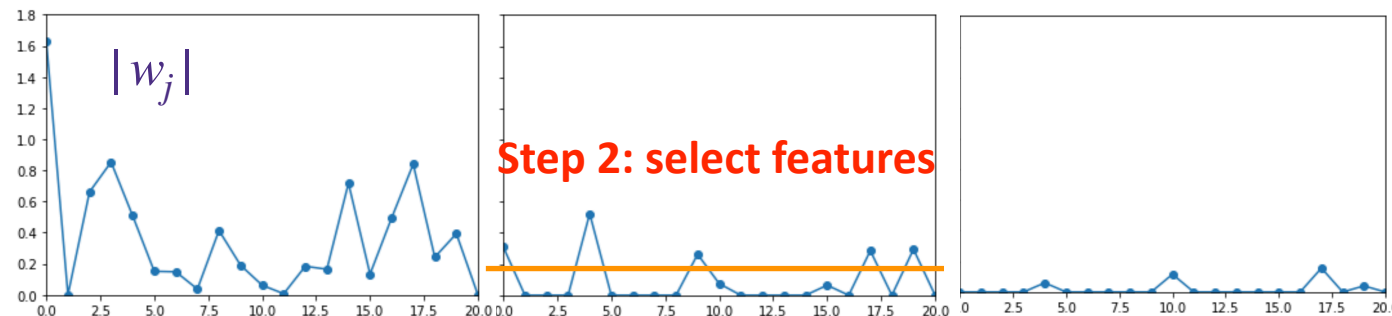
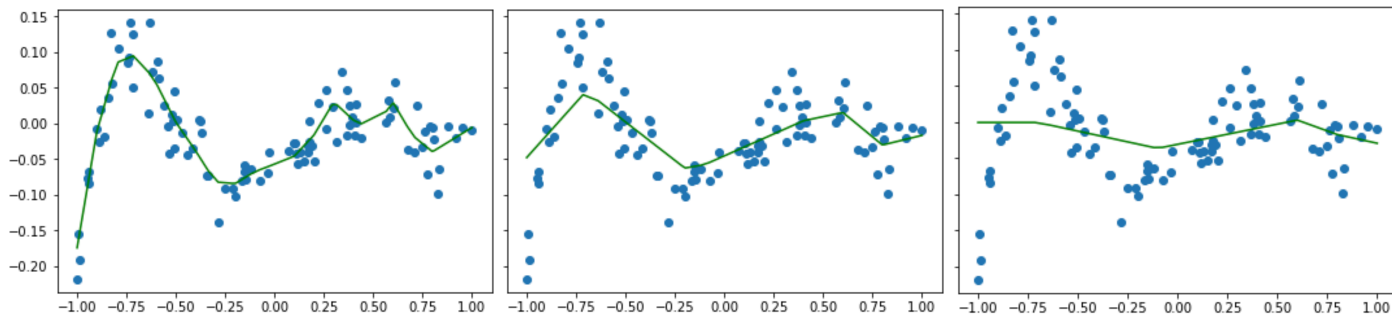
- We use Lasso on the piece-wise linear example

$$h_0(x) = 1$$

$$h_i(x) = [x + 1.1 - 0.1i]^+$$

Step 1: find optimal  $\lambda^*$

$$\text{minimize}_w \mathcal{L}(w) + \lambda \|w\|_1$$



$$\lambda = 10^{-8}$$

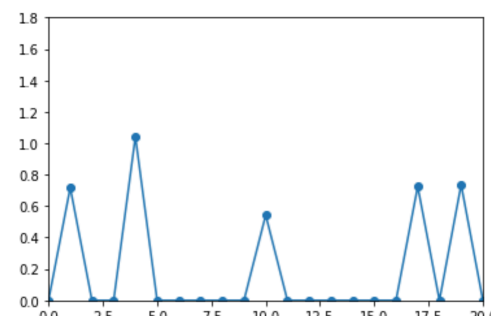
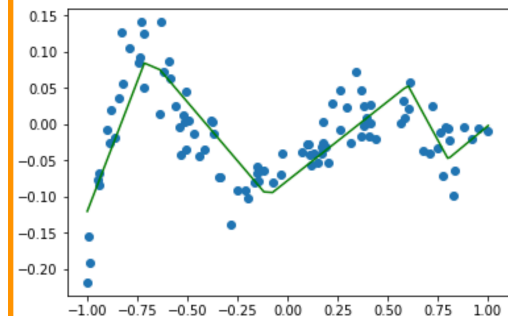
$$\lambda = 10^{-4}$$

$$\lambda = 2 \times 10^{-4}$$

- de-biasing (via re-training) is critical!

