

Computer Vision

Lecture 5: Training neural networks

Yuriy Kochura
iuriy.kochura@gmail.com
[@y_kochura](https://twitter.com/@y_kochura)

Today

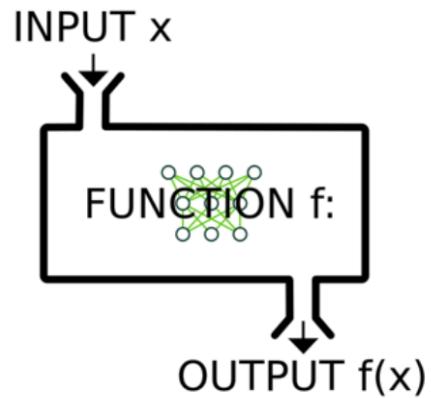
How to **optimize parameters** efficiently?

- Gradient descent
- Stochastic Gradient Descent
- Mini-batching
- Momentum
- Adaptive methods:
 - AdaGrad
 - RMSProp
 - Adam

Parameters optimization

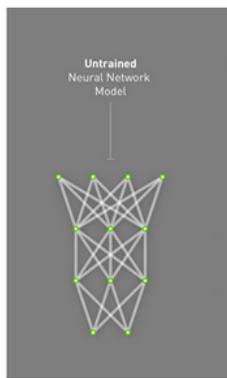
A model

While what is inside of a deep neural network may be complex, at their core, they are simply functions. They take some input and generate some output.



Components of a model

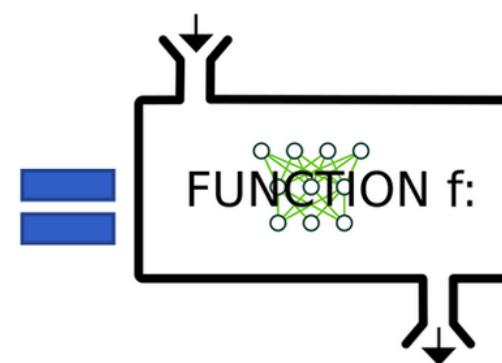
Model Architecture



Learned Weights



Model

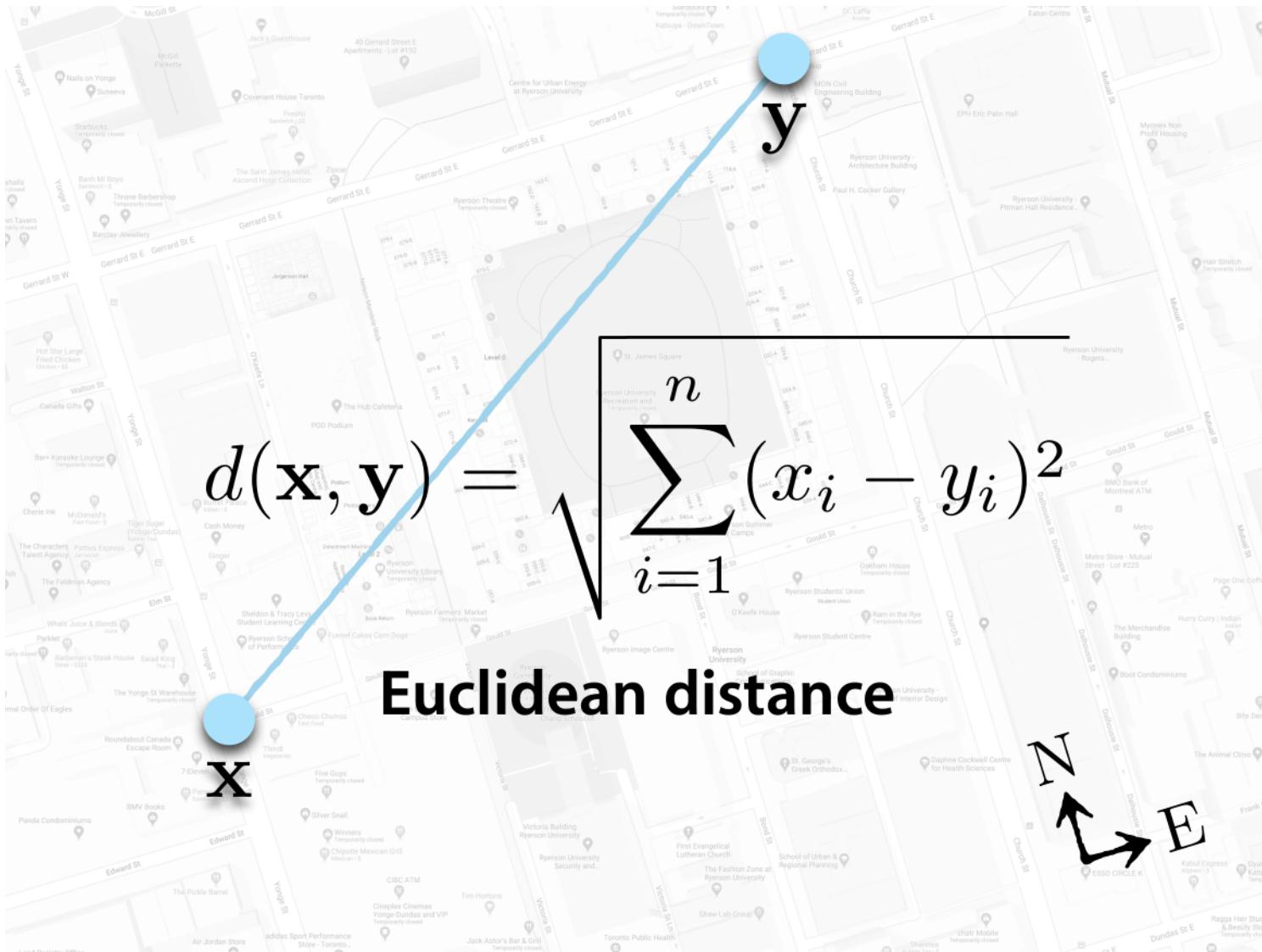


A common training process for neural networks

1. Define a task + collect data
2. Initialize the parameters
3. Choose an optimization algorithm
4. Repeat these steps:
 - 4.1. Forward propagate an input
 - 4.2 Compute the cost function
 - 4.3 Compute the gradients of the cost with respect to parameters using backpropagation
 - 4.4 Update each parameter using the gradients, according to the optimization algorithm

