

softmax

September 29, 2021

1 Softmax exercise

Complete and hand in this completed worksheet (including its outputs and any supporting code outside of the worksheet) with your assignment submission. For more details see the [assignments page](#) on the course website.

This exercise is analogous to the SVM 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

```
[1]: from __future__ import print_function
import random
import numpy as np
from cs682.data_utils import load_CIFAR10
import matplotlib.pyplot as plt

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

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/
# → autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2

[2]: def get_CIFAR10_data(num_training=49000, num_validation=1000, num_test=1000,
    → num_dev=500):
    """
    Load the CIFAR-10 dataset from disk and perform preprocessing to prepare
    it for the linear classifier. These are the same steps as we used for the
    SVM, but condensed to a single function.
```

```

"""
# Load the raw CIFAR-10 data
cifar10_dir = 'cs682/datasets/cifar-10-batches-py'

X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

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

# Normalize the data: 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 bias dimension and transform into columns
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

# Cleaning up variables to prevent loading data multiple times (which may cause
↪memory issue)
try:
    del X_train, y_train
    del X_test, y_test

```

```

    print('Clear previously loaded data.')
except:
    pass

# 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,)

```

1.1 Softmax Classifier

Your code for this section will all be written inside `cs682/classifiers/softmax.py`.

```

[3]: # First implement the naive softmax loss function with nested loops.
# Open the file cs682/classifiers/softmax.py and implement the
# softmax_loss_naive function.

from cs682.classifiers.softmax import softmax_loss_naive
import time

# Generate a random softmax weight matrix and use it to compute the loss.
W = 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('loss: %f' % loss)
print('sanity check: %f' % (-np.log(0.1)))

```

```

loss: 2.373042
sanity check: 2.302585

```

1.2 Inline Question 1:

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

Your answer: *Fill this in*

This is because we initiate our weights randomly, and we have ten different classes. Since the weights are random, it is most likely that every class would have the same score before training. This means the probabilities would be $1/\text{num_classes}$ before train, which in our case is $1/10=0.1$.

```
[4]: # 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)

# As we did for the SVM, use numeric gradient checking as a debugging tool.
# The numeric gradient should be close to the analytic gradient.
from cs682.gradient_check import grad_check_sparse
f = lambda w: softmax_loss_naive(w, X_dev, y_dev, 0.0)[0]
grad_numerical = grad_check_sparse(f, W, grad, 10)

# similar to SVM case, do another gradient check with regularization
loss, grad = softmax_loss_naive(W, X_dev, y_dev, 5e1)
f = lambda w: softmax_loss_naive(w, X_dev, y_dev, 5e1)[0]
grad_numerical = grad_check_sparse(f, W, grad, 10)
```

```
numerical: 0.287101 analytic: 0.287101, relative error: 2.864537e-07
numerical: 0.371524 analytic: 0.371524, relative error: 1.641621e-07
numerical: 0.664583 analytic: 0.664583, relative error: 3.825252e-08
numerical: -0.670274 analytic: -0.670274, relative error: 8.670363e-09
numerical: 5.662179 analytic: 5.662178, relative error: 2.007403e-08
numerical: -1.291537 analytic: -1.291537, relative error: 2.540536e-08
numerical: -1.068029 analytic: -1.068029, relative error: 9.433747e-09
numerical: 0.889227 analytic: 0.889226, relative error: 9.961711e-08
numerical: -0.699912 analytic: -0.699913, relative error: 5.477711e-08
numerical: 0.102961 analytic: 0.102961, relative error: 9.655173e-08
numerical: 0.414729 analytic: 0.414729, relative error: 7.764006e-08
numerical: -0.363680 analytic: -0.363680, relative error: 3.198507e-08
numerical: 0.465865 analytic: 0.465865, relative error: 8.059718e-08
numerical: -1.091263 analytic: -1.091263, relative error: 9.133470e-08
numerical: -1.925673 analytic: -1.925674, relative error: 1.588992e-08
numerical: -3.747741 analytic: -3.747741, relative error: 1.872167e-08
numerical: 0.269944 analytic: 0.269944, relative error: 1.539578e-07
numerical: 0.448560 analytic: 0.448560, relative error: 2.354892e-08
numerical: 1.426694 analytic: 1.426694, relative error: 4.527782e-08
numerical: -3.059476 analytic: -3.059476, relative error: 2.512482e-09
```

```
[5]: # 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, 0.000005)
toc = time.time()
print('naive loss: %e computed in %fs' % (loss_naive, toc - tic))

from cs682.classifiers.softmax import softmax_loss_vectorized
tic = time.time()
loss_vectorized, grad_vectorized = softmax_loss_vectorized(W, X_dev, y_dev, 0.
→ 000005)
toc = time.time()
print('vectorized loss: %e computed in %fs' % (loss_vectorized, toc - tic))

# As we did for the SVM, we use the Frobenius norm to compare the two versions
# of the gradient.
grad_difference = np.linalg.norm(grad_naive - grad_vectorized, ord='fro')
print('Loss difference: %f' % np.abs(loss_naive - loss_vectorized))
print('Gradient difference: %f' % grad_difference)

