

# Logistic Regression

like linear regression but used for classification  
 ↙ binary ↘ multi → etc.

## ① Estimating probability (Sigmoid)

1. it compute weighted sum of features + bias (like linear Regs.)

2. then compute logistic of this result ( $\hat{p}$ )

Note:

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

= inverse of logistic

= also called "log-odds"

= log ratio positive prediction  
Prob f neg

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

↑  
sigmoid

↘ weighted sum + Bias

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

- it has range  $[0, 1]$

- Fun fact: called sigmoid as "shape of S"

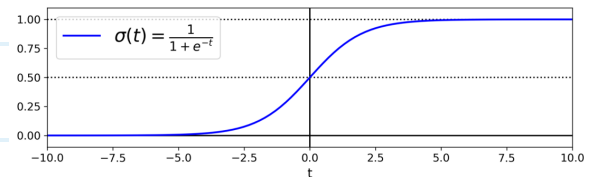


Figure 4-21. Logistic function

3. Make predication ( $\hat{y}$ )

$$\hat{y} = \begin{cases} 1 \\ 0 \end{cases}$$

$\hat{p} \geq 0.5 \rightarrow$  positive sample

$\hat{p} < 0.5 \rightarrow$  negative sample

meaning  $\hat{y} = \begin{cases} 1 & x^T \theta \geq 0 \\ 0 & x^T \theta < 0 \end{cases}$

#### 4. Calculate Loss Function

$$c(\vec{\theta}) = \begin{cases} -\log(\hat{p}) & \text{if } y=1 \\ -\log(1-\hat{p}) & \text{if } y=0 \end{cases}$$

to make sense @  $y=1 \rightarrow$  we want  $\hat{p}=1 \rightarrow -\log(\hat{p}) = -\log(1)=0$   
 $\rightarrow$  But if it approaches zero  $\hat{p} \approx 0$  it goes to " $\infty$ "

the reverse @  $y=0 \rightarrow$  we want  $\hat{p}=0 \rightarrow -\log(1-\hat{p}) = 0$

But if  $\hat{p}=1$  cost  $\approx \infty$

#### 5. Calculate Cost Function

$J(\vec{\theta}) = \text{avg of Loss}$

$$\begin{aligned} &= \frac{1}{m} \sum_{i=1}^m y^{(i)} \log(\hat{p}^{(i)}) + (1-y^{(i)}) \log(1-\hat{p}^{(i)}) \\ &= \frac{1}{m} \sum_{i=1}^m y^{(i)} \log(\hat{p}^{(i)}) + \bar{y}^{(i)} \log(\bar{p}^{(i)}) \end{aligned}$$

$\sigma(x^T \theta)$

Forward  
end  
here

#### 6. Calculate partial Derivative of Cost

$$\sigma'(x) = \frac{d}{dx} \sigma(x) = \sigma(x) (1 - \sigma(x))$$

$$\frac{d}{dx} \log(F(x)) = \frac{F'(x)}{F(x)}$$

$$\frac{d}{dx} F(g(x)) = F'(g(x)) g'(x) \Rightarrow \text{Chain Rule of derivatives}$$

$$\begin{aligned} \text{so } \frac{d}{d\theta_i} \log(\sigma(x^T \theta)) &= \frac{\sigma'(x^T \theta)}{\sigma(x^T \theta)} = \frac{\sigma(x^T \theta)(1 - \sigma(x^T \theta)) \cdot x_i}{\sigma(x^T \theta)} \\ &= (1 - \hat{p}) x_i \end{aligned}$$

$$\frac{d}{d\theta_i} \log(1 - \sigma(x^T \theta)) = \frac{-(\hat{p})(1 - \hat{p}) x_i}{1 - \hat{p}} = -\hat{p} x_i$$

$$\text{so } J(\vec{\theta}) = \frac{-1}{m} \sum_{i=1}^m y^{(i)} \log(\hat{p}) + (1 - y^{(i)}) \log(1 - \hat{p})$$

$$\text{so } \frac{d}{d\theta_j} J(\vec{\theta}) = \frac{-1}{m} \sum_{i=1}^m y^{(i)} (1 - \hat{p}) x_j^{(i)} + (1 - y^{(i)}) (-\hat{p} x_j^{(i)})$$

$$= \frac{-1}{m} \sum_{i=1}^m (y^{(i)} - \hat{p}) x_j^{(i)}$$

$$= \frac{1}{m} \sum_{i=1}^m (\sigma(\theta^T x^{(i)}) - y^{(i)}) x_j^{(i)}$$

⇒ it's like linear regression partial derivative

↳ it compute prediction error & multiply it by  $j^{\text{th}}$  feature

7- Now (Batch, Mini-Batch, Stochastic) Gradient Descent

## ② Decision Boundary

### 1-D Visualization

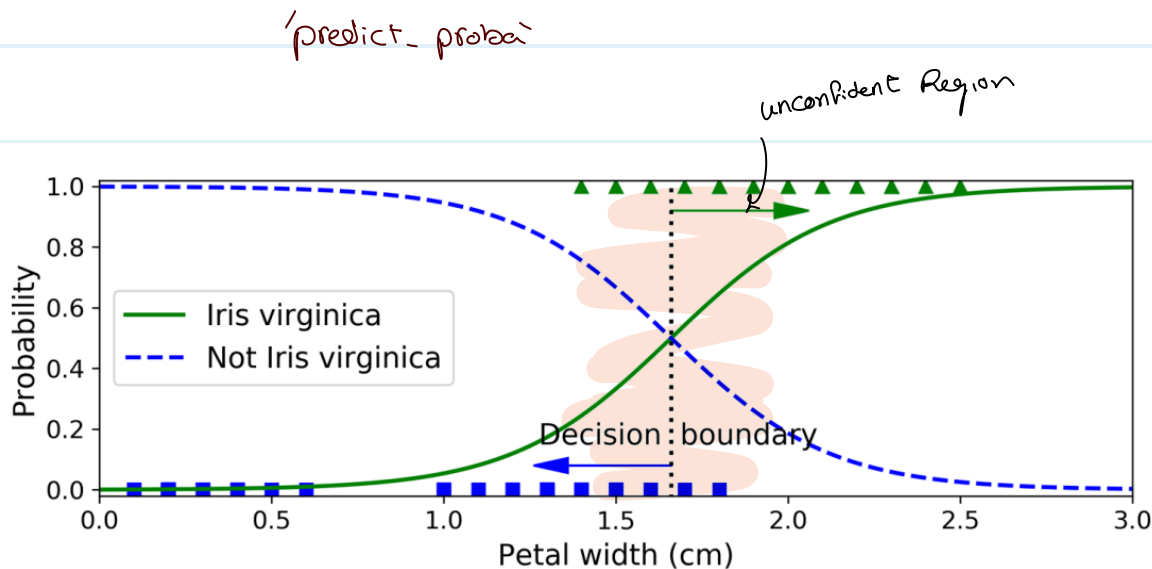
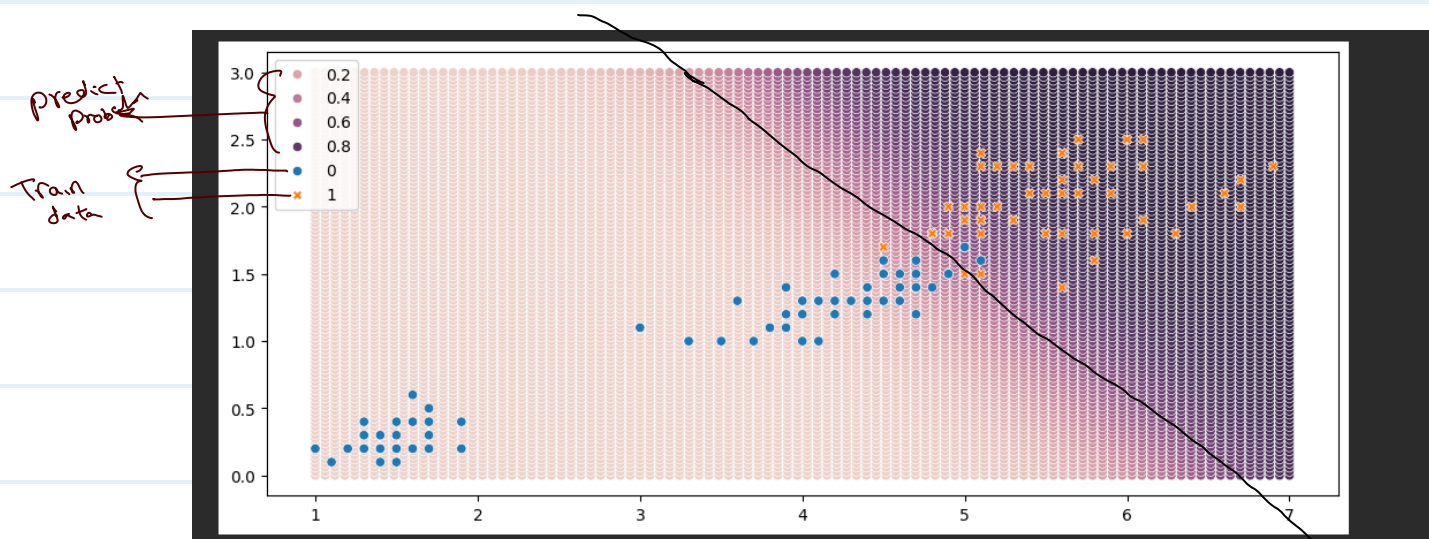