Common Loss functions

- Euclidean distance (L2 Loss)
- Mean Square Error (MSE)
- Manhattan distance (L1 Loss)
- Mean Absolute Error (MAE)
- Cross Entropy Loss



Euclidean distance (L2 Loss)

$$d(\hat{y}, y) = \sqrt{\sum_{i=1}^n (\hat{y}^{(i)} - y^{(i)})^2} = \|\hat{y} - y\|_2$$

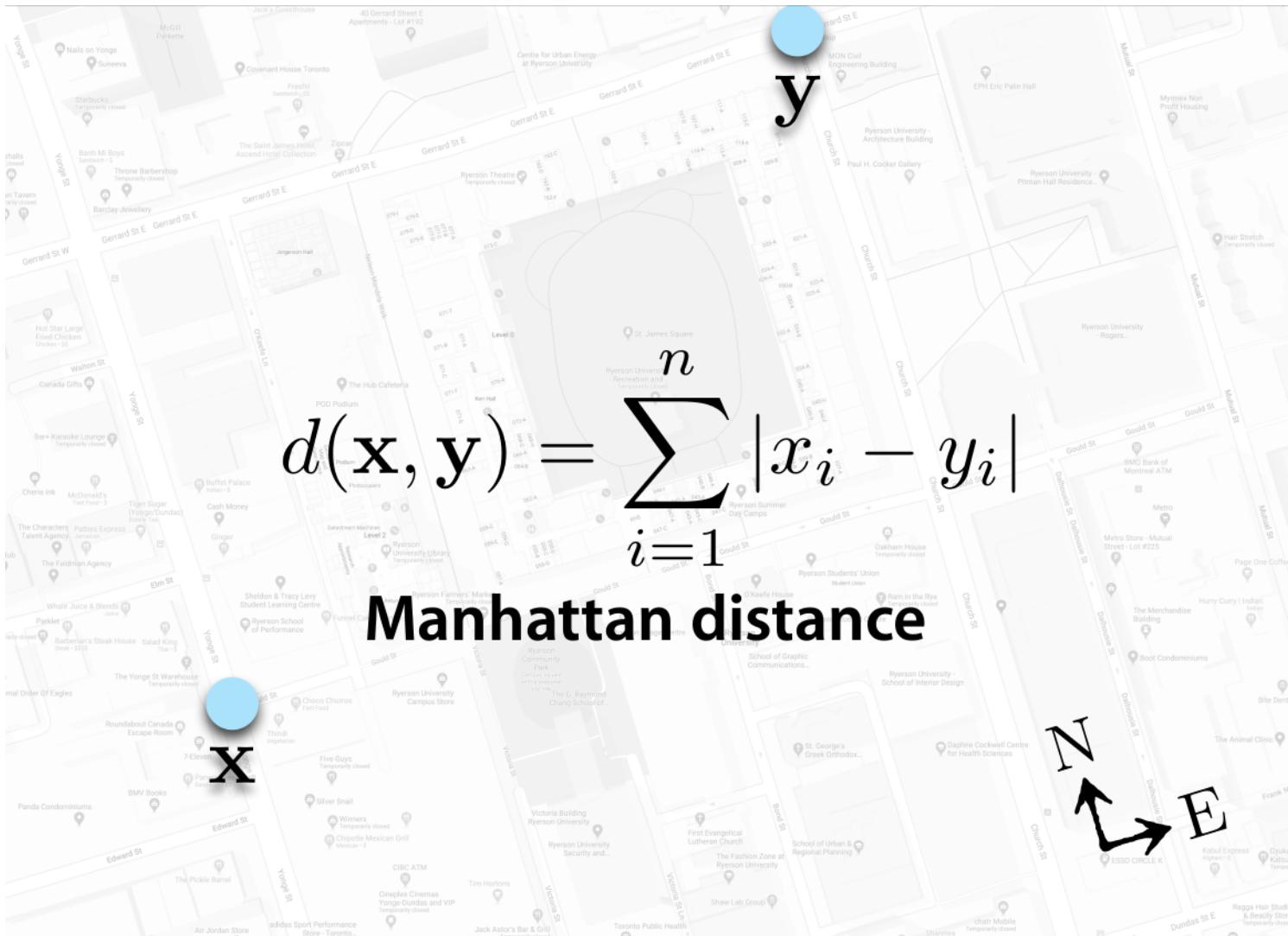
where n is a total number of training examples

