

4.

a) not true.

consider:  $f(x) = e^x$ , convex function,  
 $g(x) = -\ln x$ , convex function.

$\Rightarrow h(x) = g(f(x)) = -\ln e^x = -x$ , not convex or concave

b) function  $h(x)$  is convex if:  $h(\theta x_1 + (1-\theta)x_2) \leq \theta h(x_1) + (1-\theta)h(x_2)$ ,  
 $\forall x_1, x_2 \in \mathbb{R}, \theta \in [0, 1]$

$f(x)$  is convex  $\Rightarrow f(\theta x_1 + (1-\theta)x_2) \leq \theta f(x_1) + (1-\theta)f(x_2)$

since  $g$  is non-decreasing  $\Rightarrow g(f(\theta x_1 + (1-\theta)x_2)) \leq g(\theta f(x_1) + (1-\theta)f(x_2))$

$g(x)$  is convex  $\Rightarrow g(\theta x_1 + (1-\theta)x_2) \leq \theta g(x_1) + (1-\theta)g(x_2)$

with  $x_1 = f(x_1)$ ,  $x_2 = f(x_2) \Rightarrow g(\theta f(x_1) + (1-\theta)f(x_2)) \leq \theta g(f(x_1)) + (1-\theta)g(f(x_2))$

so:  $g(f(\theta x_1 + (1-\theta)x_2)) \leq \theta g(f(x_1)) + (1-\theta)g(f(x_2)) \quad (1-\theta)g(f(x_2))$

$\Rightarrow h(f(\theta x_1 + (1-\theta)x_2)) \leq \theta h(x_1) + (1-\theta)h(x_2)$ , proved

5.

a)  $\nabla f(x_1, x_2) = \begin{pmatrix} x_1+2 \\ x_2+1 \end{pmatrix} = 0 \Rightarrow x_1 = -2, x_2 = -1 \Rightarrow x^* = (-2, -1)$

b) step 1:  $\nabla f(x^{(0)}) = \begin{pmatrix} 0+2 \\ 0+1 \end{pmatrix} = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \Rightarrow x^{(1)} = x^{(0)} - T \nabla f(x^{(0)}) = \begin{pmatrix} 0 \\ 0 \end{pmatrix} - 1 \cdot \begin{pmatrix} 2 \\ 1 \end{pmatrix}$

step 2:  $\nabla f(x^{(1)}) = \begin{pmatrix} -2+2 \\ -1+1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \Rightarrow x^{(2)} = x^{(1)} - T \nabla f(x^{(1)}) = \begin{pmatrix} -2 \\ -1 \end{pmatrix} = \begin{pmatrix} -2 \\ -1 \end{pmatrix} \Rightarrow x^{(2)} = (-2, -1)$

c) will converge to the true minimizer  $x^*$ .

because after step 1, the gradient becomes zero, meaning  $x^{(2)}$  remains unchange and stay at the optimal point.

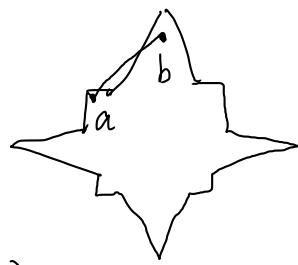
first: Hessian matrix:  $H(f) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} 0.5 & 0 \\ 0 & 2 \end{bmatrix}$

, which is positive-definite

therefore,  $f(x_1, x_2)$  is convex

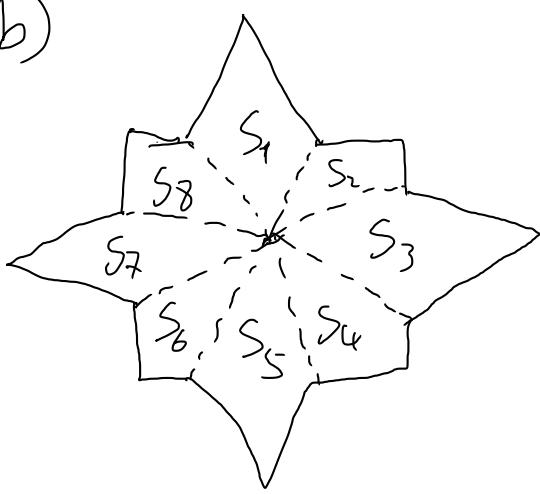
7.

a) not convex.



because the line  $ab$  is not entirely contained within the region.

b)



divide region  $S$  into smaller convex subregion  $S_1, S_2 \dots S_8$ , Apply algorithm ConvOpt to each convex subregion  $S_i$  and get local minimum of  $f$ . Finally, compare all local minimum and select the smallest local minimum as global minimum.

