

Softmax exercise

Complete and hand in this completed worksheet (including its outputs and any supporting code outside of the worksheet) with your assignment submission.

In this exercise, you will:

- implement a fully-vectorized loss function for the Softmax classifier
- implement the fully-vectorized expression for its analytic gradient
- check your implementation with numerical gradient
- use a validation set to tune the learning rate and regularization strength
- optimize the loss function with SGD
- visualize the final learned weights

Acknowledgement: This exercise is adapted from Stanford CS231n.

```
In [1]: import random
import numpy as np
from data_utils import load_CIFAR10
import matplotlib.pyplot as plt

plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

# for http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2

In [2]: def rel_error(out, correct_out):
    return np.abs(out - correct_out) / (abs(out) + abs(correct_out))

In [3]: def get_CIFAR10_data(num_training=49000, num_validation=1000, num_test=1000, num_dev=500):
    """
    Load the raw CIFAR-10 data
    """
    X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

    # subsample the data
    num_train = range(num_training, num_training + num_validation)
    X_val = X_train[num_train]
    y_val = y_train[num_train]
    mask = range(num_training)
    X_train = X_train[mask]
    y_train = y_train[mask]
    mask = range(num_test)
    X_test = X_test[mask]
    y_test = y_test[mask]

    # We will also make a development set, which is a small subset of
    # the training set
    mask = np.random.choice(num_training, num_dev, replace=False)
    X_dev = X_train[mask]
    y_dev = y_train[mask]

    # Preprocessing: reshape the image data into rows
    X_train = np.reshape(X_train, (X_train.shape[0], -1))
    X_val = np.reshape(X_val, (X_val.shape[0], -1))
    X_test = np.reshape(X_test, (X_test.shape[0], -1))
    X_dev = np.reshape(X_dev, (X_dev.shape[0], -1))

    # Normalization: subtract the mean image
    mean_image = np.mean(X_train, axis = 0)
    X_train -= mean_image
    X_val -= mean_image
    X_test -= mean_image
    X_dev -= mean_image

    # add a mean dimension at the beginning of the image
    X_train = np.hstack([X_train, np.ones((X_train.shape[0], 1))])
    X_val = np.hstack([X_val, np.ones((X_val.shape[0], 1))])
    X_test = np.hstack([X_test, np.ones((X_test.shape[0], 1))])
    X_dev = np.hstack([X_dev, np.ones((X_dev.shape[0], 1))])

    return X_train, y_train, X_val, y_val, X_test, y_test, X_dev, y_dev

# Invoke the above function to get our data.
X_train, y_train, X_val, y_val, X_test, y_test, X_dev, y_dev = get_CIFAR10_data()
print('Train data shape: ', X_train.shape)
print('Train labels shape: ', y_train.shape)
print('Validation data shape: ', X_val.shape)
print('Validation labels shape: ', y_val.shape)
print('Test data shape: ', X_test.shape)
print('Test labels shape: ', y_test.shape)
print('dev data shape: ', X_dev.shape)
print('dev labels shape: ', y_dev.shape)

Train data shape: (49000, 3073)
Train labels shape: (49000,)
Validation data shape: (1000, 3073)
Validation labels shape: (1000,)
Test data shape: (1000, 3073)
Test labels shape: (1000,)
dev data shape: (500, 3073)
dev labels shape: (500,)

In [4]: # Create one-hot vectors for label
num_class = 10
y_train_oh = np.zeros((y_train.shape[0], 10))
y_train_oh[np.arange(y_train.shape[0], y_train) == 1] = 1
y_val_oh = np.zeros((y_val.shape[0], 10))
y_val_oh[np.arange(y_val.shape[0], y_val) == 1] = 1
y_test_oh = np.zeros((y_test.shape[0], 10))
y_test_oh[np.arange(y_test.shape[0], y_test) == 1] = 1
y_dev_oh = np.zeros((y_dev.shape[0], 10))
y_dev_oh[np.arange(y_dev.shape[0], y_dev) == 1] = 1
```

Regression as classifier

The most simple and straightforward approach to learn a classifier is to map the input data (raw image values) to class label (one-hot vector). The loss function is defined as follows:

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n \|\mathbf{W}^T \mathbf{x}_i - \mathbf{y}_i\|_F^2 \quad (1)$$

Where:

- $\mathbf{W} \in \mathbb{R}^{(d+1) \times C}$: Classifier weight
- $\mathbf{X} \in \mathbb{R}^{n \times (d+1)}$: Dataset
- $\mathbf{y} \in \mathbb{R}^{n \times C}$: Class label (one-hot vector)

Optimization

Given the loss function (1), the next problem is how to solve the weight \mathbf{W} . We now discuss 2 approaches:

- Random search
- Closed-form solution

Random search

```
In [5]: bestloss = float('inf')
for num in range(100):
    W = np.random.randn(3073, 10) * 0.0001
    loss = np.linalg.norm(X_dev.dot(W) - y_dev_oh)
    if (loss < bestloss):
        bestloss = loss
        bestW = W

print('In attempt %d the loss was %f, best %f' % (num, loss, bestloss))

In attempt 0 the loss was 34.584989, best 34.584989
In attempt 1 the loss was 32.742228, best 32.742228
In attempt 2 the loss was 33.160128, best 32.742228
In attempt 3 the loss was 33.415307, best 32.742228
In attempt 4 the loss was 30.138346, best 30.138346
In attempt 5 the loss was 34.948674, best 30.138346
In attempt 6 the loss was 32.965973, best 30.138346
In attempt 7 the loss was 34.918137, best 30.138346
In attempt 8 the loss was 32.294784, best 30.138346
In attempt 9 the loss was 32.126684, best 30.138346
In attempt 10 the loss was 31.166088, best 30.138346
In attempt 11 the loss was 34.075411, best 30.138346
In attempt 12 the loss was 31.195480, best 30.138346
In attempt 13 the loss was 32.638493, best 30.138346
In attempt 14 the loss was 37.145719, best 30.138346
In attempt 15 the loss was 35.620309, best 30.138346
In attempt 16 the loss was 32.539644, best 30.138346
In attempt 17 the loss was 32.865098, best 30.138346
In attempt 18 the loss was 32.886827, best 30.138346
In attempt 19 the loss was 32.522084, best 30.138346
In attempt 20 the loss was 32.957806, best 30.138346
In attempt 21 the loss was 33.385979, best 30.138346
In attempt 22 the loss was 32.807145, best 30.138346
In attempt 23 the loss was 34.23534, best 30.138346
In attempt 24 the loss was 32.319335, best 30.138346
In attempt 25 the loss was 33.779677, best 30.138346
In attempt 26 the loss was 32.286827, best 30.138346
In attempt 27 the loss was 32.89435, best 30.138346
In attempt 28 the loss was 32.868451, best 30.138346
In attempt 29 the loss was 32.105325, best 30.138346
In attempt 30 the loss was 33.924781, best 30.138346
In attempt 31 the loss was 32.749890, best 30.138346
In attempt 32 the loss was 37.061445, best 30.138346
In attempt 33 the loss was 30.424317, best 30.138346
In attempt 34 the loss was 34.672037, best 30.138346
In attempt 35 the loss was 33.579730, best 30.138346
In attempt 36 the loss was 33.456837, best 30.138346
In attempt 37 the loss was 33.319312, best 30.138346
In attempt 38 the loss was 32.360866, best 30.138346
In attempt 39 the loss was 36.350599, best 30.138346
In attempt 40 the loss was 34.209386, best 30.138346
In attempt 41 the loss was 34.008904, best 30.138346
In attempt 42 the loss was 33.750351, best 30.138346
In attempt 43 the loss was 34.501746, best 30.138346
In attempt 44 the loss was 33.940031, best 30.138346
In attempt 45 the loss was 34.247738, best 30.138346
In attempt 46 the loss was 32.590347, best 30.138346