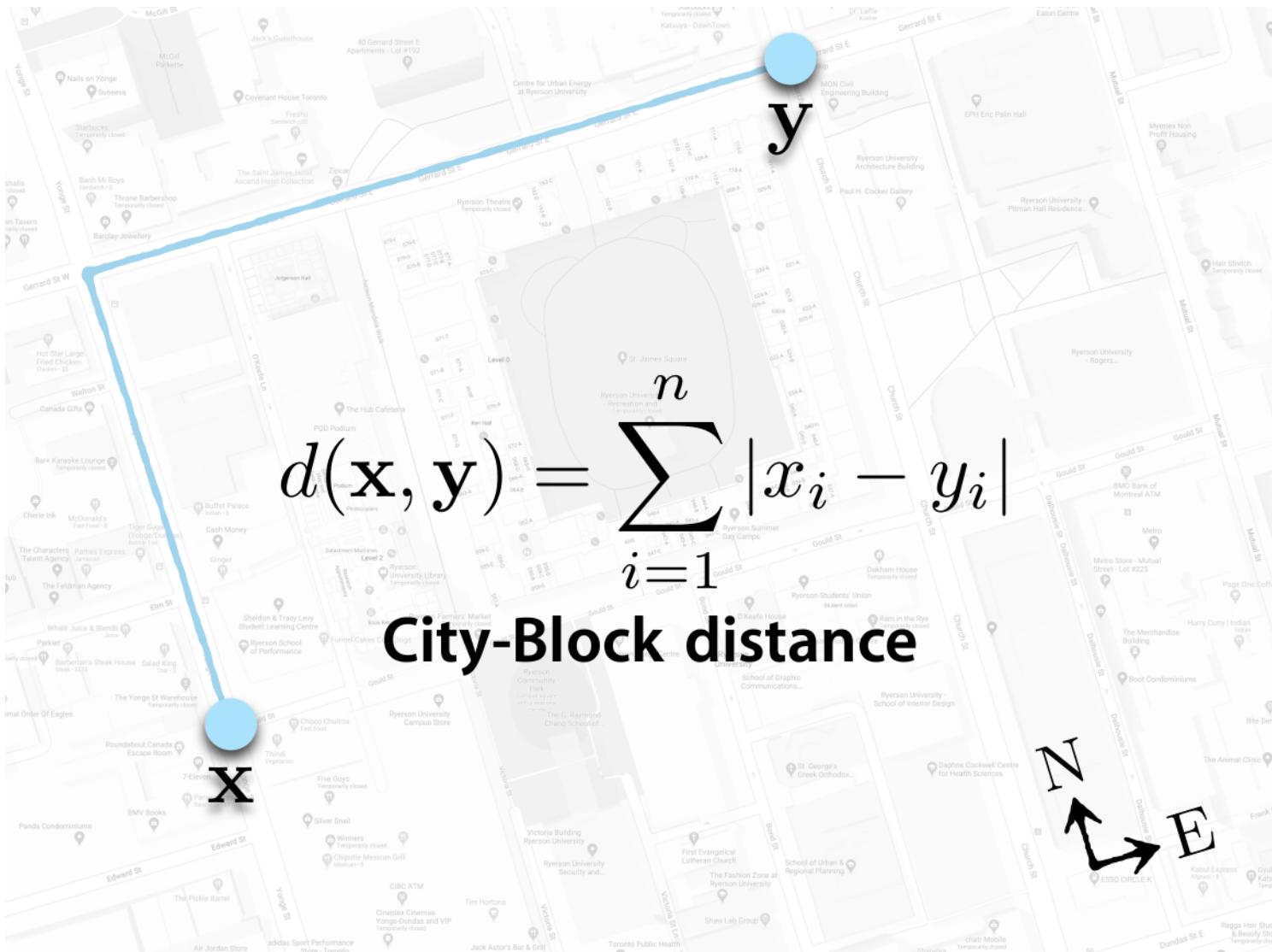
Mean Square Error (MSE)

$$\mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = (\hat{y}^{(i)} - y^{(i)})^2$$

$$\mathcal{J}(\hat{y}, y) = \frac{1}{n} \sqrt{\sum_{i=1}^n (\hat{y}^{(i)} - y^{(i)})^2} = \frac{1}{n} \|\hat{y} - y\|_2$$

where n is a total number of training examples





Manhattan distance (L1 Loss)

$$d(\hat{y}, y) = \sum_{i=1}^n |\hat{y}^{(i)} - y^{(i)}| = \|\hat{y} - y\|_1$$

where n is a total number of training examples

Mean Absolute Error (MAE)

$$\mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = |\hat{y}^{(i)} - y^{(i)}|$$

$$\mathcal{J}(\hat{y}, y) = \frac{1}{n} \sum_{i=1}^n |\hat{y}^{(i)} - y^{(i)}| = \frac{1}{n} \|\hat{y} - y\|_1$$

where n is a total number of training examples

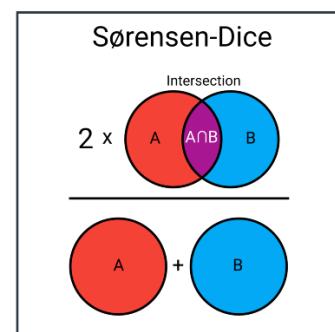
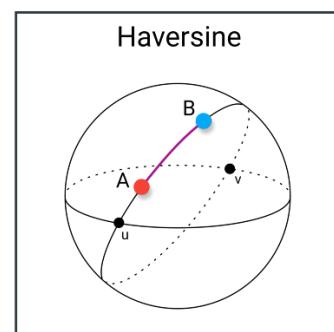
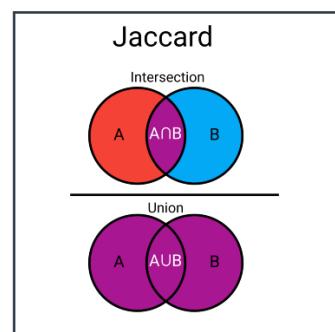
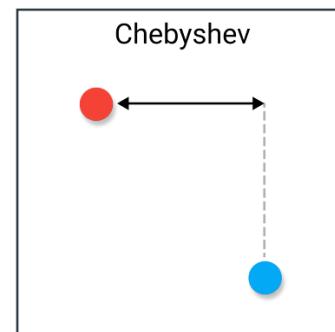
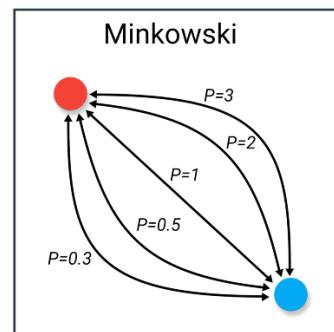
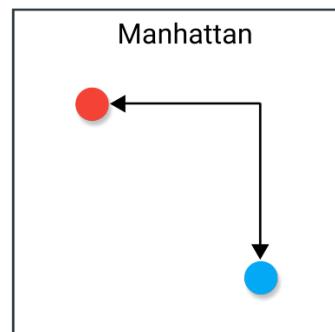
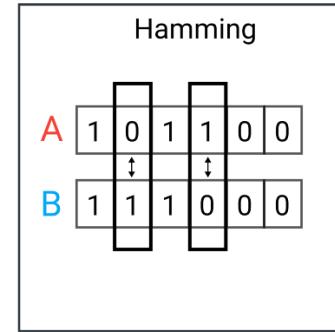
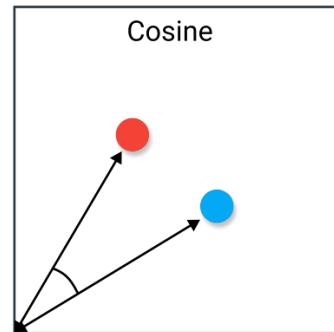
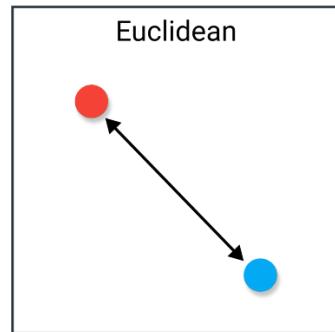
Cross Entropy Loss (binary)

