# Machine learning
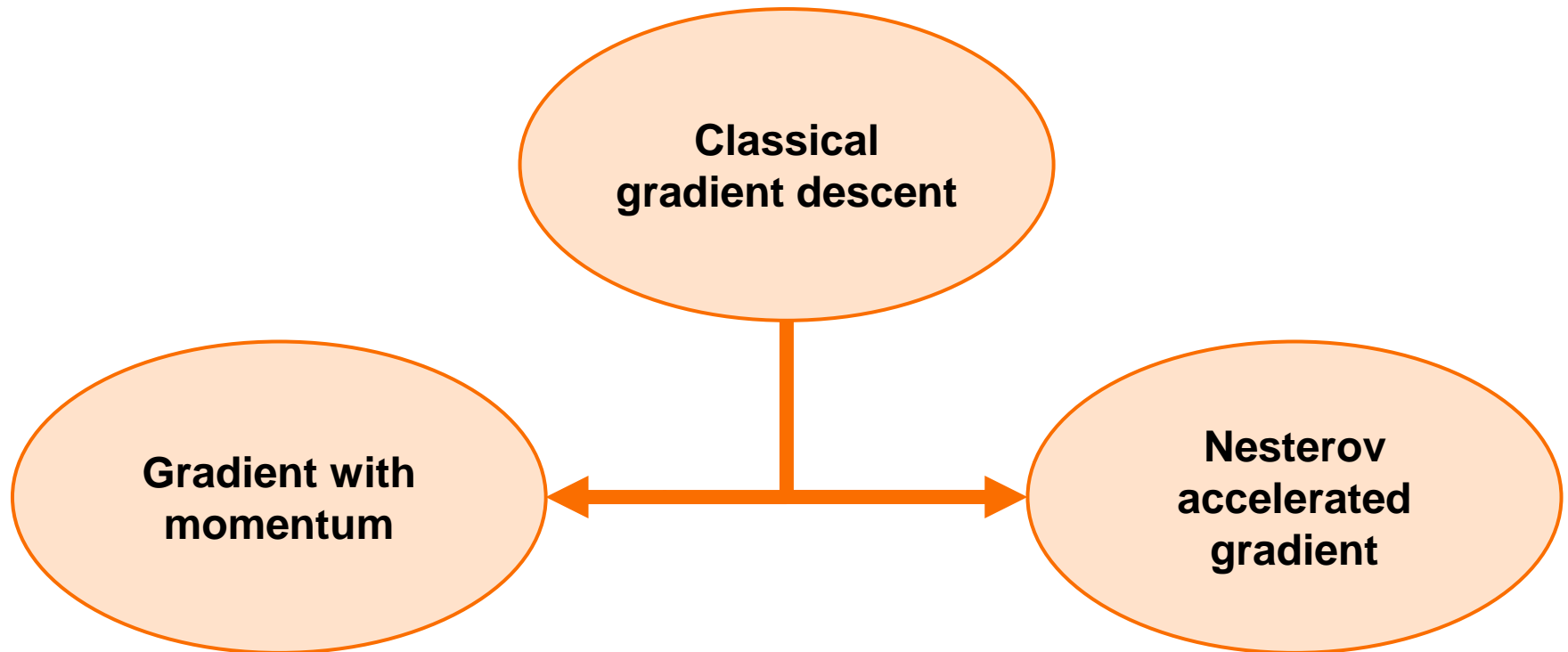
## Mini-project: Gradient Descent with Momentum and Nesterov-Accelerated Gradient

