

MACHINE LEARNING

Loss functions & training

Corso di Laurea Magistrale in Informatica

Università di Roma Tor Vergata

Giorgio Gambosi

a.a. 2023-2024



LOSS FUNCTION

- In general, the loss function $L : \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}$ measures, for any two values y, t in target space, the **cost** of referring, for any subsequent action, to t instead of the better value y
- In supervised learning, it provides a measure of the quality of the prediction returned by the prediction function h

$$\mathcal{R}(x, y) = L(h(x), y)$$

- It is a fundamental component of the empirical risk, which is just the average value of the loss function applied to all predicted value - target value pairs in the training set \mathcal{T}

$$\bar{\mathcal{R}}_{\mathcal{T}}(h) = \frac{1}{|\mathcal{T}|} \sum_{(x,t) \in \mathcal{T}} L(h(x), t)$$

- That is, it provides a measure of the quality of the predictions performed by h , at least with respect to the available data (the training set)

LOSS FUNCTION & TRAINING

- During the training phase, the empirical risk is minimized wrt the prediction function h applied, and in particular to the set of parameters θ which specifies the parametric function $h = h_\theta$
- This corresponds to minimizing the overall loss

$$\mathcal{L}(\theta; \mathcal{T}) = \sum_{i=1}^n L_i(\theta)$$

that is the sum of the loss functions $L_i = L(\theta; \mathbf{x}_i, y_i)$

LOSS FUNCTION MINIMIZATION

How to deal with loss minimization?

- we would like to compute a **global minimum**
- methods based on calculus rely on setting all derivatives to 0, that is,

$$\nabla_{\theta} \mathcal{L}(\theta; \mathcal{T}) = \mathbf{0}$$

that is

$$\frac{\partial}{\partial \theta_i} \nabla_{\theta} \mathcal{L}(\theta; \mathcal{T}) = 0 \quad \forall i$$

and solve the corresponding system of equations

Problems

- the system of equations has multiple solutions (local minima/maxima, saddle points)
- they can be hard (or impossible) to compute analytically

LOSS FUNCTION MINIMIZATION

- A local minimum of $\overline{\mathcal{R}}_{\mathcal{T}}(\boldsymbol{\theta})$ can be computed numerically, by means of iterative methods such as **gradient descent**
- Initial assignment $\boldsymbol{\theta}^{(0)} = (\theta_0^{(0)}, \theta_1^{(0)}, \dots, \theta_d^{(0)})$, with a corresponding error value

$$\overline{\mathcal{R}}_{\mathcal{T}}(\boldsymbol{\theta}^{(0)})$$

- Iteratively, the current value $\boldsymbol{\theta}^{(i-1)}$ is modified in the direction of **steepest descent** of $\overline{\mathcal{R}}_{\mathcal{T}}(\boldsymbol{\theta})$, that is the one corresponding to the negative of the gradient evaluated at $\boldsymbol{\theta}^{(i-1)}$
- At step i , $\theta_j^{(i-1)}$ is updated as follows:

$$\theta_j^{(i)} := \theta_j^{(i-1)} - \eta \left. \frac{\partial}{\partial \theta_j} \overline{\mathcal{R}}_{\mathcal{T}}(\boldsymbol{\theta}) \right|_{\boldsymbol{\theta}^{(i-1)}} = \theta_j^{(i-1)} - \frac{\eta}{|\mathcal{T}|} \sum_{(\mathbf{x}, \mathbf{t}) \in \mathcal{T}} \left. \frac{\partial}{\partial \theta_j} L(h_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{t}) \right|_{\boldsymbol{\theta}^{(i-1)}}$$

GRADIENT DESCENT

- In matrix notation:

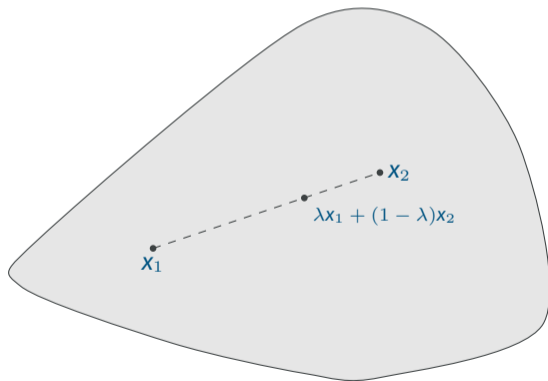
$$\boldsymbol{\theta}^{(i)} := \boldsymbol{\theta}^{(i-1)} - \eta \nabla_{\boldsymbol{\theta}} \bar{\mathcal{R}}_{\mathcal{T}}(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}^{(i-1)}} = \theta_j^{(i-1)} - \frac{\eta}{|\mathcal{T}|} \sum_{(\mathbf{x}, \mathbf{t}) \in \mathcal{T}} \nabla_{\theta} \bar{\mathcal{R}}_{\mathcal{T}}(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}^{(i-1)}}$$

- clearly this approach makes it possible to find (approximate) a local minimum, depending from the initial values; some problems
 - we are looking for a global (not simply a local) minimum
 - how to deal with saddle points?
 - how fast does the method converge?
- More on this later

CONVEXITY

A set of points $S \subset \mathbb{R}^d$ is **convex** iff for any $x_1, x_2 \in S$ and $\lambda \in (0, 1)$

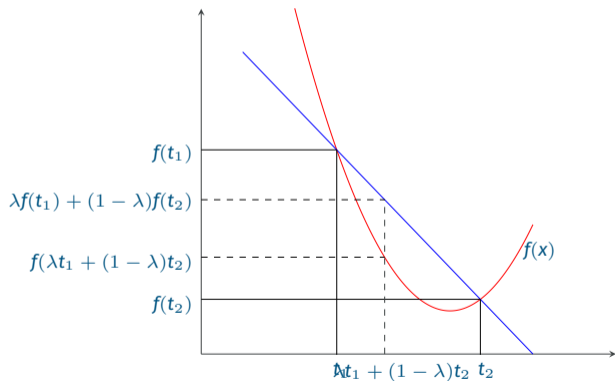
$$\lambda x_1 + (1 - \lambda)x_2 \in S$$



CONVEXITY

A function $f(\mathbf{x})$ is convex iff the set of points lying above the function is convex, that is, for all $\mathbf{x}_1, \mathbf{x}_2$ and $\lambda \in (0, 1)$,

$$f(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2) \leq \lambda f(\mathbf{x}_1) + (1 - \lambda) f(\mathbf{x}_2)$$



CONVEXITY

- Assuming $\mathcal{L}(\theta; \mathcal{T})$ is convex is a relevant simplification: if $f(\mathbf{x})$ is a convex function, then any local minimum of f is also a global minimum
- Moreover, if f is a **strictly convex** function, there exists only one local minimum for f (and it is global), that is, solving

$$\nabla_{\theta} \mathcal{L}(\theta; \mathcal{T}) = \mathbf{0}$$

provides the global minimum

- Definition: $f(\mathbf{x})$ is strictly convex iff for all $\mathbf{x}_1, \mathbf{x}_2$ and $\lambda \in (0, 1)$,

$$f(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2) < \lambda f(\mathbf{x}_1) + (1 - \lambda) f(\mathbf{x}_2)$$

- A simple but relevant case: $f(\mathbf{x})$ is quadratic. This is the case for a number of simpler ML models. Unfortunately this is not true for more complex models such as neural networks

CONVEXITY AND EMPIRICAL RISK

- convex functions properties:
 - the sum of (strictly) convex functions is (strictly) convex
 - the product of a (strictly) convex function and a constant is (strictly) convex

- since

$$\bar{\mathcal{R}}_{\mathcal{T}}(h) = \frac{1}{|\mathcal{T}|} \sum_{(\mathbf{x}, t) \in \mathcal{T}} L(h(\mathbf{x}), t) \propto \sum_{(\mathbf{x}, t) \in \mathcal{T}} L(\boldsymbol{\theta}; \mathbf{x}, t)$$

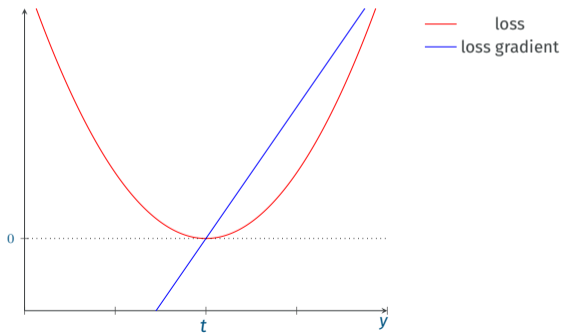
- if $L(\boldsymbol{\theta}; \mathbf{x}, t)$ is (strictly) convex then the overall cost is also (strictly) convex
- if $L(\boldsymbol{\theta}; \mathbf{x}, t)$ is convex then any local minimum of the empirical risk is also a global one
- if $L(\boldsymbol{\theta}; \mathbf{x}, t)$ is strictly convex then there exists only one minimum of the empirical risk

SOME COMMON LOSS FUNCTIONS

Let us first consider the case of **regression**.

- both y and $h(\mathbf{x})$ are real values
- loss is related to some type of point distance measure
- most common loss function for regression: **quadratic loss**

$$L(y, t) = (y - t)^2$$



QUADRATIC LOSS

- Applying quadratic loss results in the empirical risk

$$\overline{\mathcal{R}}_{\mathcal{T}}(h) = \frac{1}{|\mathcal{T}|} \sum_{(x,t) \in \mathcal{T}} (h(x) - t)^2$$

- in the common case of linear regression, the prediction is performed by means of a linear function $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$: this results into an overall loss to be minimized

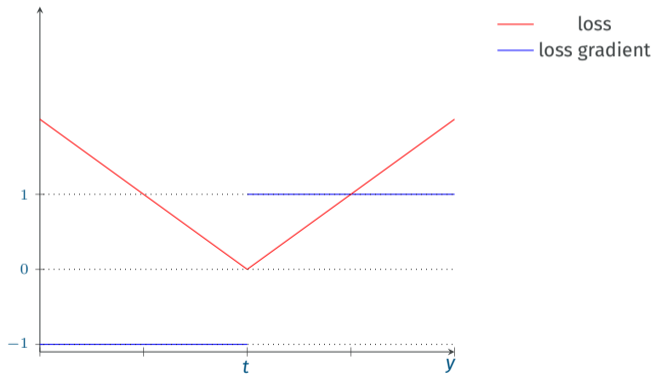
$$\mathcal{L}(\mathbf{w}, b; \mathcal{T}) = \sum_{(x,t) \in \mathcal{T}} (\mathbf{w}^T \mathbf{x} + b - t)^2$$

- since the quadratic function is strictly convex, the overall loss has only one local minimum (which is global)
- the gradient is linear

$$\frac{\partial}{\partial w_i} \mathcal{L}(\mathbf{w}, b; \mathcal{T}) = \sum_{(x,t) \in \mathcal{T}} (\mathbf{w}^T \mathbf{x} + b - t) w_i \qquad \frac{\partial}{\partial b} \mathcal{L}(\mathbf{w}, b; \mathcal{T}) = \sum_{(x,t) \in \mathcal{T}} (\mathbf{w}^T \mathbf{x} + b - t)$$

ADDITIONAL LOSS FUNCTIONS FOR REGRESSION: ABSOLUTE LOSS

- Quadratic loss is easy to deal with mathematically, but not robust to outliers, i.e. pays too much attention to outliers.
- A different loss function: **absolute loss** $L(t, y) = |t - y|$

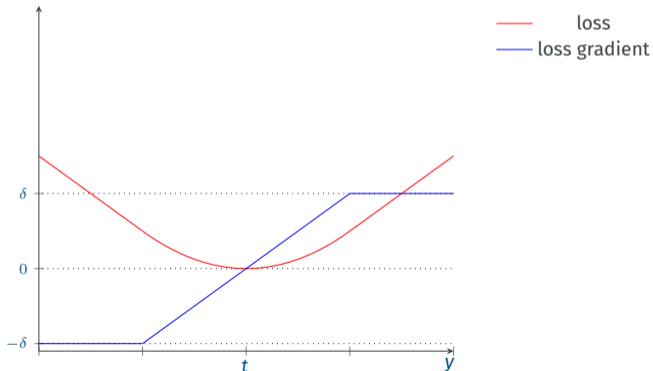


- The gradient is piecewise constant

ADDITIONAL LOSS FUNCTIONS FOR REGRESSION: HUBER LOSS

- Another different loss function: **Huber loss**

$$L(t, y) = \begin{cases} \frac{1}{2}(t - y)^2 & |t - y| \leq \delta \\ \delta(|t - y|) - \frac{\delta^2}{2} & |t - y| > \delta \end{cases}$$



LOSS FUNCTIONS FOR CLASSIFICATION

- Essentially, two approaches, depending on what we expect the prediction return:
 - prediction returns a specific class (prediction function)
 - prediction returns a probability distribution on the set of classes (prediction distribution)
- Different definition of error
 - first case: coincidence of predicted and real classes
 - second case: cumulative difference between predicted probability and 0/1 for all classes
- We consider the binary case, with two classes identified by target values -1 and 1
- Assume a real value is returned as a prediction

0/1 LOSS

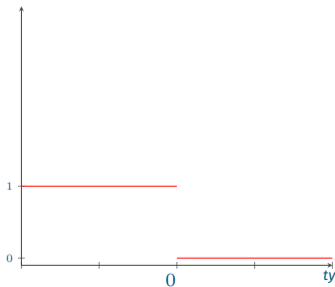
- The most “natural” loss function in classification

$$L(t, y) = \begin{cases} 1 & \text{sgn}(t) \neq y \\ 0 & \text{sgn}(t) = y \end{cases}$$

where $\text{sgn}(x)$ is 1 if $x > 0$ and -1 otherwise.

- This can be written as:

$$\mathbf{1}[ty < 0]$$



0/1 LOSS

Problem:

- not convex
- not smooth (first derivative undefined in some points or not continue)
- gradient is 0 almost everywhere (undefined at 0): gradient descent cannot be applied
- if we assume a linear prediction function

$$\bar{\mathcal{R}}_{\mathcal{T}}(h) = \frac{1}{|\mathcal{T}|} \sum_{(x,t) \in \mathcal{T}} \mathbf{1}[(\mathbf{w}^T \mathbf{x} + b)y < 0]$$

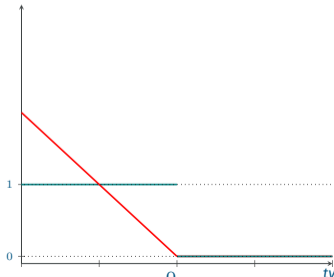
- the problem is finding the values \mathbf{w} , b which minimize the overall number of errors: this is an NP-hard, hence a computationally intractable problem.

CONVEX SURROGATE LOSS FUNCTIONS

- Approximate **from above** 0/1 loss: real 0/1 error always less than function loss
- Convex: unique local minimum = global minimum
- Smooth: may use derivatives to find minimum
- Main difference: relevance given to erroneous predictions

PERCEPTRON LOSS

- 0/1 loss assigns the same cost 1 to each error
- assume a prediction t is a real value: then, in the case of a misclassified element, the error can be measured as $-ty > 0$. That is, $L(t, y) = \max(0, -yt)$
- in the case of correctly classified element, the error is 0, while in the case of a wrong prediction, the error is equal to $|t|$
- Main difference: relevance given to erroneous predictions. The perceptron loss penalizes prediction which are largely wrong (for example a negative value $\simeq -1$ while correct class is 1)
- continuous, gradient continuous almost everywhere, convex (but not strictly convex), not surrogate

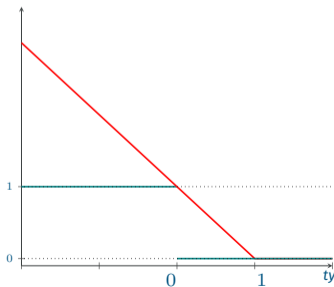


HINGE LOSS

- used in support vector machine training
- related to perceptron loss, but surrogate
- assume a prediction

$$L(t, y) = \max(0, 1 - yt)$$

- correct predictions can be penalized if “weak” (small value of t)
- continuous, gradient continuous almost everywhere, convex (but not strictly convex), surrogate



HINGE LOSS

Hinge loss $L_H(\mathbf{y}, t) = \max(0, 1 - \mathbf{y}t)$ is not differentiable wrt to \mathbf{y} at $\mathbf{y}t = 1$. The same holds for perceptron loss at $\mathbf{y}t = 0$.

For example,

$$\frac{\partial}{\partial \mathbf{y}} L_H = \begin{cases} -t & \mathbf{y}t < 1 \\ 0 & \mathbf{y}t > 1 \\ \text{undefined} & \mathbf{y}t = 1 \end{cases}$$

This is a problem if gradient descent should be applied. In this case a subgradient can be used.

SUBGRADIENT

Given a convex function (such as hinge loss) f at each differentiable point, the corresponding gradient $\nabla(x)$ provides a function which lower bounds f

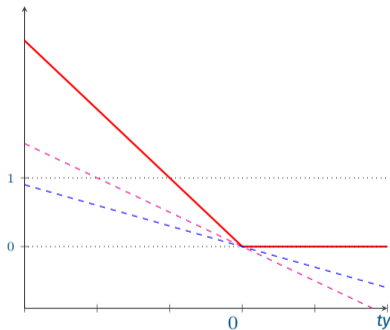
$$f(x') \geq f(x) + \nabla(x)(x - x')$$

If x is a singular point, where f is not differentiable and $\nabla(x)$ does not exist, a **subgradient** $\bar{\nabla}(x)$ is **any** function which lower bounds f

$$f(x') \geq f(x) + \bar{\nabla}(x)(x - x')$$

SUBGRADIENT AND HINGE LOSS

In the case of hinge loss, we may observe that any line whose slope in $[-t, 0]$ (if $t = 1$, in $[0, -t]$ if $t = -1$) is a subgradient



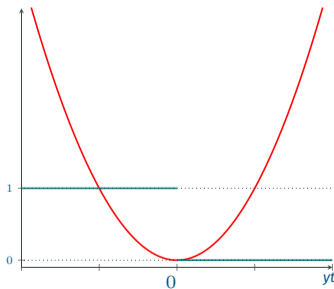
We may then choose the horizontal axis as the subgradient to use,

SQUARE LOSS

- adapted to the classification case

$$L(t, y) = (1 - yt)^2$$

- continuous, gradient continuous, convex, not surrogate
- largely wrong predictions can be too penalized
- symmetric around 0: even largely wight predictions are penalized

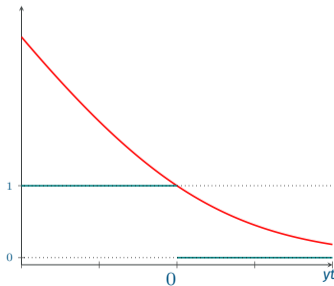


LOG LOSS (CROSS ENTROPY)

- used in logistic regression

$$L(t, y) = \frac{1}{\log 2} \log(1 + e^{-yt})$$

- a smoothed version of hinge loss
- continuous, gradient continuous, convex, surrogate
- largely wrong predictions can be too penalized
- symmetric around 0: even largely wrong predictions are penalized



WHAT IS THE RELATION WITH CROSS ENTROPY?

- given distributions p, q the cross entropy of q wrt p is defined as

$$-E_p[\log q(x)] = - \int p(x) \log q(x) dx$$

- the cross entropy is a measure of how much p and q are different
- it is related to the Kullback-Leibler divergence

$$KL(p||q) = - \int p(x) \log \frac{q(x)}{p(x)} dx = - \int p(x) \log q(x) dx + \int p(x) \log p(x) dx = -E_p[\log q(x)] - H(p)$$

where $H(p) = -E_p[\log p(x)]$ is the **entropy** of p

WHAT IS THE RELATION WITH CROSS ENTROPY?

- the entropy $H(p) = -E_p[\log p]$ denotes the expected number of bits per symbol x in a transmission channel where the distribution of symbols $p(x)$ is known
- the cross entropy $-E_p[\log q]$ denotes the total expected number of bits per symbol x in a transmission channel where the distribution of symbols $q(x)$ is used, instead of $p(x)$
- the KL divergence $KL(p||q)$ denotes the additional (with respect to the minimum) expected number of bits per symbol x in a transmission channel where the distribution of symbols $q(x)$ is used, instead of $p(x)$

WHAT IS THE RELATION WITH CROSS ENTROPY?

- consider now a classifier which predicts the probability that an element is in class C_1 and let
 - p be the probability that the element is in class C_1 : in the training set this is either 0 or 1, that is equal to the target value t
 - $y(\mathbf{x})$ be the predicted probability of the element being in class C_1
- the cross entropy $\text{CE}(\mathcal{T})$ between real and predicted probability distribution over the set of elements can be estimated as the average

$$\text{CE}(\mathcal{T}) = -\frac{1}{|\mathcal{T}|} \sum_{(\mathbf{x}, t) \in \mathcal{T}} (t \log y(\mathbf{x}) + (1-t) \log(1-y(\mathbf{x}))) = -\frac{1}{|\mathcal{T}|} \left(\sum_{(\mathbf{x}, t) \in C_1} \log y(\mathbf{x}) + \sum_{(\mathbf{x}, t) \in C_0} \log(1-y(\mathbf{x})) \right)$$

- assume now the classifier is a logistic regression, that is

$$y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + b)}}$$

then,

$$\text{CE}(\mathcal{T}) = \frac{1}{|\mathcal{T}|} \left(\sum_{(\mathbf{x}, t) \in C_1} \log(1 + e^{-(\mathbf{w}^T \mathbf{x} + b)}) + \sum_{(\mathbf{x}, t) \in C_0} \log(1 + e^{\mathbf{w}^T \mathbf{x} + b}) \right)$$

WHAT IS THE RELATION WITH CROSS ENTROPY?

- assuming now that the target encodes classes as $\bar{t} \in \{-1, 1\}$ (that is class C_0 is denoted by $\bar{t} = -1$ and class C_1 is denoted by $\bar{t} = 1$) we have

$$\text{CE}(\mathcal{T}) = \frac{1}{|\mathcal{T}|} \sum_{(x,t) \in \mathcal{T}} \log(1 + e^{-t(w^T x + b)})$$

that, apart from the constant $\log 2$ corresponds to the empirical risk if log loss is applied

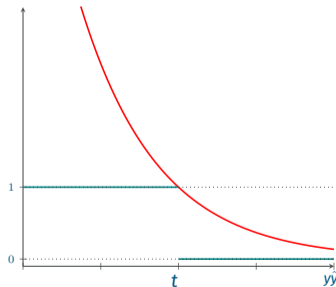
$$\bar{\mathcal{R}}_{\mathcal{T}}(h) = \frac{1}{|\mathcal{T}| \log 2} \sum_{(x,t) \in \mathcal{T}} \log(1 + e^{-t(w^T x + b)})$$

EXPONENTIAL LOSS

- used in boosting (Adaboost)

$$L(t, y) = e^{-yt}$$

- penalizes wrong predictions more than log loss: penalty grows more quickly as errors become larger
- continuous, gradient continuous, convex, surrogate



COMPUTING h^*

- In most cases, $\Theta = \mathbb{R}^d$ for some $d > 0$: in this case, the minimization of $\overline{\mathcal{R}}_{\mathcal{T}}(\mathbf{h}_{\theta})$ is unconstrained and a (at least local) minimum could be computed setting all partial derivatives to 0

$$\frac{\partial}{\partial \theta_i} \overline{\mathcal{R}}_{\mathcal{T}}(\mathbf{h}_{\theta}) = 0$$

that is, setting to zero the gradient of the empirical risk with respect to the vector of parameters θ

$$\nabla_{\theta} \overline{\mathcal{R}}_{\mathcal{T}}(\mathbf{h}_{\theta}) = \mathbf{0}$$

- The analytical solution of this set of equations is usually quite hard
- Numerical methods can be applied

GRADIENT DESCENT

- Gradient descent performs minimization of a function $J(\theta)$ through iterative updates of the current value of θ (starting from an initial value $\theta^{(0)}$) in the opposite direction to the one specified by the current value of the gradient $\nabla_{\theta}J(\theta)^{(k)} = \nabla_{\theta}J(\theta)|_{\theta=\theta^{(k)}}$

$$\theta^{(k+1)} = \theta^{(k)} - \eta \nabla_{\theta}J(\theta)^{(k)}$$

that is, for each parameter θ_i

$$\theta_i^{(k+1)} = \theta_i^{(k)} - \eta \left. \frac{\partial J(\theta)}{\partial \theta_i} \right|_{\theta^{(k)}}$$

- η is a tunable parameter, which controls the amount of update performed at each step

BATCH GRADIENT DESCENT

If minimization of the Empirical Risk is performed, gradient descent takes the form

$$\theta^{(k+1)} = \theta^{(k)} - \frac{\eta}{|\mathcal{T}|} \sum_{(\mathbf{x}, \mathbf{t}) \in \mathcal{T}} \nabla_{\theta} L(\mathbf{h}_{\theta}(\mathbf{x}), \mathbf{t})^{(k)}$$

that is,

$$\theta_i^{(k+1)} = \theta_i^{(k)} - \frac{\eta}{|\mathcal{T}|} \sum_{(\mathbf{x}, \mathbf{t}) \in \mathcal{T}} \left. \frac{\partial}{\partial \theta_i} L(\mathbf{h}_{\theta}(\mathbf{x}), \mathbf{t}) \right|_{\theta = \theta^{(k)}}$$

This is called **batch gradient descent**: observe that, at each step, all items in the training set must be considered

GRADIENT DESCENT AND LINEAR REGRESSION

For example, in the case of linear regression

$$h(\mathbf{x}) = \sum_{j=1}^d \theta_j x_j + \theta_0$$

where the loss function is usually the squared distance

$$L(h(\mathbf{x}), t) = (h(\mathbf{x}) - t)^2 = \left(\sum_{j=1}^d \theta_j x_j + \theta_0 - t \right)^2$$

the gradient is

$$\frac{\partial}{\partial \theta_i} L(h_{\theta}(\mathbf{x}), t) = 2 \left(\sum_{j=1}^d \theta_j x_j + \theta_0 - t \right) x_i \quad i = 1, \dots, d$$

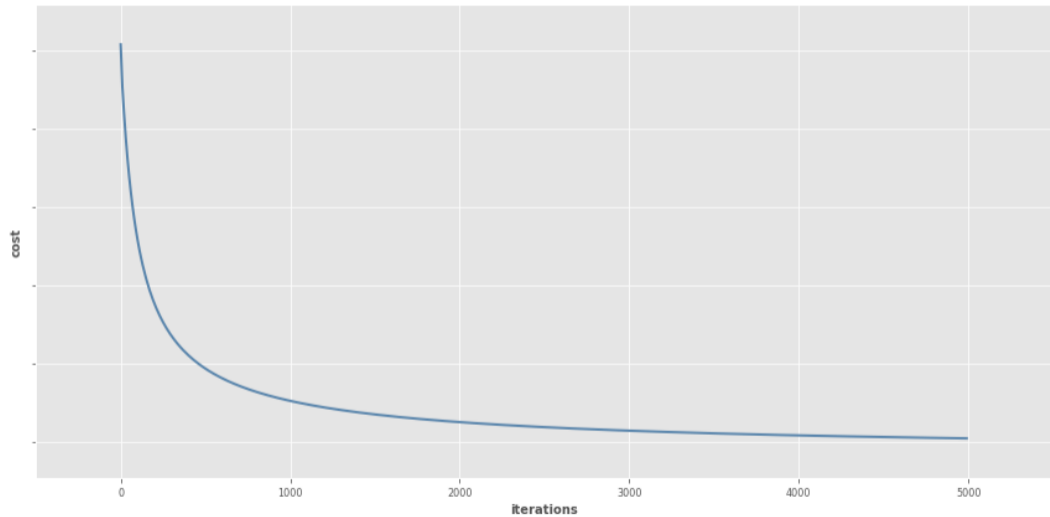
$$\frac{\partial}{\partial \theta_0} L(h_{\theta}(\mathbf{x}), t) = 2 \left(\sum_{j=1}^d \theta_j x_j + \theta_0 - t \right)$$

GRADIENT DESCENT AND LINEAR REGRESSION

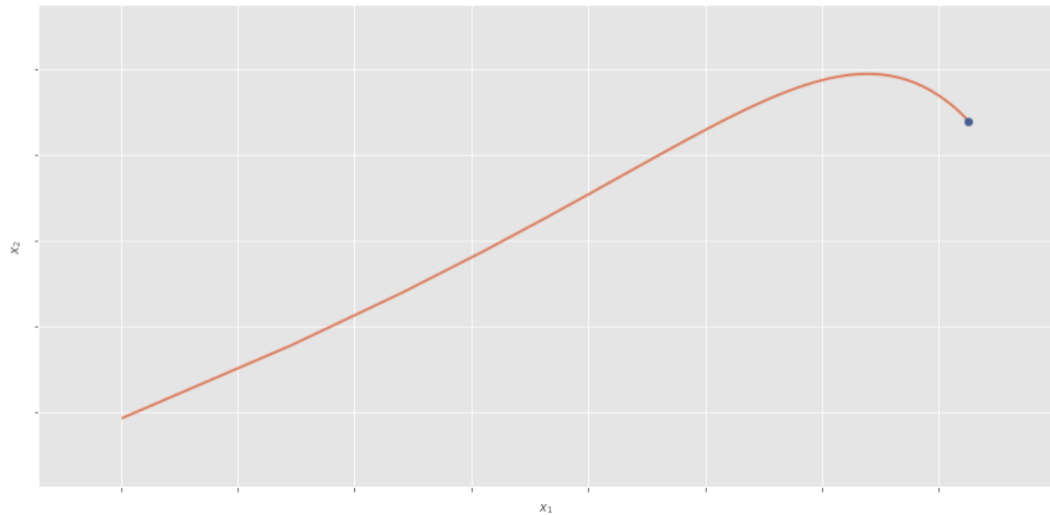
In this case, not considering the multiplicative constant 2, which may be absorbed in η , it results

$$\theta_i^{(k+1)} = \theta_i^{(k)} - \frac{\eta}{|\mathcal{T}|} \sum_{(\mathbf{x}, t) \in \mathcal{T}} \left(\sum_{j=1}^d \theta_j^{(k)} x_j + \theta_0^{(k)} - t \right) x_i \quad i = 1, \dots, d$$
$$\theta_0^{(k+1)} = \theta_0^{(k)} - \frac{\eta}{|\mathcal{T}|} \sum_{(\mathbf{x}, t) \in \mathcal{T}} \left(\sum_{j=1}^d \theta_j^{(k)} x_j + \theta_0^{(k)} - t \right)$$

GRADIENT DESCENT



GRADIENT DESCENT



GRADIENT DESCENT

As we need to calculate the gradients for the whole dataset to perform just *one* update, batch gradient descent can be very slow and is intractable for datasets that do not fit in memory. Batch gradient descent also does not allow us to update our model *online*, i.e. with new examples on-the-fly.

STOCHASTIC GRADIENT DESCENT

Batch gradient descent can be modified by performing the update, at each step, on the basis of the evaluation at a single item $\mathbf{x}_j, \mathbf{t}_j$ of the training set.

$$\theta^{(k+1)} = \theta^{(k)} - \eta \nabla_{\theta} L(h_{\theta}(\mathbf{x}_j), \mathbf{t}_j)^{(k)}$$

or

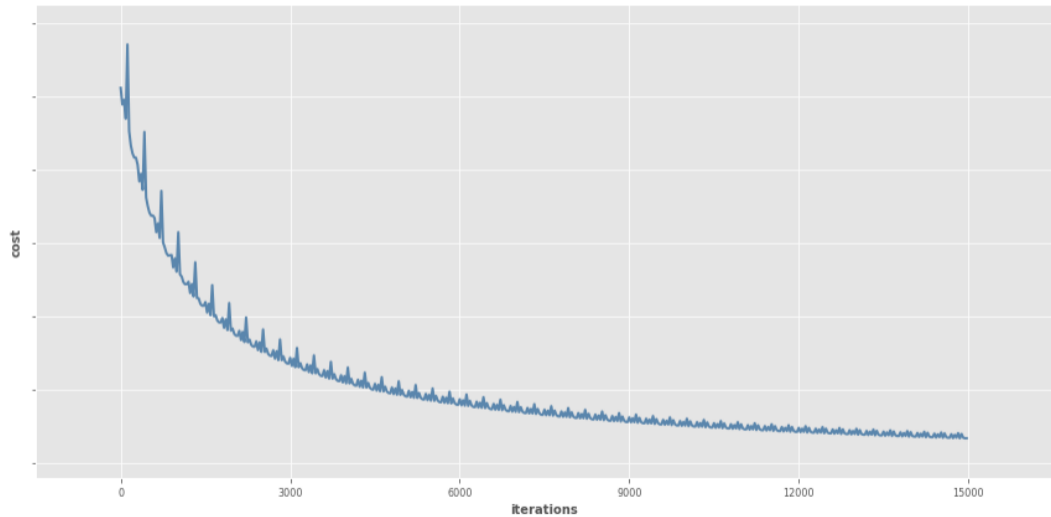
$$\theta_i^{(k+1)} = \theta_i^{(k)} - \eta \left. \frac{\partial}{\partial \theta_i} L(h_{\theta}(\mathbf{x}_j), \mathbf{t}_j) \right|_{\theta = \theta^{(k)}}$$

SGD AND LINEAR REGRESSION

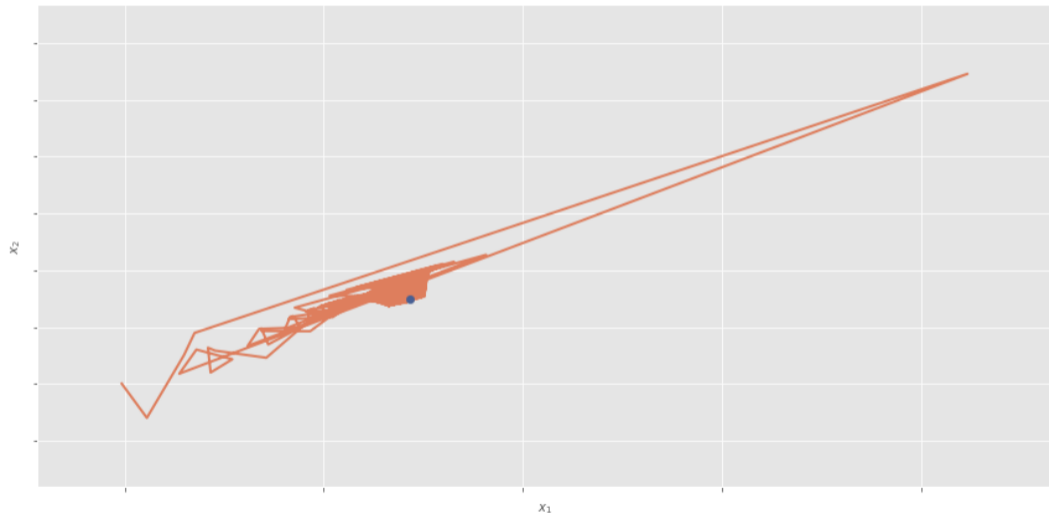
In the case of linear regression this results into

$$\theta_i^{(k+1)} = \theta_i^{(k)} - \eta \left(\sum_{r=1}^d \theta_r^{(k)} x_{jr} + \theta_0^{(k)} - t \right) x_{ji} \quad i = 1, \dots, d$$
$$\theta_0^{(k+1)} = \theta_0^{(k)} - \eta \left(\sum_{r=1}^d \theta_r^{(k)} x_{jr} + \theta_0^{(k)} - t \right)$$

STOCHASTIC GRADIENT DESCENT



STOCHASTIC GRADIENT DESCENT



MINI-BATCH GRADIENT DESCENT

An intermediate case is the one when a subset B_r of size m of the items in the training is considered at each step for gradient evaluation

$$\theta^{(k+1)} = \theta^{(k)} - \frac{\eta}{m} \sum_{(\mathbf{x}, \mathbf{t}) \in B_r} \nabla_{\theta} L(h_{\theta}(\mathbf{x}), \mathbf{t})^{(k)}$$

that is,

$$\theta_i^{(k+1)} = \theta_i^{(k)} - \frac{\eta}{m} \sum_{(\mathbf{x}, \mathbf{t}) \in B_r} \left. \frac{\partial}{\partial \theta_i} L(h_{\theta}(\mathbf{x}), \mathbf{t}) \right|_{\theta^{(k)}}$$

This is called **mini-batch gradient descent**.

MINI-BATCH GRADIENT DESCENT

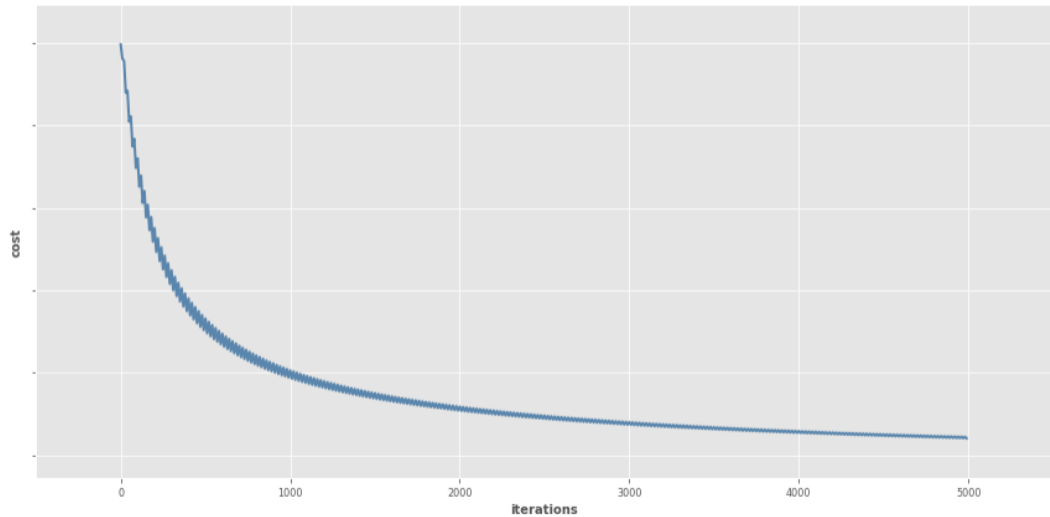
This approach

- reduces the variance of the parameter updates, which can lead to more stable convergence wrt SGD
- limits the amount of items considered for gradient evaluation before a parameter update is performed.

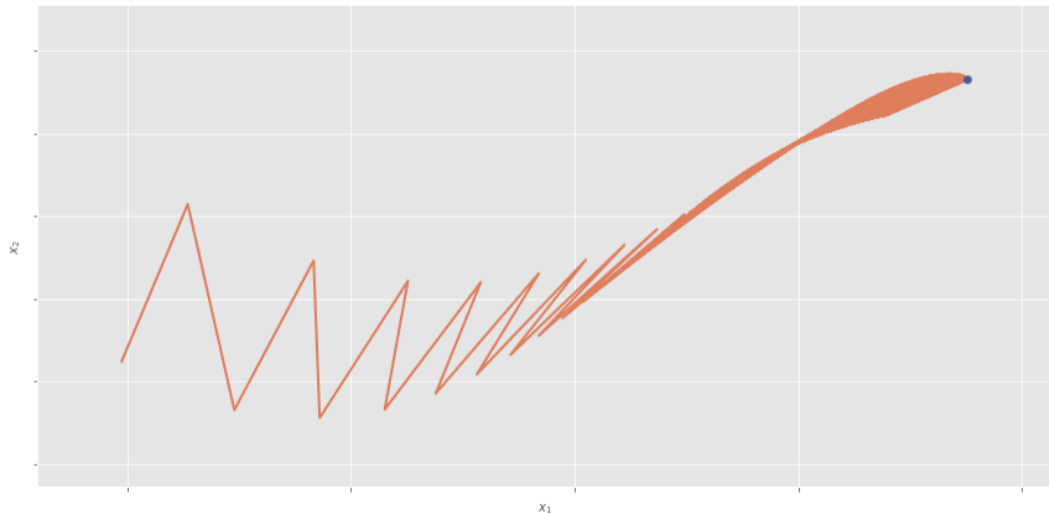
Mini-batch gradient descent is typically the algorithm of choice when training neural networks

Observe that the size $|B_r|$ of mini-batches is itself a tunable parameter

MINI-BATCH GRADIENT DESCENT



MINI-BATCH GRADIENT DESCENT



OPEN ISSUES

- Choosing a proper value for η can be difficult.
- Apply mechanisms to adjust the learning rate during training by reducing it either according to a pre-defined schedule or when the loss function decrease between epochs falls below a threshold. Both schedules and thresholds should be defined in advance.
- The same learning rate applies to updating all parameter.
- Saddle points appears in complex losses, which are usually surrounded by a plateau. Hard for simple gradient descent methods to escape, as the gradient is almost zero in all dimensions.

MOMENTUM GRADIENT DESCENT

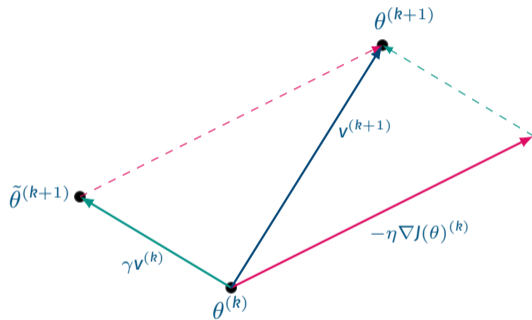
- Based on a physical interpretation of the optimization process: a body of mass $m = 1$ is moving on the surface of a cost function $J(\theta)$, with potential energy $U(\theta) = \eta J(\theta)$ and weight force (or acceleration, since $m = 1$) $F(\theta) = -\nabla U(\theta) = -\eta J'(\theta)$, at any point θ
- In gradient descent, the movement of the body is determined by the acceleration at that point, that is by the gradient $J'(\theta)$
- In momentum gradient descent, the velocity $v(\theta)$ of the body is considered: the movement of the body is determined by the velocity, that is,

$$\theta^{(k+1)} = \theta^{(k)} + v^{(k+1)}$$

with the velocity changing as determined by the acceleration

$$v^{(k+1)} = v^{(k)} - \eta J'(\theta) |_{\theta=\theta^{(k)}}$$

MOMENTUM GRADIENT DESCENT



MOMENTUM GRADIENT DESCENT

This results into

$$\begin{aligned} \mathbf{v}^{(k+1)} &= \mathbf{v}^{(k)} - \eta \sum_{(\mathbf{x}, \mathbf{t}) \in B_r} \nabla_{\theta} L(\mathbf{h}_{\theta}(\mathbf{x}), \mathbf{t})^{(k)} = \mathbf{v}^{(k-1)} - \eta \sum_{(\mathbf{x}, \mathbf{t}) \in B_r} \nabla_{\theta} L(\mathbf{h}_{\theta}(\mathbf{x}), \mathbf{t})^{(k-1)} - \eta \sum_{(\mathbf{x}, \mathbf{t}) \in B_r} \nabla_{\theta} L(\mathbf{h}_{\theta}(\mathbf{x}), \mathbf{t})^{(k)} = \dots \\ &= \mathbf{v}^{(0)} - \eta \sum_{i=0}^k \sum_{(\mathbf{x}, \mathbf{t}) \in B_r} \nabla_{\theta} L(\mathbf{h}_{\theta}(\mathbf{x}), \mathbf{t})^{(i)} \\ \theta^{(k+1)} &= \theta^{(k)} + \mathbf{v}^{(k+1)} = \theta^{(k)} \mathbf{v}^{(0)} - \eta \sum_{i=0}^k \sum_{(\mathbf{x}, \mathbf{t}) \in B_r} \nabla_{\theta} L(\mathbf{h}_{\theta}(\mathbf{x}), \mathbf{t})^{(i)} \end{aligned}$$

MOMENTUM GRADIENT DESCENT

In momentum gradient descent it is usually introduced a second parameter γ , which affects the fraction of $\mathbf{v}^{(k)}$ that is considered for the computation of $\mathbf{v}^{(k+1)}$. In terms of physical model, this corresponds to introducing an attrition coefficient. Applying the approach to the case of mini-batches, we get:

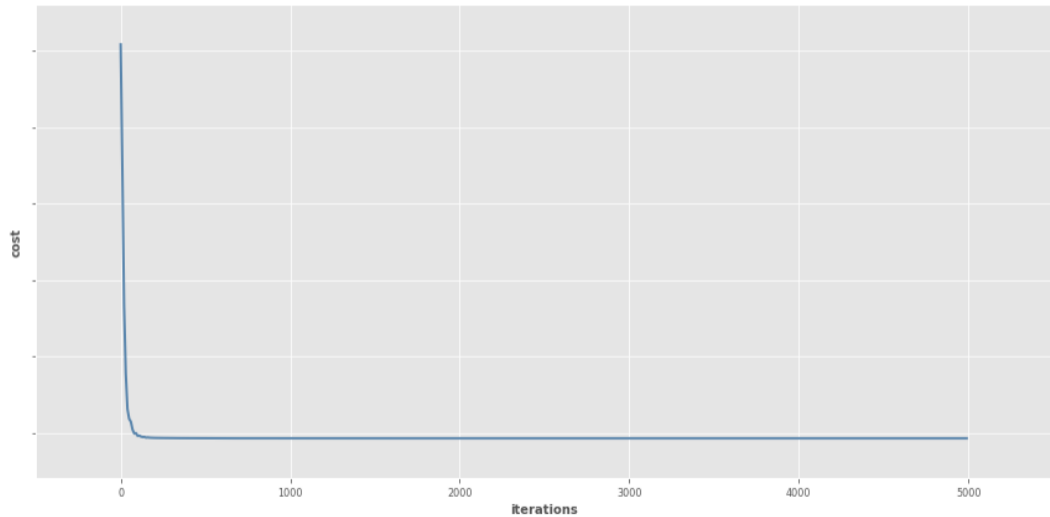
$$\mathbf{v}^{(k+1)} = \gamma \mathbf{v}^{(k)} - \eta \sum_{(\mathbf{x}, \mathbf{t}) \in B_r} \nabla_{\theta} L(h_{\theta}(\mathbf{x}), \mathbf{t})^{(k)}$$
$$\theta^{(k+1)} = \theta^{(k)} + \mathbf{v}^{(k+1)}$$

MGD AND LINEAR REGRESSION

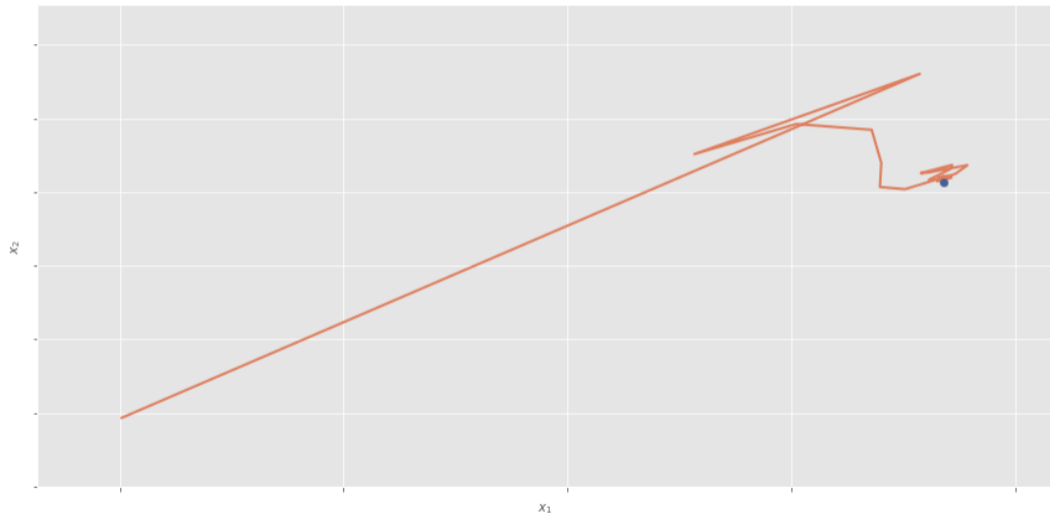
In the case of linear regression, this results into:

$$\mathbf{v}_i^{(k+1)} = \begin{cases} \gamma \mathbf{v}_i^{(k)} - \eta \sum_{(\mathbf{x}, t) \in B_r} \left(\sum_{j=1}^d \theta_j^{(k)} x_j + \theta_0^{(k)} - t \right) x_i & i = 1, \dots, d \\ \gamma \mathbf{v}_i^{(k)} - \eta \sum_{(\mathbf{x}, t) \in B_r} \left(\sum_{j=1}^d \theta_j^{(k)} x_j + \theta_0^{(k)} - t \right) & i = 0 \end{cases}$$
$$\theta_i^{(k+1)} = \theta_i^{(k)} + \mathbf{v}_i^{(k+1)}$$

MOMENTUM GRADIENT DESCENT



MOMENTUM GRADIENT DESCENT



NESTEROV GRADIENT DESCENT

In MGD, adding $\gamma \mathbf{v}^{(k)}$ to $\theta^{(k)}$ provides an approximation

$$\tilde{\theta}^{(k+1)} = \theta^{(k)} + \gamma \mathbf{v}^{(k)}$$

of the real value $\theta^{(k+1)}$

$$\tilde{\theta}^{(k+1)} = \theta^{(k)} + \gamma \mathbf{v}^{(k)}$$

$$\mathbf{v}^{(k+1)} = \gamma \mathbf{v}^{(k)} - \eta \nabla_{\theta} J(\theta)^{(k)}$$

$$\theta^{(k+1)} = \theta^{(k)} + \mathbf{v}^{(k+1)}$$

NESTEROV GRADIENT DESCENT

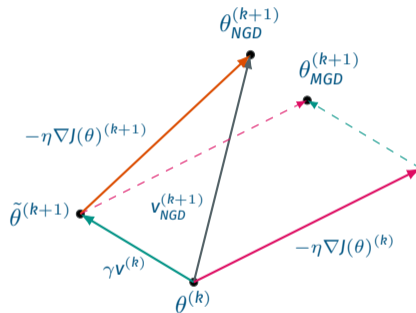
- The same approach of momentum gradient descent is applied, with the gradient estimation performed not at the current point $\theta^{(k)}$, but approximately at the next point $\theta^{(k+1)}$
- The approximation derives by considering $\tilde{\theta}^{(k)} = \theta^{(k)} + \gamma \mathbf{v}^{(k)}$ instead of $\theta^{(k+1)}$
- The updates of \mathbf{v} and θ are considered in advance with respect to momentum GD

$$\tilde{\theta}^{(k+1)} = \theta^{(k)} + \gamma \mathbf{v}^{(k)}$$

$$\mathbf{v}^{(k+1)} = \gamma \mathbf{v}^{(k)} - \eta \nabla_{\theta} J(\theta) |_{\theta = \tilde{\theta}^{(k+1)}}$$

$$\theta^{(k+1)} = \theta^{(k)} + \mathbf{v}^{(k+1)}$$

NESTEROV GRADIENT DESCENT



DYNAMICALLY UPDATING THE LEARNING RATE

Learning rate is a crucial parameter for SGD

- Too large: overshoots local minimum, loss increases
- Too small: makes very slow progress, can get stuck
- Good learning rate: makes steady progress toward local minimum

In practice: gradually decrease of the learning rate

- Step decay: periodically (every few epochs) decay η by a factor 2
- Exponential decay: $\eta^{(k)} = \eta^{(0)} e^{-\alpha k}$
- 1/t decay: $\eta^{(k)} = \frac{\eta^{(0)}}{1 + \alpha k}$

Extension: update η by monitoring the learning process

ADAGRAD

In gradient descent the update of parameter θ_j is the following

$$\theta_j^{(k+1)} = \theta_j^{(k)} - \eta \frac{\partial J(\theta)}{\partial \theta_j} \Big|_{\theta^{(k)}}$$

where the learning rate η is equal for all parameters.

We now rewrite this update in terms of the parameter update $\Delta\theta_{j,k}$, as a sequence of three steps:

$$\begin{aligned} \mathbf{g}_{j,k} &= \frac{\partial J(\theta)}{\partial \theta_j} \Big|_{\theta^{(k)}} \\ \Delta_{j,k} &= -\eta \mathbf{g}_{j,k} \\ \theta_j^{(k+1)} &= \theta_j^{(k)} + \Delta_{j,k} \end{aligned}$$

ADAGRAD

Adagrad modifies this behavior for what regards the computation of $\Delta_{j,k}$ by adapting the learning rate to the parameters, performing larger updates for infrequent and smaller updates for frequent parameters.

In Adagrad, each parameter update refers to a learning rate $\eta_j^{(k)}$, that is

$$\Delta_{j,k} = -\eta_j^{(k)} \mathbf{g}_{j,k}$$

where $\eta_j^{(k)}$ is dependent on the parameter itself and a common predefined learning rate η

ADAGRAD

In particular,

$$\eta_j^{(k)} = \frac{\eta}{\sqrt{G_{j,k} + \varepsilon}}$$

and

$$G_{j,k} = \sum_{i=0}^k g_{j,i}^2$$

is the sum of the squared derivatives of the loss function wrt to θ_j computed for all previous iterations. ε is a small smoothing constant, introduced to avoid null denominators. This results into

$$\Delta_{j,k} = -\frac{\eta}{\sqrt{G_{j,k} + \varepsilon}} g_{j,k}$$

ADAGRAD

- Learning rates decrease at each step, with the ones associated to parameters which had large gradients in the past decreasing more
- Adagrad's main weakness is its accumulation of the squared gradients in the denominator: since every added term is positive, the accumulated sum keeps growing during training. This in turn, as observed above, causes the learning rate to shrink and eventually become infinitesimally small, at which point the algorithm is no longer able to acquire additional knowledge.

RMSPROP

RMSprop seeks to reduce the aggressive, monotonically decreasing learning rate of Adagrad.

The sum over past squared gradients $G_{j,k}$ is replaced with its decaying version $\tilde{G}_{j,k}$.

This is obtained through a **decay**, obtained by applying a coefficient $0 < \gamma < 1$

$$\begin{aligned}\tilde{G}_{j,k} &= \gamma \tilde{G}_{j,k-1} + (1 - \gamma) g_{j,k}^2 \\ &= \gamma(\gamma \tilde{G}_{j,k-2} + (1 - \gamma) g_{j,k-1}^2) + (1 - \gamma) g_{j,k}^2 = \gamma^2 \tilde{G}_{j,k-2} + (1 - \gamma)(\gamma g_{j,k-1}^2 + g_{j,k}^2) \\ &= \dots \\ &= (1 - \gamma) \sum_{i=0}^k \gamma^{k-i} g_{j,i}^2\end{aligned}$$

since we assume $\tilde{G}_{j,k} = 0$ if $k < 0$.

RMSPROP

This results into the following step, at the $k + 1$ -th iteration

$$\begin{aligned}g_{j,k} &= \left. \frac{\partial J(\theta)}{\partial \theta_j} \right|_{\theta = \theta^{(k)}} \\ \tilde{G}_{j,k} &= \gamma \tilde{G}_{j,k-1} + (1 - \gamma) g_{j,k}^2 \\ \Delta_{j,k} &= - \frac{\eta}{\sqrt{\tilde{G}_{j,k} + \epsilon}} g_{j,k} \\ \theta_j^{(k+1)} &= \theta_j^{(k)} + \Delta_{j,k}\end{aligned}$$

ADADELTA

Adadelta is an extension of RMSprop in which no value η has to be arbitrarily defined: it is instead substituted by the decayed sum of previous squared updates, with the same decay γ applied for derivatives.

$$\bar{G}_{j,k} = \gamma \bar{G}_{j,k-1} + (1 - \gamma) \Delta_{j,k}^2 = (1 - \gamma) \sum_{i=0}^k \gamma^{k-i} \Delta_{j,i}^2$$

ADADELTA

The update rule is then defined as

$$\begin{aligned}g_{j,k} &= \left. \frac{\partial J(\theta)}{\partial \theta_j} \right|_{\theta^{(k)}} \\ \tilde{G}_{j,k} &= \gamma \tilde{G}_{j,k-1} + (1 - \gamma) g_{j,k}^2 \\ \Delta_{j,k} &= - \frac{\sqrt{\bar{G}_{j,k-1} + \epsilon}}{\sqrt{\tilde{G}_{j,k} + \epsilon}} g_{j,k} \\ \bar{G}_{j,k} &= \gamma \bar{G}_{j,k-1} + (1 - \gamma) \Delta_{j,k}^2 \\ \theta_j^{(k+1)} &= \theta_j^{(k)} + \Delta_{j,k}\end{aligned}$$

ADAM

In addition to storing the exponentially decaying sum $\tilde{G}_{j,k}$ of past squared derivatives $g_{j,k}^2$ like Adadelta and RMSprop (to be used in the same way as in such methods), Adam also keeps an exponentially decaying sum $\tilde{H}_{j,k}$ of past (non squared) derivatives $g_{j,k}$, as a substitute to the derivative $g_{j,k}$ in the iteration step.

$$\begin{aligned}\tilde{G}_{j,k} &= \gamma\tilde{G}_{j,k-1} + (1 - \gamma)g_{j,k}^2 \\ \tilde{H}_{j,k} &= \beta\tilde{H}_{j,k-1} + (1 - \beta)g_{j,k}\end{aligned}$$

ADAM

Since it is assumed that $\tilde{H}_{j,k} = \tilde{G}_{j,k} = 0$ if $k < 0$ and γ, β values are usually both close to 1, the methods presents a tendency (bias) to return small values of $\tilde{H}_{j,k}$ and $\tilde{G}_{j,k}$, especially during the initial time steps.

This issue is managed by applying a bias correction:

$$\hat{G}_{j,k} = \frac{\tilde{G}_{j,k}}{1 - \gamma^k}$$
$$\hat{H}_{j,k} = \frac{\tilde{H}_{j,k}}{1 - \beta^k}$$

ADAM

Parameters are updated just as we have seen in Adadelta and RMSprop, which yields the Adam update rule:

$$\begin{aligned}g_{j,k} &= \left. \frac{\partial J(\theta)}{\partial \theta_j} \right|_{\theta = \theta^{(k)}} \\ \tilde{\mathbf{G}}_{j,k} &= \gamma \tilde{\mathbf{G}}_{j,k-1} + (1 - \gamma) \mathbf{g}_{j,k}^2 \\ \tilde{\mathbf{H}}_{j,k} &= \beta \tilde{\mathbf{H}}_{j,k-1} + (1 - \beta) \mathbf{g}_{j,k} \\ \hat{\mathbf{G}}_{j,k} &= \frac{\tilde{\mathbf{G}}_{j,k}}{1 - \gamma^k} \\ \hat{\mathbf{H}}_{j,k} &= \frac{\tilde{\mathbf{H}}_{j,k}}{1 - \beta^k} \\ \Delta_{j,k} &= - \frac{\eta}{\sqrt{\hat{\mathbf{G}}_{j,k} + \epsilon}} \hat{\mathbf{H}}_{j,k} \\ \theta_j^{(k+1)} &= \theta_j^{(k)} + \Delta_{j,k}\end{aligned}$$

SECOND ORDER METHODS

Maxima (or minima) of $J(\theta)$ can be found by searching points where the gradient (all partial derivatives) is zero.

Any iterative method to compute zeros of a function (such as Newton-Raphson) can then be applied on the gradient $\nabla_{\theta}J(\theta)$

The basic idea of Newton's method is to use both the first-order derivative (gradient) and second-order derivative (Hessian matrix) to approximate the objective function with a quadratic function, and then solve the minimum optimization of the quadratic function. This process is repeated until the updated variable converges.

SECOND ORDER METHODS

$$\theta_j^{(k+1)} = \theta_j^{(k)} - \frac{J'(\theta)}{J''(\theta)} \Big|_{\theta=\theta^k}$$

More general, the high-dimensional Newton's iteration formula is

$$\theta^{(k+1)} = \theta^{(k)} - H(J(\theta))^{-1} \nabla_{\theta} J(\theta) \Big|_{\theta=\theta^k}$$

where $H(J(\theta))$ is a Hessian matrix of $J(\theta)$.

SECOND ORDER METHODS

Newton's method is an iterative algorithm that requires the computation of the inverse Hessian matrix of the objective function at each step, which makes the storage and computation very expensive.

To overcome the expensive storage and computation, approximate algorithms were considered such as **quasi-Newton methods**. The essential idea of all quasi-Newton methods is to use a positive definite matrix to approximate the inverse of the Hessian matrix, thus simplifying the complexity of the operation.