In attempt 47 the loss was 35.363975, best 30.138346
In attempt 48 the loss was 33.385979, best 30.138346
In attempt 49 the loss was 33.030771, best 30.138346
In attempt 50 the loss was 37.049055, best 30.138346
In attempt 51 the loss was 32.865468, best 30.138346
In attempt 52 the loss was 31.628220, best 30.138346
In attempt 53 the loss was 33.244617, best 30.138346
In attempt 54 the loss was 34.679795, best 30.138346
In attempt 55 the loss was 34.363134, best 30.138346
In attempt 56 the loss was 35.917236, best 30.138346
In attempt 57 the loss was 32.522759, best 30.138346
In attempt 58 the loss was 32.768395, best 30.138346
In attempt 59 the loss was 35.092037, best 30.138346
In attempt 60 the loss was 33.347134, best 30.138346
In attempt 61 the loss was 33.450936, best 30.138346
In attempt 62 the loss was 32.931332, best 30.138346
In attempt 63 the loss was 32.640132, best 30.138346
In attempt 64 the loss was 32.810023, best 30.138346
In attempt 65 the loss was 34.843407, best 30.138346
In attempt 66 the loss was 32.690375, best 30.138346
In attempt 67 the loss was 36.720599, best 30.138346
In attempt 68 the loss was 33.134931, best 30.138346
In attempt 69 the loss was 32.255037, best 30.138346
In attempt 70 the loss was 33.385979, best 30.138346
In attempt 71 the loss was 32.591734, best 30.138346
In attempt 72 the loss was 33.397416, best 30.138346
In attempt 73 the loss was 34.428359, best 30.138346
In attempt 74 the loss was 31.865540, best 30.138346
In attempt 75 the loss was 34.075230, best 30.138346
In attempt 76 the loss was 32.766628, best 30.138346
In attempt 77 the loss was 32.932035, best 30.138346
In attempt 78 the loss was 33.501741, best 30.138346
In attempt 79 the loss was 33.501285, best 30.138346
In attempt 80 the loss was 32.429037, best 30.138346
In attempt 81 the loss was 33.193304, best 30.138346
In attempt 82 the loss was 33.659201, best 30.138346
In attempt 83 the loss was 35.416439, best 30.138346
In attempt 84 the loss was 32.320320, best 30.138346
In attempt 85 the loss was 33.297539, best 30.138346
In attempt 86 the loss was 34.237956, best 30.138346
In attempt 87 the loss was 31.255983, best 30.138346
In attempt 88 the loss was 32.931332, best 30.138346
In attempt 89 the loss was 34.129418, best 30.138346
In attempt 90 the loss was 33.061938, best 30.138346
In attempt 91 the loss was 33.802207, best 30.138346
In attempt 92 the loss was 32.819316, best 30.138346
In attempt 93 the loss was 34.712956, best 30.138346
In attempt 94 the loss was 33.498338, best 30.138346
In attempt 95 the loss was 37.408749, best 30.138346
In attempt 96 the loss was 32.674211, best 30.138346
In attempt 97 the loss was 32.395267, best 30.138346
In attempt 98 the loss was 33.646270, best 30.138346
In attempt 99 the loss was 32.294125, best 30.138346

In [6]: # How bestW perform:
print('Accuracy on train set: ', np.sum(np.argmax(np.abs(1 - X_dev.dot(W)), axis=1) == y_dev).astype(np.float32)/y_train.shape[0])
print('Accuracy on test set: ', np.sum(np.argmax(np.abs(1 - X_test.dot(W)), axis=1) == y_test).astype(np.float32)/y_test.shape[0])

Accuracy on train set: 9.4
Accuracy on test set: 8.2
```