$$\mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

$$\mathcal{J}(\hat{y}, y) = -\frac{1}{n} \sum_{i=1}^n \left[y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) \right]$$

where n is a total number of training examples

9 Distance Measures in Data Science



Optimization



Optimization algorithms

Hyperparameters & Parameters

- Result of optimization is a set of parameters
- Hyperparameters vs. Parameters

Hyperparameters

- Train-test split ratio
- Learning rate in optimization algorithms
- Number of hidden layers in a NN
- Number of activation units in each layer
- Batch size

:

Parameters

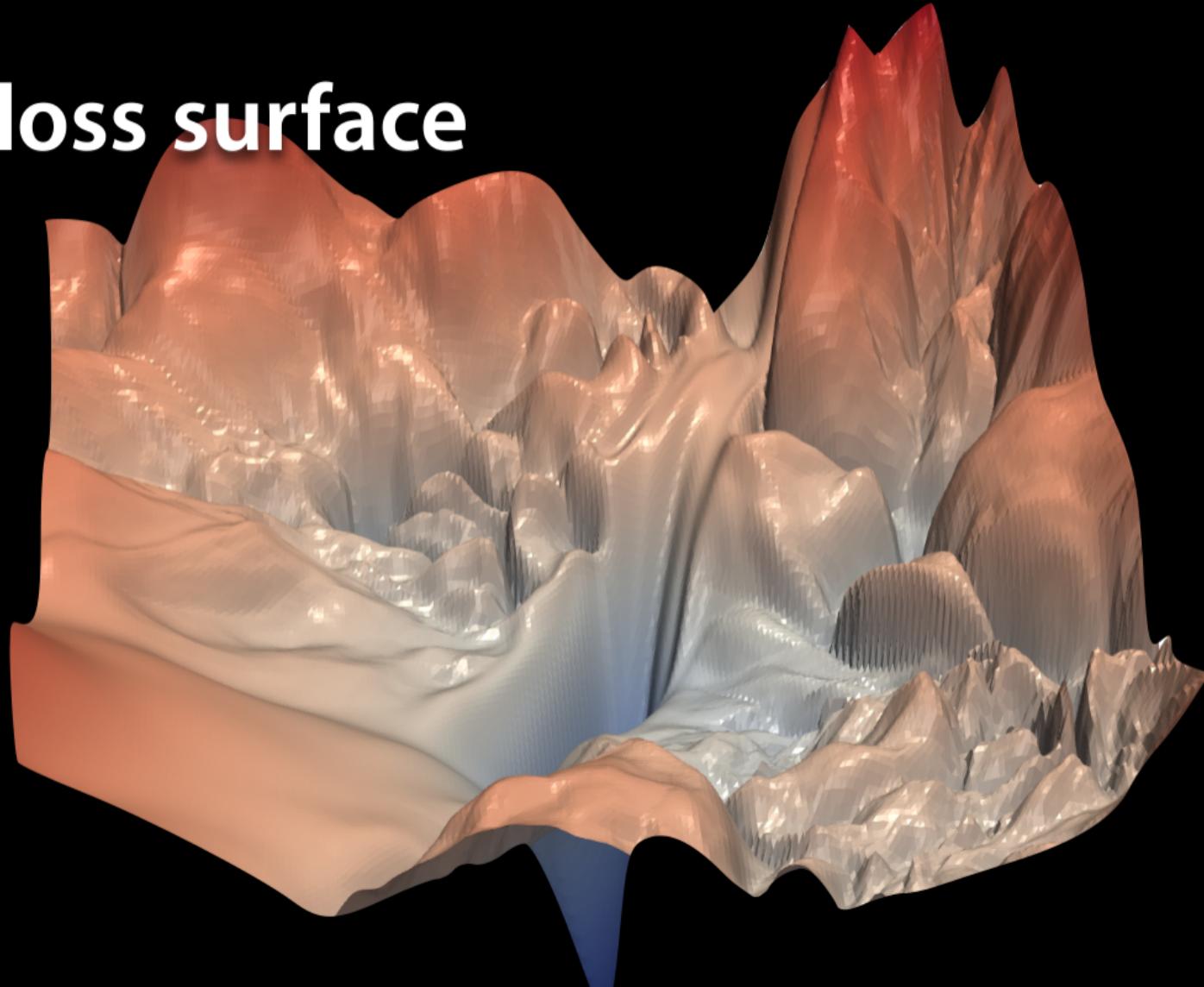
- Weights and biases of a NN
- The cluster centroids in clustering

Problem of optimization

Empirical risk minimization

$$W_*^d = \arg \min_W \mathcal{J}(W) = \frac{1}{n} \sum_{i=1}^n \mathcal{L} \left(y^{(i)}, f(\mathbf{x}^{(i)}, W) \right)$$

loss surface



A practical recommendation

Training a massive deep neural network is long, complex and sometimes confusing.

