① # of features / basis functions

② $+ \lambda \|w\|_2^2$

# Simple variable (feature/) selection: LASSO for sparse regression
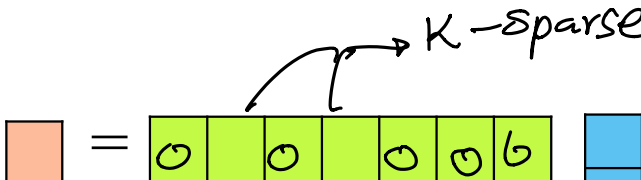
W

# Sparsity

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

- Vector $w$ is **sparse**, if many entries are zero

# Sparsity

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

- Vector $w$ is **sparse**, if many entries are zero
  - **Efficiency**: If 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 $w$

K—sparse



$\widehat{w}_{LS}^T$

$$\widehat{y}_i = \widehat{w}_{LS}^\top x_i = \sum_{j=1}^{d} x_i[j]\widehat{w}_{LS}[j]$$

$$= \sum_{j:\ \widehat{w}_{LS}\ \text{is non-zero}} x_i[j] \cdot \widehat{w}_{LS}[j]$$

: $O(K)$ computations

# Sparsity

$$\widehat{w}_{LS} = \arg\min_{w} \sum_{i=1}^{n} \left(y_i - x_i^T w\right)^2$$

- 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: **Exhaustive**

> Try all subsets of size $\overbrace{1, 2, 3, \ldots}^{d}$ and one that minimizes validation error
> Problem?

$$\sum_{k=1}^{d} \binom{d}{k} = 2^d$$

Minimum Description length.

$$Error_{CV} + \lambda \cdot \|w\|_0$$
$$= \underbrace{\phantom{xxx}}_{\# \text{ non-zero in } w.}$$

# Finding best subset: Greedy

**Forward stepwise:**
Starting from simple model and iteratively add features most useful to fit

**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.*

Forward Greedy
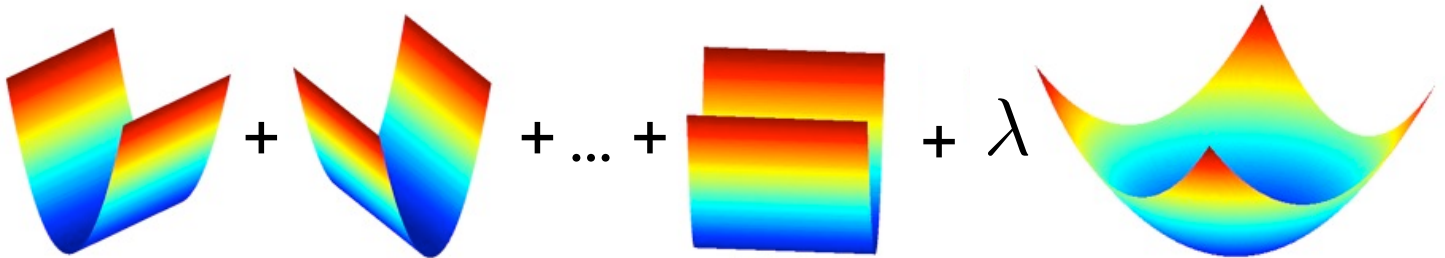
$T = \emptyset$

for $j = 1, \ldots, K$

$\quad i^* \leftarrow \arg\min_w \text{Error}_{cv}(\text{using } T \cup \{i\})$

$\quad T = T \cup \{i^*\}$

# Finding best subset: Regularize

**Ridge regression makes coefficients small**

$$\widehat{w}_{ridge} = \arg\min_{w} \underbrace{\sum_{i=1}^{n} \left(y_i - x_i^T w\right)^2}_{\text{LS-error.}} + \lambda||w||_2^2$$
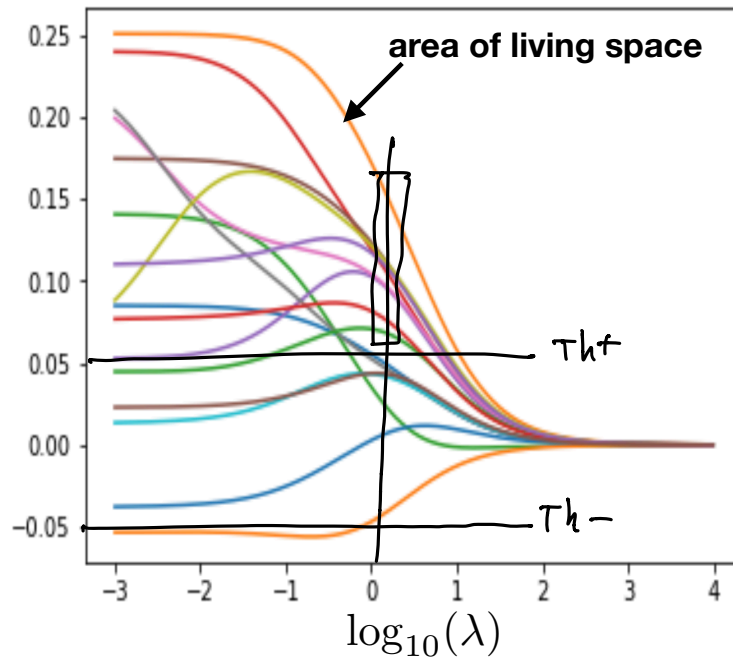
 $+$  $+ ... +$  $+ \lambda$ 