You can clearly see that the performance is very low, almost at the random level.

Closed-form solution

The closed-form solution is achieved by:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \frac{2}{n} \mathbf{X}^T (\mathbf{X}\mathbf{W} - \mathbf{y}) = 0$$
$$\Leftrightarrow \mathbf{W}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

```
In [7]: #####
# TODO:
# Implement the closed-form solution of the weight W.
#####

A = np.matmul(X_train.T, X_train)
B = np.matmul(X_train.T, y_train_oh)
W = np.matmul(np.linalg.pinv(A), B)

##### END OF YOUR CODE #####

In [8]: # Check accuracy:
print('Train set accuracy: ', np.sum(np.argmax(np.abs(1 - X_train.dot(W)), axis=1) == y_train).astype(np.float32)/y_train.shape[0])
print('Test set accuracy: ', np.sum(np.argmax(np.abs(1 - X_test.dot(W)), axis=1) == y_test).astype(np.float32)/y_test.shape[0])

Train set accuracy: 51.163265306122454
Test set accuracy: 32.195395999999996

Now, you can see that the performance is much better.
```

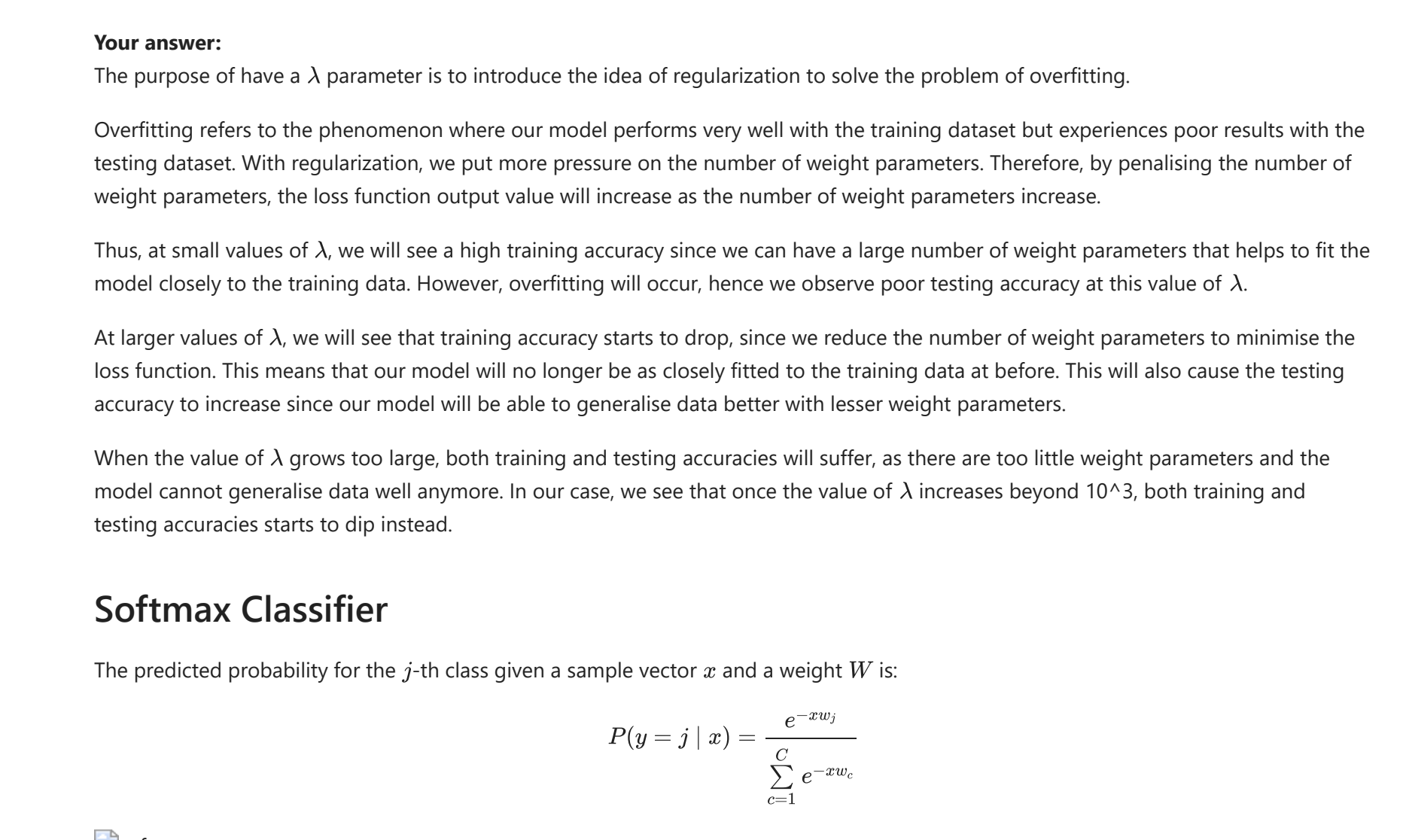
Regularization

A simple way to improve performance is to include the L2-regularization penalty.

$$\mathcal{L} = \frac{1}{n} \|\mathbf{X}\mathbf{W} - \mathbf{y}\|_F^2 + \lambda \|\mathbf{W}\|_F^2 \quad (2)$$

The closed-form solution now is:

$$\Leftrightarrow \mathbf{W}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$



Question: Try to explain why the performances on the training and test set have such behaviors as we change the value of λ .

Your answer:

The purpose of having a λ parameter is to introduce the idea of regularization to solve the problem of overfitting.

Overfitting refers to the phenomenon where our model performs very well with the training dataset but experiences poor results with the testing dataset. With regularization, we put more pressure on the number of weight parameters. Therefore, by penalising the number of weight parameters, the loss function output value will increase as the number of weight parameters increase.

Thus, at small values of λ , we will see a high training accuracy since we can have a large number of weight parameters that helps to fit the model closely to the training data. However, overfitting will occur, hence we observe poor testing accuracy at this value of λ .

