

Приемы работы с библиотекой PyTorch

Содержание

1	Вводные замечания	1
1.1	Torch Hub	3
2	Тензор	3
2.1	Типы элементов тензоров	5
2.2	Представление реальных данных с помощью тензоров	8
2.2.1	Работа с изображениями	8
3	Обобщения с помощью сверток	22
3.1	Сеть как подкласс <code>nn.Module</code>	22
3.2	Функциональные API	22
3.3	Обучение модели	23
3.4	Ширина сети	25
3.5	Регуляризация	25
	Список литературы	29

1. Вводные замечания

Чаще всего цикл обучения модели реализуют в виде обычного цикла `for` Python. Оптимизатор, доступный в модуле `torch.optim` PyTorch, который будет отвечать за обновление параметров.

По умолчанию в PyTorch используется модель немедленного выполнения (`eager mode`). Как только интерпретатор Python выполняет инструкцию, связанную с PyTorch, базовая реализация C++ или CUDA сразу же производит соответствующую операцию.

PyTorch также предоставляет возможности предварительной компиляции моделей с помощью TorchScript. Используя TorchScript, PyTorch может преобразовать модель в набор инструкций, которые можно независимо вызывать из Python, допустим, из программ на C++ или на мобильных устройствах. Это можно считать своего рода виртуальной машиной с ограниченным набором инструкций, предназначенным для операций с тензорами. Экспортировать модель можно либо в виде TorchScript для использования со средой выполнения Python, либо в стандартизированном формате ONNX (платформонезависимый формат описания моделей).

Сети среднего размера могут потребовать от нескольких часов до нескольких дней для обучения с нуля на больших реальных наборах данных на рабочих станциях с хорошим GPU [1, стр. 48]. Длительность обучения можно сократить за счет использования на одной машине нескольких GPU или даже еще сильнее – на кластере машин, оснащенных несколькими GPU.

Для примера создадим сеть AlexNet

```
# TorchVision включает несколько лучших нейросетевых архитектур для машинного зрения
from torchvision import models

alexnet = models.AlexNet()
```

Подав на вход `alexnet` данные четко определенного размера, мы выполним прямой проход (forward pass) по сети, при котором входной сигнал пройдет через первый набор нейронов, выходные сигналы которых будут поданы на вход следующего набора нейронов, и так до самого итогового выходного сигнала. На практике это означает, что при наличии объекта `input` нужного типа можно произвести прямой проход с помощью оператора `output = alexnet(input)`.

Но если мы так поступим, то получим мусор. А все потому, что сеть не была инициализирована: ее веса, числа, с которыми складываются и на которые умножаются входные сигналы, не были обучены на чем-либо, сеть сама по себе – чистый (или, точнее, сказать случайный) лист. Необходимо либо обучить ее с нуля, либо загрузить веса, полученные в результате предыдущего обучения [1, стр. 58].

В `models` названия в верхнем регистре соответствуют классам, реализующим популярные архитектуры, предназначенные для машинного зрения. С другой стороны, названия в нижнем регистре соответствуют функциям, создающим экземпляры моделей с заранее определенным количеством слоев и нейронов, а также, возможно, скачивающие и загружающие в них предобученные веса.

Для того чтобы привести входные изображения к нужному размеру, а их значения (цвета) примерно в один числовой диапазон, можно воспользоваться преобразованиями модуля `torchvision`

```
from torchvision import transforms

# это функция
preprocess = transforms.Compose([
    transforms.Resize(256),
    transforms.CenterCrop(224),
    transforms.ToTensor(),
    transforms.Normalize(
        mean=[0.485, 0.456, 0.406],
        std=[0.229, 0.224, 0.225]
    )
])
```

Здесь описана функция `preprocess`, масштабирующую входное изображение до размера 256×256 , обрезающую его до 224×224 по центру, преобразующую в тензор (многомерный массив PyTorch: в данном случае трехмерных массив, содержащий цвет, высоту и ширину) и нормализующую его компоненты RGB (красный, зеленый, синий) до заданных среднего значения и стандартного отклонения.

Если мы хотим получить от сети осмысленные ответы, все это должно соответствовать данным, полученным сетью во время обучения.

Процесс выполнения обученной модели на новых данных в сфере глубокого обучения называется *выводом* (inference). Для выполнения вывода необходимо перевести сеть в режим `eval`

```
resnet.eval()
```

Если забыть сделать это, некоторые предобученные модели, например включающие нормализацию по мини-батчам и дропаут, не дадут никаких осмысленных результатов просто по причине их внутреннего устройства. Теперь, после установки режима `eval`, можно выполнять вывод

```
out = resnet(batch_t)
# получается что-то вроде степени уверенности модели в конкретном предсказании
percentage = torch.nn.functional.softmax(out, dim=1)[0] * 100
```

Успешность работы сети во многом зависит от наличия соответствующих объектов в обучающем наборе данных. Если подать нейронной сети нечто выходящее за рамки обучающего набора данных, вполне возможно, что она достаточно уверенно вернет неправильный ответ [1, стр. 64].

Сеть представляет собой всего лишь каркас, а вся суть – в весовых коэффициентах [1, стр. 69].

1.1. Torch Hub

Автору, чтобы опубликовать модель через механизм Torch Hub, необходимо всего лишь поместить файл `hubconf.py` в корневой каталог репозитория GitHub. Структура очень проста

```
# необязательный список модулей, от которых зависит данный код
dependencies = ["torch", "math"]

# одна или несколько функций, открываемых пользователям в качестве входных точек репозитория. Эти
# функции должны инициализировать модели в соответствии с аргументами и возвращать их
def some_entry_fn(*args, **kwargs):
    model = build_some_model(*args, **kwargs)
    return model

def another_entry_fn(*args, **kwargs):
    model = build_another_model(*args, **kwargs)
    return model
```

Теперь интересные предобученные модели можно искать в репозиториях GitHub, содержащих файл `hubconf.py`, зная сразу же, что их можно будет загрузить с помощью модуля `torch.hub`.

```
import torch
from torch import hub

resnet18_model = hub.load(
    "pytorch/vision:main", # название и ветка репозитория GitHub
    "resnet18", # название точки входа
    pretrained=True # ключевой аргумент
)
```

Приведенный код скачивает копию состояния ветки `main` репозитория `pytorch/vision`, вместе с весовыми коэффициентами в локальный каталог (по умолчанию `.torch/hub` в домашнем каталоге) и выполняет функцию точки входа `resnet18`, возвращающую созданный экземпляр модели.

2. Тензор

В контексте глубокого обучения тензоры связаны с обобщением векторов и матриц на произвольную размерность. Другими словами, речь идет о многомерных массивах.

По сравнению с массивами NumPy тензоры PyTorch обладают несколькими потрясающими способностями, например возможностью чрезвычайно быстро выполнять операции на графических процессорах, умением распределять операции по нескольким устройствам или машинам, а также отслеживать породивший их граф вычислений.

Тензоры *PyTorch* и массивы *NumPy* это представления над (обычно) *непрерывными блоками памяти*, содержащими распакованные (unboxed) числовые типы данных Си, а не объекты *Python* [1, стр. 83].

Пример тензора

```
img_t = torch.randn(3, 5, 5)
batch_t = torch.tensor(2, 3, 5, 5) # [батч, каналы, строки, столбцы]
```

Иногда каналы RGB размещаются в измерении 0, а иногда – в измерении 1. Но обобщение можно производить путем отсчета с конца: каналы всегда расположены в измерении -3, третьем с конца.

```
img_gray_naive = img_t.mean(-3)
batch_gray_naive = batch_t.mean(-3)
img_gray_naive.shape, batch_gray_naive.shape # (torch.Size([5, 5]), torch.Size([2, 5, 5]))
```

PyTorch автоматически добавляет в начало измерение размером 1. Эта функция называется *транслированием* (broadcasting). *batch_t* формы (2,3,5,5) умножается на *unsqueeze_weights* формы (3,1,1), в результате чего получается тензор формы (2,3,5,5), в котором затем можно сложить третье измерение с конца (три канала).

```
weights = torch.tensor([0.2126, 0.7152, 0.0722]) # torch.Size([3])
unsqueezed_weights = weights.unsqueeze(-1).unsqueeze(-1) # torch.Size([3, 1, 1])
img_weights = (img_t * unsqueezed_weights)
batch_weights = (batch_t * unsqueezed_weights)
```

В *PyTorch* 1.3 добавилась экспериментальная возможность *именных тензоров*. У функций создания тензоров, например *tensor* и *rand*, есть аргумент *names*. В качестве аргумента *names* должна передаваться последовательность строковых значений

```
weights_named = torch.tensor([0.2126, 0.7152, 0.0722], names=["channels"])
```

При необходимости добавить названия в имеющийся тензор (не меняя существующие) можно вызвать его метод *refine_names*. Аналогично доступу по индексу с помощью многоточия можно пропускать любое количество измерений. С помощью родственного ему метода *rename* можно также переопределять или удалять (путем передачи *None*) уже существующие названия

```
img_named = img_t.refine_names(..., "channels", "rows", "columns")
```

