

Assignment 3

Supervised Learning - model training and evaluation

Kyle Bradbury

2026-02-04

Table of contents

Instructions

Instructions for all assignments can be found [here](#). Note: this assignment falls under collaboration Mode 2: Individual Assignment – Collaboration Permitted. Please refer to the syllabus for additional information. Please be sure to list the names of any students that you worked with on this assignment. Total points in the assignment add up to 90; an additional 10 points are allocated to professionalism and presentation quality.

Learning Objectives

This assignment will provide structured practice to help enable you to...

1. Understand the primary workflow in machine learning: (1) identifying a hypothesis function set of models, (2) determining a loss/cost/error/objective function to minimize, and (3) minimizing that function through gradient descent
2. Understand the inner workings of logistic regression and how linear models for classification can be developed.
3. Gain practice in implementing machine learning algorithms from the most basic building blocks to understand the math and programming behind them to achieve practical proficiency with the techniques
4. Implement batch gradient descent and become familiar with how that technique is used and its dependence on the choice of learning rate
5. Evaluate supervised learning algorithm performance through ROC curves and using cross validation
6. Apply regularization to linear models to improve model generalization performance

Exercise 1 - Classification using logistic regression: build it from the ground up

[60 points]

This exercise will walk you through the full life-cycle of a supervised machine learning classification problem. This classification problem consists of two features/predictors (e.g. petal width and petal length) and your goal is to predict one of two possible class labels (class 0 or class 1). You will build, train, and evaluate the performance of a logistic regression classifier on the data provided. Before you begin any modeling, you'll load and explore your data in Part I to familiarize yourself with it and check for any missing or erroneous data. Then, in Part II, you will explore a hypothesis set of functions to fit to the data: in this case, logistic regression models. In Part III, you will derive a cost function for the data (cross-entropy) and the specific gradient descent update equation. This will allow you to optimize your cost function to identify model parameters that minimize the cost function value of your model with respect to the training data. In Part IV, you will finally code up your model and the gradient descent training process from scratch. Through the process, you'll train your model and plot learning curves for your gradient descent implementation to verify the model learns to minimize your cost function. Lastly, in Part V, you will apply the model you designed and implemented and compare its performance to a KNN algorithm. **When complete, you will have accomplished learning objectives 1-5 above!**

A. Load, prepare, and plot your data



Note

Data for this exercise can be downloaded [here](#)

You are given some data for which you are tasked with constructing a classifier. The first step when facing any machine learning project: look at your data!

1.1 Load the data.

- In the data folder in the same directory of this notebook, you'll find the data in `A3_Q1_data.csv`. This file contains the binary class labels, y , and the features x_1 and x_2 .
- Divide your data into a training and testing set where the test set accounts for 30 percent of the data and the training set the remaining 70 percent.
- Plot the training data by class.
- Comment on the data: do the data appear separable? May logistic regression be a good choice for these data? Why or why not?

1.2 Do the data require any preprocessing due to missing values, scale differences (e.g. different ranges of values), etc.? If so, how did you handle these issues?

Next, we walk through our key steps for model fitting: choose a hypothesis set of models to train (in this case, logistic regression); identify a cost function to measure the model fit to our training data; optimize model parameters to minimize cost (in this case using gradient descent). Once we've completed model fitting, we will evaluate the performance of our model and compare performance to another approach (a KNN classifier).

B. Defining our hypothesis set: logistic regression models

Given that our data consists of two features, our logistic regression problem will be applied to a two-dimensional feature space. Recall that our logistic regression model is defined by the expression:

$$f(\mathbf{x}_i, \mathbf{w}) = \sigma(\mathbf{w}^\top \mathbf{x}_i)$$

where the sigmoid function is defined as $\sigma(x) = \frac{e^x}{1 + e^x} = \frac{1}{1 + e^{-x}}$. Also, since this is a two-dimensional problem, we define $\mathbf{w}^\top \mathbf{x}_i = w_0 x_{i,0} + w_1 x_{i,1} + w_2 x_{i,2}$ and here, $\mathbf{x}_i = [x_{i,0}, x_{i,1}, x_{i,2}]^\top$, and $x_{i,0} \triangleq 1$

