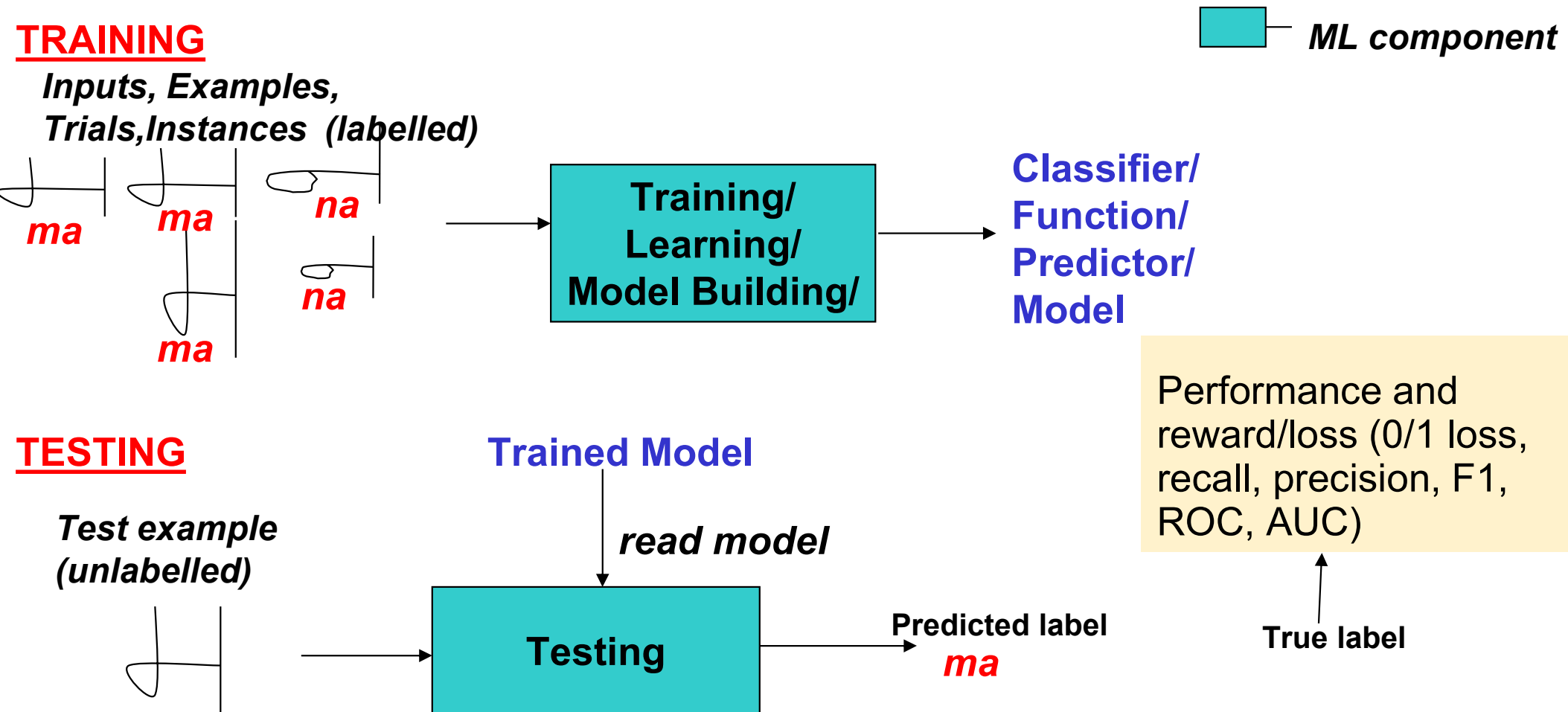


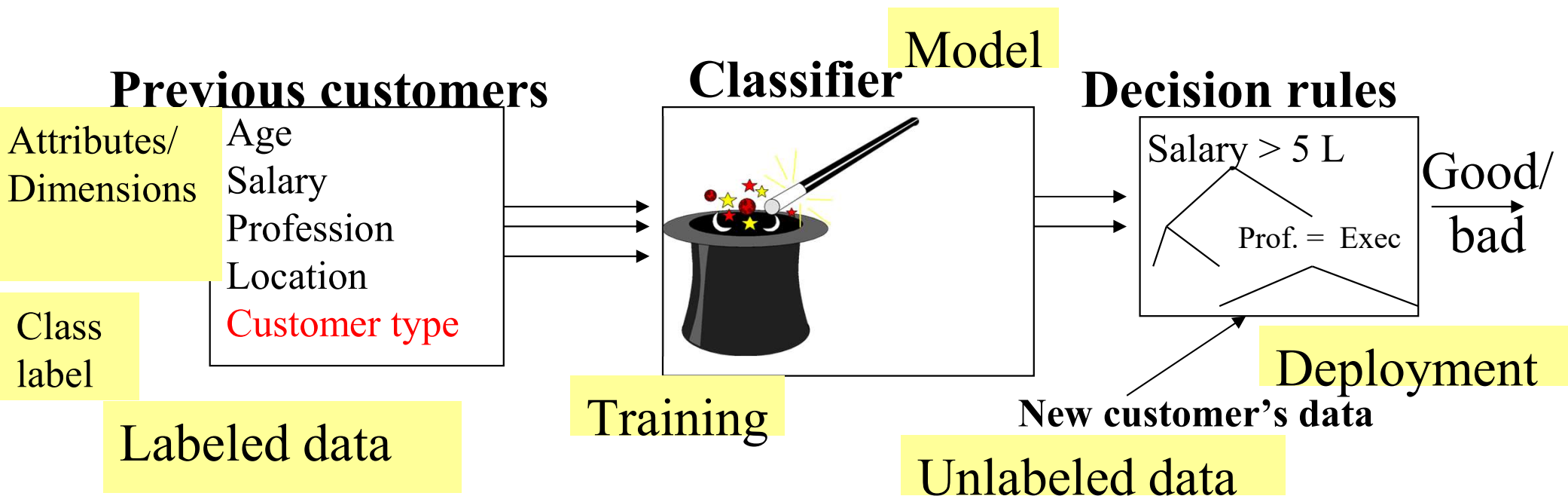
Linear Classifiers
or
Multiclass logistics regression
classifiers

Classification



Classification example

- Given old data about customers and payments, predict new applicant's loan eligibility.

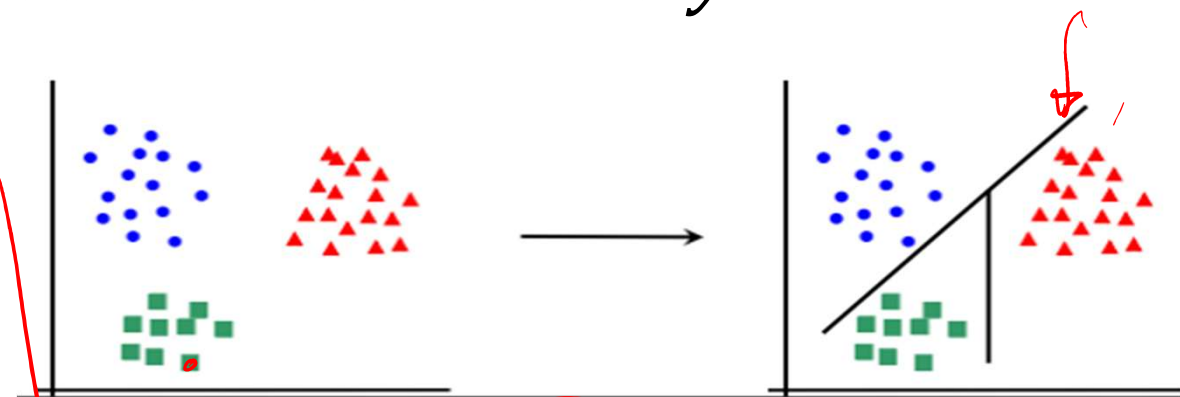


Classification: setup

- Input: Like in regression, obtained after feature engineering $x \in \mathbb{R}^d$
- Output: one of K possible class labels $y \in [1, 2, \dots, K]$
- Training data:

$$X = \begin{bmatrix} x_1^1, x_2^1, \dots, x_d^1 \\ \vdots \\ x_1^N, x_2^N, \dots, x_d^N \end{bmatrix}$$

Classifier: $C(x) \rightarrow \hat{y} \in [1, \dots, K]$



An example training dataset $d=2, K=3$

$$Y = \begin{bmatrix} 1 \\ 1 \\ 3 \\ \vdots \end{bmatrix}$$

Over-fitting in classifiers



From: <https://www.javatpoint.com/overfitting-in-machine-learning>

Types of classifiers

- Discriminative
 - Logistic regression ~~4~~
 - Decision trees
 - Neural Network
- Probabilistic
 - Generative
 - Conditional
- Kernel-based
 - Nearest neighbor classifier
 - Support Vector Machines

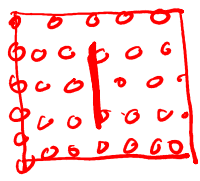
Linear classifiers

- For each class $\underline{k} \in [1, \dots, K]$, define a linear function $d+1$
 $K(d+1)$
 - $f_k(\mathbf{x}) = \underline{w_{k1}}x_1 + \dots + \underline{w_{kd}}x_d + \underline{w_{k0}}$
scoring functions, logits
- Assign predicted label \hat{y} as the class for which score is maximized.

