Regularization helps avoid Overfitting! Readings: Murphy 11.3-4

→ Regularization in Linear Regression    code: ridge_end_lasso.ipynb

- Recall Least Squares:

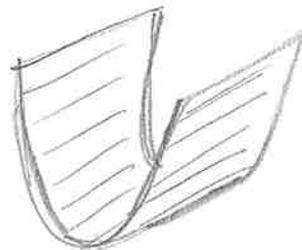$$\hat{w}_{LS} = \arg\min_w \sum_{i=1}^{n} (y_i - x_i^T w)^2$$

$$= \arg\min_w (y - Xw)^T (y - Xw)$$

$$= (X^T X)^{-1} X^T y \qquad (\text{when } (X^T X)^{-1} \text{ exists!})$$

- what if $x_i \in \mathbb{R}^d$ and $d > n$ ?

  - The objective function is flat in some directions
  - Implies optimal solution is not unique and unstable due to curvature:
    - Small changes in training data result in large changes in solution
    - often the magnitudes of $\boxed{w}$ are "very large"

= Regularization imposes "simpler" solutions by a "complexity" penalty

→ Sensitivity increases overfitting

- For a linear model:

$$y \approx b + w_1 x_1 + w_2 x_2 + \ldots + w_d x_d$$

if $|w_j|$ is large, then the prediction is [sensitive] to small changes in $x_j$

- Large sensitivity leads to overfitting and poor generalization, and equivalently, models that overfit tend to have large weights.
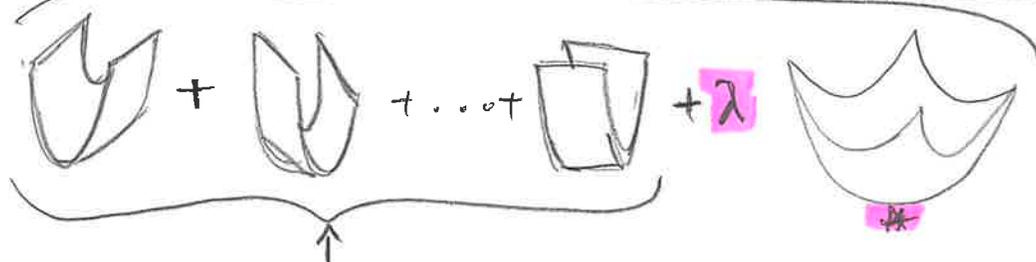
- Note: $b$ is a constant and hence there is no sensitivity for the offset $b$. Never regularize $b$

- In [Ridge Regression], we use regularizer $\|w\|_2^2$ to measure and control the [sensitivity] of the predictor.

- And, optimize for small loss and small sensitivity, by adding a regularizer in the objective (assume no offset for now)

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

regularization coefficient



Compared to old Least Squares objective:

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

$$X \in n \times d$$
$$w \in d \times 1$$

$$\hat{w}_{ridge} = \arg\min_{w} \sum_{i=1}^{n} (y_i - x_i^T w)^2 + \lambda \underbrace{\|w\|_2^2}_{②} \quad L_2 \text{ norm}$$

$$\|w\|_p \triangleq \left(|w_1|^p + \dots + |w_d|^p\right)^{\frac{1}{p}} = \|y - Xw\|^2 + \lambda \|w\|^2$$

$$= (Xw - y)^T (Xw - y) + \lambda w^T w$$

$$\nabla_w f = 2 X^T (Xw - y) + 2\lambda w = 0$$

$$(X^T X + \lambda I) w = X^T y$$

Similar to how we did in simple Linear Regression

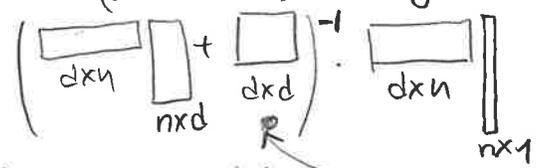$$I^d = \begin{bmatrix} 1 & 0 & \\ 0 & 1 & \ddots \\ & & 1 \\ 0 & & 1 \end{bmatrix} \in d \times d$$

$$\boxed{\hat{w}_{RIDGE} = (X^T X + \lambda I)^{-1} \cdot X^T \cdot y}$$

Used:

| Scalar Derivative | | Vector Derivative | |
|---|---|---|---|
| $f(x) \to$ | $\frac{df}{dx}$ | $f(x) \to$ | $\frac{df}{dx}$ |
| $bx \to$ | $b$ | $x^T B \to$ | $B$ |
| $bx \to$ | $b$ | $x^T b \to$ | $b$ |
| $x^2 \to$ | $2x$ | $x^T x \to$ | $2x$ |
| $bx^2 \to$ | $2bx$ | $x^T B x \to$ | $2Bx$ |

