



CAMBRIDGE

Lecture 1: Regression

Stefan Bucher

MACHINE LEARNING IN ECONOMICS
UNIVERSITY OF CAMBRIDGE

Stefan Bucher



Prince (2023, chaps. 2, 6, 8, 9).¹

1. Figures taken or adapted from Prince (2023). All rights belong to the original author and publisher. These materials are intended solely for educational purposes.



Linear Regression

Prince (2023, chap. 2)



Linear Model

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon$$

Minimizing the mean-squared error (MSE)

$$\min_{\beta} \frac{1}{n} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta)$$



Minimizing MSE Loss (2.1-2.3)



Analytical Solution

$$\min_{\beta} \frac{1}{n} (y - X\beta)^T (y - X\beta)$$

Setting the gradient with respect to β to zero

$$0 = \nabla_{\beta} [y^T y - 2\beta^T X^T y + \beta^T X^T X \beta] = -2X^T y + 2X$$

yields the BLUE estimator (Gauss-Markov)

$$\hat{\beta} = (X^T X)^{-1} X^T y$$



scikit-learn



scikit-learn

```
1 from sklearn.linear_model import LinearRegression  
2 model = LinearRegression()  
3 model.fit(X_train, y_train)  
4 print(model.coef_[2], model.intercept_)
```

546.1898661888232 152.05949984245817

Text(0, 0.5, 'Diabetes Risk')



Model Fitting & Optimization

$$\hat{\phi} = \arg \min_{\phi} L(\phi)$$

Prince (2023, chap. 6)