Метод *align_as* возвращает тензор, в котором добавлены недостающие измерения, а уже существующие переставлены в нужном порядке [1, стр. 89]

```
weights_named.shape # torch.Size([3])
img_named.shape # torch.Size([3, 5, 5])
weights_aligned = weights_named.align_as(img_named)
weights_aligned.shape # torch.Size([3, 1, 1])
```

Функции, принимающие на входе аргументы для измерений, также позволяют указывать поименованные измерения

```
(img_named * weights_aligned).sum("channels")
# то же самое
(img_named * weights_aligned).sum(0)
```

При попытке сочетать измерения с различными названиями выдается сообщение об ошибке

```
# img_named[..., :3] то же самое, что и img_named[:, :, :3]
gray_named = (img_named[..., :3] * weights_named).sum("channels") # Ошибка, т.к. размерности не
                           # совпадают
# а так можно
(img_named[:, :, :3] * weights_named).sum("channels")
```

При необходимости использовать тензоры не только в функциях, работающих с поименованными тензорами, необходимо удалить названия, установив их в `None`

```
img_named.rename(None)
```

2.1. Типы элементов тензоров

Использовать стандартные типы данных Python не рекомендуется по нескольким причинам [1, стр. 90]:

- Числовые значения в Python являются объектами. В то время как число с плавающей запятой требует для представления в компьютере только 32 бита, Python преобразует его в *полноценный объект Python* с подсчетом ссылок и т.д. Эта операция, которая называется *упаковкой* (boxing), не является проблемой при хранении небольшого количества числовых значений, но выделять память для миллионов таких объектов – совершенно нерационально.
- Списки в Python предназначены для хранения последовательных наборов объектов. В них нет операций для быстрого вычисления скалярного произведения двух векторов или их суммирования. Кроме того, *списки Python* не оптимизируют размещение своего содержимого в памяти, поскольку представляют собой *наборы указателей на объекты Python* (любые, не только числовые значения) с доступом по индексу. Наконец, списки Python одномерны, и, хотя можно создавать списки списков, это тоже нерационально.
- Интерпретатор Python работает медленно по сравнению с оптимизированным, скомпилированным кодом.

Вычисления в нейронных сетях обычно производятся над 32-битными значениями с плавающей запятой. Более высокая точность, например 64-битные значения, обычно не повышает безошибочность модели, но требует больше памяти и вычислительного времени. Нативная поддержка типа данных с половинной точностью – 16-битных значений с плавающей запятой – в стандартных CPU обычно отсутствует, зато предоставляется современными GPU. При необходимости можно перейти на *половинную точность* для снижения объема занимаемой памяти нейросетевой модели без особого влияния на степень безошибочности [1, стр. 92].

Привести результат функции создания тензора к нужному типу с помощью соответствующего метода приведения типов, можно так

```
torch.zeros(10, 2).double()
torch.ones(10, 2).short()
```

Или с помощью более удобного метода `.to()`

```
torch.zeros(10, 2).to(torch.double)
torch.ones(10, 2).to(dtype=torch.short)
```

Память под значения в тензорах выделяется непрерывными фрагментами памяти под управлением экземпляров `torch.Storage`. Хранилище представляет собой одномерный массив числовых данных, то есть непрерывный фрагмент памяти, содержащий числа заданного типа, например `float` (32-битные значения, выражающие числа с плавающей запятой) или `int64` (64-битные

значения, выражающие целые числа). Экземпляр класса `Tensor` PyTorch – это представление подобного экземпляра `Storage` с возможностью доступа к хранилищу по индексу через указание сдвига и шага по каждому измерению.

Хранилище *всегда* представляет собой *одномерный* массив вне зависимости от размерности каких-либо ссылающихся на него тензоров.

Методы, имена которых заканчиваются на символ «_», как в `zero_`, указывают, что метод работает с заменой на месте (in place), изменяя входные данные вместо того, чтобы создавать новый выходной тензор и возвращать его. Метод `zero_` обнуляет все элементы входного тензора. Все методы, в конце названия которых нет символа подчеркивания, оставляют исходный тензор неизменным и вместо этого возвращают новый.

Для транспонирования двумерных тензоров используется метод `.t()`. Но в PyTorch транспонировать можно не только матрицы. Можно транспонировать многомерный массив, и для этого достаточно указать два измерения, по которым нужно произвести транспонирование (зеркально отражая форму шага)

```
somt_t = torch.ones(3, 4, 5)
transpose_t = some_t.transpose(0, 2)
some_t.shape # torch.Size([3, 4, 5])
transpose_t.shape # torch.Size([5, 4, 3])
```

Любой из тензоров PyTorch можно перенести на (один из) GPU системы для массово-параллельных быстрых вычислений.

Помимо `dtype`, класс `Tensor` предоставляет атрибут `device`, который описывает, где на компьютере размещаются данные тензора.

```
points_gpu = torch.tensor([[4.0, 1.0], [5.0, 3.0], [2.0, 1.0]], device="cuda")
```

Вместо этого можно скопировать созданный в CPU тензор на GPU с помощью метода `to`

```
points_gpu = points.to(device="cuda")
```

При этом возвращается новый тензор с теми же числовыми данными, но хранящийся в *памяти GPU*, а не в *обычной оперативной памяти системы* [1, стр. 105].

Если на нашей машине более одного GPU, можно также указать, на каком именно GPU размещать тензор, передав отсчитываемый с нуля целочисленный номер GPU на машине, вот так

```
points_gpu = points.to(device="cuda:0")

# умножение выполняется на CPU
points = 2 * points
# умножение выполняется на GPU
points_gpu = 2 * points.to(device="cuda")
```

Отметим, что тензор `points_gpu` не передается обратно в CPU после вычисления результата. Вот что происходит в этой строке:

- Тензор `points` копируется в GPU.
- Выделяется память в GPU под новый тензор, в котором будет храниться результат умножения.
- Возвращается обращение к этому GPU-тензору.

Следовательно, если мы прибавим к результату константу

```
points_gpu = points_gpu + 4
```

операция сложения будет по-прежнему производиться в GPU и никакой информации в CPU передаваться не будет (если мы не будем выводить полученный тензор на экран или обращаться к нему). Для переноса тензора обратно в CPU необходимо указать в методе `.to()` аргумент `cpu`

```
points_gpu.to(device="cpu")
```

Можно также для получения того же результата воспользоваться сокращенными методами `.cpu()` и `.cuda()` вместо метода

```
points_gpu = points.cuda()
points_gpu = points.cuda(0)
points_gpu.cpu()
```

Стоит упомянуть, что с помощью метода `.to()` можно менять тип данных и их место размещения одновременно, указав в качестве аргументов `device` и `dtype`.

Тензоры PyTorch можно очень эффективно преобразовать в массивы NumPy и наоборот. Благодаря этому можно воспользоваться огромными объемами функциональности экосистемы Python, основанной на типах массивов NumPy. Подобная совместимость с массивами NumPy, не требующая копирования кода, возможна благодаря работы системы хранения с буферным протоколом Python.

В разреженных тензорах хранятся только ненулевые значения, а также информация об индексах.

Для сериализации объектов-тензоров PyTorch использует «за кулисами» `pickle`, а также специализированный код сериализации для хранилища. Вот как можно сохранить наш тензор `points` в файл `ourpoints.t`

```
torch.save(points, "./data/p1ch3/ourpoints.t")
```

Либо можно передать файловый дескриптор файла вместо названия

```
with open("./data/p1ch3/ourpoints.t", mode="wb") as f:
    torch.save(points, f)
```

Загрузка тензора `points` обратно также выполняется одной строкой кода

```
points = torch.load("./data/p1ch3/ourpoints.t")
```

что эквивалентно

```
with open("./data/p1ch3/ourpoints.t", mode="rb") as f:
    points = torch.load(f)
```

И хотя подобным образом можно быстро сохранять тензоры, если нужно загрузить их только в PyTorch, сам по себе формат файла не отличается совместимостью: прочитать тензор с помощью какого-либо еще ПО, помимо PyTorch, не получится.

HDF5 – переносимый, широко поддерживаемый формат представления сериализованных многомерных массивов, организованный в виде вложенного ассоциативного массива типа «ключ–значение». Python поддерживает формат HDF5 благодаря библиотеке `h5py`, принимающей и возвращающей данные в виде массивов NumPy.

```
import h5py
```



```
f = h5py.File("./ourpoints.hdf5", "w")
dset = f.create_dataset("coords", data=points.numpy())
f.close()
```

Здесь "coords" – это ключ для файла в формате HDF5. В HDF5 интересна возможность индексации набора данных на диске и обращения только к нужным нам элементам.

```
f = h5py.File("./ourpoints.hdf5", "r")
dset = f["coords"]
torch.from_numpy(dset[-2:])
f.close()
```

2.2. Представление реальных данных с помощью тензоров

2.2.1. Работа с изображениями

Изображения представляются в виде набора скалярных значений расположенных на равномерной сетке с высотой и шириной (в пикселях), например, по одному скалярному значению на каждую точку сетки (пиксель) для изображения в оттенках серого или несколько скалярных значений на каждую точку сетки для представления различных цветов.