# Finding best subset: Regularize

**Ridge regression makes coefficients small**

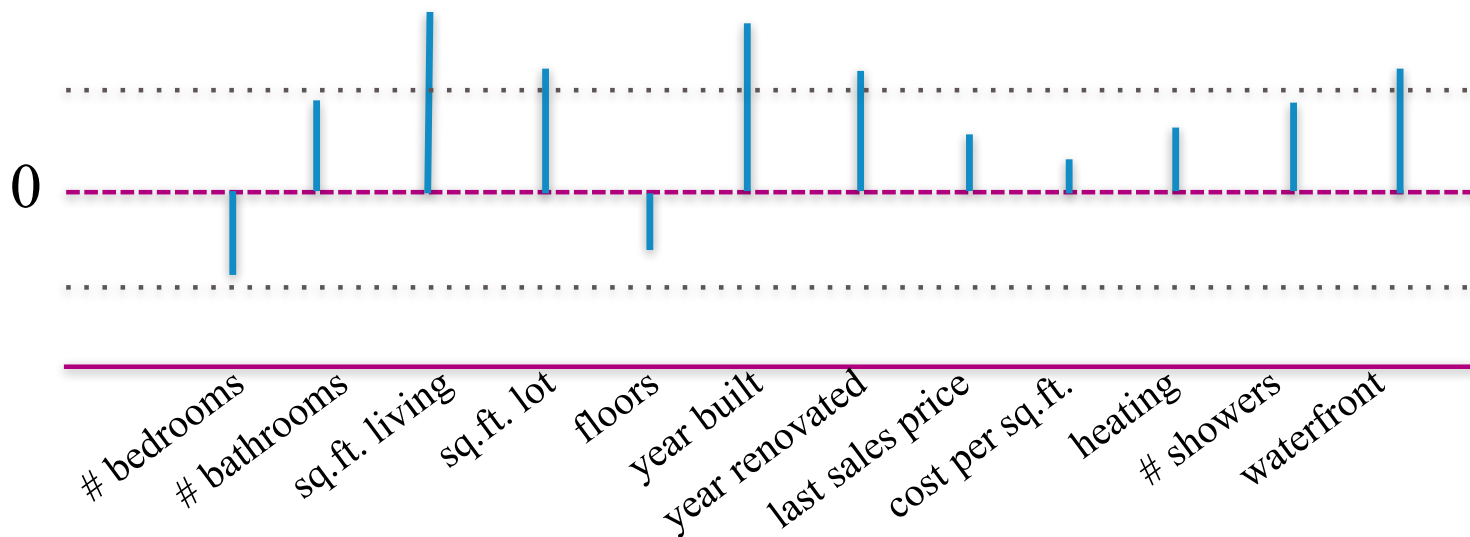$$\widehat{w}_{ridge} = \arg \min_w \sum_{i=1}^{n} \left(y_i - x_i^T w\right)^2 + \lambda||w||_2^2$$

$w_i\text{'s}$

# Thresholded Ridge Regression

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

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

# Thresholded Ridge Regression

suppose that all houses in training had same # of bathrooms as # of showers.

prediction $w_{(1)}$ $w_{bath} = 1$, $w_{shower} = 1$ → $1^2 + 1^2 = 2$ → $1 + 1 = 2$

refurbication

is same $w_{(2)}$ $w_{bath} = 2$, $w_{shower} = 0$

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

$2^2 = 4$ → $2 + 0 = 2$

Ridge $w_{(1)} > w_{(2)}$,

Lasso $w_{(1)}$ — $w_{(2)}$

Consider two related features (bathrooms, showers)