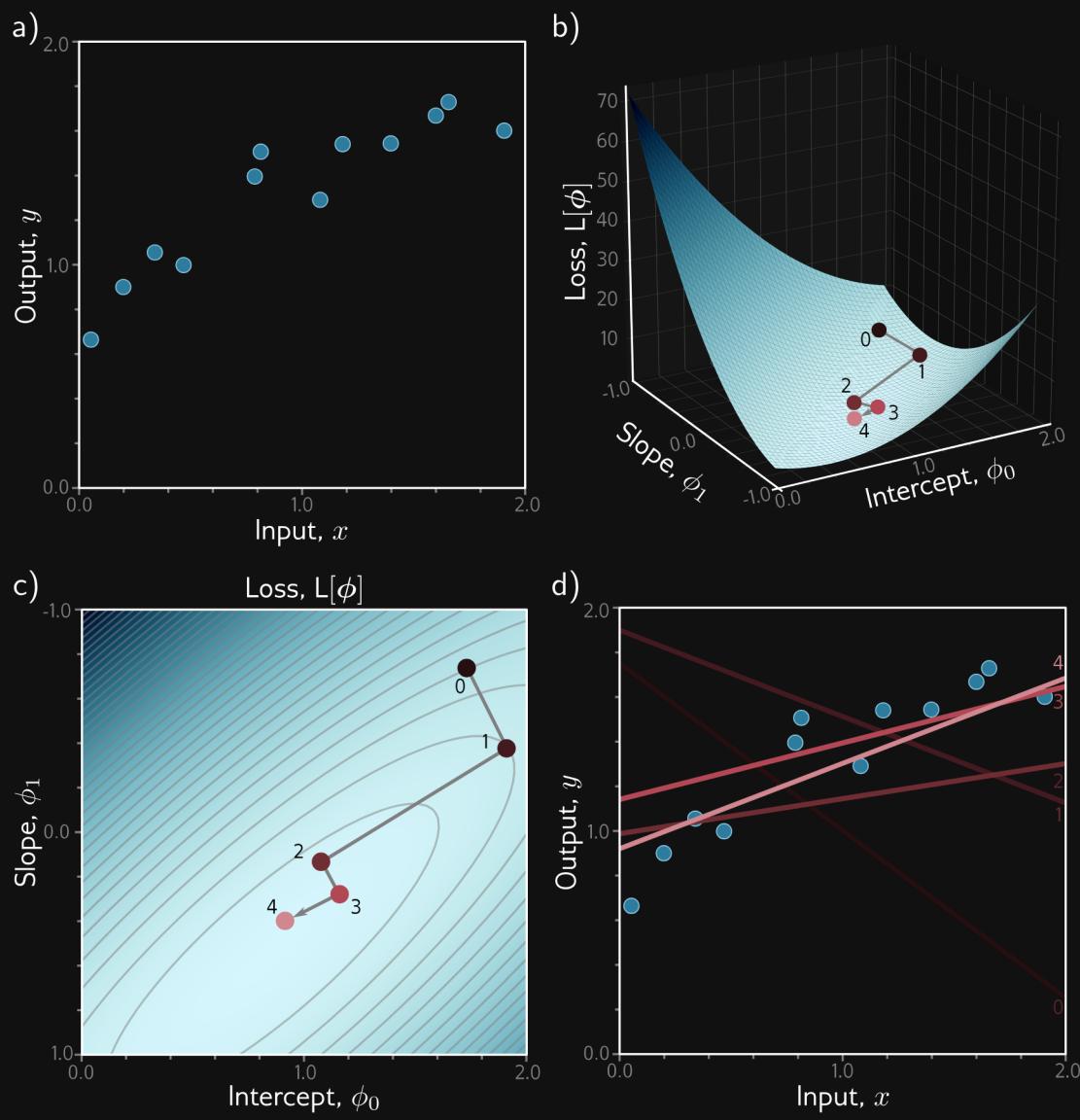


Gradient Descent¹

$$\phi \leftarrow \phi - \alpha \nabla_{\phi} L(\phi)$$

where α fixed learning rate or line search.





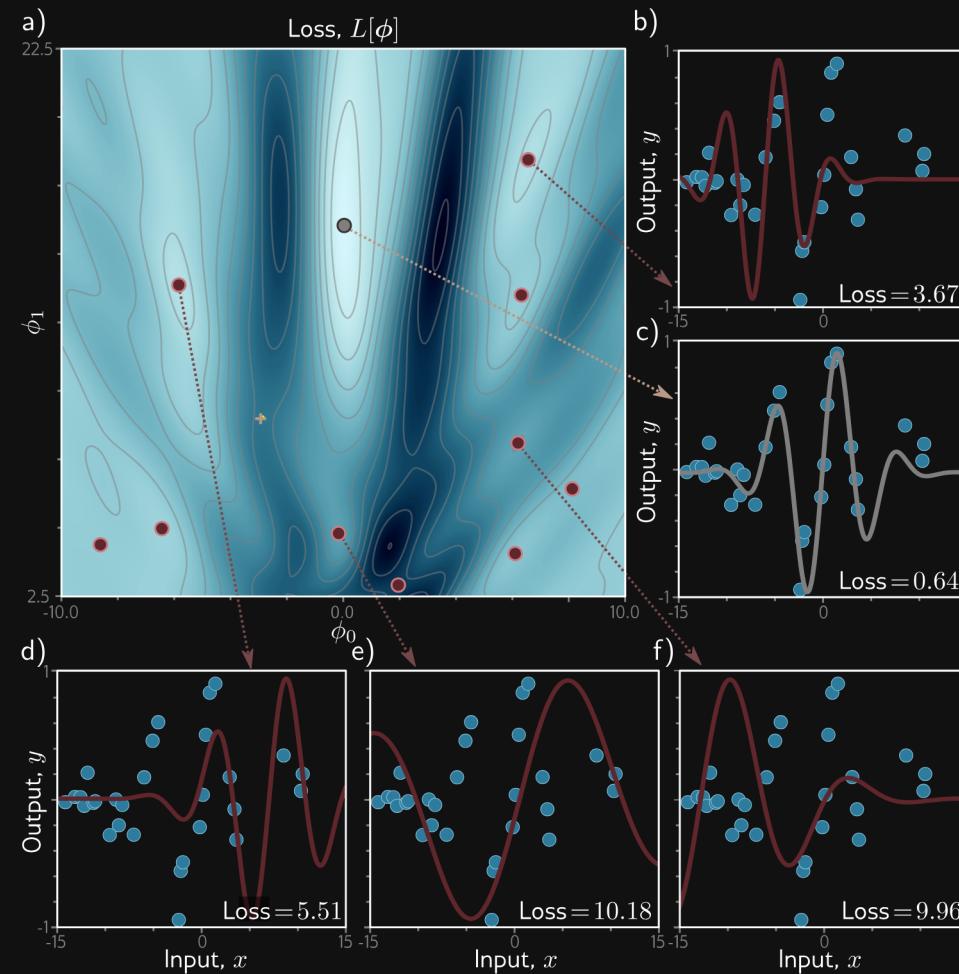
1. Cauchy (1847)