Remember from class that we impose an interpretation on our logistic regression classifier output (the confidence score) to be the conditional probability that the target variable for a given sample y_i is from class “1”, given the observed features, \mathbf{x}_i . For one sample of training data, (y_i, \mathbf{x}_i) , this probability can be written as:

$$P(Y = 1 | X = \mathbf{x}_i) = f(\mathbf{x}_i, \mathbf{w}) = \sigma(\mathbf{w}^\top \mathbf{x}_i)$$

In the context of maximizing the likelihood of our parameters given the data, we define this to be the likelihood function $L(\mathbf{w} | y_i, \mathbf{x}_i)$, corresponding to one sample from the training dataset.

Aside: the careful reader will note that this use of the term likelihood differs from the definition in Bayes' Rule. In the context of training a logistic regression model, the likelihood we are interested in is this case answers the question: “if this model was true, how plausible are the data I observed?” In this context, the likelihood function is a function of our model parameters, \mathbf{w} . It’s our goal to use this to choose the parameters to maximize the likelihood function.

No output is required for this section - just read and use this information in the later sections.

C. Find the cost function to measure how well the model parameters, \mathbf{w} , fit the model to the training data.

1.3 What is the likelihood function that corresponds to all the N samples in our training dataset that we will wish to maximize? Unlike the likelihood function written above which gives the likelihood function for a *single training data pair* (y_i, \mathbf{x}_i) , this question asks for the likelihood function for the *entire training dataset* $\{(y_1, \mathbf{x}_1), (y_2, \mathbf{x}_2), \dots, (y_N, \mathbf{x}_N)\}$.

1.4 Since a logarithm is a monotonic function, maximizing the $f(x)$ is equivalent to maximizing $\ln[f(x)]$. Express the likelihood from the last question as a cost function of the model parameters. To do so define $C(\mathbf{w})$ as the negative of the logarithm of the likelihood. Express this cost as an average cost per sample (i.e. divide your final value by N), and use this quantity going forward as the cost function to optimize.

1.5 Calculate the gradient of the cost function with respect to the model parameters $\nabla_{\mathbf{w}} C(\mathbf{w})$. Express this in terms of the partial derivatives of the cost function with respect to each of the parameters, e.g. $\nabla_{\mathbf{w}} C(\mathbf{w}) = \left[\frac{\partial C}{\partial w_0}, \frac{\partial C}{\partial w_1}, \frac{\partial C}{\partial w_2} \right]$.

To simplify notation, please use $\mathbf{w}^\top \mathbf{x}$ instead of writing out $w_0 x_{i,0} + w_1 x_{i,1} + w_2 x_{i,2}$ when it appears each time (where $x_{i,0} = 1$ for all i). You are also welcome to use $\sigma()$ to represent the sigmoid function. Lastly, this will be a function the features, $x_{i,j}$ (with the first index in the subscript representing the observation and the second the feature; targets, y_i ; and the logistic regression model parameters, w_j).

1.6 Write out the gradient descent update equation for the model parameters \mathbf{w} . This should clearly express how to update each weight from one step in gradient descent $w_j^{(k)}$ to the next $w_j^{(k+1)}$. There should be one equation for each logistic regression model parameter (or you can represent it in vectorized form). Assume that η represents the learning rate.

D. Implement gradient descent and your logistic regression algorithm

1.7 Implement your logistic regression model.

- You are provided with a template, below, for a class with key methods to help with your model development. It is modeled on the Scikit-Learn classification class convention. For this, you only need to create a version of logistic regression for the case of two feature variables (i.e. two predictors).
- Create a method called `sigmoid` that calculates the sigmoid function
- Create a method called `cost` that computes the cost function $C(\mathbf{w})$ for a given dataset and corresponding class labels. This should be the **average cost** (make sure your total cost is divided by your number of samples in the dataset).

