

Convexity

Q1:

1. Convex.

For $y \in (1, 50)$, $\log(y) > 0$; but for $x \in (-100, 100)$, $-x^2 \leq 0$, thus $\min\{-x^2, \log(y)\} = -x^2$

$$\begin{aligned} f(x, y) &= x^2 + 2y + \cos(\sin(\sqrt{\pi})) - \min\{-x^2, \log(y)\} \\ &= x^2 + 2y + \cos(\sin(\sqrt{\pi})) + x^2 \end{aligned} \quad (1)$$

The second derivatives of every element in the function are constant larger or equal than zero (namely 2,0,0,2), thus all part of function are convex.

$$f(x, y) = \underbrace{x^2}_{\text{convex}} + \underbrace{2y}_{\text{convex}} + \underbrace{\cos(\sin(\sqrt{\pi}))}_{\text{const. larger than 0}} \underbrace{-\min\{-x^2, \log(y)\}}_{\text{convex}}$$

2. Not convex in $D = (1, \infty)$.

Second derivative of $-\log(x)$ and x^3 are $\frac{1}{x^2}$ and $6x$, both are larger than zero in given domain.

$$\begin{aligned} f(x) &= \log(x) - x^3 \\ &= \overbrace{-(-\log(x) + x^3)}^{\text{concave}} \end{aligned} \quad (2)$$

convex

3. Not convex in given domain.

For positive real numbers

$$\begin{aligned} &\log(3x+1) - (-x^4 - 3x^2 + 8x - 42) \\ &= \log(3x+1) + x^4 + 3x^2 - 8x + 42 \\ &= \log(3x+1) + x^4 + 3\left(x - \frac{4}{3}\right)^2 + \frac{110}{3} > 0 \end{aligned} \quad (3)$$

Thus,

$$f(x) = -x^4 - 3x^2 + 8x - 42$$

Then analysis this function: the second derivative of x^4 , $3x^2$ and $-8x + 42$ are $12x^2$, 6 and 0, all larger or equal than zero.

$$\begin{aligned} f(x) &= -x^4 - 3x^2 + 8x - 42 \\ &= -\overbrace{\left(\underbrace{x^4}_{\text{convex}} + \underbrace{3x^2}_{\text{convex}} - \underbrace{8x + 42}_{\text{convex}}\right)}^{\text{concave}} \end{aligned} \quad (4)$$

4. Not convex.

$$\begin{aligned} \frac{\partial^2 f}{\partial x^2} &= 6yx - 2y \\ &= 2y(3x - 1) \end{aligned} \quad (5)$$

Its second derivative cross the zero at $x = \frac{1}{3}$, thus the function is not convex.

Q2:

As $\lambda \in [0, 1]$, and f_1, f_2 are convex functions

$$\begin{aligned}
 h(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) &= \max(f_1(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}), f_2(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y})) \\
 &\leq \max(\lambda f_1(\mathbf{x}) + (1 - \lambda) f_1(\mathbf{y}), \lambda f_2(\mathbf{x}) + (1 - \lambda) f_2(\mathbf{y})) \\
 &\leq \lambda \max(f_1(\mathbf{x}), f_2(\mathbf{x})) + (1 - \lambda) \max(f_1(\mathbf{y}), f_2(\mathbf{y})) \\
 &= \lambda h(\mathbf{x}) + (1 - \lambda) h(\mathbf{y})
 \end{aligned} \tag{6}$$

Convexity proved.

Q3:

As f_1, f_2 are convex functions. Assume that f_1 is a non-decreasing function.

$$\begin{aligned}
 g(\lambda x + (1 - \lambda)y) &= f_1(f_2(\lambda x + (1 - \lambda)y)) \\
 &\leq f_1(\lambda f_2(x) + (1 - \lambda)f_2(y))
 \end{aligned} \tag{7}$$

Let $f_2(x) = x', f_2(y) = y'$

$$\begin{aligned}
 g(\lambda x + (1 - \lambda)y) &= f_1(f_2(\lambda x + (1 - \lambda)y)) \\
 &\leq f_1(\lambda f_2(x) + (1 - \lambda)f_2(y)) \\
 &\leq \lambda f_1(x') + (1 - \lambda)f_1(y') \\
 &= \lambda f_1(f_2(x)) + (1 - \lambda)f_1(f_2(y)) \\
 &= \lambda g(x) + (1 - \lambda)g(y)
 \end{aligned} \tag{8}$$

Convexity proved (Convexity is preserved when f_1 is non-decreasing).

