

④ a) $h(x) = g(f(x))$, both $f(x)$ and $g(x)$ are convex over the domain

$$(g \circ f)(\lambda x + (1-\lambda)y) = g(f(\lambda x + (1-\lambda)y))$$

$$\stackrel{\textcircled{*}}{\leq} g(\lambda f(x) + (1-\lambda)f(y))$$

$$\leq \lambda g(f(x)) + (1-\lambda)g(f(y))$$

$$= \lambda (g \circ f)(x) + (1-\lambda)(g \circ f)(y)$$

} applying the convex definition

$\textcircled{*}$ this only holds assuming g preserves the monotonicity of the underlying argument, or in other words, if g is non-decreasing.

⑤ c) $f(x_1, x_2) = 0.5x_1^2 + x_2^2 + 2x_1 + x_2 + \cos(\sin(\sqrt{11}))$

$$\nabla f = (x_1 + 2, 2x_2 + 1)$$

, since the function is convex we can say that $(-2, -1/2)$ is the minimizer without extra calculations

$$\nabla f = 0 \Leftrightarrow \begin{cases} x_1 + 2 = 0 \\ 2x_2 + 1 = 0 \end{cases} \Leftrightarrow \begin{cases} x_1 = -2 \\ x_2 = -1/2 \end{cases}$$

$$\begin{aligned} \text{b) } x^{(1)} &= x^{(0)} - \tau \cdot \nabla f(x^{(0)}) \\ &= (0, 0) - 1 \cdot (2, 1) = (-2, -1) \end{aligned}$$

$$\begin{aligned} x^{(2)} &= x^{(1)} - \tau \cdot \nabla f(x^{(1)}) \\ &= (-2, -1) - (0, -1) \\ &= (-2, 0) \end{aligned}$$

c) No, because it's going to keep alternating between $(-2, -1)$ and $(-2, 0)$.

This happens because the learning rate is too high so we overshoot the actual solution to the problem and get stuck in a loop.

An obvious solution to this problem would be to use other methods of implementing gradient descent, such as adaptive grad. or other more advanced methods, or simply using a smaller learning rate such as 0.1.

- ⑦ a) The shaded region isn't convex since you can draw a line between two points in S and verify that the resulting line isn't contained in S . For example, $(1, 3.5) \rightarrow (3.5, 6)$.

Take the formal definition:

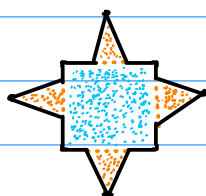
• S is convex iff $\forall x, y \in S: \lambda x + (1-\lambda)y \in S, \lambda \in [0, 1]$

\Rightarrow we can clearly confirm that taking $x = (1, 3.5); y = (3.5, 6); \lambda = 1/2$

$$1/2 \cdot (1, 3.5) + 1/2 \cdot (3.5, 6) = (0.5 + 1.75, 1.75 + 3) = (2.25, 4.75)$$

does not result in a point contained in S , which is equivalent to saying S is not convex.

- b) Since the algorithm only works in convex regions, we assume it ignores the existence of local minima. Such that the only way of using it to compute the global minimum of S , would be to break S into convex regions and applying the algorithm individually to each region S' . By comparing the solutions computed, it would be possible to compute the global minimum over S .



Example split of S resulting in only convex regions

code

November 30, 2021

1 Programming assignment 3: Optimization - Logistic Regression

```
[ ]: import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline

from sklearn.datasets import load_breast_cancer
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, f1_score
```

1.1 Your task

In this notebook code skeleton for performing logistic regression with gradient descent is given. Your task is to complete the functions where required. You are only allowed to use built-in Python functions, as well as any **numpy** functions. No other libraries / imports are allowed.

For numerical reasons, we actually minimize the following loss function

$$\mathcal{L}(\mathbf{w}) = \frac{1}{N}NLL(\mathbf{w}) + \frac{1}{2}\lambda\|\mathbf{w}\|_2^2$$

where $NLL(\mathbf{w})$ is the negative log-likelihood function, as defined in the lecture (see Eq. 33).

1.2 Exporting the results to PDF

Once you complete the assignments, export the entire notebook as PDF and attach it to your homework solutions. The best way of doing that is 1. Run all the cells of the notebook. 2. Export/download the notebook as PDF (File -> Download as -> PDF via LaTeX (.pdf)). 3. Concatenate your solutions for other tasks with the output of Step 2. On a Linux machine you can simply use **pdfunite**, there are similar tools for other platforms too. You can only upload a single PDF file to Moodle.

