

Classifying CIFAR10 images using a ResNet and Regularization techniques in PyTorch

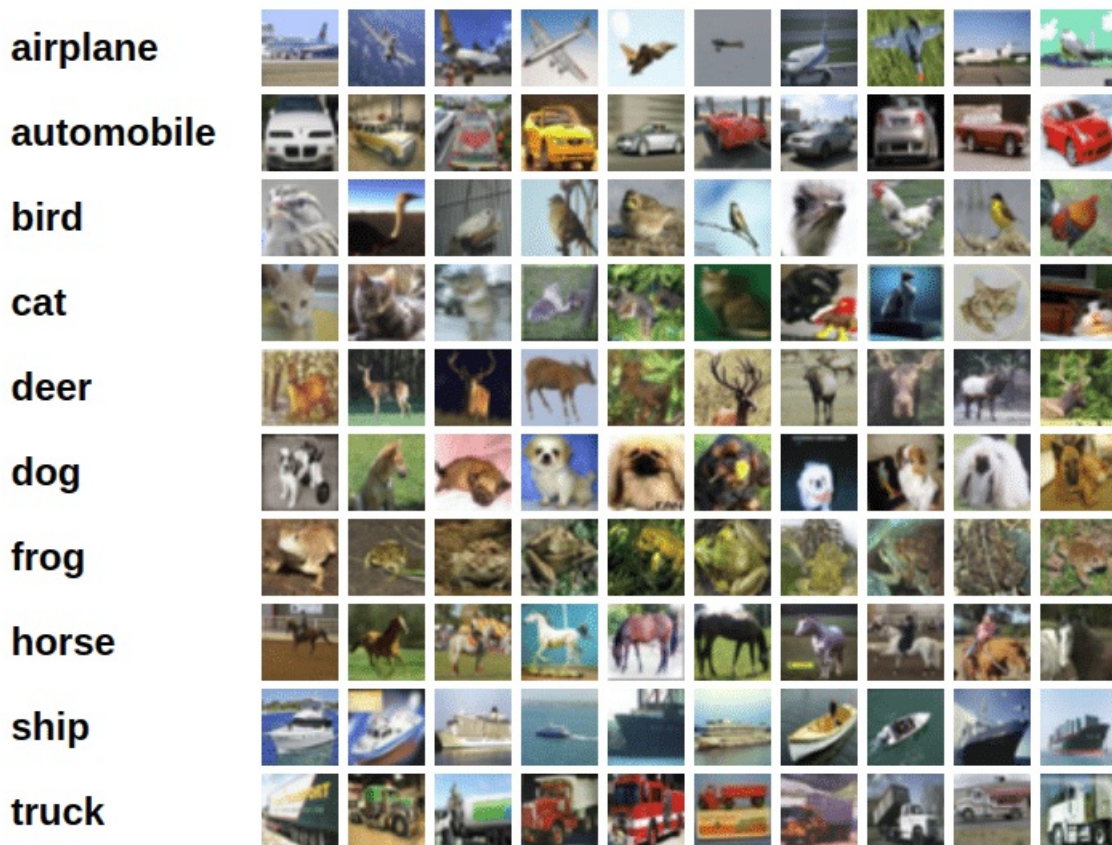
Training an image classifier from scratch to over 90% accuracy in less than 5 minutes on a single GPU

Part 6 of "PyTorch: Zero to GANs"

This post is the sixth in a series of tutorials on building deep learning models with PyTorch, an open source neural networks library. Check out the full series:

1. [PyTorch Basics: Tensors & Gradients](#)
2. [Linear Regression & Gradient Descent](#)
3. [Image Classification using Logistic Regression](#)
4. [Training Deep Neural Networks on a GPU](#)
5. [Image Classification using Convolutional Neural Networks](#)
6. [Data Augmentation, Regularization and ResNets](#)
7. [Generating Images using Generative Adversarial Networks](#)

This notebook is an extension to the tutorial [Image Classification using CNNs in PyTorch](#), where we trained a deep convolutional neural network to classify images from the CIFAR10 dataset with around 75% accuracy. Here are some images from the dataset:



In this tutorial, we'll use the following techniques to achieve over 90% accuracy in less than 5 minutes:

- Data normalization
- Data augmentation
- Residual connections
- Batch normalization
- Learning rate scheduling
- Weight Decay
- Gradient clipping
- Adam optimizer

System Setup

This notebook is hosted on [Jovian.ml](https://jovian.ml), a platform for sharing data science projects. If you want to follow along and run the code as you read, you can choose the "Run on Kaggle" option from the "Run" dropdown above. Remember select "GPU" as the accelerator and turn on internet from the sidebar within the Kaggle notebook.