- Create a method called `gradient_descent` to run **one step** of gradient descent on your training data. We'll refer to this as "batch" gradient descent since it takes into account the gradient based on all our data at each iteration of the algorithm.
- Create a method called `fit` that fits the model to the data (i.e. adjusts the model parameters to minimize cost) using your `gradient_descent` method. In doing this we'll need to make some assumptions about the following:
 - *Weight initialization.* What should you initialize the model parameters to? For this, randomly initialize the weights to a different values between 0 and 1.
 - *Learning rate.* How slow/fast should the algorithm step towards the minimum? This you will vary in a later part of this problem.
 - *Stopping criteria.* When should the algorithm be finished searching for the optimum? There are two stopping criteria: small changes in the gradient descent step size and a maximum number of iterations (to make sure the algorithm doesn't run forever). The first determines whether there was a sufficiently small change in the gradient to merit stopping the training process; this is evaluated as whether the magnitude of the step that the gradient descent algorithm takes changes by less than some threshold (we will use 10^{-6}) between iterations. Since we have a weight vector, we can compute the change in the weight by evaluating the L_2 norm (Euclidean norm) of the change in the vector between iterations. From our gradient descent update equation we know that mathematically this is $\| -\eta \nabla_{\mathbf{w}} C(\mathbf{w}) \|$. The second criterion is met if a maximum number of iterations has been reached (5,000 in this case, to prevent infinite loops from poor choices of learning rates).
 - Design your approach so that at each step in the gradient descent algorithm you evaluate the cost function for both the training and the test data for each new value for the model weights. You should be able to plot cost vs gradient descent iteration for both the training and the test data. This will allow you to plot "learning curves" that can be informative for how the model training process is proceeding.
- Create a method called `predict_proba` that predicts confidence scores (these can be thresholded into the binary predictions of the `predict` method).
- Create a method called `predict` that makes predictions based on the trained model, selecting the most probable class, given the data, as the prediction, that is class that yields the larger $P(y|\mathbf{x})$.
- (Optional, but recommended) Create a method called `learning_curve` that retrieves the cost function values that correspond to each step from a previously run gradient descent operation.
- (Optional, but recommended) Create a method called `prepare_x` which appends a column of ones as the first feature of the dataset \mathbf{X} to account for the bias term ($x_{i,1} = 1$).

This structure is strongly encouraged; however, you're welcome to adjust this to your needs (adding helper methods, modifying parameters, etc.).

```
# Logistic regression class
class Logistic_regression:
    # Class constructor
    def __init__(self):
        self.w = None      # logistic regression weights
        self.saved_w = [] # Since this is a small problem, we can save the weights
                          # at each iteration of gradient descent to build our
                          # learning curves
    # returns nothing
    pass

    # Method for calculating the sigmoid function of w^T X for an input set of weights
    def sigmoid(self, X, w):
        # returns the value of the sigmoid
        pass

    # Cost function for an input set of weights
    def cost(self, X, y, w):
        # returns the average cross entropy cost
        pass

    # Update the weights in an iteration of gradient descent
    def gradient_descent(self, X, y, lr):
        # returns a scalar of the magnitude of the Euclidean norm
        # of the change in the weights during one gradient descent step
        pass

    # Fit the logistic regression model to the data through gradient descent
    def fit(self, X, y, w_init, lr, delta_thresh=1e-6, max_iter=5000, verbose=False):
        # Note the verbose flag enables you to print out the weights at each iteration
        # (optional - but may help with one of the questions)

        # returns nothing
        pass

    # Use the trained model to predict the confidence scores (prob of positive class in this
    def predict_proba(self, X):
        # returns the confidence score for the each sample
        pass
```

```

# Use the trained model to make binary predictions
def predict(self, X, thresh=0.5):
    # returns a binary prediction for each sample
    pass

# Stores the learning curves from saved weights from gradient descent
def learning_curve(self, X, y):
    # returns the value of the cost function from each step in gradient descent
    # from the last model fitting process
    pass

# Appends a column of ones as the first feature to account for the bias term
def prepare_x(self, X):
    # returns the X with a new feature of all ones (a column that is the new column 0)
    pass

```

1.8 Choose a learning rate and fit your model. Learning curves are a plot of metrics of model performance evaluated through the process of model training to provide insights about how model training is proceeding. Show the learning curves for the gradient descent process for learning rates of $\{10^{-0}, 10^{-2}, 10^{-4}\}$. For each learning rate plot the learning curves by plotting **both the training and test data average cost** as a function of each iteration of gradient descent. You should run the model fitting process until it completes (up to 5,000 iterations of gradient descent). All of the 6 resulting curves (train and test average cost for each learning rate) should be plotted on the **same set of axes** to enable direct comparison. *Note: make sure you're using average cost per sample, not the total cost.*

- Try running this process for a really big learning rate for this problem: 10^2 . Look at the weights that the fitting process generates over the first 50 iterations and how they change. Either print these first 50 iterations as console output or plot them. What happens? How does the output compare to that corresponding to a learning rate of 10^0 and why?
- What is the impact that the different values of learning have on the speed of the process and the results?
- Of the options explored, what learning rate do you prefer and why?
- Use your chosen learning rate for the remainder of this problem.