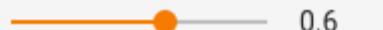
A first step towards understanding, debugging and optimizing neural networks is to make use of visualization tools for

- plotting losses and metrics,
- visualizing computational graphs,
- or showing additional data as the network is being trained.

Show data download links Ignore outliers in chart scalingTooltip sorting
method:

default

Smoothing



Horizontal Axis

STEP

RELATIVE

WALL

Runs

Write a regex to filter runs

 train eval

TOGGLE ALL RUNS

/tmp/mnist-logs

 Filter tags (regular expressions supported)

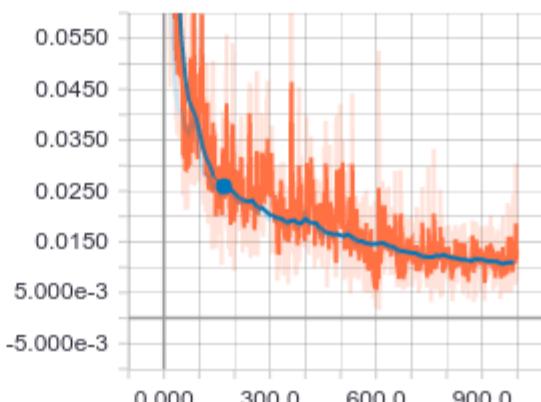
accuracy

1

cross entropy

1

cross entropy



run to download

CSV JSON

Name	Smoothed	Value	Step	Time	Relative
eval	0.02591	ep_p0.02550	ity170.0	Mon Sep 12, 15:40:41	8s
train	0.02851	0.03362	166.0	Mon Sep 12, 15:40:40	7s

mean

4

Demo

Gradient Descent

GD

Gradient descent

To minimize $\mathcal{J}(W)$, standard **batch gradient descent** (GD) consists in applying the update rule

$$g_t = \frac{1}{n} \sum_{i=1}^n \nabla_W \mathcal{L} \left(y^{(i)}, f(\mathbf{x}^{(i)}, W) \right) = \nabla_W \mathcal{J}(W)$$

$$W_{t+1} = W_t - \alpha g_t,$$

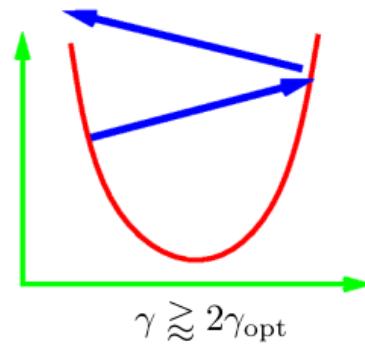
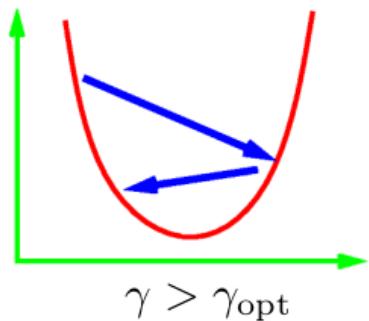
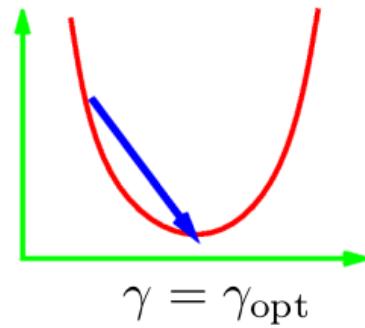
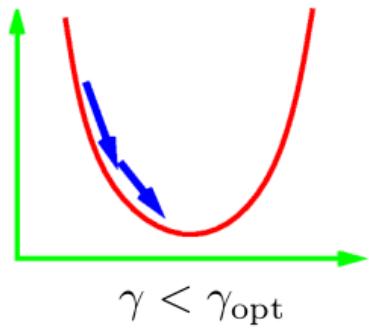
where α is the learning rate.

GD

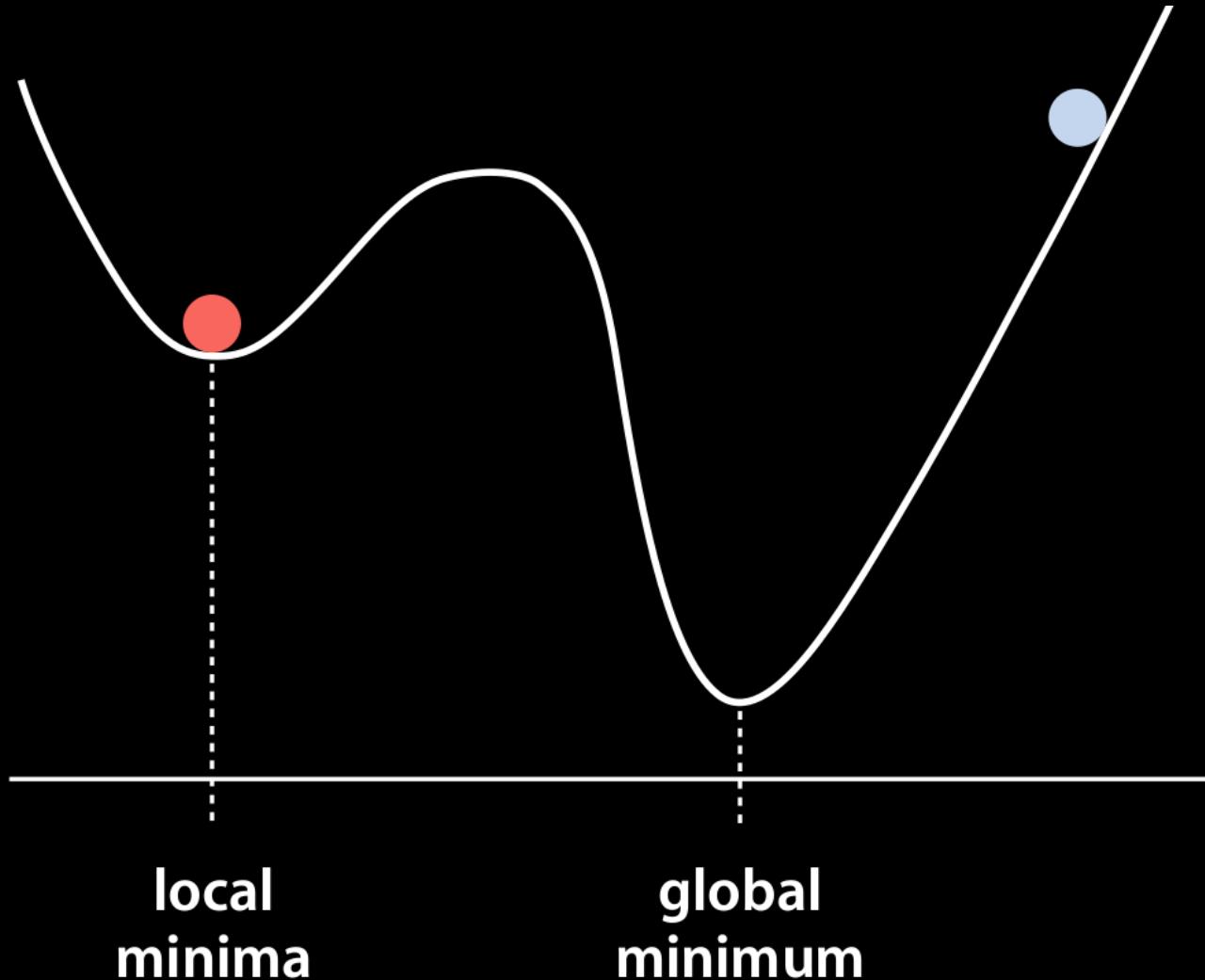
The worst optimization method in the world

