# 446 Section 4 ← (3 - η(-1))

TA: Sankar Harilal

Plans for today!

1. Reminders
2. Train/Val/Test Review
3. Gradient Descent
4. Generalized Least Squares
5. Importance of Regularization in Least Squares
6. Ridge/LASSO (if time)

# Reminders

- HW1 was due yesterday
  - Remember that you have 5 late days!
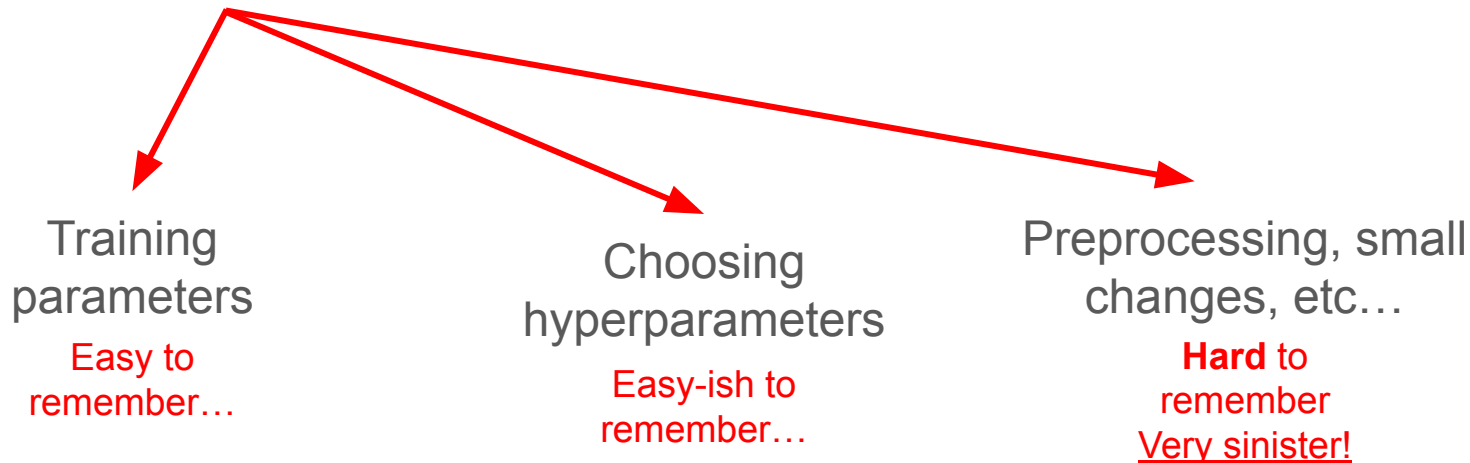- HW2 was released yesterday; due Wednesday, Feb 11

Some tips:

- Use office hours to your advantage
  - Student TA OH for homework questions
  - Professor OH for conceptual questions
- Skim the homeworks the day they are assigned and try one problem
  - Motivates you to get things done on time, starting an untouched assignment can be daunting
- Keep a tab open with the lecture slides while you do the homework for reference

# Train/val/test

# What do you never ever ever ever ever ever ever ever ever do?

## Train/tune your model on your test set!

Training parameters

Easy to remember…

Choosing hyperparameters

Easy-ish to remember…

Preprocessing, small changes, etc…

**Hard** to remember
Very sinister!

# Easiest way to combat this?

| Data |
|---|

| Data | Test |
|---|---|

| Data |
|---|

Only touch this until **everything is done**

Then you can free him

Test

# Bonus Questions:

**Q:** If you have $N$ data points, what choice of $k$ for k-fold validation would give you the same train/val splits as LOOCV?

**A:** $N$

**Q:** How would you get the **most** "pessimistic" estimate on your validation set?

**A:** Only train on 1 example and validate on the rest

# Problems 1.1, 1.2

You are given blocks of code, and something is wrong/not totally right with how they deal with the data.

*Identify them and propose solutions!*

What do you never ever ever ever ever ever ever ever do?

Train/tune your model on your test set!