Step 3: retrain

$$\text{minimize}_w \mathcal{L}(w)$$



$$\lambda = 0$$

but only use selected features

# Penalized Least Squares

---

$$\text{Ridge : } r(w) = ||w||_2^2 \qquad \text{Lasso : } r(w) = ||w||_1$$

$$\hat{w}_r = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 + \lambda r(w)$$

# Penalized Least Squares

- Regularized optimization:

$$\hat{w}_r = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 + \lambda r(w)$$

$$\text{Ridge : } r(w) = ||w||_2^2$$

$$\text{Lasso : } r(w) = ||w||_1$$

- For any  $\lambda^* \geq 0$  for which  $\hat{w}_r$  achieves the minimum, there exists a  $\mu^* \geq 0$  such that the solution of the constrained optimization,  $\hat{w}_c$ , is the same as the solution of the regularized optimization,  $\hat{w}_r$ , where

$$\hat{w}_c = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 \quad \text{subject to } r(w) \leq \mu^*$$

- so there are pairs of  $(\lambda, \mu)$  whose optimal solution  $\hat{w}_r$  are the same for the regularized optimization and constrained optimization



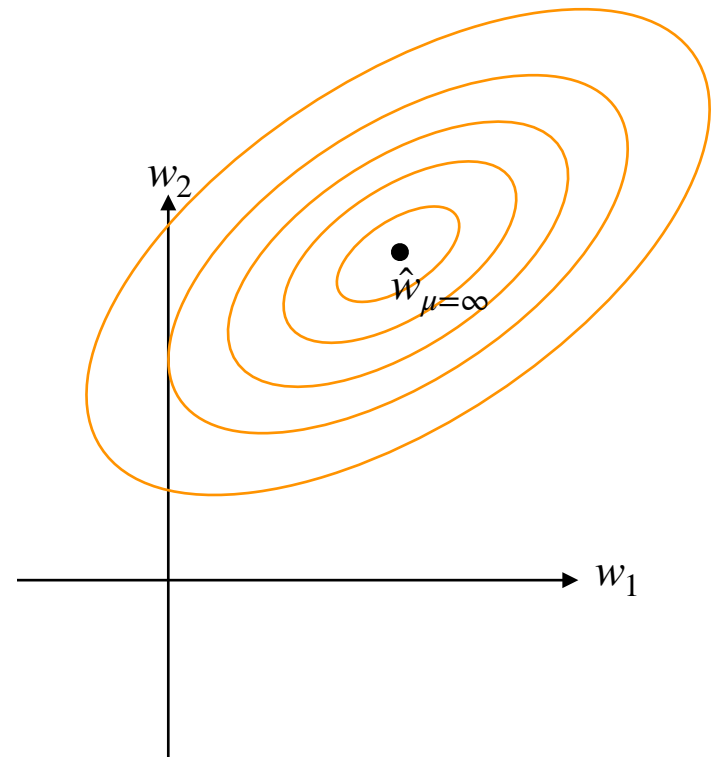
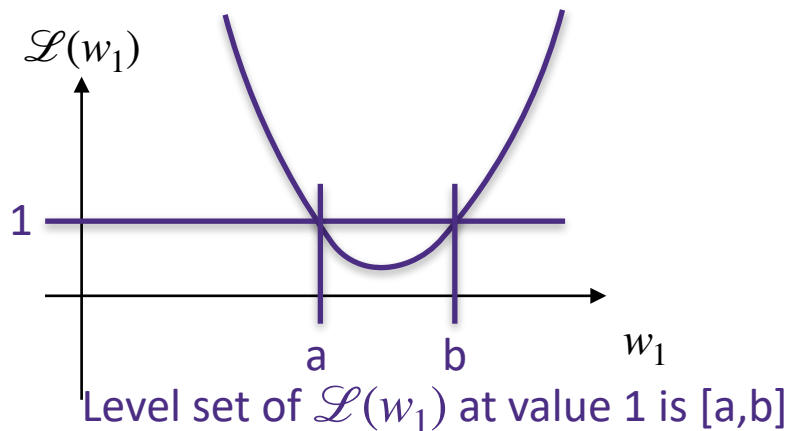
# Why does Lasso give sparse solutions?

$$\text{minimize}_w \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$\text{subject to } \|w\|_1 \leq \mu$$

