

```
[ ] 1 import numpy as np
    2 n = 20
    3 n_epochs = 1
    4 sampled_data = np.zeros(n*n_epochs)
    5 for t in range(n_epochs):
    6     for iteration in range(n):
    7         ind = np.random.randint(0, n)
    8         sampled_data[ind] = 1
    9 prob_sampled_single_epoch = np.sum(sampled_data > 0) / n
   10 prob_never_sampled_single_epoch = (1 - prob_sampled_single_epoch)
   11 print("The probability P(never) that one observation is never sampled during one epoch : ", prob_never_sampled_single_epoch)
```

The probability P(never) that one observation is never sampled during one epoch : 0.35

$$h_2 \textcircled{4} \quad \lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n = e$$

To prove  $n$  (large)  $\Pr_{\text{aver}} \approx 0.3679$

Given eqn similar to  $\lim_{n \rightarrow \infty} f(n)^{g(n)}$

for  $\Pr_{\text{aver}}$  we need to find  $\left(1 - \frac{1}{n}\right)^n$

$$\therefore \lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^n = \lim_{n \rightarrow \infty} \left(\left(1 + \frac{1}{n}\right)^n\right)^{-1}$$

$$\therefore \lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^n = e^{-1} = \frac{1}{e}$$

$$e \approx 2.718 \quad \therefore \frac{1}{e} = \frac{1}{2.718} \approx 0.3679$$

$h_2 \textcircled{2} \quad \text{Var}(\nabla f(x^{(i)}, y^{(i)})) = \sigma^2$ , variance of gradient over minibatch

$\rightarrow$  Total size =  $n$  minibatch =  $\tilde{n} \leq n$

$\therefore$  variance over minibatch  $\Rightarrow \frac{\sigma^2}{\tilde{n}}$

$h_2 \textcircled{e} \quad$  Multiply  $\tilde{n}$  by factor of  $k \rightarrow (k\tilde{n})$

$$\text{Variance becomes} \rightarrow \frac{\sigma^2}{k\tilde{n}} = \frac{1}{k} \left(\frac{\sigma^2}{\tilde{n}}\right)$$

New minibatch variance decreases by a factor of  $1/k$  compared to the original variance when we multiply it by a factor of  $k$  to minibatch size

# problem 1\_f

The authors argue: "When the minibatch size is multiplied by  $k$ , multiply the learning rate by  $k$ ."

Gradient noise decreases when we use SGD (stochastic gradient descent) while working with mini-batches as the gradient keeps the weights being updated during each (mini-batch \* epochs) iteration.

Now the author argues that when mini batch size is multiplied by a factor of  $k$  we need to multiply the learning rate by the same factor of  $k$ .

I feel the claim is incorrect because the gradient noise is proportional to the square root of the mini-batch size.

So, if you multiply the mini-batch size by a factor of  $k$ , the gradient noise will increase by a factor of  $\sqrt{k}$ .

There are a few ways to change the learning rate so that the gradient noise stays constant.

Using Learning rate schedulers

Using Exponential decay

Using Regularization techniques, such as L2 regularization, Dropout

Using Stochastic Gradient Descent (SGD) with Minibatch

Using a larger batch size

Using a more robust Optimization Algorithm



## h<sub>2</sub> ② XOR Problem

$$\text{Equation 6.2} \rightarrow J(\theta) = \frac{1}{4} \sum_{x \in X} (f^*(x) - f(x; \theta))^2$$

given  $f(x; \theta)$  where  $\theta$  consist of  $w$  &  $b$ .

$$f(x; w, b) = x^T w + b.$$

Solve for  $w \rightarrow w_1, w_2$  &  $b$

$$\frac{\partial J(\theta)}{\partial w_1} = \frac{1}{4} \sum_{x \in X} (f^*(x) - (x_1 w_1 + x_2 w_2 + b)) \cdot (2) \cdot (-x_1)$$

$$\therefore \frac{\partial J(\theta)}{\partial w_1} = -\frac{1}{2} x_1 \sum_{x \in X} (f^*(x) - (x_1 w_1 + x_2 w_2 + b)) \quad \text{--- (1)}$$

Similarly for  $w_2$  &  $b$  we get

$$\frac{\partial J(\theta)}{\partial w_2} = -\frac{1}{2} x_2 \sum_{x \in X} (f^*(x) - (x_1 w_1 + x_2 w_2 + b)) \quad \text{--- (2)}$$

$$\frac{\partial J(\theta)}{\partial b} = -\frac{1}{2} \sum_{x \in X} (f^*(x) - (x_1 w_1 + x_2 w_2 + b)) \quad \text{--- (3)}$$