Отражающие значения для различных пикселей скаляры обычно кодируются 8-битными целыми числами, как в бытовых фотоаппаратах. В медицинских, научных и промышленных приложениях нередко встречается более высокая точность, например 12- или 16-битная, для расширения диапазона или повышения чувствительности в случаях, когда пиксель отражает информацию о физическом свойстве, например о плотности костной ткани, температуре или глубине.

Загрузить изображение можно так

```
import imageio

img_arr = imageio.imread("./bobby.jpg")
img_arr.shape # (720, 1280, 3)
```

Модули PyTorch, работающие с изображениями, требуют от тензоров измерений $C \times H \times W$ (каналы, высота и ширина).

Для получения нужной нам схемы расположения можно воспользоваться методом `permute` тензора, указав в качестве параметров старые измерения для каждого из новых.

```
img = torch.from_numpy(img_arr)
out = img.permute(2, 0, 1)
```

Эта операция не копирует данные тензора, вместо этого `out` *использует то же самое хранилище*, что и `img`, только меняя информацию о размере и шаге на уровне тензора.

Несколько более эффективная альтернатива использованию для создания тензора `stack` – выделить заранее память под тензор нужного размера, а затем заполнить его загруженными из каталога изображениями следующим образом

```
batch_size = 3
batch = torch.zeros(batch_size, 3, 256, 256, dtype=torch.uint8)
```

Батч будет состоять из трех RGB-изображений по 256 пикселей высотой и 256 пикселей шириной.

Нейронные сети демонстрируют наилучшее качество обучения, когда входные данные находятся в диапазоне примерно от 0 до 1 или от -1 до 1.

Никаких принципиальных различий между тензорами, содержащими объемные пространственные данные и данные изображения, нет. Просто появляется дополнительное измерение, глубина, вслед за измерением каналов, и получается 5-мерный тензор формы $N \times C \times D \times H \times W$.

Загрузим пример КТ-снимка с помощью функции `volread` из модуля `imageio`, принимающий в качестве аргумента каталог и собирающей все файлы в формате DICOM (Digital Imaging and Communications in Medicine) в трехмерный массив NumPy

```
import imageio

dir_path = ".../2-LUNG 3.0B70f-04083"
vol_arr = imageio.volread(dir_path, "DICOM")
vol_arr.shape # (99, 512, 512)

vol = torch.from_numpy(vol_arr).float()
vol = torch.unsqueeze(vol, 0)
vol.shape # torch.Size([1, 99, 512, 512]): каналы, глубина, высота, ширина
```

При вызове метода `.view()` на тензоре возвращает новый тензор с другими размерностью и шагами *без изменения хранилища*. Это позволяет перегруппировать тензор практически без затрат, поскольку *никакие данные копировать не нужно*.

Нормализацию данных к отрезку $[0; 1]$ или $[-1; 1]$ желательно производить *для всех количественных величин*, таких как «температура» (это полезно для процесса обучения) [1, стр. 137].

По завершении обучения алгоритм способен генерировать правильные выходные сигналы при получении новых данных, *достаточно схожих* со входными данными, на которых он обучался. В случае глубокого обучения этот процесс работает даже тогда, когда входные данные и требуемые выходные сигналы *далеки* друг от друга: когда они относятся к различным предметным областям.

У нас есть модель с неизвестными значениями параметров и нужно получить оценку этих параметров, которая бы минимизировала расхождение между предсказанными выходными сигналами и измеренными значениями (ошибка). Целью процесса оптимизации должен быть поиск таких параметров модели, которые минимизировали бы функцию потерь.

Отдельная итерация обучения, во время которой обновляются параметры для всех обучающих примеров данных, называется *эпохой*.

Градиенты по параметрам модели должны быть одного порядка [1, стр. 168].

Аргумент `requires_grad` указывает PyTorch отслеживать целое семейство тензоров. Если функции дифференцируемые (как большинство операций над тензорами PyTorch), величина производной будет автоматически занесена в атрибут `.grad`.

Количество тензоров с параметром `requires_grad`, установленным в `True` аргументом, и композиции функций может быть любым. В этом случае PyTorch вычисляет производные функции потерь по всей цепочке функций (графу вычислений) и накапливает их значения в атрибутах `.grad` этих тензоров (узлы этого графа).

При вызове `.backward()` производные *накапливаются* в узлах-листьях. *Необходимо явным образом обнулять градиенты после обновления параметров на их основе*. Так что, если `backward` вызывался ранее, потери оцениваются опять, `backward` вызывается снова, после чего накапливаются градиенты во всех листьях графа, то есть суммируются с вычисленными на предыдущей итерации, в результате чего неправильное значение градиента [1, стр. 173].

Чтобы предотвратить подобное, необходимо *явным образом обнулять градиенты* на каждой итерации.

```
def training_loop(n_epochs, learning_rate, params, t_u, t_c):
```

```

for epoch in range(1, n_epochs + 1):
    if params.grad is not None:
        params.grad.zero_()

    t_p = model(t_u, *params)
    loss = loss_fn(t_p, t_c)
    loss.backward() # выполняем обратный проход и вычисляем градиенты

    with torch.no_grad():
        params -= learning_rate * params.grad

    if epoch % 500 == 0:
        print(...)

```

При вызове `loss.backward()` PyTorch обходит граф в обратном порядке, вычисляя градиенты. Контекст `torch.no_grad()` означает, что механизм автоматического вычисления градиента игнорирует внутренности блока `with`: то есть не добавляет ребра в граф прямого прохода.

У каждого оптимизатора доступны два метода: `zero_grad` и `step`. Метод `zero_grad` обнуляет атрибут `grad` всех передаваемых оптимизатору параметров при его создании. А метод `step` обновляет значения параметров в соответствии с реализуемой конкретным оптимизатором стратегией оптимизации.

```

def training_loop(n_epochs, optimizer, params, t_u, t_c):
    for epoch in range(1, n_epochs + 1):
        t_p = model(t_u, *params)
        loss = loss_fn(t_p, t_c)

        # обнуляем градиент, потому как в противном случае
        # производные накапливались бы в узлах-листьях графа вычислений
        optimizer.zero_grad()
        loss.backward() # обратный проход по графу; вычисляем градиенты
        optimizer.step() # обновляем параметры модели

        if epoch % 500 == 0:
            print(...)

    return params

```

Очень гибкая модель с большим количеством параметров стремится к минимизации функции потерь в точках данных и нет никаких гарантий, что она будет себя вести нужным образом *вдали* или *между* точками данных [1, стр. 180].

Потери на обучающем наборе данных показывают, можно ли вообще подогнать нашу модель к этому обучающему набору данных - другими словами, достаточны ли *разрешающие возможности* (capacity) этой модели для обработки содержащейся в данных информации [1, стр. 182].

Глубокая нейронная сеть потенциально может аппроксимировать очень сложные функции при условии достаточно большого числа нейронов, а значит, и параметров. Чем меньше параметров, тем проще должна быть форма функции, чтобы наша сеть смогла ее аппроксимировать. Итак, правило 1: если потери на обучающем наборе данных не уменьшаются, вероятно, модель слишком проста *для имеющихся данных*.

Что ж, если вычисленная на проверочном наборе данных функция потерь не убывает вместе с обучающим набором, значит, наша модель обучается лучше аппроксимировать полученные во время обучения примеры данных, но не *обобщается* на примеры данных, которые не входят в

этот конкретный набор. Правило 2: если потери на обучающем и проверочном наборах данных расходятся – модель переобучена.

С интуитивной точки зрения более простая модель может описывать обучающие данные не так хорошо, как более сложная, но зато, вероятно, будет вести себя более равномерным образом между точками данных.

Следовательно, процесс выбора правильного размера нейросетевой модели в смысле количества параметров основан на двух шагах:

- увеличение размера до тех пор, пока модель не будет хорошо подогнана к данным,
- а затем уменьшение, пока не будет устранено переобучение.

Разбиение набора данных. Перетасовка элементов тензора эквивалентна перестановке его индексов – как раз то, что делает функция `randperm`

```
n_samples = t_u.shape[0]
n_val = int(0.2 * n_samples)

shuffled_indices = torch.randperm(n_samples)
train_indices = shuffled_indices[:-n_val]
test_indices = shuffled_indices[-n_val:]

train_t_u = t_u[train_indices]
train_t_c = t_c[train_indices]

val_t_u = t_u[val_indices]
val_t_c = t_c[val_indices]
```

Осталось только дополнительно вычислять потери на проверочном наборе данных *на каждой эпохе*, чтобы заметить переобучение

```
def training_loop(n_epochs, optimizer, params, train_t_u, val_t_u, train_t_c, val_t_c):
    for epoch in range(1, n_epochs + 1):
        train_t_p = model(train_t_u, *params)
        train_loss = loss_fn(train_t_p, train_t_c)

        val_t_p = model(val_t_u, *params)
        val_loss = loss_fn(val_t_p, val_t_c)

        optimizer.zero_grad()
        # здесь только train_loss.backward(), поскольку мы не хотим обучать
        # модель на проверочном наборе данных
        train_loss.backward()
        optimizer.step()

        if epoch <= 3 or epoch % 500 == 0:
            print(...)

    return params
```

Запуск

```
params = torch.tensor([1.0, 0.0], requires_grad=True)
learning_rate = 1e-2
optimizer = optim.SGD([params], lr=learning_rate)

training_loop(
    n_epochs = 3000,
    optimizer = optimizer,
    params = params,
```

```

train_t_u = train_t_un,
val_t_u = val_t_un,
train_t_c = train_t_c,
val_t_c = val_t_c,
)

```

Наша цель – убедиться, что убывают как потери на обучающем наборе данных, *так* и потери на проверочном [1, стр. 186].