Otherwise, to run the code on your machine, you can clone the notebook, install the required dependencies using [conda](https://conda.io/), and start Jupyter by running the following commands:

```
pip install jovian --upgrade          # Install the jovian library
jovian clone aakashns/05b-cifar10-resnet  # Download notebook & dependencies
cd 05b-cifar10-resnet                # Enter the created directory
conda create -n 05b-cifar10-resnet     # Create virtual env
conda activate 05b-cifar10-resnet      # Activate virtual env
conda install jupyter                 # Install Jupyter
jupyter notebook                      # Start Jupyter
```

For a more detailed explanation of the above steps, check out the System setup section in the [first notebook](#). Before you start executing the code below, you may want to clear the cell outputs by selecting "Kernel > Restart and Clear Output" from the Jupyter notebook menu bar, to avoid confusion.

We begin by importing the required modules & libraries.

In [1]:

```
# Uncomment and run the commands below if imports fail
# !conda install numpy pandas pytorch torchvision cpuonly -c pytorch -y
# !pip install matplotlib --upgrade --quiet
```

In [27]:

```
import os
import torch
import torchvision
import tarfile
import torch.nn as nn
import numpy as np
import torch.nn.functional as F
from torchvision.datasets.utils import download_url
from torchvision.datasets import ImageFolder
from torch.utils.data import DataLoader
import torchvision.transforms as tt
from torch.utils.data import random_split
from torchvision.utils import make_grid
import matplotlib.pyplot as plt
%matplotlib inline
```

In [3]:

```
project_name='05b-cifar10-resnet'
```

Preparing the Data

Let's begin by downloading the dataset and creating PyTorch datasets to load the data, just as we did in the previous tutorial.

In [4]:

```
# Download the dataset
dataset_url = "http://files.fast.ai/data/cifar10.tgz"
download_url(dataset_url, '.')

# Extract from archive
with tarfile.open('./cifar10.tgz', 'r:gz') as tar:
    tar.extractall(path='./data')

# Look into the data directory
data_dir = './data/cifar10'
print(os.listdir(data_dir))
classes = os.listdir(data_dir + "/train")
```

```
print(classes)
```

Downloading <http://files.fast.ai/data/cifar10.tgz> to ./cifar10.tgz

```
['labels.txt', 'test', 'train']  
['automobile', 'dog', 'horse', 'frog', 'cat', 'deer', 'bird', 'ship', 'airplane', 'truck']
```

There are a few important changes we'll make while creating the PyTorch datasets:

1. **Use test set for validation:** Instead of setting aside a fraction (e.g. 10%) of the data from the training set for validation, we'll simply use the test set as our validation set. This just gives a little more data to train with. In general, once you have picked the best model architecture & hyperparameters using a fixed validation set, it is a good idea to retrain the same model on the entire dataset just to give it a small final boost in performance.
2. **Channel-wise data normalization:** We will normalize the image tensors by subtracting the mean and dividing by the standard deviation across each channel. As a result, the mean of the data across each channel is 0, and standard deviation is 1. Normalizing the data prevents the values from any one channel from disproportionately affecting the losses and gradients while training, simply by having a higher or wider range of values than others.
3. **Randomized data augmentations:** We will apply randomly chosen transformations while loading images from the training dataset. Specifically, we will pad each image by 4 pixels, and then take a random crop of size 32 x 32 pixels, and then flip the image horizontally with a 50% probability. Since the transformation will be applied randomly and dynamically each time a particular image is loaded, the model sees slightly different images in each epoch of training, which allows it to generalize better.