Solving eq (1) for XOR Inputs

when  $x_1 = 0, x_2 = 0$   $f^*(x) = 0$  we get 0

when  $x_1 = 0, x_2 = 1$   $f^*(x) = 1$  we get 0

when  $x_1 = 1, x_2 = 0$   $f^*(x) = 1$  we get

$$\begin{aligned} & -\frac{1}{2} (f^*(x) - (x_1 w_1 + x_2 w_2 + b)) \\ & = -\frac{1}{2} (1) [1 - (w_1 + b)] = -(1 - (w_1 + b)) \end{aligned}$$

When  $x_1=1, x_2=1, f^*(x)=0$  we get

$$-\frac{1}{2}(1) [0 - (w_1 + w_2 + b)] = \frac{(w_1 + w_2 + b)}{2}$$

Summation to zero yields

$$0 + 0 + -\frac{(1 - (w_1 + b))}{2} + \frac{1}{2}(w_1 + w_2 + b) = 0$$

$$-1 + w_1 + b + w_1 + w_2 + b = 0$$

$$2b + 2w_1 + w_2 = 1 \quad \dots (4)$$

Similarly solving for eqn (2)

When  $x_2=0, x_1=0, f^*(x)=0$

we get 0

When  $x_2=0, x_1=1, f^*(x)=1$

we get 0

When  $x_2=0, x_1=0, f^*(x)=1$  we get

$$-\frac{1}{2}(1) [1 - (w_2 + b)] = -\frac{1}{2}[1 - w_2 - b]$$

When  $x_2=1, x_1=1, f^*(x)=0$  we get

$$-\frac{1}{2}(1) [0 - (w_1 + w_2 + b)] = \frac{-1}{2}[-w_1 - w_2 - b]$$

Summation to Zero yields

$$0 = 0 + 0 + \left(-\frac{1}{2}\right)[1 - w_2 - b] + \left(-\frac{1}{2}\right)(-w_1 - w_2 - b)$$

$$0 = -1 + w_2 + b + w_1 + w_2 + b$$

$$0 = 2b + 2w_2 + w_1 = 1$$

$$\therefore 2b + 2w_2 + w_1 = 1 \quad \dots (5)$$



Similarly solving for eq<sup>m</sup> (3)

When  $x_1 = 0, x_2 = 0, f^*(x) = 0$  we get

$$-\frac{1}{2} [0 - (b)] = +\frac{b}{2}$$

When  $x_1 = 0, x_2 = 1, f^*(x) = 2$  we get

$$-\frac{1}{2} [1 - (w_2 + b)] = -\frac{1}{2} (1 - w_2 - b)$$

When  $x_1 = 1, x_2 = 0, f^*(x) = 2$  we get

$$-\frac{1}{2} [1 - (w_1 + b)] = -\frac{1}{2} (1 - w_1 - b)$$

When  $x_1 = 1, x_2 = 1, f^*(x) = 0$  we get

$$-\frac{1}{2} [0 - (w_1 + w_2 + b)] = -\frac{1}{2} [-w_1 - w_2 - b]$$

Summation to Zero yields

$$\left(+\frac{b}{2}\right) + \left(-\frac{1}{2} (1 - w_2 - b)\right) + \left(-\frac{1}{2} (1 - w_1 - b)\right) + \left(-\frac{1}{2} (-w_1 - w_2 - b)\right) = 0$$

$$+\frac{b}{2} - \frac{1}{2} + \frac{w_2}{2} + \frac{b}{2} - \frac{1}{2} + \frac{w_1}{2} + \frac{b}{2} + \frac{w_1}{2} + \frac{w_2}{2} + \frac{b}{2} = 0$$

$$4b - 2 + 2w_1 + 2w_2 = 0$$

Divide by 2 on both sides

$$2b - 1 + w_1 + w_2 = 0$$

$$\Rightarrow 2b = (1 - w_1 - w_2) \dots (6)$$

Substituting eqn (6) in eqn (4) we get

$$(1 - w_1 - w_2) + 2w_1 + w_2 = 1$$

$$w_1 = 0 \quad \therefore \text{--- (7)}$$

Substituting eqn (6) in eqn (5) we get

$$(1 - w_1 - w_2) + 2w_2 + w_1 = 1$$

$$w_2 = 0 \quad \therefore \text{--- (8)}$$

Substituting eqn (7) & eqn (8) in eqn (6)