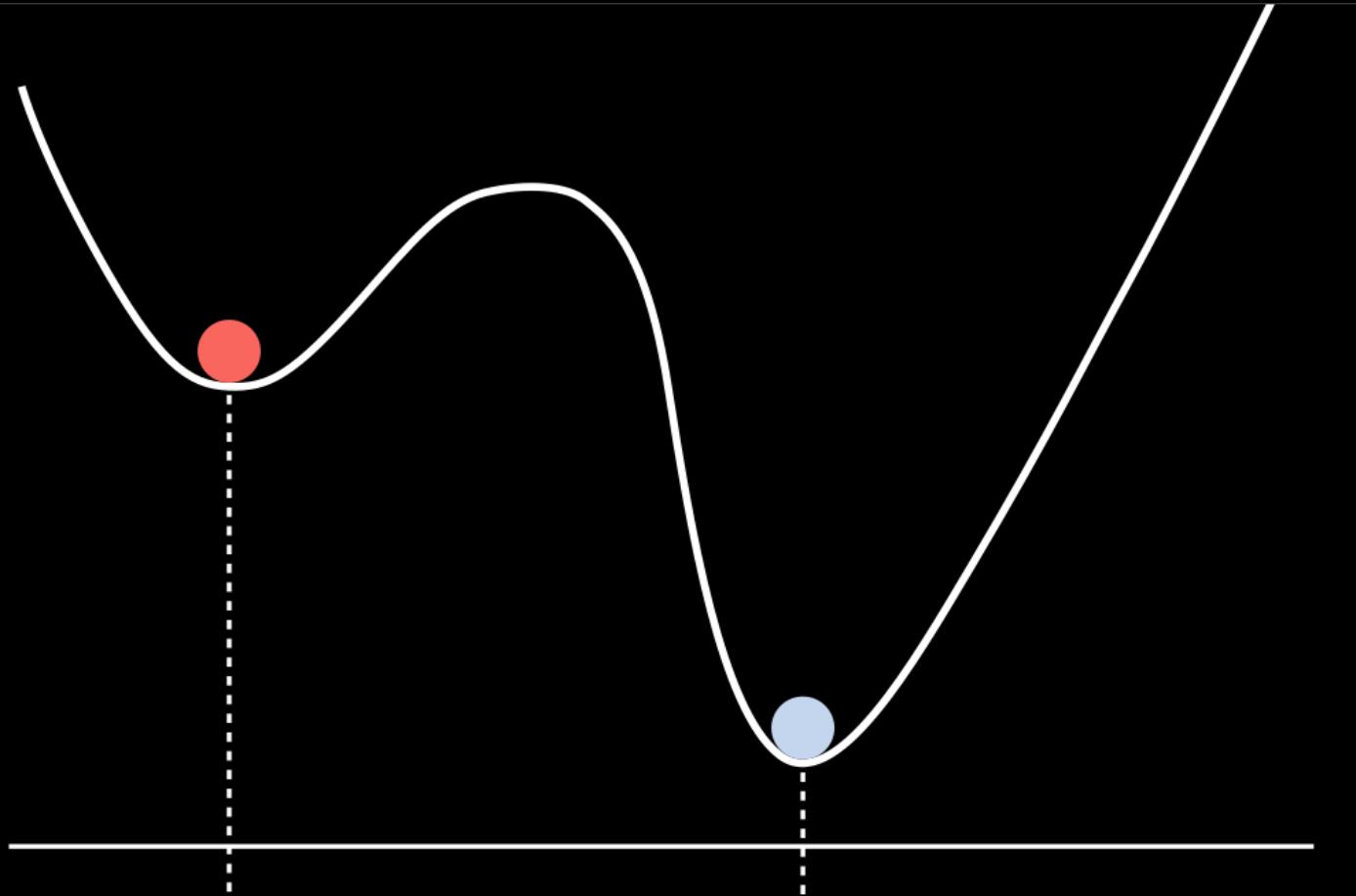
Side note: You should probably never use gradient descent directly, consider it a building block for other methods.