Noah Gollnick

Wirtschaftsinformatik

**Main goal: Accelerate the machine learning process**

# Learning-rate-Dilemma in classic Gradient descent

**High learning rate** ⟷ **Low learning rate**

- Faster training process, smaller number of iterations for best result

- Higher probability of overshooting and oscillating around the minimum

- Slower training process, higher number of iterations for best result

- Smaller probability of overshooting, higher chance to find the perfect minimum

Technische
Universität
Braunschweig

igp

# Approach of using momentum

**Possible Approach: Dynamic learning rate**

**Assumption 1:**
**The risk of overshooting is higher at the end of the training process**

**Assumption 2:**
**The loss is getting generally better over time, so the steps are getting smaller**

**Learning algorithm adjusts to the size of the last gradient descent step**

Technische
Universität
Braunschweig

igp

# Concept and Computation of Momentum

Momentum: Adding a velocity vector to the changes made in one machine learning step. The velocity equals the difference between the parameters of the last 2 steps:

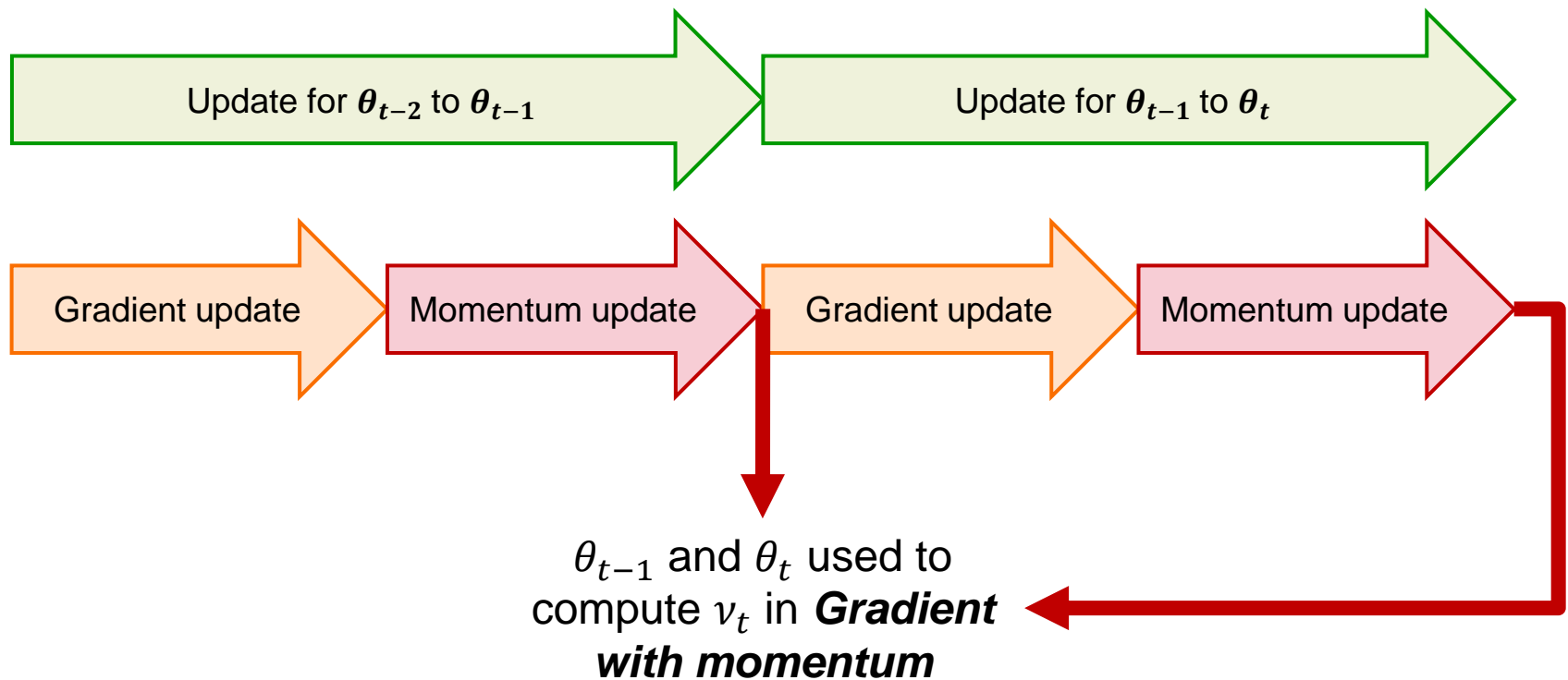$$\boldsymbol{v_t} = \boldsymbol{\theta_t} - \boldsymbol{\theta_{t-1}}$$