Minimization of convex functions**Q4:**

Suppose there exists at least one point $\mathbf{m} \in \mathbb{R}^N$, such that $f(\mathbf{m}) < f(\boldsymbol{\theta}^*)$.

Because the convexity of f ,

$$\begin{aligned}
 f(\lambda \boldsymbol{\theta}^* + (1 - \lambda)f(\mathbf{m})) &\leq \lambda f(\boldsymbol{\theta}^*) + (1 - \lambda)f(\mathbf{m}) \\
 &< \lambda f(\boldsymbol{\theta}^*) + (1 - \lambda)f(\boldsymbol{\theta}^*) \\
 &= f(\boldsymbol{\theta}^*)
 \end{aligned} \tag{9}$$

However, when $\lambda \rightarrow 1$, $f(\lambda \boldsymbol{\theta}^* + (1 - \lambda)f(\mathbf{m})) \rightarrow f(\boldsymbol{\theta}^*)$, this means that in a sufficiently small region, there exist some points that the function value at those points are smaller than the local minimum, which of course is not correct. Thus, the statement that a local minimum of a convex function is also the global minimum has been proven.

Gradient Descent**Q5:**

Programming assignment 5: Optimization: Logistic regression

```
In [12]: 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

L(w) = 1/N * NLL(w) + 1/2 * lambda ||w||_2^2

where *NLL*(w) is the negative log-likelihood function, as defined in the lecture (Eq. 33)

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. Download the notebook in HTML (click File > Download as > .html)
- 3. Convert the HTML to PDF using e.g. <https://www.sejda.com/html-to-pdf> or `wkhtmltopdf` for Linux ([tutorial](#))
- 4. Concatenate your solutions for other tasks with the output of Step 3. 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.

This way is preferred to using `nbconvert`, since `nbconvert` clips lines that exceed page width and makes your code harder to grade.

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.

```
In [13]: 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 [14]: 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 [15]: 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
    nll = 0
    for i in range(len(y)):
        nll = nll + (y[i]*np.log(sigmoid(np.dot(X[i],w))) + (1-y[i])*np.log(1-sigmoid(np.dot(X[i],w))))
    nll = -nll
    return nll
```

Computing the loss function L(w) (nothing to do here)

```
In [16]: 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 * np.linalg.norm(w[1:])**2
```

Task 3: Implement the gradient ∇_wL(w)

Make sure that you compute the gradient of the loss function *L*(w) (not simply the NLL!)

```
In [84]: def get_gradient(X, y, w, mini_batch_indices, lmbda):
    """
    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.
    t.
        This includes the full batch gradient as well, if mini_batch_indices = np.arange(n_train).
    lmbda: float
        Regularization strength. lmbda = 0 means having no regularization.

    Returns
    -----
    dw : array, shape [D]
        Gradient w.r.t. w.
    """
    # TODO
    N = len(mini_batch_indices)
    dw = np.dot(X[mini_batch_indices].T, (y[mini_batch_indices] - sigmoid(np.dot(X[mini_batch_indices],w))))

    dw = dw/(-N)

    np.insert(dw,0,0)
    dw = dw + lmbda*w
    return dw
```

Train the logistic regression model (nothing to do here)

```
In [85]: def logistic_regression(X, y, num_steps, learning_rate, mini_batch_size, lmbda, 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

    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 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, 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 [86]: 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
    y_pred = []
    for i in range(len(X)):
        y_pred.append((sigmoid(np.dot(w.T,X[i]))>=0.5).astype(np.int))
    return y_pred
```

Full batch gradient descent

```
In [87]: # Change this to True if you want to see loss values over iterations.
verbose = False

In [88]: 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,
                                         lmbda=0.1,
                                         verbose=verbose)

In [89]: 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 [90]: 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: 0.9357, f1_score: 0.9488

In [63]: 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 L(w)')
plt.legend()
plt.show()
```