- the **level set** of a function  $\mathcal{L}(w_1, w_2)$  is defined as the set of points  $(w_1, w_2)$  that have the same function value
- the level set of a quadratic function is an oval
- the center of the oval is the least squares solution  $\hat{w}_{\mu=\infty} = \hat{w}_{\text{LS}}$

1-D example with quadratic loss



# Why does Lasso give sparse solutions?

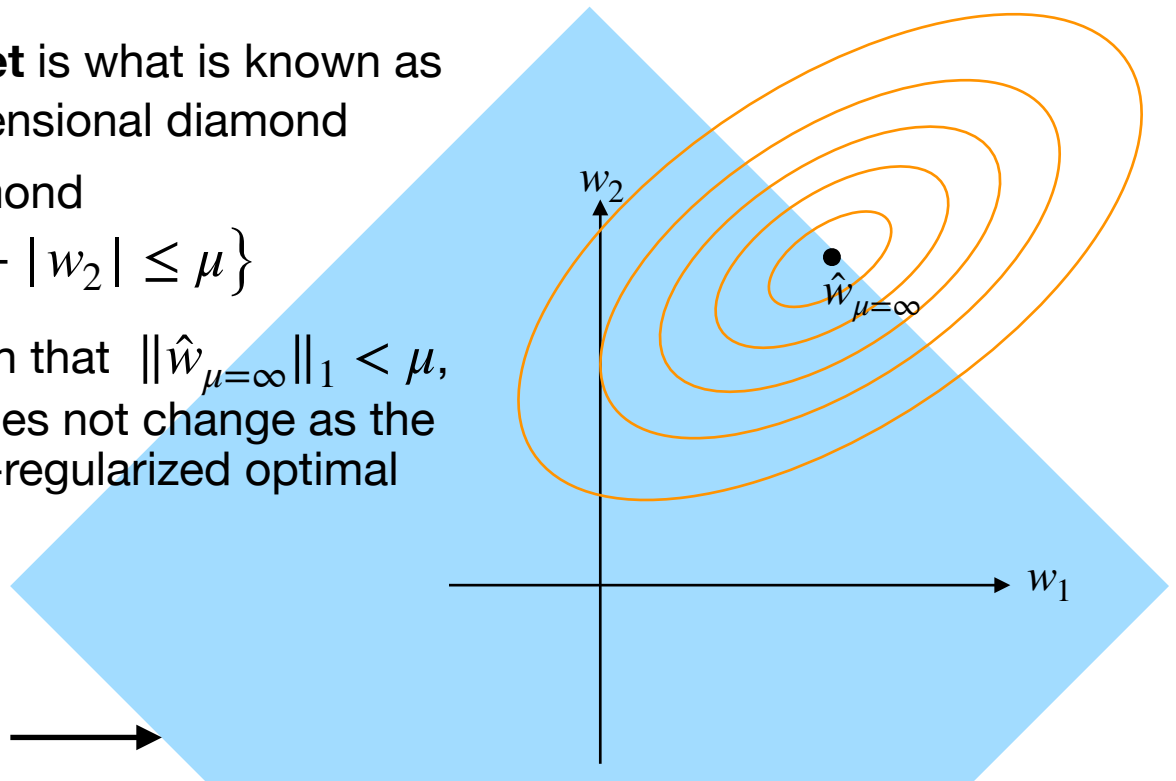
$$\text{minimize}_w \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$\text{subject to } \|w\|_1 \leq \mu$$

- as we decrease  $\mu$  from infinity, the feasible set becomes smaller
- the shape of the **feasible set** is what is known as  $L_1$  ball, which is a high dimensional diamond
- In 2-dimensions, it is a diamond

$$\{(w_1, w_2) \mid |w_1| + |w_2| \leq \mu\}$$

- when  $\mu$  is large enough such that  $\|\hat{w}_{\mu=\infty}\|_1 < \mu$ , then the optimal solution does not change as the feasible set includes the un-regularized optimal solution