$0$

# bedrooms   # bathrooms   sq. ft. living   sq. ft. lot   floors   year built   year renovated   last sales price   cost per sq. ft.   heating   # showers   waterfront
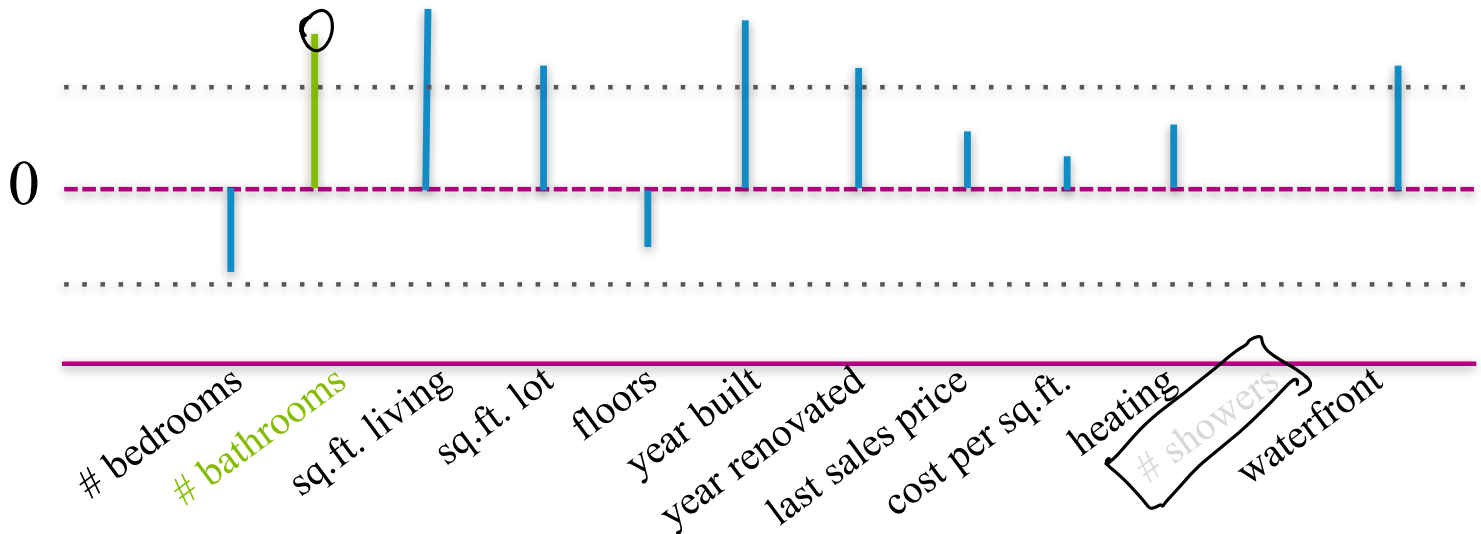
# Thresholded Ridge Regression

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

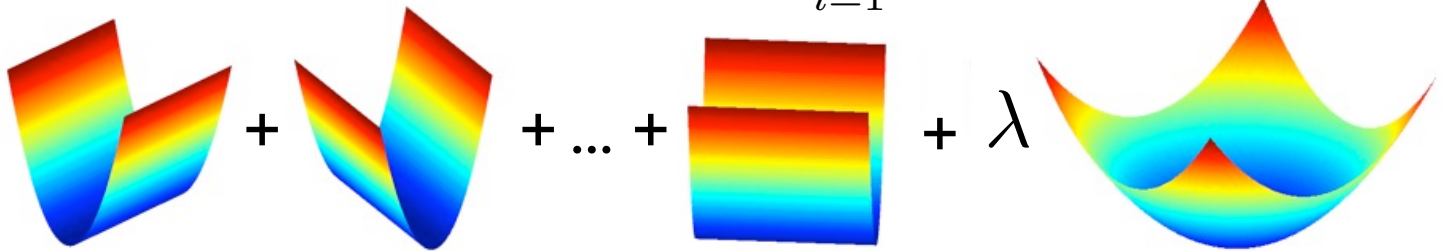What if we **didn't** include showers? Weight on bathrooms increases!



**Can another regularizer perform selection automatically?**
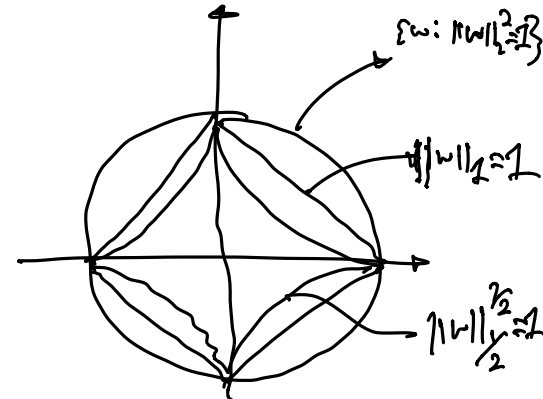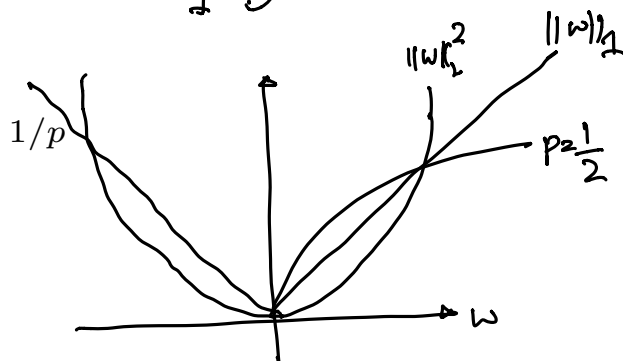
# Recall Ridge Regression