Вопрос: раз мы никогда не вызываем `backward()` для `val_loss`, зачем вообще формировать граф вычислений? Можно просто вызывать `model` и `loss_fn` как обычные функции, без отслеживания истории вычислений.

```

def training_loop(n_epochs, optimizer, params, train_t_u, val_t_u, train_t_c, val_t_c):
    for epoch in range(1, n_epochs + 1):
        train_t_p = model(train_t_u, *params)
        train_loss = loss_fn(train_t_p, train_t_c)

        with torch.no_grad():
            val_t_p = model(val_t_u, *params)
            val_loss = loss_fn(val_t_p, val_t_c)
            # Убеждаемся, что для нашего вывода аргумент requires_grad == False
            assert val_loss.requires_grad == False

        optimizer.zero_grad()
        train_loss.backward()
        optimizer.step()

```

Нейрон – по сути представляет собой линейное преобразование входного сигнала (например, умножение входного сигнала на какое-либо число (*вес*) и прибавление к нему константы *смещения*) с последующим применением фиксированной нелинейной функции (*функции активации*).

Функция активации играет две выжные роли [1, стр. 195]:

- Во внутренних частях модели благодаря функции активации возможны различные наклоны графика выходного сигнала в разных значениях – нечто, по определению *недоступное для линейной функции*. Искусно сочетая эти по-разному наклоненные участки для различных выходных сигналов, нейронные сети могут аппроксимировать любые функции.
- На последнем слое сети она локализует выходные сигналы предыдущей линейной операции в заданном интервале.

Нейрон – это просто линейная функция с последующей функцией активации.

Функции активации по определению [1, стр. 199]:

- *нелинейны* – сколько ни применяй преобразование вида $w \cdot x + b$ без функции активации, все равно получится функция той же самой (аффинной линейной) формы. *Нелинейность позволяет сети в целом аппроксимировать более сложные функции*,
- *дифференцируемы*, что дает возможность вычисления градиентов. Точечные разрывы (Hardtanh, ReLU etc.), допустимы.

В отсутствие этих характеристик сеть либо превратиться обратно в линейную модель, либо с трудом будет поддаваться обучению.

Для функций активации справедливо следующее:

- Имеется по крайней мере один диапазон *чувствительности*, внутри которого нетривиальные изменения входного сигнала приводят к соответствующим нетривиальным изменениям выходного. Необходимо для обучения.

- У многих из них есть также диапазон *нечувствительности* (*насыщения*), в котором изменения входного сигнала практически не приводят к изменениям выходного.

В целом получается механизм, обладающий большими возможностями: при получении на входе различных данных в сети, составленной из *линейных* и *активационных блоков*:

- различные нейроны могут возвращать для одних входных сигналов результаты, относящиеся к различным диапазонам,
- соответствующие этим входным сигналам ошибки в основном влияют на нейроны, работающие в диапазоне чувствительности, а на остальные блоки процесс обучения практически не влияет.

В результате объединения множества линейных и активационных блоков параллельно и последовательно получается математический объект, способный аппроксимировать сложные функции. Различные сочетания нейронов реагируют в различных диапазонах на входные сигналы, причем параметры этих блоков можно довольно легко оптимизировать посредством градиентного спуска, поскольку процесс обучения напоминает обучение линейной функции, вплоть до момента насыщения выходного сигнала.

В PyTorch есть отдельный подмодуль, посвященный нейронным сетям, – `torch.nn`. Он включает «кирпичики», необходимые для создания всех видов нейросетевых архитектур. В терминологии PyTorch эти «кирпичики» называются *модулями* (в других фреймворках подобные стандартные блоки часто называются *слоями* (layers)). Модуль PyTorch – это класс Python, наследующий базовый класс `nn.Module`.

Подмодули должны быть атрибутами верхнего уровня, а не быть закопаны внутри экземпляров `list` или `dict`! В противном случае оптимизатор не сможет их найти (а значит, и их параметры). На случай, если модели потребуется список или ассоциативный массив подмодулей, в PyTorch есть классы `nn.ModuleList` и `nn.ModuleDict`. У `nn.Module` есть подкласс `nn.Linear`, применяющий ко входным сигналам аффинное преобразование.

У всех подклассов `nn.Module` в PyTorch есть метод `__call__`, позволяющий создавать экземпляры `nn.Linear` и вызывать их как функции следующим образом

```
import torch.nn as nn

linear_model = nn.Linear(1, 1)
linear_model(t_un_val)
```

Вызов экземпляра `nn.Module` с набором инструментов приводит к вызову метода `forward` с теми же аргументами, который реализует *прямой проход* вычислений, в то время как `__call__` выполняет другие немаловажные операции до и после вызова `forward`. Так что формально можно вызвать `forward` напрямую, и он вернет тот же результат, что и `__call__`, но делать это из пользовательского кода не рекомендуется

```
y = model(x)    # Правильно!
y = model.forward(x)  # НЕправильно! Так делать не надо!!!
```

Конструктор `nn.Linear` принимает три аргумента: число входных признаков, число выходных признаков и булево значение, указывающее, включает линейная модель смещение или нет

```
import torch.nn as nn

linear_model = nn.Linear(1, 1)
linear_model(t_un_val)
```

Все модули в `nn` ориентированы на генерацию выходных сигналов сразу для батча из нескольких входных сигналов. Следовательно, если нам нужно выполнить `nn.Linear` для десяти примеров данных, можно создать входной тензор размеров $B \times N_{\text{вх}}$, где B – размер батча, а $N_{\text{вх}}$ – число входных признаков, и пропустить его один раз через модель.

Причины для организации данных по батчам многогранны. Одна из них – желание полноценно загрузить вычислениями имеющиеся вычислительные ресурсы. В частности, GPU позволяют сильно распараллеливать вычисления, так что при одиночном входном сигнале для маленькой модели большинство вычислительных элементов будет простаивать.

Еще одно преимущество в том, что некоторые развитые модели способны использовать статистическую информацию по целому батчу, и эти статистические показатели будут точнее при большом размере батча.

```
linear_model = nn.Linear(1, 1)
optimizer = optim.SGD(
    linear_model.parameters(),
    lr=1e-02,
)
```

При вызове метода `training_loss.backward()` в листьях графа вычислений накапливаются градиенты. При вызове `optimizer.step()` программа проходит по всем объектам `Parameter` и меняет их на соответствующую содержимому атрибута `grad` долю.

```
def training_loop(
    n_epochs,
    optimizer,
    model,
    loss_fn,
    t_u_train,
    t_u_val,
    t_c_train,
    t_c_val,
):
    for epoch in range(1, n_epochs + 1):
        t_p_train = model(t_u_train)
        loss_train = loss_fn(t_p_train, t_c_train)

        t_p_val = model(t_u_val)
        loss_val = loss_fn(t_p_val, t_c_val)

        # требуется обязательно обнулять градиент на каждой итерации
        optimizer.zero_grad()
        # выполняем проход в обратном направлении и вычисляем градиенты
        loss_train.backward()
        # обновляем параметры модели
        optimizer.step()
```

Модуль `nn` включает несколько распространенных функций потерь, одна из которых – `nn.MSELoss`.

Модуль `nn` предоставляет удобный способ соединения модулей цепочкой с помощью контейнера `nn.Sequential`

```
seq_model = nn.Sequential(
    nn.Linear(1, 13),
    nn.Tanh(),
    nn.Linear(13, 1)
)
```

Модель переходит от одного входного признака до 13 скрытых признаков, пропускает их через функцию активации `Tanh` и, наконец, объединяет получившиеся 13 чисел в один выходной признак.