Training parameters

**Easy to remember…**

Choosing hyperparameters

**Easy-ish to remember…**

Preprocessing, small changes, etc…

**Hard** to remember
Very sinister!

## 1.1. Program 1

```
1   # Given dataset of 1000-by-50 feature
2   # matrix X, and 1000-by-1 labels vector
3
4   mu = np.mean(X, axis=0)
5   X = X - mu
6
7   idx = np.random.permutation(1000)
8   TRAIN = idx[0:900]
9   TEST = idx[900::]
10
11  ytrain = y[TRAIN]
12  Xtrain = X[TRAIN, :]
13
14  # solve for argmin_w ||Xtrain*w - ytrain||_2
15  w = np.linalg.solve(np.dot(Xtrain.T, Xtrain), np.dot(Xtrain.T, ytrain))
16
17  b = np.mean(ytrain)
18
19  ytest = y[TEST]
20  Xtest = X[TEST, :]
21
22  train_error = np.dot(np.dot(Xtrain, w)+b - ytrain,
23                np.dot(Xtrain, w)+b - ytrain ) / len(TRAIN)
24  test_error = np.dot(np.dot(Xtest, w)+b - ytest,
25                np.dot(Xtest, w)+b - ytest ) / len(TEST)
26
27  print('Train error = ', train_error)
28  print('Test error = ', test_error)
```

*mu* is calculated from the **entire** data (train + test), intertwining them!

*This is bad!*

**Correct procedure:**
- **Split into train and test**
- **Compute the mean of the train data ($\mu_{train}$)**
- **De-mean both train and test data using $\mu_{train}$**

## 1.2. Program 2

```python
1   # We are given: 1) dataset X with n=1000 samples and 50 features and 2) a vector y of length 1000 with labels.
2   # Consider the following code to train a model, using cross validation to perform hyperparameter tuning.
3
4   def fit(Xin, Yin, _lambda):
5       w = np.linalg.solve(np.dot(Xin.T, Xin) + _lambda * np.eye(Xin.shape[1]), np.dot(Xin.T, Yin))
6       b = np.mean(Yin) - np.dot(w, mu)
7       return w, b
8
9   def predict(w, b, Xin):
10      return np.dot(Xin, w) + b
11
12  idx = np.random.permutation(1000)
13  TRAIN = idx[0:800]
14  VAL = idx[800:900]
15  TEST = idx[900::]
16
17  ytrain = y[TRAIN]
18  Xtrain = X[TRAIN, :]
19  yval = y[VAL]
20  Xval = X[VAL, :]
21
22  # demean data
23  mu = np.mean(Xtrain, axis=0)
24  Xtrain = Xtrain - mu
25  Xval = Xval - mu
26
27  # use validation set to pick the best hyper-parameter to use
28  lambdas = [10 ** -5, 10 ** -4, 10 ** -3, 10 ** -2]
29  err = np.zeros(len(lambdas))
30
31  for idx, _lambda in enumerate(lambdas):
32      w, b = fit(Xtrain, ytrain, _lambda)
33      yval_hat = predict(w, b, Xval)
34      err[idx] = np.mean((yval_hat - yval)**2)
35
36  lambda_best = lambdas[np.argmin(err)]
37
38  Xtot = np.concatenate((Xtrain, Xval), axis=0)
39  ytot = np.concatenate((ytrain, yval), axis=0)
40
41  w, b = fit(Xtot, ytot, lambda_best)
42
43  ytest = y[TEST]
44  Xtest = X[TEST, :]
45
46  # demean data
47  Xtest = Xtest - mu
48
49  ytot_hat = predict(w, b, Xtot)
50  train_error = np.mean((ytot_hat - ytot) **2)
51  ytest_hat = predict(w, b, Xtest)
52  test_error = np.mean((ytest_hat - ytest) **2)
53
54  print('Train error = ', train_error)
55  print('Test error = ', test_error)
```

The final model is trained on BOTH the training and validation sets.