Learning Rate



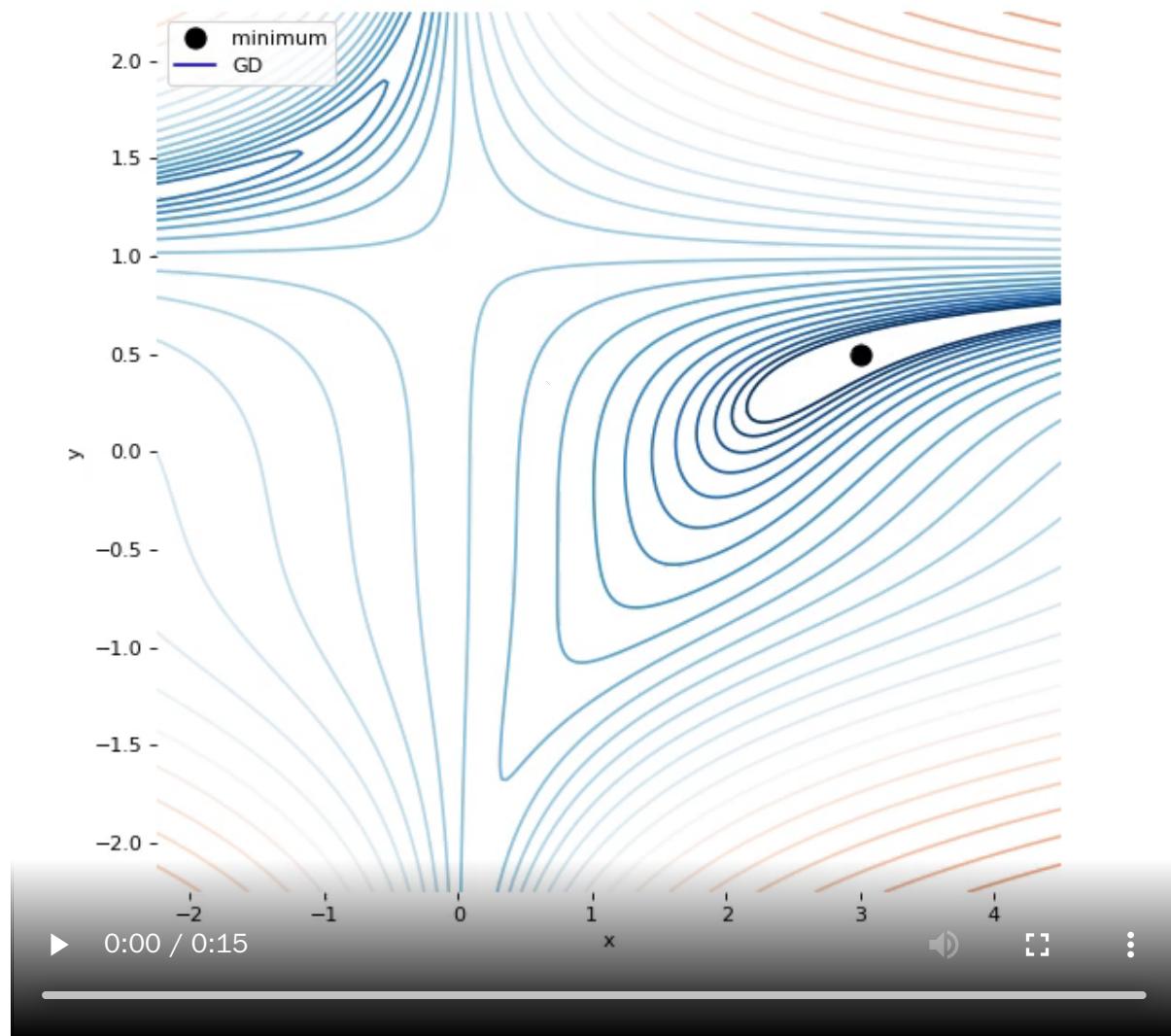
Here $\gamma = \alpha$ is the learning rate.





local
minima

global
minimum



Stochastic Gradient Descent

SGD

SGD

To reduce the computational complexity, **stochastic gradient descent** (SGD) consists in updating the parameters after every sample

$$\ell^{(i)} = \mathcal{L} \left(y^{(i)}, f(\mathbf{x}^{(i)}, W) \right)$$

$$g_t^{(i)} = \nabla_W \ell^{(i)}$$

$$W_{t+1} = W_t - \alpha g_t^{(i)}$$

Advantages of SGD

- Stochastic gradients are drastically cheaper to compute (proportional to your dataset size), so you can often take thousands of SGD steps for the cost of one GD step.
- On average, the stochastic gradient is a good estimate of the gradient.
- The noise can prevent the optimizing converging to bad local minima, a phenomena called annealing.

Mini-batching

Mini-batching

Visiting the samples in mini-batches and updating the parameters each time

$$\begin{aligned} g_t^{(k)} &= \frac{1}{B} \sum_{i=1}^B \nabla_W \mathcal{L} \left(y_k^{(i)}, f(\mathbf{x}_k^{(i)}, W) \right) \\ W_{t+1} &= W_t - \alpha g_t^{(k)}, \end{aligned}$$

where k is an index of mini-batch.

- Increasing the batch size B reduces the variance of the gradient estimates and enables the speed-up of batch processing.
- The interplay between B and α is still unclear.

