

SOFTMAX

GENERALIZATION OF LOGISTIC REGRESSION

Multi-Class Classification

LOGISTIC REGRESSION → LEARNING FRAMEWORK

$$x^{(i)} \in \mathbb{R}^m \quad y^{(i)} \in \{0, 1\}$$

$x_0 = 1$

$$h_{\vartheta}(x) = \frac{1}{1 + e^{-\vartheta^T x}}$$

$$\vartheta^T = [\vartheta_0, \vartheta_1, \vartheta_2, \dots, \vartheta_m] \in \mathbb{R}^{m+1}$$

$$x = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_m \end{bmatrix}$$

TO DISTINGUISH
AMONG TWO CLASSES.

$$\text{TRAINING SET} = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

M SAMPLES

GOAL: MINIMIZE THE COST FUNCTION $J(\vartheta)$

$$J(\vartheta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h_{\vartheta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\vartheta}(x^{(i)})) \right]$$

CROSS-ENTROPY LOSS

REPEAT {

$$\vartheta_j \leftarrow \vartheta_j - \alpha \frac{\partial J(\vartheta)}{\partial \vartheta_j}$$

GRADIENT DESCENT IS USED TO FIND PARAMETERS ϑ

FOR ALL j
SIMULTANEOUSLY }

LEARNING RATE: 0.1, 0.03, 0.01, 0.003, 0.001, ...

TO PREVENT OVERFITTING REGULARIZATION IS USED

$$J(\vartheta) = -\frac{1}{m} \left[\dots \right] + \frac{\lambda}{2m} \sum_{j=1}^m \vartheta_j^2$$

1 VS ALL
APPROACH
TO PERFORM
MULTI-CLASS
CLASSIFICATION

SOFTMAX → LEARNING FRAMEWORK TO PERFORM

MULTI-CLASS CLASSIFICATION.

$x^{(i)} \in \mathbb{R}^m$ $x_0^{(i)} = 1$ IT IS AN $m \times m$ EXTENSION OF THE LOGISTIC REGRESSION.

$$y^{(i)} \in \{0, 1, \dots, K-1\}$$

$\underbrace{\text{K CLASSES}}_{\text{ONE HOT ENCODING}}$

$$\text{TRAINING SET} = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

↑
m SAMPLES

WE WANT TO ESTIMATE $P_{\theta}(y=k|x)$ $\forall k=0, 1, \dots, K-1$

⇒ OUR HYPOTHESES HAVE TO ESTIMATE K VECTORS

SOFTMAX MODEL

$$h_{\theta}(x) = \begin{bmatrix} P_{\theta}(0)(y=0|x) \\ P_{\theta}(1)(y=1|x) \\ P_{\theta}(2)(y=2|x) \\ \vdots \\ P_{\theta}(K-1)(y=K-1|x) \end{bmatrix} = \frac{1}{\sum_{k=0}^{K-1} e^{\theta^{(k)\top} x}}$$

\uparrow
NORMALIZE THE DISTRIBUTION

$e^{\theta^{(0)\top} x}$
 $e^{\theta^{(1)\top} x}$
 $e^{\theta^{(2)\top} x}$
 \vdots
 $e^{\theta^{(K-1)\top} x}$

NOTE THAT
 WHEN PARAMETERS
 ARE FIXED
 (LEARNED)

$$\sum_{k=0}^{K-1} P(y=k|x) = 1$$

$v \in \mathbb{R}^{(m+1) \times K}$
 $v = \begin{bmatrix} v_0^{(0)} & v_1^{(0)} & \dots & v_{m-1}^{(0)} \\ v_0^{(1)} & v_1^{(1)} & \dots & v_{m-1}^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ v_0^{(K-1)} & v_1^{(K-1)} & \dots & v_{m-1}^{(K-1)} \end{bmatrix}$
 IS A MATRIX

SOFTMAX COST FUNCTION

- IT USES CROSS-ENTROPY LOSS FUNCTION

- LET $\mathbb{1}\{\cdot\}$ THE INDICATOR FUNCTION

$$\mathbb{1}\{\text{PROPOSITION}\} = \begin{cases} 1 & \text{IF PROPOSITION IS TRUE} \\ 0 & \text{IF PROPOSITION IS FALSE} \end{cases}$$

$$\bullet \text{LET } h_{\psi}^{(c)}(x) = P(y=c|x) = \frac{e^{\psi^{(c)T}x}}{\sum_{k=0}^{K-1} e^{\psi^{(k)T}x}}$$

$$J(\psi) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{c=0}^{K-1} \mathbb{1}\{y^{(i)}=c\} \log \left(h_{\psi}^{(c)}(x^{(i)}) \right) \right] + \frac{\lambda}{2m} \sum_{c=0}^{K-1} \sum_{j=1}^m \psi_j^{(c)2}$$

REGULARIZATION TERM

LEARNING GOAL IS TO FIND THE PARAMETERS

$$\boldsymbol{\vartheta} = \begin{bmatrix} \vartheta_0^{(0)} & \vartheta_0^{(1)} & \vartheta_0^{(2)} & \dots & \vartheta_0^{(K-1)} \\ \vartheta_1^{(0)} & \vartheta_1^{(1)} & \vartheta_1^{(2)} & \dots & \vartheta_1^{(K-1)} \\ \vdots & \vdots & \vdots & & \vdots \\ \vartheta_m^{(0)} & \vartheta_m^{(1)} & \vartheta_m^{(2)} & \dots & \vartheta_m^{(K-1)} \end{bmatrix}$$

