

Lecture 3⁻¹⁻

Learning the whole dataset

So far we train neural networks to model one (single) training example (a data set)

(input₀, goal₀).

First consider the simple network with one input and one output

In fact the learning process is reduced to memorizing data:
the weight₀ (w_0) is selected to fit the prediction₀, pred₀ to the given goal₀.

$$\text{weight}_0 = \text{goal}_0 / \text{input}_0,$$

$$\Rightarrow \text{pred}_0 := \text{input}_0 * \text{weight}_0 = \text{goal}_0$$

-2-

Next time when the network gets input $:=$ input₀ it returns the correct predict value goal₀.

Still if this network see new input value inp, its prediction can be not accurate at all.

A real training (learning) is based on the larger (whole) dataset (info we have from experiments)

input = [input₀, input₁, ..., input_n]

goal = [goal₀, goal₁, ..., goal_n]

The weight is obtained by solving the error equation (the loss function)

$$w_0 = \arg \min error(w) \quad (*)$$

where

$$error(w) = \frac{1}{2} \sum_{j=0}^{M-1} (inp_j \times w - goal_j)^2$$

The optimal weight w_0 is a solution of the equation

$$\frac{\partial error(w)}{\partial w} = 0$$

i.e.

$$\sum_{j=0}^M (inp_j \times w - goal_j) \times inp_j = 0$$

$$\Rightarrow w_0 = \left(\sum_{j=0}^M goal_j \times inp_j \right) / \sum_{j=0}^M (inp_j)^2 \quad (2)$$

-4-

In Machine Learning (ML)
an equivalent and more general
relations are defined

$$\text{pred}_j = \alpha + w \times \text{inp}_j \quad (\alpha - \text{bias})$$

This relation is called
a linear regression model.

In the deep learning method
more general networks are used.
Then instead formula (2) the
optimal value of the weight
is determined by using the
gradient descent algorithm
(and efficient modifications
of it!).

Multiple inputs

Let the network has N inputs,
we denote the j -th data set

$$\text{input}^j = [\text{inp}_0^j, \text{inp}_1^j, \dots, \text{inp}_{N-1}^j]$$

$j = 0, 1, \dots, M-1,$

$$\text{goal} = [\text{goal}^0, \text{goal}^1, \dots, \text{goal}^{M-1}]$$

$$\text{weight} = [w_0, w_1, \dots, w_{N-1}]$$

$$\text{predict} = [\text{pred}^0, \text{pred}^1, \dots, \text{pred}^{M-1}]$$

The pred^j is calculated as

$$\text{pred}^j = \sum_{n=0}^{N-1} \text{inp}_n^j \times w_n.$$

The loss function (error)

$$\text{error} = \frac{1}{2} \sum_{j=0}^{M-1} (\text{pred}^j - \text{goal}^j) \times \times 2.$$

$$\text{delta}^j = \text{pred}^j - \text{goal}^j$$

If bias coefficient is included into prediction formula

$$\text{pred}^j = \sum_{n=0}^{N-1} \text{inp}_n^j \times W_n + d,$$

then the loss function depends also on one more parameter

$$\text{parameters} = [\text{weight}, d].$$

In gradient descent step we

use also the derivative with respect to d (∇ error depends also on $\frac{\partial \text{error}}{\partial d}$)

$$\frac{\partial \text{error}}{\partial d} = \sum_{j=0}^{M-1} (\text{pred}^j - \text{goal}^j) \times \frac{\partial \text{pred}^j}{\partial d}$$