- Ridge Regression objective:

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



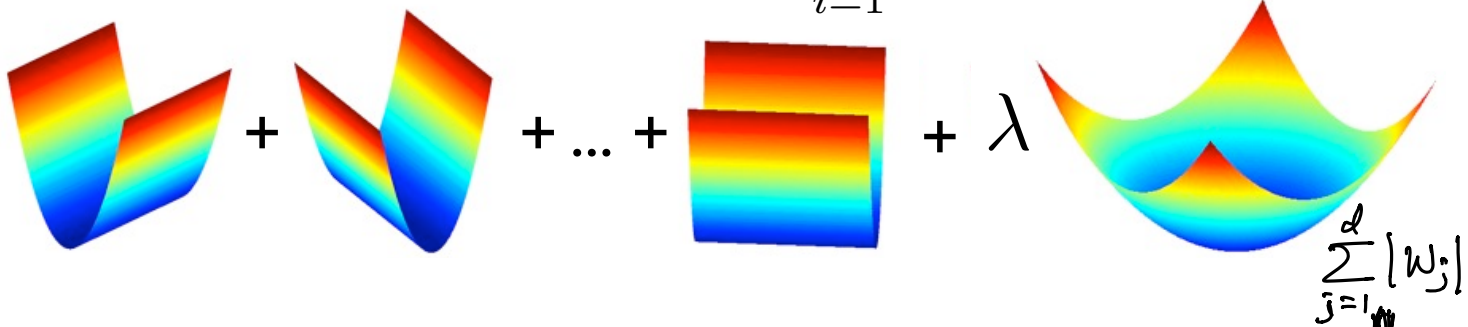1-D          2-D

$$||w||_p = \left(\sum_{i=1}^{d} |u_i|^p\right)^{1/p}$$

$||w||_2^2$     $||w||_1$

$p = \frac{1}{2}$

$w$

$||w||_1 = |w_1| + |w_2|$

$\{w : ||w||_2^2 \leq 1\}$

$||w||_1 \leq 1$

$||w||_2^{1/2} \leq 1$

# Ridge vs. Lasso Regression

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

- Ridge Regression objective:

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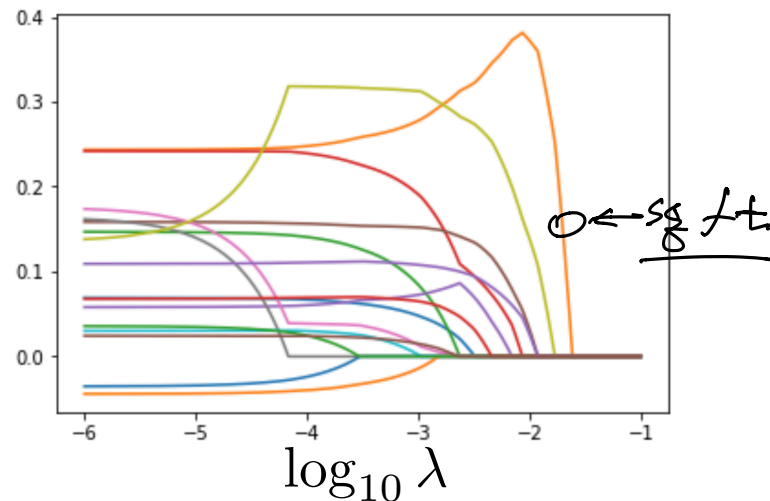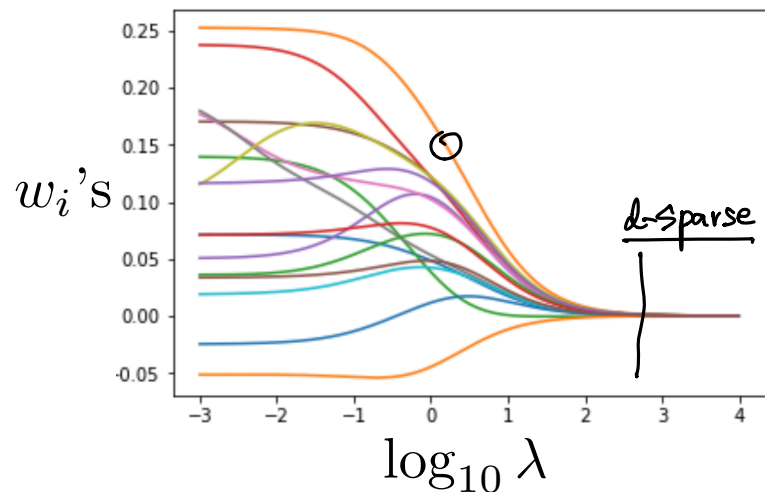


$+$ ... $+$ ... $+$ ... $+ \lambda$

$$\sum_{j=1}^{d} |w_j|$$

- <u>Lasso</u> objective:

$$\widehat{w}_{lasso} = \arg\min_{w} \sum_{i=1}^{n} \left( y_i - x_i^T w \right)^2 + \boxed{\lambda \|w\|_1}$$