Stefan Bucher

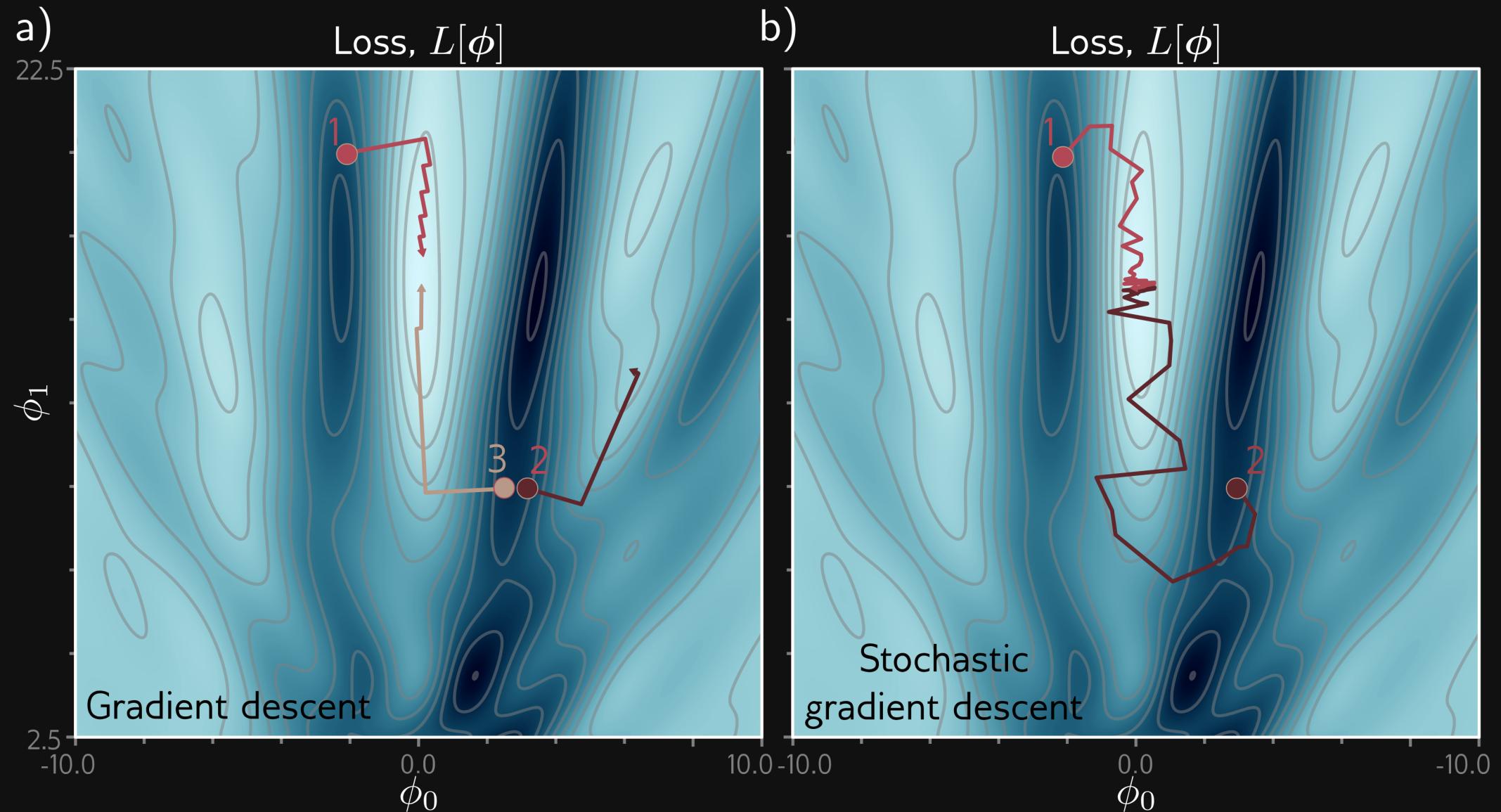


Nonlinear Model

Loss functions generally non-convex, possibly multiple local minima.