→ **Shrinkage Properties**

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

$$= (X^T X + \lambda I)^{-1} \cdot X^T \cdot y$$

$$\left( \underbrace{\Box}_{\substack{d \times n \\ n \times d}} + \underbrace{\Box}_{\substack{d \times d \\ R}} \right)^{-1} \cdot \underbrace{\Box}_{d \times n} \left| \right|_{n \times 1}$$
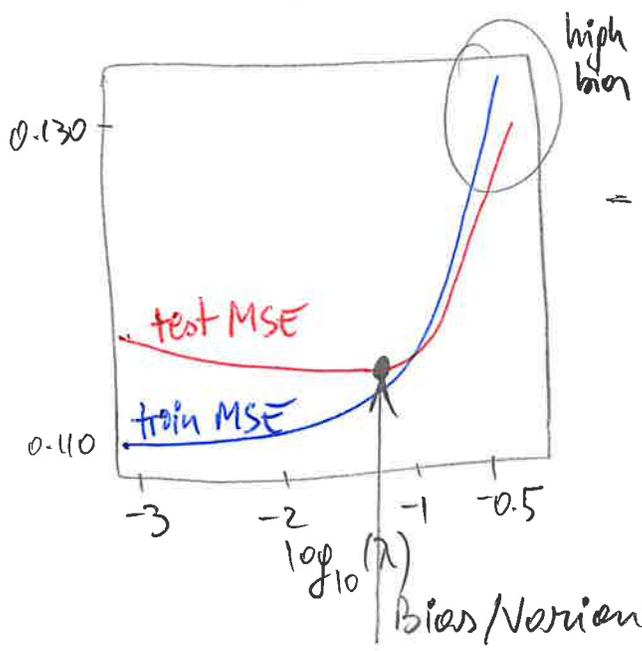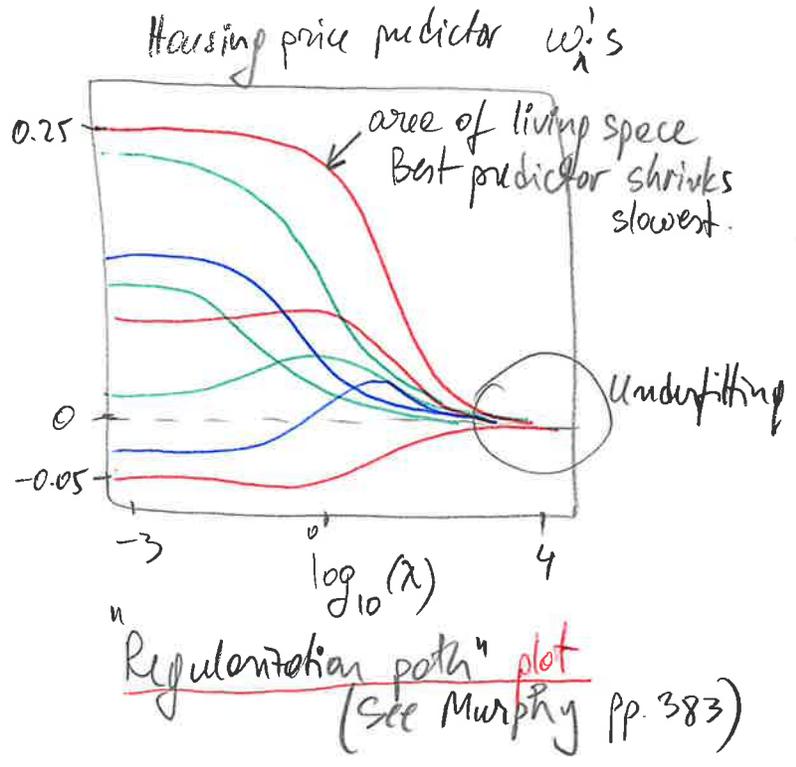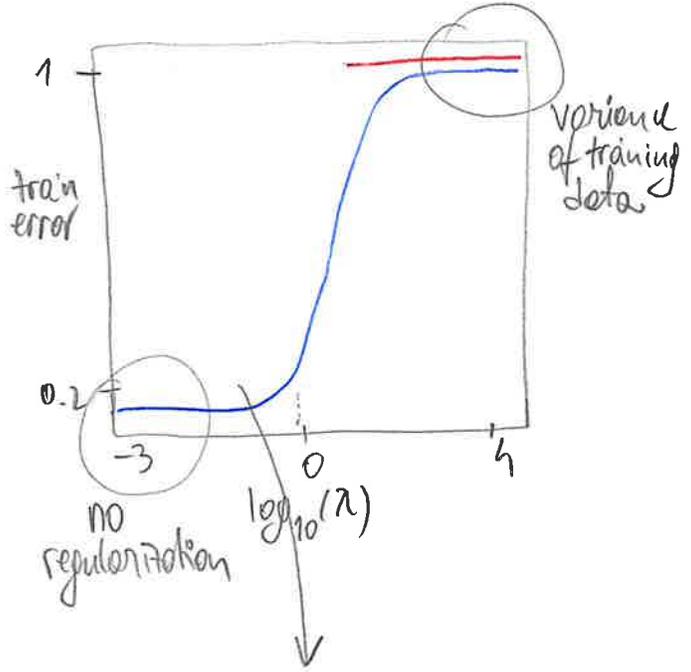
Larger $\lambda$ "divides by" larger term! (makes $w$ smaller)