**feasible set:**  $\{w \in \mathbb{R}^2 \mid \|w\|_1 \leq \mu\}$  →

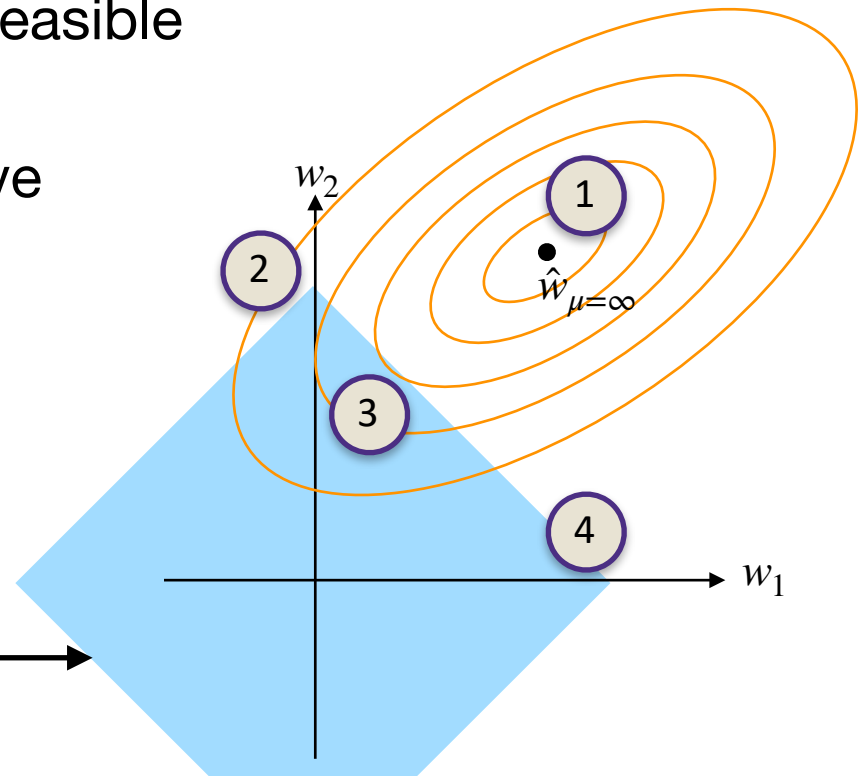
# Why does Lasso give sparse solutions?

$$\text{minimize}_w \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$\text{subject to } \|w\|_1 \leq \mu$$

- As  $\mu$  decreases (which is equivalent to increasing regularization  $\lambda$ ) the feasible set (blue diamond) shrinks
- The optimal solution of the above optimization is ?