*This is… eh…*

**Your hyperparameters selected on just the train data may not hold for train + val:**
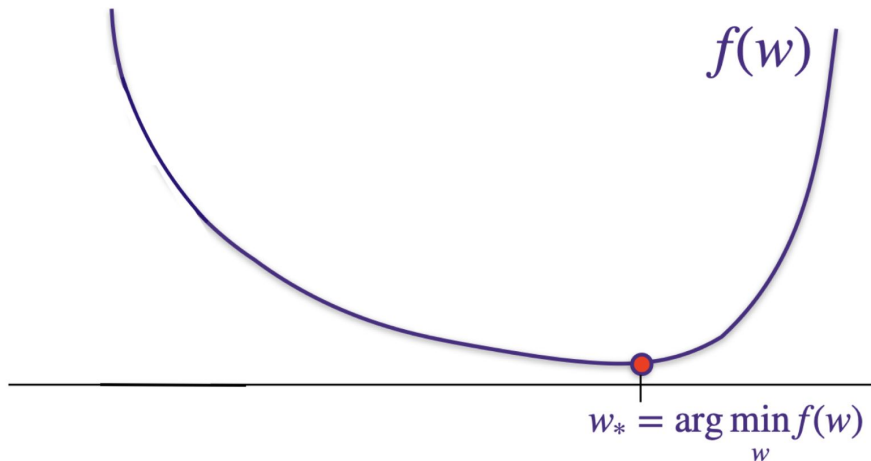-   **More data is good but you should ensure that the hyperparameters you tuned do not depend on the number of elements.**
-   **Tradeoff between more data and better test error estimate**

# Gradient Descent

# Gradient Descent

Purpose of this exercise: Understanding how gradient descent relates to approximations, and why it works.

Consider some function $f(w)$, which has some $w_*$ for which $w_* = \arg\min_w f(w)$:



$f(w)$

$w_* = \arg\min_w f(w)$
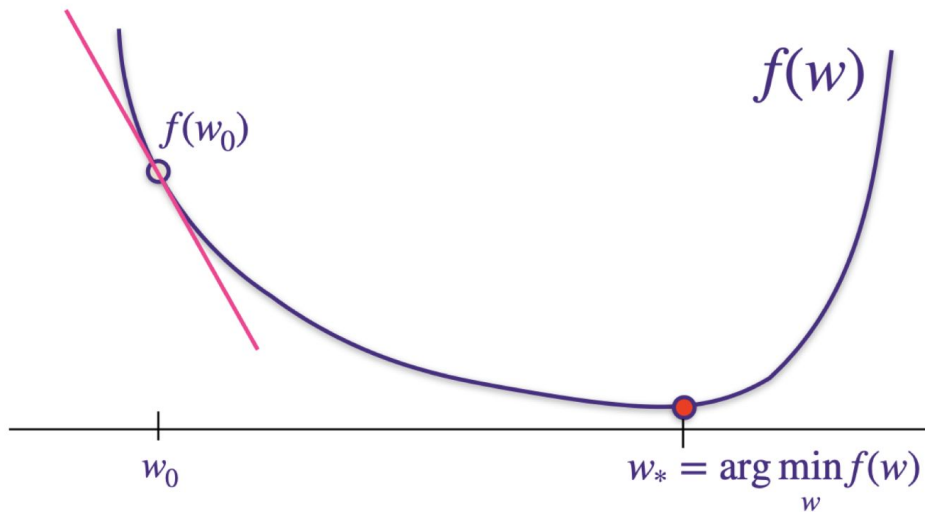
# Question 2a

Let $w_0$ be some initial guess for the minimum of $f(w)$. Gradient descent will allow us to improve this solution.

(a) For some $w$ that is very close to $w_0$, give the Taylor series approximation for $f(w)$ starting at $f(w_0)$.
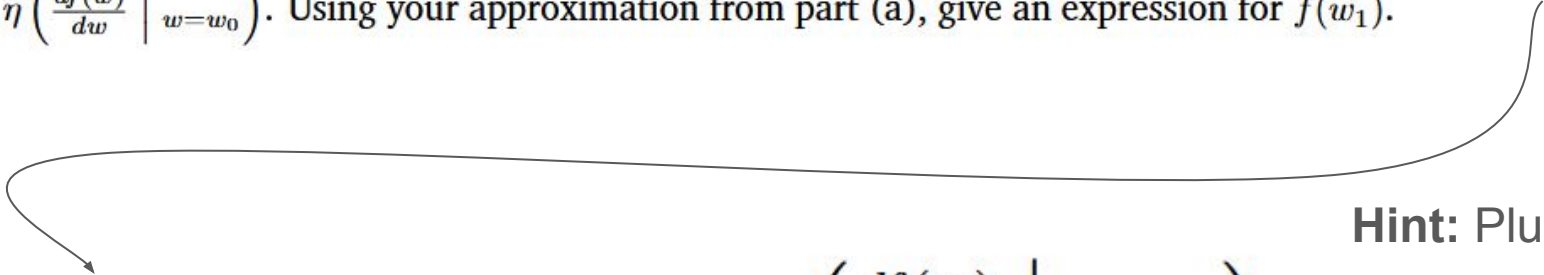
Remember Taylor expansion?