Momentum

Momentum

SGD + Momentum = Stochastic **heavy ball** method

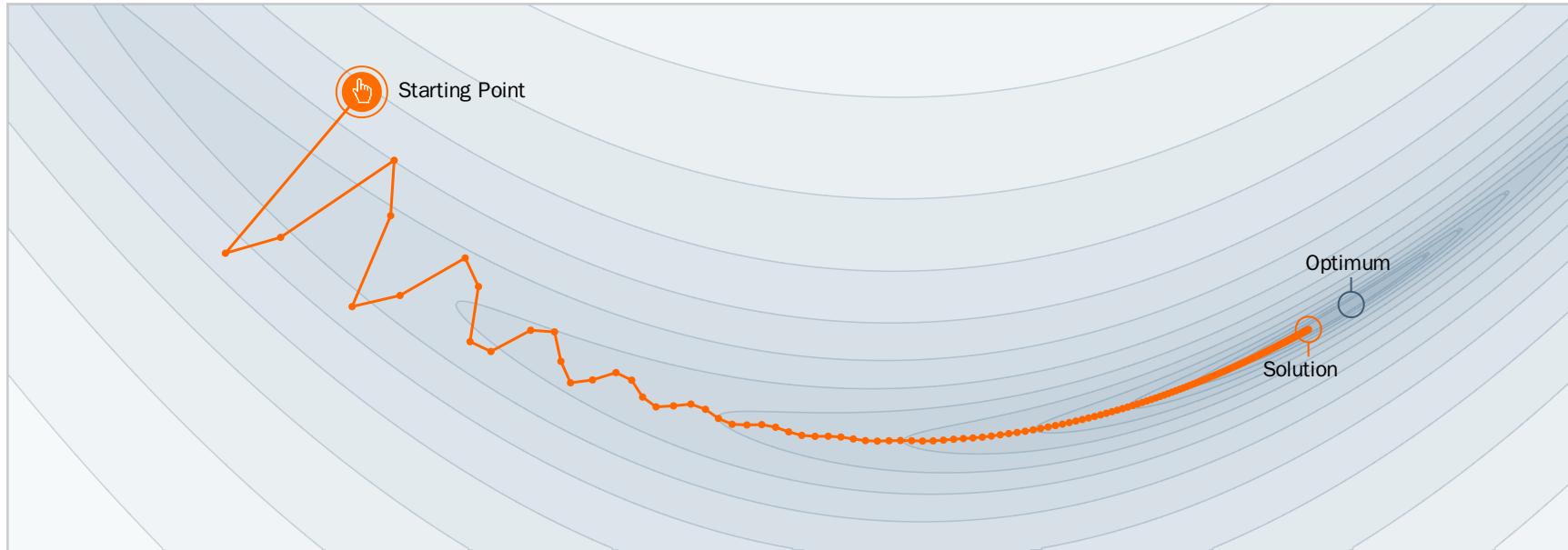
$$p_{t+1} = \beta_t p_t + \nabla \ell^{(i)}(W_t)$$

$$W_{t+1} = W_t - \alpha_t p_{t+1}$$

Rule for updating:

$$W_{t+1} = W_t - \alpha_t \nabla \ell^{(i)}(W_t) + \beta_t (W_t - W_{t-1})$$

Key idea: The next step becomes a combination of the previous step's direction and the new negative gradient.



Step-size $\alpha = 0.02$



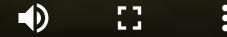
Momentum $\beta = 0.99$



We often think of Momentum as a means of dampening oscillations, speeding up the iterations, leading to faster convergence. But interesting behavior. It allows a larger range of step-sizes to be used, but creates its own oscillations. What is going on?

SGD with Momentum

▶ 0:00 / 0:46



Advantages

SGD with momentum has **three** nice properties:

- it can go through local barriers
- it accelerates if the gradient does not change much
- it dampens oscillations in narrow valleys

Practical Aspects of momentum

It's basically “**free lunch**”, in almost all situations, **SGD + momentum** is better than SGD, and very rarely worse!

Recommended Parameters:

$\beta = 0.9$ or 0.99 almost always works well. Sometimes slight gains can be had by tuning it.

The step size parameter (α) usually needs to be decreased when the momentum parameter is increased to maintain convergence.

Demo

Adaptive methods

Adaptive methods

The magnitude of the gradients often varies highly between layers due to, so a global learning rate may not work well.

General IDEA: Instead of using the same learning rate for every weight in our network, **maintain an estimate of a better rate separately for each weight.**

The exact way of adapting to the learning rates varies between algorithms, but most methods either **adapt to the variance of the weights**, or to the **local curvature** of the problem.

AdaGrad

Per-parameter downscale by square-root of sum of squares of all its historical values.

$$\begin{aligned} r_t &= r_{t-1} + g_t \odot g_t \\ W_{t+1} &= W_t - \frac{\alpha}{\varepsilon + \sqrt{r_t}} \odot g_t \end{aligned}$$

- AdaGrad eliminates the need to manually tune the learning rate. Most implementation use $\alpha = 0.01$ as default.
- It is good when the objective is convex.
- r_t grows unboundedly during training, which may cause the step size to shrink and eventually become infinitesimally small.
- δ is an additive constant that ensures that we do not divide by 0.

RMSProp

Same as AdaGrad but accumulate an exponentially decaying average of the gradient.

Key IDEA: normalize by the root-mean-square of the gradient

$$r_t = \rho r_{t-1} + (1 - \rho) g_t \odot g_t$$
$$W_{t+1} = W_t - \frac{\alpha}{\varepsilon + \sqrt{r_t}} \odot g_t$$

- Perform better in non-convex settings.

Adam: RMSprop with a kind of momentum

“Adaptive Moment Estimation”

Similar to RMSProp with momentum, but with bias correction terms for the first and second moments.

$$p_t = \rho_1 p_{t-1} + (1 - \rho_1) g_t$$

$$\hat{p}_t = \frac{p_t}{1 - \rho_1^t}$$

$$r_t = \rho_2 r_{t-1} + (1 - \rho_2) g_t \odot g_t$$

$$\hat{r}_t = \frac{r_t}{1 - \rho_2^t}$$

$$W_{t+1} = W_t - \alpha \frac{\hat{p}_t}{\varepsilon + \sqrt{\hat{r}_t}}$$

- Good defaults are $\rho_1 = 0.9$ and $\rho_2 = 0.999$.
- Just as momentum improves SGD, it improves RMSProp as well.

Practical side

For poorly conditioned problems, Adam is often much better than SGD.

Use Adam over RMSprop due to the clear advantages of momentum.

BUT, Adam is poorly understood theoretically, and has known disadvantages:

- Does not converge at all on some simple example problems!
- Gives worse generalization error on many computer vision problems (i.e. ImageNet)
- Requires more memory than SGD
- Has 2 momentum parameters, so some tuning may be needed

Demo

demo - losslandscape

The end