The velocity vector is added to the next step (multiplicated with a scalar factor) while subtracting the gradient:

$$\boldsymbol{\theta_{t+1}} = \boldsymbol{\theta_t} - \boldsymbol{\varepsilon \Delta f(\theta_t)} + \boldsymbol{\mu v_t}$$

Technische
Universität
Braunschweig

igp

# Difference between GDM and NAG

Difference between classical Gradient with momentum and Nesterov accelerated gradient:
NAG splits up gradient update step and momentum update step

Update for $\boldsymbol{\theta}_{t-2}$ to $\boldsymbol{\theta}_{t-1}$

Update for $\boldsymbol{\theta}_{t-1}$ to $\boldsymbol{\theta}_t$

Gradient update

Momentum update

Gradient update

Momentum update

$\theta_{t-1}$ and $\theta_t$ used to compute $v_t$ in *Gradient with momentum*

Technische
Universität
Braunschweig

igp

# Difference between GDM and NAG

Difference between classical Gradient with momentum and Nesterov accelerated gradient:
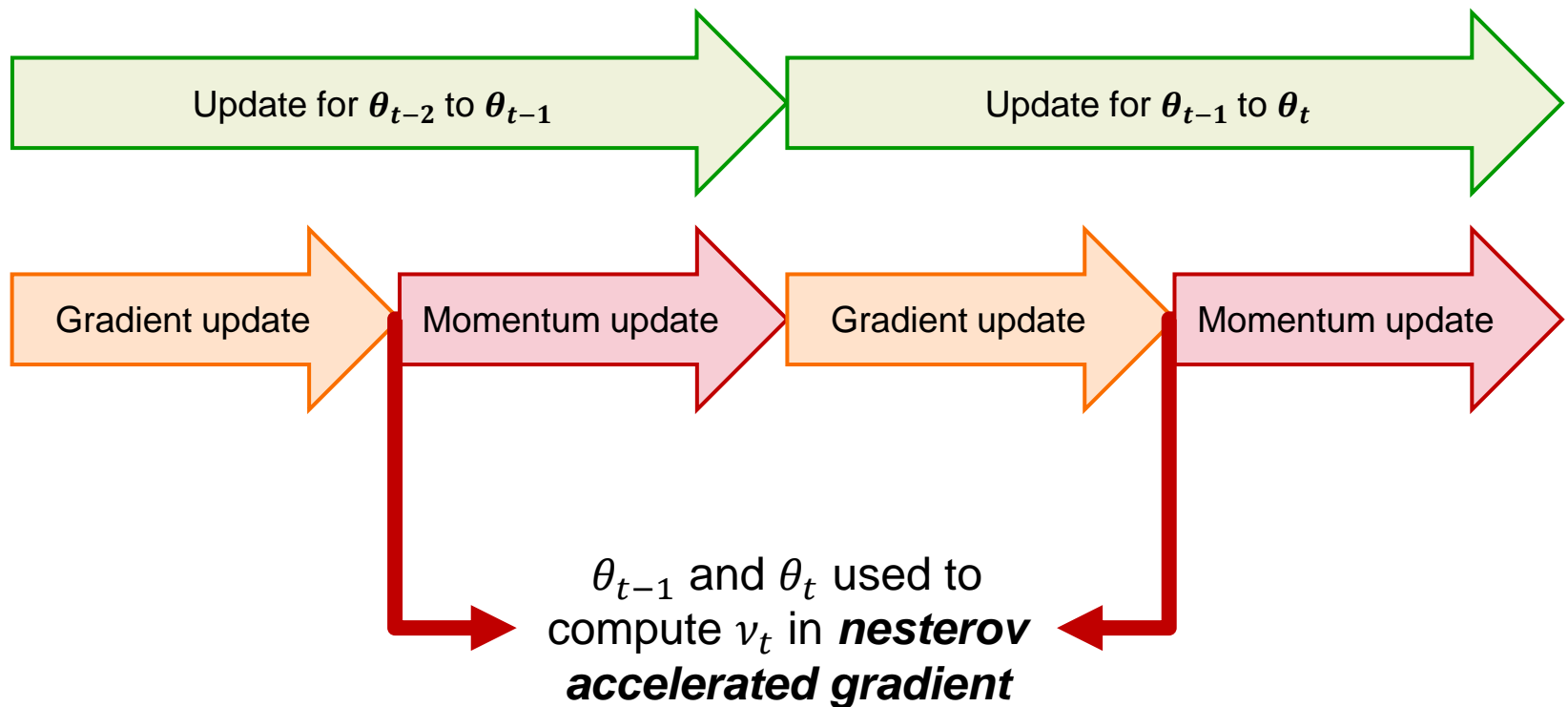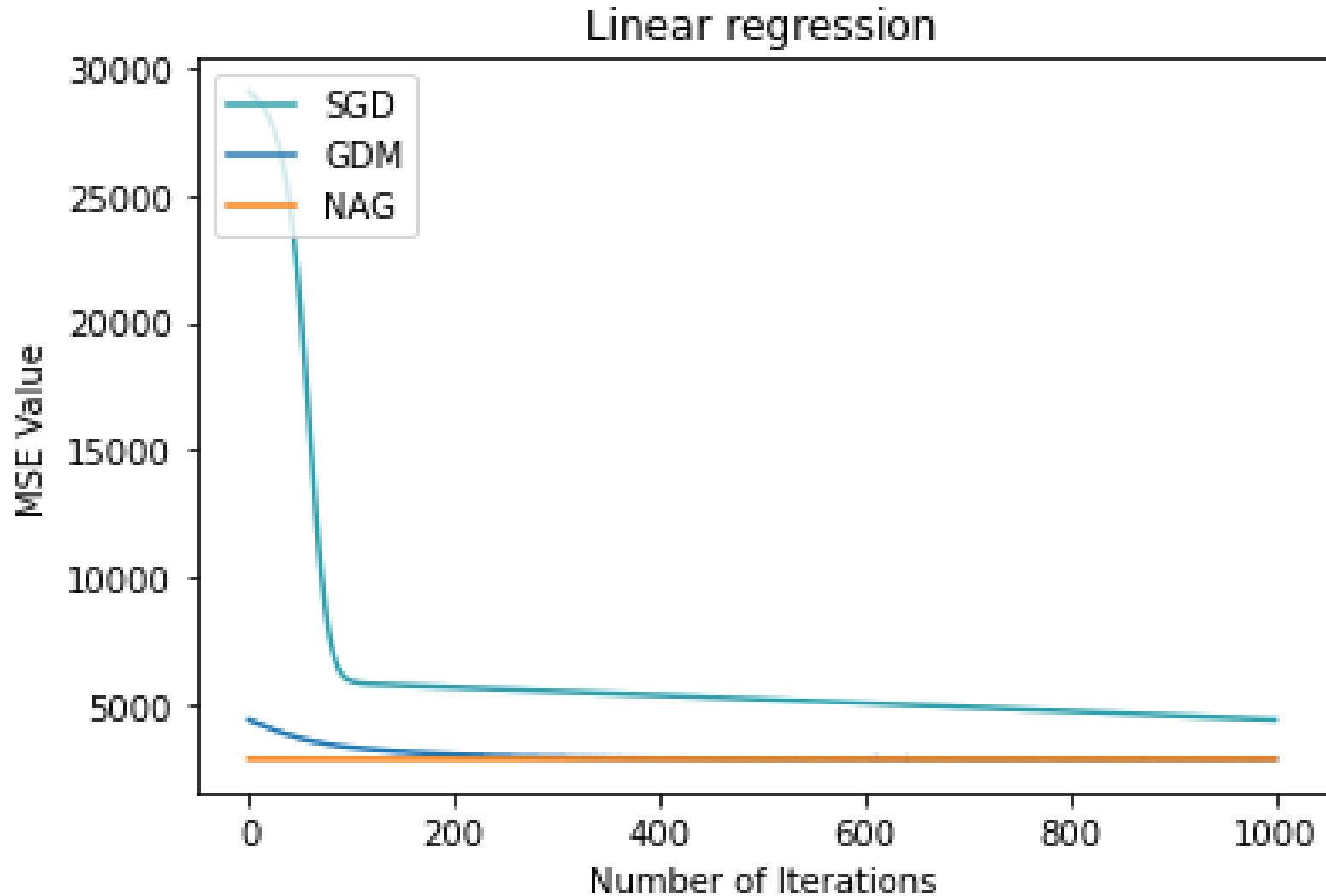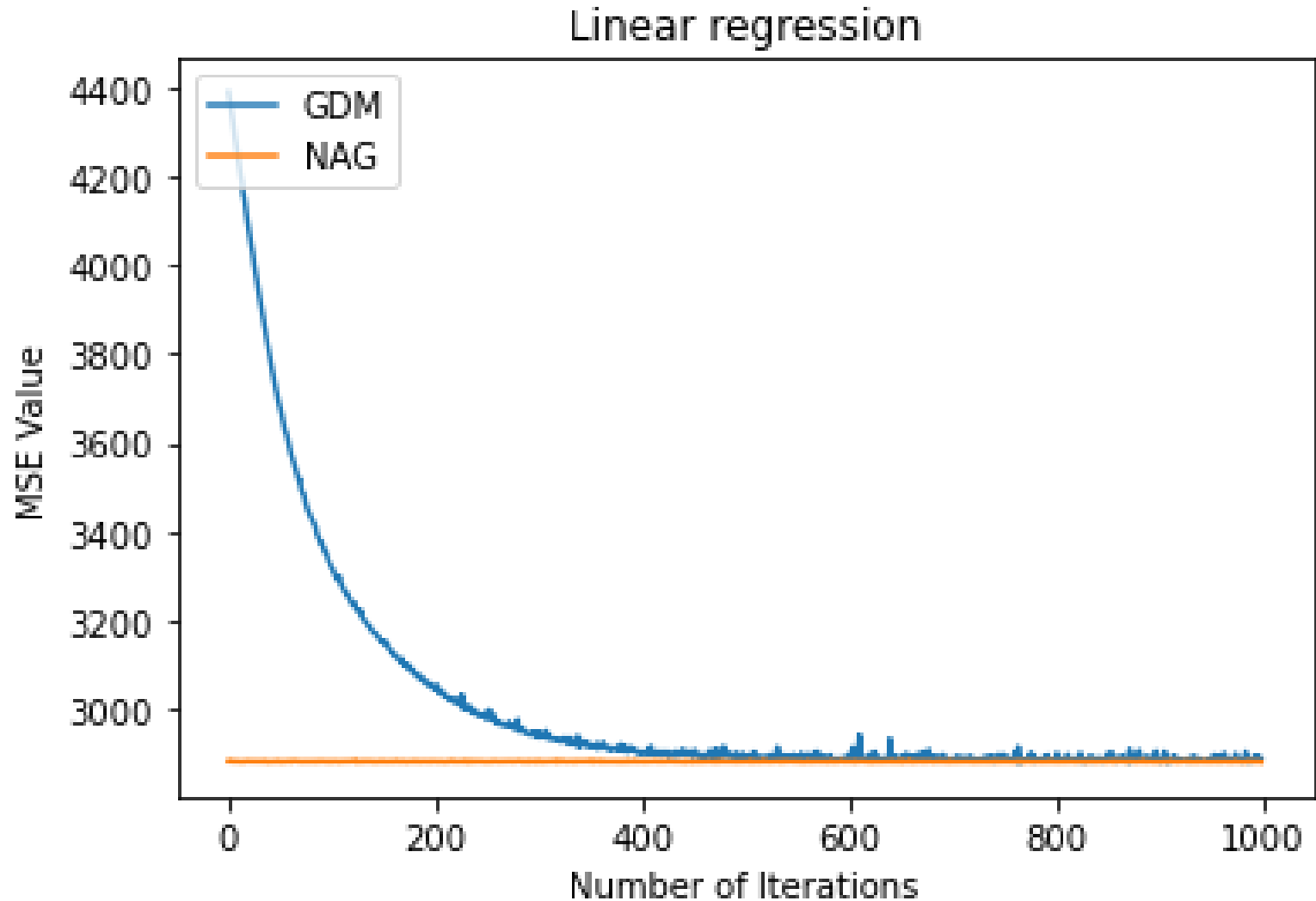NAG splits up gradient update step and momentum update step

Update for $\boldsymbol{\theta}_{t-2}$ to $\boldsymbol{\theta}_{t-1}$

Update for $\boldsymbol{\theta}_{t-1}$ to $\boldsymbol{\theta}_t$

Gradient update

Momentum update

Gradient update

Momentum update

$\theta_{t-1}$ and $\theta_t$ used to compute $v_t$ in ***nesterov accelerated gradient***

# Comparing Methods for Linear Regression

| Tested Model: | Linear Regression |
|---|---|
| Dataset: | Diabetes Dataset (Scikit-Learn) |
| Number of Samples: | 422 |
| Number of Features | 10 |
| Loss function: | Mean Squared Error |
| Learning Rate: | 0.00001 |
| Momentum rate (if Momentum): | 0.9 |
| Number of iterations: | 1000 |
| Random Seed: | 57 |

Technische Universität Braunschweig

igp

# Comparing Methods for Linear Regression



Linear regression

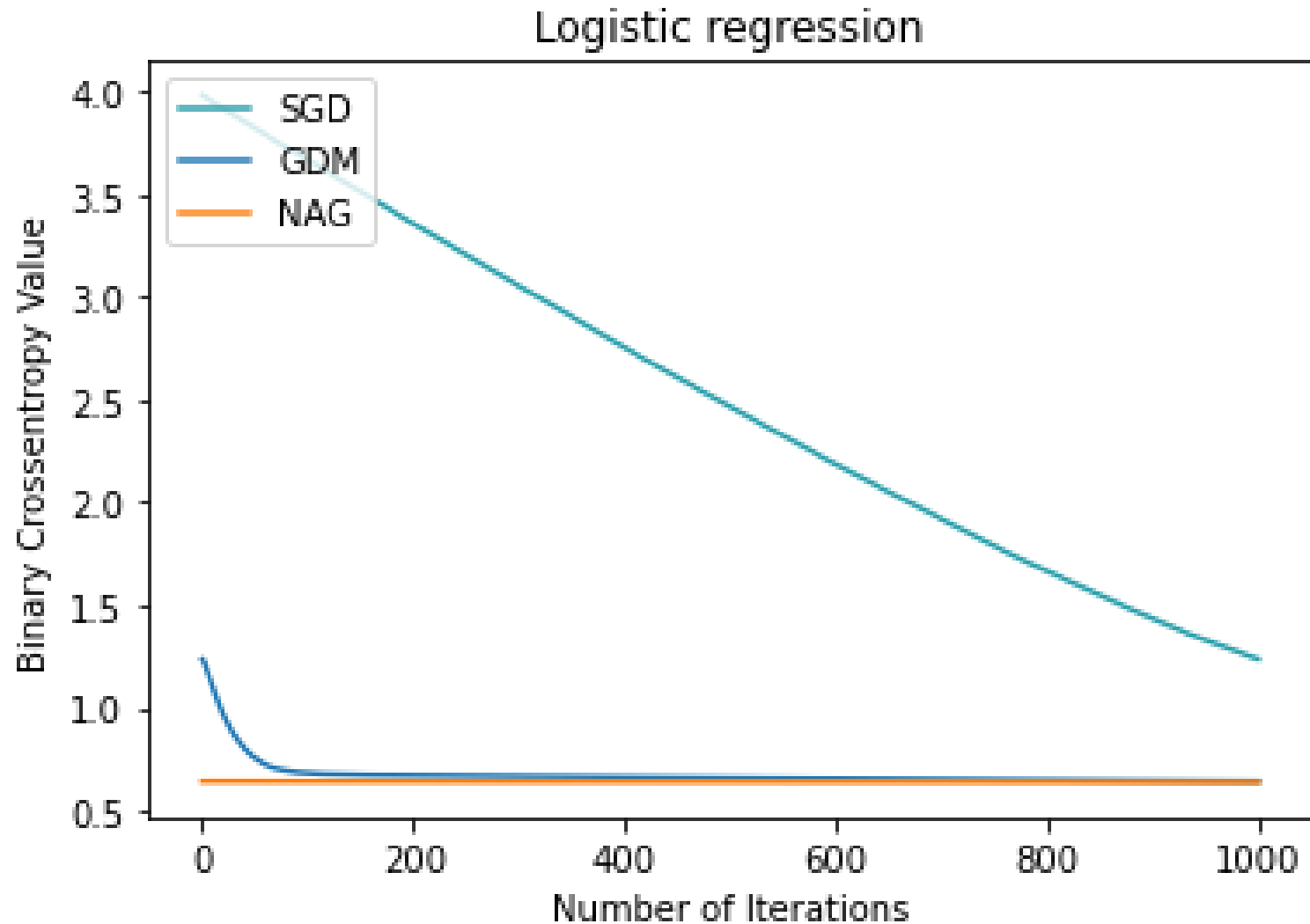# Comparing Methods for Linear Regression



Linear regression

# Comparing Methods for Linear Regression

| Tested Model: | Logistic Regression |
|---|---|
| Dataset: | Iris Dataset (Scikit-Learn) |
| Number of Samples: | 100 |
| Number of Features | 4 |
| Loss function: | Binary Crossentropy |
| Learning Rate: | 0.001 |
| Momentum rate (if Momentum): | 0.9 |
| Number of iterations: | 1000 |
| Random Seed: | 57 |

Technische
Universität
Braunschweig

igp

# Comparing Methods for Logistic Regression



Logistic regression

# Comparing Methods for Logistic Regression



Logistic regression

# Library Implementations for Sklearn/Pytorch/Keras

**Scikit-learn:**      No Explicit Bulit-in Implementation

**Theano/Lasagne:**

- Function sgd for gradient descent
- Function apply_momentum for momentum
- Function apply_nesterov_momentum for NAG

**Pytorch:**

- Class SGD for gradient descent
- momentum as float argument
- nesterov as boolean argument

**Tensorflow/Keras:**

- Class SGD for gradient descent
- momentum as float argument
- nesterov as boolean argument

# Momentum and NAG does not give you better results, but it give you a good result much faster!

**You might use it when:**

- …the prediction just have to be good, but not perfect
- …fast Training process is required

**For Example: model is trained live while usage multiple times**

**You might not use it when:**

- …there is a high requirement for the best accuracy possible
- …time is no important matter

**For example: model is only trained once without gathering new data**

# References

http://www.cs.utoronto.ca/~ilya/pubs/ilya_sutskever_phd_thesis.pdf

https://jlmelville.github.io/mize/nesterov.html

https://scikit-learn.org/stable/index.html

https://github.com/Lasagne/Lasagne

https://pytorch.org/

https://keras.io/