$+$ ... $+$ ... $+$ ... $+ \lambda$

# Example: house price with 16 features

test error is red and train error is blue



error

$w_i$'s

$\log_{10} \lambda$

Ridge regression

$\log_{10} \lambda$

Lasso regression

# Lasso regression naturally gives sparse features

- **feature selection** with Lasso regression

    1. choose $\lambda$ based on cross validation error
    2. keep only those features with non-zero (or not-too-small) parameters in $w$ at optimal $\lambda$ $\longrightarrow$ feature selection
    3. **retrain** with the sparse model and $\lambda = 0$
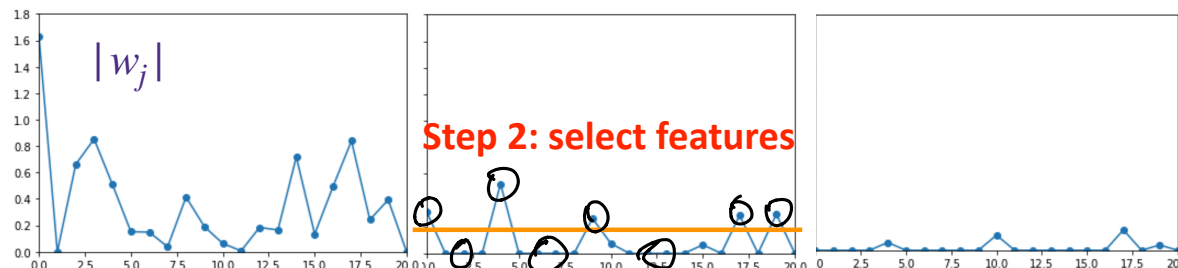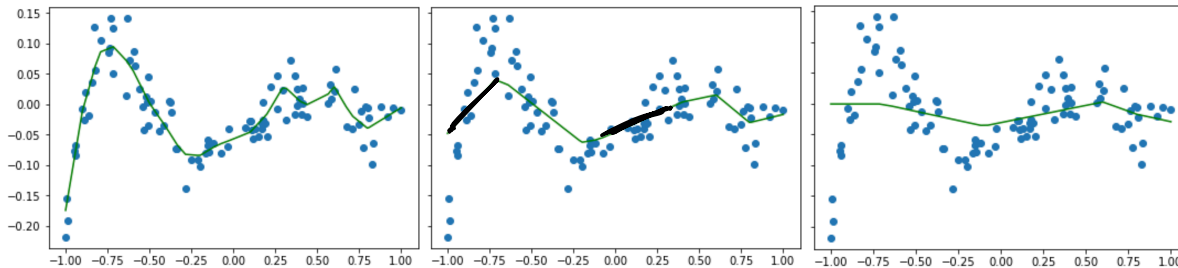
    No Regularization.

# Example: piecewise-linear fit

$$h_0(x) = 1$$

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

- We use Lasso on the piece-wise linear example

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

$$\text{minimize}_w \quad \mathscr{L}(w)$$
s.t. sparsity step-2



$|w_j|$

**Step 2: select features**

$$\lambda = 10^{-8} \qquad \lambda = 10^{-4} \qquad \lambda = 2 \times 10^{-4} \qquad \lambda = 0$$

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

but only use selected features

# Penalized Least Squares

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

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

# Penalized Least Squares

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

$$\widehat{w}_r = \arg\min_w \sum_{i=1}^{n} \left(y_i - x_i^T w\right)^2 + \lambda r(w) \quad \leftarrow \text{Penalized}$$
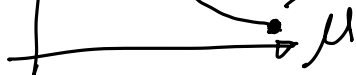
$$\widehat{w}_\lambda = \widehat{w}_\mu$$

For any $\lambda \geq 0$ for which $\hat{w}_r$ achieves the minimum, there exists a $\mu \geq 0$ such that

$$\widehat{w}_r = \arg\min_w \sum_{i=1}^{n} \left(y_i - x_i^T w\right)^2 \qquad \text{subject to } r(w) \leq \mu$$

Constrained

$\lambda$ $(\mu = 0, \lambda = \infty) \rightarrow w = 0$

$(\mu = \infty, \lambda = 0) \rightarrow w = \widehat{w}_{LS}$
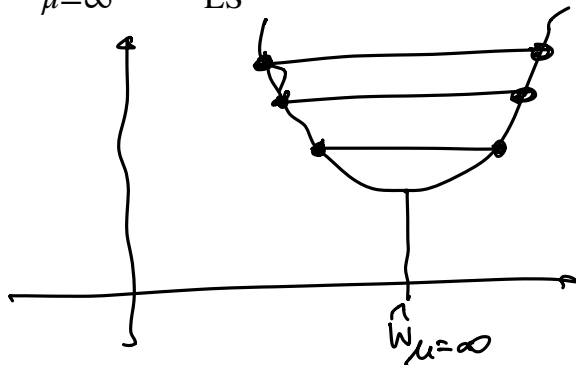
$\mu$