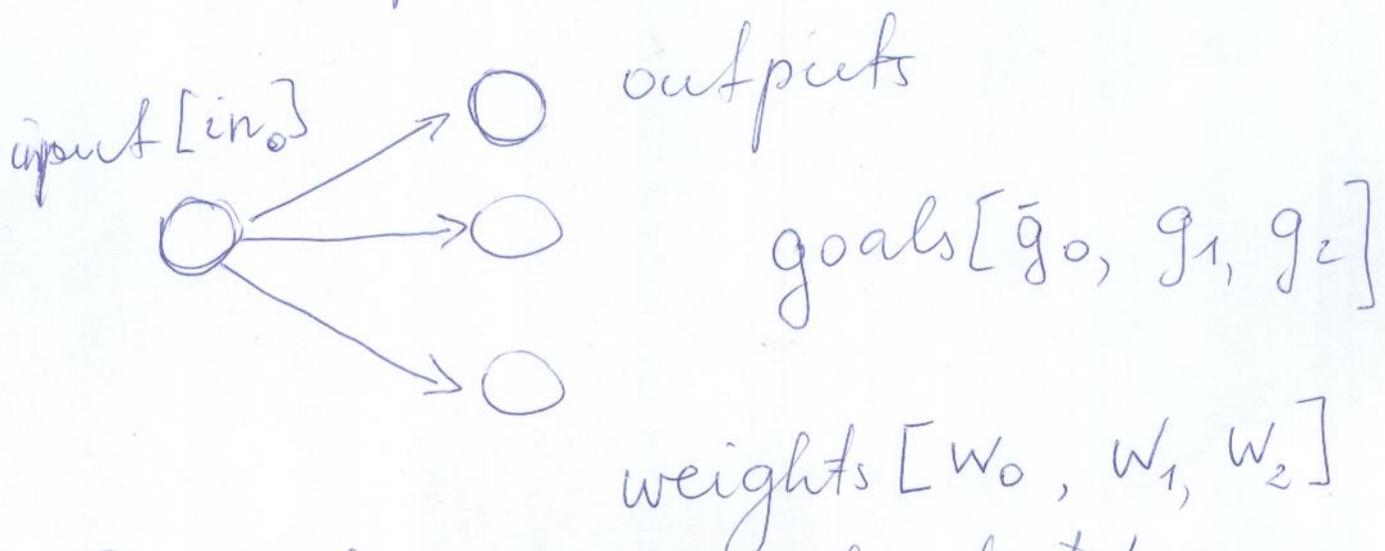
$$= \sum_{j=0}^{M-1} (\text{pred}^j - \text{goal}^j)$$

$$= \sum_{j=0}^{M-1} \text{delta}^j$$

Multiple outputs

For simplicity we start with a network fore

- one single input data
- multiple data for outputs



Predictions are calculated as

$$\boxed{\text{pred}_j = w_j \times in_0} \quad (1)$$

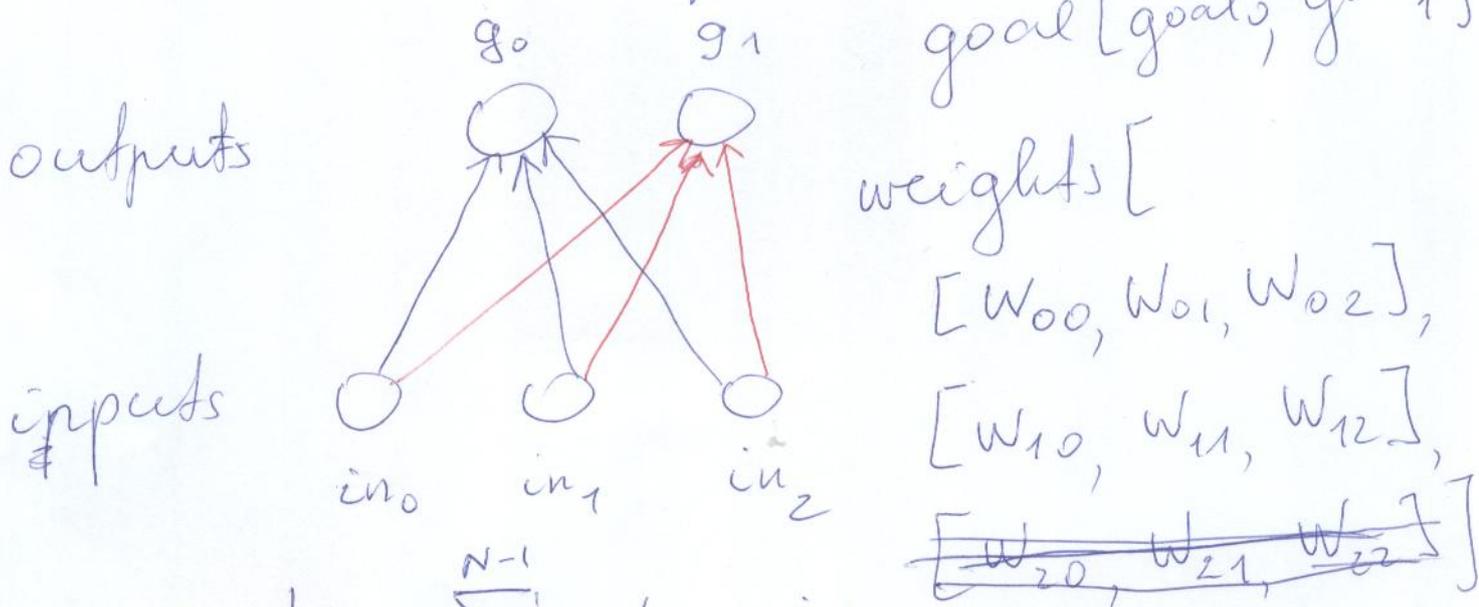
$$\text{delta}_j = \text{pred}_j - g_j$$

We make two main notes:

- Three networks are independent
- Gradient descent method can be applied in parallel for each network.

