

Stochastic Gradient Descent (SGD)

SGD is:

- an iterative optimization method, stochastic approximation of gradient descent
- the fundamental idea behind many machine learning algorithms such as linear support vector machines, logistic regression, graphical models
- the de facto standard training algorithm for neural networks

Given an objective function $Q_i(\mathbf{w})$ which computes some sort of error for sample i , e.g.

$$Q_i(\mathbf{w}) = (\hat{y}_i - y_i)^2$$

Classical SGD:

Initialize parameter vector \mathbf{w} and set learning rate η

```

while termination criterion not satisfied do
  randomly shuffle samples in training set
  for  $i \leftarrow 1, n$  do
     $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla Q_i(\mathbf{w})$ 
  end for
end while
  
```

SGD Extensions

Extension 1: Batch updates

Issue: Each iteration of classical SGD evaluates the gradient at a single point x_i only

- ⇒ often results in noisy estimates of the true gradient
- ⇒ computationally expensive (blocking vectorization)

Idea: Combine multiple gradient estimates of a so-called “mini-batch”, then do a single update

- ⇒ smoother convergence thanks to better estimates
- ⇒ speed-up owing to vectorization

Extension 2: Implicit updates (ISGD)

Issues:

- Classical SGD sensitive to choice of η and may easily diverge
- Better convergence with some decay schedule of η over iterations, $\eta = \eta(\text{iter})$

Idea: Evaluate stochastic gradient at the next iteration rather than the current one:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla Q_i(\mathbf{w}_{k+1})$$

To see the effect, assuming an ordinary least squares objective

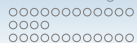
$$Q_i(\mathbf{w}) = \frac{1}{2}(\hat{y}_i - y_i)^2$$

in the usual samples $\{(x_i, y_i)\}_{i=1}^N$, compare the resulting update rules:

$$\text{Classical SGD: } \mathbf{w}_{k+1} = \mathbf{w}_k + \eta(y_i - \mathbf{x}_i^\top \mathbf{w}_k) \mathbf{x}_i$$

$$\text{Implicit SGD (closed-form!): } \mathbf{w}_{k+1} = \mathbf{w}_k + \frac{\eta}{1 + \eta \|\mathbf{x}_i\|^2} (y_i - \mathbf{x}_i^\top \mathbf{w}_k) \mathbf{x}_i$$

Note the normalizing effect on η , resulting in numerically stable procedure for virtually all η



Extension 3: Momentum

Issue: Classical SGD frequently shows oscillations, hampering convergence

Idea: Introduce some “memory” $\Delta \mathbf{w}_k$ of the current direction, i.e., do not take “hard turns” at every step:

$$\Delta \mathbf{w}_{k+1} = -\eta \nabla Q_i(\mathbf{w}_k) + \alpha \Delta \mathbf{w}_k$$

resulting in the update rule

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \Delta \mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla Q_i(\mathbf{w}_k) + \alpha \Delta \mathbf{w}_k$$

with $\alpha \in [0, 1)$ an exponential decay factor (hyperparameter)