$\hookrightarrow$ To approximate a function around a point $\underline{a}$

$$f(x) \approx f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f'''(a)}{3!}(x-a)^3 \cdots$$

$\uparrow$
Exact at $\underline{a}$, close around $\underline{a}$

Better and better approximations

# Question 2a

Let $w_0$ be some initial guess for the minimum of $f(w)$. Gradient descent will allow us to improve this solution.

(a) For some $w$ that is very close to $w_0$, give the Taylor series approximation for $f(w)$ starting at $f(w_0)$.

For $w$ very close to $w_0$, we see that $f(w) \approx f(w_0) + (w - w_0) \left( \frac{df(w)}{dw} \Big|_{w=w_0} \right)$.

# Question 2b

(b) Now, let us choose some $\eta > 0$ that is *very small*. With this very small $\eta$, let's assume that $w_1 = w_0 - \eta \left( \frac{df(w)}{dw} \Big|_{w=w_0} \right)$. Using your approximation from part (a), give an expression for $f(w_1)$.

**Hint:** Plug in here

$$f(w) \approx f(w_0) + (w - w_0) \left( \frac{df(w)}{dw} \Big|_{w=w_0} \right).$$

Fancy way of saying f'(w₀)

(Derivative of f(w) at w₀)

# Question 2b

$$w_1 = w_0 - \eta \left( \frac{df(w)}{dw} \bigg|_{w=w_0} \right) \quad \leftarrow \text{Given}$$

$$f(w_1) \approx f(w_0) + (w_1 - w_0) \left( \frac{df(w)}{dw} \bigg|_{w=w_0} \right)$$

plug in here

$$= f(w_0) + \left( w_0 - \eta \left( \frac{df(w)}{dw} \bigg|_{w=w_0} \right) - w_0 \right) \left( \frac{df(w)}{dw} \bigg|_{w=w_0} \right)$$

$$= f(w_0) - \eta \left( \frac{df(w)}{dw} \bigg|_{w=w_0} \right)^2$$

# Question 2c

(c) Given your expression for $f(w_1)$ from part (b), explain why, if $\eta$ is small enough and if the function approximation is a good enough approximation, we are guaranteed to move in the "right" direction closer to the minimum $w_*$.

**Remember:**
We want to
minimize this

**Hint:** Why
would this be
good?

$$f(w_1) \approx f(w_0) - \eta \left( \frac{df(w)}{dw} \bigg|_{w=w_0} \right)^2$$

# Question 2c

Note that in part (b), the derivative is squared and will always be a nonnegative value. Therefore, $f(w_1) < f(w_0)$.

$$f(w_1) \approx f(w_0) - \eta \left( \frac{df(w)}{dw} \bigg|_{w=w_0} \right)^2$$

**In English:** The loss function after a weight update will always evaluate to be smaller than before the weight update
- If the step size is small enough
- If the approximation is good enough

# Question 2d

(d) Building from your answer in part (c), write a general form for the gradient descent algorithm.