- when $\lambda = 0 \Rightarrow$ least squares model
- this defines a family of models parameterized by $\boxed{\lambda}$
- Large $\lambda$ means more regularization and simpler model
- Small $\lambda$ means less regularization and more complex model.

Ridge Regression : minimize $\sum_{i=1}^{n} (w^T x_i - y_i)^2 + \lambda \|w\|_2^2$

Training MSE $\frac{1}{n} \sum_{i=1}^{n} (y_i - x_i^T \hat{w}_{ridge}^{(\lambda)})^2$



train error

1

0.2

−3          0          4

$\log_{10}(\lambda)$

no regularization

Variance of training data

Housing price predictor $w_{\lambda}^i$'s



0.25

0

−0.05

−3          0          4

$\log_{10}(\lambda)$

area of living space
Best predictor shrinks slowest.

Underfitting

"Regularization path" plot
(see Murphy pp. 383)



0.130

0.110

test MSE

train MSE

high bias

−3    −2    −1    −0.5

$\log_{10}(\lambda)$

Bias/Variance trade-off

We are not changing the model class, we are just changing $\boxed{\lambda}$.

= Gain in test MSE comes from shrinking $w$'s to get a less sensitive predictor;

− which in turn reduces the variance

− This is the role of regularizer: reduce sensitivity/variance !

→ **Bias-Variance Properties**

- $\hat{\omega}_{ridge} = (X^T X + \lambda I)^{-1} X^T y$

- To analyze bias-variance tradeoff, we need to assume probabilistic generative model:

$$x_i \sim P_X, \quad y = Xw + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I)$$

for some ground truth model parameter $\boxed{w}$

- The **true error** at a sample with feature $x$ is:

$$\mathbb{E}_{y, \mathcal{D}_{train} | x}\left[(y - x^T \hat{\omega}_{ridge})^2 \mid x\right] =$$

↗ test sample    ↖ fit to a given set $X, y$

→ $\eta(x)$ same as before

$$= \underbrace{\mathbb{E}_{y|x}\left[(y - \mathbb{E}[y|x])^2 \mid x\right]}_{\text{irreducible error}} + \underbrace{\mathbb{E}_{\mathcal{D}_{train}}\left[(\mathbb{E}[y|x] - x^T \hat{\omega}_{ridge})^2 \mid x\right]}_{\text{Learning error}}$$

$$= \mathbb{E}_{y|x}\left[(y - x^T w)^2 \mid x\right] + \mathbb{E}_{\mathcal{D}_{train}}\left[(\underset{\uparrow}{x^T w} - x^T \hat{\omega}_{ridge})^2 \mid x\right]$$

ground truth

$$= \underbrace{\sigma^2}_{\substack{\text{irreducible} \\ \text{error}}} + \underbrace{\left(x^T w - \mathbb{E}_{\mathcal{D}_{train}}[x^T \hat{\omega}_{ridge} | x]\right)^2}_{\text{Bias-squared}} + \underbrace{\mathbb{E}_{\mathcal{D}_{train}}\left[\left(\mathbb{E}_{\mathcal{D}_{train}}[x^T \hat{\omega}_{ridge} | x] - x^T \hat{\omega}_{ridge}\right)^2 | x\right]}_{\text{Variance}}$$

