

Gradient Descent

Johanni Brea

Introduction à l'apprentissage automatique

GYMINF 2021

Table of Contents

1. Gradient Descent

2. Stochastic Gradient Descent

3. Adaptive Learning Rates and Momentum

4. Early Stopping

Gradient Descent

1. Input: loss function L , initial guess $\beta^{(0)} = (\beta_0^{(0)}, \dots, \beta_p^{(0)})$
learning rate η , maximal number of steps T .
2. For $t = 1, \dots, T$
 - ▶ $\delta_i = \eta \frac{\partial L}{\partial \beta_i} (\beta^{(t-1)})$
 - ▶ $\beta_i^{(t)} = \beta_i^{(t-1)} - \delta_i$
3. Return $\beta^{(T)}$

Gradient Descent

1. Input: loss function L , initial guess $\beta^{(0)} = (\beta_0^{(0)}, \dots, \beta_p^{(0)})$
learning rate η , maximal number of steps T .
2. For $t = 1, \dots, T$
 - ▶ $\delta_i = \eta \frac{\partial L}{\partial \beta_i} (\beta^{(t-1)})$
 - ▶ $\beta_i^{(t)} = \beta_i^{(t-1)} - \delta_i$
3. Return $\beta^{(T)}$

Automatic Differentiation software uses the chain rule and symbolic derivatives for primitive functions, to compute the derivative of almost any code we write.

Table of Contents

1. Gradient Descent

2. Stochastic Gradient Descent

3. Adaptive Learning Rates and Momentum

4. Early Stopping

Stochastic Gradient Descent (SGD)

Computing the loss over all samples $1, \dots, n$ can be computationally costly.
A subset of the training data may be sufficient to estimate the gradient direction.

Stochastic Gradient Descent (SGD)

Computing the loss over all samples $1, \dots, n$ can be computationally costly.
A subset of the training data may be sufficient to estimate the gradient direction.

1. Input: loss function L , initial guess

$$\beta^{(0)} = (\beta_0^{(0)}, \dots, \beta_p^{(0)})$$

learning rate η , maximal number of steps T ,
tolerance Δ , **batch size** B .

where $L(\beta; \mathcal{I})$ is the loss function
evaluated on the training samples
with indices in \mathcal{I} , e.g.

2. For $t = 1, \dots, T$

- ▶ **Determine batch of training indices \mathcal{I}**

- ▶ $\delta_i = \eta \frac{\partial L}{\partial \beta_i} (\beta^{(t-1)}; \mathcal{I})$

- ▶ $\beta_i^{(t)} = \beta_i^{(t-1)} - \delta_i$

3. Return $\beta^{(T)}$

$$L(\beta; \mathcal{I}) = \frac{1}{B} \sum_{i \in \mathcal{I}} (y_i - x_i^T \beta)^2$$

Stochastic Gradient Descent (SGD)

Computing the loss over all samples $1, \dots, n$ can be computationally costly.
A subset of the training data may be sufficient to estimate the gradient direction.

1. Input: loss function L , initial guess

$$\beta^{(0)} = (\beta_0^{(0)}, \dots, \beta_p^{(0)})$$

learning rate η , maximal number of steps T ,
tolerance Δ , **batch size** B .

2. For $t = 1, \dots, T$

- ▶ **Determine batch of training indices \mathcal{I}**

- ▶ $\delta_i = \eta \frac{\partial L}{\partial \beta_i} (\beta^{(t-1)}; \mathcal{I})$

- ▶ $\beta_i^{(t)} = \beta_i^{(t-1)} - \delta_i$

3. Return $\beta^{(T)}$

where $L(\beta; \mathcal{I})$ is the loss function evaluated on the training samples with indices in \mathcal{I} , e.g.

$$L(\beta; \mathcal{I}) = \frac{1}{B} \sum_{i \in \mathcal{I}} (y_i - x_i^T \beta)^2$$

Example $B = 6$

	batch 1	batch 2	batch 3	...
\mathcal{I}	1 8 3 13 93	9 14 2 26 31	...	

Table of Contents

1. Gradient Descent

2. Stochastic Gradient Descent

3. Adaptive Learning Rates and Momentum

4. Early Stopping

Adaptive Learning Rates and Momentum

Momentum

- ▶ $\delta_i = \eta \frac{\partial L}{\partial \beta_i} \left(\beta^{(t-1)} \right)$
- ▶ $v_i^{(t)} = \mu v_i^{(t-1)} + (1 - \mu) \delta_i$
- ▶ $\beta_i^{(t)} = \beta_i^{(t-1)} - v_i^{(t)}$

Adaptive Learning Rates and Momentum

Momentum

- ▶ $\delta_i = \eta \frac{\partial L}{\partial \beta_i} \left(\beta^{(t-1)} \right)$
- ▶ $v_i^{(t)} = \mu v_i^{(t-1)} + (1 - \mu) \delta_i$
- ▶ $\beta_i^{(t)} = \beta_i^{(t-1)} - v_i^{(t)}$

Adaptive Learning Rates

For every parameter a different learning rate η_i can be chosen. It can also change over time.

<https://doi.org/10.1371/journal.pcbi.1007640>

Adaptive Learning Rates and Momentum

Momentum

- ▶ $\delta_i = \eta \frac{\partial L}{\partial \beta_i} \left(\beta^{(t-1)} \right)$
- ▶ $v_i^{(t)} = \mu v_i^{(t-1)} + (1 - \mu) \delta_i$
- ▶ $\beta_i^{(t)} = \beta_i^{(t-1)} - v_i^{(t)}$

Adaptive Learning Rates

For every parameter a different learning rate η_i can be chosen. It can also change over time.

<https://doi.org/10.1371/journal.pcbi.1007640>

Modern methods like ADAM(W) include momentum and automatically adapting learning rates for the different parameters.

Table of Contents

1. Gradient Descent

2. Stochastic Gradient Descent

3. Adaptive Learning Rates and Momentum

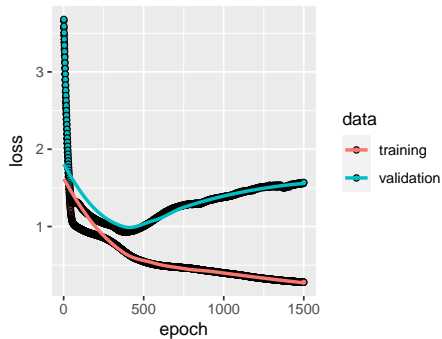
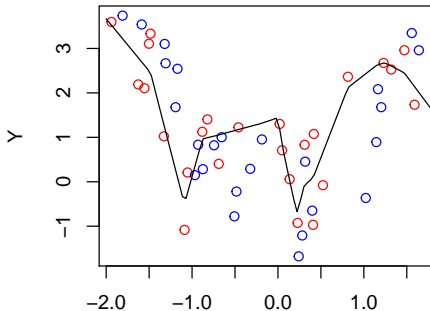
4. Early Stopping

Early Stopping

Start with small weights and stop gradient descent when validation loss starts to increase.

training data validation data

black line: a flexible neural network trained
with gradient descent



Quiz

- ▶ If we choose a small learning rate larger than zero, the training loss in gradient descent is decreasing in every step.

Quiz

- ▶ If we choose a small learning rate larger than zero, the training loss in gradient descent is decreasing in every step.
- ▶ If we choose a small learning rate larger than zero, the training loss in stochastic gradient descent is decreasing in every step.

Quiz

- ▶ If we choose a small learning rate larger than zero, the training loss in gradient descent is decreasing in every step.
- ▶ If we choose a small learning rate larger than zero, the training loss in stochastic gradient descent is decreasing in every step.
- ▶ Stochastic gradient descent requires less computation than full gradient descent for each update of the parameters.