

$$F = G \frac{m_1 m_2}{d^2}$$

# Deep Learning for Particle Physicists

Lewis Tunstall | AEC Graduate Seminar | May 6th 2022

$$\frac{\partial^2 u}{\partial t^2} = c^2 \frac{\partial^2 u}{\partial x^2}$$

$$\frac{df}{dt} = \lim_{h \rightarrow 0} \frac{f(t+h) - f(t)}{h}$$



# Housekeeping



## Timetable

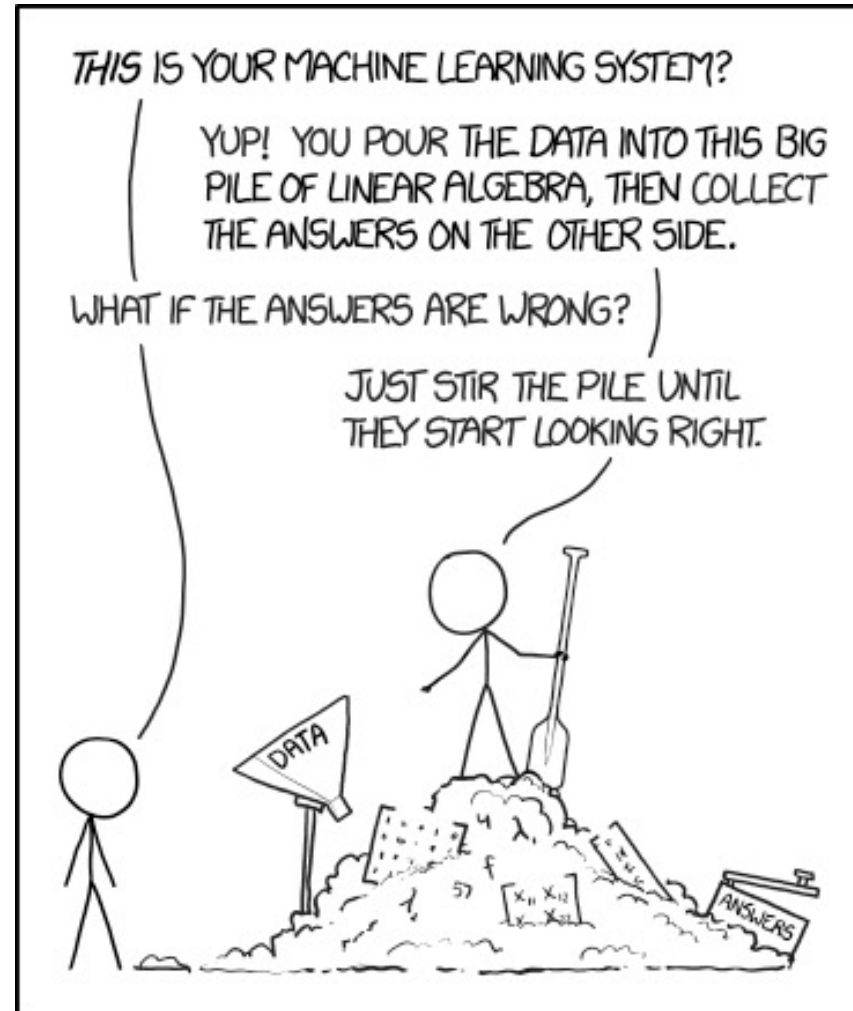
14:15-16:00 every Friday in Room 119 from April 29 to **June 10**