Stochastic Gradient Descent



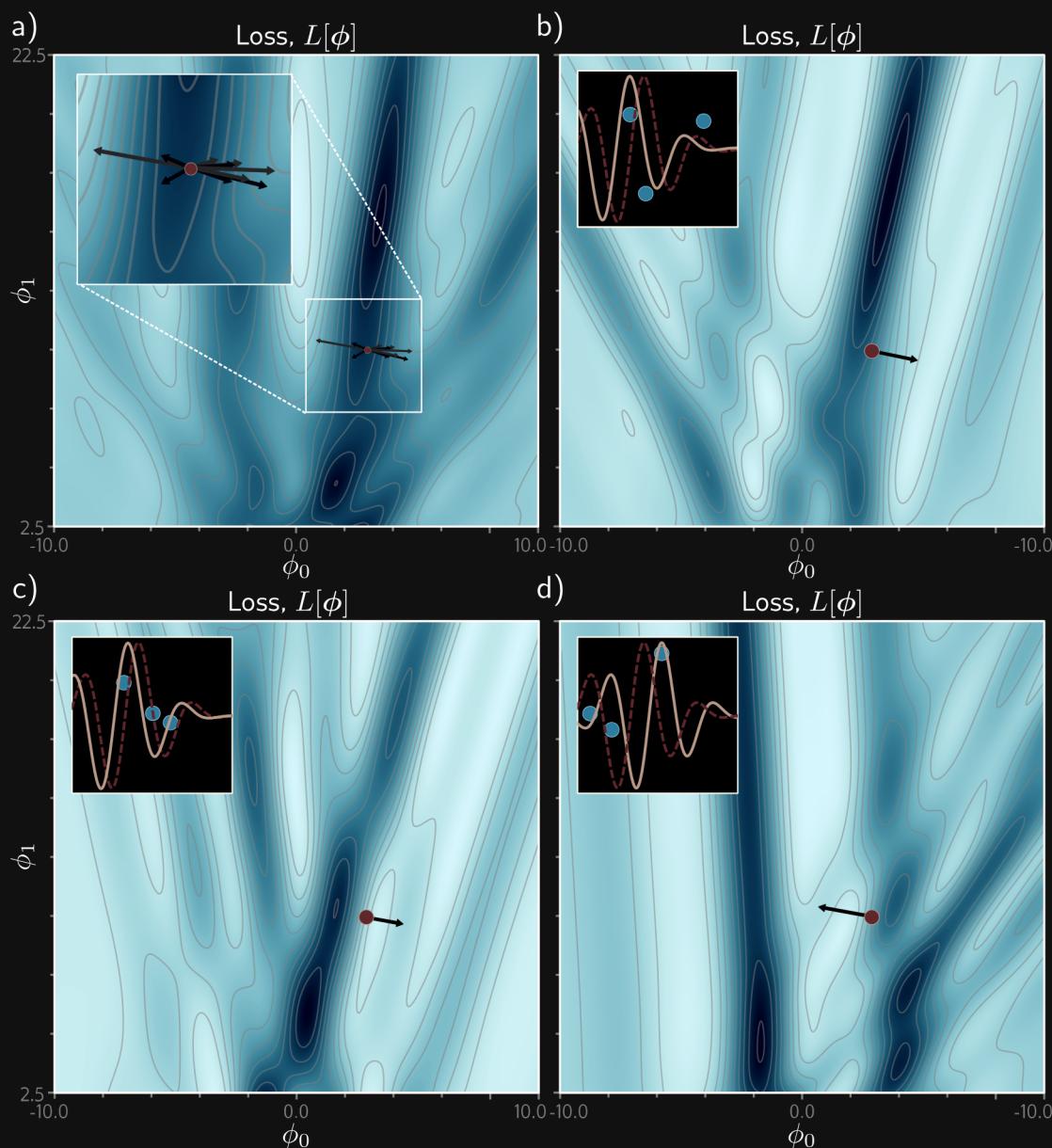
Stochastic Gradient Descent¹

Compute gradient on a mini-batch B_t of data points

$$\phi_{t+1} \leftarrow \phi_t - \alpha \sum_{i \in B_t} \nabla_{\phi} l_i(\phi_t; x_i, y_i)$$

- Sampled without replacement within each epoch (sweep of training data).
- Learning rate often decreasing over time.