# Programming assignment 3: Optimization - Logistic Regression

```
In [8]: 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
```

## 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 Slide 39).

## Load and preprocess the data

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

```
In [9]: 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)
```

## Task 1: Implement the sigmoid function

```
In [10]: 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.
    """
    # TODO
    t_sigmoid = 1 / (1 + np.exp(-t))

    return t_sigmoid
```

## Task 2: Implement the negative log likelihood

As defined in Eq. 33

```
In [11]: 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.
    """
    # TODO
    z = X @ w

    predictions = 1 / (1 + np.exp(-z))

    nll = -np.sum(y * np.log(predictions) + (1 - y) * np.log(1 - predictions))

    return nll
```

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

```
In [12]: def compute_loss(X, y, w, lmbda):
    """
        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).
    lmbda : 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) + lmbda * 0.5 * np.linalg.n
```

## Task 3: Implement the gradient $\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w})$

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

```
In [13]: def get_gradient(X, y, w, mini_batch_indices, lmbda):
    """
    Calculates the gradient (full or mini-batch) of the negative log likelihood.

    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.
        This includes the full batch gradient as well, if mini_batch_indices = np.arange(N).
    lmbda: float
        Regularization strength. lmbda = 0 means having no regularization.

    Returns
    -----
    dw : array, shape [D]
        Gradient w.r.t. w.
    """
    # TODO
    X_batch = X[mini_batch_indices]
    y_batch = y[mini_batch_indices]

    z = X_batch @ w

    predictions = 1 / (1 + np.exp(-z))

    error = predictions - y_batch

    dw = (X_batch.T @ error) / len(y_batch)

    dw[1:] += lmbda * w[1:]

    return dw
```

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

```
In [14]: def logistic_regression(X, y, num_steps, learning_rate, mini_batch_size, lmbda, v
    """
        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
    n_train = X.shape[0] # number of training instances

    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 parameters
        # if mini_batch_size = n_train we perform full batch GD and this loop runs
        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, lmbda)
            trace.append(loss)
            # print loss to monitor the progress
            if verbose:
                print('Step {0}, loss = {1:.4f}'.format(step, loss))
```

```
return w, trace
```

## **Task 4: Implement the function to obtain the predictions**

```
In [15]: 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.
    """
    # TODO
    z = X @ w

    predictions = 1 / (1 + np.exp(-z))

    y_pred = (predictions >= 0.5).astype(int)

    return y_pred
```

## Full batch gradient descent

```
In [16]: # Change this to True if you want to see loss values over iterations.  
verbose = False
```

```
In [18]: 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,
                                                    lmbda=0.1,
                                                    verbose=verbose)
```

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

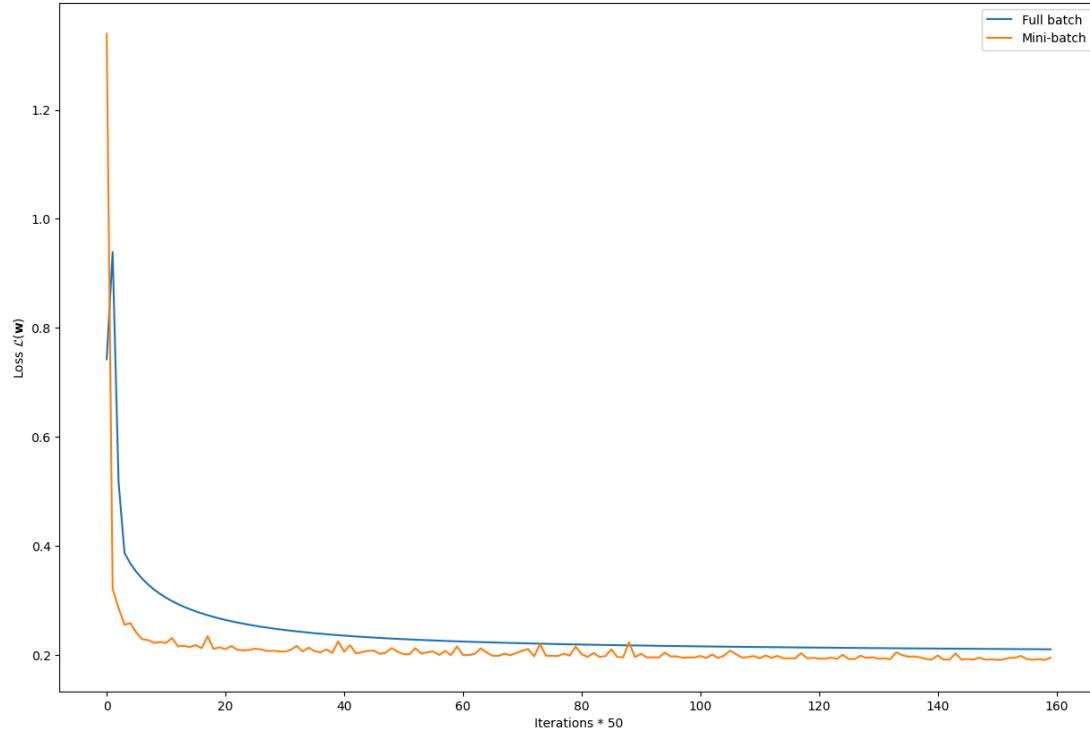
```
Full batch: accuracy: 0.9240, f1_score: 0.9384
Mini-batch: accuracy: 0.9415, f1_score: 0.9533
```

```
In [19]: 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_m))
```

```
Full batch: accuracy: 0.9240, f1_score: 0.9384
Mini-batch: accuracy: 0.9415, f1_score: 0.9533
```

```
In [20]: 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()
```



```
In [ ]:
```