( Sample data from Independent Gaussians

Suppose $X^T X = nI$, then $\hat{\omega}_{ridge} = (X^T X + \lambda I)^{-1} X^T (Xw + \varepsilon)$

$$= \frac{n}{n+\lambda} \underset{\uparrow}{\textcircled{w}} + \frac{1}{n+\lambda} X^T \textcircled{\varepsilon}$$

ground truth     ↖ trades off weight vs. noise. )

Hw assignment to verify

$$\nabla^2 + \underbrace{\frac{\lambda^2}{(n+\lambda)^2}(\omega^T x)^2}_{\text{Bias-squared}} + \underbrace{\frac{\nabla^2_n}{(n+\lambda)^2}\|x\|_2^2}_{\text{Variance}}$$

$\underbrace{\qquad}_{\substack{\text{irreducible}\\\text{error}}}$

$-$ $\lambda$ trades off
Bias vs. Variance.

$-$ Larger $\lambda$ $\Rightarrow$ smaller variance.
Larger Bias.

## True Error:

$$E_{y,\mathcal{D}_{train}|x}\left[(y - x^T \hat{\omega}_{ridge})^2 | x\right] = \nabla^2 + \underbrace{\frac{\lambda^2}{(n+\lambda)^2}(\omega^T x)^2}_{\text{Bias-squared}} + \underbrace{\frac{\nabla^2_n}{(n+\lambda)^2}\|x\|_2^2}_{\text{Variance}}$$



$\lambda \to 0$

$\hat{\omega}_{ridge} \to \hat{\omega}_{LS}$

$\lambda \to \infty$

$\hat{\omega}_{ridge} \to 0$

## Takeaways:

- Regularization: penalizes complex models towards preferred, simpler models.

- Ridge Regression:
    - L2 penalized least-square regression.
    - Regularization parameter trades off model complexity w/ training error.
    - Never regularize the offset!

- We learned to measure $\boxed{\text{sensitivity}}$ by the size of weights $\|\omega\|_2^2$

$$\hat{\omega}_{LS} = \arg\min_{\omega} \sum_{i=1}^{n}(y_i - x_i^T \omega)^2$$
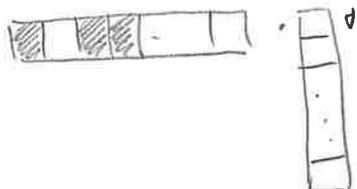
# Sparsity & the LASSO

— how to make the model compact and interpretable.

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

— vector $\boxed{\omega}$ is $\boxed{\text{sparse}}$ if many entries are $\emptyset$

= A vector $\omega$ is said to be $\boxed{k\text{-sparse}}$ if at most $k$ entries are non-zero.

— we are interested in $k$-sparse $\omega$, with $k << d$, because:
- computationally more efficient
- get rid of redundant/spurious features
- explanability/interpretability

$\boxed{\text{Efficiency}}$ — if $size(\omega) = 100$ billion, each prediction $\omega^T x$ is expensive.

— if $\omega$ sparse, prediction computation depends only on non-zeros in $\omega$.

$$\hat{y}_i = \hat{\omega}_{LS}^T \cdot x_i = \sum_{j=1}^{d} \hat{\omega}_{LS}[j] \times x_i[j] =$$

$$\Box = \boxed{\phantom{xx}} \cdot \boxed{\phantom{x}} = \sum_{j: \hat{\omega}_{LS}[j] \neq 0}^{d} \hat{\omega}_{LS}[j] \times x_i[j]$$

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

$\boxed{\text{interpretability}}$ — what are the relevant features to make a prediction?

— How do we find the "best" subset of features useful for predicting the price, among all combinations?