Make sure you are using **nbconvert** Version 5.5 or later by running **jupyter nbconvert --version**. Older versions clip lines that exceed page width, which makes your code harder to grade.

1.3 Load and preprocess the data

In this assignment we will work with the UCI ML Breast Cancer Wisconsin (Diagnostic) dataset <https://goo.gl/U2Uwz2>.

Features are computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. They describe characteristics of the cell nuclei present in the image. There are 212 malignant examples and 357 benign examples.

```
[ ]: X, y = load_breast_cancer(return_X_y=True)

# Add a vector of ones to the data matrix to absorb the bias term
X = np.hstack([np.ones([X.shape[0], 1]), X])

# Set the random seed so that we have reproducible experiments
np.random.seed(123)

# Split into train and test
test_size = 0.3
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=test_size)
```

1.4 Task 1: Implement the sigmoid function

```
[ ]: def sigmoid(t):
    """
    Applies the sigmoid function elementwise to the input data.

    Parameters
    -----
    t : array, arbitrary shape
        Input data.

    Returns
    -----
    t_sigmoid : array, arbitrary shape.
        Data after applying the sigmoid function.
    """

    res = 1/(1+np.exp(-t))
    return res
```

1.5 Task 2: Implement the negative log likelihood

As defined in Eq. 33

```
[ ]: def negative_log_likelihood(X, y, w):
    """
    Negative Log Likelihood of the Logistic Regression.

    Parameters
    -----
    X : array, shape [N, D]
        (Augmented) feature matrix.
```

```

y : array, shape [N]
    Classification targets.
w : array, shape [D]
    Regression coefficients (w[0] is the bias term).

Returns
-----
nll : float
    The negative log likelihood.
"""

# since sigmoid overflows to 0 we update values to eps
sig = sigmoid(X @ w)
sig = np.where(sig == 0, np.finfo(np.float64).eps, sig)

neg_sig = 1-sig
neg_sig = np.where(neg_sig == 0, np.finfo(np.float64).eps, neg_sig)

x1 = np.sum(y @ np.log(sig))
x2 = np.sum((1-y) @ np.log(neg_sig))

return -(x1+x2)

```

1.5.1 Computing the loss function $\mathcal{L}(w)$ (nothing to do here)

```

[ ]: def compute_loss(X, y, w, lambda):
    """
    Negative Log Likelihood of the Logistic Regression.

    Parameters
    -----
    X : array, shape [N, D]
        (Augmented) feature matrix.
    y : array, shape [N]
        Classification targets.
    w : array, shape [D]
        Regression coefficients (w[0] is the bias term).
    lambda : float
        L2 regularization strength.

    Returns
    -----
    loss : float
        Loss of the regularized logistic regression model.
    """
    # The bias term w[0] is not regularized by convention

```

```

    return negative_log_likelihood(X, y, w) / len(y) + lambda * 0.5 * np.linalg.
↪norm(w[1:])**2

```

1.6 Task 3: Implement the gradient $\nabla_w \mathcal{L}(w)$

Make sure that you compute the gradient of the loss function $\mathcal{L}(w)$ (not simply the NLL!)

```

[ ]: def get_gradient(X, y, w, mini_batch_indices, lambda):
    """
    Calculates the gradient (full or mini-batch) of the negative log_
↪likelihood w.r.t. w.

    Parameters
    -----
    X : array, shape [N, D]
        (Augmented) feature matrix.
    y : array, shape [N]
        Classification targets.
    w : array, shape [D]
        Regression coefficients (w[0] is the bias term).
    mini_batch_indices: array, shape [mini_batch_size]
        The indices of the data points to be included in the (stochastic)_
↪calculation of the gradient.
        This includes the full batch gradient as well, if mini_batch_indices =_
↪np.arange(n_train).
    lambda: float
        Regularization strength. lambda = 0 means having no regularization.

    Returns
    -----
    dw : array, shape [D]
        Gradient w.r.t. w.
    """

    sig = sigmoid(X[mini_batch_indices] @ w)
    sum = (y[mini_batch_indices]-sig) @ X[mini_batch_indices]

    return -sum/len(mini_batch_indices) + lambda*w