Названия модулей в `Sequential` представляют собой просто порядковые номера модулей в списке аргументов. Что любопытно, `Sequential` также принимает на входе `OrderedDict`, в котором можно указать название каждого из передаваемых `Sequential` модулей

```
from collections import OrderedDict

seq_model = nn.Sequential(OrderDict([
    ("hidden_linear", nn.Linear(1, 8)),
    ("hidden_activation", nn.Tanh()),
    ("output_linear", nn.Linear(8, 1))
]))

for name, param in seq_model.named_parameters():
    print(name, param.shape)
# output
hidden_linear.weight torch.Size([8, 1])
hidden_linear.bias torch.Size([8])
output_linear.weight torch.Size([1, 8])
output_linear.bias torch.Size([1])
```

Обращаться к конкретным объектам `Parameter` можно путем указания подмодулей в качестве атрибутов

```
seq_model.output_linear.bias
# output
Parameter containing:
tensor([0.1402], requires_grad=True)
```

Можно посмотреть *градиенты параметра weight* линейной части скрытого слоя. Запускаем цикл обучения для новой модели нейронной сети, после чего смотрим на получившиеся градиенты после последней эпохи

```
seq_model.hidden_linear.weight.grad
```

Функции активации, в дополнение к линейным преобразованиям, позволяют нейронным сетям аппроксимировать сильно нелинейные функции, оставляя их при этом достаточно простыми для оптимизации [1, стр. 215].

Чтобы преобразовать изображение PIL в тензор PyTorch, можно воспользоваться модулем `torchvision.transforms`. Есть преобразование `ToTensor`, превращающее массивы NumPy и изображения PIL в тензоры. Оно также располагает измерения выходного тензора в порядке $C \times H \times W$ (каналы, высота, ширина).

```
from torchvision import transforms, datasets

cifar10_train = datasets.CIFAR10(data_path, train=True, download=True)
img, label = cifar10_train[99]

to_tensor = transforms.ToTensor()
img_t = to_tensor(img)
img_t.shape # torch.Size([3, 32, 32])
```

Изображение `img` превратилось в тензор формы $3 \times 32 \times 32$, то есть в изображение размером 32×32 с тремя цветовыми каналами (RGB).

Это преобразование можно передать непосредственно в виде аргумента `datasets.CIFAR10`

```
tensor_cifar10 = datasets.CIFAR10(
    data_path,
    train=True,
    download=False,
    transform=transforms.ToTensor(),
)
```

Рекомендуемая практика – нормализовать набор данных до нулевого среднего значения и единичного стандартного отклонения по каждому из каналов.

При выборе функций активации, линейных около нуля (± 1 или ± 2), ограничение данных тем же диапазоном повышает вероятность *ненулевых градиентов нейронов*, а значит, и *ускоряет обучение*. Кроме того, нормализация каналов к одинаковому распределению гарантирует смещение и обновление информации из разных каналов (посредством градиентного спуска) с одинаковой скоростью обучения.

Чтобы обеспечить нулевое среднее и единичное стандартное отклонение по каждому из каналов, необходимо вычислить среднее значение и стандартное отклонение каждого из каналов набора данных и применить следующее преобразование $v_n[c] = (v[c] - \text{mean}[c]) / \text{stdev}[c]$.

Поскольку набор данных CIFAR-10 невелик, можно работать с ним полностью в оперативной памяти. Разместим все возвращаемые объектом `Dataset` тензоры последовательно в дополнительном измерении

```
imgs = torch.stack([img_t for img_t, _ in tensor_cifar_10], dim=3)
imgs.shape # torch.Size([3, 32, 32, 50 000])
```

Теперь можно легко вычислить поканальные средние значения

```
imgs.view(3, -1) # torch.Size([3, 5120000])
imgs.view(3, -1).mean(dim=1)
imgs.view(3, -1).std(dim=1)
```

Теперь

```
transformed_cifar10 = datasets.CIFAR10(
    data_path,
    train=True,
    download=False,
    transform=transforms.Compose([
        transforms.ToTensor(),
        transforms.Normalize(
            mean=imgs.view(3, -1).mean(dim=1),
            std=imgs.view(3, -1).std(dim=1),
        )
    ])
)
```

Сколько признаков содержит каждый пример данных? Так, $3 \times 32 \times 32$ равняется 3072 входных признака на каждый пример. Получаем новую модель `nn.Linear` с 3072 входными признаками и некоторым количеством скрытых признаков, за которым следует функция активации, а затем еще один `nn.Linear`, сокращающий модель до соответствующего количества выходных признаков (в данном случае 2)

```
import torch.nn as nn

n_out = 2
```

```

model = nn.Sequential(
    nn.Linear(
        3072, # входные признаки
        512, # размер скрытого слоя
    ),
    nn.Tanh(),
    nn.Linear(
        512, # размер скрытого слоя
        n_out, # выходные признаки
    )
)

```

Нейронной сети требуется по крайней мере один скрытый слой (активации, поэтому два модуля) с *нелинейностью* между слоями, чтобы сеть могла *усваивать произвольные функции*, в противном случае модуль будет просто *линейной*.

Здесь необходимо понять, что выходной сигнал носит категориальный характер: птица или самолет. Для представления категориальной величины следует воспользоваться унитарным кодированием, например, [1, 0] для самолета и [0, 1] для птицы (порядок выбран произвольно). Такая схема подходит и в случае 10 классов, как в полном наборе данных CIFAR-10; просто вектор будет длиной 10.

В нашем частном случае бинарной классификации (птица или самолет; по сути птица или не птица) два значения – избыточно, поскольку одно всегда равно 1 минус второе. И действительно, PyTorch позволяет выдавать на выходе одно значение *вероятности*, получая вероятность путем использования в конце модели функции активации `nn.Sigmoid` и функции потерь на основе *бинарной перекрестной энтропии* `nn.BCELoss`. Существует также `nn.BCELossWithLogits`, объединяющая эти два шага [1, стр. 229].

В идеальном случае сеть должна выдавать на выходе `torch.tensor([1.0, 0.0])` для самолета и `torch.tensor([0.0, 1.0])` – для птицы. На практике же, поскольку наш классификатор не будет идеален, следует ожидать от сети неких промежуточных значений. Главное в этом случае, что мы можем интерпретировать выходные сигналы как вероятности: первая запись – вероятность класса "airplane", а вторая – "bird".

Чтобы развернуть изображение формы $3 \times 32 \times 32$, а затем преобразовать его в вектор-строку делаем так

```

img_batch = img.view(-1).unsqueeze(0)
img_batch.shape # torch.Size([1, 3072]); 3 * 32 * 32 = 3072

```

Что касается функции потерь для задач классификации, то здесь необходимо накладывать штраф на ошибки классификации, а не кропотливо штрафовать все, что не равно в точности 0.0 или 1.0, поэтому в задачах классификации плохо работает среднеквадратическая функция потерь.

В этом случае *необходимо максимизировать вероятность*, соответствующую истинному классу.

Другими словами, нам нужна функция потерь, принимающая очень высокие значения, когда правдоподобие (истинность параметров нашей модели при имеющихся данных) низко: настолько низко, что вероятности альтернативных вариантов выше. И наоборот, потери должны быть низкими, когда правдоподобие данного варианта выше, чем у альтернатив, и мы не хотим заикливаться на доведении вероятности до 1 [1, стр. 234].

Действующая подобным образом функция потерь существует и называется *отрицательной логарифмической функцией правдоподобия* (negative log likelihood, NNL).

Для каждого примера данных в батче мы делаем следующее:

1. Производим прямой проход и получаем выходные значения из последнего (линейного) слоя.
2. Вычисляем для них многомерную логистическую функцию и получаем вероятности.
3. Извлекаем предсказанную вероятность для истинного класса (правдоподобие параметров).

Отметим, что истинный класс известен, поскольку обучение производится с учителем, – это наши эталонные данные.

4. Вычисляем ее логарифм, ставим перед ним знак «минус» и прибавляем к потерям.

Функция `nn.LogSoftmax()` обеспечивает численную устойчивость вычислений

```
model = nn.Sequential(
    nn.Linear(3072, 512),
    nn.Tanh(),
    nn.Linear(512, 2),
    nn.LogSoftmax(dim=1),
)
```

Среднеквадратическая функция потерь (MSE) насыщается намного раньше и – что принципиально – для совершенно неправильных предсказаний. Основная причина состоит в том, что уклон MSE слишком мал, чтобы компенсировать пологость многомерной логистической функции активации для неправильных предсказаний. **Поэтому MSE для вероятностей плохо подходит для задач классификации** [1, стр. 237].

```
import torch
import torch.nn as nn

model = nn.Sequential(
    nn.Linear(3072, 512),
    nn.Tanh(),
    nn.Linear(512, 2),
    nn.LogSoftmax(dim=1),
)

learning_rate = 1e-2

optimizer = optim.SGD(model.parameters(), lr=learning_rate)

n_epochs = 100

for epoch in range(n_epochs):
    for img, label in cifar2:
        out = model(img.view(-1).unsqueeze(0))
        loss = loss_fn(out, torch.tensor([label]))

        optimizer.zero_grad()
        loss.backward()
        optimizer.step()

    print(...)
```

Мы поняли, что обработка всех 10 000 изображений одним батчем – это перебор, так что решили создать внутренний цикл, чтобы *обрабатывать по одному примеру данных за раз* и производить *обратное распространение ошибки по этому одному примеру*.

Перетасовывая примеры данных на каждой эпохе и вычисляя градиент по одному или (что желательно из соображений устойчивости) нескольким примерам данных за раз, мы фактически вносим элемент случайности в алгоритм градиентного спуска. Оказывается, что следование градиентам, вычисленным по мини-батчам, которые представляют собой лишь слабые аппроксимации градиентов, вычисленных по всему набору данных, улучшает сходимость и предотвращает «застывание» процесса оптимизации во встреченных по пути локальных минимумах [1, стр. 238].

