

## Supervised Learning

Definition of *Supervised Learning*:

(Goodfellow et.al. *Deep Learning Book*, p. 105)

Supervised learning involves observing several examples  $X = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  of a random vector  $\mathbf{x}$  and associated values  $Y = \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(m)}\}$  of a random vector  $\mathbf{y}$ , and learning to predict  $\mathbf{y}$  from  $\mathbf{x}$ , usually by estimating  $p(\mathbf{y}|\mathbf{x})$ .

## Cost Function

In order to train the desired behavior of a machine learning model with a set of parameters  $\theta$  it is important to define the right *cost function*, as the gradient descent algorithm will minimize this function. The cost function  $J(\theta)$  computes a *cost value*  $c$  dependent on the model parameters  $\theta$ :

$$J(\theta) = c \quad (1)$$

In case of a *classification task*, the network has to assign an  $n$ -dimensional input vector  $\mathbf{x} \in \mathbb{R}^n$  to a certain class  $i$  of  $k$  classes  $i \in \{1, \dots, k\}$ . One possibility to achieve this is to train the network to compute a *probability distribution* over all classes  $k$  for a given  $\mathbf{x}$ . The cost function could be the *mean square error (MSE)* of the output of the model compared to the desired output of the model.

## Gradient Descent

Learning of a parameterized model is to optimize the parameters of the model in a way to minimize a *cost function* (also called *objective function*, *loss function* or *error function*).

As typically the optimal values of the parameters cannot be calculated directly, an iterative optimization approach is used.

If we assume that  $J(\theta)$  is the cost function providing a cost value  $c$  for a parameter set  $\theta$ . We want to find the optimal value for  $\theta$  so that  $J(\theta)$  is minimal. We use the derivative  $J'(\theta)$  which gives us the slope at point  $\theta$ . If the slope  $J'(\theta) > 0$ , decreasing  $\theta$  will decrease  $J(\theta)$ . If the slope  $J'(\theta) < 0$ , increasing  $\theta$  will decrease  $J(\theta)$ . By iteratively calculating new values for  $\theta$  with:

$$\theta^{new} = \theta - \epsilon J'(\theta) \quad (2)$$

we can find at least a local minimum for  $J(\theta)$  if  $\epsilon$  is small enough.  $\epsilon$  is called the *learning rate* and is a positive small number (usually  $\epsilon \ll 1$ ).

As  $\theta$  is an  $n$ -dimensional vector, the derivative is also a vector called the *gradient*  $\nabla_{\theta} J(\theta)$ . Element  $i$  of the gradient is the partial derivative of  $J$  with respect to  $\theta_i$ . The iterative process of formula (2) is written:

$$\theta^{new} = \theta - \epsilon \nabla_{\theta} J(\theta) \quad (3)$$

This iterative technique is called *gradient descent* and is generally attributed to *Augustin-Louis Cauchy*, who first suggested it in 1847.

# Stochastic Gradient Descent (SGD)

(Goodfellow et.al. Deep Learning Book, p. 150)

Nearly all *deep learning* algorithms are working with a particular version of gradient descent: *stochastic gradient descent (SGD)*.

We have a set of several examples  $X = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  and  $Y = \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(m)}\}$  of a random vector  $\mathbf{x}$  and an associated value or vector  $\mathbf{y}$ , and we are going to learn to predict  $\mathbf{y}$  from  $\mathbf{x}$  with gradient descent. We define the negative conditional log-likelihood (NLL) as our cost function  $J(\theta)$  of a set of parameter  $\theta$ :

$$J(\theta) = E_{\mathbf{x}, \mathbf{y} \sim \hat{P}_{data}} [L(\mathbf{x}, \mathbf{y}, \theta)] = \frac{1}{m} \sum_{i=1}^m L(\mathbf{x}^{(i)}, y^{(i)}, \theta) \quad (4)$$

$L$  is the per-example loss:

$$L(\mathbf{x}, \mathbf{y}, \theta) = -\log p(\mathbf{y}|\mathbf{x}, \theta) \quad (5)$$

For this additive cost function, the gradient descent requires the computing of all per-example losses:

$$\nabla_{\theta} J(\theta) = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} L(\mathbf{x}^{(i)}, y^{(i)}, \theta) \quad (6)$$

When the training size  $m$  is large, this is computational expensive or even impractical.

The idea of stochastic gradient descent is to see the gradient as an *expectation* (like in formula (4)). This expectation can be approximately estimated using a smaller set of examples, a *minibatch* of examples  $B_X = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m')}\}$  and  $B_Y = \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(m')}\}$  drawn uniformly from the training set. The size of the minibatch  $m'$  is typically chosen to be a small number ranging between 1 and a few hundred.

The estimate of the gradient  $\mathbf{g}$  is calculated:

$$\mathbf{g} = \frac{1}{m'} \nabla_{\theta} \sum_{i=1}^{m'} L(\mathbf{x}^{(i)}, y^{(i)}, \theta) \quad (7)$$

using examples  $\mathbf{x}^{(i)}$  and  $\mathbf{y}^{(i)}$  from the minibatch  $B_X$  and  $B_Y$ . Analog to formula (3) the parameters  $\theta$  are changed along the negative estimate of the gradient  $\mathbf{g}$  multiplied by the learning rate  $\epsilon$ :

$$\theta^{new} = \theta - \epsilon \mathbf{g} \quad (8)$$