```

1.6.1 Train the logistic regression model (nothing to do here)

```

[ ]: def logistic_regression(X, y, num_steps, learning_rate, mini_batch_size, lambda,
↪verbose):
    """
    Performs logistic regression with (stochastic) gradient descent.

    Parameters

```

```

-----
X : array, shape [N, D]
    (Augmented) feature matrix.
y : array, shape [N]
    Classification targets.
num_steps : int
    Number of steps of gradient descent to perform.
learning_rate: float
    The learning rate to use when updating the parameters w.
mini_batch_size: int
    The number of examples in each mini-batch.
    If mini_batch_size=n_train we perform full batch gradient descent.
lmbda: float
    Regularization strength. lmbda = 0 means having no regularization.
verbose : bool
    Whether to print the loss during optimization.

Returns
-----
w : array, shape [D]
    Optimal regression coefficients (w[0] is the bias term).
trace: list
    Trace of the loss function after each step of gradient descent.
"""

trace = [] # saves the value of loss every 50 iterations to be able to plot
→it later
n_train = X.shape[0] # number of training instances-1 :D

w = np.zeros(X.shape[1]) # initialize the parameters to zeros

# run gradient descent for a given number of steps
for step in range(num_steps):
    permuted_idx = np.random.permutation(n_train) # shuffle the data

    # go over each mini-batch and update the paramters
    # if mini_batch_size = n_train we perform full batch GD and this loop
→runs only once
    for idx in range(0, n_train, mini_batch_size):
        # get the random indices to be included in the mini batch
        mini_batch_indices = permuted_idx[idx:idx+mini_batch_size]
        gradient = get_gradient(X, y, w, mini_batch_indices, lmbda)

        # update the parameters
        w = w - learning_rate * gradient

```

```

        # calculate and save the current loss value every 50 iterations
    if step % 50 == 0:
        loss = compute_loss(X, y, w, lambda)
        trace.append(loss)
        # print loss to monitor the progress
        if verbose:
            print('Step {0}, loss = {1:.4f}'.format(step, loss))
    return w, trace

```

1.7 Task 4: Implement the function to obtain the predictions

```

[ ]: def predict(X, w):
    """
    Parameters
    -----
    X : array, shape [N_test, D]
        (Augmented) feature matrix.
    w : array, shape [D]
        Regression coefficients (w[0] is the bias term).

    Returns
    -----
    y_pred : array, shape [N_test]
        A binary array of predictions.
    """

    return np rint(sigmoid(X@w))

```

1.7.1 Full batch gradient descent

```

[ ]: # Change this to True if you want to see loss values over iterations.
    verbose = False

```

```

[ ]: n_train = X_train.shape[0]
    w_full, trace_full = logistic_regression(X_train,
                                             y_train,
                                             num_steps=8000,
                                             learning_rate=1e-5,
                                             mini_batch_size=n_train,
                                             lambda=0.1,
                                             verbose=verbose)

```

```

[ ]: n_train = X_train.shape[0]
    w_minibatch, trace_minibatch = logistic_regression(X_train,
                                                        y_train,
                                                        num_steps=8000,
                                                        learning_rate=1e-5,

```



```
mini_batch_size=50,  
lambda=0.1,  
verbose=verbose)
```

Our reference solution produces, but don't worry if yours is not exactly the same.

Full batch: accuracy: 0.9240, f1_score: 0.9384

Mini-batch: accuracy: 0.9415, f1_score: 0.9533

```
[ ]: y_pred_full = predict(X_test, w_full)  
     y_pred_minibatch = predict(X_test, w_minibatch)  
  
     print('Full batch: accuracy: {:.4f}, f1_score: {:.4f}'  
           .format(accuracy_score(y_test, y_pred_full), f1_score(y_test,   
↪y_pred_full)))  
     print('Mini-batch: accuracy: {:.4f}, f1_score: {:.4f}'  
           .format(accuracy_score(y_test, y_pred_minibatch), f1_score(y_test,   
↪y_pred_minibatch)))
```

Full batch: accuracy: 0.9240, f1_score: 0.9384

Mini-batch: accuracy: 0.9415, f1_score: 0.9533

```
[ ]: plt.figure(figsize=[15, 10])  
     plt.plot(trace_full, label='Full batch')  
     plt.plot(trace_minibatch, label='Mini-batch')  
     plt.xlabel('Iterations * 50')  
     plt.ylabel('Loss  $\mathcal{L}(\mathbf{w})$ ')  
     plt.legend()  
     plt.show()
```

