

# CS5242 : Neural Networks and Deep Learning

## Lecture 4: Vanilla Neural Networks Loss and Optimization

Semester 1 2021/22

Xavier Bresson

<https://twitter.com/xbresson>

Department of Computer Science  
National University of Singapore (NUS)



# Outline

- Quality of neural network
- Maximizing the quality by gradient ascent
- Learning rate strategy
- Neural network loss
- Minimizing the loss by gradient descent
- Mini-batch/stochastic gradient descent
- Summary

# Outline

- **Quality of neural network**
- Maximizing the quality by gradient ascent
- Learning rate strategy
- Neural network loss
- Minimizing the loss by gradient descent
- Mini-batch/stochastic gradient descent
- Summary

# Neural network quality

- Quality of neural networks :
  - Evaluate the confidence of the net to predict the correct class over the whole training set.
- First, assume for simplicity that the training set consists of
  - 4 classes : Horse, Cat, Dog, Deer
  - 5 labelled pictures



HORSE



CAT



DOG

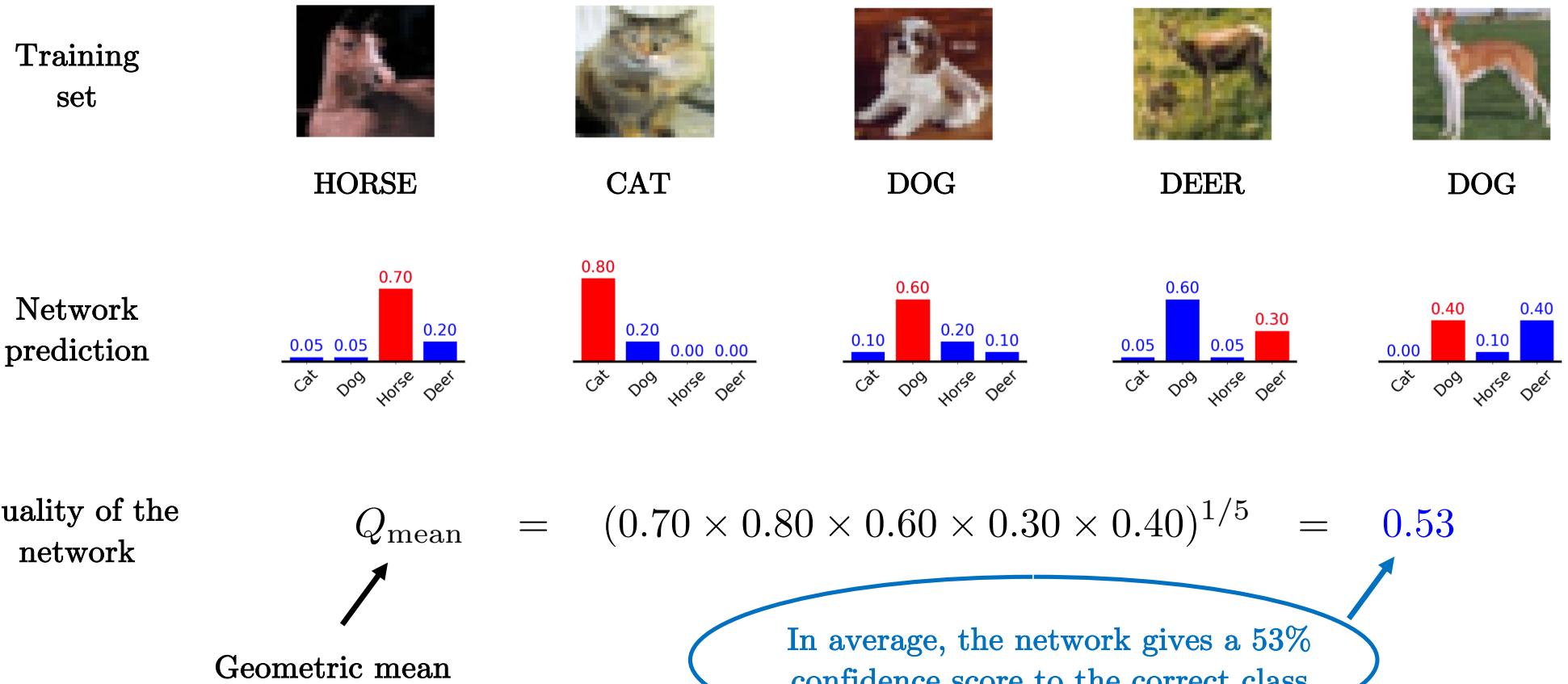


DEER

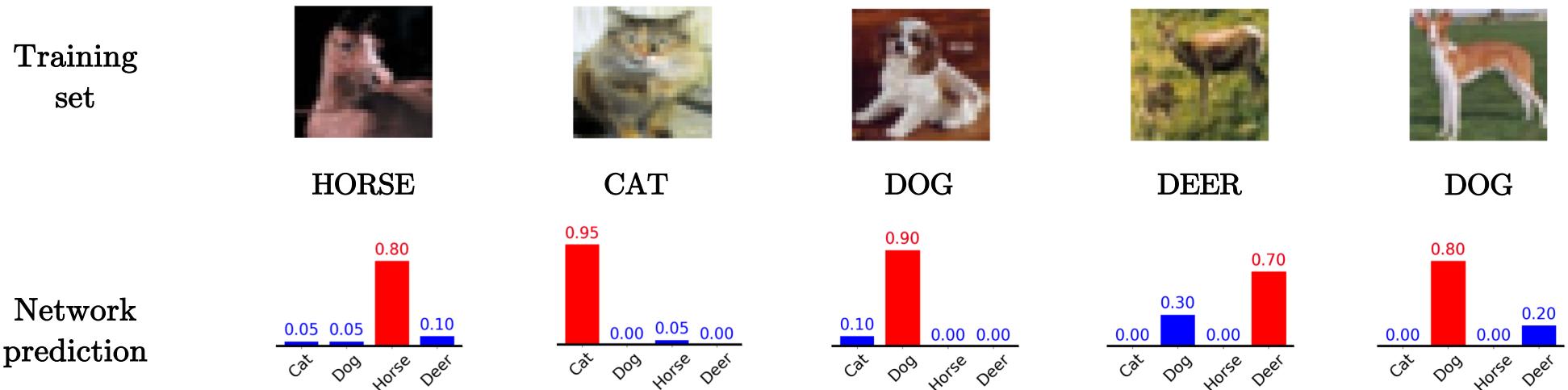


DOG

# Neural network quality



## Good NN quality



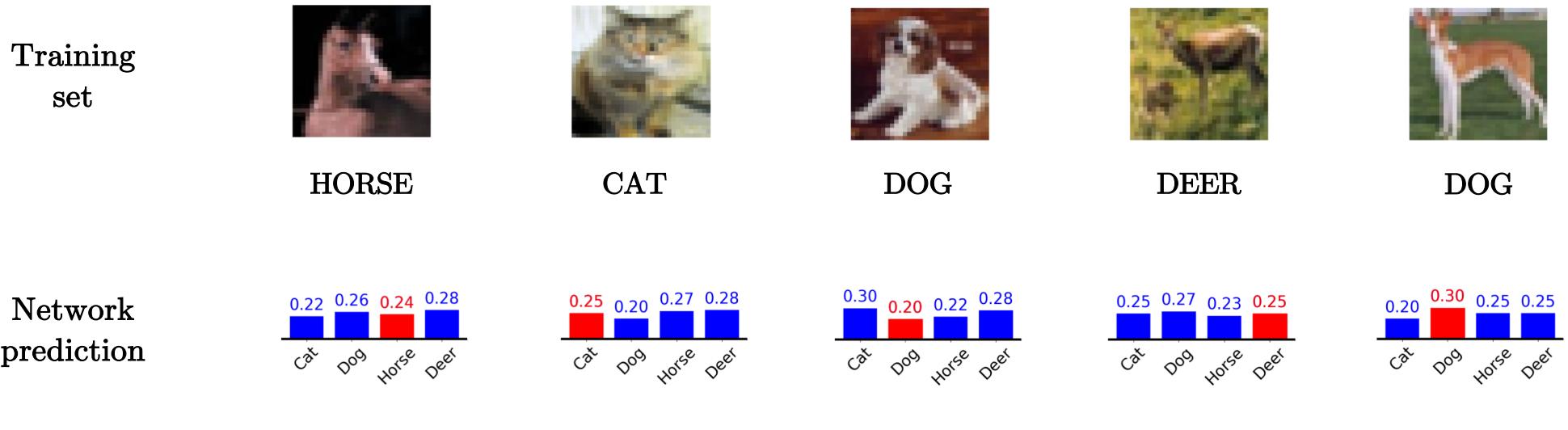
Quality of the network

$$Q_{\text{mean}} = (0.80 \times 0.95 \times 0.90 \times 0.70 \times 0.80)^{1/5} = 0.83$$

In average, the network gives a 83% confidence score to the correct class

Optimal classification:  $Q_{\text{mean}} \approx 1.00$

## Bad NN quality



Quality of the network

$$Q_{\text{mean}} = (0.24 \times 0.25 \times 0.20 \times 0.25 \times 0.30)^{1/5} = 0.24$$