In [5]:

```
# Data transforms (normalization & data augmentation)  
stats = ((0.4914, 0.4822, 0.4465), (0.2023, 0.1994, 0.2010))  
train_tfms = tt.Compose([tt.RandomCrop(32, padding=4, padding_mode='reflect'),  
                        tt.RandomHorizontalFlip(),  
                        tt.ToTensor(),  
                        tt.Normalize(*stats, inplace=True)])  
valid_tfms = tt.Compose([tt.ToTensor(), tt.Normalize(*stats)])
```

In [6]:

```
# PyTorch datasets  
train_ds = ImageFolder(data_dir+'train', train_tfms)  
valid_ds = ImageFolder(data_dir+'test', valid_tfms)
```

Next, we can create data loaders for retrieving images in batches. We'll use a relatively large batch size of 500 to utilize a larger portion of the GPU RAM. You can try reducing the batch size & restarting the kernel if you face an "out of memory" error.

In [7]:

```
batch_size = 400
```

In [8]:

```
# PyTorch data loaders  
train_dl = DataLoader(train_ds, batch_size, shuffle=True, num_workers=3, pin_memory=True)  
valid_dl = DataLoader(valid_ds, batch_size*2, num_workers=3, pin_memory=True)
```

Let's take a look at some sample images from the training dataloader.

In [9]:

```
def show_batch(dl):
    for images, labels in dl:
        fig, ax = plt.subplots(figsize=(12, 12))
        ax.set_xticks([]); ax.set_yticks([])
        ax.imshow(make_grid(images[:64], nrow=8).permute(1, 2, 0))
        break
```

In [10]:

```
show_batch(train_dl)
```



The colors seem out of place because of the normalization. Note that normalization is also applied during inference. If you look closely, you can see the cropping and reflection padding in some of the images. Horizontal flip is a bit difficult to detect from visual inspection.

Using a GPU

To seamlessly use a GPU, if one is available, we define a couple of helper functions (`get_default_device` & `to_device`) and a helper class `DeviceDataLoader` to move our model & data to the GPU as required. These are described in more detail in a [previous tutorial](#).

In [11]:

```
def get_default_device():
    """Pick GPU if available, else CPU"""
    if torch.cuda.is_available():
```



```

-- torch.cuda.is_available(),
    return torch.device('cuda')
else:
    return torch.device('cpu')

def to_device(data, device):
    """Move tensor(s) to chosen device"""
    if isinstance(data, (list,tuple)):
        return [to_device(x, device) for x in data]
    return data.to(device, non_blocking=True)

class DeviceDataLoader():
    """Wrap a dataloader to move data to a device"""
    def __init__(self, dl, device):
        self.dl = dl
        self.device = device

    def __iter__(self):
        """Yield a batch of data after moving it to device"""
        for b in self.dl:
            yield to_device(b, self.device)

    def __len__(self):
        """Number of batches"""
        return len(self.dl)

```

Based on where you're running this notebook, your default device could be a CPU (`torch.device('cpu')`) or a GPU (`torch.device('cuda')`)

In [12]:

```

device = get_default_device()
device

```

Out[12]:

```

device(type='cuda')

```

We can now wrap our training and validation data loaders using `DeviceDataLoader` for automatically transferring batches of data to the GPU (if available).

In [13]:

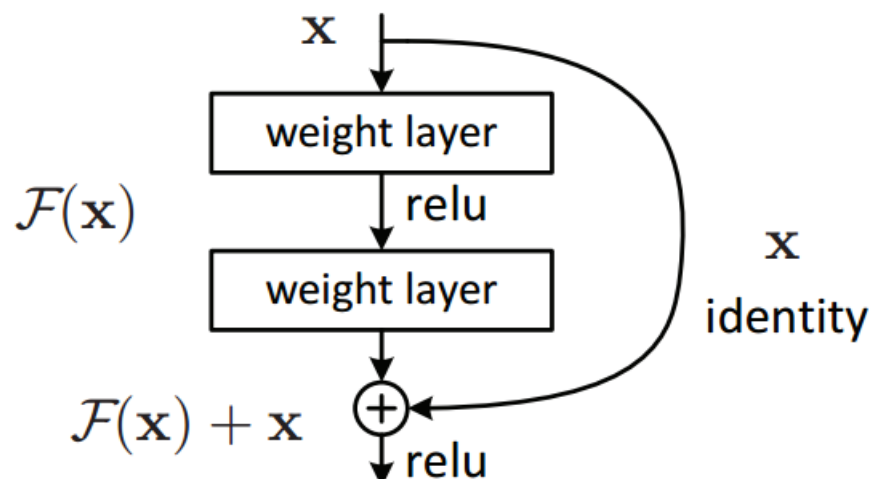
```

train_dl = DeviceDataLoader(train_dl, device)
valid_dl = DeviceDataLoader(valid_dl, device)