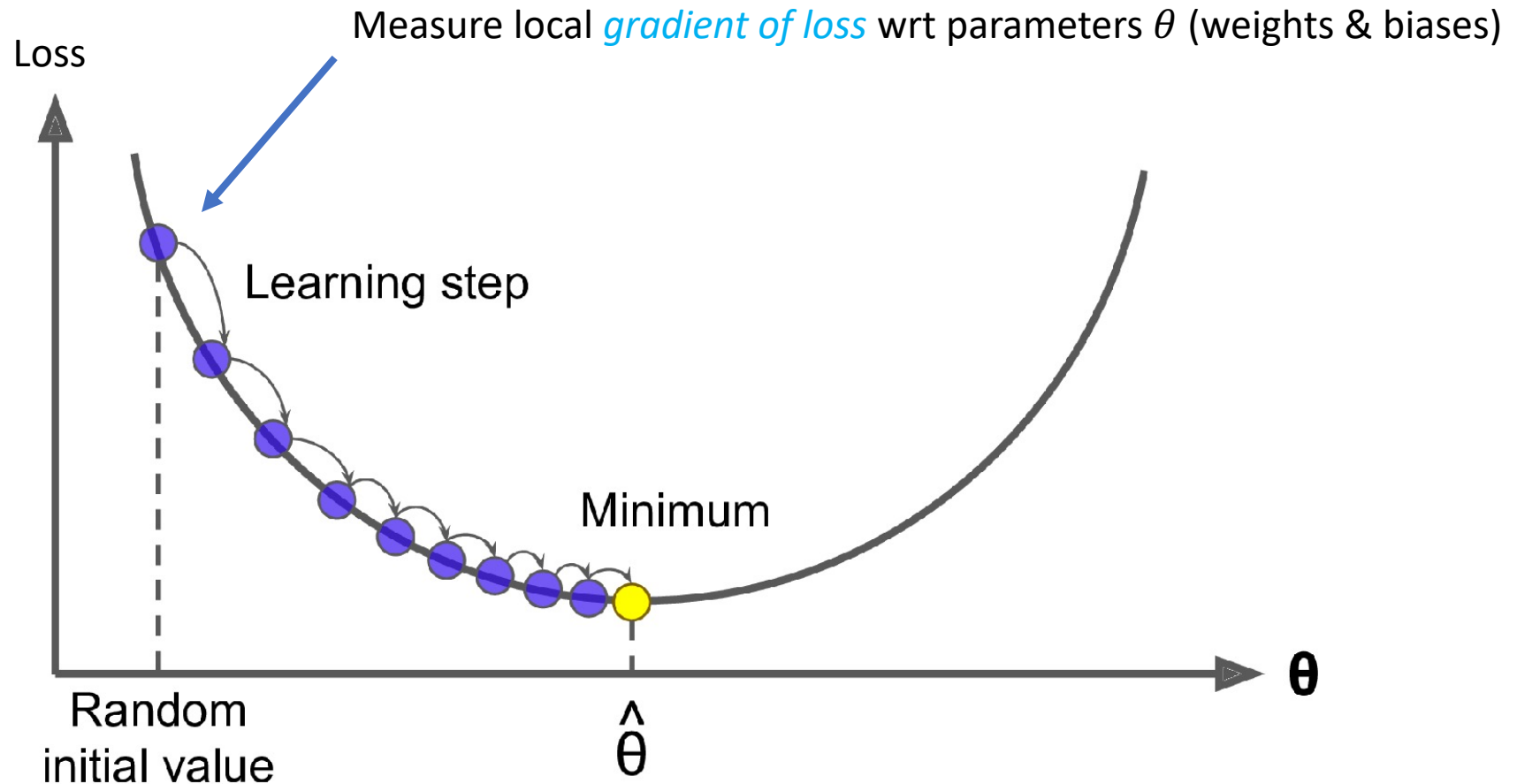
## Hugging Face Hub

- <https://huggingface.co/join>
- <https://huggingface.co/dl4phys> (organisation for this class)

# Training models



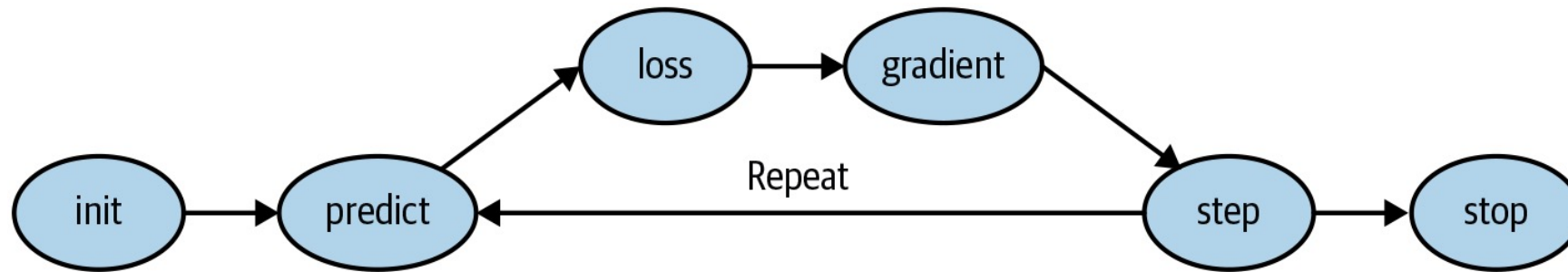
# Gradient descent



[Hands-On Machine Learning with Scikit-Learn & TensorFlow, A. Geron](#)

**Basic idea:** tweak parameters *iteratively* to minimise a *loss function*

# Gradient descent




Deep Learning for Coders with fastai and PyTorch, J. Howard & S. Gugger

7 main steps for training ML models and deep NNs

# Common loss functions

model  
prediction

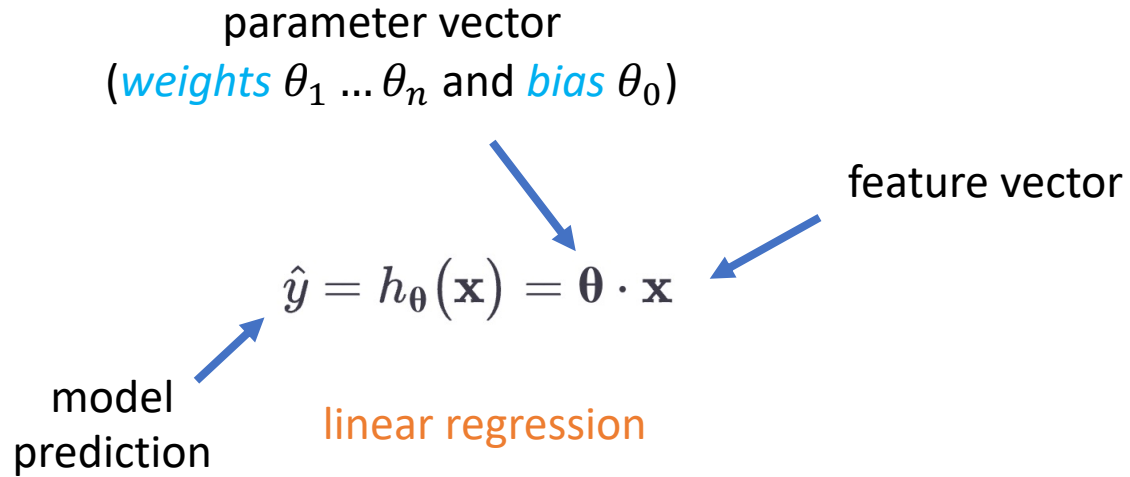

$$\hat{y} = h_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta} \cdot \mathbf{x}$$