In average, the network gives a 24% confidence score to the correct class

Random classification:  $Q_{\text{mean}} \approx 0.25 \approx 1/K$

# Neural network quality

The  $i^{\text{th}}$   
picture of the  
training set

$$\mathbf{x}^{(i)}$$



$$i = 1$$



$$i = 2$$



$$i = 3$$



$$i = 4$$



$$i = 5$$

The class  
index of the  
 $i^{\text{th}}$  picture

$$\text{cl}(i)$$

$$3$$

$$1$$

$$2$$

$$4$$

$$2$$

Network  
output for the  
 $i^{\text{th}}$  picture

$$Q_{\text{mean}} = (0.70 \times 0.80 \times 0.60 \times 0.30 \times 0.40)^{1/5}$$

$$= \left( (\text{entry } \text{cl}(1) \text{ of } \mathbf{p}^{(1)}) \times (\text{entry } \text{cl}(2) \text{ of } \mathbf{p}^{(2)}) \times \cdots \times (\text{entry } \text{cl}(5) \text{ of } \mathbf{p}^{(5)}) \right)^{1/5}$$

$$= \left( \prod_{i=1}^5 \text{entry } \text{cl}(i) \text{ of vector } \mathbf{p}^{(i)} \right)^{1/5}$$



# Neural network quality

- With a training set containing  $N$  labeled pictures :

$$Q_{\text{mean}} = \left( \prod_{i=1}^N \text{entry } \text{cl}(i) \text{ of vector } \mathbf{p}^{(i)} \right)^{1/N}$$

The class index of the  $i^{\text{th}}$  picture (integer)

How confident the network is that the  $i^{\text{th}}$  picture belongs to the class  $\text{cl}(i)$

Network output for the  $i^{\text{th}}$  picture  $\mathbf{x}^{(i)}$  :

$$\mathbf{p}^{(i)} = \text{Softmax}(W\mathbf{x}^{(i)} + \mathbf{b})$$

Vector with 4 entries

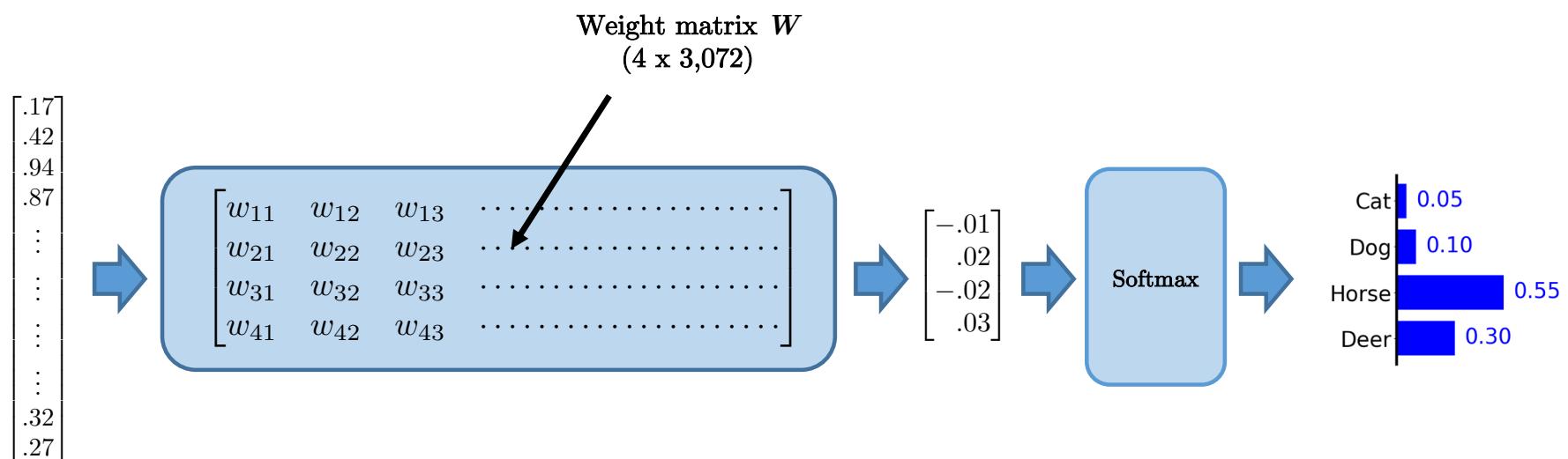
Cat	0.10
Dog	0.10
Horse	0.60
Deer	0.20

# Outline

- Quality of neural network
- **Maximizing the quality by gradient ascent**
- Learning rate strategy
- Neural network loss
- Minimizing the loss by gradient descent
- Mini-batch/stochastic gradient descent
- Summary

# Maximizing the quality of neural networks

- Learning = Maximizing the quality  $Q_{mean}$  of the network over its parameters  $W$ .
- Maximizing  $Q_{mean}$  will be done by gradient ascent.
- Let us consider a vanilla one-layer NN :



## Quality function

- The weight matrix  $W$  contains the 4 templates of size 3,072 :

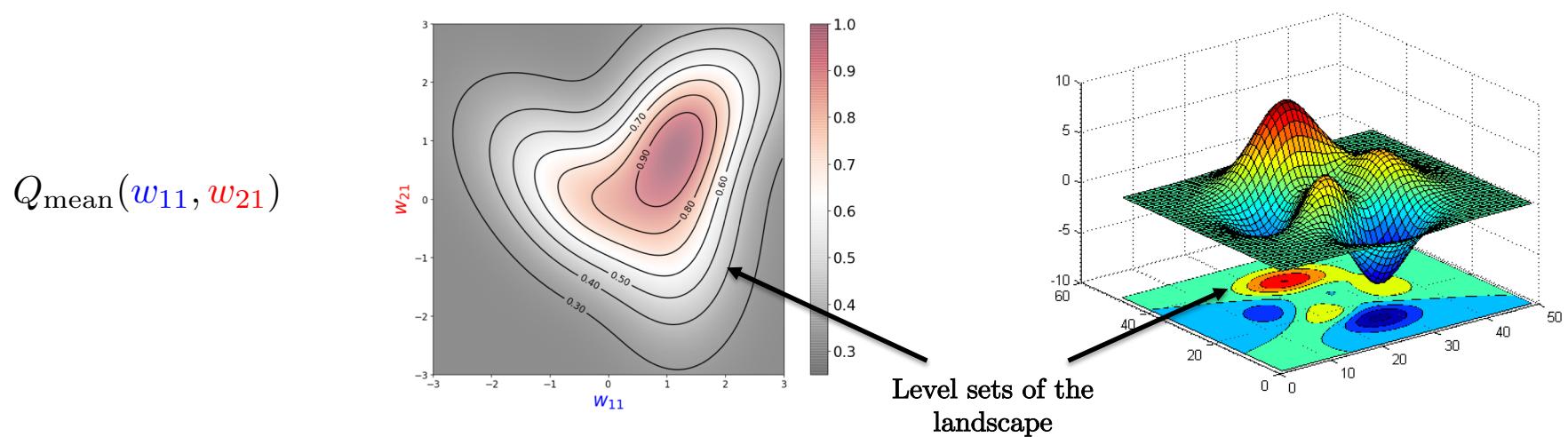
$$W = \begin{bmatrix} w_{11} & w_{12} & w_{13} & \dots \\ w_{21} & w_{22} & w_{23} & \dots \\ w_{31} & w_{32} & w_{33} & \dots \\ w_{41} & w_{42} & w_{43} & \dots \end{bmatrix}$$

- The Quality  $Q_{\text{mean}}$  is a function of  $\underbrace{12,288}_{4 \times 3,072}$  parameters/weights :

$$Q_{\text{mean}}(w_{11}, w_{12}, w_{13}, \dots ; w_{21}, w_{22}, w_{23}, \dots ; w_{31}, w_{32}, w_{33}, \dots ; w_{41}, w_{42}, w_{43}, \dots)$$
$$Q_{\text{mean}}(W)$$

# Landscape of quality function

- For each configuration  $W$  of the network, there exists a value  $Q_{mean}(W)$ .
  - All the possible configurations  $W$  provide the landscape of  $Q_{mean}$ .
- As  $Q_{mean}$  has 12k variables, it is not possible to visualize this landscape.
- Let's assume that  $Q_{mean}$  has only 2 parameters for visualization purpose :