Hint: how could we generalize this equation from part b?

$$w_1 = w_0 - \eta \left( \frac{df(w)}{dw} \bigg|_{w=w_0} \right)$$
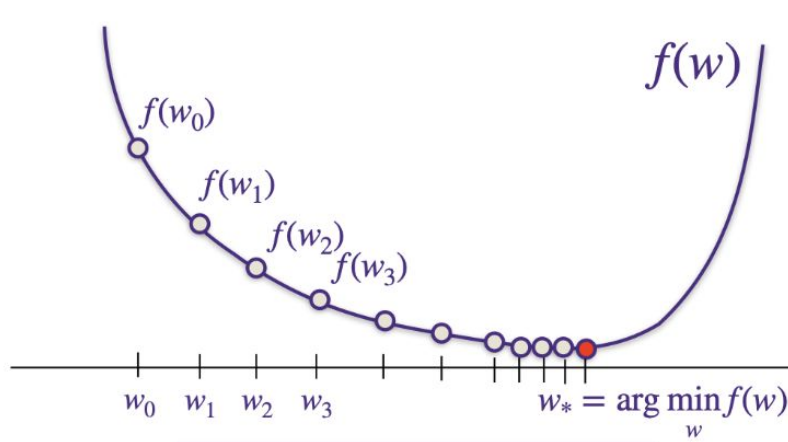
# Question 2d

Gradient descent is written as:

For $k = 0, 1, 2, 3, ...,$ $w_{k+1} = w_k - \eta \left( \frac{df(w)}{dw} \Big|_{w=w_k} \right).$

Note that as $k \to \infty$, $\left( \frac{df(w)}{dw} \Big|_{w=w_k} \right) \to 0.$

We visualize as:

**Convergence guarantees iff convex!**

# Generalized Least Squares

# Least Squares Proof(s)

Should look familiar…

$$\widehat{\omega}_{\text{general}} = (X^\top X + \lambda D)^{-1} X^\top y$$

Has shown up…

- In lecture (Lecture 2)
- On your homework (A5 Ridge Regression proof)
- And now here!

You can look at the generalized proof in your own time.

$$\widehat{\omega}_{\text{general}} = \left( \sum_{i=1}^{n} x_i x_i^\top + \lambda D \right)^{-1} \left( \sum_{i=1}^{n} x_i y_i \right)$$

# Question 3.2a

$$\widehat{\omega}_{\text{general}} = (X^\top X + \lambda D)^{-1} X^\top y$$

(a) In the simple least squares case ($\lambda = 0$ above), what happens to the resulting $\widehat{\omega}$ if we double all the values of $y_i$?

**Solution:**

As can be seen from the formula $\widehat{\omega} = (X^\top X)^{-1} X^\top y$, doubling $y$ doubles $\omega$ as well. This makes sense intuitively as well because if the observations are scaled up, the model should also be.

# Question 3.2b

$$\widehat{\omega}_{\text{general}} = (X^\top X + \lambda D)^{-1} X^\top y$$

(b) In the simple least squares case ($\lambda = 0$ above), what happens to the resulting $\widehat{\omega}$ if we double the data matrix $X \in \mathbb{R}^{n \times d}$?

**Solution:**

As can be seen from the formula $\widehat{\omega} = (X^\top X)^{-1} X^\top y$, doubling $X$ halves $\omega$. This also makes sense intuitively because the error we are trying to minimize is $\|X\omega - y\|_2^2$, and if the $X$ has doubled, while $y$ has remained unchanged, then $\omega$ must compensate for it by reducing by a factor of 2.

# Importance of Regularization in Least Squares

# Question 3.2c

$$\widehat{\omega}_{\text{general}} = (X^\top X + \lambda D)^{-1} X^\top y$$

(c) Suppose $D = I$ (that is, it is the identity matrix). That is, this is the *ridge* regression setting. Explain why $\lambda > 0$ ensures that the solution exists and the matrix can be inverted.

# 3.2c setup

Let's do a linear algebra refresher so that we can show off an interesting and actually useful result about the utility of regularization!

$A : \mathbb{R}^d \to \mathbb{R}^n, \ \text{if} \ d \gg n$

$$A : \begin{bmatrix} \vdots & \vdots & \vdots \\ \ & \ & \ \end{bmatrix} n$$

must have a non-empty nullspace (can show using rank-nullity theorem)

This is <u>bad</u> for invertibility

<u>Note</u>: $Null(A) = Null(A^T A)$

$\hookrightarrow$ proof in Section 02 handout

Invertible means $(A^T A)^{-1}(A^T A) = I$, so $(A^T A)^{-1}(A^T A)x = x$

---

Way to think about nullspaces

$A \in \mathbb{R}^{n \times d} \quad x \in \mathbb{R}^d$

Nullspace: Subspace of $\mathbb{R}^d$, contains all solutions to $\underline{Ax = 0}$

