

# Developing an AI application

Going forward, AI algorithms will be incorporated into more and more everyday applications. For example, you might want to include an image classifier in a smart phone app. To do this, you'd use a deep learning model trained on hundreds of thousands of images as part of the overall application architecture. A large part of software development in the future will be using these types of models as common parts of applications.

In this project, you'll train an image classifier to recognize different species of flowers. You can imagine using something like this in a phone app that tells you the name of the flower your camera is looking at. In practice you'd train this classifier, then export it for use in your application. We'll be using [this dataset](http://www.robots.ox.ac.uk/~vgg/data/flowers/102/index.html) (<http://www.robots.ox.ac.uk/~vgg/data/flowers/102/index.html>) of 102 flower categories, you can see a few examples below.

hard-leaved  
pocket orchid



cautleya spicata



orange dahlia



The project is broken down into multiple steps:

- Load and preprocess the image dataset
- Train the image classifier on your dataset
- Use the trained classifier to predict image content

We'll lead you through each part which you'll implement in Python.

When you've completed this project, you'll have an application that can be trained on any set of labeled images. Here your network will be learning about flowers and end up as a command line application. But, what you do with your new skills depends on your imagination and effort in

building a dataset. For example, imagine an app where you take a picture of a car, it tells you what the make and model is, then looks up information about it. Go build your own dataset and make something new.

First up is importing the packages you'll need. It's good practice to keep all the imports at the beginning of your code. As you work through this notebook and find you need to import a package, make sure to add the import up here.

```
In [1]: # Import the necessary packages
get_ipython().run_line_magic('matplotlib', 'inline')
get_ipython().run_line_magic('config', "InlineBackend.figure_format = 're

import matplotlib.pyplot as plt

import torch
import numpy as np
from torch import nn
from torch import optim
import torch.nn.functional as F
from torchvision import datasets, transforms, models
import json
from PIL import Image

device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
```

## Load the data

Here you'll use `torchvision` to load the data ([documentation \(http://pytorch.org/docs/0.3.0/torchvision/index.html\)](http://pytorch.org/docs/0.3.0/torchvision/index.html)). The data should be included alongside this notebook, otherwise you can [download it here \(https://s3.amazonaws.com/content.udacity-data.com/nd089/flower\\_data.tar.gz\)](https://s3.amazonaws.com/content.udacity-data.com/nd089/flower_data.tar.gz). The dataset is split into three parts, training, validation, and testing. For the training, you'll want to apply transformations such as random scaling, cropping, and flipping. This will help the network generalize leading to better performance. You'll also need to make sure the input data is resized to 224x224 pixels as required by the pre-trained networks.

The validation and testing sets are used to measure the model's performance on data it hasn't seen yet. For this you don't want any scaling or rotation transformations, but you'll need to resize then crop the images to the appropriate size.

The pre-trained networks you'll use were trained on the ImageNet dataset where each color channel was normalized separately. For all three sets you'll need to normalize the means and standard deviations of the images to what the network expects. For the means, it's `[0.485, 0.456, 0.406]` and for the standard deviations `[0.229, 0.224, 0.225]`, calculated from the ImageNet images. These values will shift each color channel to be centered at 0 and range from -1 to 1.

```
In [2]: # Load the relevant datasets
data_dir = 'flowers'
train_dir = data_dir + '/train'
valid_dir = data_dir + '/valid'
test_dir = data_dir + '/test'
```

```
In [3]: # Define the transforms for the training, validation, and testing sets
train_transforms = transforms.Compose([transforms.RandomRotation(30),
                                       transforms.RandomResizedCrop(224),
                                       transforms.RandomHorizontalFlip(),
                                       transforms.ToTensor(),
                                       transforms.Normalize([0.485, 0.456, 0.229], [0.224, 0.224, 0.224])])

data_transforms = transforms.Compose([transforms.Resize(256),
                                       transforms.CenterCrop(224),
                                       transforms.ToTensor(),
                                       transforms.Normalize([0.485, 0.456, 0.229], [0.224, 0.224, 0.224])])

# Load the datasets with ImageFolder
train_datasets = datasets.ImageFolder(train_dir, transform=train_transforms)
valid_datasets = datasets.ImageFolder(valid_dir, transform=data_transforms)
test_datasets = datasets.ImageFolder(test_dir, transform=data_transforms)

# Using the image datasets and the trainforms, define the dataloaders
trainloader = torch.utils.data.DataLoader(train_datasets, batch_size=64,
                                           shuffle=True)
validloader = torch.utils.data.DataLoader(valid_datasets, batch_size=32,
                                           shuffle=False)
testloader = torch.utils.data.DataLoader(test_datasets, batch_size=32,
                                           shuffle=False)
```

## Label mapping

You'll also need to load in a mapping from category label to category name. You can find this in the file `cat_to_name.json`. It's a JSON object which you can read in with the `json` module (<https://docs.python.org/2/library/json.html>). This will give you a dictionary mapping the integer encoded categories to the actual names of the flowers.

```
In [4]: with open('cat_to_name.json', 'r') as f:
        cat_to_name = json.load(f)
```

# Building and training the classifier

Now that the data is ready, it's time to build and train the classifier. As usual, you should use one of the pretrained models from `torchvision.models` to get the image features. Build and train a new feed-forward classifier using those features.

We're going to leave this part up to you. If you want to talk through it with someone, chat with your fellow students! You can also ask questions on the forums or join the instructors in office hours.

Refer to [the rubric \(https://review.udacity.com/#!/rubrics/1663/view\)](https://review.udacity.com/#!/rubrics/1663/view) for guidance on successfully completing this section. Things you'll need to do:

- Load a [pre-trained network \(http://pytorch.org/docs/master/torchvision/models.html\)](http://pytorch.org/docs/master/torchvision/models.html) (If you need a starting point, the VGG networks work great and are straightforward to use)
- Define a new, untrained feed-forward network as a classifier, using ReLU activations and dropout
- Train the classifier layers using backpropagation using the pre-trained network to get the features
- Track the loss and accuracy on the validation set to determine the best hyperparameters

We've left a cell open for you below, but use as many as you need. Our advice is to break the problem up into smaller parts you can run separately. Check that each part is doing what you expect, then move on to the next. You'll likely find that as you work through each part, you'll need to go back and modify your previous code. This is totally normal!

When training make sure you're updating only the weights of the feed-forward network. You should be able to get the validation accuracy above 70% if you build everything right. Make sure to try different hyperparameters (learning rate, units in the classifier, epochs, etc) to find the best model. Save those hyperparameters to use as default values in the next part of the project.

```
In [5]: # Choose a pretrained network model
model = models.densenet121(pretrained=True)

# Freeze parameters so the program doesn't backprop through them
for param in model.parameters():
    param.requires_grad = False

# Build a feed-forward classifier
input_size = 1024
output_size = 102
hidden_layer_units = 516
classifier = nn.Sequential(nn.Linear(input_size, hidden_layer_units),
                           nn.ReLU(),
                           nn.Dropout(p=0.5),
                           nn.Linear(hidden_layer_units, output_size),
                           nn.LogSoftmax(dim=1))

model.classifier = classifier

# Define the criterion and optimizer
criterion = nn.NLLLoss()
optimizer = optim.Adam(model.classifier.parameters(), lr=0.003)
```

/opt/conda/lib/python3.6/site-packages/torchvision-0.2.1-py3.6.egg/torchvision/models/densenet.py:212: UserWarning: nn.init.kaiming\_normal is now deprecated in favor of nn.init.kaiming\_normal\_.

```
In [6]: # Implement a function for the validation pass
def validation(model, validloader, criterion):
    valid_loss = 0
    accuracy = 0
    for images, labels in validloader:

        if torch.cuda.is_available():
            images, labels = images.cuda(), labels.cuda()

        output = model.forward(images)
        valid_loss += criterion(output, labels).data.item()

        ps = torch.exp(output).data
        equality = (labels.data == ps.max(1)[1])
        accuracy += equality.type_as(torch.FloatTensor()).mean()

    return valid_loss, accuracy
```

```
In [7]: # Train the network with a cross-validation pass
if torch.cuda.is_available():
    model.cuda()
```

```

epochs = 3
steps = 0
running_loss = 0
print_every = 40

for e in range(epochs):
    model.train()

    for images, labels in iter(trainloader):
        steps += 1

        if torch.cuda.is_available():
            images, labels = images.cuda(), labels.cuda()

        optimizer.zero_grad()

        output = model.forward(images)
        loss = criterion(output, labels)
        loss.backward()
        optimizer.step()

        running_loss += loss.data.item()

    if steps % print_every == 0:
        # Make sure network is in eval mode for inference
        model.eval()

        # Turn off gradients for validation, to save memory and computation
        with torch.no_grad():
            valid_loss, accuracy = validation(model, validloader, criterion)

        print("Epoch: {}/{}.. ".format(e+1, epochs),
              "Training Loss: {:.3f}.. ".format(running_loss/print_every),
              "Validation Loss: {:.3f}.. ".format(valid_loss/len(validloader)),
              "Validation Accuracy: {:.3f}".format(accuracy/len(validloader)))

        running_loss = 0

        # Make sure training is back on
        model.train()

```

```

Epoch: 1/3.. Training Loss: 4.160.. Validation Loss: 3.174.. Validation Accuracy: 0.262
Epoch: 1/3.. Training Loss: 2.946.. Validation Loss: 1.786.. Validation Accuracy: 0.636
Epoch: 2/3.. Training Loss: 2.219.. Validation Loss: 1.213.. Validation Accuracy: 0.719
Epoch: 2/3.. Training Loss: 1.860.. Validation Loss: 0.956.. Validation Accuracy: 0.793
Epoch: 2/3.. Training Loss: 1.711.. Validation Loss: 0.824.. Validation Accuracy: 0.824

```

```

tion Accuracy: 0.779
Epoch: 3/3.. Training Loss: 1.512.. Validation Loss: 0.672.. Valida
tion Accuracy: 0.844
Epoch: 3/3.. Training Loss: 1.477.. Validation Loss: 0.625.. Valida
tion Accuracy: 0.842

```

## Testing your network

It's good practice to test your trained network on test data, images the network has never seen either in training or validation. This will give you a good estimate for the model's performance on completely new images. Run the test images through the network and measure the accuracy, the same way you did validation. You should be able to reach around 70% accuracy on the test set if the model has been trained well.

```

In [8]: # Implement a function for the test pass
def testing(model, testloader, criterion):
    test_loss = 0
    accuracy = 0
    for images, labels in testloader:

        if torch.cuda.is_available():
            images, labels = images.cuda(), labels.cuda()

        output = model.forward(images)
        test_loss += criterion(output, labels).data.item()

        ps = torch.exp(output).data
        equality = (labels.data == ps.max(1)[1])
        accuracy += equality.type_as(torch.FloatTensor()).mean()

    return test_loss, accuracy

```

```

In [9]: # Do validation on the test set
model.eval()

# Turn off gradients for validation, to save memory and computations
with torch.no_grad():
    test_loss, accuracy = testing(model, testloader, criterion)

print("Test Loss: {:.3f}.. ".format(test_loss/len(testloader)),
      "Test Accuracy: {:.3f}".format(accuracy/len(testloader)))

```

```

Test Loss: 0.670.. Test Accuracy: 0.818

```



## Save the checkpoint

Now that your network is trained, save the model so you can load it later for making predictions. You probably want to save other things such as the mapping of classes to indices which you get from one of the image datasets: `image_datasets['train'].class_to_idx`. You can attach this to the model as an attribute which makes inference easier later on.

```
model.class_to_idx = image_datasets['train'].class_to_idx
```

Remember that you'll want to completely rebuild the model later so you can use it for inference. Make sure to include any information you need in the checkpoint. If you want to load the model and keep training, you'll want to save the number of epochs as well as the optimizer state, `optimizer.state_dict`. You'll likely want to use this trained model in the next part of the project, so best to save it now.

```
In [10]: checkpoint = {'input_size': input_size,
                       'output_size': output_size,
                       'hidden_layer_units': hidden_layer_units,
                       'class_to_idx': train_datasets.class_to_idx,
                       'optimizer_dict': optimizer.state_dict(),
                       'classifier': model.classifier,
                       'state_dict': model.state_dict()}

torch.save(checkpoint, 'checkpoint.pth')
```

## Loading the checkpoint

At this point it's good to write a function that can load a checkpoint and rebuild the model. That way you can come back to this project and keep working on it without having to retrain the network.

```
In [11]: # Load a checkpoint and rebuild the model
def load_checkpoint(filepath):
    pretrained_model = models.densenet121(pretrained=True)
    if torch.cuda.is_available():
        loaded_model = torch.load(filepath)
    else:
        loaded_model = torch.load(filepath, map_location=lambda storage,
pretrained_model.class_idx = loaded_model['class_to_idx']
pretrained_model.classifier = loaded_model['classifier']
pretrained_model.load_state_dict(loaded_model['state_dict'])

    for param in pretrained_model.parameters():
        param.requires_grad = False

    if torch.cuda.is_available():
        pretrained_model.cuda()

    return pretrained_model
```

```
In [12]: model = load_checkpoint('checkpoint.pth')

/opt/conda/lib/python3.6/site-packages/torchvision-0.2.1-py3.6.egg/tor
chvision/models/densenet.py:212: UserWarning: nn.init.kaiming_normal i
s now deprecated in favor of nn.init.kaiming_normal_.
```

# Inference for classification

Now you'll write a function to use a trained network for inference. That is, you'll pass an image into the network and predict the class of the flower in the image. Write a function called `predict` that takes an image and a model, then returns the top  $K$  most likely classes along with the probabilities. It should look like

```
probs, classes = predict(image_path, model)
print(probs)
print(classes)
> [ 0.01558163  0.01541934  0.01452626  0.01443549  0.01407339 ]
> [ '70', '3', '45', '62', '55' ]
```

First you'll need to handle processing the input image such that it can be used in your network.

## Image Preprocessing

You'll want to use `PIL` to load the image ([documentation \(https://pillow.readthedocs.io/en/latest/reference/Image.html\)](https://pillow.readthedocs.io/en/latest/reference/Image.html)). It's best to write a function that preprocesses the image so it can be used as input for the model. This function should process the images in the same manner used for training.

First, resize the images where the shortest side is 256 pixels, keeping the aspect ratio. This can be done with the `thumbnail` (<http://pillow.readthedocs.io/en/3.1.x/reference/Image.html#PIL.Image.Image.thumbnail>) or `resize` (<http://pillow.readthedocs.io/en/3.1.x/reference/Image.html#PIL.Image.Image.thumbnail>) methods. Then you'll need to crop out the center 224x224 portion of the image.

Color channels of images are typically encoded as integers 0-255, but the model expects floats 0-1. You'll need to convert the values. It's easiest with a Numpy array, which you can get from a PIL image like so `np_image = np.array(pil_image)`.

As before, the network expects the images to be normalized in a specific way. For the means, it's `[0.485, 0.456, 0.406]` and for the standard deviations `[0.229, 0.224, 0.225]`. You'll want to subtract the means from each color channel, then divide by the standard deviation.

And finally, PyTorch expects the color channel to be the first dimension but it's the third dimension in the PIL image and Numpy array. You can reorder dimensions using `ndarray.transpose` (<https://docs.scipy.org/doc/numpy-1.13.0/reference/generated/numpy.ndarray.transpose.html>). The color channel needs to be first and retain the order of the other two dimensions.

```
In [13]: def process_image(image):  
    ''' Scales, crops, and normalizes a PIL image for a PyTorch model,  
        returns an Numpy array  
    '''  
  
    # Resize the image and make the shorter side 256 pixels  
    width, height = image.size  
    ratio = width / height  
    if width < height:  
        width = 256  
        height = width / ratio  
    elif width > height:  
        height = 256  
        width = ratio * height  
    else:  
        width, height = 256  
    image = image.resize((round(width), round(height)))  
  
    # Crop the image and convert it to an nparray  
    np_image = np.array(image.crop((16, 16, 240, 240)))  
  
    # Normalize the image  
    means = np.array([0.485, 0.456, 0.406])  
    stds = np.array([0.229, 0.224, 0.225])  
    np_image = (np_image / 255 - means) / stds  
  
    # Change the color channel to meet PyTorch expectations  
    return np_image.transpose(2, 0, 1)
```

To check your work, the function below converts a PyTorch tensor and displays it in the notebook. If your `process_image` function works, running the output through this function should return the original image (except for the cropped out portions).

```
In [14]: def imshow(image, ax=None, title=None):
    """Imshow for Tensor."""
    if ax is None:
        fig, ax = plt.subplots()

    # PyTorch tensors assume the color channel is the first dimension,
    # but matplotlib assumes it is the third dimension
    image = image.transpose((1, 2, 0))

    # Undo preprocessing
    mean = np.array([0.485, 0.456, 0.406])
    std = np.array([0.229, 0.224, 0.225])
    image = std * image + mean

    # Image needs to be clipped between 0 and 1 or it looks like noise wh
    image = np.clip(image, 0, 1)

    ax.imshow(image)

    return ax
```

## Class Prediction

Once you can get images in the correct format, it's time to write a function for making predictions with your model. A common practice is to predict the top 5 or so (usually called top- $K$ ) most probable classes. You'll want to calculate the class probabilities then find the  $K$  largest values.

To get the top  $K$  largest values in a tensor use `x.topk(k)`. (<http://pytorch.org/docs/master/torch.html#torch.topk>). This method returns both the highest  $k$  probabilities and the indices of those probabilities corresponding to the classes. You need to convert from these indices to the actual class labels using `class_to_idx` which hopefully you added to the model or from an `ImageFolder` you used to load the data ([see here](#)). Make sure to invert the dictionary so you get a mapping from index to class as well.

Again, this method should take a path to an image and a model checkpoint, then return the probabilities and classes.

```
probs, classes = predict(image_path, model)
print(probs)
print(classes)
> [ 0.01558163  0.01541934  0.01452626  0.01443549  0.01407339]
> ['70', '3', '45', '62', '55']
```

```
In [15]: def predict(image_path, model, topk=5):
    ''' Predict the class (or classes) of an image using a trained deep l
    ...

    model.eval()
    # Implement the code to predict the class from an image file
    image = Image.open(image_path)
    image_tensor = torch.from_numpy(process_image(image))
    image_tensor.unsqueeze_(0)
    model = model.double()

    # Swap the dictionary keys and values to obtain the numerical classes
    idx_to_class = {v: k for k, v in model.class_idx.items()}

    model.to(device)
    output = model.forward(image_tensor.to(device))
    largest = torch.exp(output).data.topk(topk)
    probs, classes = largest[0].tolist()[0], list(map(lambda i:idx_to_cla
                                                    largest[1].tolist())

    return probs, classes

probs, classes = predict('flowers/train/24/image_06826.jpg', model)
print(probs)
print(classes)
```

```
[0.7514867028398498, 0.06913552856499615, 0.03214981142966451, 0.02968
5848140125815, 0.02435186864020786]
['24', '74', '88', '43', '99']
```

```
In [16]: def predict_names(image_path, model, cat_to_name, topk=5):
    probs, classes = predict(image_path, model, topk)
    names = []
    for c in classes:
        names.append(cat_to_name[c])
    return probs, names

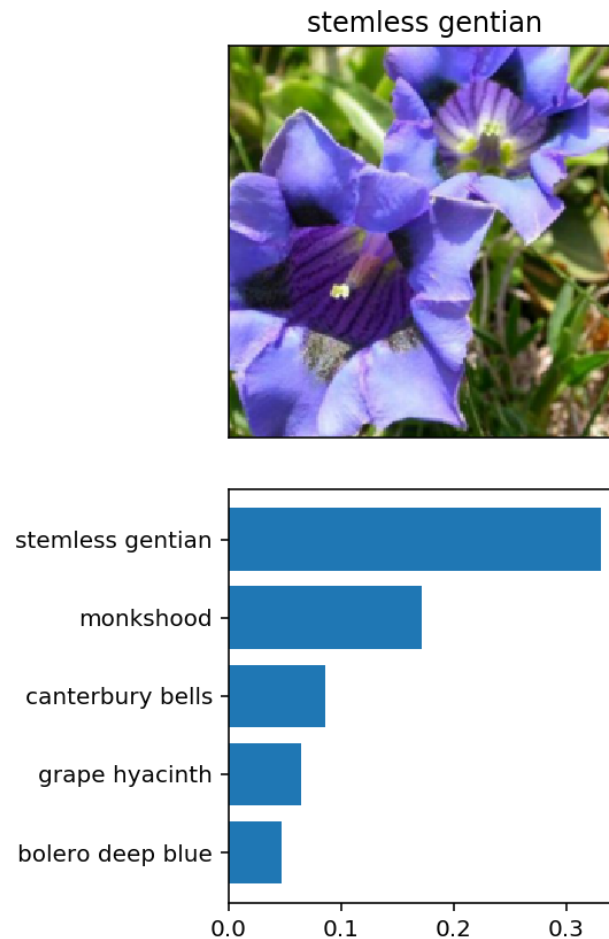
probs, names = predict_names('flowers/train/24/image_06826.jpg', model, c

print(probs)
print(names)
```

```
[0.7514867028398498, 0.06913552856499615, 0.03214981142966451, 0.02968
5848140125815, 0.02435186864020786]
['red ginger', 'rose', 'cyclamen', 'sword lily', 'bromelia']
```