Обычно размер мини-батча представляет собой константу, задаваемую до обучения, аналогично скорости обучения.

Модуль `torch.utils.data` включает класс, помогающий с перетасовкой и организацией данных по мини-батчам: `Dataloader`. Задача загрузчика данных состоит в выборе мини-батчей из набора данных с гибкими возможностями использования различных стратегий выборки. Одна из самых распространенных стратегий: равномерная выборка после *перетасовки* данных в каждой эпохе.

```
train_loader = torch.utils.data.DataLoader(cifar2, batch_size=64, shuffle=True)
```

Пример

```
import torch
import torch.nn as nn

train_loader = torch.utils.data.DataLoader(
    cifar2, batch_size=64,
    shuffle=True,
)

model = nn.Sequential(
    nn.Linear(3072, 512),
    nn.Tanh(),
    nn.Linear(512, 2)
    nn.LogSoftmax(dim=1)
)

learning_rate = 1e-2

optimizer = optim.SGD(model.parameters(), lr=learning_rate)

# отрицательный логарифм правдоподобия
loss_fn = nn.NLLLoss()

n_epochs = 100

for epoch in range(n_epochs):
    for imgs, labels in train_loader:
        batch_size = imgs.shape[0]
        outputs = model(imgs.view(batch_size, -1))
        loss = loss_fn(outputs, labels)

        optimizer.zero_grad()
        loss.backward()
        optimizer.step()
```

На каждой итерации внутреннего цикла `imgs` представляет собой мини-батч 64 RGB-изображений (размером 32×32), а `labels` – тензор размером 64 с индексами меток.

Наша цель в том, чтобы правильно присвоить изображениям метки классов, причем желательно на независимом наборе данных

```
val_loader = torch.utils.data.DataLoader(
    cifar2_val,
    batch_size=64,
    shuffle=False, # NB
)

correct = 0
total = 0

with torch.no_grad(): # NB!
    for imgs, labels in val_loader:
        batch_size = imgs.shape[0]
        outputs = model(imgs.view(batch_size, -1))
        _, predicted = torch.max(outputs, dim=1)
        total += labels.shape[0]
        correct += int((predicted == labels).sum())

print(correct / total)
```

Достаточно часто последний слой `nn.LogSoftmax` не включают в сеть, используя в качестве функции потерь `nn.CrossEntropyLoss`.

Полносвязная сеть не является *инвариантной относительно сдвига*. Это значит, что сеть, обученная распознавать «Спитфайер» (самолет), начинающийся с позиции 4,4, не сможет распознать *тот же самый* «Спитфайер», начинающийся с позиции 8,8. Нам пришлось бы дополнять (augment) набор данных, то есть применять случайные сдвиги к изображениям во время обучения, чтобы сеть могла заметить «Спитфайер» в любом месте изображения, причем это пришлось бы делать для всех изображений в наборе [1, стр. 246].

Учтите, что *весовые коэффициенты ядра заранее не известны*, а инициализируются *случайным* образом и обновляются посредством обратного распространения ошибки. Отметим также, что для всего изображения используется одно и то же ядро, а это значит, что и весовые коэффициенты ядра [1, стр. 252]. В производную функции потерь по сверточным весам вносят свой вклад все изображение.

Свертка эквивалентна нескольким линейным операциям, весовые коэффициенты которых равны нулю практически везде, кроме окрестностей отдельных пикселей, и получают одинаковые обновления во время обучения.

Очень часто применяют ядра, размеры которых одинаковы по всем измерениям, поэтому в PyTorch есть сокращенная форма записи для них: `kernel_size=3` для двумерной свертки означает форму 3×3 , для трехмерной свертки – форму $3 \times 3 \times 3$.

```
conv = nn.Conv2d(3, 16, kernel_size=3)
# или так
conv = nn.Conv2d(3, 16, kernel_size=(3, 3))
```

Здесь 3 признака на один пиксель (3 канала; RGB) и 16 выходных каналов.

В результате прохода двумерной свертки получается двумерное изображение, пиксели которого представляют собой взвешенную сумму значений по локальным окрестностям входного изображения.

Как начальные значения весовых коэффициентов ядра `conv.weight`, так и смещения задаются случайным образом, так что выходное изображение особого смысла не имеет.

Лучше придерживаться ядер нечетных размеров; ядра четного размера встречаются редко [1, стр. 256].

Задача сверточной сети состоит в оценке ядра набора фильтров в последовательных слоях, преобразующих многоканальное изображение в другое многоканальное изображение, в котором различные каналы соответствуют разным признакам (например, один канал – для среднего значения, другой – для вертикальных краев и т.д.).

Первый набор ядер работает с маленькими окрестностями низкоуровневых признаков первого порядка, а второй набор фактически работает с более широкими окрестностями, *генерируя признаки, представляющие собой композицию предыдущих признаков*. Благодаря этому замечательному механизму сверточные нейронные сети способны анализировать очень сложные кадры.

```
# свертка + нелинейность + пулинг
model = nn.Sequential(
    nn.Conv2d(3, 16, kernel_size=3, padding=1), # свертка
    nn.Tanh(), # нелинейность
    nn.MaxPool2d(2), # пулинг
    nn.Conv2d(16, 8, kernel_size=3, padding=1),
    nn.Tanh(),
    nn.MaxPool2d(2),
    ...
)
```

Первая операция свертки превращает три канала RGB в 16, благодаря чему у сети появляется возможность генерировать *16 независимых признаков* (16 каналов). Далее мы применяем функцию активации *Tanh*. Полученное 16-канальное изображение 32×32 субдискретизируется первым слоем *nn.MaxPool2d* до 16-канального изображения 16×16 .

Теперь субдискретизированное изображение подвергается еще одной операции свертки, выдающей на выходе 8-канальный выходной сигнал 16×16 . Если повезет, это выходное изображение будет состоять из высокоуровневых признаков. И опять же мы применяем функцию активации *Tanh*, после чего производим субдискретизацию до 8-канального выходного изображения 8×8 .

После уменьшения входного изображения до набора 8×8 признаков можно надеяться вернуть из сети значения вероятностей, подходящих для подачи на вход отрицательной логарифмической функции правдоподобия. Нужно преобразовать 8-канальное изображение 8×8 в одномерный вектор и завершить нашу сеть набором полносвязанных слоев

```
model = nn.Sequential(
    nn.Conv2d(3, 16, kernel_size=3, padding=1),
    nn.Tanh(),
    nn.MaxPool2d(2),
    nn.Conv2d(16, 8, kernel_size=3, padding=1),
    nn.Tanh(),
    nn.MaxPool2d(2),
    ... # Пропущен важный момент
    nn.Linear(8 * 8 * 8, 32),
    nn.Tanh(),
    nn.Linear(32, 2),
)
```

Для повышения разрешающих возможностей модели можно *увеличить количество выходных каналов* сверточных слоев (то есть число *признаков*, генерируемых каждым из сверточных слоев), в результате чего увеличится и размер линейного слоя.

PyTorch позволяет производить в модели любые вычисления путем создания подклассов `nn.Module`.

Готовые или пользовательские свертки как подмодули обычно включаются в программу посредством описания в конструкторе `__init__` и присваивания их `self` для использования в функции `forward`. Их параметры в то же время хранятся в них на протяжении всего жизненного цикла нашего модуля. Обратите внимание, что перед этим необходимо вызывать `super().__init__()`.

3. Обобщения с помощью сверток

3.1. Сеть как подкласс `nn.Module`

Задача классификационных сетей обычно заключается в сжатии информации в том смысле, что мы начинаем с изображения, содержащего значительное количество пикселей, и сжимаем его в вектор вероятностей классов.

Для создания подкласса `nn.Module` как минимум необходимо описать функцию `forward`, принимающую входные сигналы модуля и возвращающую выходной.

Подмодули должны быть атрибутами верхнего уровня, а не быть «закопаны» внутри экземпляров `list` или `dict`! В противном случае оптимизатор не сможет их (а значит, и их параметры) найти. На случай, если модели потребуется список или ассоциативный массив подмодулей, в PyTorch есть классы `nn.ModuleList` и `nn.ModuleDict` [1, стр. 269].