```

Model with Residual Blocks and Batch Normalization

One of the key changes to our CNN model this time is the addition of the residual block, which adds the original input back to the output feature map obtained by passing the input through one or more convolutional layers.



Here is a very simply Residual block:

In [14]:

```
class SimpleResidualBlock(nn.Module):
    def __init__(self):
        super().__init__()
        self.conv1 = nn.Conv2d(in_channels=3, out_channels=3, kernel_size=3, stride=1, padding=1)
        self.relu1 = nn.ReLU()
        self.conv2 = nn.Conv2d(in_channels=3, out_channels=3, kernel_size=3, stride=1, padding=1)
        self.relu2 = nn.ReLU()

    def forward(self, x):
        out = self.conv1(x)
        out = self.relu1(out)
        out = self.conv2(out)
        return self.relu2(out) + x # ReLU can be applied before or after adding the input
```

In [15]:

```
simple_resnet = to_device(SimpleResidualBlock(), device)
```

```
for images, labels in train_dl:
    out = simple_resnet(images)
    print(out.shape)
    break
```

```
del simple_resnet, images, labels
torch.cuda.empty_cache()
```

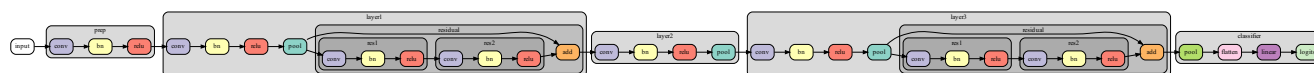
```
torch.Size([400, 3, 32, 32])
```

This seeming small change produces a drastic improvement in the performance of the model. Also, after each convolutional layer, we'll add a batch normalization layer, which normalizes the outputs of the previous layer.

Go through the following blog posts to learn more:

- Why and how residual blocks work: <https://towardsdatascience.com/residual-blocks-building-blocks-of-resnet-fd90ca15d6ec>
- Batch normalization and dropout explained: <https://towardsdatascience.com/batch-normalization-and-dropout-in-neural-networks-explained-with-pytorch-47d7a8459bcd>

We will use the ResNet9 architecture, as described in [this blog series](#) :



In [16]:

```
def accuracy(outputs, labels):
    _, preds = torch.max(outputs, dim=1)
    return torch.tensor(torch.sum(preds == labels).item() / len(preds))

class ImageClassificationBase(nn.Module):
    def training_step(self, batch):
        images, labels = batch
        out = self(images) # Generate predictions
        loss = F.cross_entropy(out, labels) # Calculate loss
        return loss

    def validation_step(self, batch):
        images, labels = batch
        out = self(images) # Generate predictions
        loss = F.cross_entropy(out, labels) # Calculate loss
        acc = accuracy(out, labels) # Calculate accuracy
        return {'val_loss': loss.detach(), 'val_acc': acc}

    def validation_epoch_end(self, outputs):
        batch_losses = [x['val_loss'] for x in outputs]
        epoch_loss = torch.stack(batch_losses).mean() # Combine losses
        batch_accs = [x['val_acc'] for x in outputs]
```

```

epoch_acc = torch.stack(batch_accs).mean() # Combine accuracies
return {'val_loss': epoch_loss.item(), 'val_acc': epoch_acc.item()}

def epoch_end(self, epoch, result):
    print("Epoch [{}], last_lr: {:.5f}, train_loss: {:.4f}, val_loss: {:.4f}, val_acc: {:.4f}".format(
        epoch, result['lrs'][-1], result['train_loss'], result['val_loss'], result['val_acc']))