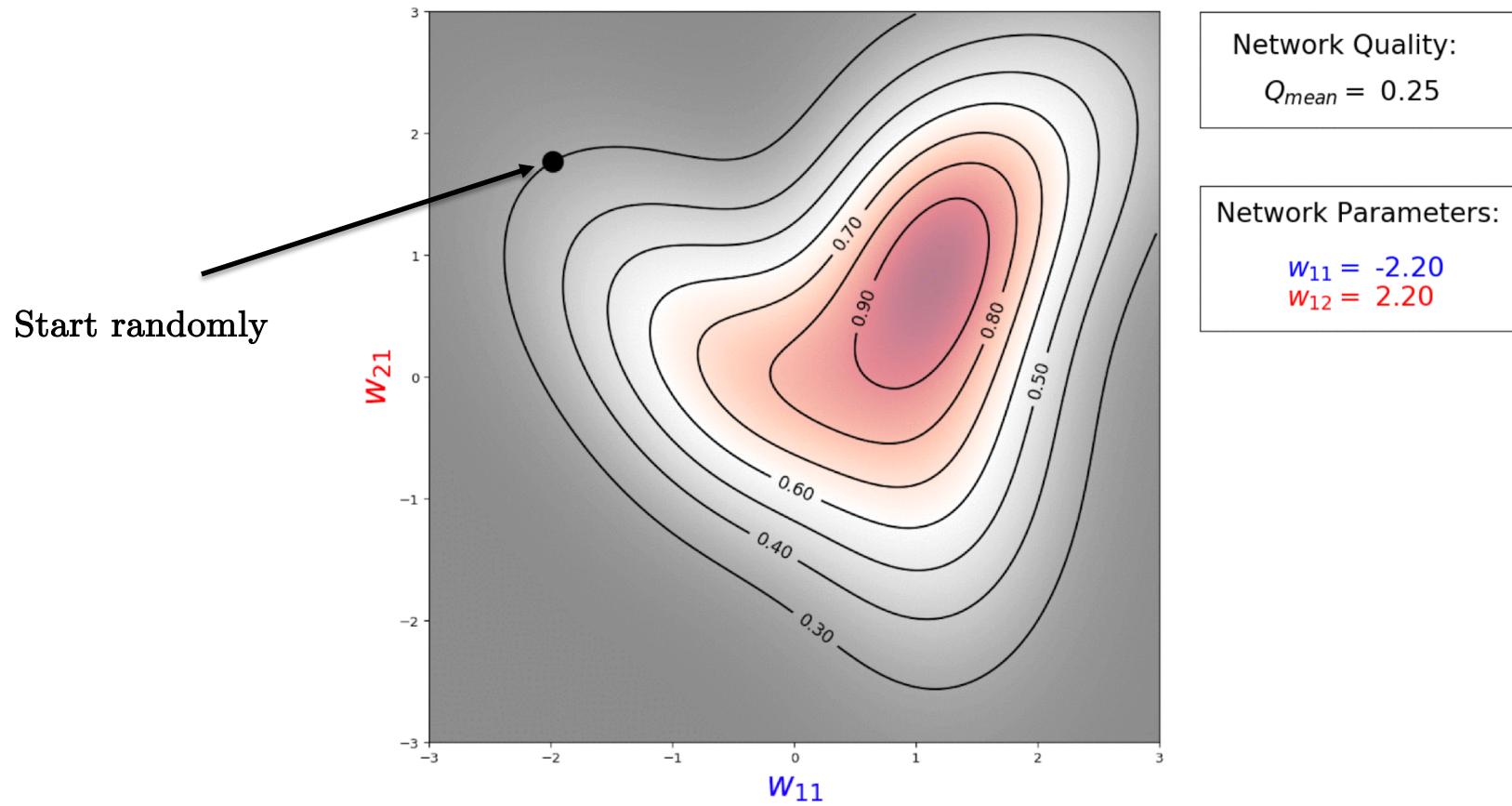


# Maximizing the quality

- Gradient ascent algorithm:
  - Start randomly.
  - Move in the direction of the steepest ascent.
  - Stop when reaching the maximum.



# Initialization

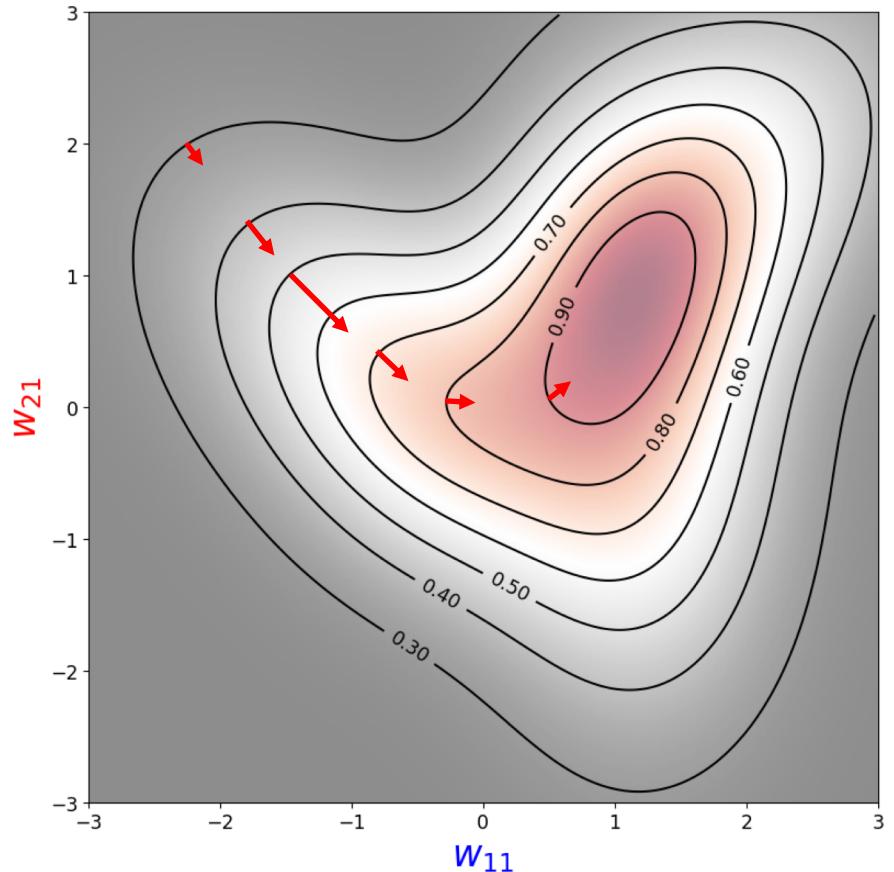


# Update scheme

- Gradient ascent:  
Follow the gradient!
- The parameters are updated according to the equation:

$$\begin{bmatrix} w_{11} \\ w_{21} \end{bmatrix} = \begin{bmatrix} w_{11} \\ w_{21} \end{bmatrix} + lr \begin{bmatrix} \frac{\partial Q}{\partial w_{11}} \\ \frac{\partial Q}{\partial w_{21}} \end{bmatrix}$$

- So the displacements are proportional to the gradient.
- The gradient is:
  - Perpendicular to the level sets.
  - Its magnitude is proportional to the slope of the landscape.



# Maximum

- At the maximum:

- We reach the top of the mountain, which is perfectly flat.
- We stop moving since the magnitude of the gradient is zero.

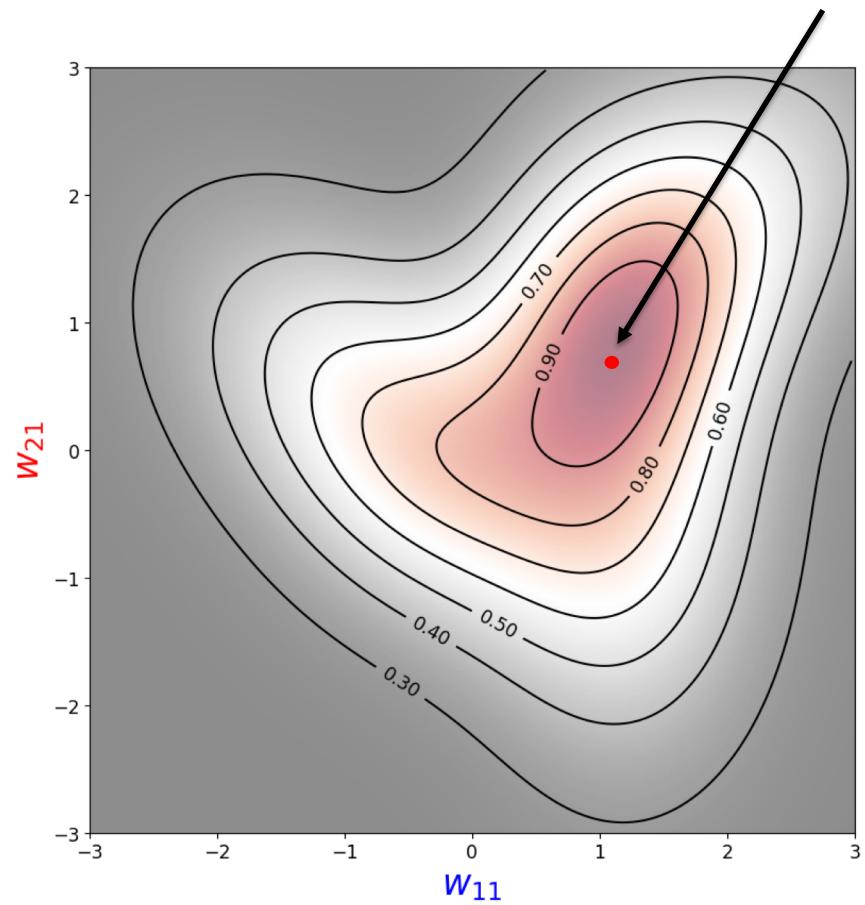
At the maximum location:

$$\begin{aligned} w_{11} &= 1.1 \\ w_{21} &= 0.8 \end{aligned}$$

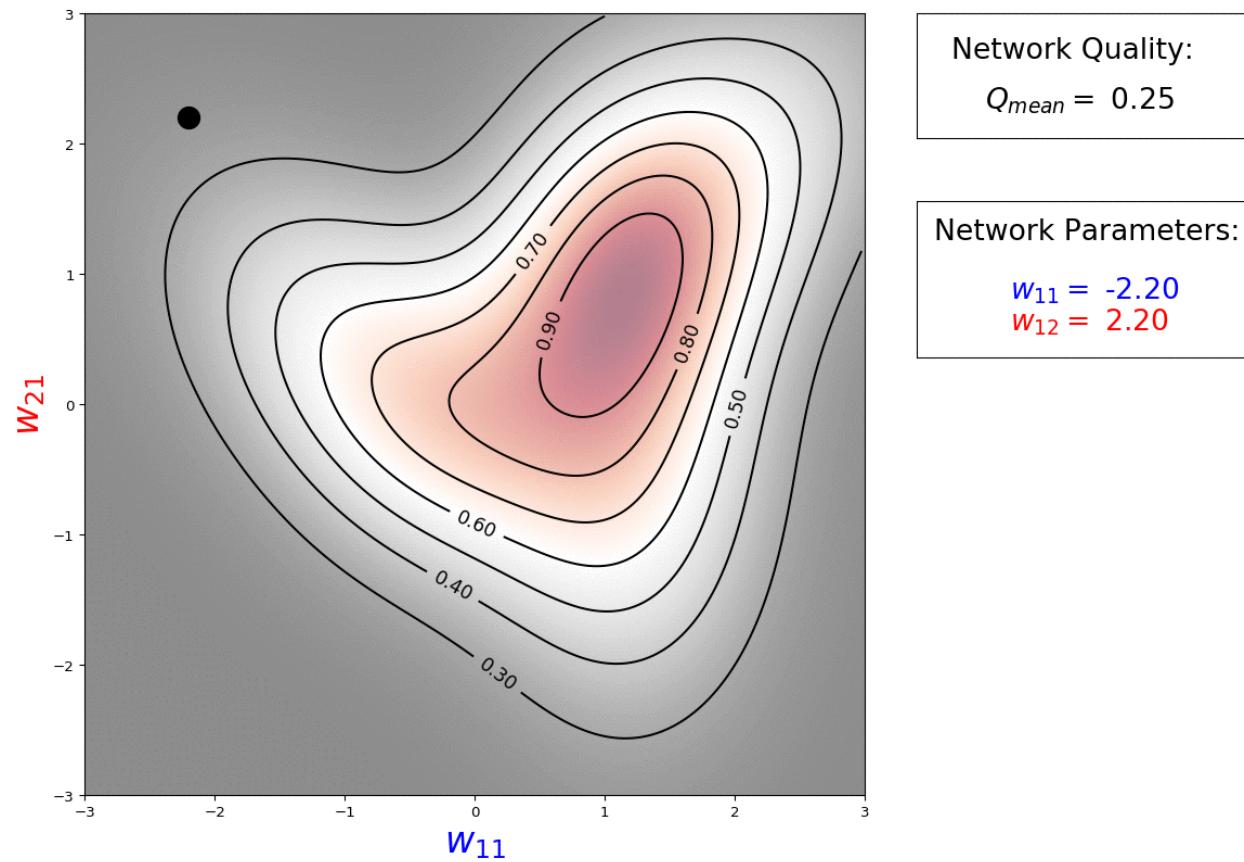
We have:

$$\begin{bmatrix} \frac{\partial Q}{\partial w_{11}} \\ \frac{\partial Q}{\partial w_{21}} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Maximum of the landscape



# Gradient ascent in action



# High-dimensional quality function

- The previous  $Q_{mean}$  function has 2 variables.
  - The gradient vector  $\nabla Q_{mean}$  has thus 2 elements.

$$Q_{\text{mean}}(\color{blue}{w_{11}}, \color{red}{w_{21}}) \quad \nabla Q_{\text{mean}} = \begin{bmatrix} \frac{\partial Q}{\partial w_{11}} \\ \frac{\partial Q}{\partial w_{21}} \end{bmatrix}$$

- But the original quality  $Q_{mean}$  is a **function of 12,288 parameters/weights**.
  - The gradient vector  $\nabla Q_{mean}$  has consequently 12,288 elements :

$$Q_{\text{mean}}(w_{11}, w_{12}, w_{13}, \dots ; \color{blue}{w_{21}}, \color{red}{w_{22}}, \color{green}{w_{23}}, \dots ; \color{teal}{w_{31}}, w_{32}, w_{33}, \dots ; \color{magenta}{w_{41}}, w_{42}, w_{43}, \dots)$$

$$\nabla Q_{\text{mean}} = \begin{bmatrix} \frac{\partial Q}{\partial w_{11}} \\ \frac{\partial Q}{\partial w_{12}} \\ \frac{\partial Q}{\partial w_{13}} \\ \vdots \\ \frac{\partial Q}{\partial w_{21}} \\ \frac{\partial Q}{\partial w_{22}} \\ \frac{\partial Q}{\partial w_{23}} \\ \vdots \\ \frac{\partial Q}{\partial w_{31}} \\ \frac{\partial Q}{\partial w_{32}} \\ \frac{\partial Q}{\partial w_{33}} \\ \vdots \\ \frac{\partial Q}{\partial w_{41}} \\ \frac{\partial Q}{\partial w_{42}} \\ \frac{\partial Q}{\partial w_{43}} \\ \vdots \end{bmatrix}$$

## High-dimensional update formula

$$w_{ij} = w_{ij} + \text{lr} \frac{\partial Q}{\partial w_{ij}} \quad \Rightarrow$$

$$\begin{bmatrix} w_{11} \\ w_{12} \\ w_{13} \\ \vdots \\ w_{21} \\ w_{22} \\ w_{23} \\ \vdots \\ w_{31} \\ w_{32} \\ w_{33} \\ \vdots \\ w_{41} \\ w_{42} \\ w_{43} \\ \vdots \end{bmatrix} = \begin{bmatrix} w_{11} \\ w_{12} \\ w_{13} \\ \vdots \\ w_{21} \\ w_{22} \\ w_{23} \\ \vdots \\ w_{31} \\ w_{32} \\ w_{33} \\ \vdots \\ w_{41} \\ w_{42} \\ w_{43} \\ \vdots \end{bmatrix} + \text{lr} \begin{bmatrix} \partial Q / \partial w_{11} \\ \partial Q / \partial w_{12} \\ \partial Q / \partial w_{13} \\ \vdots \\ \partial Q / \partial w_{21} \\ \partial Q / \partial w_{22} \\ \partial Q / \partial w_{23} \\ \vdots \\ \partial Q / \partial w_{31} \\ \partial Q / \partial w_{32} \\ \partial Q / \partial w_{33} \\ \vdots \\ \partial Q / \partial w_{41} \\ \partial Q / \partial w_{42} \\ \partial Q / \partial w_{43} \\ \vdots \end{bmatrix}$$

# Matrix representation

- **Scalar form** of the gradient ascent equation:

$$w_{ij} = w_{ij} + \text{lr} \frac{\partial Q}{\partial w_{ij}}$$

- **Matrix form** of the gradient ascent equation:

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} & \dots \\ w_{21} & w_{22} & w_{23} & \dots \\ w_{31} & w_{32} & w_{33} & \dots \\ w_{41} & w_{42} & w_{43} & \dots \end{bmatrix} = W$$
$$\begin{bmatrix} w_{11} & w_{12} & w_{13} & \dots \\ w_{21} & w_{22} & w_{23} & \dots \\ w_{31} & w_{32} & w_{33} & \dots \\ w_{41} & w_{42} & w_{43} & \dots \end{bmatrix} + \text{lr} \begin{bmatrix} \frac{\partial Q}{\partial w_{11}} & \frac{\partial Q}{\partial w_{12}} & \frac{\partial Q}{\partial w_{13}} & \dots \\ \frac{\partial Q}{\partial w_{21}} & \frac{\partial Q}{\partial w_{22}} & \frac{\partial Q}{\partial w_{23}} & \dots \\ \frac{\partial Q}{\partial w_{31}} & \frac{\partial Q}{\partial w_{32}} & \frac{\partial Q}{\partial w_{33}} & \dots \\ \frac{\partial Q}{\partial w_{41}} & \frac{\partial Q}{\partial w_{42}} & \frac{\partial Q}{\partial w_{43}} & \dots \end{bmatrix}$$

Gradient reshaped  
as a matrix

$$\frac{\partial Q}{\partial W}$$

# Outline

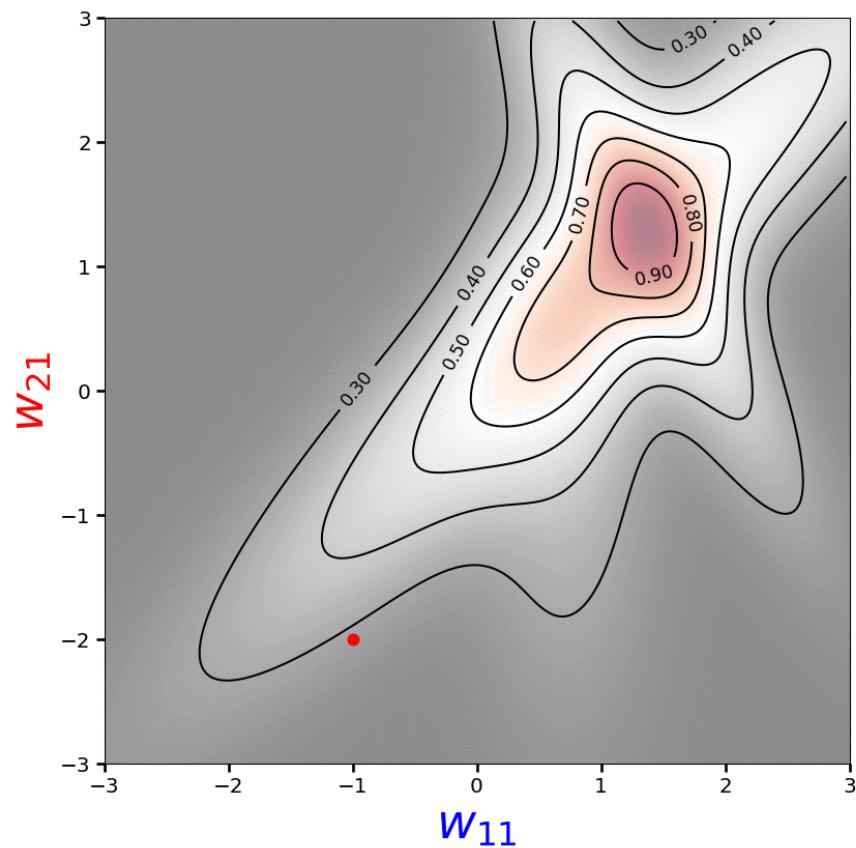
- Quality of neural network
- Maximizing the quality by gradient ascent
- **Learning rate strategy**
- Neural network loss
- Minimizing the loss by gradient descent
- Mini-batch/stochastic gradient descent
- Summary

# Learning rate

- It is important to carefully choose the learning rate:
  - If the learning rate is **too big**, the algorithm will **diverge**.
  - If the learning rate is **too small**, it will **take too long to reach the maximum**.