$$2b = (1 - w_1 - w_2)$$

$$= (1 - 0 - 0)$$

$$= 1$$

$$2b = 1$$

$$\therefore b = \frac{1}{2} = 0.5$$

```

1  # Training Linear Regression Model with user defined params + grid search
2
3  import numpy as np
4
5  X_tr = np.reshape(np.load("age_regression_Xtr.npy"), (-1, 48*48))
6  y_tr = np.load("age_regression_ytr.npy")
7  X_te = np.reshape(np.load("age_regression_Xte.npy"), (-1, 48*48))
8  y_te = np.load("age_regression_yte.npy")
9  print(X_tr.shape, y_tr.shape, X_te.shape, y_te.shape)
10
11  split_80_20: int = int(X_tr.shape[0]*0.8)
12  print('80 / 20 Split : ', split_80_20)
13  X_tr_80, y_tr_80 = X_tr[:split_80_20,:], y_tr[:split_80_20]
14  X_vl_20, y_vl_20 = X_tr[split_80_20:,:], y_tr[split_80_20:]
15  del X_tr, y_tr
16  print(X_tr_80.shape, y_tr_80.shape, X_te.shape, y_te.shape, X_vl_20.shape, y_vl_20.shape)
17
18  # custom grid params
19  grid_search_mini_batch_size : list = [160, 200, 400, 640]
20  grid_search_learning_rate : list = [3e-8, 3e-6, 3e-5, 3e-3]
21  grid_search_no_of_epochs : list = [16, 32, 48, 64]
22  grid_search_regularization_val : list = [1e-7, 1e-5, 1e-3, 1e-1]
23
24  # custom helper functions
25  def get_initialize_weights_and_bias_liner_reg(input_dims:int) -> tuple:
26      # return np.random.randn(input_dims), 0.1
27      return np.zeros(input_dims), 0.1
28
29  def get_y_hat_pred_func(x: np.ndarray, w: np.ndarray, b: np.ndarray) -> np.ndarray:
30      return np.dot(x, w) + b
31
32  def get_unregularized_loss_func(y_pred: np.ndarray, y_true: np.ndarray) -> float:
33      return 0.5 * np.mean((y_pred - y_true)**2)
34
35  def get_regularized_loss_func(y_pred: np.ndarray, y_true: np.ndarray, w: np.ndarray, alpha: float) -> float:
36      return get_unregularized_loss_func(y_pred, y_true) + 0.5 * alpha * np.dot(w.T, w)
37
38  # hyperparameter tuning
39  best_params: dict = {
40      'batch_size': 0,
41      'learning_rate': 0,
42      'num_epochs': 0,
43      'regularization_strength': 0,
44      'weights': 0,
45      'bias': 0,
46  }
47  best_loss: float = X_tr_80.shape[0]

```



```

47 best_loss: float = X_tr_80.shape[0]
48
49 # save history (log of each step)
50 result_params: list = list()
51
52 # Shuffle training data
53 shuffling_indices: np.ndarray = np.arange(X_tr_80.shape[0])
54 np.random.shuffle(shuffling_indices)
55 X_train_shuffled: np.ndarray = X_tr_80[shuffling_indices]
56 y_train_shuffled: np.ndarray = y_tr_80[shuffling_indices]
57
58 for initial_mini_batch_size in grid_search_mini_batch_size:
59     for initial_no_of_epochs in grid_search_no_of_epochs:
60         for initial_learning_rate in grid_search_learning_rate:
61             for initial_regularization_strength in grid_search_regularization_val:
62
63                 # Initialize weights and bias
64                 w, b = get_initialize_weights_and_bias_liner_reg(input_dims=X_tr_80.shape[1])
65
66                 # Training Loop
67                 for epoch in range(initial_no_of_epochs):
68
69                     for batch_start in range(0, X_tr_80.shape[0], initial_mini_batch_size):
70
71                         # taking mini batches from training data
72                         X_mini_batch = X_train_shuffled[batch_start:batch_start+initial_mini_batch_size,:]
73                         y_mini_batch = y_train_shuffled[batch_start:batch_start+initial_mini_batch_size]
74
75                         # we compute the y hat predict
76                         y_pred_mini = get_y_hat_pred_func(X_mini_batch, w, b)
77
78                         # we compute gradient of w
79                         grad_w = np.dot(X_mini_batch.T, (y_pred_mini - y_mini_batch)) / X_mini_batch.shape[0]
80
81                         # apply L2 regularization to learning weights
82                         grad_w += initial_regularization_strength * w
83
84                         # update learning weights
85                         w -= initial_learning_rate * grad_w
86
87                 # compute loss on validation set
88                 y_val_pred = get_y_hat_pred_func(X_vl_20, w, b)
89
90                 # compute and print the loss for monitoring training progress
91                 unreg_val_loss: float = round(get_unregularized_loss_func(y_val_pred, y_vl_20),4)
92                 reg_val_loss: float = round(get_regularized_loss_func(y_val_pred, y_vl_20, w, initial_regularization_strength),4)