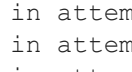
At larger values of λ , we will see that training accuracy starts to drop, since we reduce the number of weight parameters to minimise the loss function. This means that our model will no longer be able to generalise data better to the training data at before. This will also cause the testing accuracy to increase since our model will be able to classify data better with lesser weight parameters.

When the value of λ grows too large, both training and testing accuracies will suffer, as there are too little weight parameters and the model cannot generalise data well anymore. In our case, we see that once the value of λ increases beyond 10^{-1} , both training and testing accuracies starts to dip instead.

Softmax Classifier

The predicted probability for the j -th class given a sample vector \mathbf{x} and a weight \mathbf{W} is:

$$P(y = j | \mathbf{x}) = \frac{e^{-\mathbf{w}_j^T \mathbf{x}}}{\sum_{c=1}^C e^{-\mathbf{w}_c^T \mathbf{x}}}$$



Your code for this section will all be written inside `classifiers/softmax.py`.

```
In [11]: # First implement the naive softmax loss function with nested loops.
# Open the file classifiers/softmax.py to implement the
# softmax_loss_naive function.

from classifiers.softmax import softmax_loss_naive
import time

# Generate a random softmax weight matrix and use it to compute the loss.
np.random.randn(3073, 10) * 0.0001
loss, grad = softmax_loss_naive(W, X_dev, y_dev, 0.0)

# As a rough sanity check, our loss should be something close to -log(0.1).
print('sanity check: %f' % (-np.log(0.1)))

loss: 2.337964
sanity check: 2.302585

Question: Why do we expect our loss to be close to -log(0.1)? Explain briefly.**

Your answer:
Since we are calculating our loss based on random weights (i.e., we haven't started the "learning" process yet), we expect that the initial loss has to be close to -log(0.1) because initially all the classes are equally likely to be chosen. Since CIFAR-10 consists of samples which belong to one of ten classes, the probability of the correct class will be 1/10 = 0.1. The softmax loss is the negative log probability of the correct class, therefore it is -log(0.1).

Optimization

Random search

In [12]: bestloss = float('inf')
for num in range(100):
    W = np.random.randn(3073, 10) * 0.0001
    loss, _ = softmax_loss_naive(W, X_dev, y_dev, 0.0)
    if (loss < bestloss):
        bestloss = loss
        bestW = W

print('In attempt %d the loss was %f, best %f' % (num, loss, bestloss))

In attempt 0 the loss was 2.358189, best 2.358189
In attempt 1 the loss was 2.358189, best 2.358189
In attempt 2 the loss was 2.310783, best 2.310783
In attempt 3 the loss was 2.353139, best 2.310783
In attempt 4 the loss was 2.392929, best 2.310783
In attempt 5 the loss was 2.362636, best 2.310783
In attempt 6 the loss was 2.365318, best 2.310783
In attempt 7 the loss was 2.416753, best 2.310783
In attempt 8 the loss was 2.364832, best 2.310783
In attempt 9 the loss was 2.395973, best 2.310783
In attempt 10 the loss was 2.353233, best 2.310783
In attempt 11 the loss was 2.400715, best 2.310783
In attempt 12 the loss was 2.346051, best 2.310783
In attempt 13 the loss was 2.363616, best 2.310783
In attempt 14 the loss was 2.359163, best 2.310783
In attempt 15 the loss was 2.339410, best 2.310783
In attempt 16 the loss was 2.406397, best 2.310783
In attempt 17 the loss was 2.363616, best 2.310783
In attempt 18 the loss was 2.359163, best 2.310783
In attempt 19 the loss was 2.365997, best 2.310783
In attempt 20 the loss was 2.388597, best 2.310783
In attempt 21 the loss was 2.392862, best 2.310783
In attempt 22 the loss was 2.313808, best 2.310783
In attempt 23 the loss was 2.423105, best 2.310783
In attempt 24 the loss was 2.358605, best 2.310783
In attempt 25 the loss was 2.364470, best 2.310783
In attempt 26 the loss was 2.374066, best 2.310783
In attempt 27 the loss was 2.375507, best 2.310783
In attempt 28 the loss was 2.359737, best 2.310783
In attempt 29 the loss was 2.388236, best 2.310783
In attempt 30 the loss was 2.356644, best 2.310783