Generalization

1. The ANN can be considered for multiple inputs and multiple outputs



$$pred_k = \sum_{n=0}^{N-1} W_{kn} * in_n$$

inputs

$k = 0, 1$ - outputs

A real training is based on the larger datasets.

- Let us assume that we have one set of inputs (dataset)

$$\text{input} = [\text{inp}_0, \text{inp}_1, \dots, \text{inp}_{N-1}]$$

$N = \#$ of inputs

- We want to predict K outputs.

The training is done on a dataset of size $M = \#$ size of the dataset

$$\text{goals} = [\text{goal}^0, \text{goal}^1, \dots, \text{goal}^{M-1}]$$

$$\text{goal}^m = [g_0^m, g_1^m, \dots, g_{K-1}^m]$$

$$m = 0, \dots, M-1.$$

Predictions

$$\text{predictions} = [\text{predict}_0, \dots, \text{predict}_{K-1}]$$

$$\text{predict}_k = \sum_{n=0}^{N-1} w_{k,n} \text{inp}_n, \quad k=0, \dots, K-1.$$

Our task is to make a training and to find weights of ANN

$$W = [W_0, \dots, W_{k-1}] \quad K = \# \text{ outputs}$$

$$W_k = [w_{k0}, w_{k1}, \dots, w_{k,N-1}] \quad N = \# \text{ inputs}$$

We have K independent ANN for each output and the training of each ANN can be done in parallel.

The loss function error_k(W_k):

$$\text{error}_k(W_k) = \frac{1}{2} \sum_{m=0}^{M-1} \left(\sum_{n=0}^{N-1} \text{inp}_n^m \times w_{kn} - \text{goal}_k^m \right)^2$$

The optimal weights W_k are calculated by using gradient descent algorithm.

if inputs are also different for each m

$$\text{error}_k(W_k) = \frac{1}{2} \sum_{m=0}^{M-1} \left(\sum_{n=0}^{N-1} \text{inp}_n^m \times w_{kn} - \text{goal}_k^m \right)^2$$


Full, batch and stochastic gradient descent

It was noted that the case of multiple outputs can be resolved as K independent trainings of separate ANN.

Thus we restrict to the case of multiple N inputs and one output

The size of the dataset is equal to M :

$$\text{goals} = [\text{goal}^0, \text{goal}^1, \dots, \text{goal}^{M-1}]$$

$$\text{inputs} = [\text{input}^0, \text{input}^1, \dots, \text{input}^{M-1}]$$

$$\text{input}^m = [\text{inp}_0^m, \text{inp}_1^m, \dots, \text{inp}_{N-1}^m]$$

$$m = 0, \dots, M-1.$$

We want to find weights

$$W = [w_0, w_1, \dots, w_{N-1}]$$

They are used to compute predictions

$$\text{pred}^m = \sum_{n=0}^{N-1} \text{inp}_n^m * w_n$$

The loss function error(W) is defined as

$$\text{error}(W) = \frac{1}{2} \sum_{m=0}^{M-1} (\text{pred}^m - \text{goal}^m)^2$$

We have investigated the

Full gradient descent version

An entire dataset is used at each iteration of iteration: the average weight-delta is calculated over the entire data set and only then the weights are changed

Batch gradient descent

This modification of GDM that takes an approach laying in between of two previous modifications

Instead of updating the weights after just one example (as in *stochastic gradient descent*) or after the entire dataset (after all M examples of data) we choose a batch size (between 8 and 256) of examples, after which the weights are updated

- 13 -

Weight updates:

$$W^{s+1} = W^s - \alpha \Delta W^s,$$

where weight-delta is calculated as

$$\Delta W^s = \left(\sum_{m=0}^{M-1} \left(\sum_{n=0}^{N-1} \text{inp}_n^m \times W_n^s - \text{goal}^m \right) \text{inp}_n^m \right)$$
$$= \sum_{m=0}^{M-1} \underbrace{\text{delta}^m \times \text{inp}^m}_{\text{dot product}}$$

Stochastic gradient descent

It performs a prediction and weight update for each training example separately (for each $m=0, \dots, M-1$).