# Why does Lasso give sparse solutions?

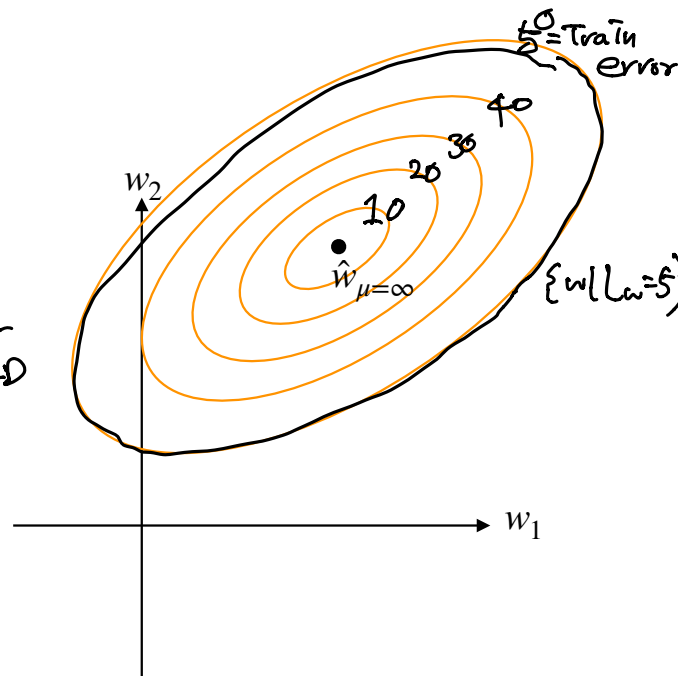$$\text{minimize}_w \quad \sum_{i=1}^{n} (w^T x_i - y_i)^2$$

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

- the **level set** of a function $\mathscr{L}(w_1, w_2)$ 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 $\hat{w}_{\mu=\infty} = \hat{w}_{LS}$

# Why does Lasso give sparse solutions?

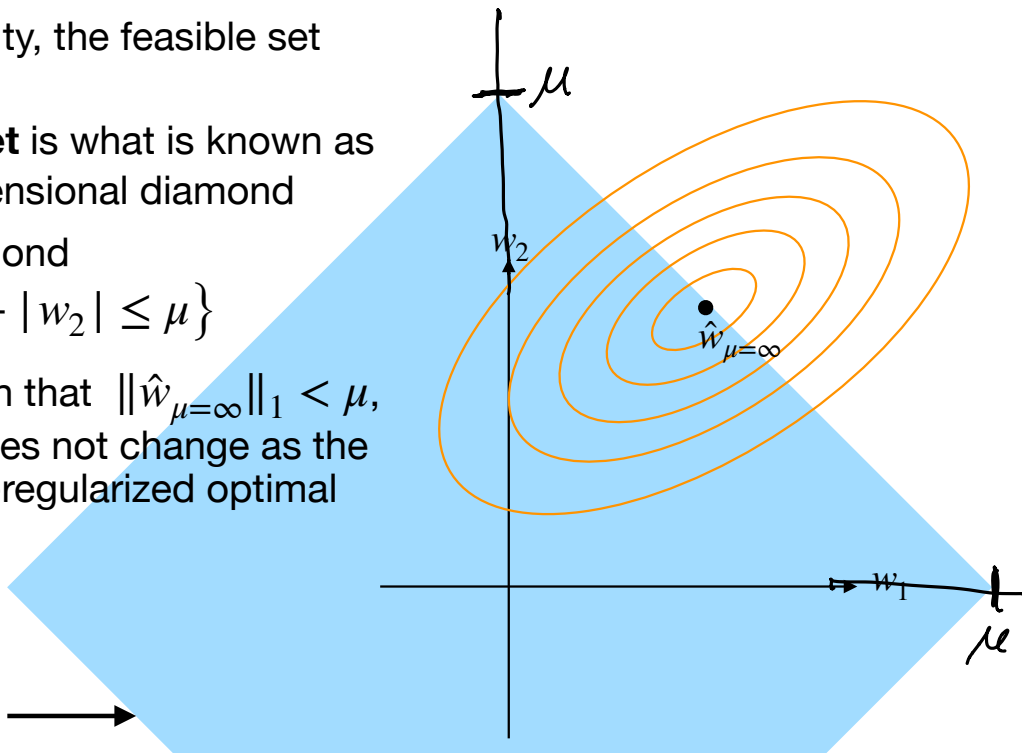$$\text{minimize}_w \quad \sum_{i=1}^{n} (w^T x_i - y_i)^2$$