In attempt 31 the loss was 2.417594, best 2.310783
In attempt 32 the loss was 2.373376, best 2.310783
In attempt 33 the loss was 2.387114, best 2.310783
In attempt 34 the loss was 2.362240, best 2.310783
In attempt 35 the loss was 2.317753, best 2.310783
In attempt 36 the loss was 2.365536, best 2.310783
In attempt 37 the loss was 2.367833, best 2.310783
In attempt 38 the loss was 2.378711, best 2.310783
In attempt 39 the loss was 2.355449, best 2.310783
In attempt 40 the loss was 2.390439, best 2.310783
In attempt 41 the loss was 2.357839, best 2.310783
In attempt 42 the loss was 2.324688, best 2.310783
In attempt 43 the loss was 2.326473, best 2.310783
In attempt 44 the loss was 2.364281, best 2.310783
In attempt 45 the loss was 2.394717, best 2.310783
In attempt 46 the loss was 2.323231, best 2.310783
In attempt 47 the loss was 2.400970, best 2.310783
In attempt 48 the loss was 2.359326, best 2.310783
In attempt 49 the loss was 2.276516, best 2.272516
In attempt 50 the loss was 2.31057, best 2.272516
In attempt 51 the loss was 2.365610, best 2.272516
In attempt 52 the loss was 2.365215, best 2.272516
In attempt 53 the loss was 2.381307, best 2.272516
In attempt 54 the loss was 2.401309, best 2.272516
In attempt 55 the loss was 2.387564, best 2.272516
In attempt 56 the loss was 2.322308, best 2.272516
In attempt 57 the loss was 2.366056, best 2.272516
In attempt 58 the loss was 2.310527, best 2.272516
In attempt 59 the loss was 2.383944, best 2.272516
In attempt 60 the loss was 2.347397, best 2.272516
In attempt 61 the loss was 2.317753, best 2.272516
In attempt 62 the loss was 2.330848, best 2.272516
In attempt 63 the loss was 2.328594, best 2.272516
In attempt 64 the loss was 2.414387, best 2.272516
In attempt 65 the loss was 2.307981, best 2.272516
In attempt 66 the loss was 2.383025, best 2.272516
In attempt 67 the loss was 2.378787, best 2.272516
In attempt 68 the loss was 2.366056, best 2.272516
In attempt 69 the loss was 2.377702, best 2.272516
In attempt 70 the loss was 2.391061, best 2.272516
In attempt 71 the loss was 2.349717, best 2.272516
In attempt 72 the loss was 2.351057, best 2.272516
In attempt 73 the loss was 2.345127, best 2.272516
In attempt 74 the loss was 2.384438, best 2.272516
In attempt 75 the loss was 2.400660, best 2.272516
In attempt 76 the loss was 2.330442, best 2.272516
In attempt 77 the loss was 2.344338, best 2.272516
In attempt 78 the loss was 2.330186, best 2.272516
In attempt 79 the loss was 2.331407, best 2.272516
In attempt 80 the loss was 2.365672, best 2.272516
In attempt 81 the loss was 2.377288, best 2.272516
In attempt 82 the loss was 2.343781, best 2.272516
In attempt 83 the loss was 2.346305, best 2.272516
In attempt 84 the loss was 2.425035, best 2.272516
In attempt 85 the loss was 2.362370, best 2.272516
In attempt 86 the loss was 2.423469, best 2.272516
In attempt 87 the loss was 2.342048, best 2.272516
In attempt 88 the loss was 2.414161, best 2.272516
In attempt 89 the loss was 2.388981, best 2.272516
In attempt 90 the loss was 2.313597, best 2.272516
In attempt 91 the loss was 2.385751, best 2.272516
In attempt 92 the loss was 2.341823, best 2.272516
In attempt 93 the loss was 2.425128, best 2.272516
In attempt 94 the loss was 2.384035, best 2.272516
In attempt 95 the loss was 2.351891, best 2.272516
In attempt 96 the loss was 2.343007, best 2.272516
In attempt 97 the loss was 2.326912, best 2.272516
In attempt 98 the loss was 2.345184, best 2.272516
In attempt 99 the loss was 2.345184, best 2.272516

In [13]: # How bestW perform on trainset
scores = X_train.dot(bestW)
y_pred = np.argmax(scores, axis=1)
print('Accuracy on train set %f' % np.mean(y_pred == y_train))

# evaluate performance of test set
scores = X_test.dot(bestW)
y_pred = np.argmax(scores, axis=1)
print('Accuracy on test set %f' % np.mean(y_pred == y_test))

Accuracy on train set: 0.132959
Accuracy on test set: 0.158000

Compare the performance when using random search with regression classifier and softmax classifier. You can see how much useful the softmax classifiers.
```