In other words, all the vectors are "annihilated" by $A$

Invertible means $(A^\top A)^{-1}(A^\top A) = I$, so $(A^\top A)^{-1}(A^\top A)x = x$

If $A^\top A$ has a non-empty nullspace, then

$$\exists x \text{ s.t. } (A^\top A)x = 0$$

$\downarrow$

Makes $(A^\top A)^{-1}\underbrace{(A^\top A)x}= x$ impossible!
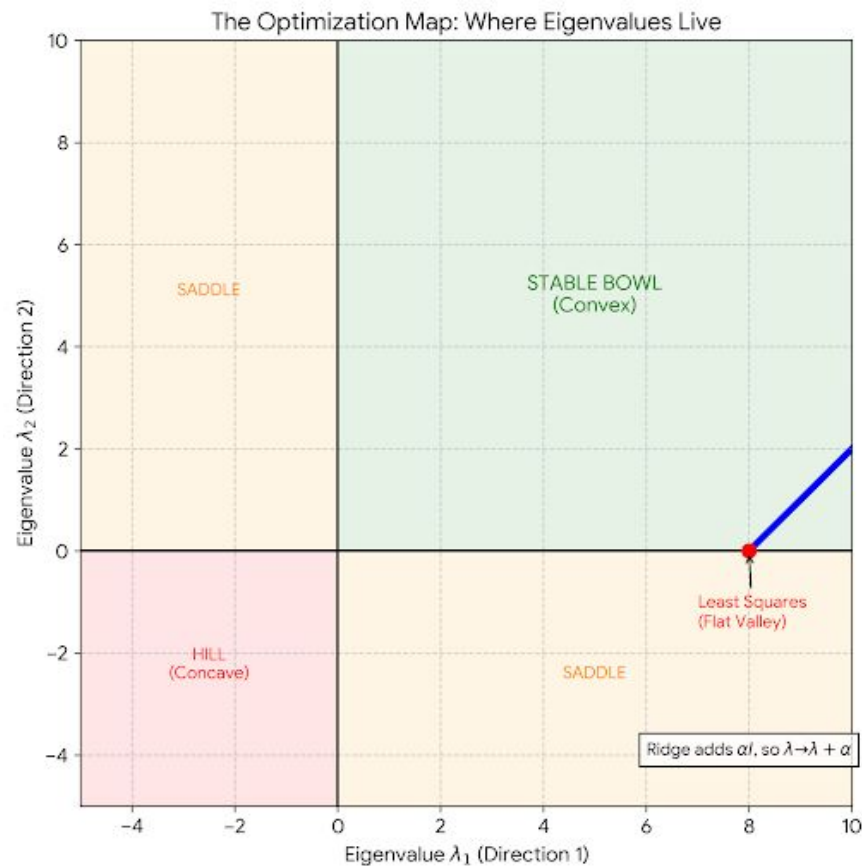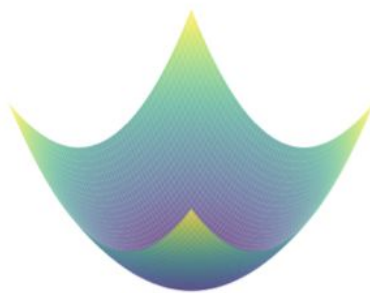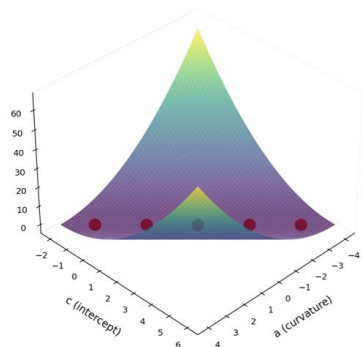
If this $= 0$, no way to recover $x$!

Main idea: If $X \in \mathbb{R}^{n \times d}, d \gg n,$ then

$\text{Null}(X)$ and $\text{Null}(X^\top X)$ are non-empty

This means $X^\top X$ has no inverse

An issue because... $\hat{w} = (X^\top X)^{-1} X^\top y$

Let's not give up!

