

446 Section 4 $\leftarrow (3 - \eta(-1))$

TA: Varun Ananth

Plans for today!

1. This
2. Reminders
3. Train/Val/Test Problems
4. Gradient Descent
5. Generalized Least Squares
6. Ridge/LASSO (if time)

Reminders

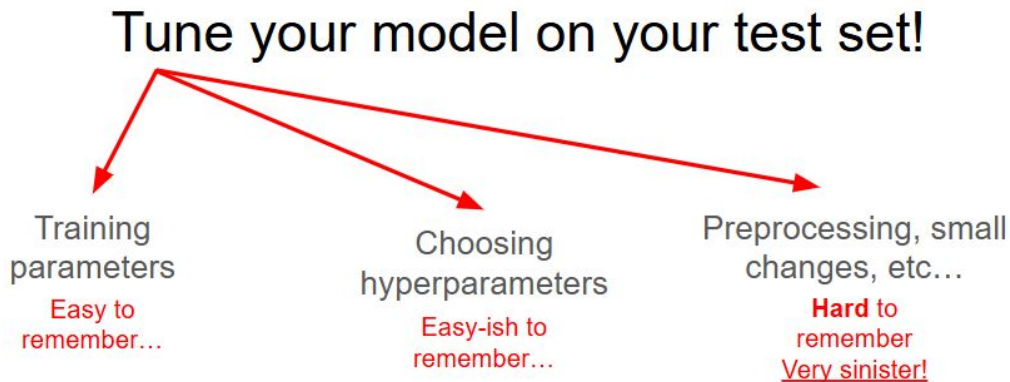
- HW1 was due yesterday
 - Remember the late day policy!
- HW2 is released
- Midterm in a week...
 - February 7th, Friday

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 do?



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!

Calculate a mean just on the train data, and use this to de-mean both the train and test datasets

1.2. Program 2

```
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, lambda_best)
50 train_error = np.mean((ytot_hat - ytot)**2)
51 ytest_hat = predict(w, b, Xtest, lambda_best)
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

- Tradeoff between more data and better test error estimate

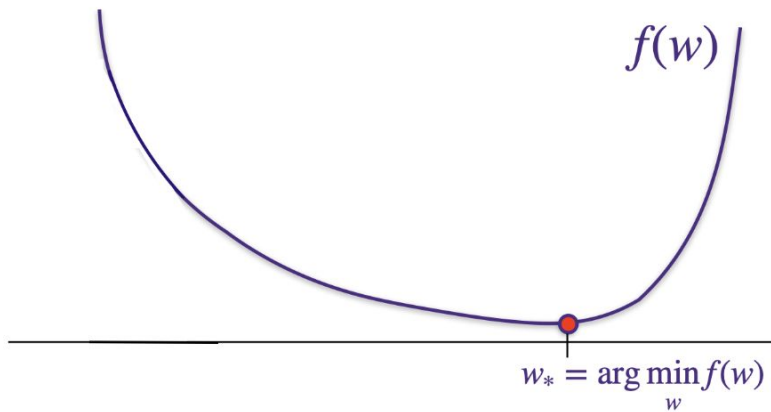
Gradient Descent

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

Gradient Descent

Purpose of this
exercise:

Understanding how
gradient descent
relates to
approximations, and
why it works.



2a

(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?

↳ To approximate a function around a point a

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

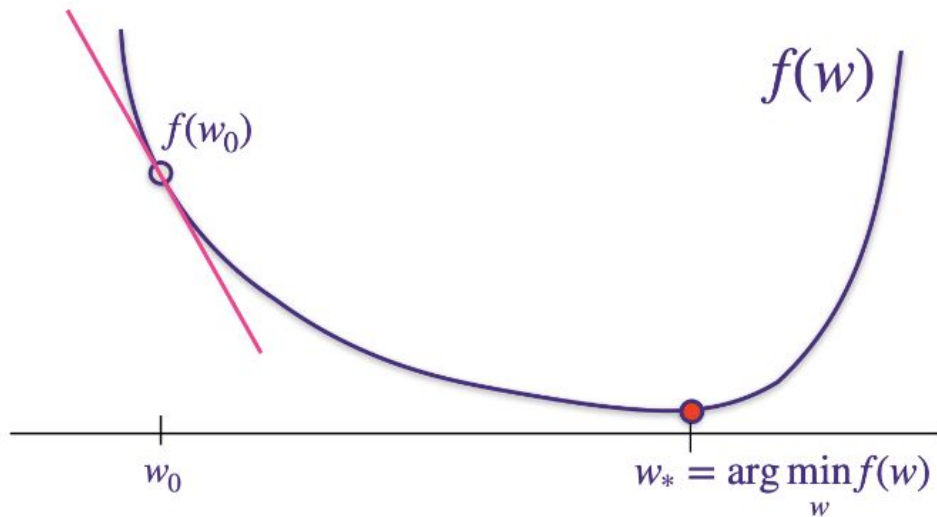
Exact at a , close around a

Better and better approximations

2a (answer)

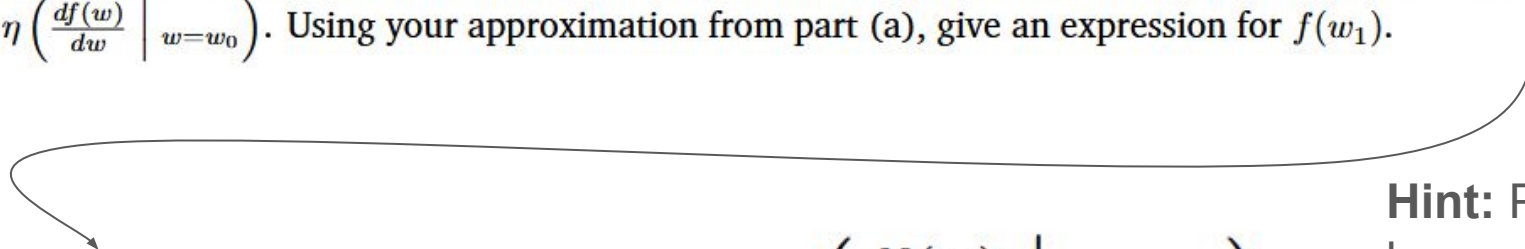
(a) For some w that is very close to w_0 , give the Taylor series approximation for $f(w)$ starting at $f(w_0)$.

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



2b

(b) Now, let us choose some $\eta > 0$ that is *very small*. With this very small η , 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 $f(w_1)$.


$$f(w) \approx f(w_0) + (w - w_0) \underbrace{\left(\frac{df(w)}{dw} \Big|_{w=w_0} \right)}_{\text{Fancy way of saying } f'(w_0)}.$$

Hint: Plug in here

Fancy way of saying $f'(w_0)$

(Derivative of $f(w)$ at w_0)

2b (answer)

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

$$\begin{aligned} f(w_1) &\approx f(w_0) + (w_1 - w_0) \left(\frac{df(w)}{dw} \Big|_{w=w_0} \right) \\ &= f(w_0) + \left(w_0 - \eta \left(\frac{df(w)}{dw} \Big|_{w=w_0} \right) - w_0 \right) \left(\frac{df(w)}{dw} \Big|_{w=w_0} \right) \\ &= f(w_0) - \eta \left(\frac{df(w)}{dw} \Big|_{w=w_0} \right)^2 \end{aligned}$$

2c

- (c) Given your expression for $f(w_1)$ from part (b), explain why, if η 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