## Sanity Checking

Now that you can use a trained model for predictions, check to make sure it makes sense. Even if the testing accuracy is high, it's always good to check that there aren't obvious bugs. Use `matplotlib` to plot the probabilities for the top 5 classes as a bar graph, along with the input image. It should look like this:



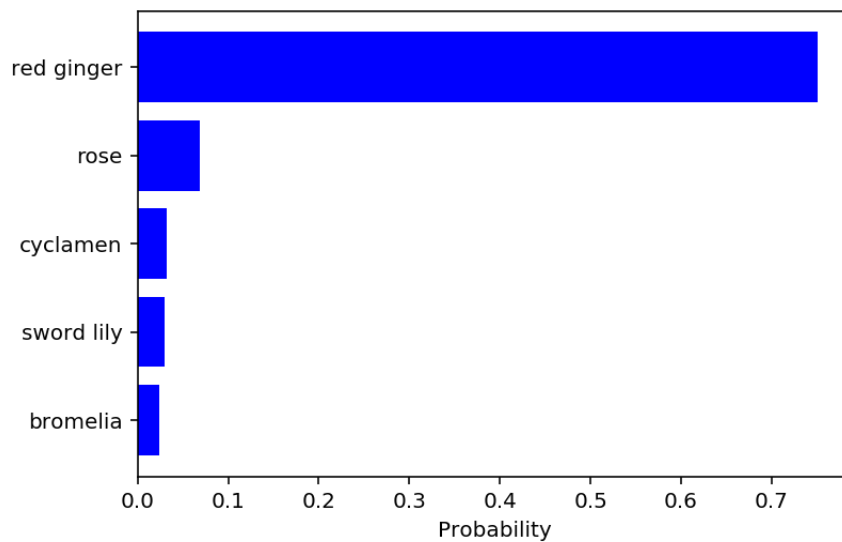
You can convert from the class integer encoding to actual flower names with the `cat_to_name.json` file (should have been loaded earlier in the notebook). To show a PyTorch tensor as an image, use the `imshow` function defined above.

```
In [17]: test_flower = 'flowers/train/24/image_06826.jpg'
probs, names = predict_names(test_flower, model, cat_to_name)
sanity_check = process_image(Image.open(test_flower))

# Display the test flower image
flower_ax = imshow(sanity_check)
flower_ax.set_yticks(list())
flower_ax.set_xticks(list())
flower_ax.set_title(names[0])
```


```
# Plot the predictions and related probabilities
fig, ax = plt.subplots()
y_pos = np.arange(len(names))
ax.barh(y_pos, probs, align='center', color='blue')
ax.set_yticks(y_pos)
ax.set_yticklabels(names)
ax.invert_yaxis()
ax.set_xlabel('Probability')
plt.show()
```


red ginger





In [ ]:



 Present

 Slides

 Themes

 Help