3.2. Функциональные API

В PyTorch есть *функциональные* аналоги для всех модулей `nn`. Под функциональными здесь подразумевается «без внутреннего состояния» – другими словами, «выходное значение которых целиком и полностью определяется значениями входных аргументов».

```
import torch.nn.functional as F

class Net(nn.Module):
    def __init__(self):
        super().__init__()
        self.conv1 = nn.Conv2d(3, 16, kernel_size=3, padding=1)
        self.conv2 = nn.Conv2d(16, 8, kernel_size=3, padding=1)
        self.fc1 = nn.Linear(8 * 8 * 8, 32)
        self.fc2 = nn.Linear(32, 2)

    def forward(self, x):
        out = F.max_pool2d( # пулинг
            torch.tanh( # нелинейность
                self.conv1(x) # свертка
            ), 2)
        out = F.max_pool2d(
            torch.tanh(
                self.conv2(x)
            ), 2)
        out = out.view(-1, 8 * 8 * 8)
        out = torch.tanh(self.fc1(out))
        out = self.fc2(out)

        return out
```


3.3. Обучение модели

В основе сверточной сети лежат два вложенных цикла: внешний – по *эпохам*, а внутренний – на основе объекта `DataLoader`, генерирующего батчи из объекта `Dataset`. На каждой итерации цикла необходимо [1, стр. 273]:

- Пропустить входные сигналы через модель (прямой проход).
- Вычислить функцию потерь (также часть прямого прохода).
- Обнулить все старые градиенты.
- Вызвать `loss.backward()` для вычисления градиентов функции потерь относительно каждого из параметров (обратный проход).
- Оптимизировать в сторону уменьшения потерь.

```
import datetime

def training_loop(n_epochs, optimizer, model, loss_fn, train_loader):
    for epoch in range(1, n_epochs + 1):
        loss_train = 0.0
        for imgs, labels in train_loader:
            outputs = model(imgs)
            loss = loss_fn(outputs, labels)

            # избавляемся от градиентов с предыдущих итераций
            optimizer.zero_grad()
            # выполняем обратный проход; то есть вычисляем градиенты
            # по всем обучаемым параметрам сети
            loss.backward()
            # обновляем модель
            optimizer.step()
            # суммируем потери за эту эпоху
            loss_train += loss.item()

        if epoch == 1 or epoch % 10 == 0:
            # получаем средние потери на батч
            print("{} Epoch {}, Training loss {}".format(
                datetime.datetime.now(), epoch,
                loss_train / len(train_loader))
            )
```

Обучение в течение 100 эпох

```
# Объект DataLoader организует примеры данных из нашего набора по батчам.
# Перетасовка обеспечивает случайный порядок примеров данных из набора
train_loader = torch.utils.data.DataLoader(cifar2, batch_size=64, shuffle=False)

model = Net()
optimizer = optim.SGD(model.parameters(), lr=1e-02)
loss_fn = nn.CrossEntropyLoss()

training_loop(
    n_epochs=100,
    optimizer=optimizer,
    model=model,
    loss_fn=loss_fn,
    train_loader=train_loader,
)
```

Сохранить модель можно так

```
# сохраняются только веса
torch.save(model.state_dict(), data_path + "birds_vs_airplanes.pt")
```

Файл `birds_vs_airplanes.pt` теперь содержит все параметры объекта `model`: весовые коэффициенты и смещения для двух модулей свертки и двух линейных модулей. Никакой структуры, только весовые коэффициенты. Это значит, что при развертывании модели в реальных условиях нам понадобится описание класса `model`

```
loaded_model = Net()
loaded_model.load_state_dict(torch.load(data_path + "birds_vs_airplanes.pt"))
```

В `nn.Module` есть реализована функция `.to`, перемещающая все параметры в GPU (или приводящая тип данных, если передать ей аргумент `dtype`).

Между `Module.to` и `Tensor.to` существует тонкое различие. `Module.to` производит операции с заменой на месте, то есть изменяет экземпляр модуля. А `Tensor.to` – нет, возвращая *новый тензор*.

Рекомендуемой практикой является создание экземпляра `Optimizer` *после* перемещения всех параметров на нужное устройство [1, стр. 275].

Перенос вычислений на GPU при его наличии считается хорошим стилем программирования. Неплохим паттерном программирования будет установка значения переменной `device` в зависимости от `torch.cuda.is_available`:

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

```
def training_loop(
    n_epochs, optimizer,
    model, loss_fn, train_loader
):
    for epoch in range(1, n_epochs + 1):
        loss_train = 0.0
        for imgs, labels in train_loader:
            imgs = imgs.to(device=device) # NB
            labels = labels.to(device=device) # NB
            outputs = model(imgs)
            loss = loss_fn(outputs, labels)
```

Можно создать экземпляры модели и перенести ее на `device`

```
train_loader = torch.utils.data.DataLoader(cifar2, batch_size=64, shuffle=True)

# перенести модель (все ее параметры) на GPU
model = Net().to(device=device)
optimizer = optim.SGD(model.parameters(), lr=1e-02)
loss_fn = nn.CrossEntropyLoss()

training_loop(
    n_epochs=100,
    optimizer=optimizer,
    model=model,
    loss_fn=loss_fn,
    train_loader=train_loader,
)
```

Если забыть перенести саму модель или входные данные на GPU, вы получите сообщения об ошибках, указывающие, что тензоры располагаются на различных устройствах, поскольку операторы PyTorch не поддерживают смеси входных данных на GPU и CPU.

PyTorch попытается загрузить веса на то же устройство, с которого они были сохранены, то есть весовые коэффициенты с GPU будут восстановлены на GPU. Лаконичным вариантом будет потребовать от PyTorch переопределить информацию об устройстве при загрузке весовых коэффициентов

```
loaded_model = Net().to(device=device)
loaded_model.load_state_dict(
    torch.load(
        data_path + "birds_vs_airplanes.pt",
        map_location=device
    )
)
```

3.4. Ширина сети

Ширина сети это количество нейронов в слое или каналов на каждую операцию свертки. Расширить модель в PyTorch очень легко. Необходимо просто указать большее количество выходных каналов в первой свертке и увеличивать следующие слои соответствующим образом, не забывая менять функцию `forward` так, чтобы отразить увеличившуюся длину вектора при переходе на полносвязные слои.

Чем больше разрешающих возможностей модели, тем с большей степенью изменчивости входных сигналов сможет справиться модель, но в то же время тем выше вероятность переобучения, поскольку модель сможет воспользоваться дополнительными параметрами для запоминания несущественных аспектов входных данных.

3.5. Регуляризация

Основные способы регуляризации:

- штрафы на весовые коэффициенты (L_1 , L_2 -норма),
- дропаут,
- нормализация по батчам (альтернатива дропауту).

Первый способ достижения устойчивости обобщения: добавления члена регуляризации в формулу потерь. Этот дополнительный член ограничивает рост весовых коэффициентов модели в процессе обучения: он устроен так, что они стремятся оставаться маленькими. Другими словами, он налагает штраф на большие значения весов. В результате форма функции потерь становится более гладкой, и для модели нет особого смысла подстраиваться под отдельные примеры данных.

Наиболее популярные виды членов регуляризации: L_2 -регуляризация и L_1 -регуляризация. L_2 -регуляризацию также называют *затуханием весов* (weight decay). Прибавление к функции потерь члена L_2 -регуляризации эквивалентно уменьшению каждого весового коэффициента пропорционально его текущему значению во время шага оптимизации (отсюда и название «затухание веса»). Обратите внимание, что затухание веса относится ко всем параметрам сети, в том числе и к смещениям.

```
def training_loop_l2reg(n_epochs, optimizer, model, loss_fn, train_loader):
    for epoch in range(1, n_epochs + 1):
        loss_train = 0.0
```

```

for imgs, labels in train_loader:
    imgs = imgs.to(device)
    labels = labels.to(device=device)
    outputs = model(imgs)
    loss = loss_fn(outputs, labels)

    l2_lambda = 0.001
    l2_norm = sum(p.pow(2.0).sum() for p in model.parameters())
    loss = loss + l2_lambda * l2_norm

    optimizer.zero_grad()
    loss.backward()
    optimizer.step()

    loss_train += loss.item()

```

Впрочем, в оптимизаторе SGD в PyTorch уже есть параметр `weight_decay`, соответствующий $2 * \lambda$, который напрямую осуществляет затухание весов во время их обновления. Он полностью эквивалентен прибавлению L_2 -нормы весовых коэффициентов к функции потерь без необходимости накопления в функции потерь и вовлечения автоматического вычисления градиентов.

Идея дропаута действительно проста: обнуляем случайную часть выходных сигналов нейронов по сети, причем этот случайный выбор производится на каждой итерации обучения.

Фактически в результате этой процедуры на каждой итерации формируются слегка отличающиеся модели с различными топологиями нейронов, уменьшая шансы нейронов модели скоординироваться в процессе запоминания, что происходит при переобучении.