```

```

naive loss: 2.373042e+00 computed in 0.093381s
vectorized loss: 2.373042e+00 computed in 0.004147s
Loss difference: 0.000000
Gradient difference: 0.000000

```

```

[6]: # Use the validation set to tune hyperparameters (regularization strength and
# learning rate). You should experiment with different ranges for the learning
# rates and regularization strengths; if you are careful you should be able to
# get a classification accuracy of over 0.35 on the validation set.
from cs682.classifiers import Softmax
results = {}
best_val = -1
best_softmax = None
learning_rates = [8e-8, 9e-8, 1e-7]
regularization_strengths = [9e3, 1e4, 1.5e4]

#####
# TODO:
# Use the validation set to set the learning rate and regularization strength. #
# This should be identical to the validation that you did for the SVM; save   #
# the best trained softmax classifier in best_softmax.                         #
#####
# Your code
for lr in learning_rates:
    for reg in regularization_strengths:
        softmax = Softmax()

```

```

softmax.train(X_train, y_train, lr, reg, num_iters=3000)
y_train_pred = softmax.predict(X_train)
train_accuracy = np.mean(y_train == y_train_pred)
y_val_pred = softmax.predict(X_val)
val_accuracy = np.mean(y_val == y_val_pred)
results[(lr, reg)] = (train_accuracy, val_accuracy)

if val_accuracy > best_val:
    best_softmax = softmax
    best_val = val_accuracy

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

# Print out results.
for lr, reg in sorted(results):
    train_accuracy, val_accuracy = results[(lr, reg)]
    print('lr %e reg %e train accuracy: %f val accuracy: %f' % (
        lr, reg, train_accuracy, val_accuracy))

print('best validation accuracy achieved during cross-validation: %f' %
      best_val)

```

```

lr 8.000000e-08 reg 9.000000e+03 train accuracy: 0.359837 val accuracy: 0.370000
lr 8.000000e-08 reg 1.000000e+04 train accuracy: 0.357898 val accuracy: 0.368000
lr 8.000000e-08 reg 1.500000e+04 train accuracy: 0.346408 val accuracy: 0.363000
lr 9.000000e-08 reg 9.000000e+03 train accuracy: 0.357245 val accuracy: 0.377000
lr 9.000000e-08 reg 1.000000e+04 train accuracy: 0.358612 val accuracy: 0.372000
lr 9.000000e-08 reg 1.500000e+04 train accuracy: 0.337245 val accuracy: 0.351000
lr 1.000000e-07 reg 9.000000e+03 train accuracy: 0.361286 val accuracy: 0.372000
lr 1.000000e-07 reg 1.000000e+04 train accuracy: 0.358367 val accuracy: 0.371000
lr 1.000000e-07 reg 1.500000e+04 train accuracy: 0.344653 val accuracy: 0.355000
best validation accuracy achieved during cross-validation: 0.377000

```

```

[7]: # evaluate on test set
# Evaluate the best softmax on test set
y_test_pred = best_softmax.predict(X_test)
test_accuracy = np.mean(y_test == y_test_pred)
print('softmax on raw pixels final test set accuracy: %f' % (test_accuracy, ))

```

softmax on raw pixels final test set accuracy: 0.370000

Inline Question - True or False

It's possible to add a new datapoint to a training set that would leave the SVM loss unchanged, but this is not the case with the Softmax classifier loss.

Your answer: True

Your explanation:

We could find data that when in SVM it's score for all other classes would lead to $S_j - S_{yi} < -1$, which then the loss would only be 0. However, when in softmax, since e^x could never be zero, $e^{(s_j)}$ could never be zero which means $e^{(s_{yi})}/\sum_j e^{(s_j)}$ would always be < 1 , therefore $-\log(e^{(s_{yi})}/\sum_j e^{(s_j)})$ would never be 0.

```
[8]: # Visualize the learned weights for each class
w = best_softmax.W[:-1,:] # strip out the bias
w = w.reshape(32, 32, 3, 10)

w_min, w_max = np.min(w), np.max(w)

classes = ['plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 'ship', 'truck']
for i in range(10):
    plt.subplot(2, 5, i + 1)

    # Rescale the weights to be between 0 and 255
    wimg = 255.0 * (w[:, :, :, i].squeeze() - w_min) / (w_max - w_min)
    plt.imshow(wimg.astype('uint8'))
    plt.axis('off')
    plt.title(classes[i])
```