# Visualized

- If $X$ has a non-empty null space, matrix $X^T X$ has an eigenvalue of 0. This corresponds to a "flat valley" – no unique solution
- By adding $I$, we shift all the eigenvalues, moving it into a "stable bowl"

Let's add in $\lambda I : \hat{w} = (X^T X + \lambda I)^{-1} X^T Y$

Is $X^T X + \lambda I$ always invertible for $\lambda > 0$?

A matrix $A$ is positive semi-definite if $x^\top A x \geq 0$ and positive definite if $x^\top A x > 0$
**Positive Definite (PD):** All eigenvalues are strictly positive
**Positive Semi-Definite (PSD):** All eigenvalues are 0 or positive

We want to show that $u(X^T X + \lambda I)u > 0 \quad \forall u \in \mathbb{R}^d$

$\hookrightarrow$ Show matrix is positive definite, meaning it must have an inverse

We want to show that $u(X^T X + \lambda I)u > 0 \quad \forall u \in \mathbb{R}^d$

$$u(X^T X + \lambda I)u = \underbrace{u(X^T X)u}_{\text{Is this always } > 0?} + u(\lambda I)u$$

L2 norm

$$u(X^T X)u = u X^T X u = \|Xu\|_2^2 \geq 0 \quad \text{Yes!}$$

$$u(X^T X + \lambda I)u = u(X^T X)u + u(\lambda I)u$$

$$\geq u(\lambda I)u \leftarrow \text{this is PD, } \therefore \; > 0$$

$$> 0$$

We have shown $X^T X + \lambda I$ is PD and therefore always invertible if $\lambda > 0$!

$\hookrightarrow$ $\underline{\text{Even if } d \gg n!}$

# Question 3.2c

$$\widehat{\omega}_{\text{general}} = (X^\top X + \lambda D)^{-1} X^\top y$$

(c) Suppose $D = I$ (that is, it is the identity matrix). That is, this is the *ridge* regression setting. Explain why $\lambda > 0$ ensures that the solution exists and the matrix can be inverted.
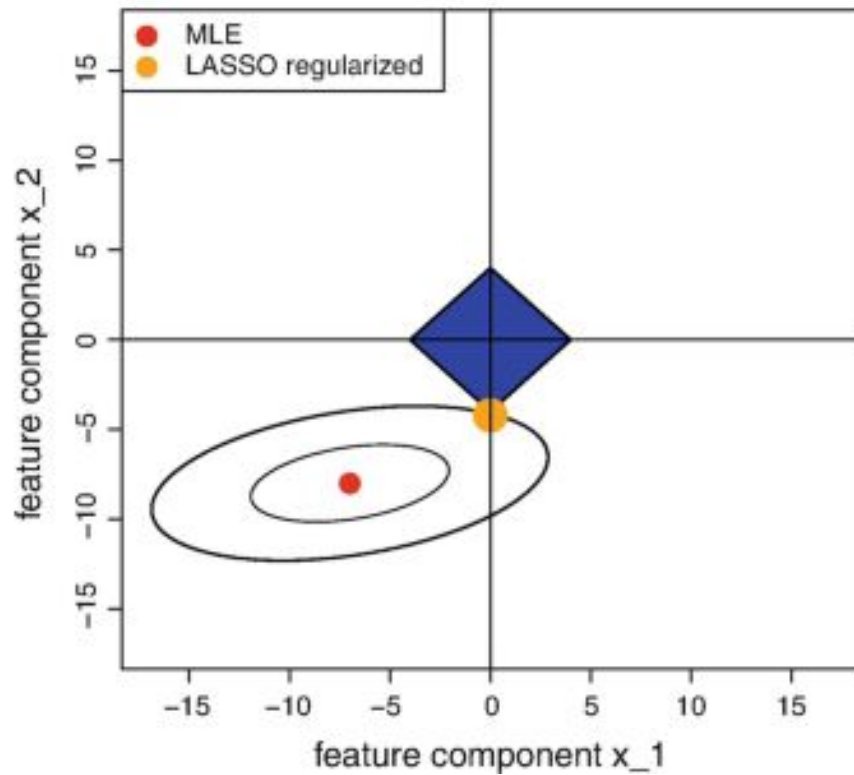
**Solution:**