## Learning rate too large

- lr = 5.0



Iteration 1

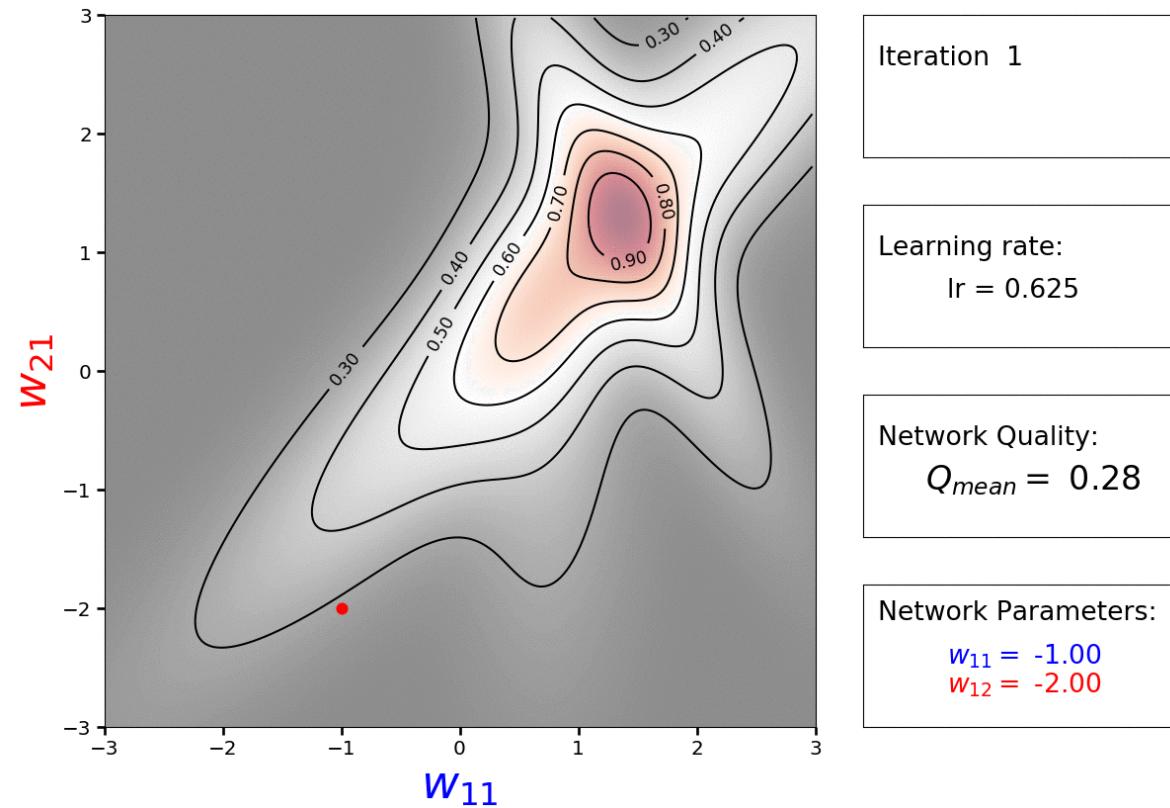
Learning rate:  
lr = 5.0

Network Quality:  
 $Q_{mean} = 0.28$

Network Parameters:  
 $w_{11} = -1.00$   
 $w_{12} = -2.00$

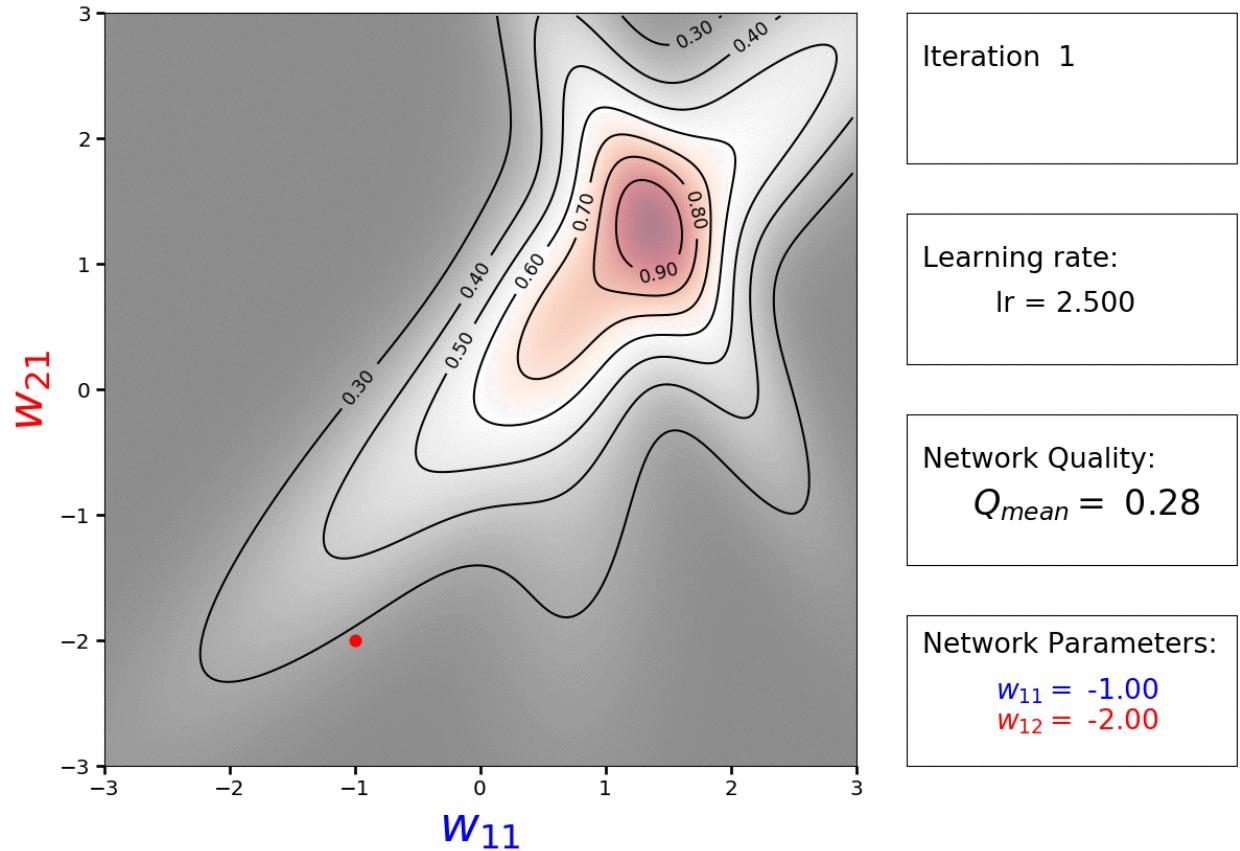
## Learning rate too small

- $lr=0.625$ 
  - 8 times smaller than  $lr=5.0$



# Learning rate schedule

- Best way to go !
  - Make a learning rate schedule :
    - Start with  $lr=2.5$ .
    - Divide by 2 every 10 iterations.



# Outline

- Quality of neural network
- Maximizing the quality by gradient ascent
- Learning rate strategy
- **Neural network loss**
- Minimizing the loss by gradient descent
- Mini-batch/stochastic gradient descent
- Summary

## Optimizing product is slow

- Quality function is a product of functions:

$$Q_{\text{mean}} = \left( \prod_{i=1}^N \text{entry } \text{cl}(i) \text{ of vector } \mathbf{p}^{(i)} \right)^{1/N}$$

- Gradient of a product is computationally expensive (e.g. derivate of a product of 50,000 terms)

$$w_{ij} = w_{ij} + \text{lr} \frac{\partial Q}{\partial w_{ij}}$$

- It is much cheaper to compute the gradient of a sum.

- How? The **log** function over a product provides a **sum**.

## New loss

- News loss = - log of neural network quality

$$\mathcal{L} = - \log Q_{\text{mean}} = - \log \left( \prod_{i=1}^N \text{entry } \text{cl}(i) \text{ of vector } \mathbf{p}^{(i)} \right)^{1/N}$$

$$= - \frac{1}{N} \log \left( \prod_{i=1}^N \text{entry } \text{cl}(i) \text{ of vector } \mathbf{p}^{(i)} \right)$$

$$= - \frac{1}{N} \sum_{i=1}^N \log \left( \text{entry } \text{cl}(i) \text{ of vector } \mathbf{p}^{(i)} \right)$$

This function is much easier to differentiate.

# Cross-entropy loss

- The new loss is called cross-entropy loss.
- Cross-entropy is a standard function in information theory.
- Cross-entropy is a function of  $W$ :

$$\begin{aligned} \mathcal{L} &= \frac{1}{N} \sum_{i=1}^N -\log \left( \text{entry } \text{cl}(i) \text{ of vector } \mathbf{p}^{(i)} \right) \\ \mathcal{L} &= \frac{1}{N} \sum_{i=1}^N -\text{entry } \text{cl}(i) \text{ of vector } \log \left( \mathbf{p}^{(i)} \right) \\ \mathcal{L}(w_{11}, w_{12}, \dots) &= \frac{1}{N} \sum_{i=1}^N -\text{entry } \text{cl}(i) \text{ of vector } \text{LogSoftmax}(W \mathbf{x}^{(i)}) \end{aligned}$$

Function of  
12,288 parameters.

log  $\begin{pmatrix} p_1^{(i)} \\ p_2^{(i)} \\ p_3^{(i)} \\ p_4^{(i)} \end{pmatrix}$  =  $\log(p_2^{(i)})$

## Cross-entropy loss

$$\mathcal{L}(w_{11}, w_{12}, \dots) = \overbrace{\frac{1}{N} \sum_{i=1}^N \underbrace{-\text{entry } \text{cl}(i) \text{ of vector } \log(\mathbf{p}^{(i)})}_{\ell^{(i)}(w_{11}, w_{12}, \dots)}}^{\text{Average loss over all data points}}$$

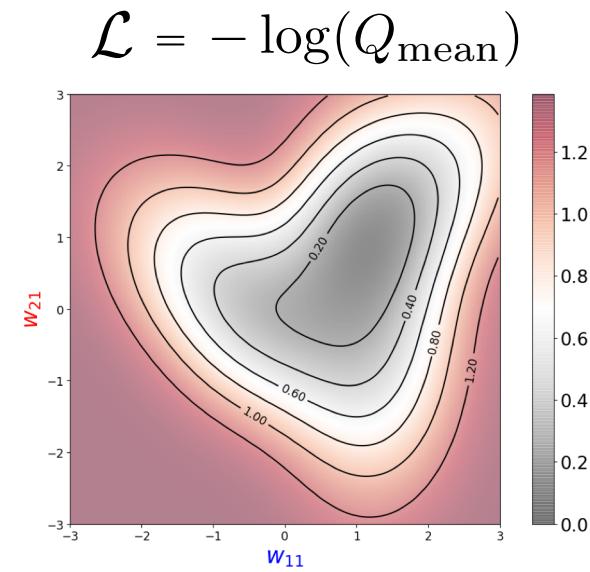
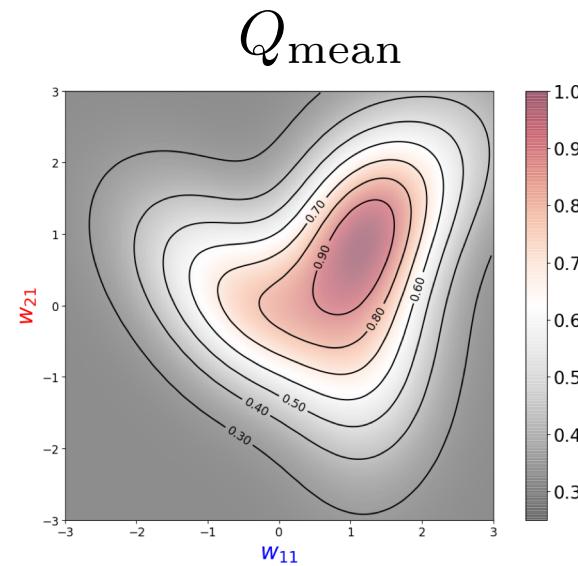
Loss of the  $i^{\text{th}}$  data point =  
- log of the probability given to the correct class

# Cross-entropy loss

	$i = 1$	$i = 2$	$i = 3$	$i = 4$	$i = 5$
$\mathbf{x}^{(i)}$					
$\text{cl}^{(i)}$	3	1	2	4	2
$\mathbf{p}^{(i)}$	$\begin{bmatrix} 0.05 \\ 0.05 \\ \textcolor{red}{0.70} \\ 0.20 \end{bmatrix}$	$\begin{bmatrix} \textcolor{red}{0.80} \\ 0.20 \\ 0.00 \\ 0.00 \end{bmatrix}$	$\begin{bmatrix} 0.05 \\ \textcolor{red}{0.60} \\ 0.20 \\ 0.10 \end{bmatrix}$	$\begin{bmatrix} 0.05 \\ 0.60 \\ 0.05 \\ \textcolor{red}{0.30} \end{bmatrix}$	$\begin{bmatrix} 0.00 \\ \textcolor{red}{0.40} \\ 0.10 \\ 0.40 \end{bmatrix}$

$$\begin{aligned}
 \text{Average Loss} &= \frac{1}{5} \sum_1^5 -\text{entry } \text{cl}(i) \text{ of vector } \log(\mathbf{p}^{(i)}) \\
 &= \frac{1}{5} \left( -\log(\textcolor{red}{0.70}) - \log(\textcolor{red}{0.80}) - \log(\textcolor{red}{0.60}) - \log(\textcolor{red}{0.30}) - \log(\textcolor{red}{0.40}) \right) \\
 &= \frac{1}{5} \left( -0.37 + 0.22 + 0.51 + 1.20 + 0.92 \right) = \textcolor{red}{0.68}
 \end{aligned}$$

## Range of values

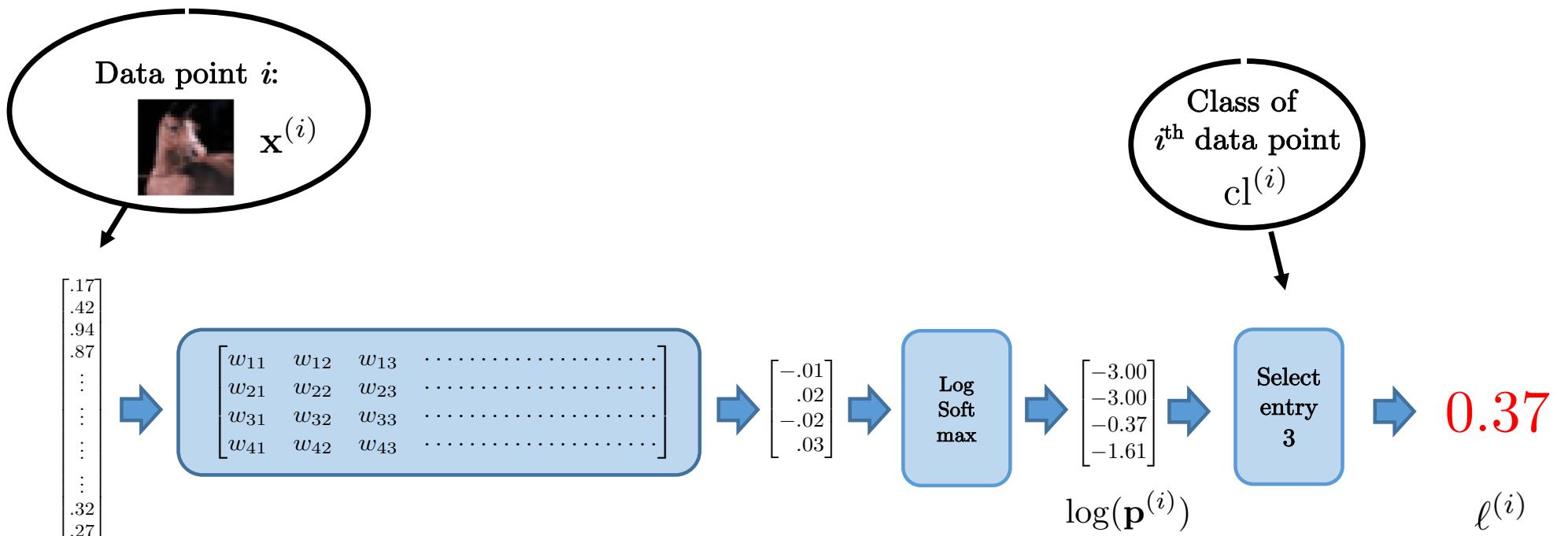


$$0.25 \leq Q_{\text{mean}} \leq 1$$

$0 \leq \mathcal{L} \leq 1.4$

-  $\log(1)$       We will want  $L$  to be as close to 0 as possible.      -  $\log(0.25)$

## Representing the network and the loss



# Outline

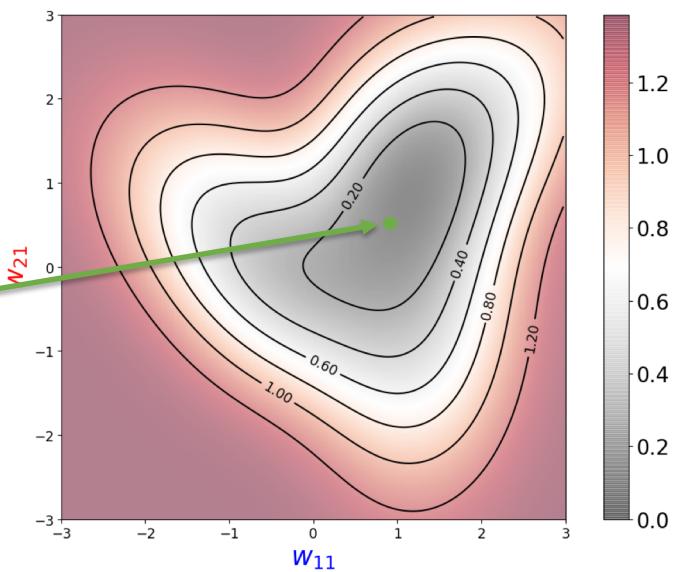
- Quality of neural network
- Maximizing the quality by gradient ascent
- Learning rate strategy
- Neural network loss
- **Minimizing the loss by gradient descent**
- Mini-batch/stochastic gradient descent
- Summary

# Loss

- The total loss over the whole training set is:

$$\mathcal{L}(w_{11}, w_{12}, \dots) = \frac{1}{N} \sum_{i=1}^N \ell^{(i)}(w_{11}, w_{12}, \dots)$$

We want to tune the internal parameters of the network in order to minimize the total loss  
(ideally make it equal to 0)



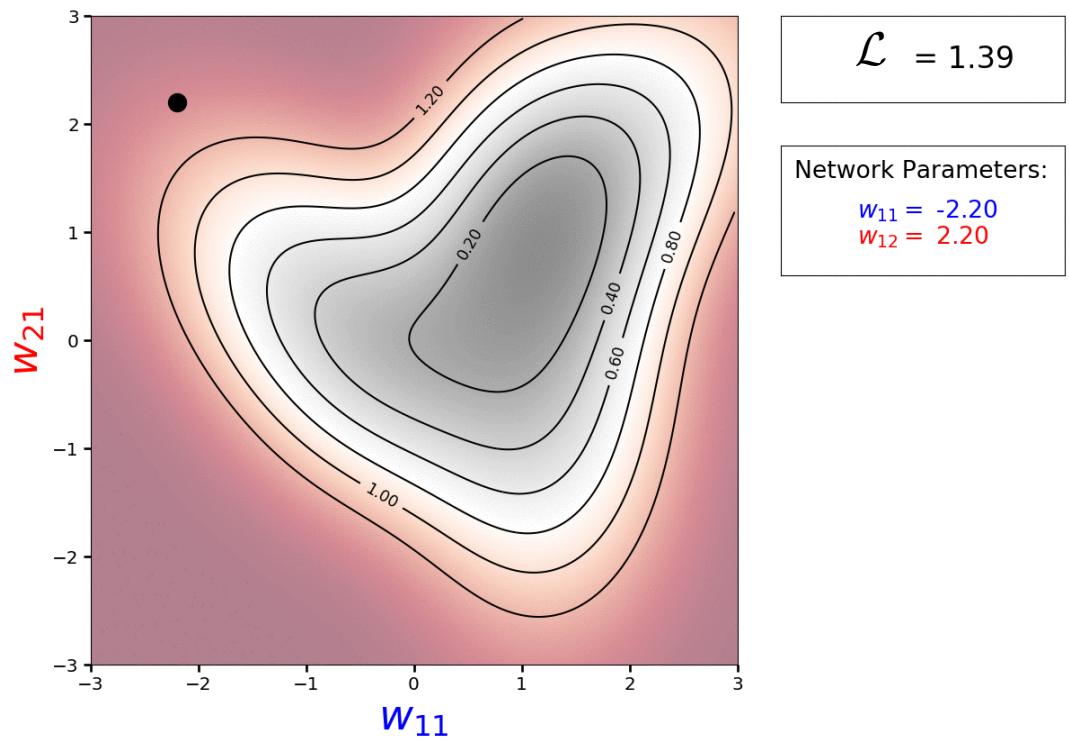
# Loss minimization

- The total loss is **minimized** by gradient descent:

$$w_{rs} = w_{rs} - \text{lr} \frac{\partial \mathcal{L}}{\partial w_{rs}}$$

Minimization

$$\begin{bmatrix} w_{11} \\ w_{12} \end{bmatrix} = \begin{bmatrix} w_{11} \\ w_{12} \end{bmatrix} - \text{lr} \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial w_{11}} \\ \frac{\partial \mathcal{L}}{\partial w_{12}} \end{bmatrix}$$



## Matrix representation

- Matrix form of the gradient descent equation:

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} & \dots \\ w_{21} & w_{22} & w_{23} & \dots \\ w_{31} & w_{32} & w_{33} & \dots \\ w_{41} & w_{42} & w_{43} & \dots \end{bmatrix} =$$

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} & \dots \\ w_{21} & w_{22} & w_{23} & \dots \\ w_{31} & w_{32} & w_{33} & \dots \\ w_{41} & w_{42} & w_{43} & \dots \end{bmatrix} - \text{lr} \begin{bmatrix} \partial \mathcal{L}/\partial w_{11} & \partial \mathcal{L}/\partial w_{12} & \partial \mathcal{L}/\partial w_{13} & \dots \\ \partial \mathcal{L}/\partial w_{21} & \partial \mathcal{L}/\partial w_{22} & \partial \mathcal{L}/\partial w_{23} & \dots \\ \partial \mathcal{L}/\partial w_{31} & \partial \mathcal{L}/\partial w_{32} & \partial \mathcal{L}/\partial w_{33} & \dots \\ \partial \mathcal{L}/\partial w_{41} & \partial \mathcal{L}/\partial w_{42} & \partial \mathcal{L}/\partial w_{43} & \dots \end{bmatrix}$$

$$W = W - \text{lr} \frac{\partial \mathcal{L}}{\partial W}$$

# Outline

- Quality of neural network
- Maximizing the quality by gradient ascent
- Learning rate strategy
- Neural network loss
- Minimizing the loss by gradient descent
- **Mini-batch/stochastic gradient descent**
- Summary

# Gradient descent

- Loss:

$$\mathcal{L}(w_{11}, w_{12}, \dots) = \frac{1}{N} \sum_{i=1}^N \ell^{(i)}(w_{11}, w_{12}, \dots)$$

(Average) loss                                      Loss of data point  $\mathbf{x}^{(i)}$

- (Full batch) gradient descent:

$$W = W - \text{lr} \frac{\partial \mathcal{L}}{\partial W}$$

$$\frac{\partial \mathcal{L}}{\partial W} = \underbrace{\frac{1}{N} \sum_{i=1}^N \frac{\partial \ell^{(i)}}{\partial W}}$$

The gradient matrix  $\frac{\partial \mathcal{L}}{\partial W}$  is the average of 50,000 gradient matrices  $\frac{\partial \ell^{(i)}}{\partial W}$  !

## Mini-batch gradient descent

- We can be more efficient by computing the average of only  $m$  gradient matrices.
- Mini-batch gradient descent:

$$W = W - \text{lr} \frac{\partial \mathcal{L}}{\partial W}$$
$$\frac{\partial \mathcal{L}}{\partial W} \approx \underbrace{\frac{1}{m} \sum_{i \in B} \frac{\partial \ell^{(i)}}{\partial W}}$$

Random approximation of the gradient  $\frac{\partial \mathcal{L}}{\partial W}$  using a batch  $B$  of  $m$  indices chosen at random:

$$B = \{i_1, i_2, i_3, \dots, i_m\}$$

# Mini-batch gradient descent

- **Advantages** of mini-batch gradient descent:

- **The gradient is quickly computed**: No need to go through the whole training set (which can be huge) to update the weights.
- **Datasets can be redundant** : Many data are similar, a subset of the data will give almost the same gradient than using all data points.

$$W = W - \text{lr} \frac{\partial \mathcal{L}}{\partial W}$$
$$\left\{ \begin{array}{lcl} \frac{\partial \mathcal{L}}{\partial W} & = & \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell^{(i)}}{\partial W} & \text{Mini-batch gradient} \\ \frac{\partial \mathcal{L}}{\partial W} & \approx & \frac{1}{m} \sum_{i \in B} \frac{\partial \ell^{(i)}}{\partial W} & \text{Full batch gradient} \end{array} \right.$$

# Mini-batch gradient descent

- Batch size :
  - No batch (single data): batch size is 1. This is called stochastic gradient descent (SGD). It approximates the gradient by averaging over 1 data point. Best for generalization but not as fast as mini-batch GD.
  - Mini-batch: batch size is 200. This is called mini-batch GD. It approximates the gradient by averaging over 200 data point. Good for generalization and fastest GD.
  - Full batch (all training data): batch size is 60,000. This is called gradient descent (GD). It provides the exact gradient by averaging over all data points. Weak for generalization and slowest GD.