linear regression

$$\text{MSE}(\mathbf{X}, h_{\boldsymbol{\theta}}) = \frac{1}{m} \sum_{i=1}^m \left( \boldsymbol{\theta}^{\top} \mathbf{x}^{(i)} - y^{(i)} \right)^2$$

mean squared error

# Common loss functions

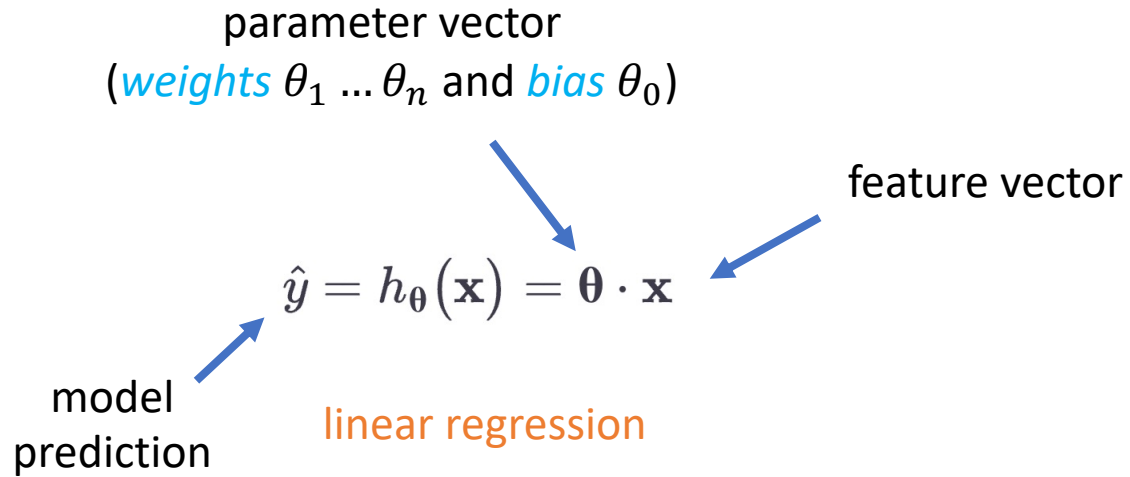


$$\text{MSE}(\mathbf{X}, h_{\theta}) = \frac{1}{m} \sum_{i=1}^m \left( \boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)} \right)^2$$

mean squared error

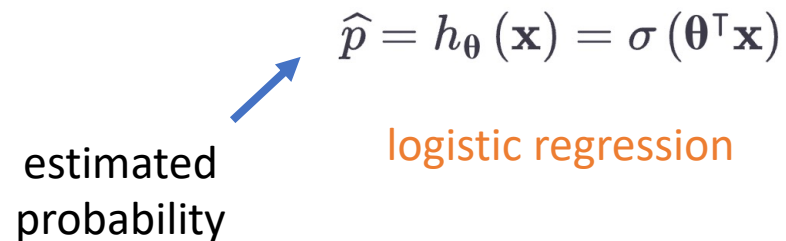


# Common loss functions



$$\text{MSE}(\mathbf{X}, h_{\theta}) = \frac{1}{m} \sum_{i=1}^m \left( \boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)} \right)^2$$

mean squared error



$$J(\boldsymbol{\theta}) = -\frac{1}{m} \sum_{i=1}^m \left[ y^{(i)} \log(\hat{p}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{p}^{(i)}) \right]$$

binary cross entropy

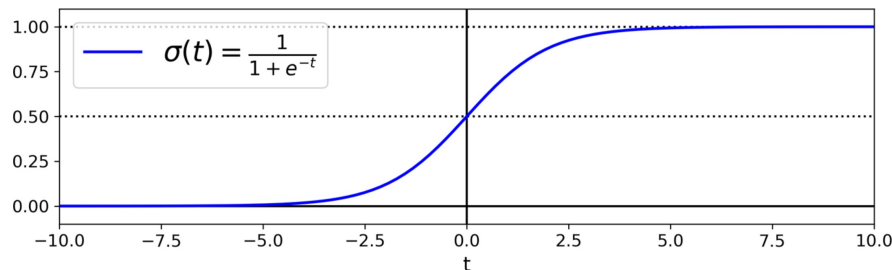
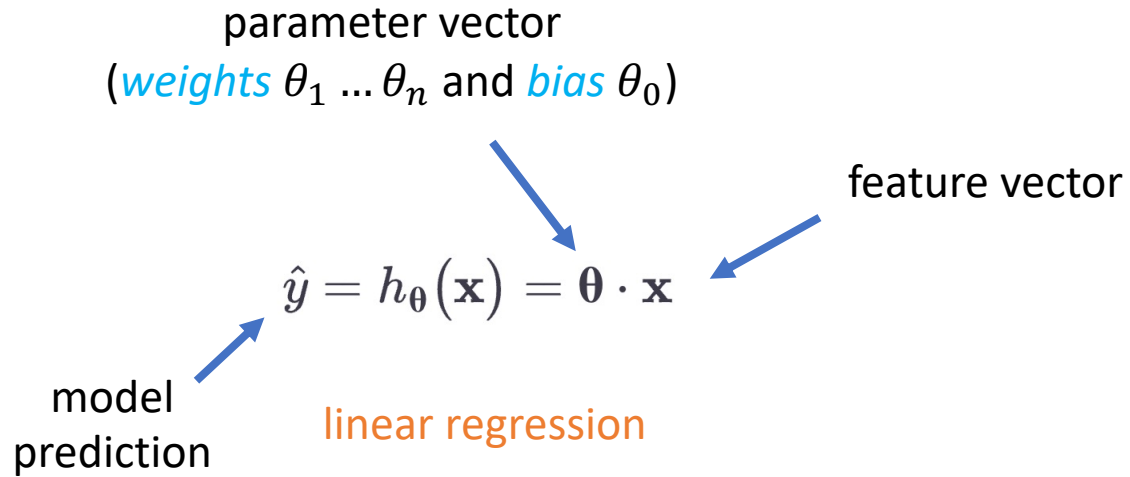


Figure 4-21. Logistic function



# Common loss functions



$$\text{MSE}(\mathbf{X}, h_{\theta}) = \frac{1}{m} \sum_{i=1}^m \left( \boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)} \right)^2$$

mean squared error

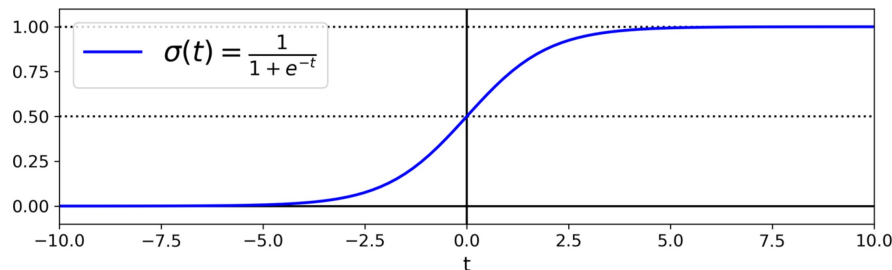
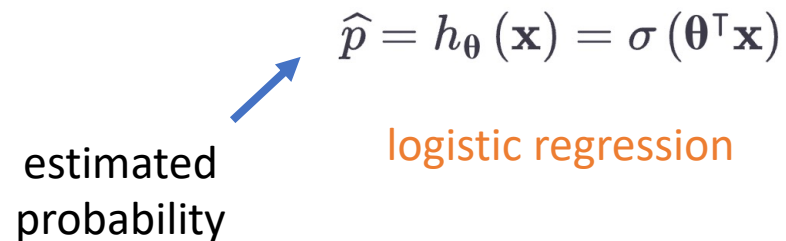


Figure 4-21. Logistic function

$$J(\boldsymbol{\theta}) = -\frac{1}{m} \sum_{i=1}^m \left[ y^{(i)} \log(\hat{p}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{p}^{(i)}) \right]$$