```

In [17]:

```

def conv_block(in_channels, out_channels, pool=False):
    layers = [nn.Conv2d(in_channels, out_channels, kernel_size=3, padding=1),
              nn.BatchNorm2d(out_channels),
              nn.ReLU(inplace=True)]
    if pool: layers.append(nn.MaxPool2d(2))
    return nn.Sequential(*layers)

class ResNet9(ImageClassificationBase):
    def __init__(self, in_channels, num_classes):
        super().__init__()

        self.conv1 = conv_block(in_channels, 64)
        self.conv2 = conv_block(64, 128, pool=True)
        self.res1 = nn.Sequential(conv_block(128, 128), conv_block(128, 128))

        self.conv3 = conv_block(128, 256, pool=True)
        self.conv4 = conv_block(256, 512, pool=True)
        self.res2 = nn.Sequential(conv_block(512, 512), conv_block(512, 512))

        self.classifier = nn.Sequential(nn.MaxPool2d(4),
                                         nn.Flatten(),
                                         nn.Linear(512, num_classes))

    def forward(self, xb):
        out = self.conv1(xb)
        out = self.conv2(out)
        out = self.res1(out) + out
        out = self.conv3(out)
        out = self.conv4(out)
        out = self.res2(out) + out
        out = self.classifier(out)
        return out

```

In [18]:

```

model = to_device(ResNet9(3, 10), device)
model

```

Out[18]:

```

ResNet9(
  (conv1): Sequential(
    (0): Conv2d(3, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (2): ReLU(inplace=True)
  )
  (conv2): Sequential(
    (0): Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (2): ReLU(inplace=True)
    (3): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
  )
  (res1): Sequential(
    (0): Sequential(
      (0): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
      (1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      (2): ReLU(inplace=True)
    )
    (1): Sequential(
      (0): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
      (1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      (2): ReLU(inplace=True)
    )
  )
)

```



```

for epoch in range(epochs):
    # Training Phase
    model.train()
    train_losses = []
    lrs = []
    for batch in train_loader:
        loss = model.training_step(batch)
        train_losses.append(loss)
        loss.backward()

    # Gradient clipping
    if grad_clip:
        nn.utils.clip_grad_value_(model.parameters(), grad_clip)

    optimizer.step()
    optimizer.zero_grad()

    # Record & update learning rate
    lrs.append(get_lr(optimizer))
    sched.step()

    # Validation phase
    result = evaluate(model, val_loader)
    result['train_loss'] = torch.stack(train_losses).mean().item()
    result['lrs'] = lrs
    model.epoch_end(epoch, result)
    history.append(result)
return history

```

In [21]:

```

history = [evaluate(model, valid_dl)]
history

```

Out[21]:

```
[{'val_loss': 2.307021379470825, 'val_acc': 0.07249999791383743}]
```

We're now ready to train our model. Instead of SGD (stochastic gradient descent), we'll use the Adam optimizer which uses techniques like momentum and adaptive learning rates for faster training. You can learn more about optimizers here:

<https://ruder.io/optimizing-gradient-descent/index.html>

In [62]:

```

epochs = 8
max_lr = 0.01
grad_clip = 0.1
weight_decay = 1e-4
opt_func = torch.optim.Adam

```

In [22]:

```

%%time
history += fit_one_cycle(epochs, max_lr, model, train_dl, valid_dl,
                        grad_clip=grad_clip,
                        weight_decay=weight_decay,
                        opt_func=opt_func)

```

```

Epoch [0], last_lr: 0.00393, train_loss: 1.2882, val_loss: 1.3784, val_acc: 0.5857
Epoch [1], last_lr: 0.00935, train_loss: 1.0240, val_loss: 1.7711, val_acc: 0.5280
Epoch [2], last_lr: 0.00972, train_loss: 0.8086, val_loss: 0.9093, val_acc: 0.7018
Epoch [3], last_lr: 0.00812, train_loss: 0.5827, val_loss: 0.6836, val_acc: 0.7687
Epoch [4], last_lr: 0.00556, train_loss: 0.4600, val_loss: 0.4701, val_acc: 0.8425
Epoch [5], last_lr: 0.00283, train_loss: 0.3473, val_loss: 0.3804, val_acc: 0.8715
Epoch [6], last_lr: 0.00077, train_loss: 0.2535, val_loss: 0.3279, val_acc: 0.8894
Epoch [7], last_lr: 0.00000, train_loss: 0.1893, val_loss: 0.2739, val_acc: 0.9089
CPU times: user 42.3 s, sys: 9.38 s, total: 51.7 s
Wall time: 4min 7s