# Connection to template matching

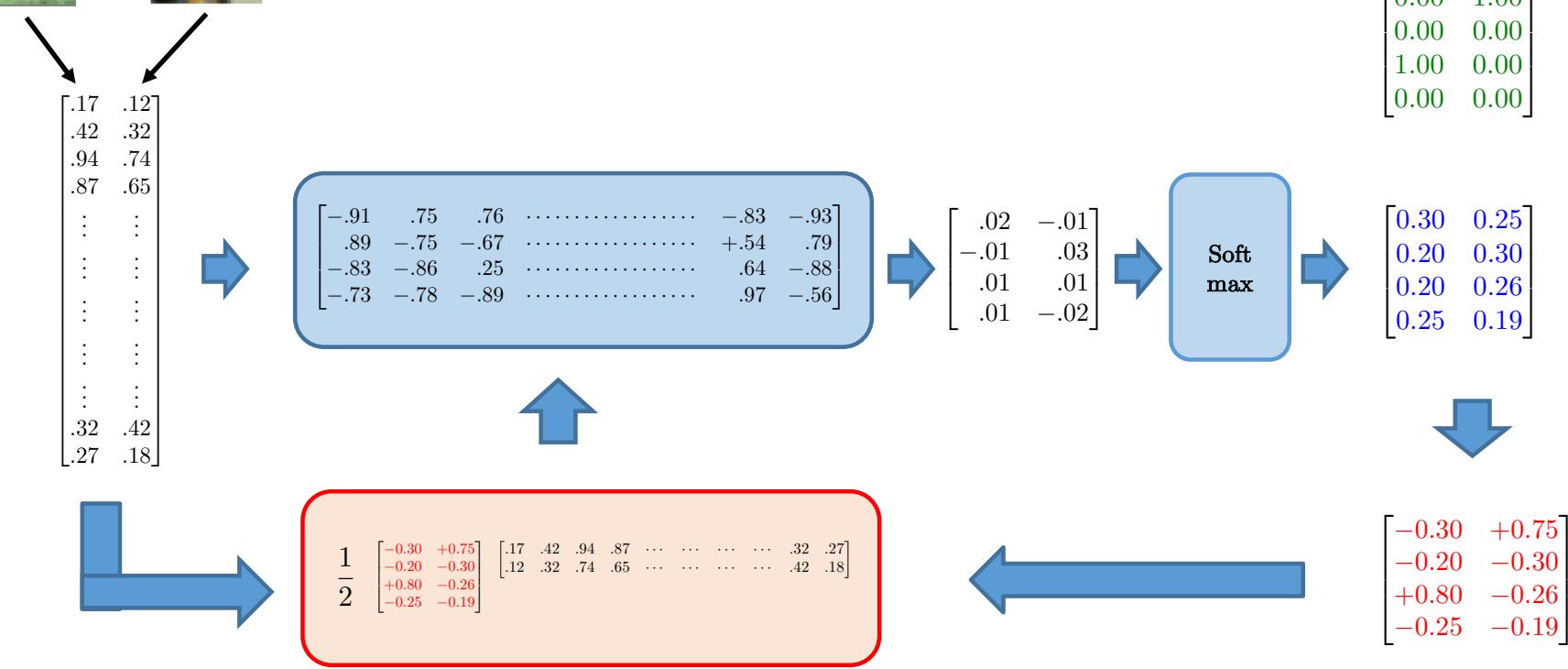
- Mini-batch GD = template update formula ! (No proof given)
  - The mini-batch gradient of the cross-entropy loss turns out to be exactly the same formula for template matching. (No coincidence)

$$\frac{\partial \mathcal{L}^B}{\partial W} = -\frac{1}{m} EX^T$$

- The template update scheme reduces to Mini-batch GD. (See next slide)
- When we did template matching update, we actually minimized the cross-entropy loss by Mini-batch GD.



## Connection to template matching



$$W = W - \text{lr} \circ \frac{\partial \mathcal{L}^B}{\partial W}$$

$$\frac{\partial \mathcal{L}^B}{\partial W} = -\frac{1}{m} E X^T$$

# Outline

- Quality of neural network
- Maximizing the quality by gradient ascent
- Learning rate strategy
- Neural network loss
- Minimizing the loss by gradient descent
- Mini-batch/stochastic gradient descent
- **Summary**

# Summary

Cross-entropy loss :  $\mathcal{L}(w_{11}, w_{12}, \dots) = \frac{1}{N} \sum_{i=1}^N \underbrace{\ell^{(i)}(w_{11}, w_{12}, \dots)}$

Each data point has a loss describing how well it has been classified by the network.

Gradient Descent :

$$\frac{\partial \mathcal{L}}{\partial W} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell^{(i)}}{\partial W}$$

Mini-batch Gradient Descent :

$$\frac{\partial \mathcal{L}^B}{\partial W} = \underbrace{\frac{1}{m} \sum_{i \in B} \frac{\partial \ell^{(i)}}{\partial W}}$$

Random approximation of the gradient using a batch  $B$   $B = \{i_1, i_2, i_3, \dots, i_m\}$

Doing **Mini-batch GD** is exactly the same thing than doing **template matching with mini-batch**.

Doing **GD** is exactly the same thing than doing **template matching with full batch**.

# Lab 01

- Cross-entropy loss:

jupyter cross\_entropy\_demo Last Checkpoint: a few seconds ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 Logout

In [1]:  
import torch  
import torch.nn as nn  
import utils

Make a Cross Entropy Criterion and call it mycrit

In [2]:  
mycrit=nn.CrossEntropyLoss()  
print(mycrit)  
CrossEntropyLoss()

Make a batch of labels

In [3]:  
labels=torch.LongTensor([2,3])  
print(labels)  
tensor([2, 3])

Make a batch of scores

jupyter cross\_entropy\_exercise Last Checkpoint: a few seconds ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 Logout

In [ ]:  
import torch  
import torch.nn as nn  
import utils

Make a Cross Entropy Criterion and call it criterion. The command is nn.CrossEntropyLoss().

In [ ]:  
criterion= # COMPLETE HERE  
print(criterion)

Suppose that there only two classes (class 0 and class 1).

Suppose we have a batch of three data points:

x<sup>(0)</sup> belongs to class 0

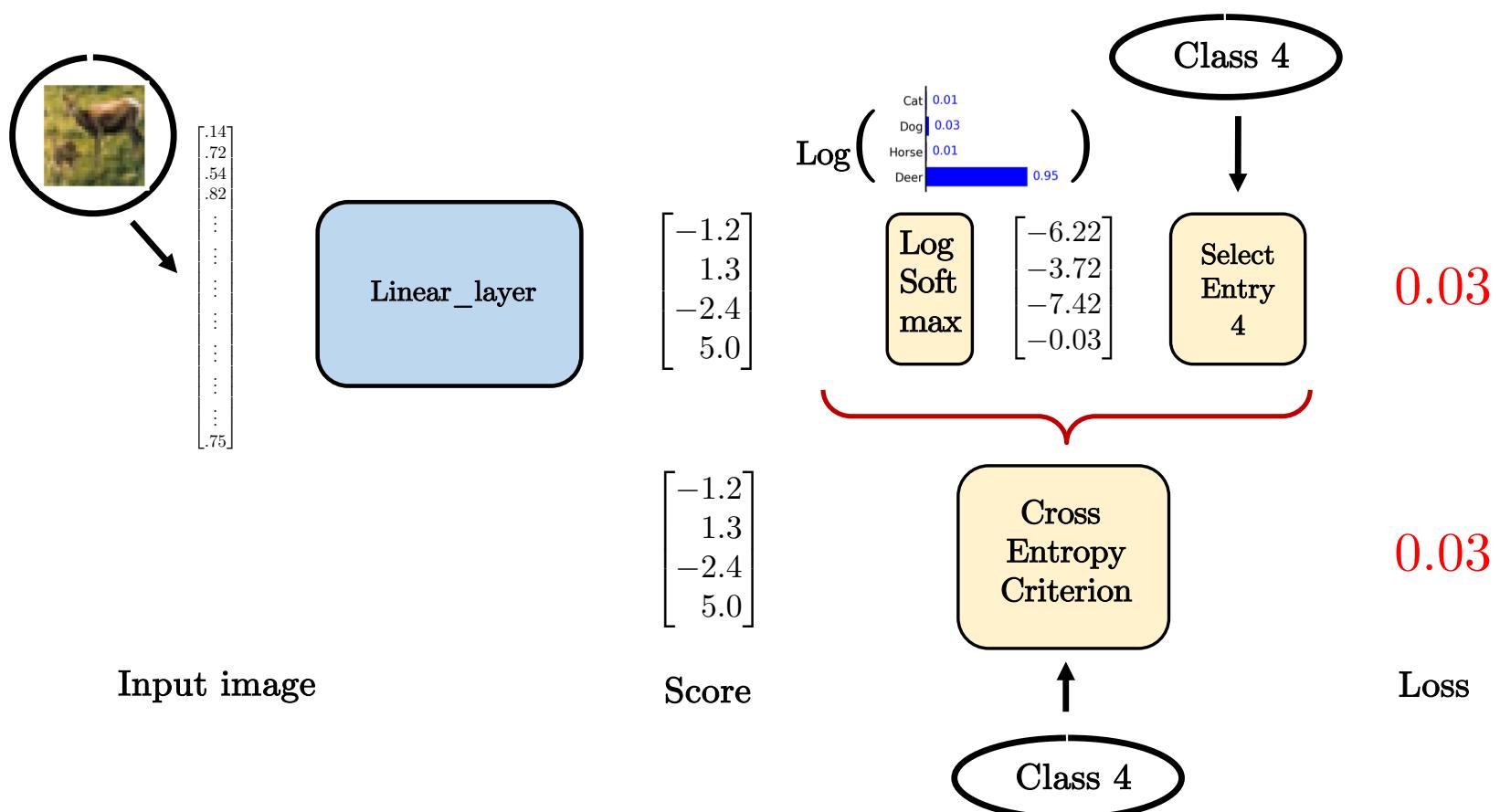
x<sup>(1)</sup> belongs to class 1

x<sup>(2)</sup> belongs to class 1

Put the labels of each of these point a LongTensor:

# Lab 01

- Function `nn.CrossEntropyLoss()` :





Questions?