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

Hint: Why
would this be
good?

2c (answer)

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

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

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

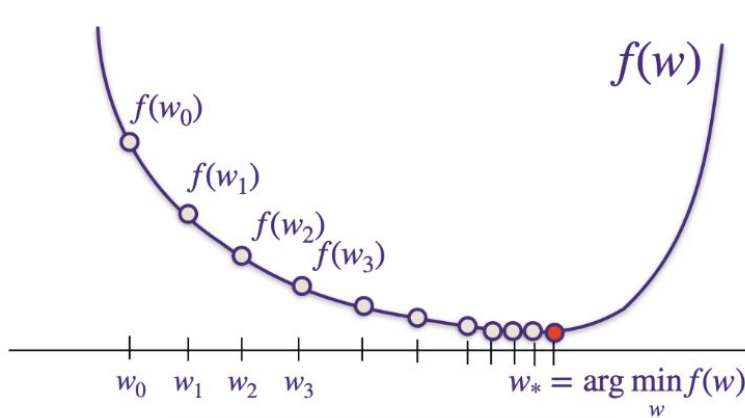
2d (answer)

Gradient descent is written as:

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

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

We visualize as:



**Convergence
guarantees iff
convex!**

Generalized Least Squares

Least Squares Proof(s)

Should look familiar...

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.

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

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

3.2a

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

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

3.2a (answer)

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

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

Solution:

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

3.2b

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

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

3.2b (answer)

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

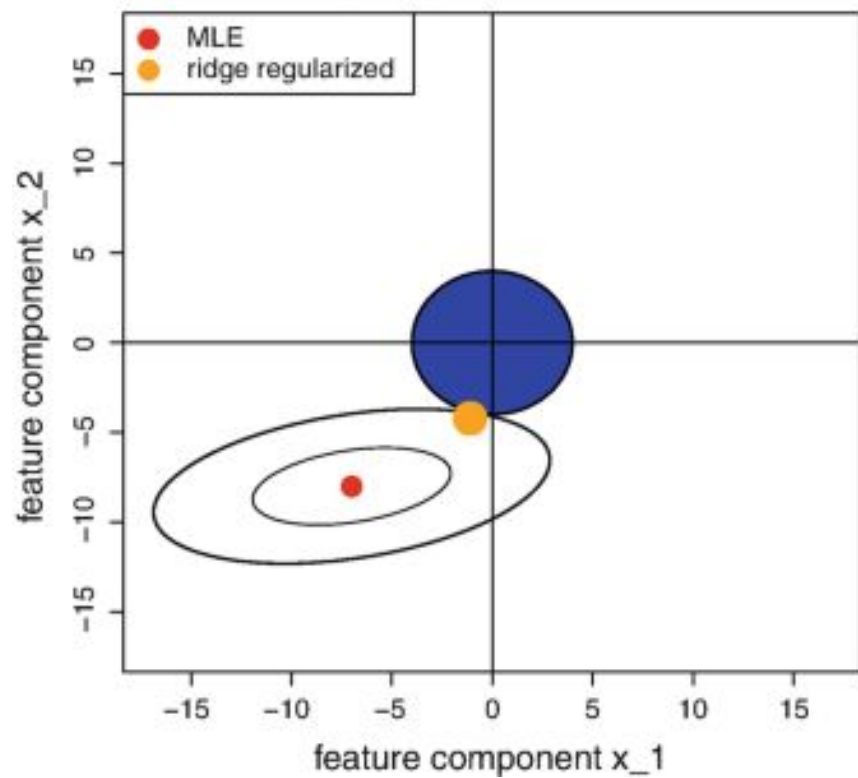
- (b) In the simple least squares case ($\lambda = 0$ above), what happens to the resulting $\hat{\omega}$ if we double the data matrix $X \in \mathbb{R}^{n \times d}$?

Solution:

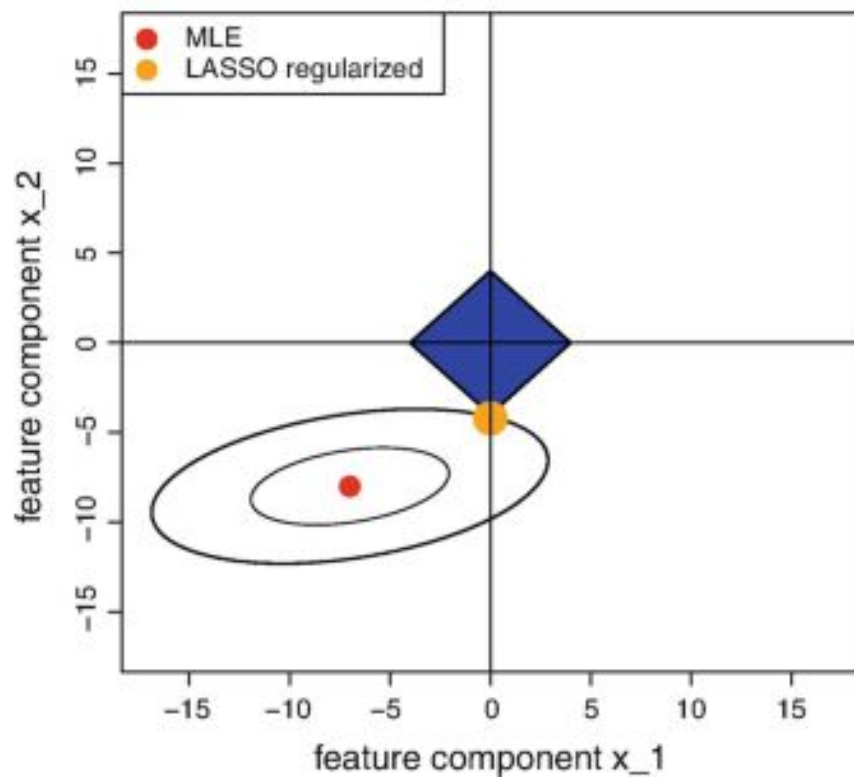
As can be seen from the formula $\hat{\omega} = (X^{\top} X)^{-1} X^{\top} y$, doubling X halves ω . This also makes sense intuitively because the error we are trying to minimize is $\|X\omega - y\|_2^2$, and if the X has doubled, while y has remained unchanged, then ω must compensate for it by reducing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

Ridge vs. LASSO

ridge regularization (L2)



LASSO regularization (L1)



Questions/Chat Time!