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

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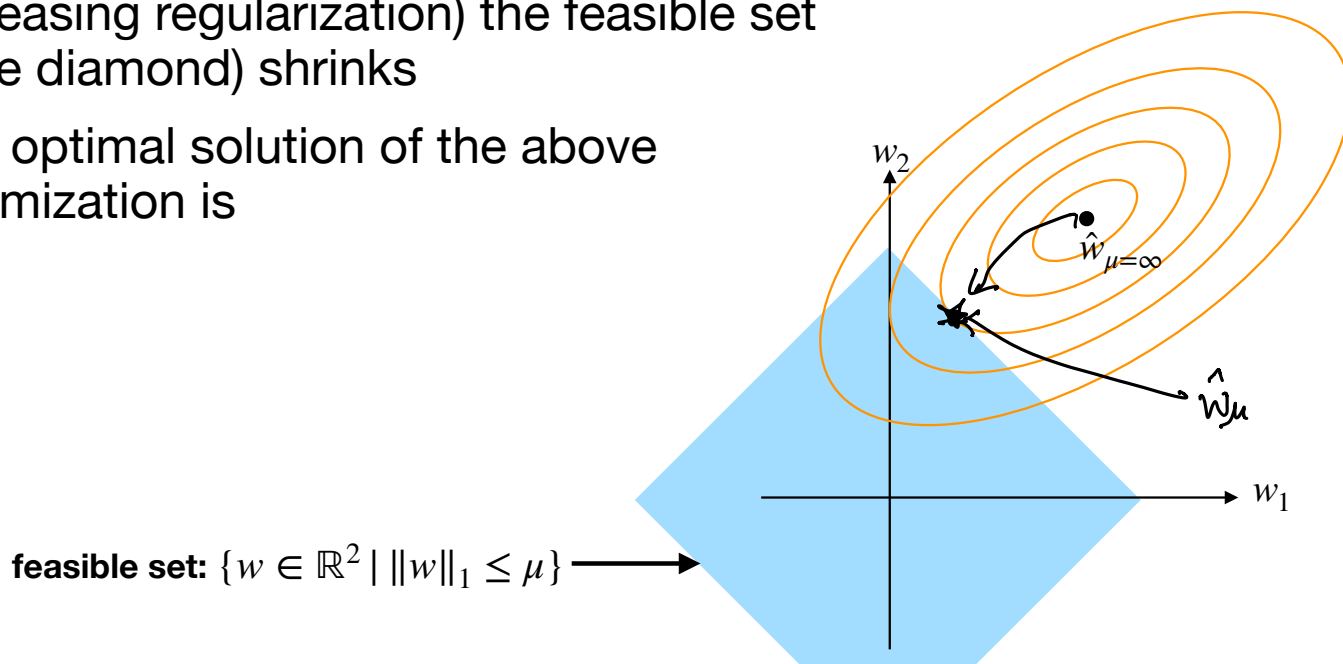
feasible set: $\{w \in \mathbb{R}^2 \mid \|w\|_1 \leq \mu\}$ $\longrightarrow$

# Why does Lasso give sparse solutions?

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

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

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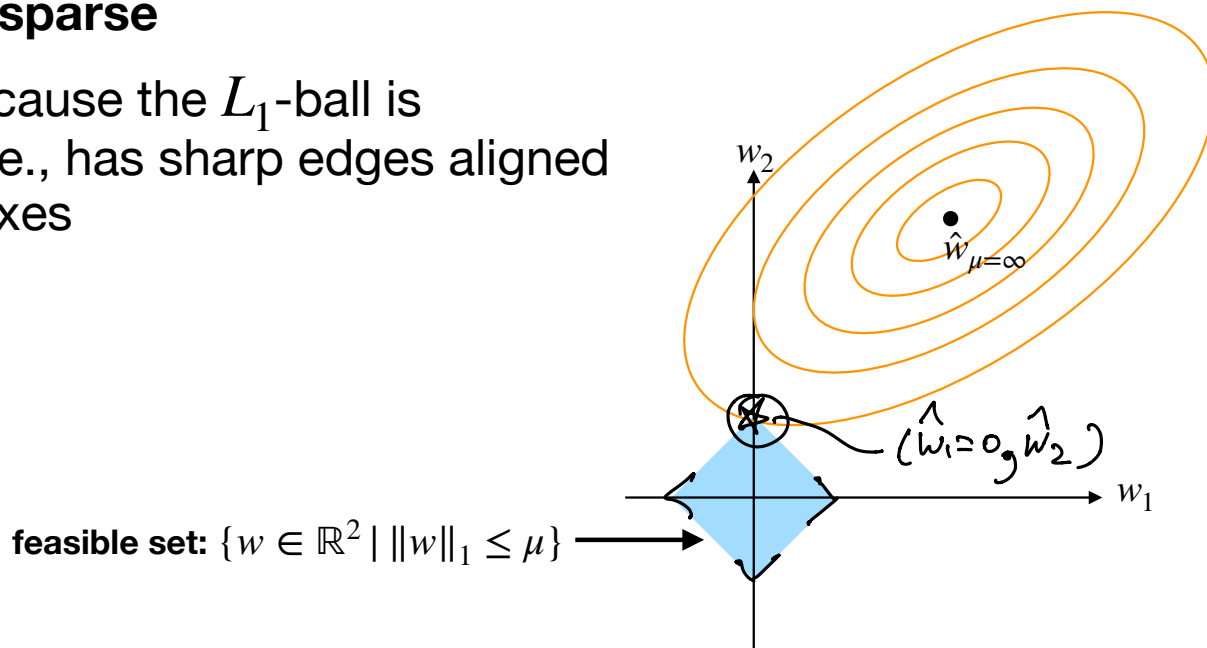
- The optimal solution of the above optimization is