Figure 4-23. Estimated probabilities and decision boundary

### 2-D Visualization



Note: Logistic Regression  $\rightarrow$  use "C" which inverse " $\alpha$ " to control regularization

decision Boundary

"implementation" in Note Book

### ③ Softmax Regression (multinomial Regression)

Logistic Regression  $\Rightarrow$  it's used for multi-class / not multi-output  
 $\Rightarrow$  without train binary classifier for each class

Steps =

1 - Calculate Score function for each class ( $s_k(\vec{x})$ ) for each class  $k$

$$s_k(\vec{x}) = \vec{x}^T \vec{\theta}^{(k)}$$

there's parameters for each class  
 stored as  $\Theta = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_k \end{pmatrix}$

2 - Estimate probability of each class

2.1 - By taking exponential to score

2.2 - then normalize it

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

estimated probability to class "k"

$K = \# \text{ classes}$

$\vec{s}(\vec{x}) = \text{vector of scores per each class}$

3 - Predict class which has  $\Rightarrow$  highest probability  
 $=$  highest score

$$\hat{y} = \underset{k}{\operatorname{argmax}} (\hat{p}_k) = \underset{k}{\operatorname{argmax}} (s_k(\vec{x})) = \underset{k}{\operatorname{argmax}} (\theta_k^T \cdot \vec{x})$$

Backward  
step

#### 4. Cross-Entropy loss

$$\begin{aligned}
 J(\Theta) &= -\sum_{i=1}^m H(p^{(i)}, \hat{p}^{(i)}) && \text{Where } p^{(i)} = [0 \ 0 \ 0 \ 1 \ 0] \\
 &= -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(\hat{p}_k^{(i)}) && \text{for } y=4 \text{ for 5-classes} \\
 &= -\frac{1}{m} \sum_{i=1}^m \log(\hat{p}_{(y)}^{(i)}) \\
 &\quad \uparrow \\
 &\quad \text{estimated probability of true class}
 \end{aligned}$$

Note : 2.1 Note contain Cross Entropy

#### 5. let's calculate partial derivative

$$\nabla_{\theta_{(k)}} J(\Theta) = \frac{1}{m} \sum_{i=1}^m (\hat{p}_k^{(i)} - y_k^{(i)}) \vec{x}^{(i)}$$

#### 6. Update parameters

then use partial derivative to update  $\theta_{(k)}$ , for each class

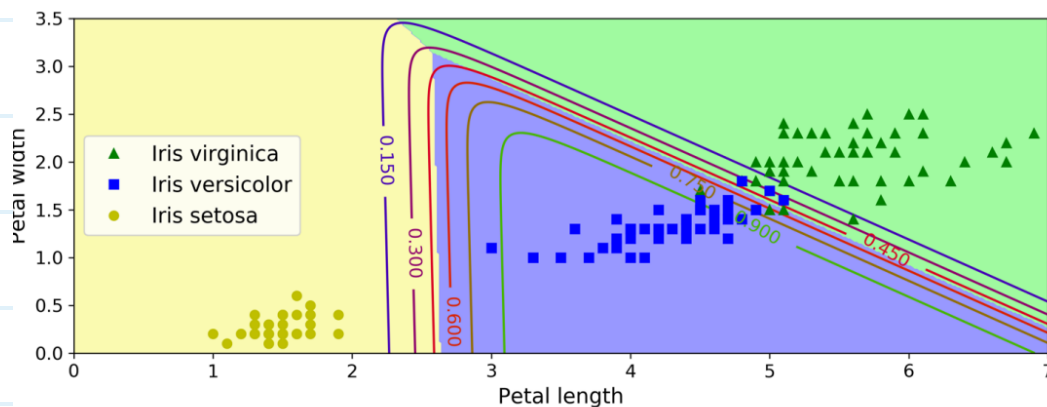


Figure 4-25. Softmax Regression decision boundaries

Platz  
gratis