

# Training Models

Using Machine Learning Tools

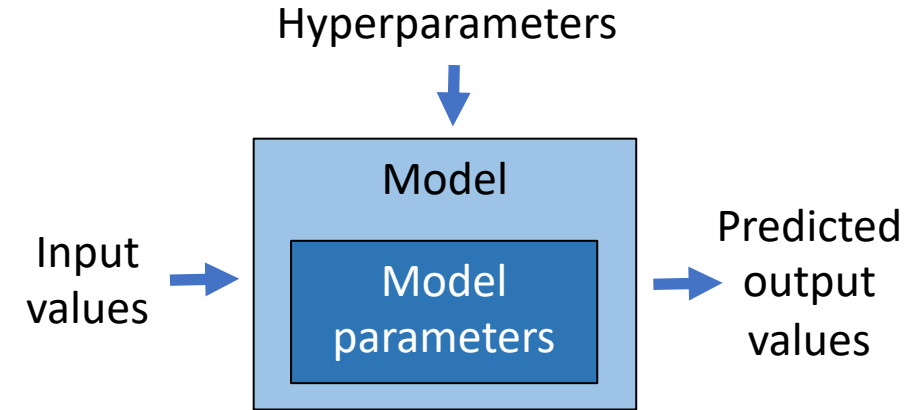
Geron Chapter 4

# Last Time ...

- Classification vs. regression
- Range of classifier approaches, white box vs. black box
- Several performance metrics
  - Accuracy, Confusion matrix
  - Precision, Recall
  - ROC, AUC
- Look at the data at all stages:
  - scatter plot
  - model parameters / decision boundaries
- Today we will look at training (i.e. optimisation of model parameters) in more detail

# Model

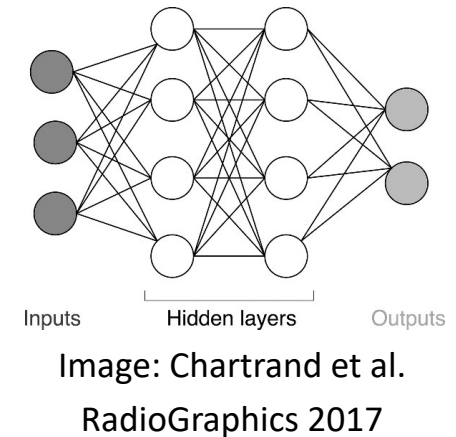
- A mapping of input values to predicted outcome values
- Flexible due to model parameters
- Constrained by fixed hyper-parameters
- Example: Linear model



$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n = \boldsymbol{\theta}^T \cdot \mathbf{x} = h_{\boldsymbol{\theta}}(\mathbf{x})$$

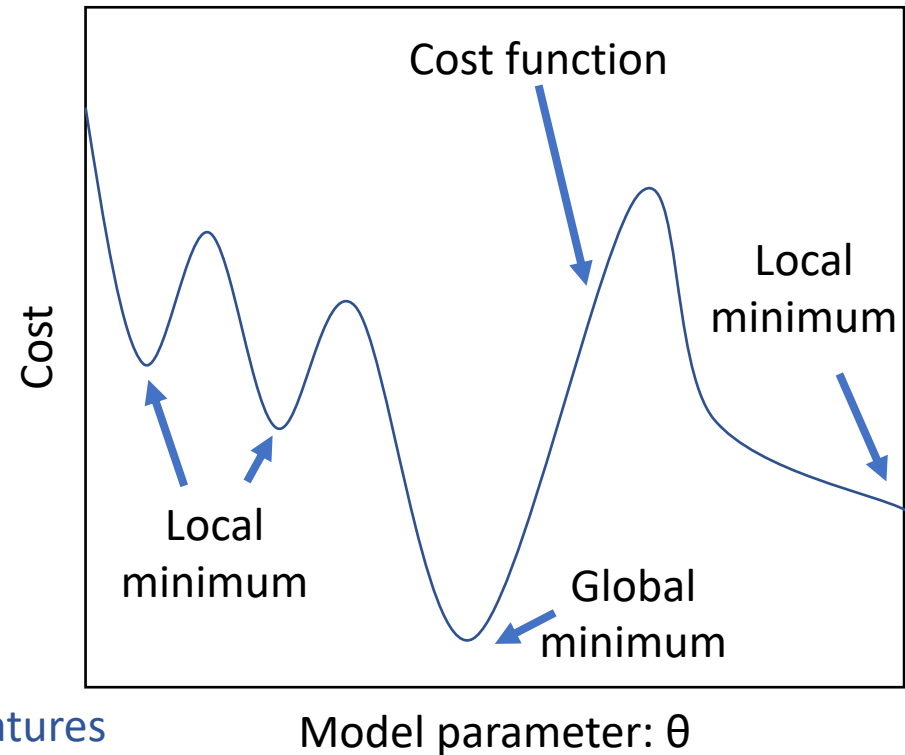
Diagram illustrating the linear model equation with annotations:

- $\hat{y}$ : Predicted output value
- $\theta_0, \theta_1, \theta_2, \dots, \theta_n$ : Model parameters
- $x_1, x_2, \dots, x_n$ : Input values for *one sample*
- $\boldsymbol{\theta}^T \cdot \mathbf{x}$ : The dot product of the parameter vector and the input vector
- $h_{\boldsymbol{\theta}}(\mathbf{x})$ : The hypothesis function



# Cost Function = Error = Loss Function

- Measures errors or differences between predicted and target values
- Want to minimise it for training data
- Global minimum: smallest value overall
- Local minimum: smallest value in some region
- Example: Mean square error (MSE)



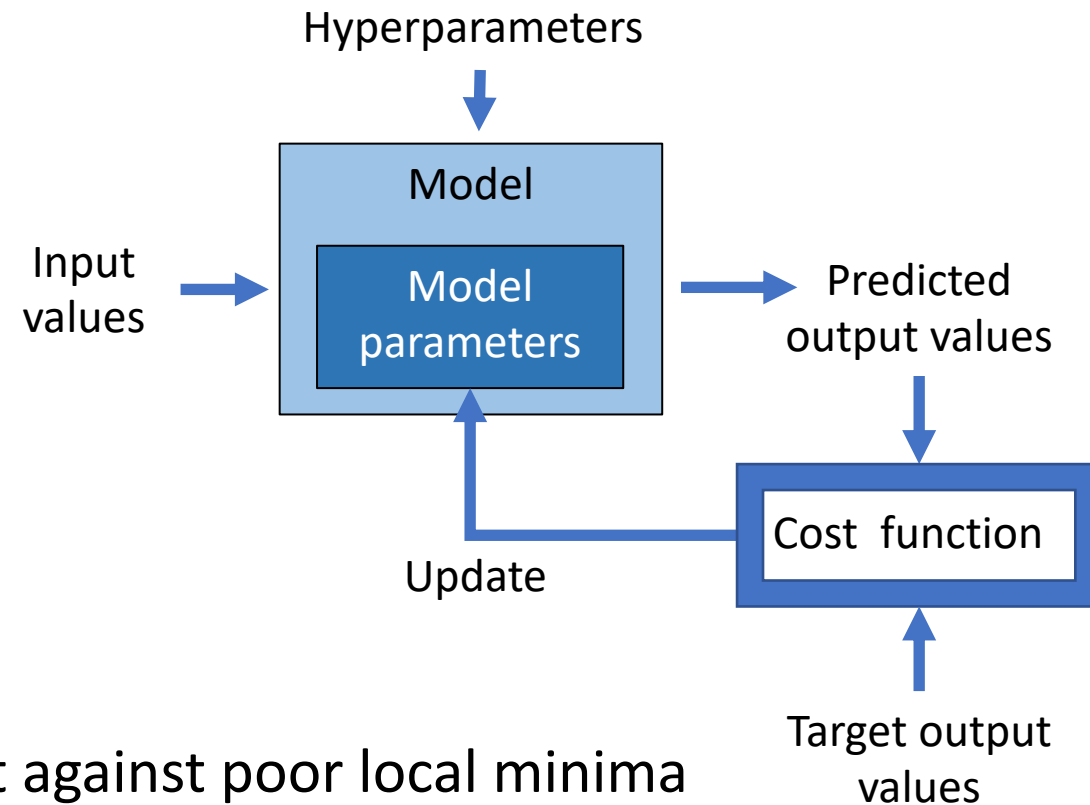
$$MSE(\mathbf{X}, h_{\theta}) = \frac{1}{M} \sum_{i=1}^M (\theta^T \mathbf{x}^{(i)} - y^{(i)})^2$$

Annotations for the equation:

- Input features for  $i^{th}$  sample (points to  $\mathbf{x}^{(i)}$ )
- Targets for  $i^{th}$  sample (points to  $y^{(i)}$ )
- Parameters (same for all samples) (points to  $\theta$ )

# Training = Fitting = Optimisation

- Minimise cost function by adjusting model parameters
- Start with initial guess of model parameters
- Iteratively change model parameters & evaluate cost function
- Ideal algorithm: fast, but robust against poor local minima



# Gradient Descent

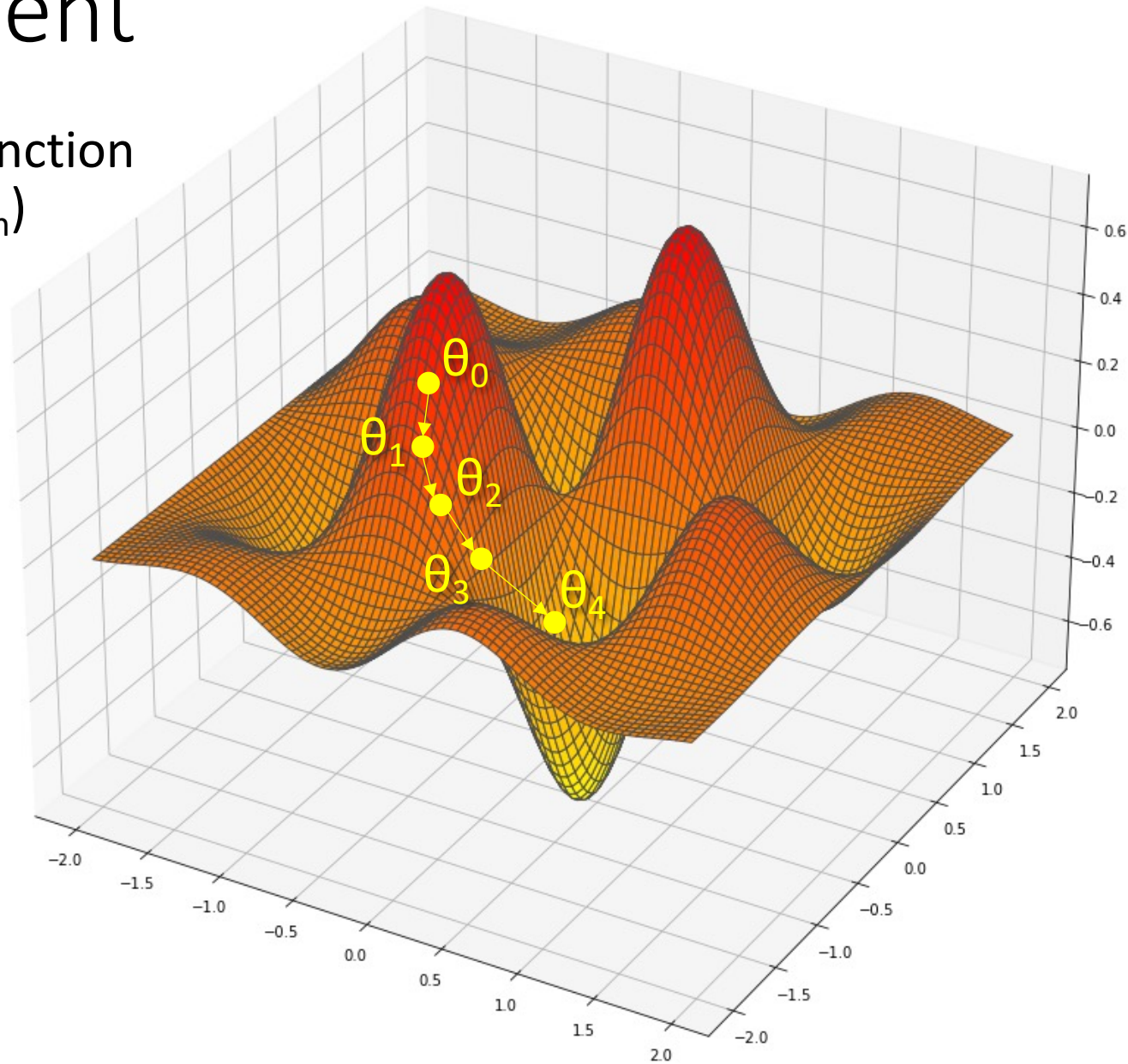
- Partial derivatives of cost function  
= Local gradient =  $\nabla_{\theta} \text{MSE}(\theta_n)$

- Iteratively step downhill  
(negative gradient)

$$\theta_{n+1} = \theta_n - \eta \nabla_{\theta} \text{MSE}(\theta_n)$$

Learning rate “eta”  $\eta$   
= sets size of step

This is an important  
hyper-parameter!

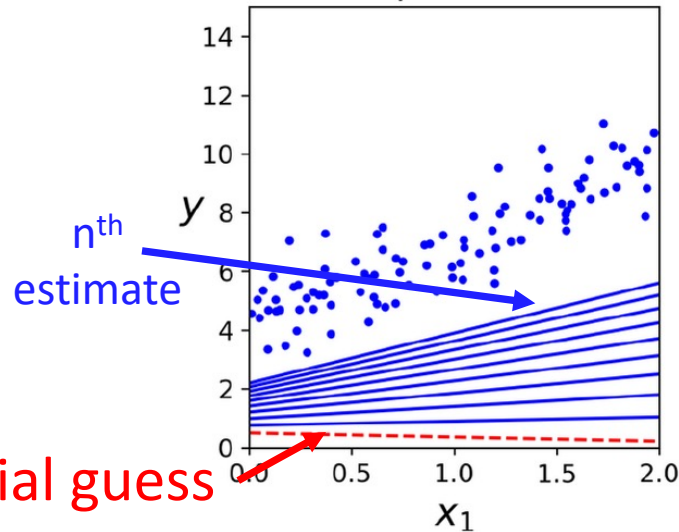


# Learning rate $\eta$

Example: trying to  
fit line to points

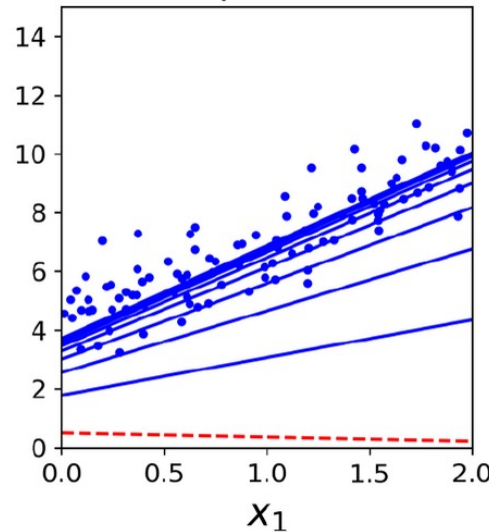
Too small

$\eta = 0.02$



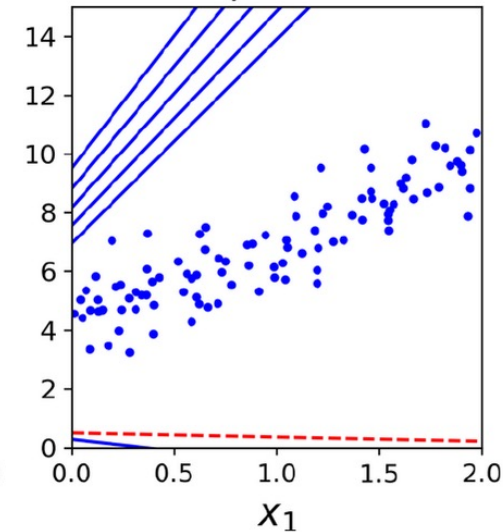
Good

$\eta = 0.1$

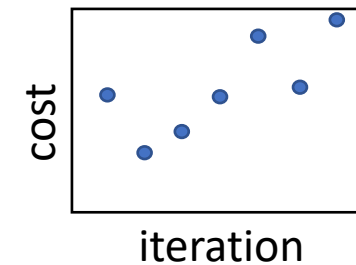
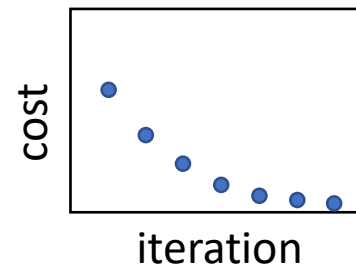
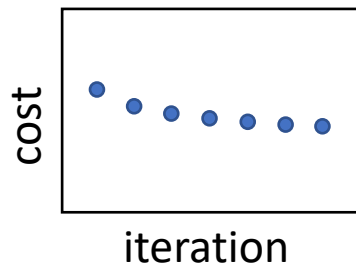


Too large

$\eta = 0.5$



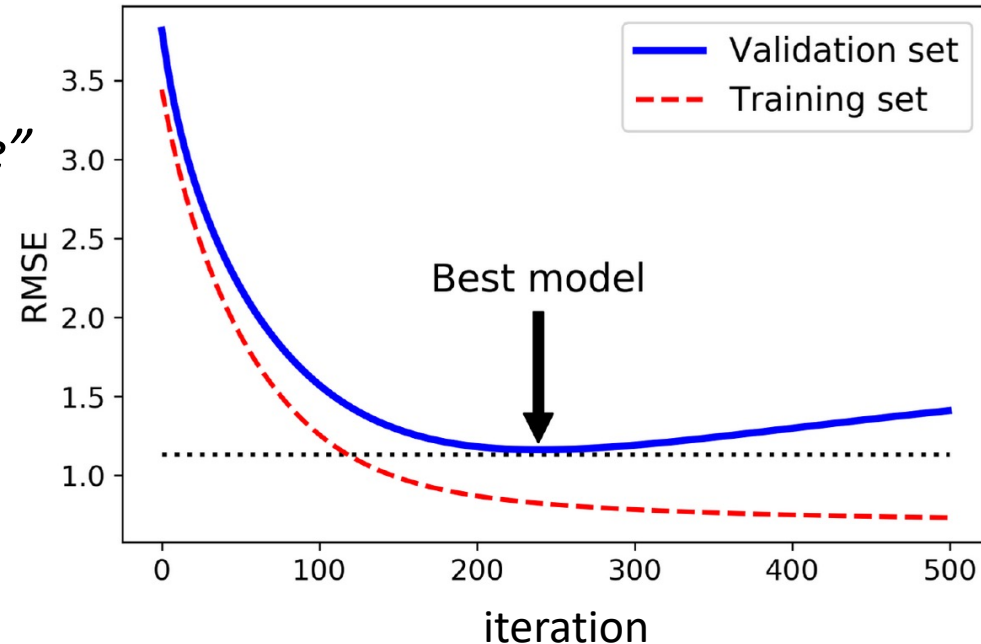
- Small  $\eta$  values take a long time to change, but go in the right direction
- Large  $\eta$  values overstep and can easily become unstable
- Can monitor the behaviour of the cost function values over iterations to spot these



# Stopping Criteria and Final Result

Stop if:

- No further improvement
  - e.g. 5 iterations in a row show “no change”
  - `early_stopping`: turn on/off
  - `n_iter_no_change`: number of iterations
  - `tol`: if cost difference between steps is less than  $\epsilon$  (tolerance) then treat it as *no change*
- Maximum number of iterations reached
  - `max_iter`: maximum number of iterations
- Final result: best model across the training process
  - this might *not* be the final model



Images: Geron, Hands On ML



# Stochastic Gradient Descent

- Pick one random sample
- Calculate the cost function gradient only from that sample

Better algorithm:

- Shuffle instances of the training set
- Use one instance after the other
- Adjust the learning rate  $\eta$
- Reshuffle and repeat

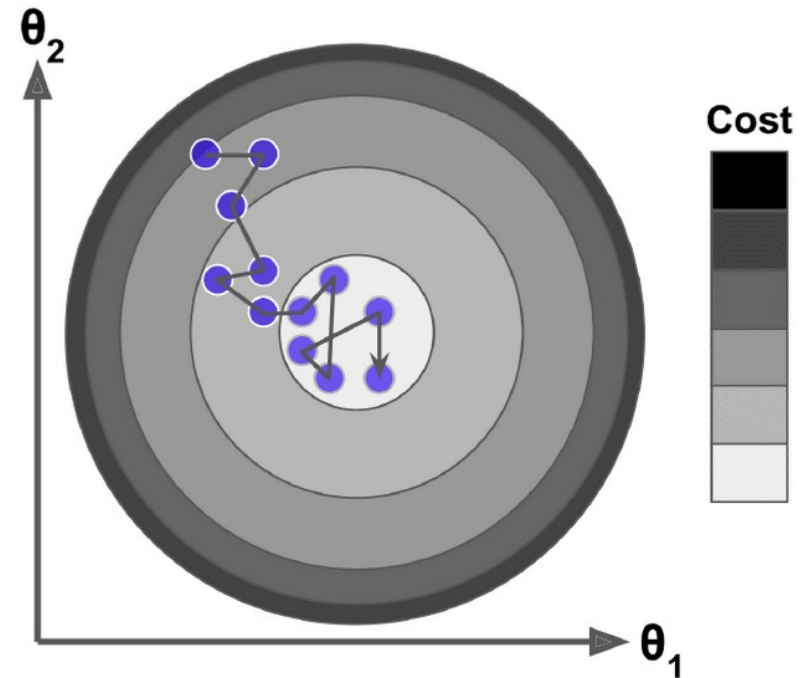
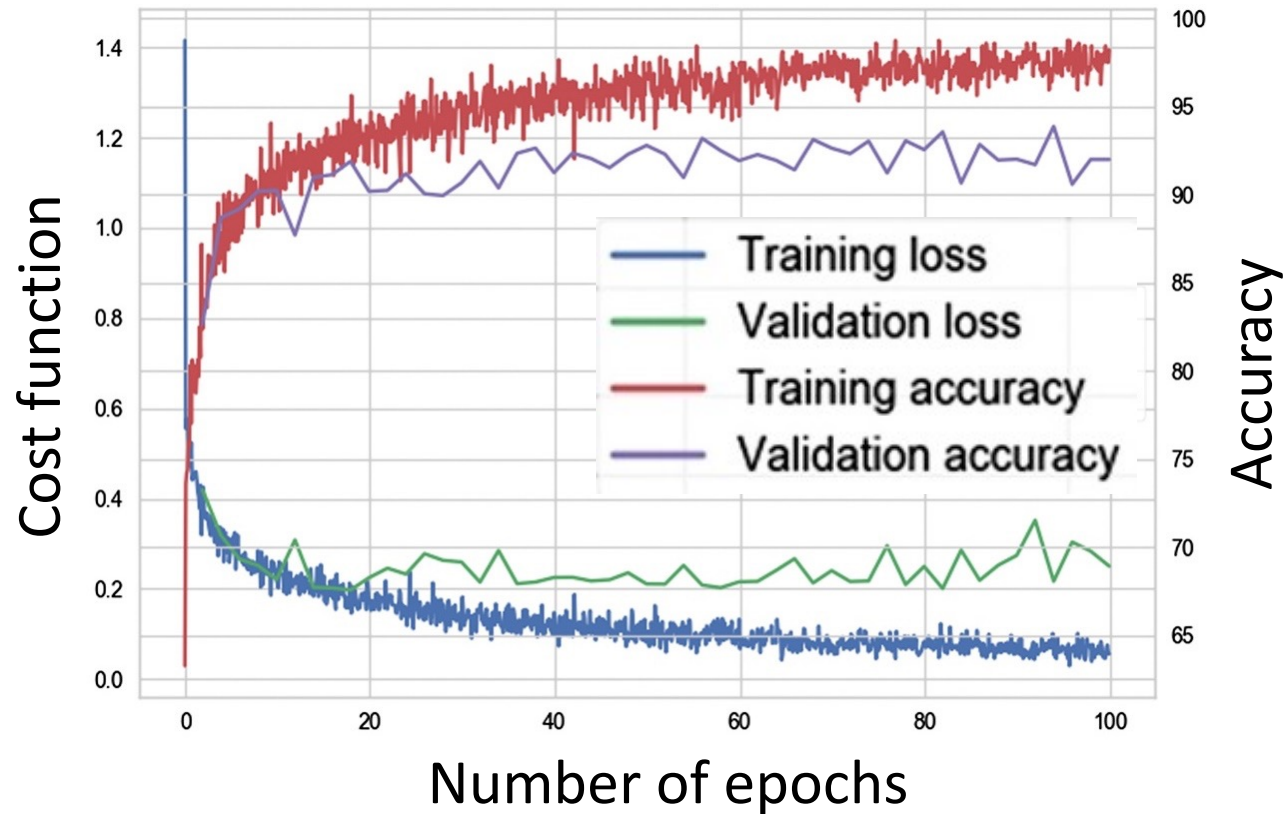


Image: Geron, Hands On ML

Pros: Fast, low memory, randomisation can help escape local minimum

Cons: Very noisy & no guarantee that minimum is reached

# Learning Curve = Training Curve



- Epoch = one pass through whole dataset  $\cong$  iteration
- Shows both training **and** validation performance
  - allows both underfitting and overfitting to be seen

Image: Chartrand et al.  
RadioGraphics 2017

# Sources of Generalisation Error

- **Variance**
  - *Irreducible error*
    - Due to randomness in the data itself
  - *Overfitting* leads to over-sensitivity to small variations in the data
    - Too many model parameters
- **Bias** or systematic error
  - *Sub-optimal model choice or hyperparameter choice*
    - Especially underfitting
  - *Representativeness of data*
    - Lack of data coverage (model extrapolations are usually bad)
    - Bias in the data (e.g. due to limitations or bias in sampling)
    - Imbalances in the data (e.g. due to nature of problem)
      - disease vs healthy ; suspicious vs normal transactions

# Regularisation

- Add a term to the cost function that tries to prevent overfitting
  - usually controlled by an adjustable weight  $\alpha$

$$\text{Cost} = \text{data\_term} + \alpha * \text{regularisation\_term}$$

- Purpose
  - Prevent overfitting by penalising large parameter values or lack of smoothness in outputs
  - Add a-priori knowledge (desired properties) to an underdetermined problem
- A form of multi-objective optimisation

# L2 (Ridge/Tikhonov) Regularisation

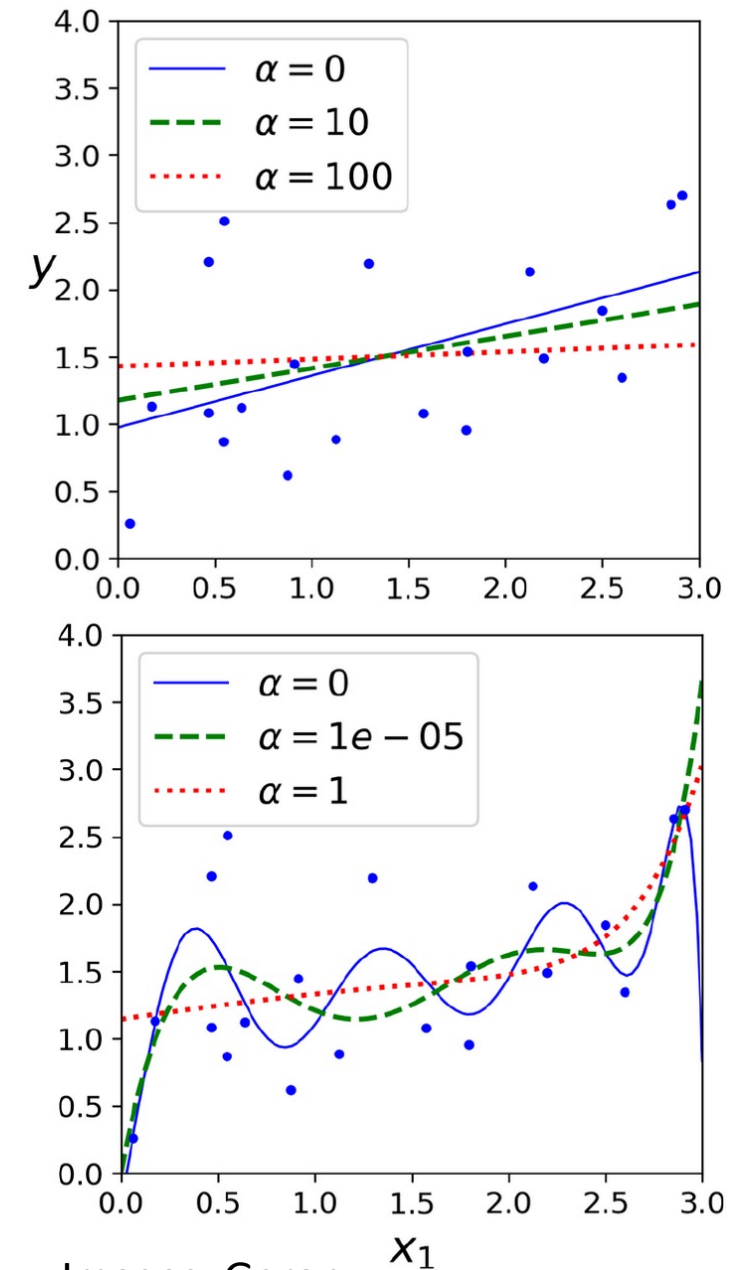
- Effect: Keep model parameters small

$$J(\boldsymbol{\theta}) = \text{MSE}(\boldsymbol{\theta}) + \alpha \frac{1}{2} \sum_{i=1}^n \theta_i^2$$

Diagram illustrating the components of the L2 regularisation cost function  $J(\boldsymbol{\theta})$ :

- Cost**: The overall function  $J(\boldsymbol{\theta})$ .
- Data term**: The Mean Squared Error  $\text{MSE}(\boldsymbol{\theta})$ .
- Regularisation term**: The penalty term  $\alpha \frac{1}{2} \sum_{i=1}^n \theta_i^2$ .
- Regularisation strength  $\alpha$** : The coefficient  $\alpha$  that scales the regularisation term.
- L2 Norm**: The term  $\frac{1}{2} \sum_{i=1}^n \theta_i^2$ , which is the squared L2 norm of the model parameters.
- Model parameters, in this case except  $\theta_0$** : The parameters  $\theta_i$  for  $i = 1, \dots, n$ .

- Scaling of data important for setting  $\alpha$
- Scikit learn: `penalty` parameter “l2”