```

In [52]:

In [5]:

```
train_time='4:07'
```

Our model trained to over **90% accuracy in just 5 minutes!** Try playing around with the data augmentations, network architecture & hyperparameters to achieve the following results:

1. 94% accuracy in under 10 minutes (easy)
2. 90% accuracy in under 2.5 minutes (intermediate)
3. 94% accuracy in under 5 minutes (hard)

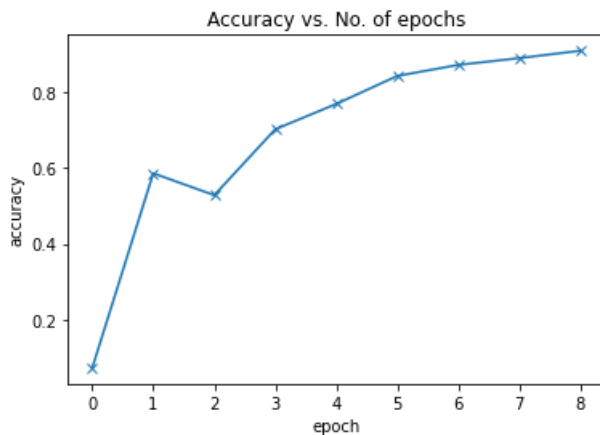
Let's plot the validation set accuracies to study how the model improves over time.

In [23]:

```
def plot_accuracies(history):  
    accuracies = [x['val_acc'] for x in history]  
    plt.plot(accuracies, '-x')  
    plt.xlabel('epoch')  
    plt.ylabel('accuracy')  
    plt.title('Accuracy vs. No. of epochs');
```

In [24]:

```
plot_accuracies(history)
```



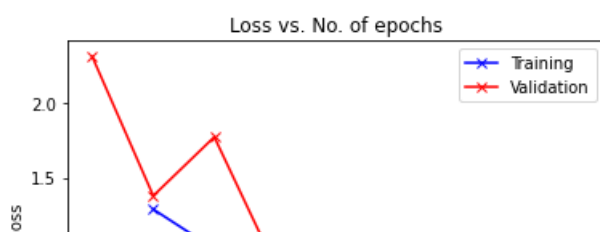
We can also plot the training and validation losses to study the trend.

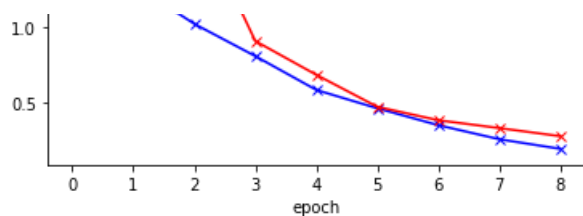
In [25]:

```
def plot_losses(history):  
    train_losses = [x.get('train_loss') for x in history]  
    val_losses = [x['val_loss'] for x in history]  
    plt.plot(train_losses, '-bx')  
    plt.plot(val_losses, '-rx')  
    plt.xlabel('epoch')  
    plt.ylabel('loss')  
    plt.legend(['Training', 'Validation'])  
    plt.title('Loss vs. No. of epochs');
```

In [26]:

```
plot_losses(history)
```





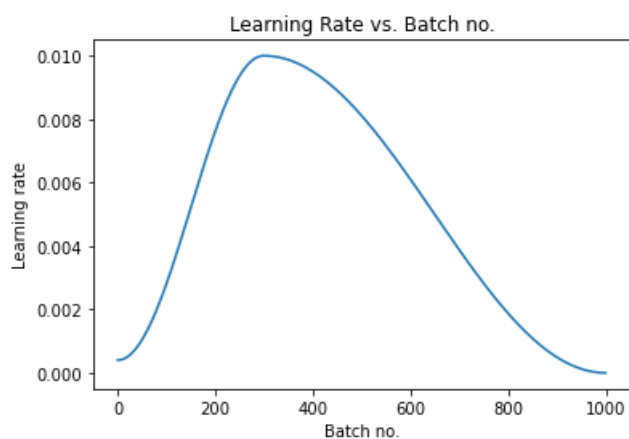
It's clear from the trend that our model isn't overfitting to the training data just yet. Finally, let's visualize how the learning rate changed over time, batch-by-batch over all the epochs.