E. Evaluate your model performance through cross validation

1.9 Test the performance of your trained classifier using K-folds cross validation resampling technique. The scikit-learn package [StratifiedKFolds](#) may be helpful.

- Train your logistic regression model and a K-Nearest Neighbor classification model with $k = 7$ nearest neighbors.
- Using the trained models, make four plots: two for logistic regression and two for KNN. For each model have one plot showing the training data used for fitting the model, and the other showing the test data. On each plot, include the decision boundary resulting from your trained classifier.
- Produce a Receiver Operating Characteristic curve (ROC curve) that represents the performance from cross validated performance evaluation for each classifier (your logistic regression model and the KNN model, with $k = 7$ nearest neighbors). For the cross validation, use $k = 10$ folds.
 - Plot these curves on the same set of axes to compare them. You should not plot one curve for each fold of k-folds; instead, you should plot one ROC curve for Logistic Regression and one for KNN (each should incorporate all 10 folds of validation data).
 - On the ROC curve plot, also include the chance diagonal for reference (this represents the performance of the worst possible classifier). This is represented as a line from $(0, 0)$ to $(1, 1)$.
 - Calculate the Area Under the Curve for each model and include this measure in the legend of the ROC plot along with each model's name.
- Comment on the following:
 - What is the purpose of using cross validation for this problem?
 - How do the models compare in terms of performance (both ROC curves and decision boundaries) and which model (logistic regression or KNN) would you select to use on previously unseen data for this problem and why? How confident are you in your choice and what information would increase your confidence in your decision?

Exercise 2 - Digits classification

An exploration of regularization, imbalanced classes, ROC and PR curves

[30 points]

The goal of this exercise is to apply your supervised learning skills on a very different dataset: in this case, image data; MNIST: a collection of images of handwritten digits. Your goal is to train a classifier that is able to distinguish the number “3” from all possible numbers and to do so as accurately as possible. You will first explore your data (this should always be your starting point to gain domain knowledge about the problem.). Since the feature space in this problem is 784-dimensional, overfitting is possible. To avoid overfitting you will investigate the impact of regularization on generalization performance (test accuracy) and compare regularized and

unregularized logistic regression model test error against other classification techniques such as naive Bayes and random forests and draw conclusions about the best-performing model.

Start by loading your dataset from the [MNIST dataset](#) of handwritten digits, using the code provided below. MNIST has a training set of 60,000 examples, and a test set of 10,000 examples. The digits have been size-normalized and centered in a fixed-size image.

Your goal is to classify whether or not an example digit is a 3. Your binary classifier should predict $y = 1$ if the digit is a 3, and $y = 0$ otherwise. Create your dataset by transforming your labels into a binary format (3's are class 1, and all other digits are class 0).

2.1 Plot 10 examples of each class (i.e. class $y = 0$, which are not 3's and class $y = 1$ which are 3's), from the training dataset.

Note that the data are composed of samples of length 784. These represent 28 x 28 images, but have been reshaped for storage convenience. To plot digit examples, you'll need to reshape the data to be 28 x 28 (which can be done with numpy `reshape`).

2.2 How many examples are present in each class? Show a plot of samples by class (bar plot). What fraction of samples are positive? What issues might this cause?

2.3 Identify the value of the regularization parameter that optimizes model performance on out-of-sample data. Using a logistic regression classifier, apply lasso regularization and retrain the model and evaluate its performance on the test set over a range of values on the regularization coefficient. You can implement this using the [LogisticRegression](#) module and activating the 'l1' penalty; the parameter C is the inverse of the regularization strength. Vary the value of C logarithmically from 10^{-4} to 10^4 (and make your x-axes logarithmic in scale) and evaluate it at least 20 different values of C . As you vary the regularization coefficient, Plot the following four quantities (this should result in 4 separate plots)...

