



## Dive into DeepLearning - Chapters 1 → 5

### Introduction

#### Supervised Problem

1. The data we can learn from
2. A model of how to transform the data
3. An objective function that quantifies how well (or badly) the model is doing
4. An algorithm to adjust the model's parameters to optimize the objective function

### Preliminaries

#### → Data manipulation

`torch.arange(12)`

`x.shape`

`torch.numel(M)` = total number of elements in M

`x.sum(dim=1, keepdim=True)` = keeps dimension of vector unitário

`x.numpy()`

`torch.from_numpy(M)`

`torch.mv(A, x)` = matrix vector multiplication

`torch.mm(A, B)` = matrix matrix multiplication

`torch.norm()`

#### → Linear Algebra

`B = A.clone()`

#### → Probability

`from torch.distributions import multinomial`

`My playlist today: #`

Meow  
Meow



## Linear Regression

$$g = \mathbf{x}w + b \text{ linear model}$$

Loss function = a measure of fitness

$$L(w, b) = \frac{1}{2} (\hat{y} - y)^2$$

When training we want to find the parameters  $(w^*, b^*)$  that minimize total loss across all training examples

$$w^*, b^* = \underset{w, b}{\operatorname{argmin}} L(w, b)$$

- \*Analytic solution: too much restrictive, computational cost
- \*Minibatch

The most naive application of gradient descent consist of taking the derivative of the loss function, which is an average of all losses computed of every ~~single~~ update  $\rightarrow$  too slow often sample a minibatch

$$w = w - \eta \nabla L(w)$$

Minibatch  $\beta$

$$(w, b) \leftarrow (w, b) - \eta \sum_{i \in \beta} \nabla_{(w, b)} L^{(i)}(w, b)$$

Concise implementation

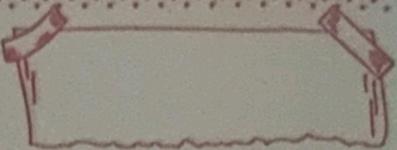
from torch import nn

nn.Sequential = pipeline of nn networks

FORONI: nn.Linear(2, 1)

in dimension  $\rightarrow$  out dimension

We can do it!



`net[i]` → layer i

`net[0].weight`, `net[0].bias`

or `MSELoss()`

trainer = `torch.optim.SGD`

## Softmax regression

Regression = how much / how many?

Classification = which one

classification ↗ hard assignment: categories only

↪ soft assignment: also class probs

classification treated as one-hot encoding

$$y \in \{(1,0,0), (0,1,0), (0,0,1)\}$$

dog      cat      chicken  
 $y_1$        $y_2$        $y_3$

Three logits  $O_1, O_2, O_3$  for each input

$$O_1 = x_1 w_1 + \dots + x_n w_n + b_1$$

$$O_2 = x_1 w_1 + \dots + x_n w_n + b_2$$

$$O_3 = x_1 w_1 + \dots + x_n w_n + b_3$$

$$O = w \cdot x + b$$

Just as linear regression softmax is a single layer neural network and fully connected

You might be tempted to suggest that we interpret the logits "O" directly as an output



of interest. However, there are some problems with directly interpreting the output of linear layers as probabilities; they also don't sum up to 1, can be negative. This violates the axioms of probability.

$$\hat{y} = \text{softmax}(o) \text{ where } \hat{y}_j = \frac{\exp(o_j)}{\sum_k \exp(o_k)}$$

$$\operatorname{argmax}_{\hat{y}_j} \hat{y}_j = \operatorname{argmax}_j o_j$$

Note that the softmax operation does not change the ordering among the logits  $o$ , which are simply the pre-softmax values that determine the probabilities assigned to each class.

Vectorization of minibatch  $X$

$$O = XW + b$$

$$\hat{Y} = \text{softmax}(O)$$

Loss function = log likelihood

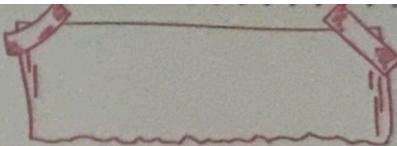
$$P(Y, X) = \prod_{i=1}^n P(y^{(i)} | x^{(i)})$$

Maximum likelihood = minimize negative likelihood

$$l(y, \hat{y}) = - \sum_{j=1}^n y_j \log \hat{y}_j$$

The image classification dataset

One of the widely used for image classification is the



MNIST dataset. While it had a good run as a benchmark dataset, even simple models by today's standards achieve classification accuracy over 95%, making it unsuitable for distinguishing between stronger models and weaker ones. Today, MNIST serves as more of aanity check than benchmark.

## Multilayer Perceptron

Linearity implies the weaker assumption of monotonicity  
↑ feature ↑ model output

We can overcome those limitations of linear models and handle a more general class of functions by incorporating one or more hidden layers.

> Multilayer Perceptron - MLP

Stack many fully connected layers one top of each other

$$H = \sigma(XW^{(1)} + b^{(1)}) \rightarrow \text{activation function}$$
$$O = HW^{(2)} + b^{(2)}$$

1. ReLU: Rectified Linear Unit

$$\text{ReLU}(x) = \max(x, 0)$$

2. sigmoid( $x$ ) =  $\frac{1}{1 + \exp(-x)}$

Meow

Meow



$$3. \tanh(x) = \frac{1 - \exp(-2x)}{1 + \exp(-2x)}$$

or. Parameter

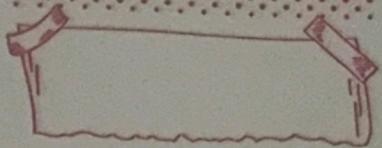
`Torch.zeros_like(matrix)`

underfit  $\downarrow$  train eval  $\downarrow$  test eval  
overfit  $\uparrow$  train eval  $\downarrow$  test eval

Weight decay

Limiting the number of the features is a popular technique to reduce overfitting.

Weight decay (commonly called L<sub>2</sub> regularization)



## Dropout

Faced with more features than examples, linear models tend to overfit. But given more examples than features, we can generally count on linear models not overfit.

Unfortunately, the readability with which linear models generalize comes at cost. Naively applied, linear models do not take into account interactions among features. For every feature, a linear model must assign either a positive or a negative weight, ignoring context.

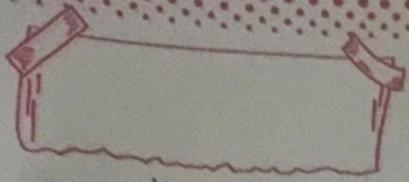
In traditional texts, this fundamental tension between generalizability and flexibility is described as bias-variance trade-off. Linear models have high bias, they can only represent a small class of functions. However, these models have low variance: they give similar results across different samples of data.

Deep neural networks inhabit the opposite end of the bias-variance spectrum. Unlike linear models, neural network are not confined to looking each feature individually. Even when we have far more examples than features, neural network are capable of overfitting.

Another useful notion of simplicity is smoothness, i.e., that function should not be sensitive to small change to its input.

My playlist today: #





smoothness  $\rightarrow$  dropout = support noise in examples  
in Dropout

## Propagation and computational graphs

Forward propagation (or forward pass) refers to the calculation and storage of intermediate (including outputs) for a neural network in order from the input layer to the output layer.

Plotting computational graphs help us visualize the dependencies of operations and variables within calculations.

Backpropagation refers to the method of calculating the gradient of neural network parameters. In short, the method traverses the network in reverse order from the output to the input layer, according to the chain rule in calculus.

## Layers and blocks

in state dict

in init.constant(matrix, a)

in init.zeros'(matrix)