**feasible set:**  $\{w \in \mathbb{R}^2 \mid \|w\|_1 \leq \mu\}$   $\longrightarrow$

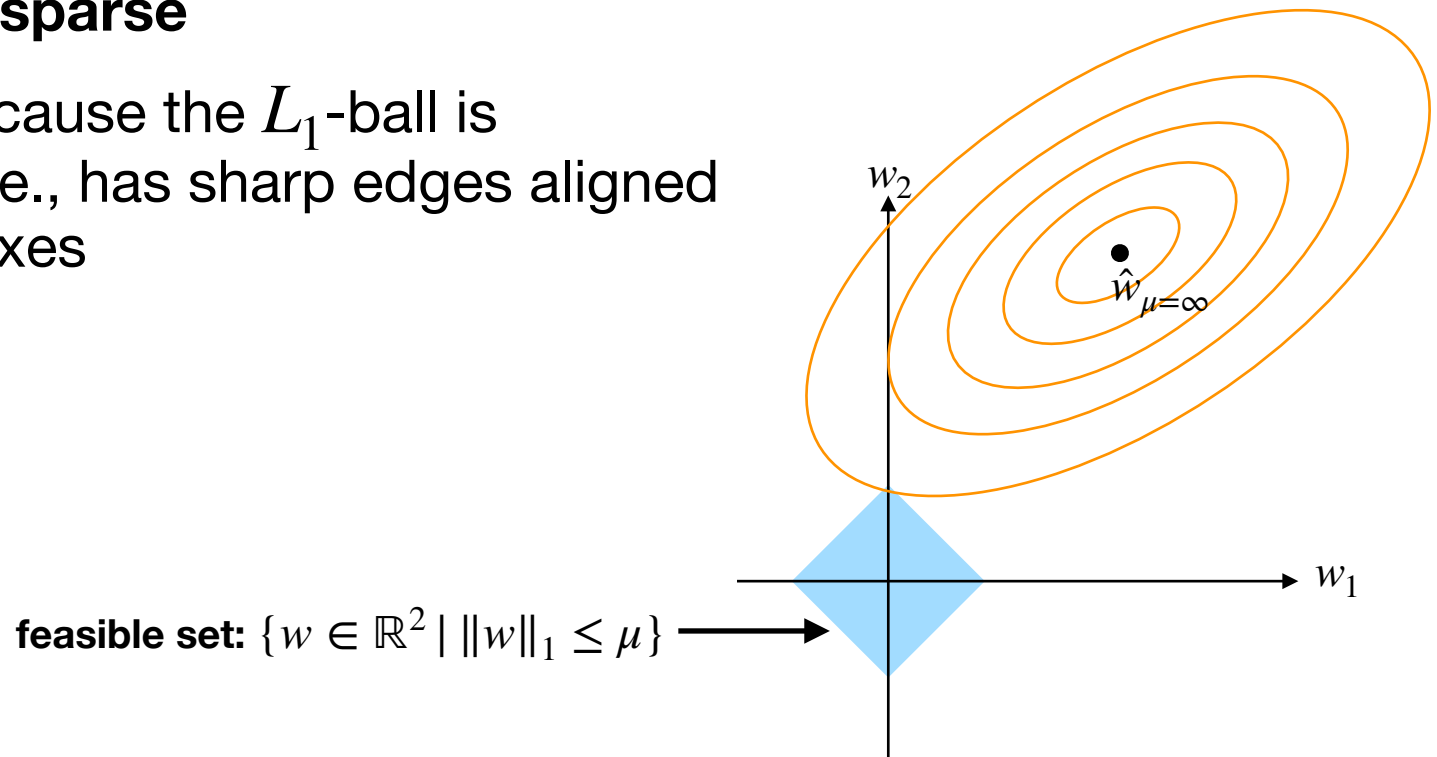


# Why does Lasso give sparse solutions?

$$\text{minimize}_w \sum_{i=1}^n (w^T x_i - y_i)^2$$

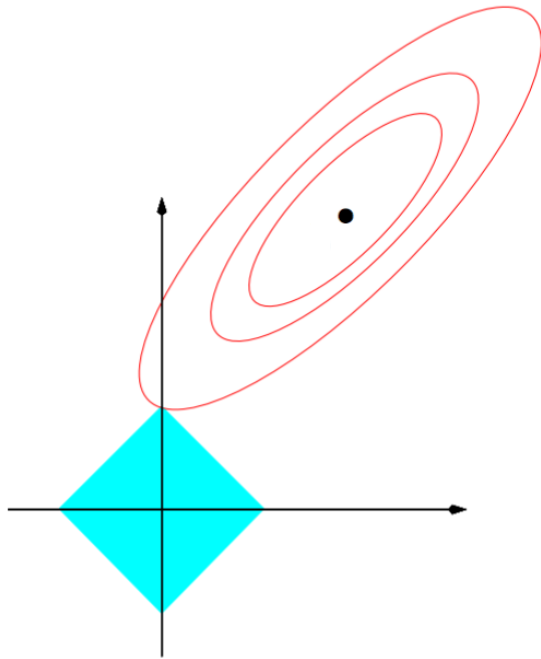
$$\text{subject to } \|w\|_1 \leq \mu$$

- For small enough  $\mu$ , the optimal solution becomes **sparse**
- This is because the  $L_1$ -ball is “pointy”, i.e., has sharp edges aligned with the axes



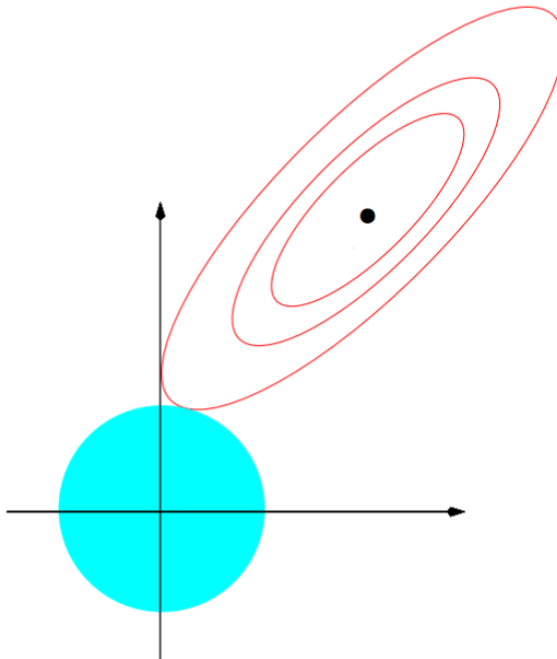
# Penalized Least Squares

- Lasso regression finds sparse solutions, as  $L_1$ -ball is “pointy”
- Ridge regression finds dense solutions, as  $L_2$ -ball is “smooth”



$$\text{minimize}_w \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$\text{subject to } \|w\|_1 \leq \mu$$



$$\text{minimize}_w \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$\text{subject to } \|w\|_2^2 \leq \mu$$

# Questions?

---