- The number of model parameters that are estimated to be nonzero (in the logistic regression model, one attribute is `coef_`, which gives you access to the model parameters for a trained model)
- The cross entropy loss (which can be evaluated with the Scikit Learn `log_loss` function)
- Area under the ROC curve (AUC)
- The F_1 -score (assuming a threshold of 0.5 on the predicted confidence scores, that is, scores above 0.5 are predicted as Class 1, otherwise Class 0). Scikit Learn also has a `f1_score` function which may be useful. -Which value of C seems best for this problem? Please select the closest power of 10. You will use this in the next part of this exercise.

2.4 Train and test a (1) logistic regression classifier with minimal regularization (using the Scikit Learn package, set `penalty='l1'`, $C=1e100$ to approximate this), (2) a logistic regression classifier with the best value of the regularization parameter from the last section, (3) a Gradient Boosting classifier, and (4) a Random Forest (RF) classifier (using default parameters for the RF classifier).

- Compare your classifiers' performance using ROC and Precision Recall (PR) curves. For the ROC curves, all your curves should be plotted on the same set of axes so that you can directly compare them. Please do the same with the PR curves.
- Plot the line that represents randomly guessing the class (50% of the time a “3”, 50% not a “3”). You SHOULD NOT actually create random guesses. Instead, you should think through the theory behind how ROC and PR curves work and plot the appropriate lines. It's a good practice to include these in ROC and PR curve plots as a reference point.
- For PR curves, an excellent resource on how to correctly plot them can be found [here](#) (ignore the section on “non-linear interpolation between two points”). This describes how a random classifier is represented in PR curves and demonstrates that it should provide a lower bound on performance.
- When training your logistic regression model, it's recommended that you use solver=“liblinear”; otherwise, your results may not converge.
- Describe the performance of the classifiers you compared. Did the regularization of the logistic regression model make much difference here? Which classifier you would select for application to unseen data.

2.5 Incorporating model uncertainty. So far, we've only looked at point estimates of our model performance. However, we know that limitations in either training data or test data such as small size, noise, or biases can add uncertainty to our model fit (based on the training data) and our performance evaluation metrics (based on the test data). Ideally, we could evaluate the uncertainty of both the training sample and the test sample. However, whenever we change the training set, we need to retrain the model which can be quite computationally expensive and often we want to know, specifically, how good a single trained model is. In that case, we focus on the uncertainty in the test set and address the question, “given our trained model, how uncertain is its measured performance?”

In this case, your goal is to compare each of the models you trained in 2.4 while including a measure of uncertainty by bootstrap sampling your test dataset. To do this you will first create 1,000 bootstrapped sample “test” datasets that will be used for each model (**your test data bootstrap samples must be identical for each model or the comparison will be meaningless**). Evaluate the [average precision](#) (which is the area under the precision recall curve) for each model for each of the 1,000 bootstrap samples. Create a box and whisker plot where each model is represented with its mean and the 2.5th and 97.5th percentiles are represented by the “whiskers” and the box represents the 25th and 75th percentiles.

Comment on your results and their implications for the model comparison.

```
# Load the MNIST Data
from sklearn.datasets import fetch_openml
from sklearn.model_selection import train_test_split
import numpy as np
import matplotlib.pyplot as plt
import pickle
```

```

# Set this to True to download the data for the first time and False after the first time
#   so that you just load the data locally instead
download_data = True

if download_data:
    # Load data from https://www.openml.org/d/554
    X, y = fetch_openml('mnist_784', return_X_y=True, as_frame=False)

    # Adjust the labels to be '1' if y==3, and '0' otherwise
    y[y!='3'] = 0
    y[y=='3'] = 1
    y = y.astype('int')

    # Divide the data into a training and test split
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=1/7, random_state=88)

    file = open('tmpdata', 'wb')
    pickle.dump((X_train, X_test, y_train, y_test), file)
    file.close()
else:
    file = open('tmpdata', 'rb')
    X_train, X_test, y_train, y_test = pickle.load(file)
    file.close()

```