

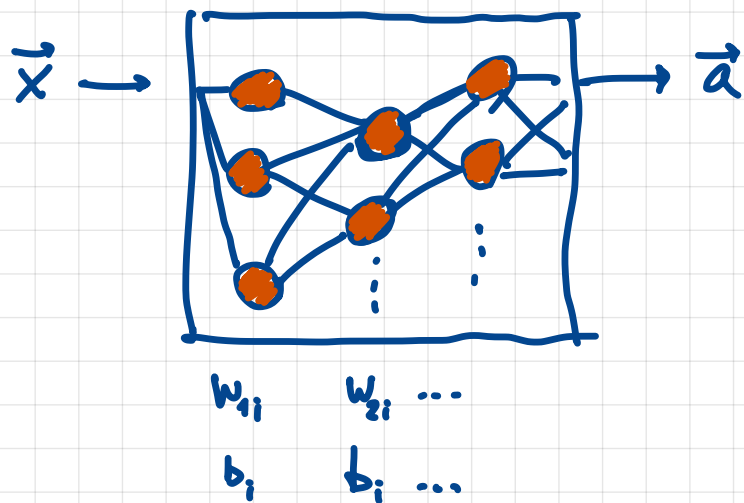
Short Takes 331

Machine learning:
Stochastic
gradient descent



More ML: Stochastic gradient descent

- A neural network (last video!) can be regarded as a complicated mathematical expression connecting the input vector \vec{x} with the output layer activation \vec{a} .



- Our job is to optimize the weights and biases such that the output matches a desired result for a given input.
 \vec{a} \vec{y} \vec{x}

- Suppose we are given data samples
 $\{(\vec{x}, \vec{y})_1, (\vec{x}, \vec{y})_2, \dots, (\vec{x}, \vec{y})_{N_{\text{samples}}}\}$

You'll want to divide these N_{samples} into a training set and a test set
 N_T N_V

- To optimize the network one chooses a cost function $C(w, b)$.

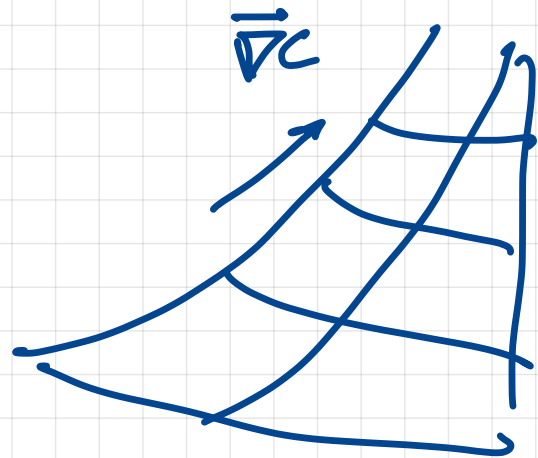
Normally, to optimize in this multivariable space, one uses methods like gradient descent (GD)...

Idea: To minimize $C(w, b)$, calculate the gradient vector $\vec{\nabla} C = \begin{pmatrix} \frac{\partial C}{\partial w} \\ \frac{\partial C}{\partial b} \end{pmatrix}$ and move in the opposite direction...

• start search at $\begin{pmatrix} w_{(0)} \\ b_{(0)} \end{pmatrix}$

• update to $\begin{pmatrix} w_{(1)} \\ b_{(1)} \end{pmatrix} = \begin{pmatrix} w_{(0)} \\ b_{(0)} \end{pmatrix} - \eta \cdot \bar{\nabla} C$

free parameter
(learning rate)



• Assuming $C(w, b)$ does have a minimum, this process will (at least intuitively) take us there.

• In practice,

• $C(w, b) = \frac{1}{N_T} \sum_{i=1}^{N_T} C_i(w, b)$,

where $C_i(w, b)$ is calculated based on sample $(\bar{x}_i, \bar{y}_i)_i$ only.