In [37]:

```
def plot_lrs(history):
    lrs = np.concatenate([x.get('lrs', []) for x in history])
    plt.plot(lrs)
    plt.xlabel('Batch no.')
    plt.ylabel('Learning rate')
    plt.title('Learning Rate vs. Batch no.');
```

In [39]:

```
plot_lrs(history)
```



As expected, the learning rate starts at a low value, and gradually increases for 30% of the iterations to a maximum value of 0.01, and then gradually decreases to a very small value.

Save and Commit

Let's save the weights of the model, record the hyperparameters, and commit our experiment to Jovian. As you try different ideas, make sure to record every experiment so you can look back and analyze the results.

In [42]:

```
torch.save(model.state_dict(), 'cifar10-resnet9.pth')
```

In [43]:

```
!pip install jovian --upgrade --quiet
```

In [44]:

```
import jovian
```

In [66]:

```
jovian.reset()
jovian.log_hyperparams({'arch': 'resnet9'}
```

```
jovian.log_hyperparams(train=100000,
                      epochs=epochs,
                      lr=max_lr,
                      scheduler='one-cycle',
                      weight_decay=weight_decay,
                      grad_clip=grad_clip,
                      opt=opt_func.__name__)
```

[jovian] Hyperparams logged.

In [67]:

```
jovian.log_metrics(val_loss=history[-1]['val_loss'],
                  val_acc=history[-1]['val_acc'],
                  train_loss=history[-1]['train_loss'],
                  time=train_time)
```

[jovian] Metrics logged.

In []:

```
jovian.commit(project=project_name, environment=None, outputs=['cifar10-resnet9.pth'])
```

[jovian] Attempting to save notebook..

[jovian] Detected Kaggle notebook...

[jovian] Uploading notebook to https://jovian.ml/aakashns/05b-cifar10-resnet

Summary and Further Reading

Here's a summary of the different techniques used in this tutorial to improve our model performance and reduce the training time:

- **Data normalization:** We normalized the image tensors by subtracting the mean and dividing by the standard deviation of pixels across each channel. Normalizing the data prevents the pixel values from any one channel from disproportionately affecting the losses and gradients. [Learn more](#)
- **Data augmentation:** We applied random transformations while loading images from the training dataset. Specifically, we will pad each image by 4 pixels, and then take a random crop of size 32 x 32 pixels, and then flip the image horizontally with a 50% probability. [Learn more](#)
- **Residual connections:** One of the key changes to our CNN model was the addition of the residual block, which adds the original input back to the output feature map obtained by passing the input through one or more convolutional layers. We used the ResNet9 architecture [Learn more](#).
- **Batch normalization:** After each convolutional layer, we added a batch normalization layer, which normalizes the outputs of the previous layer. This is somewhat similar to data normalization, except it's applied to the outputs of a layer, and the mean and standard deviation are learned parameters. [Learn more](#)
- **Learning rate scheduling:** Instead of using a fixed learning rate, we will use a learning rate scheduler, which will change the learning rate after every batch of training. There are [many strategies](#) for varying the learning rate during training, and we used the "One Cycle Learning Rate Policy". [Learn more](#)
- **Weight Decay:** We added weight decay to the optimizer, yet another regularization technique which prevents the weights from becoming too large by adding an additional term to the loss function. [Learn more](#)
- **Gradient clipping:** We also added gradient clipping, which helps limit the values of gradients to a small range to prevent undesirable changes in model parameters due to large gradient values during training. [Learn more](#).
- **Adam optimizer:** Instead of SGD (stochastic gradient descent), we used the Adam optimizer which uses techniques like momentum and adaptive learning rates for faster training. There are many other optimizers to choose from and experiment with. [Learn more](#).

As an exercise, you should try applying each technique independently and see how much each one affects the performance and training time. As you try different experiments, you will start to cultivate the intuition for picking the right architectures, data augmentation & regularization techniques.

In []: