

# Logistic Regression and Support Vector Machines

## Using Machine Learning Tools

Geron Chapter 4 and 5

# Last Time ...

Training: Minimise cost function by adjusting model parameters

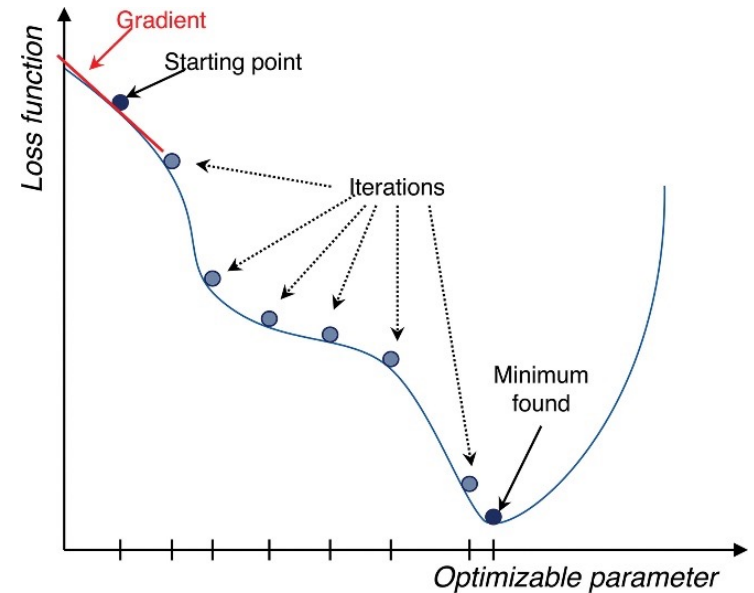
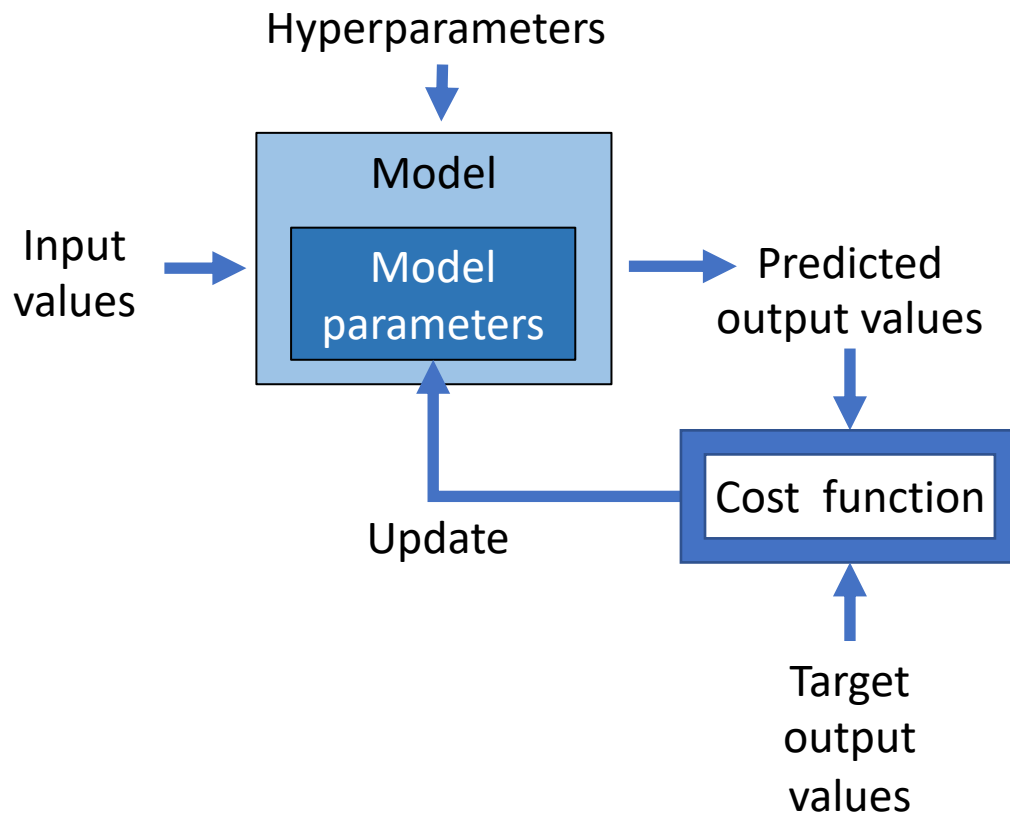


Image: Chartrand et al. RadioGraphics 2017

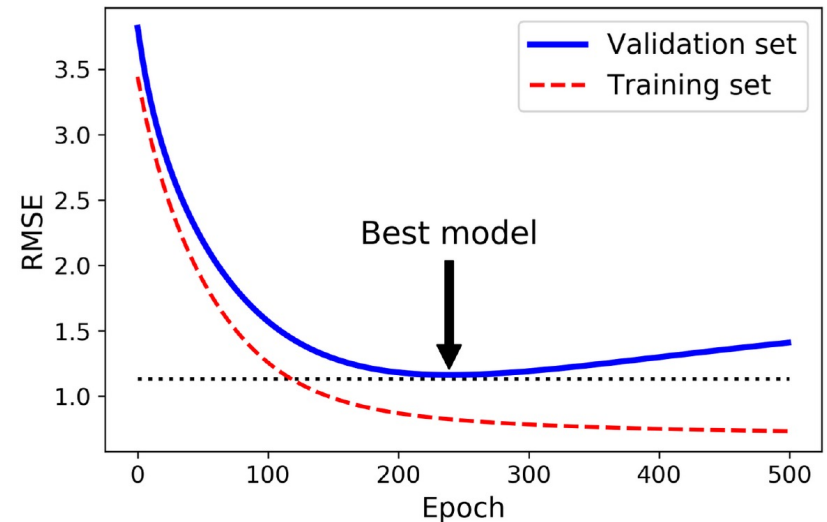


Image: Geron, Hands On ML

# Logistic Regression = Logit Regression

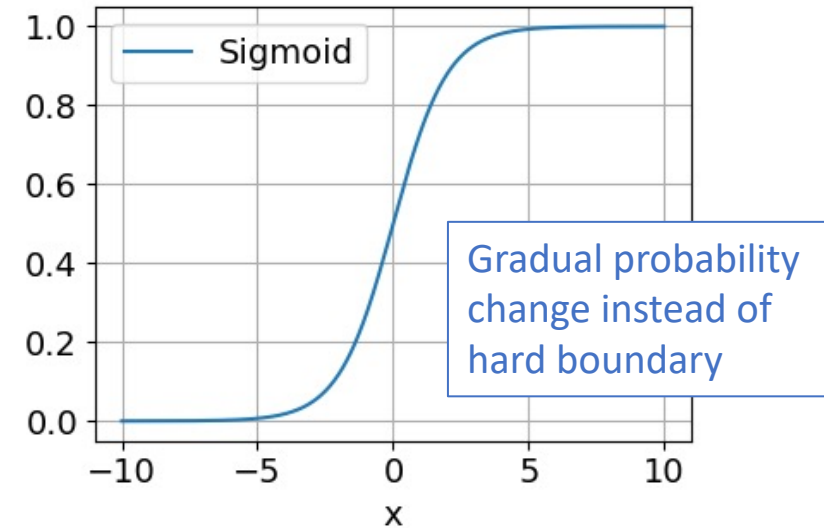
Linear model

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n = \mathbf{x}^T \cdot \boldsymbol{\theta}$$

+ Logistic function (sigmoid)

$$\sigma(t) = \frac{1}{1 + \exp(-t)}$$

= Logistic model (regression)



Logistic fct. between 0 and 1

Negative inputs  $\sigma(t) < 0.5$

Positive inputs  $\sigma(t) > 0.5$

# Logistic Regression = Logit Regression

Linear model

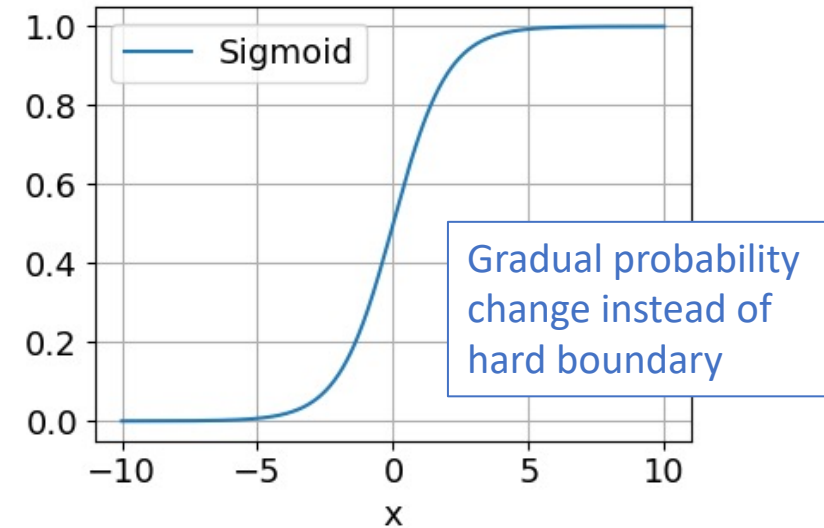
$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n = \mathbf{x}^T \cdot \boldsymbol{\theta}$$

+ Logistic function (sigmoid)

$$\sigma(t) = \frac{1}{1 + \exp(-t)}$$

= Logistic model (regression)

$$\hat{p} = h_{\boldsymbol{\theta}}(\mathbf{x}) = \sigma(\mathbf{x}^T \boldsymbol{\theta})$$



Logistic fct. between 0 and 1

Negative inputs  $\sigma(t) < 0.5$

Positive inputs  $\sigma(t) > 0.5$

Logit = inverse of logistic fct.

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right)$$

# Logistic Regression = Logit Regression

Linear model

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n = \mathbf{x}^T \cdot \boldsymbol{\theta}$$

+ Logistic function (sigmoid)

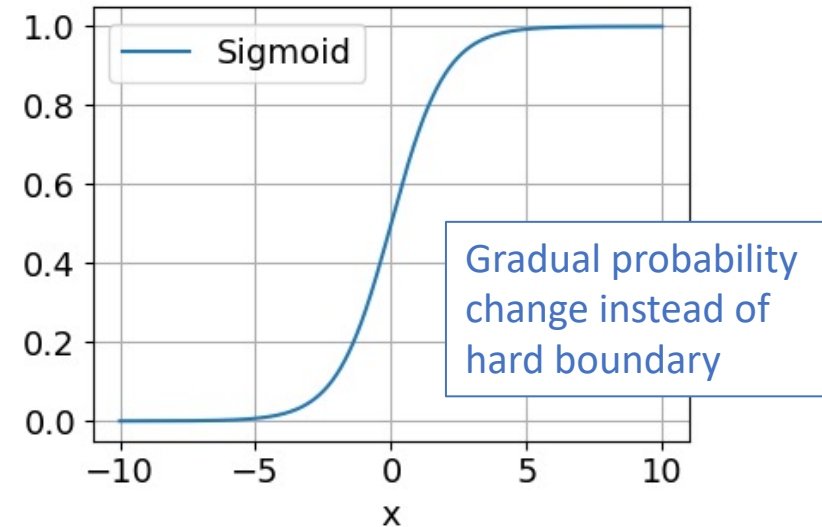
$$\sigma(t) = \frac{1}{1 + \exp(-t)}$$

= Logistic model (regression)

$$\hat{p} = h_{\boldsymbol{\theta}}(\mathbf{x}) = \sigma(\mathbf{x}^T \boldsymbol{\theta})$$

+ Threshold = binary classifier

$$\hat{y} = \begin{cases} 0 & \text{if } \hat{p} < 0.5 \\ 1 & \text{if } \hat{p} \geq 0.5 \end{cases}$$



Logistic fct. between 0 and 1

Negative inputs  $\sigma(t) < 0.5$

Positive inputs  $\sigma(t) > 0.5$

Logit = inverse of logistic fct.

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right)$$

# Training a Logistic Model

- We are maximising probability of observing our targets  $y$ , given data  $x$
- Probability of observing binary data  $y$  if the real probability is  $p$ :

$$P(y) = \begin{cases} p & \text{if } y = 1 \\ 1 - p & \text{if } y = 0 \end{cases} \quad \text{or} \quad P(y) = p^y (1 - p)^{1-y}$$

# Training a Logistic Model

- We are maximising probability of observing our targets  $y$ , given data  $x$
- Probability of observing binary data  $y$  if the real probability is  $p$ :

$$P(y) = \begin{cases} p & \text{if } y = 1 \\ 1 - p & \text{if } y = 0 \end{cases} \quad \text{or} \quad P(y) = p^y (1 - p)^{1-y}$$

- For a set of observations:  $P(\mathbf{y}) = \prod_i \hat{p}_i^{y_i} (1 - \hat{p}_i)^{1-y_i}$

$$\text{or } \log(P(\mathbf{y})) = \sum_i y_i \log(\hat{p}_i) + (1 - y_i) \log(1 - \hat{p}_i)$$

# Training a Logistic Model

- We are maximising probability of observing our targets  $y$ , given data  $x$
- Probability of observing binary data  $y$  if the real probability is  $p$ :

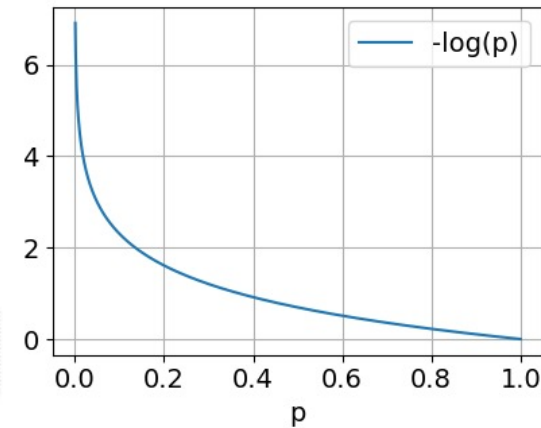
$$P(y) = \begin{cases} p & \text{if } y = 1 \\ 1 - p & \text{if } y = 0 \end{cases} \quad \text{or} \quad P(y) = p^y (1 - p)^{1-y}$$

- For a set of observations:  $P(\mathbf{y}) = \prod_i \hat{p}_i^{y_i} (1 - \hat{p}_i)^{1-y_i}$

$$\text{or } \log(P(\mathbf{y})) = \sum_i y_i \log(\hat{p}_i) + (1 - y_i) \log(1 - \hat{p}_i)$$

- Where the probability is the estimated probability output by the regression  $\hat{p} = h_{\theta}(\mathbf{x}) = \sigma(\mathbf{x}^T \theta)$

- Minimise this *negative* log probability: called **log loss** or **binary cross-entropy**  $J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[ y^{(i)} \log(\hat{p}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{p}^{(i)}) \right]$



- Cost function is convex  $\rightarrow$  good for optimisation, as they have no local minima



# Decision Boundaries

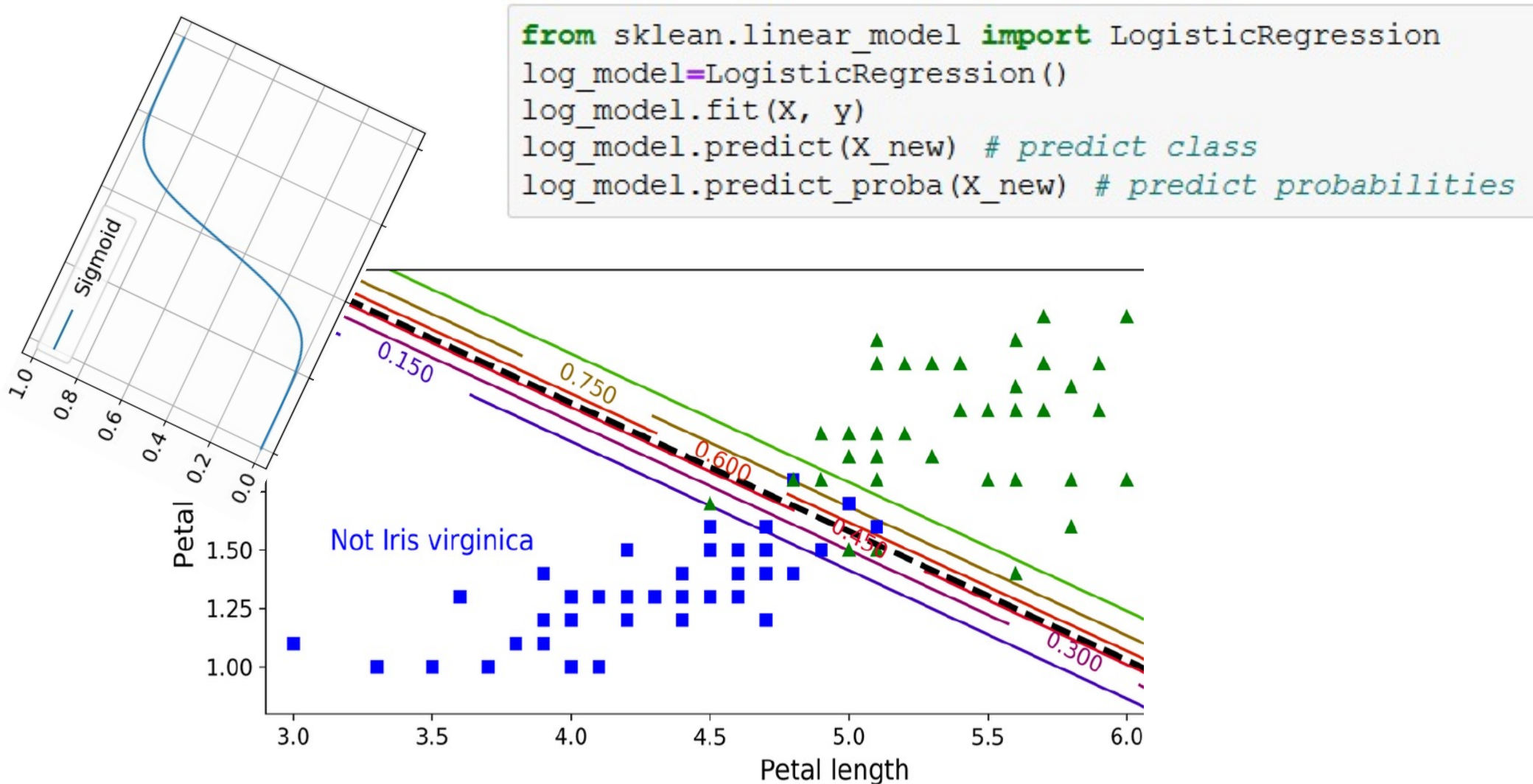


Image: Geron, Hands On ML

# Softmax Regression

- Extend logistic regression to multiple classes
- Fit a linear model for each class  $k$   $s_k(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\theta}^{(k)}$  ← Called “logits”
- Map the class scores into a probability distribution using:

- *Softmax* function:

$$\hat{p}_k = \sigma(\mathbf{s}(\mathbf{x}))_k = \frac{\exp(s_k(\mathbf{x}))}{\sum_{j=1}^K \exp(s_j(\mathbf{x}))}$$

Outputs in range (0,1)

Monotonic

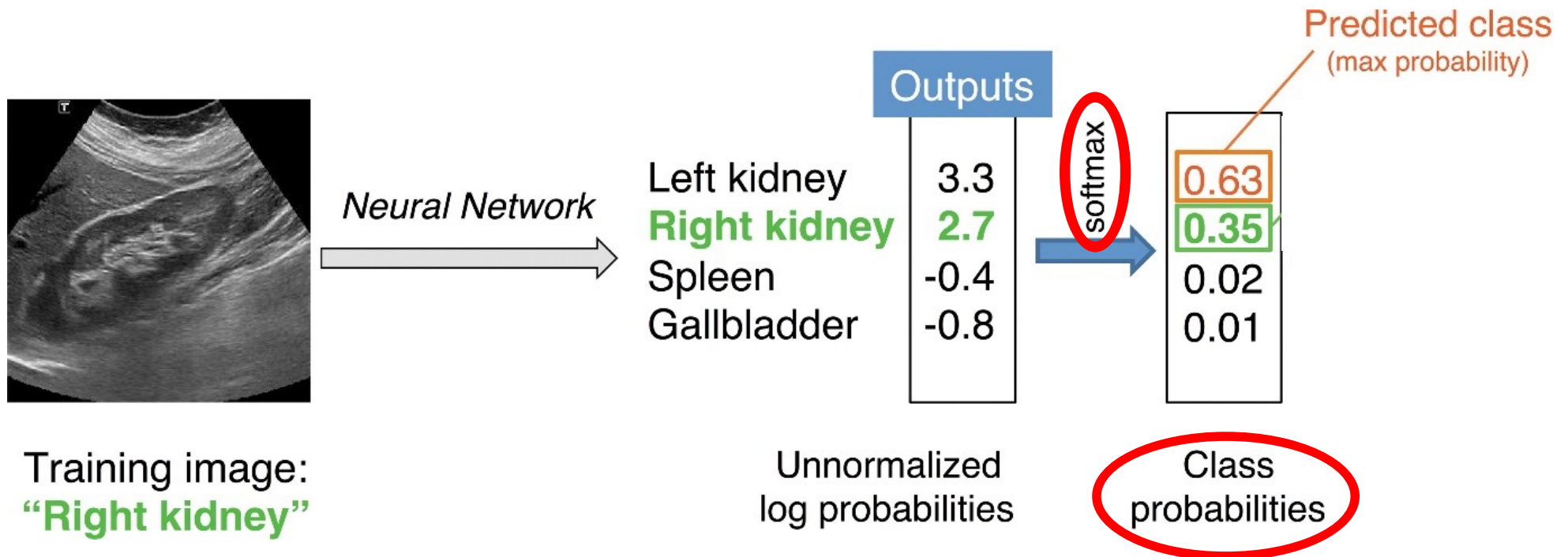
$$\sum_k \hat{p}_k = 1$$

- Cost function: *Cross entropy*

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

Differentiable, so use  
Gradient Descent

# Softmax Normalisation in Neural Networks



$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \text{ for } i = 1, \dots, K \text{ and } \mathbf{z} = (z_1, \dots, z_K) \in \mathbb{R}^K$$

(Chartrand et al.  
RadioGraphics 2017)

# Part 2

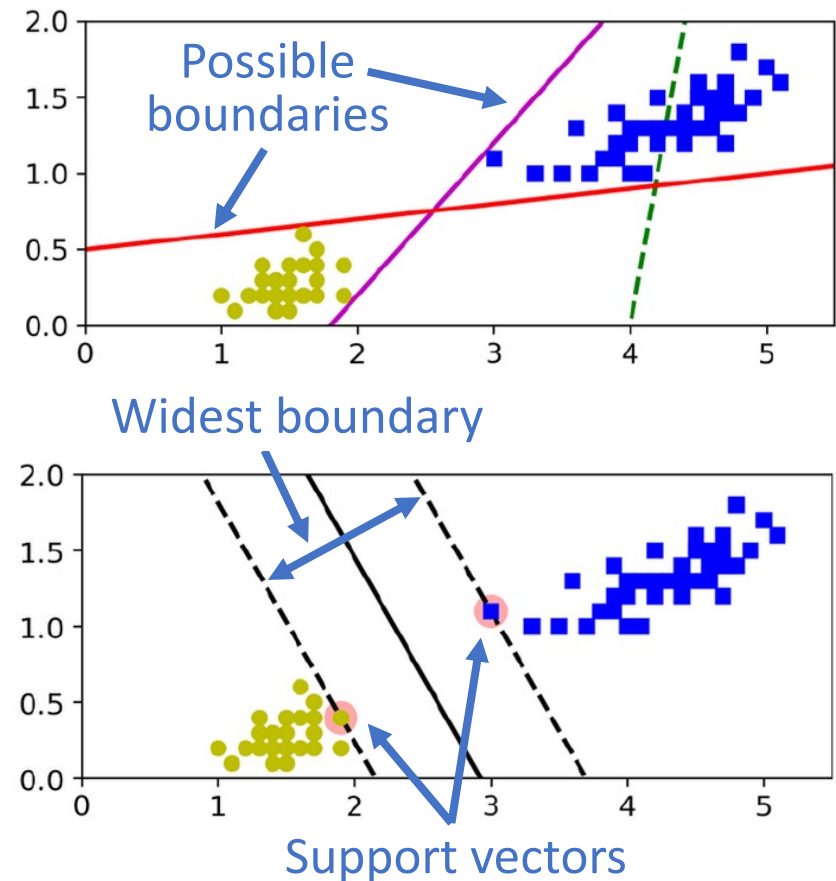
## Support Vector Machines

# Support Vector Machine (SVM)

- Key idea: fit the **widest possible boundary** between two classes while limiting boundary violations
- Turns out that the boundary is defined by **a small set of samples** at the boundary edge → **support vectors**
- Good for complex boundaries & small to medium datasets (becomes quite slow for large sets)
- Can be linear or non-linear (boundary shape)
- Internally, uses a smart mathematical approach, called the “kernel trick”

# Linear SVM - Hard Margin

- Fit the widest possible boundary between two classes **with zero** boundary violations
- As for any method based on *distances* an SVM is **sensitive to the scale** of inputs, especially relative scaling
  - Make sure pre-processing does scaling!



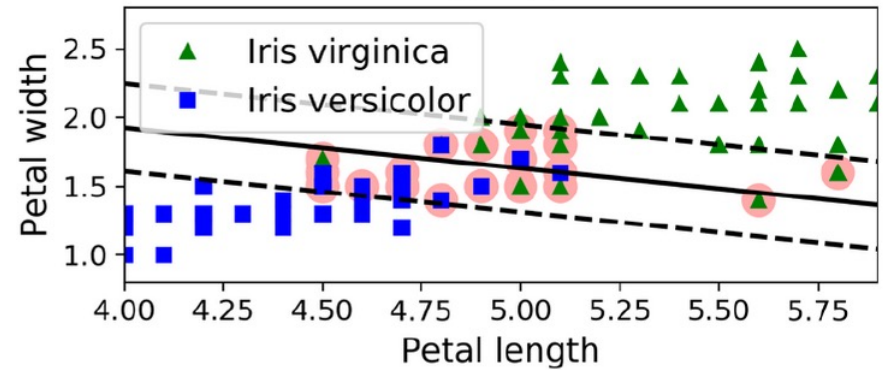
Images: Geron, Hands On ML

# Linear SVM – Soft Margin

- Fit the *widest possible boundary* between two classes **while limiting boundary violations**
  - Trade-off between these
- Parameter **C** controls how important boundary violations are in the loss function
  - C near zero → Hard Margin where boundary width is most important
  - Large C → boundary width less important compared to boundary violations

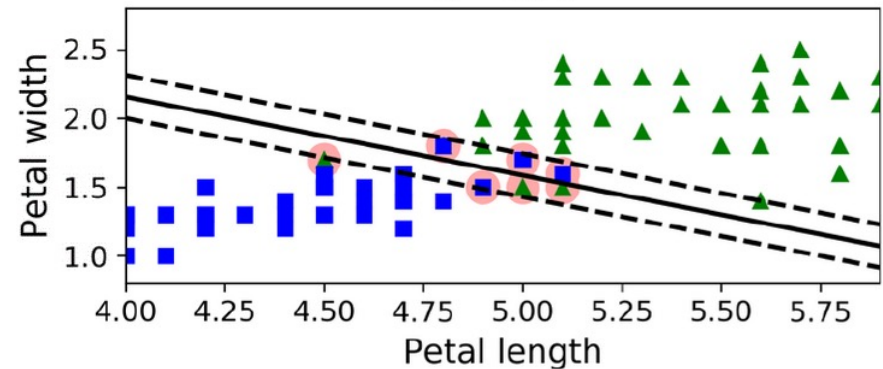
Wider boundary, more violations

$C = 1$



Less violations, narrower boundary

$C = 100$

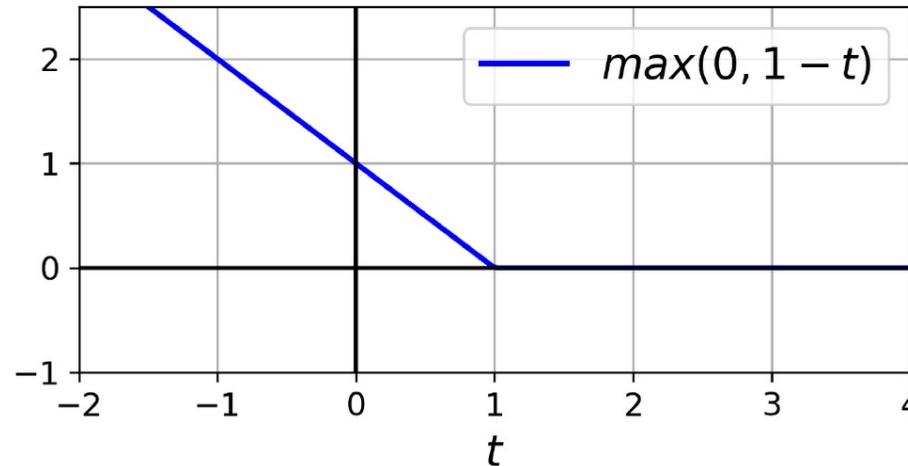


Images: Geron, Hands On ML

# Hinge loss function & SVM

- Can write the SVM loss using the Hinge loss function:

$$\max(0, 1 - t)$$



- You will see this as an option in other machine learning methods



# Non-Linear Boundaries

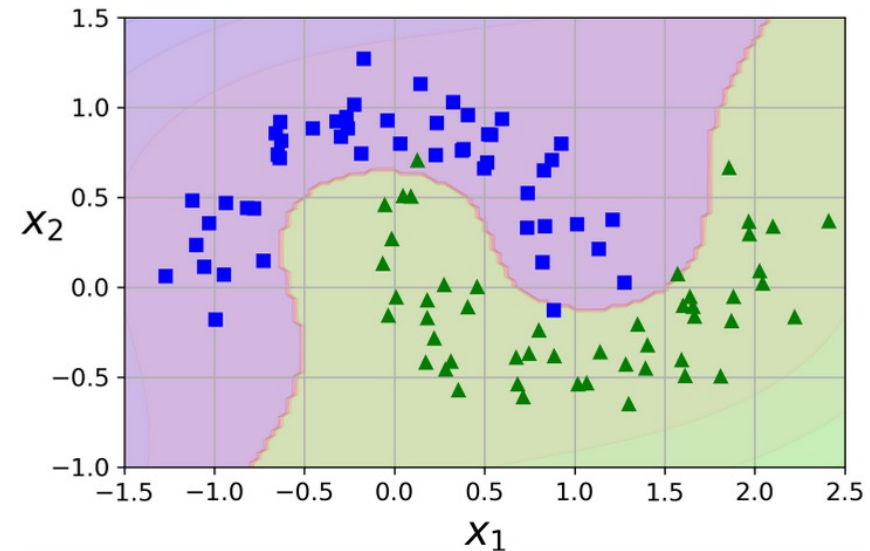
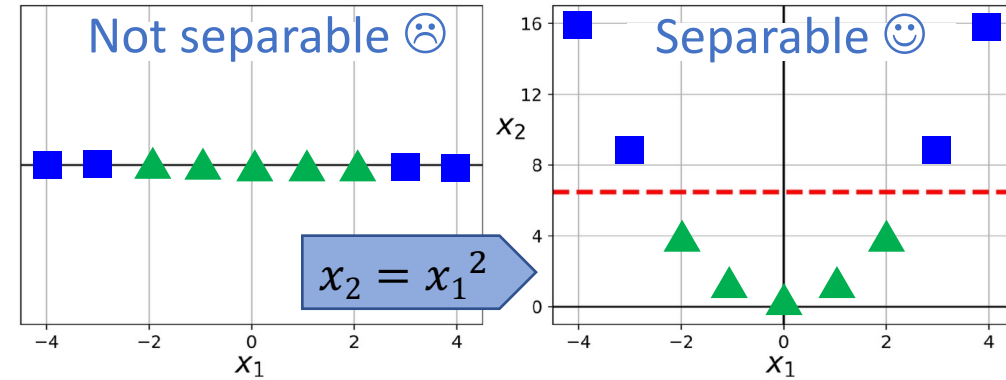
- Problem: Many decision boundaries are not linear
- Idea: Add non-linear functions of the input data as extra features and use a linear boundary in this higher dimensional space
  - Longer computation time
- e.g. Gaussian radial basis function (RBF)

Distance between samples

$$\phi_{\gamma}(\mathbf{x}, \ell) = \exp \left( -\gamma \|\mathbf{x} - \ell\|^2 \right)$$

New non-linear feature

Hyperparameter Gamma



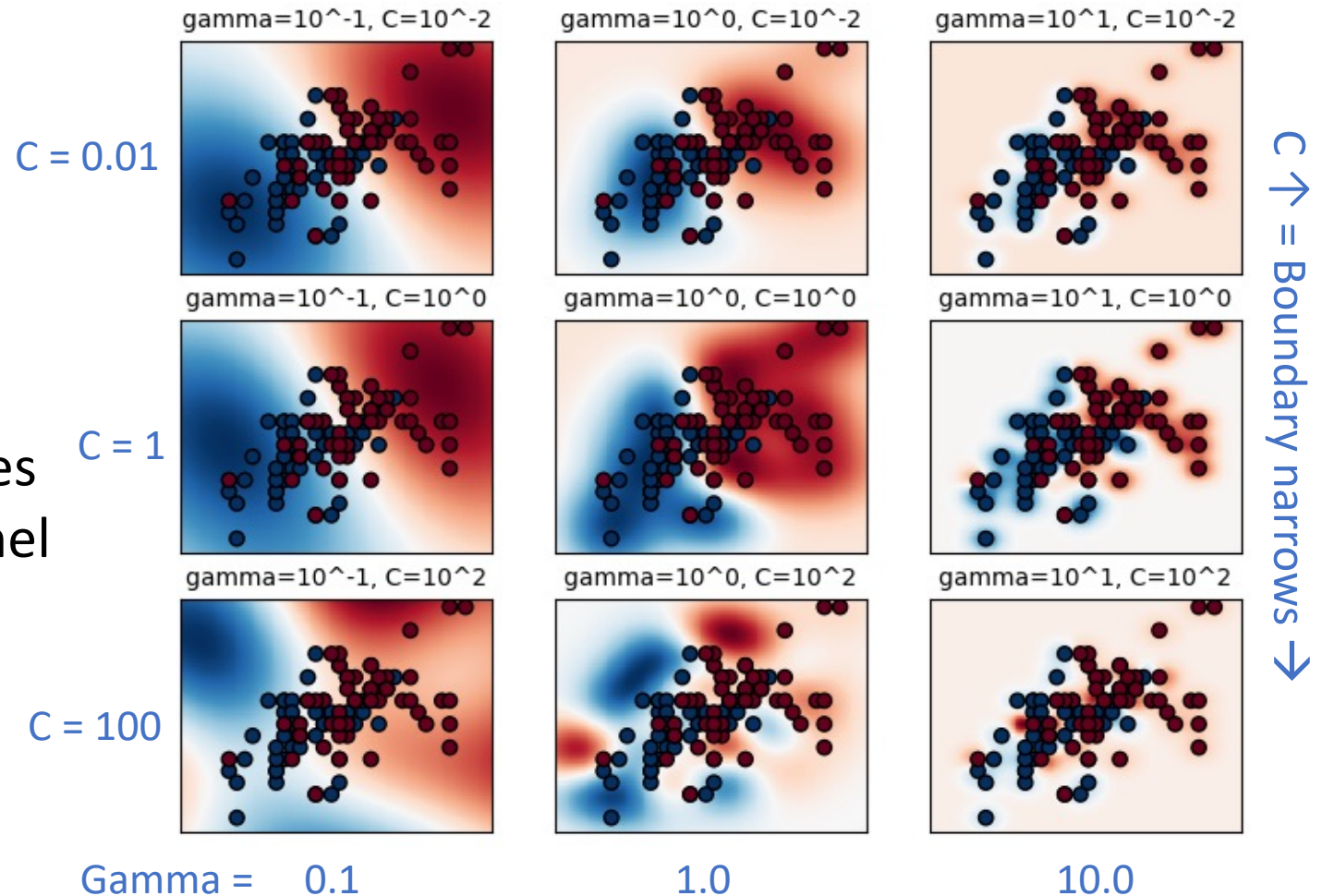
Images: Geron, Hands On ML

# SVM Hyperparameters: C & Gamma

Gamma  $\uparrow$  = More local (less smooth)  $\rightarrow$

An example:

- 2 classes
- 50 samples each
- 2 (original) features
- Gaussian RBF kernel



# Kernel “trick”

- Problem: we have introduced many extra features
  - More powerful classifier, but SLOW!
- Turns out that SVM optimisation only needs to know the *dot product* between each sample’s new high dimensional feature set (i.e. vector)
- Store these dot products in a matrix:  $N_{\text{samp}} \times N_{\text{samp}}$
- **Kernel** = a function  $K(\mathbf{a}, \mathbf{b})$  that specifies the dot product
- Kernel trick: use kernel functions to get dot products without ever actually calculating the extra features
- Common kernels (best practice to start simple)

Linear:	$K(\mathbf{a}, \mathbf{b}) = \mathbf{a}^\top \mathbf{b}$	Linear classifier
Polynomial:	$K(\mathbf{a}, \mathbf{b}) = (\gamma \mathbf{a}^\top \mathbf{b} + r)^d$	Use polynomial feature combinations without having to compute them
Gaussian RBF:	$K(\mathbf{a}, \mathbf{b}) = \exp(-\gamma \ \mathbf{a} - \mathbf{b}\ ^2)$	Use similarity to support vectors with Gaussian drop off
Sigmoid:	$K(\mathbf{a}, \mathbf{b}) = \tanh(\gamma \mathbf{a}^\top \mathbf{b} + r)$	

# Summary

- Logistic regression
  - Probability-based method
  - Uses sigmoid (logistic) function to get outputs in  $[0,1]$  range
  - *Sigmoid & Softmax* functions used in many deep learning models
- Support Vector Machines
  - Fit the *widest possible boundary* between two classes while limiting boundary violations
  - Create extra non-linear features → high dimensional space where they can be separated with a linear boundary in this space
  - Best to try simple kernels first, compare different kernels
  - Hyperparameters also need to be optimised
  - Data needs to be scaled uniformly (internally it is distance-based)