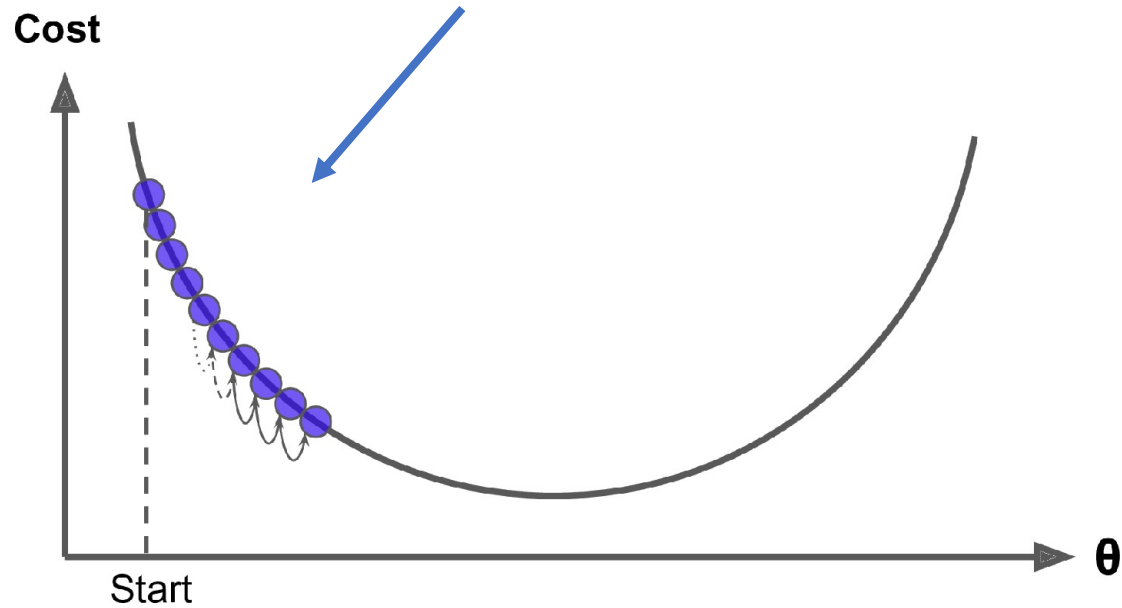
binary cross entropy

generalises to *cross entropy*  
for  $K > 2$  classes

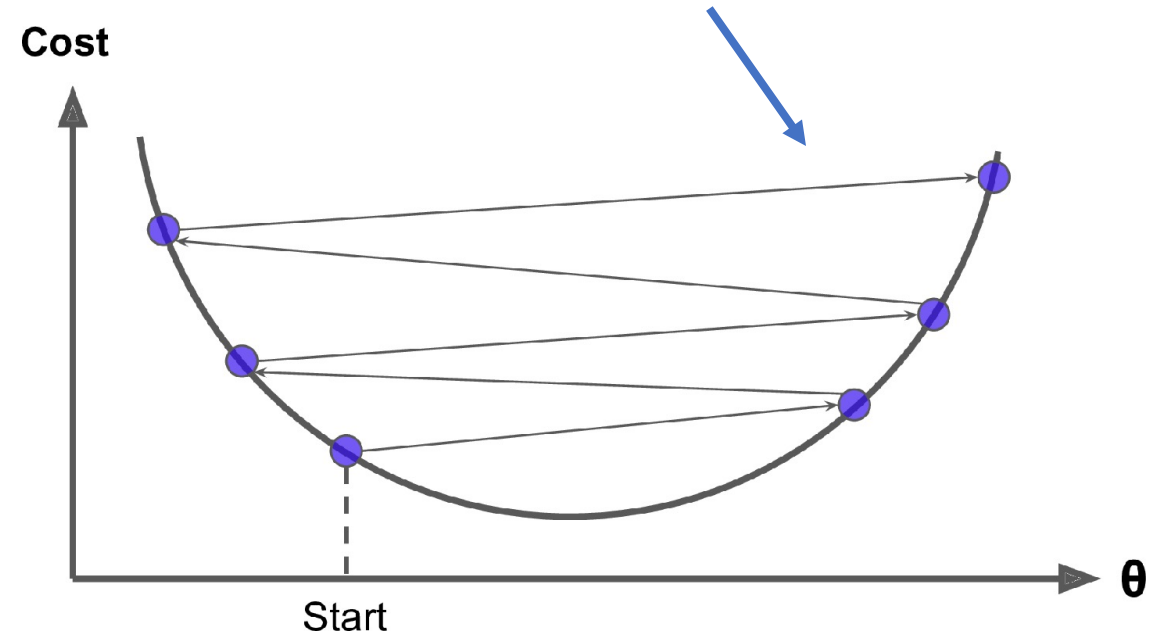
$$J(\Theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(\hat{p}_k^{(i)})$$

# Gradient descent

Takes many steps to converge



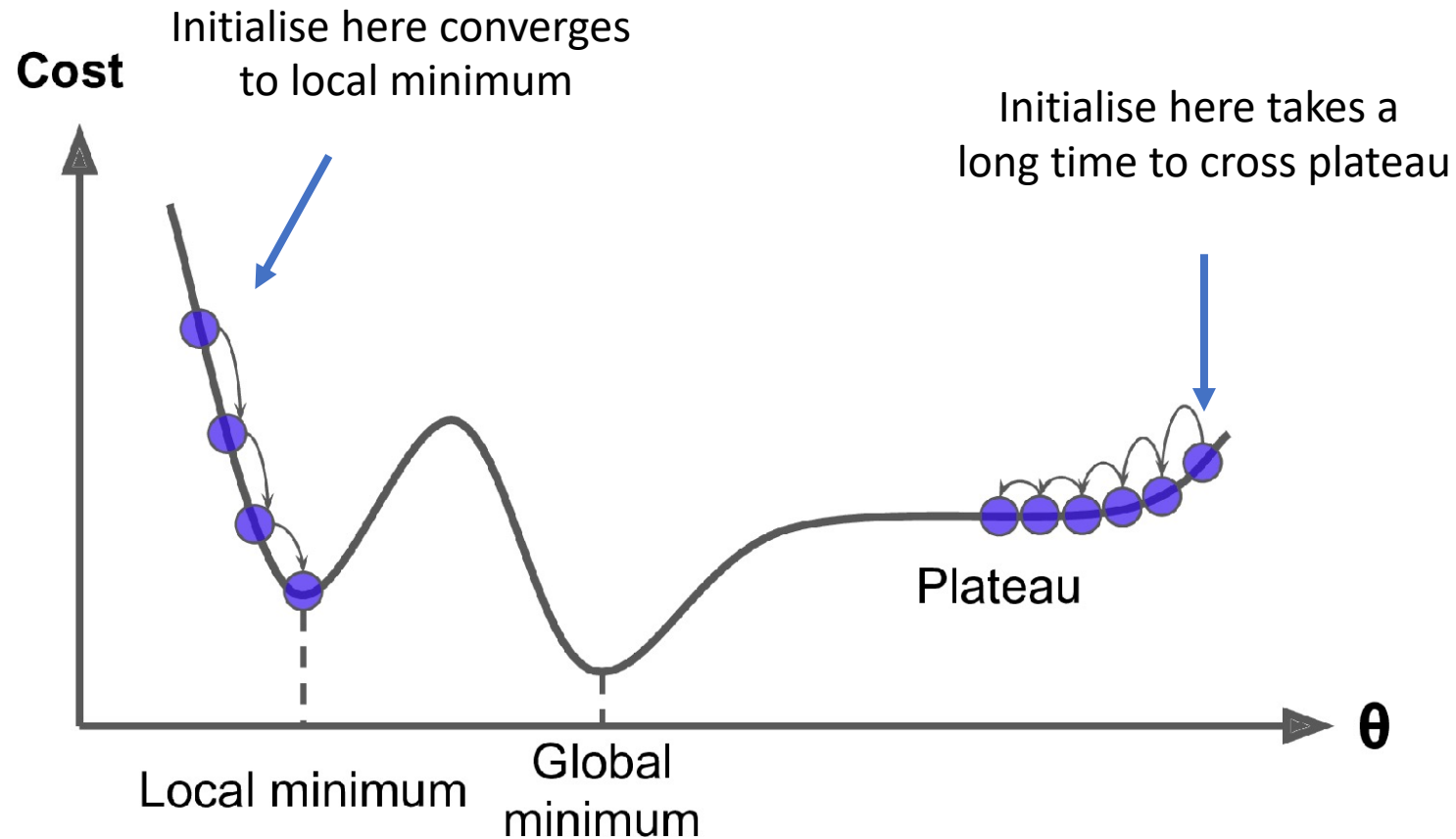
Overshoots the minimum & possibly diverges



Hands-On Machine Learning with Scikit-Learn & TensorFlow, A. Geron

*Small* vs *large* learning rates

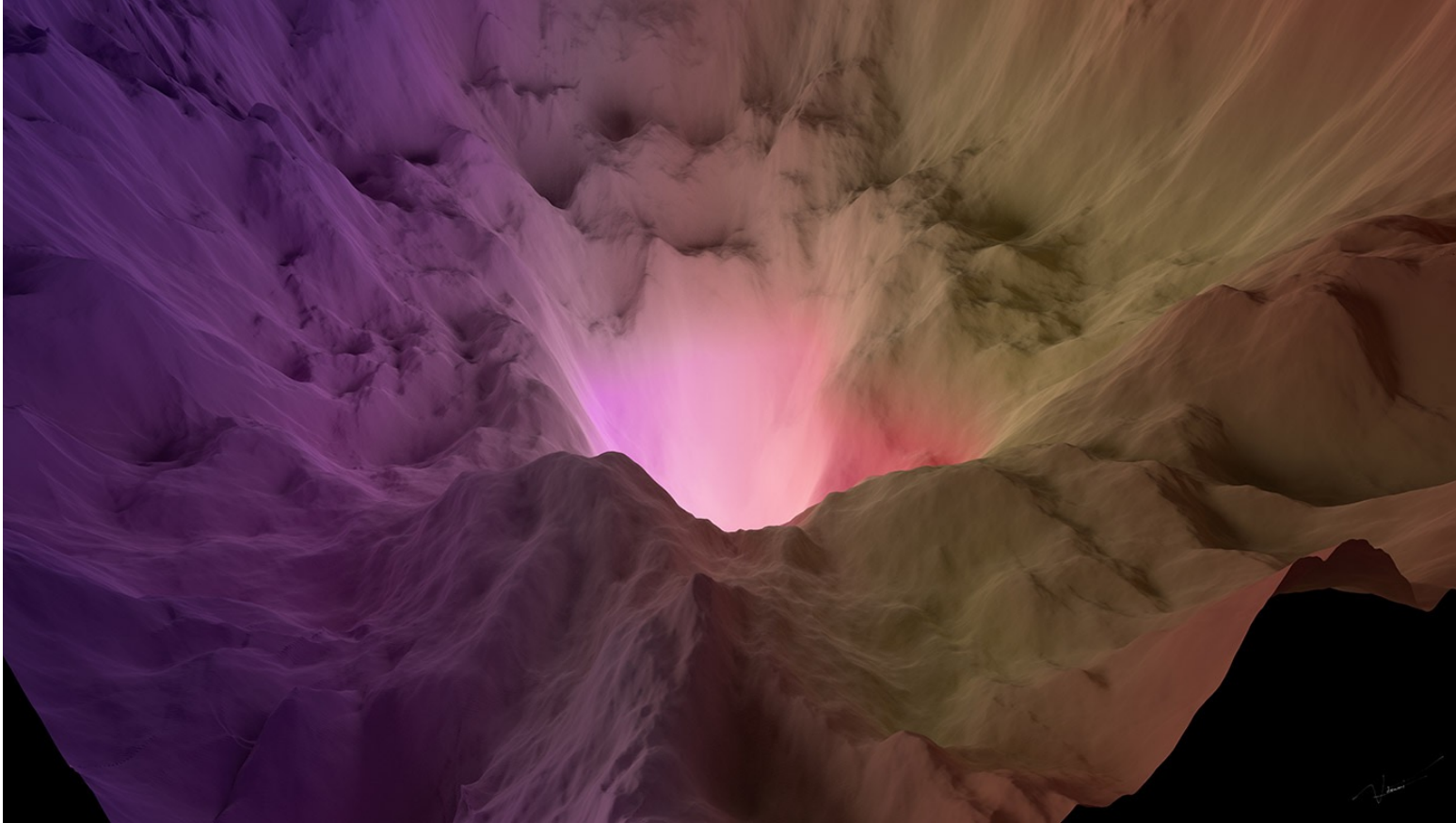
# Gradient descent



Hands-On Machine Learning with Scikit-Learn & TensorFlow, A. Geron

Loss functions have multiple minima, plateaus, ridges etc

# Gradient descent

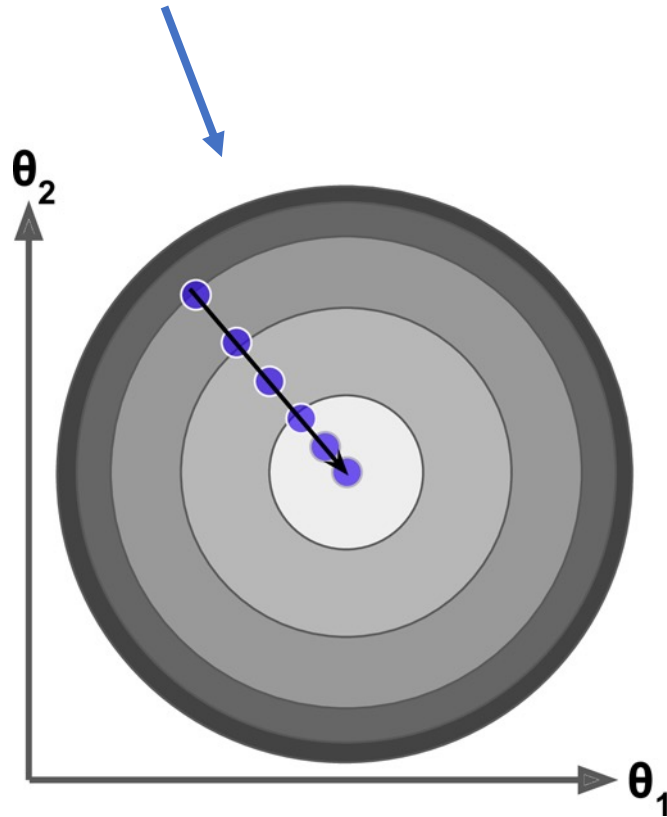


<https://losslandscape.com/gallery/>

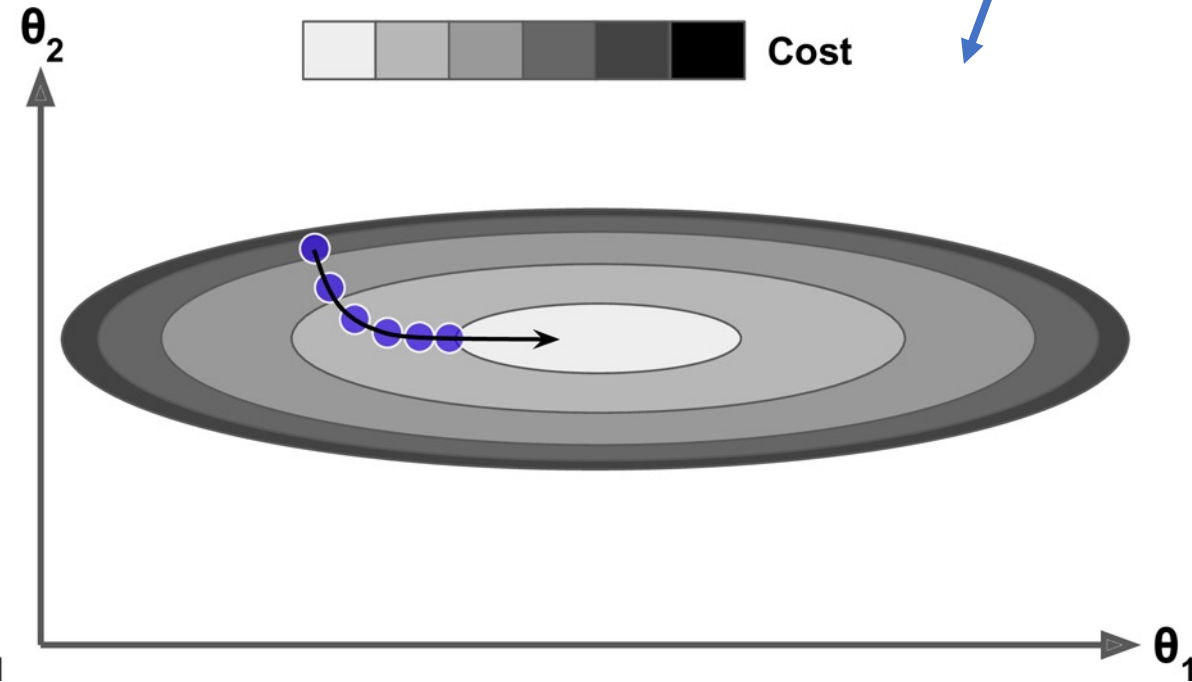
Loss functions have multiple minima, plateaus, ridges etc

# Gradient descent

Similar feature scales  
⇒ quick to find minimum



Longer time to converge when  
features scales are very different



Hands-On Machine Learning with Scikit-Learn & TensorFlow, A. Geron

## Feature scaling

# Batch gradient descent

$$\text{MSE}(\mathbf{X}, h_{\boldsymbol{\theta}}) = \frac{1}{m} \sum_{i=1}^m \left( \boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)} \right)^2$$



the loss we wish to minimise

$$\frac{\partial}{\partial \theta_j} \text{MSE}(\boldsymbol{\theta}) = \frac{2}{m} \sum_{i=1}^m \left( \boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)} \right) x_j^{(i)}$$



the gradients wrt each parameter

$$\nabla_{\boldsymbol{\theta}} \text{MSE}(\boldsymbol{\theta}) = \begin{pmatrix} \frac{\partial}{\partial \theta_0} \text{MSE}(\boldsymbol{\theta}) \\ \frac{\partial}{\partial \theta_1} \text{MSE}(\boldsymbol{\theta}) \\ \vdots \\ \frac{\partial}{\partial \theta_n} \text{MSE}(\boldsymbol{\theta}) \end{pmatrix} = \frac{2}{m} \mathbf{X}^T (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})$$