```

```

93
94     # print(f"Epoch (epoch + 1)/{initial_no_of_epochs}, UnRegularized Val Loss: {unreg_val_loss}%, Regularized Val Loss: {reg_val_loss}%")
95     result_params.append({
96         'batch_size': initial_mini_batch_size,
97         'learning_rate': initial_learning_rate,
98         'num_epochs': initial_no_of_epochs,
99         'regularization_strength': initial_regularization_strength,
100         'bias': b,
101         'unregularized_loss': unreg_val_loss,
102         'regularized_loss': reg_val_loss,
103     })
104
105     # Check if this set of hyperparameters is the best so far
106     if unreg_val_loss < best_loss:
107         best_loss = unreg_val_loss
108         best_params = {
109             'batch_size': initial_mini_batch_size,
110             'learning_rate': initial_learning_rate,
111             'num_epochs': initial_no_of_epochs,
112             'regularization_strength': initial_regularization_strength,
113             'weights': w,
114             'bias': b,
115         }
116     else:
117         continue
118
119     # the best hyperparameters and corresponding loss
120     print(f"\nBest Hyperparameters: {best_params}, \nBest Validation Loss: {best_loss}")
121
122     # compute loss on validation set
123     y_te_pred = get_y_hat_pred_func(X_te, best_params['weights'], best_params['bias'])
124
125     # compute and display the loss for monitoring training progress
126     unreg_te_loss: float = round(get_unregularized_loss_func(y_te_pred, y_te), 4)
127     reg_te_loss: float = round(get_regularized_loss_func(y_te_pred, y_te, best_params['weights'], best_params['regularization_strength']), 4)
128     print(f"\nComputing Loss on Test Data Set\nUnRegularized Test Loss: {unreg_te_loss}%, Regularized Test Loss: {reg_te_loss}%")

```

(4000, 2304) (4000,) (2500, 2304) (2500,)
80 / 20 Split : 3200
(3200, 2304) (3200,) (2500, 2304) (2500,) (800, 2304) (800,)

Best Hyperparameters: {'batch\_size': 160, 'learning\_rate': 0.003, 'num\_epochs': 64, 'regularization\_strength': 1e-07, 'weights': array([ 0.88437179, 0.77671802, 0.41219991, ..., -0.28188384, -0.22051565, -0.25063656]), 'bias': 0.1},
Best Validation Loss: 114.3292

Computing Loss on Test Data Set
UnRegularized Test Loss: 112.385%, Regularized Test Loss: 112.385%

hw2 (4) Linear Regression with MAE minimization via Gradient Descent

$$f_{MAE}(w, b) = \frac{1}{n} \sum_{i=1}^n |(x^{(i)})^T w + b - y^{(i)}|$$

$$t = \left[ (x^{(i)})^T w + b - y^{(i)} \right]$$

(A) For prediction that overestimates  
 $\hat{y}^{(i)} < [(x^{(i)})^T w + b]$

(B) For predictions that underestimates  
 $\hat{y}^{(i)} > [(x^{(i)})^T w + b]$



Consider overestimation case  
When  $t > 0$   $\frac{\partial |t|}{\partial t} = 1$

When  $t < 0$   $\frac{\partial |t|}{\partial t} = \frac{\partial (-t)}{\partial t} = -1$

Considering the summation, yields

$$\frac{\partial f_{MAE}(w, b)}{\partial b} = \frac{1}{n} \sum_{i=1}^{n^+} \frac{\partial |t|}{\partial t}$$

$$= \frac{1}{n} [\text{no. of cases } t > 0 - \text{no. of cases } t < 0]$$

$$= \frac{n^+}{n} \quad \text{--- (1)}$$

Consider underestimation case

When  $t > 0$   $\frac{\partial |t|}{\partial t} = 1$

When  $t < 0$   $\frac{\partial |t|}{\partial t} = \frac{\partial (-t)}{\partial t} = -1$

Similarly we get  $\left(\frac{-n^-}{n}\right)$

$$\begin{aligned} \therefore \nabla_b f_{MAE}(w, b) &= \frac{n^+}{n} - \frac{n^-}{n} \\ &= \frac{1}{n} (n^+ - n^-) \end{aligned}$$

hw (4b) Expression for  $\nabla_w f_{MAE}(w, b)$

$$t = ((x^{(p)})^T w + b - y^{(p)})$$

$$\text{When } t > 0 \quad \frac{\partial |t|}{\partial w} = x^{(p)}$$

$$\text{When } t < 0 \quad \frac{\partial |t|}{\partial w} = -x^{(p)}$$

Case where predictions overestimate  
and underestimate we get

$$\frac{\partial f_{MAE}(w, b)}{\partial w} = \frac{1}{n} \left[ \text{no. of cases } t > 0 \cdot x^{(p)} - \text{no. of cases } t < 0 \cdot x^{(p)} \right]$$

$$\nabla_w f_{MAE}(w, b) = \frac{1}{n} \sum_{i=1}^n \left[ (n^+ \cdot x^{(p)}) - (n^- \cdot x^{(p)}) \right]$$

```
1 # problem 4_c
2
3 import numpy as np
4
5 def get_update_b_with_gradient_descent(X: np.ndarray, y: np.ndarray, w: np.ndarray, b: float, lr: float) -> float:
6
7     # Compute the predicted values
8     y_hat_pred: np.ndarray = np.dot(X, w) + b
9
10    # we compute the gradient b for using MAE equation
11    grad_b: np.ndarray = (1/X.shape[0]) * np.sum(np.sign(y_hat_pred-y))
12
13    # update parameter b
14    b_new: float = b - lr * grad_b
15
16    return b_new
```

```
1 print(b)
2 print(get_update_b_with_gradient_descent(X=X_te, y=y_te, w=w, b=b, lr=initial_learning_rate))
```

```
0.1
0.1024096
```



$$\text{hw}_2(5) \quad p(y|x) = \mathcal{N}(\mu = x^T \omega, \sigma^2)$$

$$= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y - x^T \omega)^2}{2\sigma^2}\right)$$

to prove:

$$\omega = \left( \sum_{i=1}^n x^{(i)} (x^{(i)})^T \right)^{-1} \left( \sum_{i=1}^n x^{(i)} y^{(i)} \right)$$

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x^{(i)T} \omega - y^{(i)})^2$$

$$P(D | w, \sigma^2) = \prod_{i=1}^n P(y^{(i)} | x^{(i)}, w, \sigma^2)$$

Taking log of likelihood

$$\log P(D | w, \sigma^2) = \sum_{i=1}^n \log P(y^{(i)} | x^{(i)}, w, \sigma^2)$$

$$\text{Now, } \log P(y^{(i)} | x^{(i)}, w, \sigma^2) = -\frac{1}{2} \log (2\pi\sigma^2) - \frac{(y^{(i)} - (x^{(i)})^T w)^2}{2\sigma^2}$$

$$= -\frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - (x^{(i)})^T w)^2$$

Taking  $\nabla$  wrt  $w$  we get after setting to 0

$$\frac{1}{\sigma^2} \sum_{i=1}^n (x^{(i)} (x^{(i)})^T w - x^{(i)} y^{(i)}) = 0$$

$$\sum_{i=1}^n x^{(i)} x^{(i)T} w = \sum_{i=1}^n x^{(i)} y^{(i)}$$

$$\therefore w = \left( \sum_{i=1}^n x^{(i)} x^{(i)T} \right)^{-1} \left( \sum_{i=1}^n x^{(i)} y^{(i)} \right)$$

Now for  $\sigma^2$

$$\frac{\partial}{\partial \sigma^2} P(D | w, \sigma^2) = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (y^{(i)} - (x^{(i)})^T w)^2$$

setting  $\frac{\partial}{\partial \sigma^2} P(D|\omega, \sigma^2) = 0$

$$-\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (y^{(i)} - (x^{(i)})^T \omega)^2 = 0$$

multiply both sides by  $2\sigma^4$

$$-n\sigma^2 + \sum_{i=1}^n (y^{(i)} - (x^{(i)})^T \omega)^2 = 0$$

$$n\sigma^2 = \sum_{i=1}^n (y^{(i)} - (x^{(i)})^T \omega)^2$$

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n ((x^{(i)})^T \omega - y^{(i)})^2$$