Eq. $C_i(w, b) = (\bar{y}_i - \bar{a}(\bar{x}_i, w, b))^2$

"quadratic cost"

desired output (pointing to \bar{y}_i)

network output (pointing to $\bar{a}(\bar{x}_i, w, b)$)

Another choice: "cross entropy"

$$C_i(w, b) = - \left[\bar{y}_i \cdot \ln \bar{a}(\bar{x}_i, w, b) + (1 - \bar{y}_i) \cdot \ln (1 - \bar{a}(\bar{x}_i, w, b)) \right]$$

note dot prod.

(we assume here $0 \leq y \leq 1$)

• Calculating $\bar{\nabla} C$ for all the samples is prohibitively costly.

We take mini-batches to approximate ...

$$N_T = m \cdot M$$

Size
of mini-batch

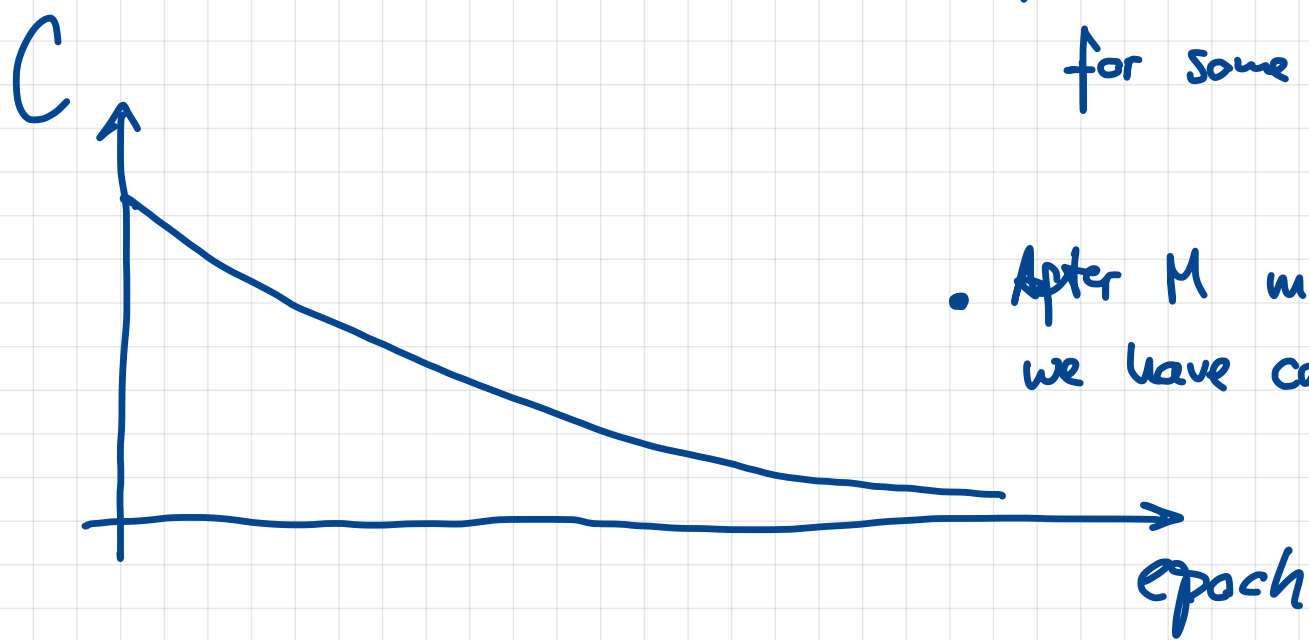
mini-batches

of randomly chosen
samples

The "stochastic" part.

$$\vec{\nabla} C = \frac{1}{N_T} \cdot \sum_{i=1}^{N_T} \vec{\nabla} C_i \approx \frac{1}{M} \sum_{i=1}^M \vec{\nabla} C_i$$

use in gradient descent
for some number of GD iterations.



- After M mini-batch rounds, we have completed an "epoch".

For NN's, $\vec{\nabla} C_i$ can be calculated quickly via "back propagation" ...
next time!

————— x —————