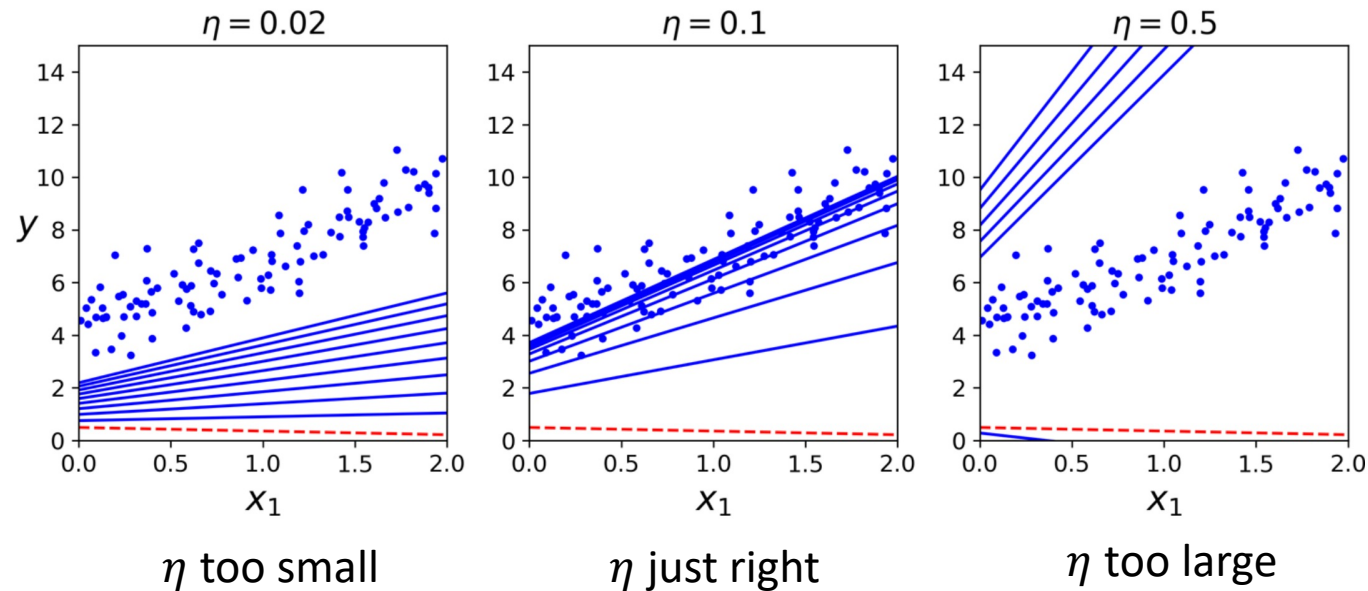
gradients as a *matrix product*  
computed over full training set for each step  
=> *batch gradient descent*

Calculating gradients for regression tasks

# Batch gradient descent

define size of step via *learning rate*

$$\boldsymbol{\theta}^{(\text{next step})} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} \text{MSE}(\boldsymbol{\theta})$$

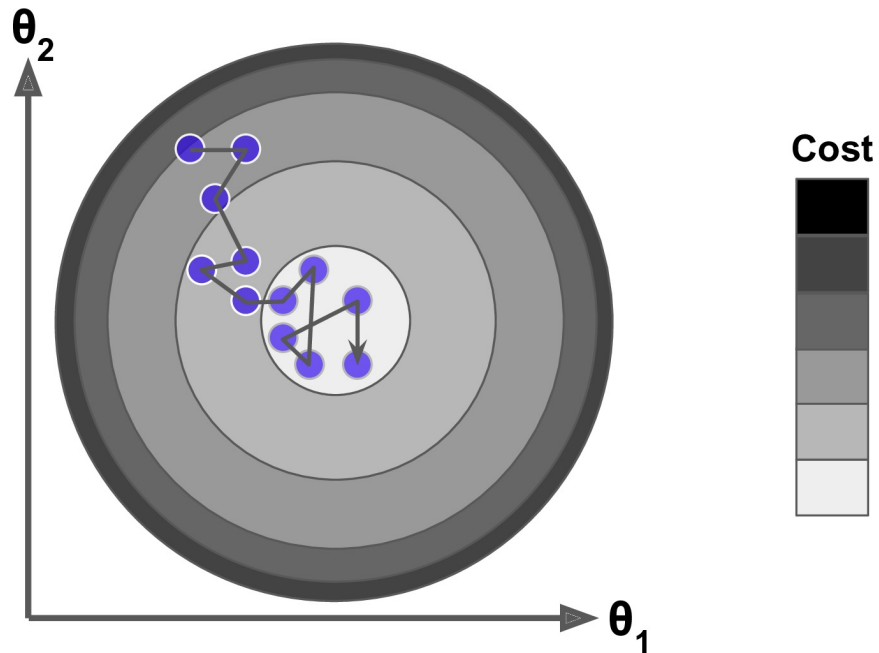


Hands-On Machine Learning with Scikit-Learn & TensorFlow, A. Geron

## The gradient descent step



# Stochastic gradient descent

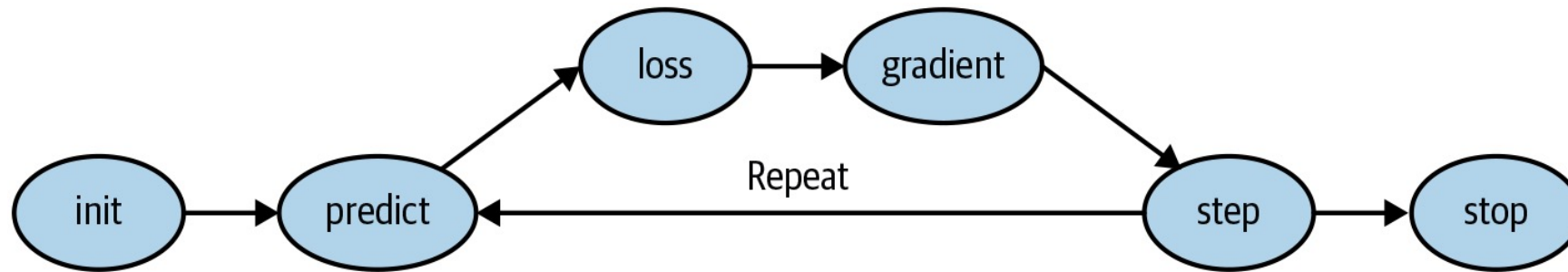


- ✓ much *faster* than batch gradient descent
- ✓ randomness allows algorithm to jump out of local minima
- ✗ randomness implies hard to settle in minimum  
(use a *learning rate schedule* to decay learning rate)

Hands-On Machine Learning with Scikit-Learn & TensorFlow, A. Geron

Sample & compute gradients on *single instance*

# Gradient descent summary



Deep Learning for Coders with fastai and PyTorch, J. Howard & S. Gugger