В PyTorch можно реализовать дропаут в модели с помощью добавления модуля `nn.Dropout` между нелинейной функцией активации и линейным или сверточным модулем последующего слоя. В качестве аргумента необходимо указать вероятность, с которой будут обнуляться входные сигналы.

```

class NetDropout(nn.Module):
    def __init__(self, n_chans1=32):
        super().__init__()
        self.n_chans1 = n_chans1
        self.conv1 = nn.Conv2d(3, n_chans1, kernel_size=3, padding=1)
        self.conv1_dropout = nn.Dropout2d(p=0.4)
        self.conv2 = nn.Conv2d(n_chans1, n_chans1 // 2, kernel_size=3, padding=1)
        self.conv2_dropout = nn.Dropout2d(p=0.4)
        self.fc1 = nn.Linear(8 * 8 * n_chans1 // 2, 32)
        self.fc2 = nn.Linear(32, 2)

    def forward(self, x):
        out = F.max_pool2d(torch.tanh(self.conv1(x)), 2)
        out = self.conv1_dropout(out) # Dropout
        out = F.max_pool2d(torch.tanh(self.conv2(out)), 2)
        out = self.conv2_dropout(out) # Dropout
        out = self.view(-1, 8 * 8 * self.n_chans1 // 2)
        out = torch.tanh(self.fc1(out))
        out = self.fc2(out)

    return out

```

Дропаут обычно происходит во время обучения, в то время как во время использования обученной модели в реальных условиях модуль дропаута обходят или, что эквивалентно, присваи-

вают равную нулю вероятность. Этот процесс контролируется свойством `train` модуля `Dropout` [1, стр. 283].

PyTorch позволяет переключаться между двумя режимами, вызывая `model.train()` или `model.eval()` для любого подкласса `nn.Module`.

В статье Сергея Йоффе «Нормализация по батчам: ускорение обучения нейронных сетей путем сокращения внутреннего ковариантного сдвига» (Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift) описывается методика, которая позволяет *повысить скорость обучения и снизить зависимость обучения от начальных значений*, а также играет роль *регуляризатора*, тем самым представляя альтернативу дропауту [1, стр. 283].

Основная идея *нормализации по батчам* состоит в нормализации входных сигналов функций активации сети так, чтобы получить определенное желательное распределение для мини-батчей. Если вспомнить внутренние механизмы обучения и роль нелинейных функций активации, становится ясно, что это помогает *избежать чрезмерного углубления входных сигналов функций активации в область насыщения*, что гасит градиенты и замедляет обучение [1, стр. 283].

На практике нормализация по батчам сдвигает и масштабирует промежуточные входные сигналы на основе среднего значения и стандартного отклонения, вычисленных в этой промежуточной точке по примерам данных мини-батча. Эффект от регуляризации основан на том, что отдельные примеры данных и следующие далее по конвейеру функции активации всегда рассматриваются моделью как сдвинутые и нормализованные, в зависимости от статистических показателей выделенного случайным образом мини-батча. Авторы статьи полагают, что нормализация по батчам исключает или по крайней мере сокращает необходимость в дропауте.

Поскольку цель нормализации по батчам – масштабировать входные сигналы функции активации, логично будет производить ее после линейного преобразования (свертки а данном случае), то есть перед функцией активации

```
class NetBatchNorm(nn.Module):
    def __init__(self, n_chans1=32):
        super().__init__()
        self.n_chans1 = n_chans1
        self.conv1 = nn.Conv2d(3, n_chans1, kernel_size=3, padding=1)
        self.conv1_batchnorm = nn.BatchNorm2d(num_features=n_chans1)
        self.conv2 = nn.Conv2d(n_chans1, n_chans1 // 2, kernel_size=3, padding=1)
        self.conv2_batchnorm = nn.BatchNorm2d(num_features=n_chans1 // 2)
        self.fc1 = nn.Linear(8 * 8 * n_chans1 // 2, 32)
        self.fc2 = nn.Linear(32, 2)

    def forward(self, x):
        out = self.conv1_batchnorm(self.conv1(x)) # batch norm
        out = F.max_pool2d(torch.tanh(out))
        out = self.conv2_batchnorm(self.conv2(out)) # batch norm
        out = out.view(-1, 8 * 8 * self.n_chans1 // 2)
        out = torch.tanh(self.fc1(out))
        out = self.fc2(out)

    return out
```

Как и дропаут, нормализация по батчам должна вести себя по-разному во время обучения и во время выполнения вывода. На самом деле *во время выполнения вывода* желательно, чтобы выходной сигнал для конкретного входного сигнала не зависел от прочих входных сигналов, подаваемых на вход модели. А это значит, что необходим способ нормализовать данные, но при этом *раз и навсегда зафиксировать параметры нормализации*.

При обработке мини-батчей, помимо оценки среднего значения и стандартного отклонения для текущего мини-батча, PyTorch также обновляет в качестве приближения *скользящие оценки среднего значения и стандартного отклонения*, репрезентативные для всего набора данных.

Таким образом, если пользователь указывает `model.eval()` и модель содержит модуль нормализации по батчам, скользящие оценки фиксируются и используются для нормализации. Для возврата к использованию статистических показателей мини-батчей мы вызываем `model.train()` точно так же, как для дропаута.

Глубина обеспечивает возможность работы сети с *иерархической информацией*, когда для анализа какого-либо входного сигнала необходимо понимать *контекст* [1, стр. 285].

Углубление сети обычно ухудшает сходимость. Производные функции потерь относительно параметров, особенно в первых слоях, приходится умножать на множество других чисел, генерируемых цепочкой операций между функцией потерь и параметром. Эти множители могут быть маленькими, приводя в результате к еще меньшим числам, или большими, поглощая маленькие числа из-за приближенности операций с плавающей запятой. В сухом остатке мы получаем, что в результате длинной цепочки операций умножение *вклад отдельного параметра в градиент исчезает*, и это ведет к *неэффективному обучению* данного слоя, поскольку ни этот, ни другие параметры не будут обновляться должным образом.

Добавление в модель обходной связи (skip connection) наподобие ResNet сводится к прибавлению выходного сигнала какого-то слоя к входному сигналу другого слоя

```
class ResNet(nn.Module):
    def __init__(self, n_chans1=32):
        super().__init__()
        self.n_chans1 = n_chans1
        self.conv1 = nn.Conv2d(3, n_chans1, kernel_size=3, padding=1)
        self.conv2 = nn.Conv2d(n_chans1, n_chans1 // 2, kernel_size=3, padding=1)
        self.conv3 = nn.Conv2d(n_chans1 // 2, n_chans1 // 2, kernel_size=3, padding=1)
        self.fc1 = nn.Linear(4 * 4 * n_chans1 // 2, 32)
        self.fc2 = nn.Linear(32, 2)

    def forward(self, x):
        out = F.max_pool2d(torch.relu(self.conv1(x)), 2)
        out = F.max_pool2d(torch.relu(self.conv2(out)), 2)
        out1 = out
        out = F.max_pool2d(torch.relu(self.conv3(out)) + out1, 2) # skip connection
        out = out.view(-1, 4 * 4 * self.n_chans1 // 2)
        out = self.fc2(out)

        return out
```

Обходная связь создает прямой путь от расположенных глубоко параметров к функции потерь, благодаря чему они вносят более непосредственный вклад в градиент функции потерь, ведь частные производные функции потерь по этим параметрам теперь получают шанс не умножаться на коэффициент в длинной цепочке прочих операций.

Отмечается, что обходные связи благотворно влияют на сходимость, особенно на начальных этапах обучения. Кроме того, поверхность функции потерь глубоких остаточных сетей намного глаже, чем у сетей прямого распространения той же глубины и ширины.

Применение обходных связей в ResNet дало возможность успешно обучать модели глубиной более 100 слоев.

`nn.Sequential` гарантирует, что выходной сигнал одного блока будет использован как входной сигнал следующего, а также что все параметры блока видимы `Net`.

```
class ResBlock(nn.Module):
    def __init__(self, n_chans):
        super().__init__()
        # Слои BatchNorm свел бы на нет эффект смещения, так что смещение обычно опускают
        self.conv = nn.Conv2d(n_chans, n_chans, kernel_size=3, padding=1, bias=False)
        self.batch_norm = nn.BatchNorm(n_num_features=n_chans)
        torch.nn.init.kaiming_normal_(self.conv.weight, nonlinearity="relu")
        torch.nn.init.constant_(self.batch_norm.weight, 0.5)
        torch.nn.init.zeros_(self.batch_norm.bias)

    def forward(self, x):
        out = self.conv(x)
        out = self.batch_norm(out)
        out = torch.relu(out)

        return out + x

class NetResDeep(nn.Module):
    def __init__(self, n_chans1=32, n_blocks=10):
        super().__init__()
        self.n_chans1 = n_chans1
        self.conv1 = nn.Conv2d(3, n_chans1, kernel_size=3, padding=1)
        self.resblocks = nn.Sequential(
            *(n_blocks * [ResBlock(n_chans=n_chans1)]))
        self.fc1 = nn.Linear(8 * 8 * n_chans1, 32)
        self.fc2 = nn.Linear(32, 2)

    def forward(self, x):
        out = F.max_pool2d(torch.relu(self.conv1(x)), 2)
        out = self.resblocks(out)
        out = F.max_pool2d(out, 2)
        out = out.view(-1, 8 * 8 * self.n_chans1)
        out = torch.relu(self.fc1(out))
        out = self.fc2(out)

        return out
```

У регуляризации весов и дропаута статистическая интерпретация в качестве методов регуляризации более строгая, чем нормализация по батчам. Нормализация по батчам предназначена скорее для *улучшения сходимости* [1, стр. 292].

Список литературы

1. *Стивенс Э.* PyTorch. Освещаая глубокое обучение. – СПб.: Питер, 2022. – 576 с.