1. Robbins & Monro (1951)



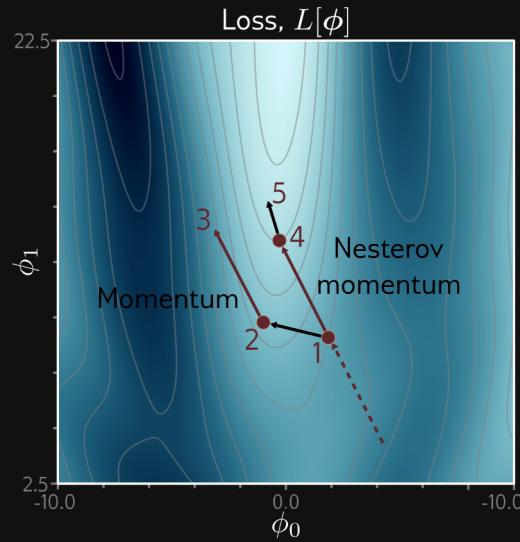
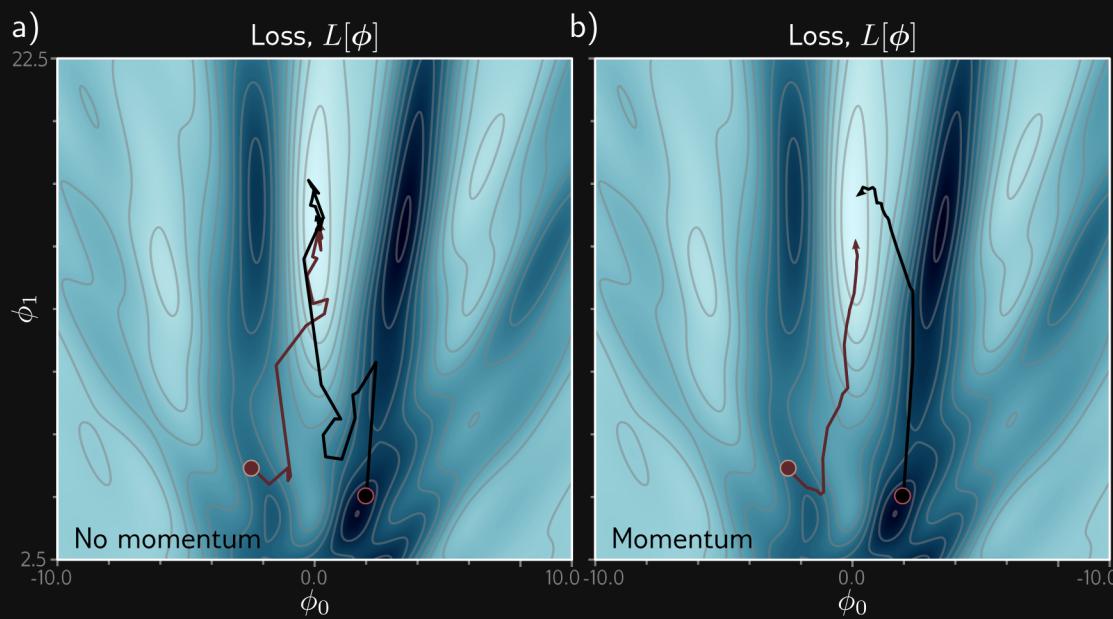
Nesterov Accelerated Momentum¹

$$m_{t+1} \leftarrow \beta m_t + (1 - \beta) \sum_{i \in B_t} \nabla_{\phi} l_i(\phi_t - \alpha \beta m_t; x_i, y_i)$$

$$\phi_{t+1} \leftarrow \phi_t - \alpha m_{t+1}$$

adds momentum (m_t) *before* evaluating gradient.





1. Nesterov (1983)

Stefan Bucher



Adaptive Moment Estimation (Adam)¹

Normalizes the gradient while relying on momentum to avoid convergence issues.



$$m_{t+1} \leftarrow \beta m_t + (1 - \beta) \sum_{i \in B_t} \nabla_{\phi} l_i(\phi_t; x_i, y_i)$$