Stochastic Gradient descent

Even though it is possible to achieve closed-form solution with softmax classifier, it would be more complicated. In fact, we could achieve very good results with gradient descent approach. Additionally, in case of very large dataset, it is impossible to load the whole dataset into the memory. Gradient descent can help to optimize the loss function in batch.

$$\mathbf{W}^{t+1} = \mathbf{W}^t - \alpha \frac{\partial \mathcal{L}(\mathbf{x}; \mathbf{W})}{\partial \mathbf{W}^t}$$

Where α is the learning rate, \mathcal{L} is a loss function, and \mathbf{x} is a batch of training dataset.

```
In [14]: # Complete the implementation of softmax_loss_naive and implement a (naive)
# version of the gradient that uses nested loops.
loss, grad = softmax_loss_naive(W, X_dev, y_dev, 0.0)

# Use numeric gradient checking as a debugging tool.
# The numeric gradient should be close to the analytic gradient.
def grad_check_sparse(f, W, grad, 10):
    f = lambda w: softmax_loss_naive(W, X_dev, y_dev, 0.0) [0]
    grad_numerical = grad_check_sparse(f, W, grad, 10)

# gradient check with regularization
loss, grad = softmax_loss_naive(W, X_dev, y_dev, 1e2)
f = lambda w: softmax_loss_naive(W, X_dev, y_dev, 1e2) [0]
grad_numerical = grad_check_sparse(f, W, grad, 10)

numerical: -1.268820 analytic: -1.268820, relative error: 2.673774e-09
numerical: -1.498243 analytic: -1.498243, relative error: 4.54630e-09
numerical: -1.091292 analytic: -1.091292, relative error: 2.083992e-08
numerical: -1.282992 analytic: -1.282992, relative error: 2.987137e-08
numerical: -2.215568 analytic: -2.215568, relative error: 9.234415e-09
numerical: -2.071038 analytic: -2.071038, relative error: 1.66383e-08
numerical: -2.138951 analytic: -2.138951, relative error: 2.953210e-08
numerical: -6.087443 analytic: -6.087443, relative error: 2.667196e-09
numerical: -1.571848 analytic: -1.571848, relative error: 9.234415e-09
numerical: -2.643592 analytic: -2.643592, relative error: 8.611738e-09
numerical: 3.104153 analytic: 3.104153, relative error: 1.037942e-08
numerical: -3.464932 analytic: -3.464933, relative error: 2.571873e-08
numerical: -1.950068 analytic: -1.950069, relative error: 4.475809e-09
numerical: -3.118712 analytic: -3.118712, relative error: 3.211463e-08
numerical: -0.523393 analytic: -0.523393, relative error: 1.889307e-08
numerical: -0.842392 analytic: -0.842392, relative error: 1.160366e-07
numerical: 0.510741 analytic: 0.510743, relative error: 1.861022e-08
numerical: 2.391333 analytic: 2.391333, relative error: 3.636120e-08
numerical: 1.400442 analytic: 1.400442, relative error: 2.312027e-08
numerical: -0.925175 analytic: -0.925175, relative error: 4.953093e-08

In [15]: # Now that we have a naive implementation of the softmax loss function and its gradient,
# implement a vectorized version in softmax_loss_vectorized.
# The two versions should compute the same results, but the vectorized version should be
# much faster.
tic = time.time()
loss_naive, grad_naive = softmax_loss_naive(W, X_dev, y_dev, 0.00001)
toc = time.time()
print('naive loss: %e computed in %fs' % (loss_naive, toc - tic))

from classifiers.softmax import softmax_loss_vectorized
tic = time.time()
loss_vectorized, grad_vectorized = softmax_loss_vectorized(W, X_dev, y_dev, 0.00001)
toc = time.time()
print('vectorized loss: %e computed in %fs' % (loss_vectorized, toc - tic))

# We use the Frobenius norm to compare the two versions
of the gradient.
loss_difference = np.linalg.norm(grad_naive - grad_vectorized, ord='fro')
print('loss difference: %f' % np.abs(loss_naive - loss_vector
```