Images: Geron,  
Hands On ML

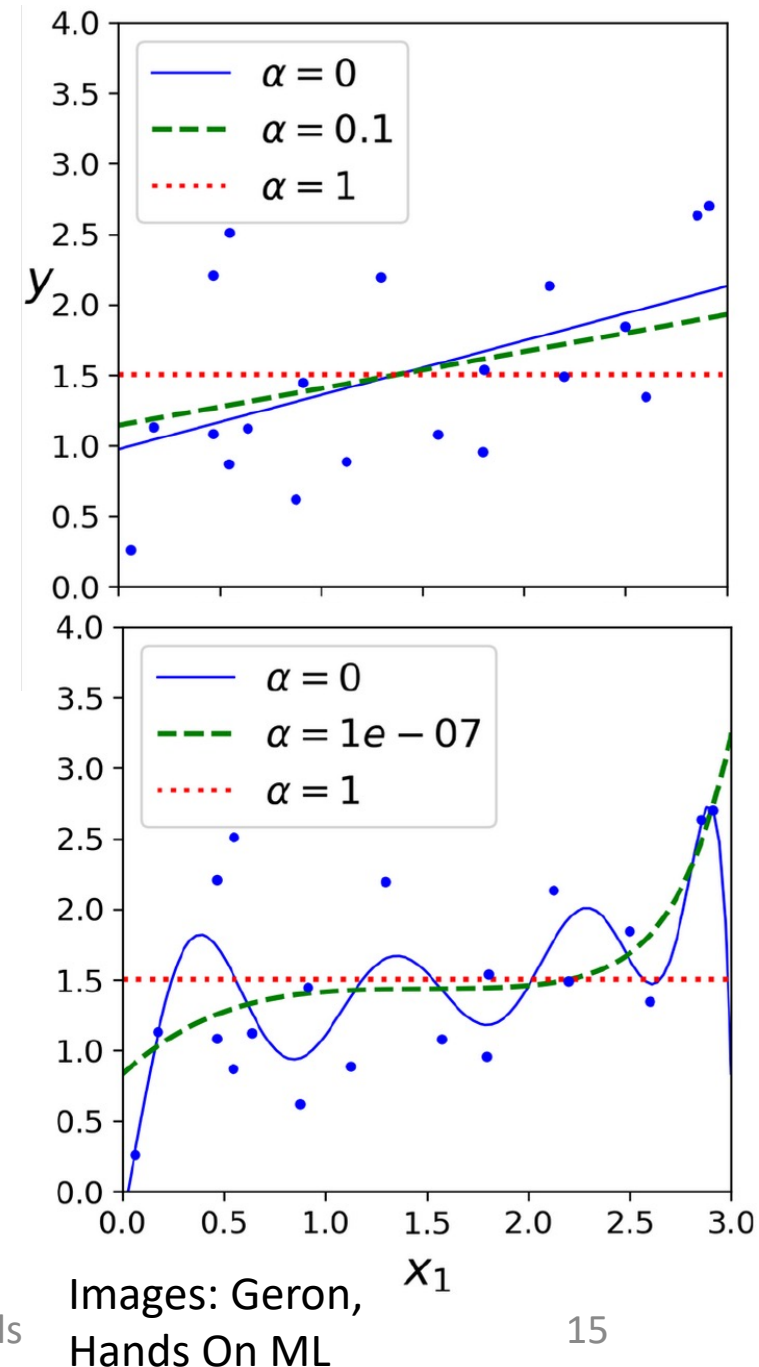
# L1 (Lasso) Regularisation

- LASSO = Least Absolute Shrinkage and Selection Operator
- Effect: Keep model parameters small
- Also tends to eliminate least important features (i.e. sets some  $\theta_i = 0$ )

$$J(\boldsymbol{\theta}) = \text{MSE}(\boldsymbol{\theta}) + \alpha \sum_{i=1}^n |\theta_i|$$

L1 Norm

- Scaling of data important for setting  $\alpha$
- Scikit learn: `penalty` parameter “l1”
- Not differentiable at 0, but most optimisers can cope with this



# “Elastic Net”

- Just a mixture of L2 and L1 regularisation

The diagram illustrates the Elastic Net cost function  $J(\boldsymbol{\theta})$  with several annotations. Above the equation, three blue brackets categorize the terms: 'Cost' for  $J(\boldsymbol{\theta})$ , 'Data term' for  $\text{MSE}(\boldsymbol{\theta})$ , and 'Regularisation term' for the entire regularisation part. The equation is  $J(\boldsymbol{\theta}) = \text{MSE}(\boldsymbol{\theta}) + r\alpha \sum_{i=1}^n |\theta_i| + \frac{1-r}{2}\alpha \sum_{i=1}^n \theta_i^2$ . Below the equation, two blue brackets identify the regularisation components: 'L1' for  $r\alpha \sum_{i=1}^n |\theta_i|$  and 'L2' for  $\frac{1-r}{2}\alpha \sum_{i=1}^n \theta_i^2$ . Two blue curved arrows originate from the text 'Mixture ratio r' at the bottom and point to the coefficients  $r$  and  $\frac{1-r}{2}$  in the equation.

$$J(\boldsymbol{\theta}) = \text{MSE}(\boldsymbol{\theta}) + r\alpha \sum_{i=1}^n |\theta_i| + \frac{1-r}{2}\alpha \sum_{i=1}^n \theta_i^2$$

Cost      Data term      Regularisation term

L1      L2

Mixture ratio  $r$

# Summary

- Training: Minimise cost function by adjusting model parameters
- Regularisation: Penalise large parameters via an extra term in the cost function
- It can be helpful to understand how these work to get some intuition for good settings to use and how to diagnose problems
- Critically evaluate the implementation at hand and the defaults
  - Example: SGD classifier by default uses L2 (Ridge) regularisation