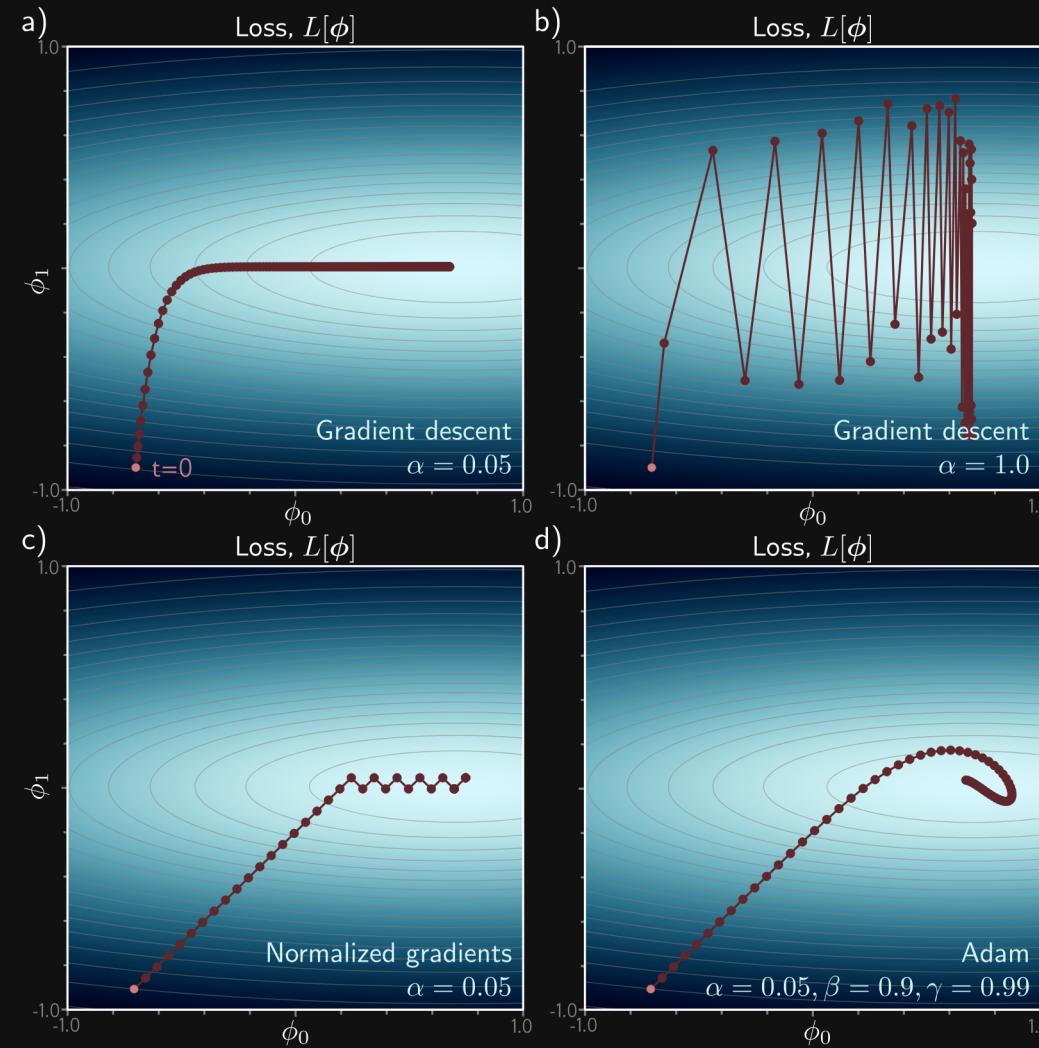
$$v_{t+1} \leftarrow \gamma v_t + (1 - \gamma) \left(\sum_{i \in B_t} \nabla_{\phi} l_i(\phi_t; x_i, y_i) \right)^2$$

$$\phi_{t+1} \leftarrow \phi_t - \alpha \frac{\frac{m_{t+1}}{1-\beta^{t+1}}}{\sqrt{\frac{v_{t+1}}{1-\gamma^{t+1}}} + \epsilon}$$

1. Kingma & Ba (2015)



Adaptive Moment Estimation (Adam)



PyTorch



PyTorch



Setup

```
1 import torch
2 from torch import nn
3 from torchvision import datasets
4 from torch.utils.data import DataLoader
5 from torchvision.transforms import ToTensor, Compose, Grayscale
6 import random
7
8 if torch.backends.mps.is_available():
9     DEVICE = torch.device("mps") # Apple Silicon Metal Performance Shaders
10 elif torch.cuda.is_available():
11     DEVICE = torch.device("cuda") # GPU
12 else:
13     DEVICE = torch.device("cpu") # CPU
```