House price dataset

Lot size
Single family
Year Built
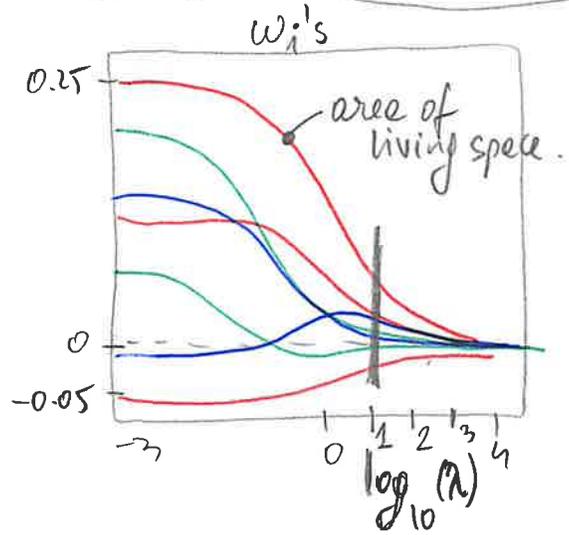Last sold price
Finished sqft
Finished basement
Parking type
Cooling

Heating
Exterior materials
Roof type
Structure type
...

==Finding Best subset:==

→ ==Exhaustive== - time consuming

→ ==Greedy== - Forward stepwise - start from simple model; add iteratively features __most__ useful to fit.

- Backward stepwise - start w/ full model; iteratively remove features __least__ useful to fit.
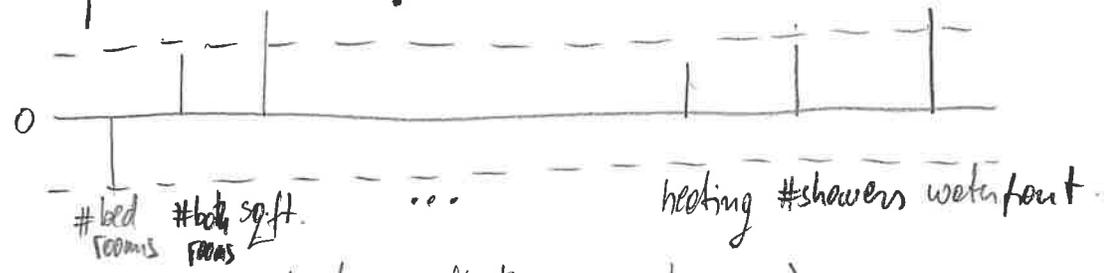
→ ==Regularize== - Recall that Ridge Regression makes coefficients small

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



Threshold Ridge Regression: - just set small ridge coefficient to $\emptyset$

- How to pick threshold?



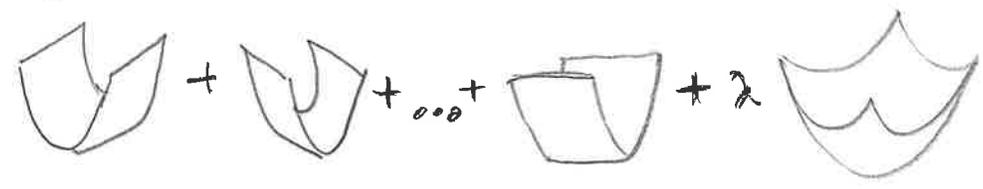#bed rooms  #bath rooms  sq.ft.  ...  heating #showers waterfront

- Challenge: related features (bathrooms, showers)
  - if did not include showers, weight on bathrooms increases!
  - we want a feature selection scheme to select one of the two __automatically__!

- Can we do this automatic selection w/ another regularizer?
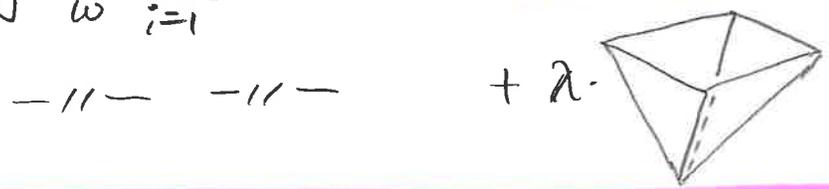
# Ridge vs. Lasso Regression

- **Ridge** Regression objective:

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



- **Lasso** objective:

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

$$-//- \quad -//- \quad + \lambda \cdot$$



**LASSO** = Least Absolute Shrinkage and Selection Operator

- Sensitivity of a model $\omega$ is measured in L1-norm:

$$\|\omega\|_1 = \sum_{i=1}^{d} \underbrace{|\omega[j]|}_{\text{absolute value}}$$

$$\left( \begin{array}{l} \ell_p\text{-norm of a vector } \omega \in \mathbb{R}^d \text{ is:} \\[2mm] \|\omega\|_p \overset{\Delta}{=} \left( \sum_{j=1}^{d} |\omega[j]|^p \right)^{1/p} \end{array} \right)$$

Example: house price w/ 16 features



Most relevant features last largest in reg. path

Ridge regression

Lasso regression

==Lasso regression== – naturally gives ==sparse== features.

1 > model selection – choose $\lambda$ based on (cross) validation error
2 > feature selection – keep only features with non-zero (or not-too-small) parameters in $w$ at optimal $\lambda$

3 > retrain with the sparse model and $\lambda = 0$.

- ==Regularized Least Squares==
  Optimization:
  $$\hat{w}_R = \arg\min_w \sum_{i=1}^{n} (y_i - x_i^T w)^2 + \lambda \cdot r(w)$$

  Ridge : $r(w) = \|w\|_2^2$
  Lasso : $r(w) = \|w\|_1$

Theorem:

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 <u>constrained</u> optimization, $\hat{w}_c$, is the same as the solution of the <u>regularized</u> 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 points $(\lambda, \mu)$ whose optimal solution $\hat{w}_R$ are the same for the ==regularized optimization== and <u>constrained optimization</u>!

Why does Lasso give sparse solutions?

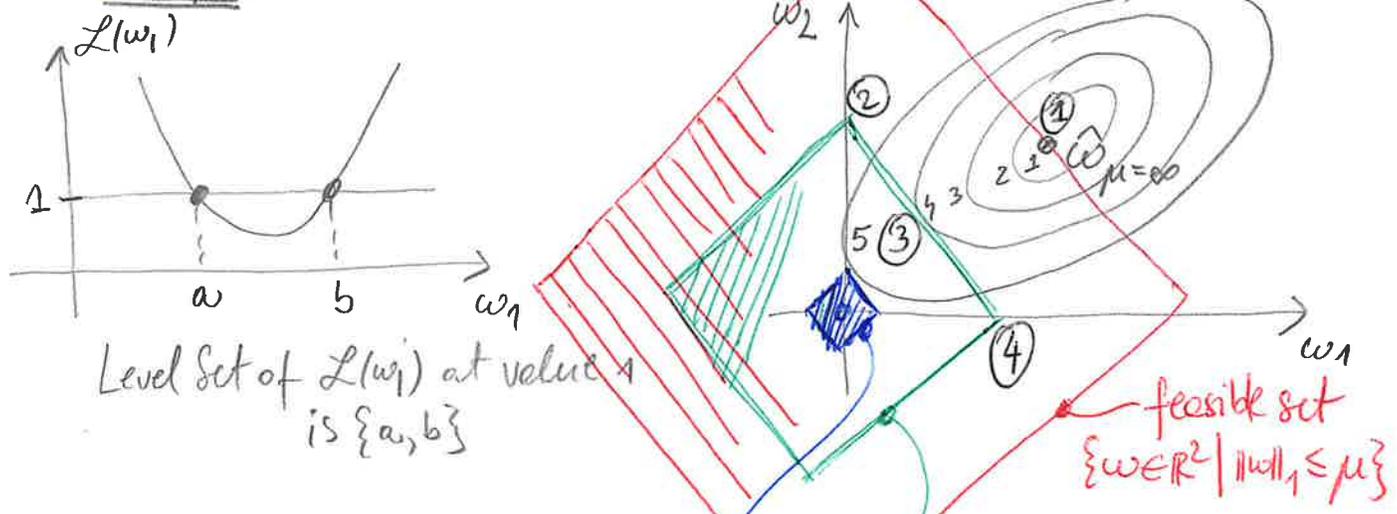$$\underset{w}{\text{minimize}} \sum_{i=1}^{n} (w^T x_i - y_i)^2$$

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

- "Level set" of a function $\mathcal{L}(w_1, w_2)$ is defined as set of points $(w_1, w_2)$ that have the same function value.
- Level set of a quadratic function is an oval.
- Center of oval is the Least Squares solution $\boxed{\hat{w}_{\mu=\infty} = \hat{w}_{LS}}$
- Examples:



Level Set of $\mathcal{L}(w_1)$ at value $1$ is $\{a, b\}$

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

- As we decrease $\mu$ from $\infty$, the feasible set becomes smaller.
- The shape of the feasible set is known as L1 ball, which is a diamond
- In 2-dimensions, diamond is: $\{(w_1, w_2) \mid |w_1| + |w_2| \leq \mu\}$

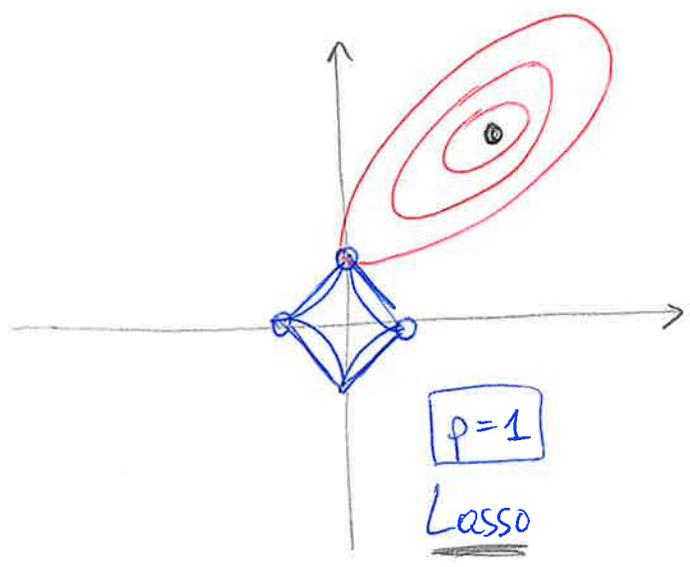◆ ⟶ when $\mu$ is large enough such that $\|\hat{w}_{\mu=\infty}\|_1 < \mu$, then, the optimal solution does not change as the feasible set includes the un-regularized optimal solution

◆ — As $\mu$ decreases (which is equivalent to increasing regularization $\lambda$) the feasible set (green diamonds in figure above) shrinks.
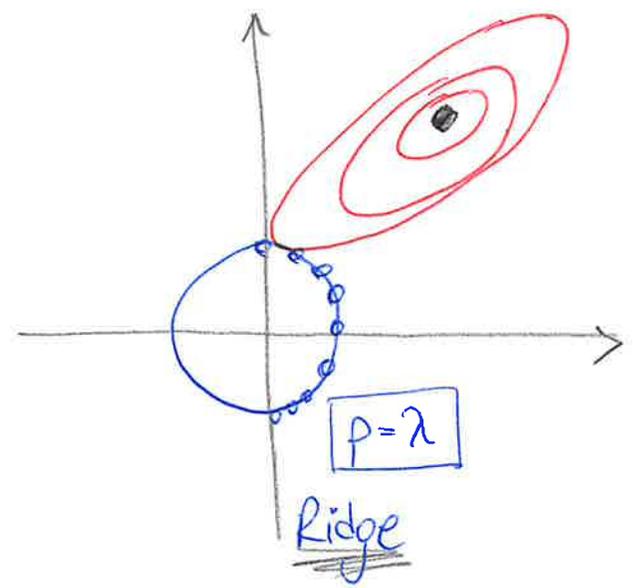
◆ — For small enough $\mu$, the optimal solution becomes [sparse]; because the L1-ball is "pointy" (i.e., has sharp edges aligned with the axes)

## Constrained Least Squares

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



$p=1$

Lasso

$$\begin{cases} \underset{\omega}{\text{minimize}} \ \sum_{i=1}^{n} (\omega^T x_i - y_i)^2 \\ \text{subject to:} \ \|\omega\|_1 \le \mu \end{cases}$$



$p=2$

Ridge

$$\begin{cases} \underset{\omega}{\text{minimize}} \ \sum_{i=1}^{n} (\omega^T x_i - y_i)^2 \\ \text{subject to:} \ \|\omega\|_2^2 \le \mu. \end{cases}$$

## L1-Ball in Higher Dimensions

L1-ball in 3 dimensions



Corners: 1-sparse

Edges: 2-sparse.