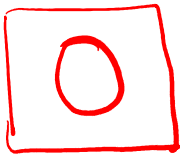
$$\hat{y} = \underset{k \in [1 \dots K]}{\operatorname{argmax}} f_k(\mathbf{x})$$

Example

Letter recognition ϕ Features:



$x^1 = 4$



x^2



x^3



x^4

$x_1 \equiv \# \text{ of vertical lines}$

$x_2 \equiv \# \text{ of circles}$

$x_3 \equiv \# \text{ of horizontal lines}$

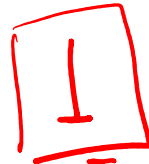
$x_4 \equiv \# \text{ of slanted lines}$

$x^1: [1, 0, 0, 0]$ $x^2: [0, 1, 0, 0]$

$x^3: [1, 0, 1, 1]$

$x^4: [2, 0, 1, 0]$

$x^5: [1, 0, 1, 0]$



x^5

$w_{11}, w_{12}, w_{13}, w_{14} = [10, -1, -1, -1], w_{10} = 0$

$w_{21}, w_{22}, w_{23}, w_{24} = [-1, 10, -1, -1], w_{20} = 0$

$w_{31}, w_{32}, w_{33}, w_{34} = [5, 0, 10, 5], w_{30} = 0$

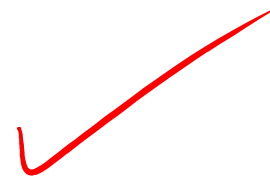
$$f_1(x^4) \equiv 18$$

$$f_2(x^4) \equiv -1$$

$$f_3(x^4) \equiv 20$$

$$\operatorname{argmax}_{k \in \{1, 2, 3\}} (f_k(x)) = \hat{y} = 3$$

$$y = 3$$



Special case of binary classifiers

$$k = 2$$

$$f_1(x) = w_{11}x_1 + w_{12}x_2 + \dots + w_{1d}x_d + w_{10}$$

$$f_2(x) = w_{21}x_1 + w_{22}x_2 + \dots + w_{2d}x_d + w_{20}$$

$$\arg \max_{k=(1,2)} f_k(x) \equiv \arg \max(\underline{f_1(x)}, \underline{f_2(x)}) \equiv \arg \max(f_1(x) - f_2(x), 0)$$

$$\begin{aligned} f_1(x) &\equiv (w_{11} - w_{21})x_1 + (w_{12} - w_{22})x_2 + \dots + (w_{10} - w_{20}) \\ &\equiv \underbrace{(\vec{w}_1 - \vec{w}_2)}_{\vec{w}} \vec{x} + (w_{10} - w_{20}) \end{aligned}$$

$\vec{w} = \vec{w}_1 - \vec{w}_2$

$$f(x) = w \cdot x + w_0$$

Training objective

- Find parameters such that predicted class match actual class in training data
- Consider training instance: x^i, y^i

$\hat{y} = \underset{y}{\operatorname{argmax}} f_y(x^i; w)$ ← this needs to match true y^i

~~\max_w~~ $|f_{y^i}(x^i; w) - \max_y f_y(x^i; w)| \neq 0$ for correct prediction

→ $E(x^i, y^i, w) \equiv -f_{y^i}(x^i; w) + \max_y f_y(x^i; w) \geq 0 < 0$ for incorrect classification.

Training objective: minimize $w \sum_{i=1}^N E(x^i, y^i, w)$ not differentiable.

Differentiable rewrite with softmax

- Maximum of K real values:

- $\max(o_1, o_2, \dots, o_K)$ $\approx \log(e^{o_1} + e^{o_2} + \dots + e^{o_K})$

HW

o_1 was the maximum

$$\begin{aligned} &= o_1 + \log\left(e^{o_1 - o_1} + \underbrace{e^{o_2 - o_1}}_{\leq 1} + \underbrace{e^{o_K - o_1}}_{\leq 1}\right) \\ &= o_1 + \log(1 + 0.0001 + \dots) \\ &\approx o_1 \end{aligned}$$

Final objective

minimize $\sum_{i=1}^N \text{Loss}(x^i, y^i, w)$

$\rightarrow w$

$\text{loss}(x, y, w)$
as the surrogate
for $E(x, y, w)$

$$= \sum_{i=1}^N -f_y(x^i, w) + \log \sum_{y=1}^K e^{f_y(x^i, w)}$$

Demo

- <https://colab.research.google.com/drive/1MtvbOIHgmkV05GJXfefisim6TiClEXts?usp=sharing>