Key components:

- tensors are like numpy arrays but GPU-accelerated
- autograd



Tensors

Tensors “live” on a device

```
1 x = torch.randn(2, 2, device="cpu")
2 print(x.device)
3 x = x.to(DEVICE)
4 print(x.device)
```

cpu
mps:0



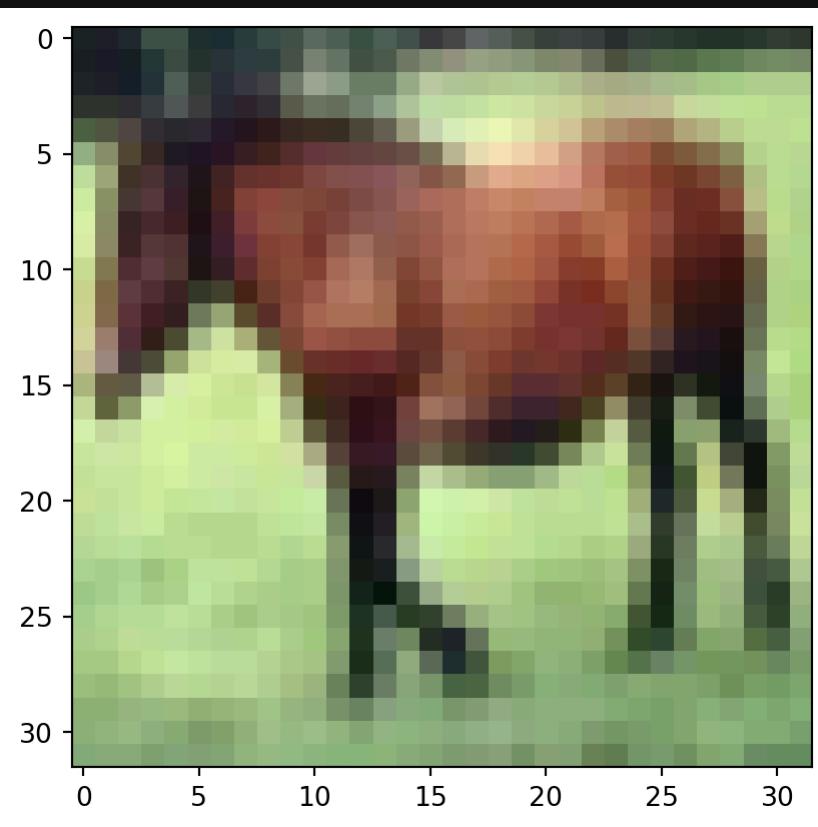
Datasets

```
1 training_data = datasets.CIFAR10(  
2     root="data",    # path where the images will be stored  
3     train=True,  
4     download=True,  
5     transform=ToTensor()  
6     )  
7 test_data = datasets.CIFAR10(  
8     root="data",  
9     train=False,  
10    download=True,  
11    transform=ToTensor()  
12    )  
13 image, label = training_data[7]  
14 print(f"Label: {training_data.classes[label]}")  
15 print(f"Image size: {image.shape}")  
16 plt.imshow(image.permute(1, 2, 0)) # imshow expects channel order Height x  
17 plt.show()
```

Label: horse

Image size: torch.Size([3, 32, 32])





Stefan Bucher



Data Loader

```
1 # Data Loader for Batching
2 def seed_worker(worker_id):
3     worker_seed = torch.initial_seed() % 2**32
4     numpy.random.seed(worker_seed)
5     random.seed(worker_seed)
6     g_seed = torch.Generator()
7     g_seed.manual_seed(torch.initial_seed() % 2**32)
8
9 train_dataloader = DataLoader(training_data, batch_size=64, shuffle=True, n
10 test_dataloader = DataLoader(test_data,   batch_size=64, shuffle=True, num_w
11
12 # Load the next batch
13 batch_images, batch_labels = next(iter(train_dataloader))
14 print('Batch size:', batch_images.shape)
15 plt.imshow(batch_images[0].permute(1, 2, 0))
16 plt.show()
```





References

Prince, Simon J. D. 2023. *Understanding Deep Learning*. Cambridge, Massachusetts: The MIT Press.