The solution is $\widehat{\omega} = (X^\top X + \lambda I)^{-1} X^\top y$. We already saw in a previous part that $X^\top X$ is always positive semidefinite, that is, its eigenvalues are at least zero. Adding $\lambda I$, where $\lambda > 0$, ensures that $X^\top X + \lambda I$ is in fact positive *definite*. This helps us have a good condition number.
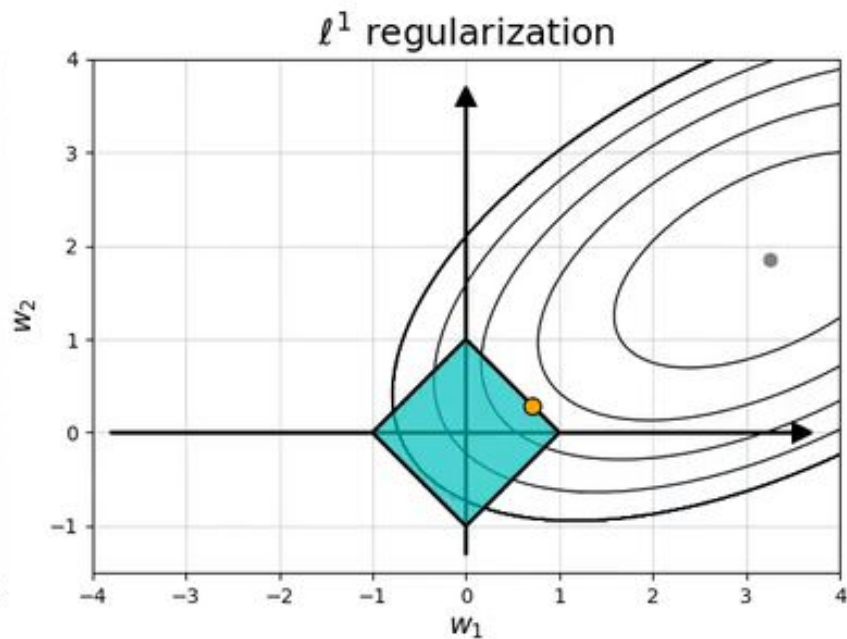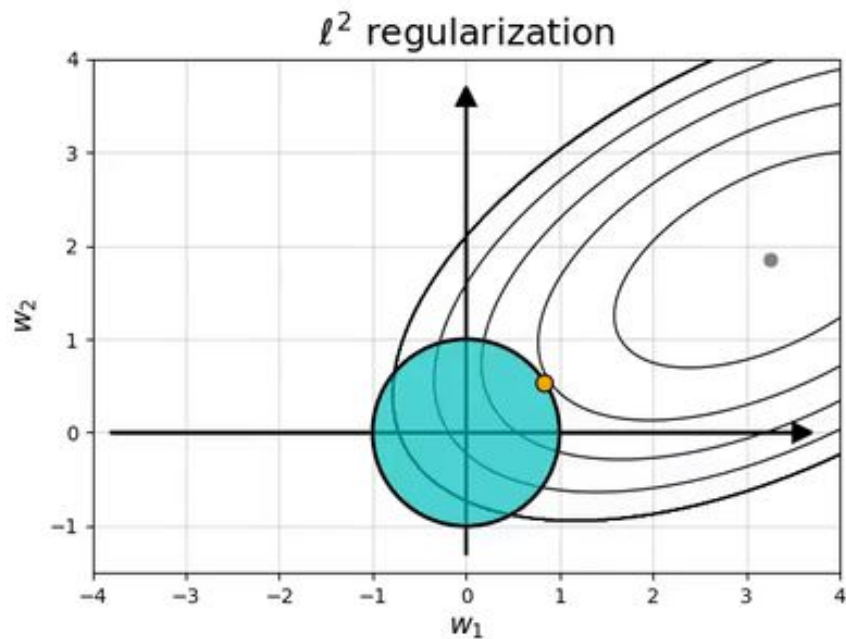
# Ridge vs. LASSO

ridge regularization (L2)

LASSO regularization (L1)

$\ell^1$ induces sparse solutions for least squares

$\ell^2$ regularization

$\ell^1$ regularization

by @itayevron

# Questions/Chat Time!