THIS IS DONE
BY FINDING
THE MINIMUM
OF THE DERIVATIVE
OF THE COST
FUNCTION
WITH GRADIENT
DESCENT (OR WITH
STOCHASTIC GRADIENT
DESCENT)

REPEAT {

$$\hat{\vartheta}_j^{(c)} \leftarrow \hat{\vartheta}_j^{(c)} - \alpha \frac{\partial J(\boldsymbol{\vartheta})}{\partial \hat{\vartheta}_j^{(c)}}$$

FOR ALL $j=0, \dots, m$ AND $c=0, \dots, K-1$

SIMULTANEOUSLY UNTIL CONVERGENCE {

$$\frac{\partial J(\boldsymbol{\vartheta})}{\partial \hat{\vartheta}_j^{(c)}} = -\frac{1}{m} \sum_{i=1}^m \left[\left(\{y^{(i)}=c\} - h_{\boldsymbol{\vartheta}}^{(c)}(x^{(i)}) \right) x_j^{(i)} \right] + \frac{\lambda}{m} \hat{\vartheta}_j^{(c)}$$

WITH REGULARIZATION

To PREVENT
OVERFITTING

PROPERTY OF SOFTMAX

SOFTMAX MODEL IS "OVERPARAMETERIZED", WHICH MEANS

THAT FOR ANY HYPOTHESES WHICH MIGHT FIT THE

DATA, THERE ARE MULTIPLE PARAMETERS SETTINGS

THAT GIVE EXACTLY THE SAME SOLUTION TO MAP

DATA $x \in \mathbb{R}^m$ TO PROBABILITYS $y \in \{0, \dots, K-1\}$

$$h_{\varphi^{(c)}}(x^{(i)}) = P_{\varphi^{(c)}}(y^{(i)} = c \mid x^{(i)}) = \frac{\ell^{(\varphi^{(c)} - \psi)^T x^{(i)}}}{\sum_{k=0}^{K-1} \ell^{(\varphi^{(k)} - \psi)^T x^{(i)}}}$$

$$= \frac{\ell^{(\varphi^{(c)} - \psi)^T x^{(i)}}}{\ell^{(\varphi^{(0)} - \psi)^T x^{(i)}} + \sum_{k=0}^{K-1} \ell^{(\varphi^{(k)} - \psi)^T x^{(i)}}} = \frac{\ell^{(\varphi^{(c)} - \psi)^T x^{(i)}}}{\sum_{k=0}^{K-1} \ell^{(\varphi^{(k)} - \psi)^T x^{(i)}}}$$

SUBTRACTING
 $\psi \in \mathbb{R}^{m+1}$
 TO EACH
 PARAMETER
 $\varphi^{(c)} \quad c=0 \dots K-1$
 DO NOT AFFECT
 THE HYPOTHESIS

h_{φ}

$$\psi = \begin{bmatrix} \psi_0 \\ \vdots \\ \psi_m \end{bmatrix}$$

SOFTMAX COST FUNCTION ALONG WITH THE LOGISTIC
REGRESSION COST FUNCTION:

FOR $K=2$ (NOT CONSIDERING REGULARIZATION)

$$\begin{aligned}
 J(\boldsymbol{\varphi}) &= -\frac{1}{m} \left[\sum_{i=1}^m \sum_{c=0}^1 \mathbb{1}\{y^{(i)} = c\} \log(h_{\boldsymbol{\varphi}}^{(c)}(\mathbf{x}^{(i)})) \right] = \\
 &= -\frac{1}{m} \left[\sum_{i=1}^m (1-y^{(i)}) \log(h_{\boldsymbol{\varphi}}^{(0)}(\mathbf{x}^{(i)})) + y^{(i)} \log(h_{\boldsymbol{\varphi}}^{(1)}(\mathbf{x}^{(i)})) \right] = \\
 &= -\frac{1}{m} \left[\sum_{i=1}^m (1-y^{(i)}) \log(1-h_{\boldsymbol{\varphi}}^{(1)}(\mathbf{x}^{(i)})) + y^{(i)} \log(h_{\boldsymbol{\varphi}}^{(1)}(\mathbf{x}^{(i)})) \right]
 \end{aligned}$$

LOGISTIC REGRESSION
IS A SPECIAL CASE OF SOFTMAX

$$h_{\boldsymbol{\varphi}}^{(0)} + h_{\boldsymbol{\varphi}}^{(1)} = 1 \Rightarrow h_{\boldsymbol{\varphi}}^{(0)} = 1 - h_{\boldsymbol{\varphi}}^{(1)}$$

RELATIONSHIP BETWEEN LOGISTIC REGRESSION AND SOFTMAX

By using the property of "rank parametrization" it

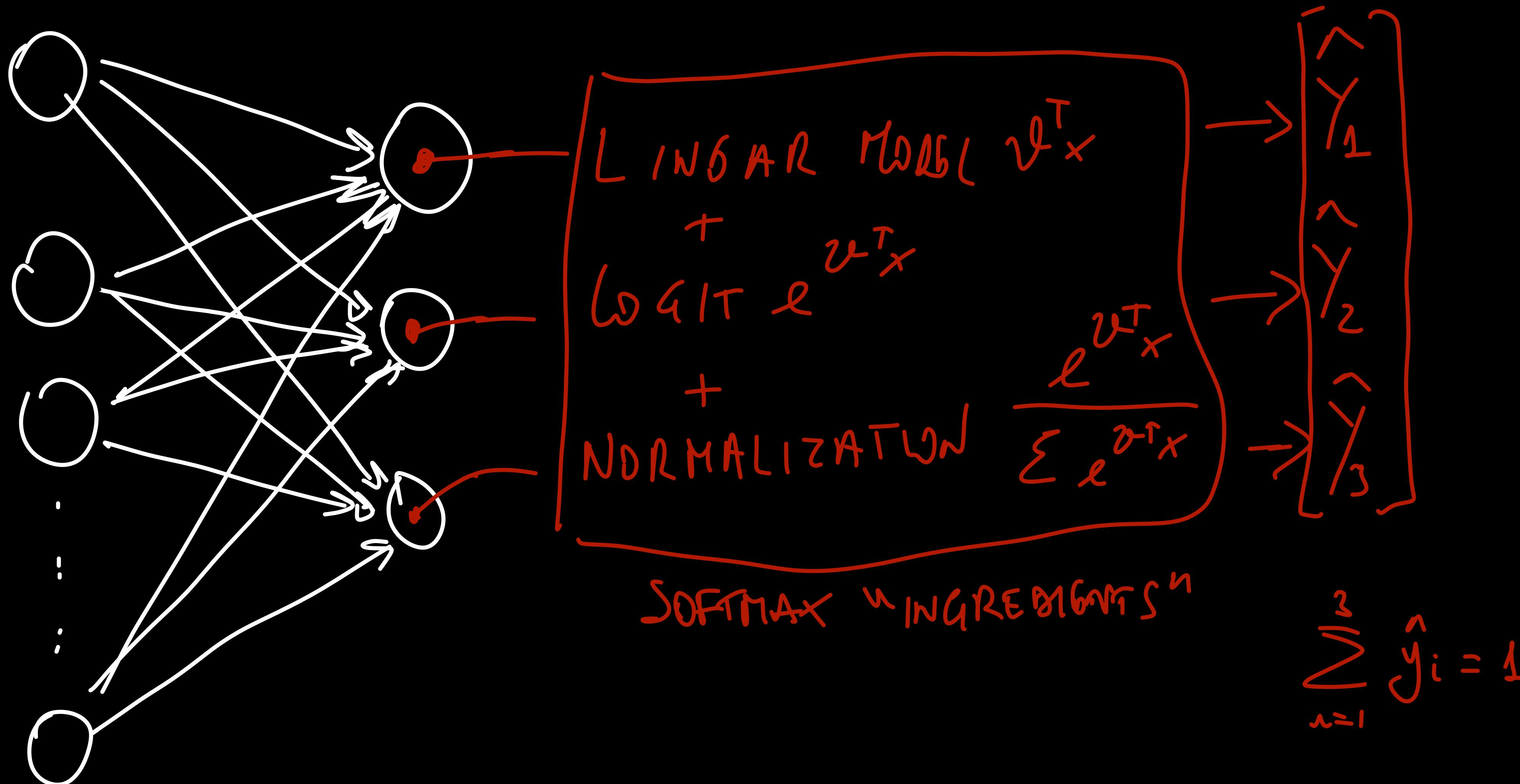
is simple to show that in the special case of $K=2$

the softmax regression reduces to logistic regression

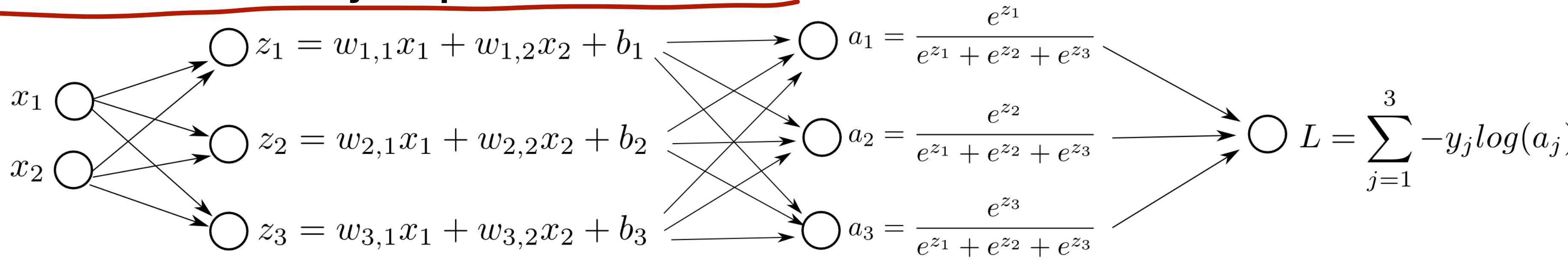
$$\begin{aligned}
 h_{\theta}(x) &= \left[\frac{e^{v^{(0)}T x}}{e^{v^{(0)}T x} + e^{v^{(1)}T x}} \quad \frac{e^{v^{(1)}T x}}{e^{v^{(0)}T x} + e^{v^{(1)}T x}} \right] = \left[\frac{(v^{(0)} - v^{(1)})^T x}{e^{(v^{(0)} - v^{(1)})^T x} + e^{(\vec{0})^T x}} \quad \frac{e^{(\vec{0})^T x}}{e^{(v^{(0)} - v^{(1)})^T x} + e^{(\vec{0})^T x}} \right] \\
 &= \left[\frac{(v^{(0)} - v^{(1)})^T x}{1 + e^{(v^{(0)} - v^{(1)})^T x}} \quad \frac{1}{1 + e^{(v^{(0)} - v^{(1)})^T x}} \right]^T = \begin{bmatrix} 1 - & \\ & 1 \end{bmatrix}^T
 \end{aligned}$$

$v^{(0)} - v^{(1)}$ can be replaced by a single parameter ϑ

SOFTMAX "NETWORK"



The softmax is a very simple Neural Network



$$b_j = v_{j,0}$$

Model PARAMETERS

$$v^T = \begin{bmatrix} b_1 & w_{1,1} & w_{1,2} \\ b_2 & w_{2,1} & w_{2,2} \\ b_3 & w_{3,1} & w_{3,2} \end{bmatrix}$$

NOTATION
CONSIDERANDO
SLIDE
PREVIOUSNTI

$$v_j^{(c)} = w_{c,j} \quad \forall j=1,2,3 \quad \forall c=1,2,3$$

$$v_0^{(c)} = b_c \quad \forall c=1,2,3$$

INPUT

$$X = \begin{bmatrix} x_0 = 1 \\ x_1 \\ x_2 \end{bmatrix}$$

OUTPUT \equiv CLASS LABEL

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

$$y_j \in \{0, 1\} \quad \forall j=1,2,3$$

$$\sum_j y_j = 1$$

PARAMETERS

$$v^T = \begin{bmatrix} -1 & 1 & 0 \\ 1 & 0.5 & 1 \\ 2 & 1 & 1.5 \end{bmatrix}$$

INPUT

$$X = \begin{bmatrix} 1 \\ 0.2 \\ 0.1 \end{bmatrix}$$

LINEAR MODEL OUTPUT

$$z = \begin{bmatrix} -0.8 \\ 1.2 \\ 2.35 \end{bmatrix}$$

LOGIT OUTPUT

$$e^z = \begin{bmatrix} 0.45 \\ 3.32 \\ 10.48 \end{bmatrix}$$

NORMALIZATION (SOFTMAX OUTPUT)

$$a_j = \frac{e^{z_j}}{\sum_i e^{z_i}}$$

COMPUTING LOSS L

1-HOT ENCODING
LABEL CLASS

$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

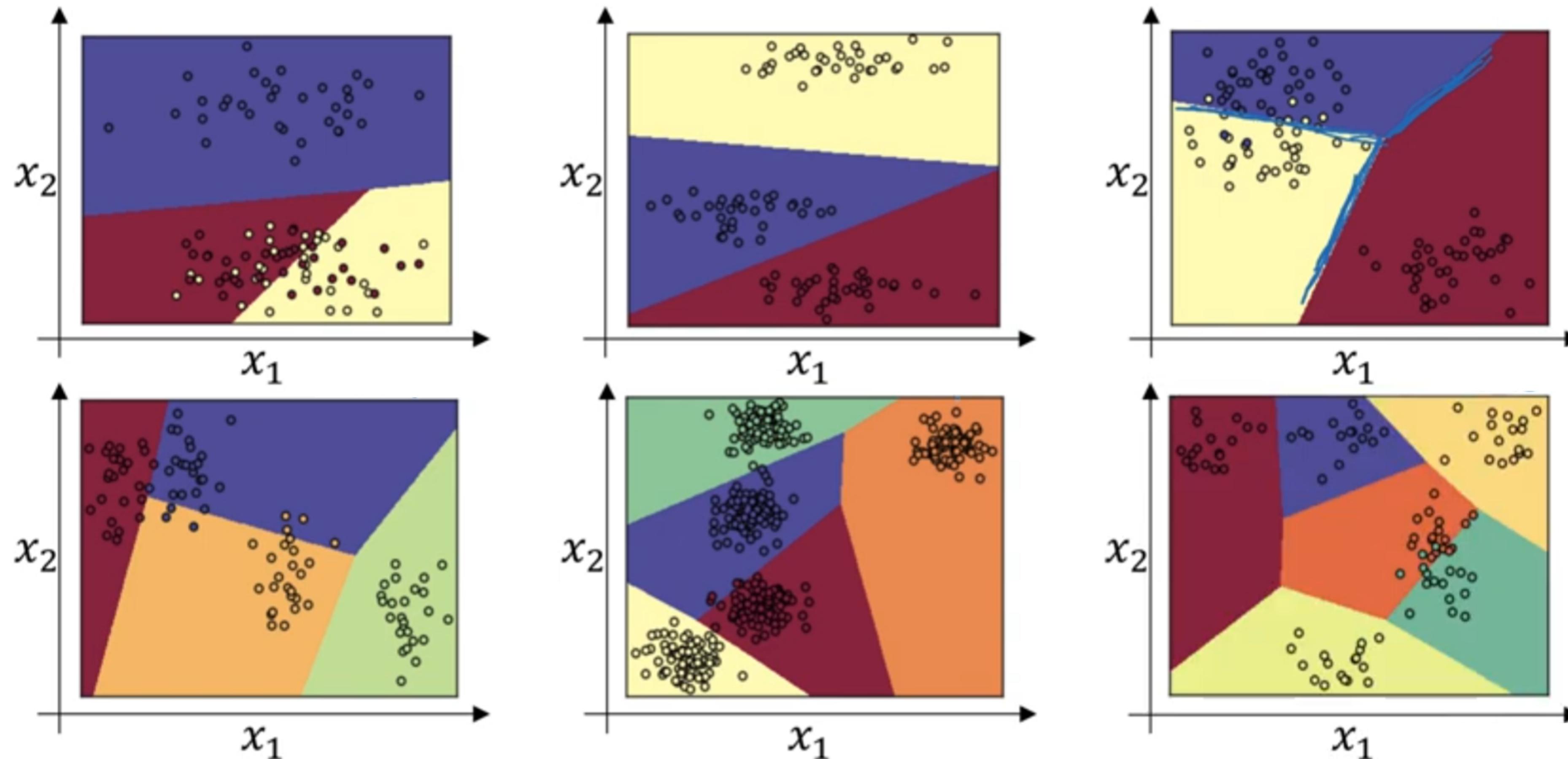
$$\sum_j e^{z_j} = 14.25$$

$$-1 \cdot \log(0.03)$$

LOSS SU UN CAMPIONE
FARMO MEDIA SUTUTTI I CAMPIONI

$$\sum_{i=1}^m l^{(i)}$$

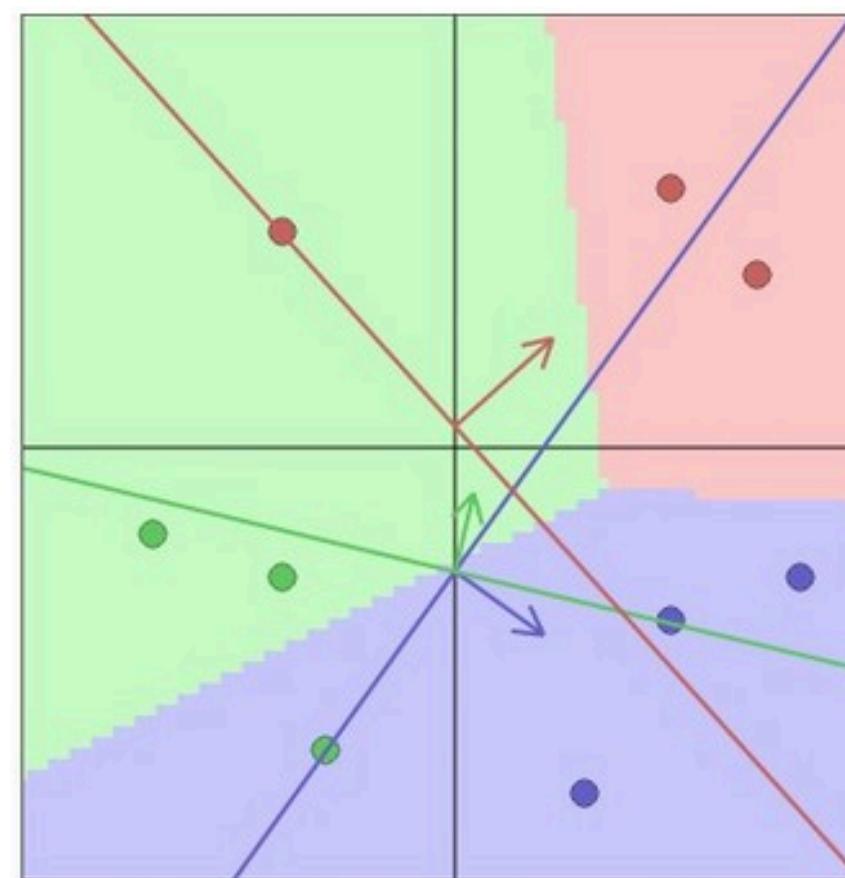
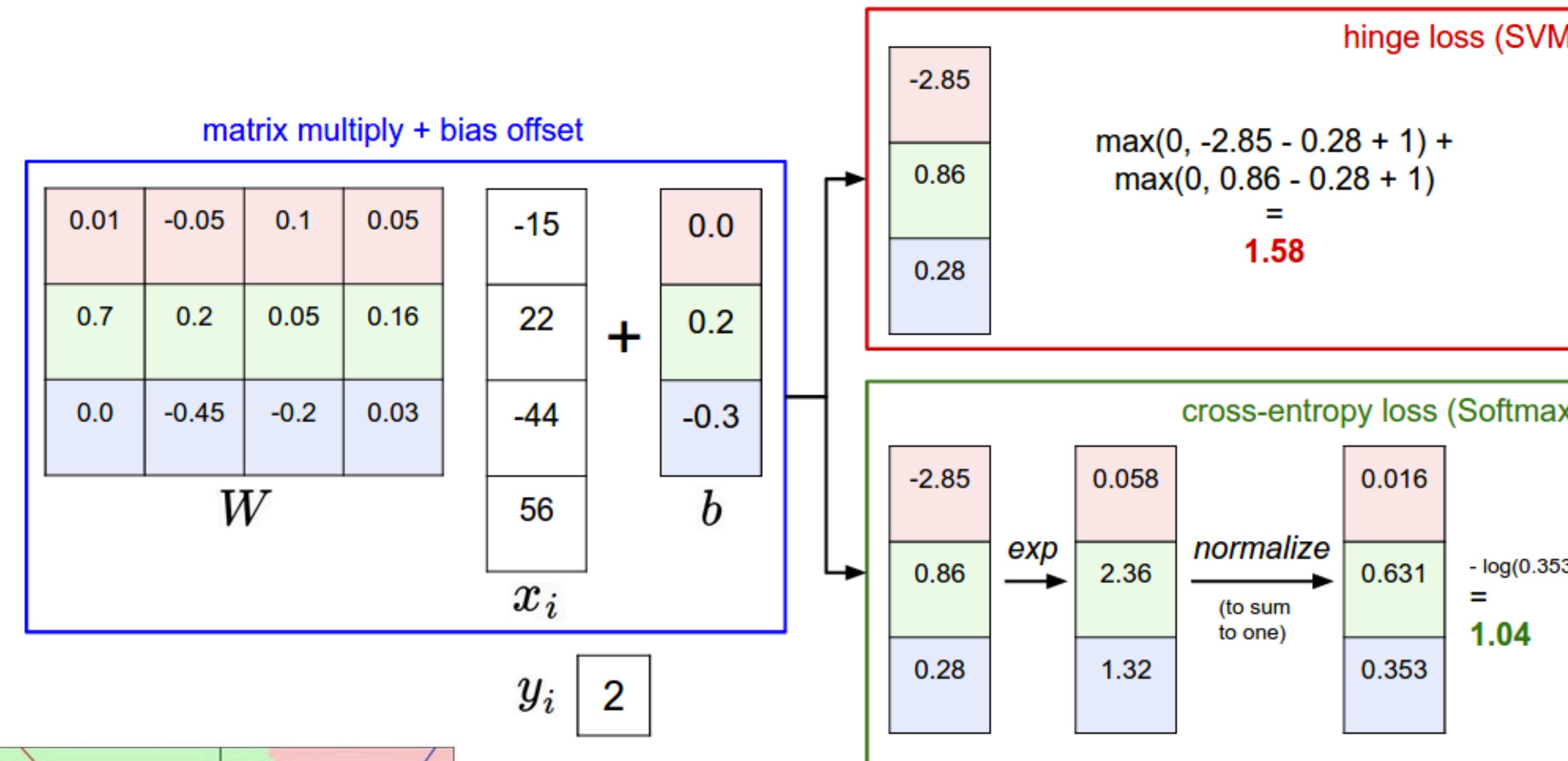
Softmax Classifier - Esempi



LIN6AR BOUNDARIES AND WE CAN USE $\phi(x)$ K6LN6L
(OUTPUT OF AN6WORK) TO REPRE6NT DATA
AND HAVING NON-LIN6AR BOUNDARIES

Possiamo usare
AD ESEMPIO UN
ME6ICO POLINOMIALE

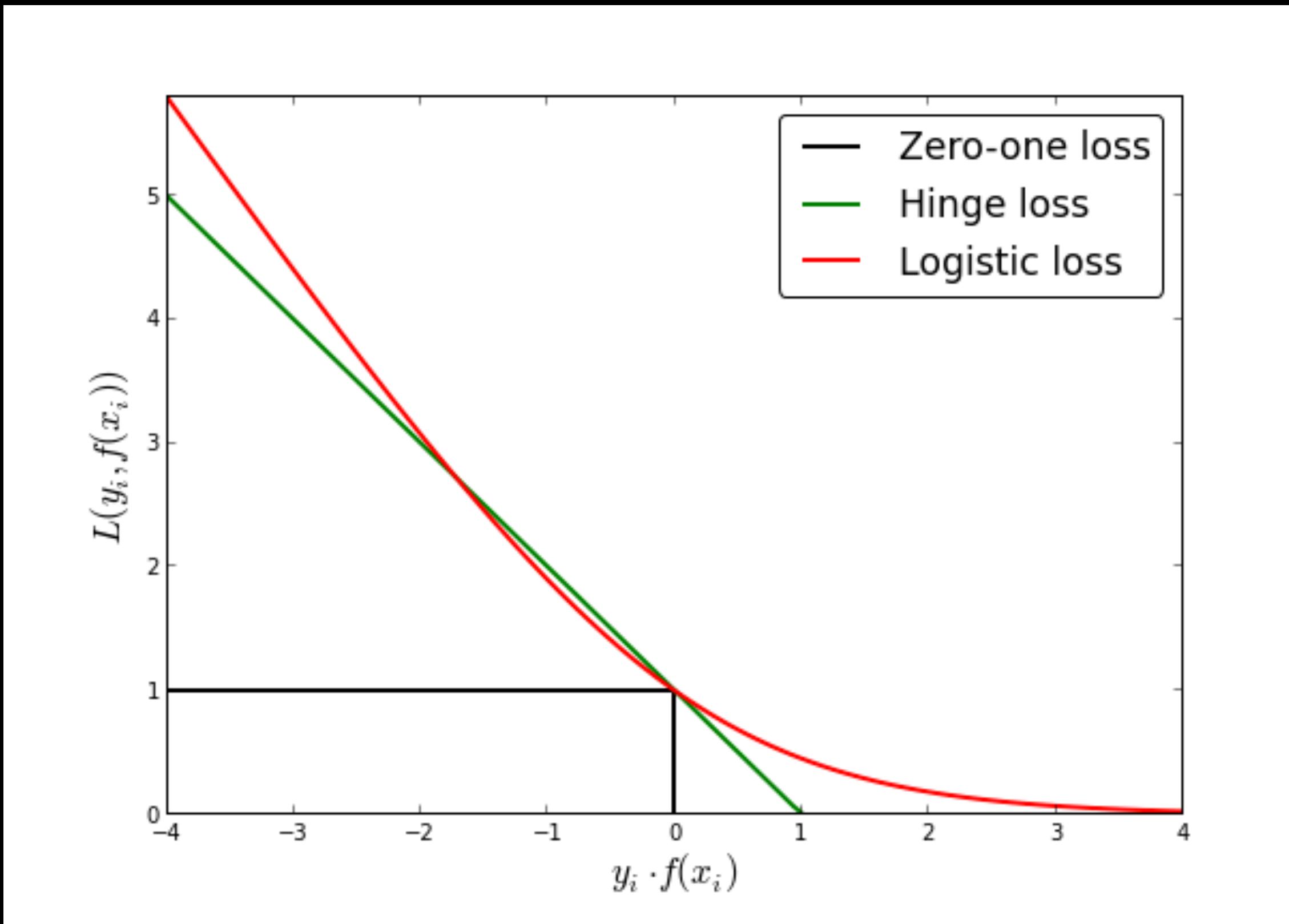
Changing the Loss Function in this linear machine framework you can obtain the Multi-Class SVM



$$L = \underbrace{\frac{1}{N} \sum_i \sum_{j \neq y_i} \max(0, f_j - f_{y_i} + 1)}_{\text{data loss}} + \lambda \underbrace{\sum_k \sum_l W_{k,l}^2}_{\text{regularization loss}}$$

SAME FRAMEWORK
OF SOFTMAX
BUT A DIFFERENT
LOSS FUNCTION
IS USED

Loss Functions for Classification Some Examples



$$L(x, y) = \sum_{i=1}^m \text{loss}\left(f(x^{(i)}), y^{(i)}\right)$$

ZERO-ONE LOSS

$$\mathbb{1}_{\{f(x^{(i)}) = y^{(i)}\}}$$

INTEGRATOR
FUNCTION

HINGE LOSS

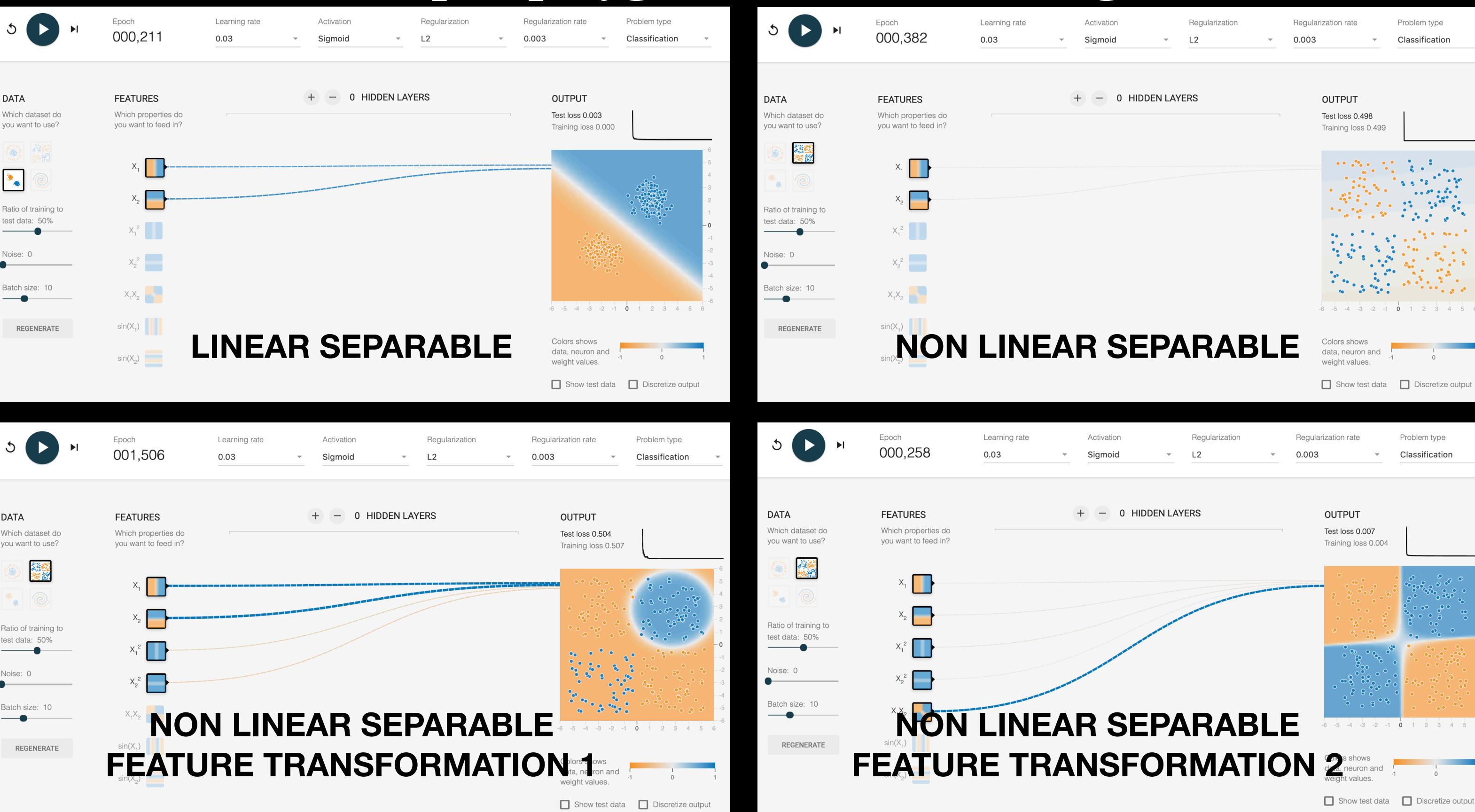
$$\max(0, 1 - f(x^{(i)})y^{(i)})$$

LOGISTIC

$$\log(1 + \exp(f(x^{(i)})y^{(i)}))$$

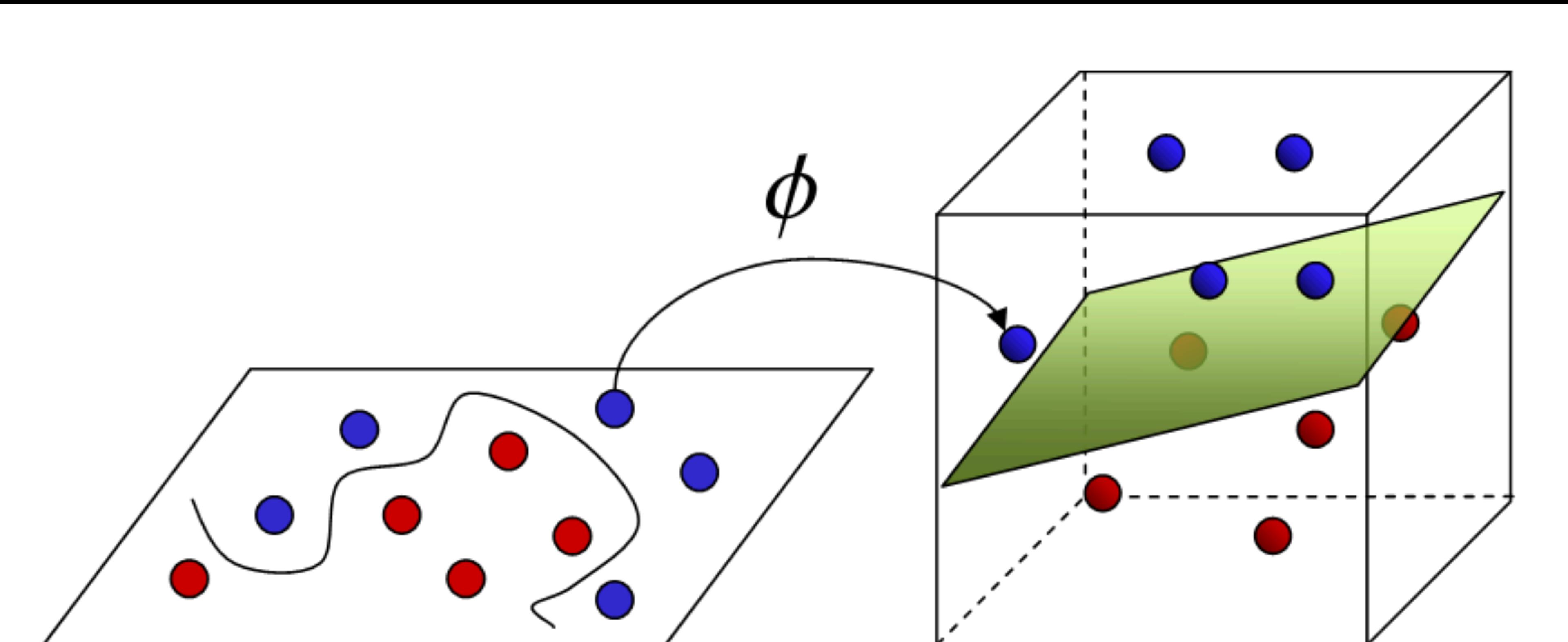
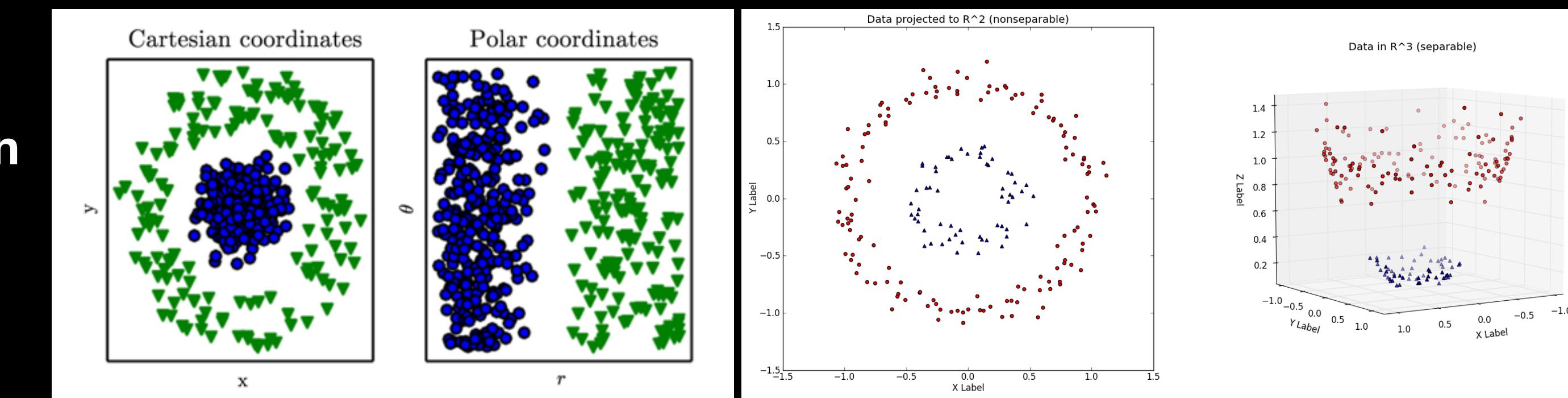
BINARY CLASSIFICATION

<https://playground.tensorflow.org>



Role of Representation

How to define?



Input Space

Feature Space

• REPRESENTATION LEARNING INTUITION

- WE HAVE A NEW FRAMEWORK FOR MULTI-CLASS CLASSIFICATION = SOFTMAX
- WE KNOW THAT IF WE HAVE A GOOD KERNEL FUNCTION $\phi(x)$ TO MAP FEATURES IN HIGH DIMENSIONAL NON LINEAR FEATURES SPACE WE CAN SOLVE COMPLEX PROBLEM
- HOW TO CHOOSE $\phi(x)$ IS NOT SIMPLE
CAN WE LEARN A FUNCTION $\phi(x)$ FOR REPRESENTATION?
 \Rightarrow NEURAL NET & DEEP LEARNING.

Network of Neurons - Learning Representation

<https://playground.tensorflow.org>