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

$w_2$

$\hat{w}_{\mu=\infty}$

$\hat{w}_\mu$

$w_1$

# Why does Lasso give sparse solutions?

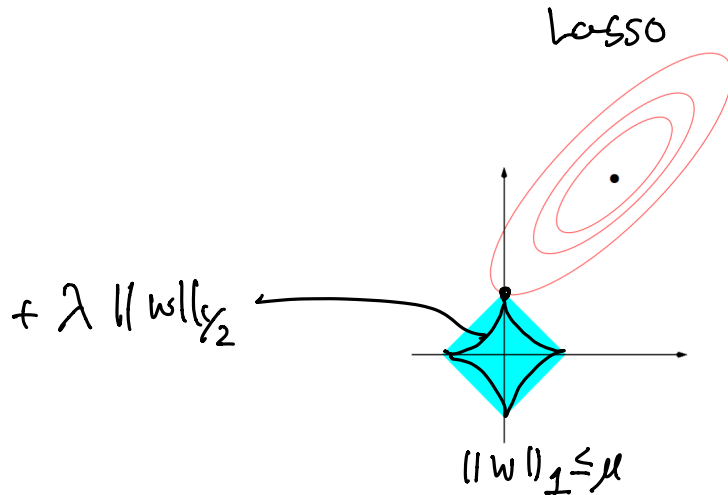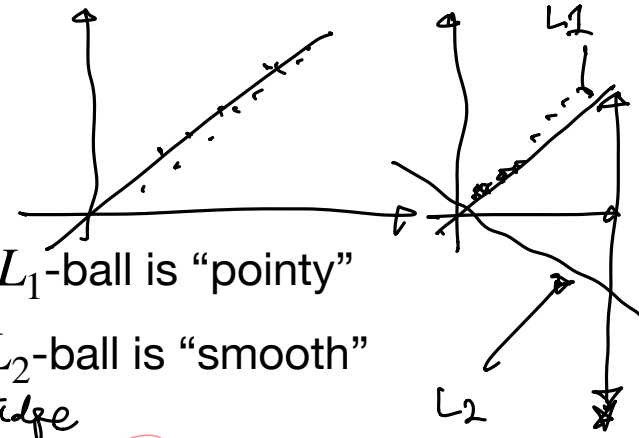$$\text{minimize}_w \quad \sum_{i=1}^{n} (w^T x_i - y_i)^2$$

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

- 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



$w_2$

$\hat{w}_{\mu=\infty}$

$(\hat{w}_1 = 0, \hat{w}_2)$

$w_1$

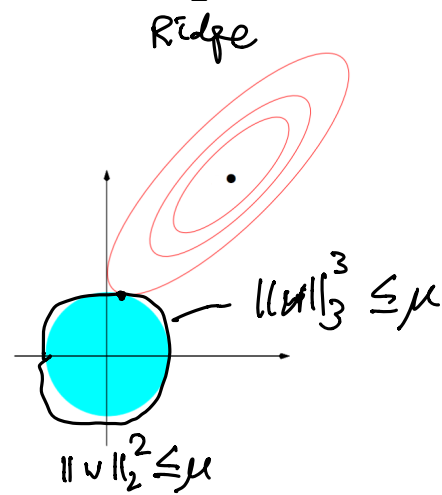**feasible set:** $\{w \in \mathbb{R}^2 \mid \|w\|_1 \leq \mu\}$ ⟶

# 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"

L1

L2

Lasso

Ridge

$+ \lambda \|w\|_{1/2}$

$\|w\|_1 \leq \mu$

$\|w\|_2^2 \leq \mu$

$\|w\|_3^3 \leq \mu$

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

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

subject to $\|w\|_1 \leq \mu$

subject to $\|w\|_2^2 \leq \mu$

# Questions?

vs.

$$\min \quad h_w(x, r)$$

$$s.t. \quad \| w \|_1 \leq 1$$

less Sparse.

$$\hat{w}_1$$

$$\min \quad h_w(x, r)$$

$$s.t. \quad \| w \|_{1/2} \leq 1$$

more sparse

$$\hat{w}_{1/2}$$

more zeros.

k-sparse ⇌ k non-zero

more sparse ⇌ k↓ ⇌ more zero.