



UNIVERSIDADE
CATÓLICA
PORTUGUESA

BRAGA

Deep Learning

Session 5

Training Neural Networks

Applied Data Science

2024/2025

Objective Function: Analogous to Training

Babies to not throw food on the floor



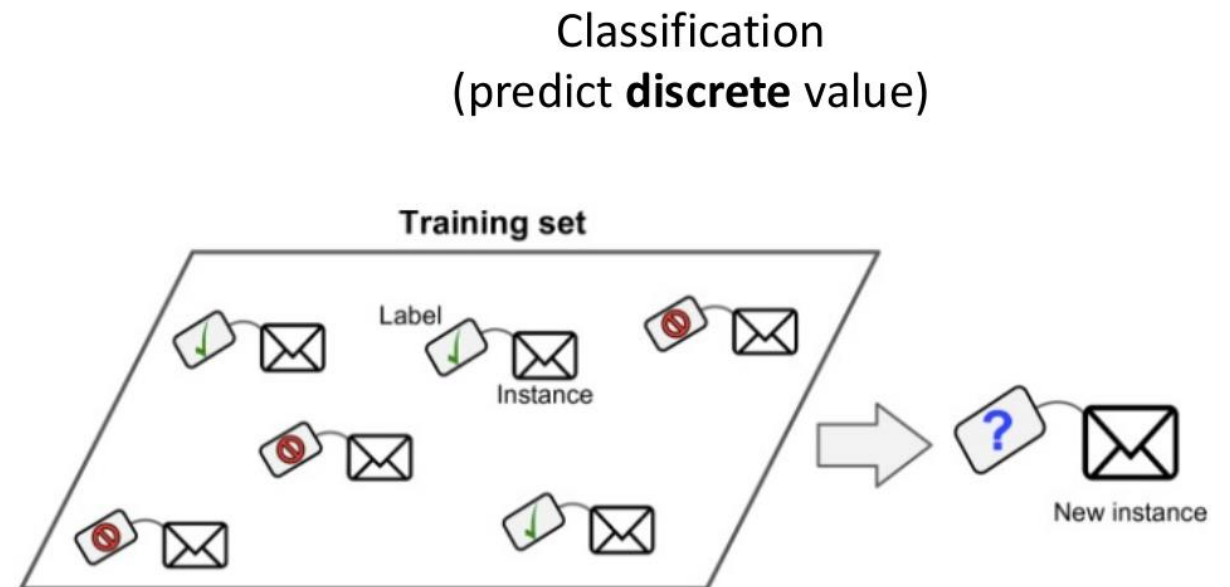
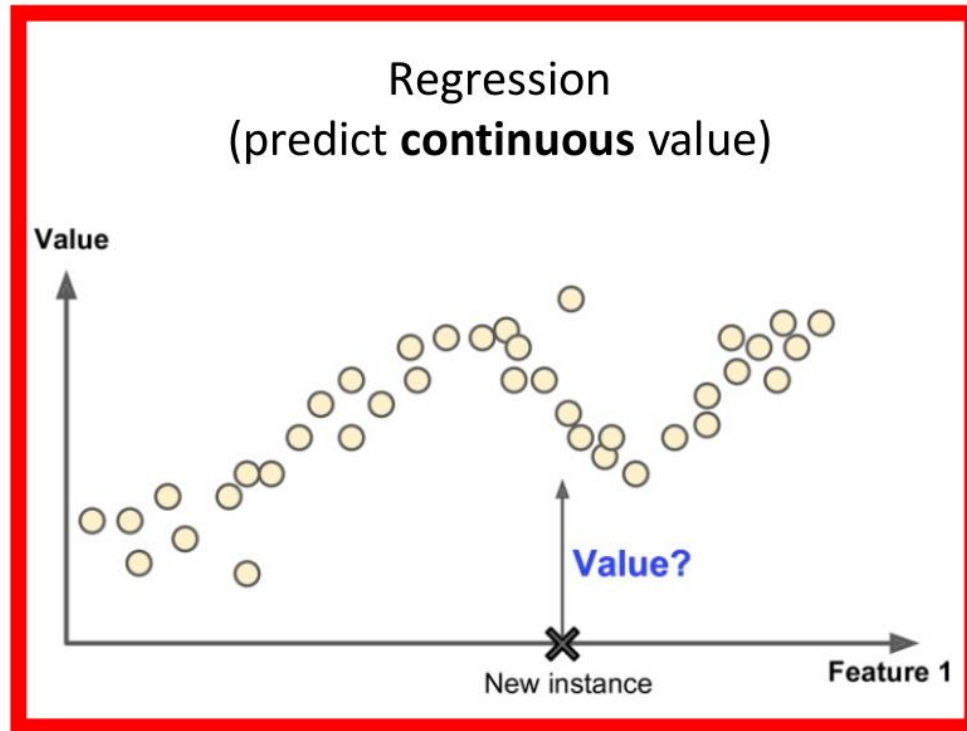
<https://www.youtube.com/watch?v=58Pr-QrVNqU>

Dogs to learn to sit

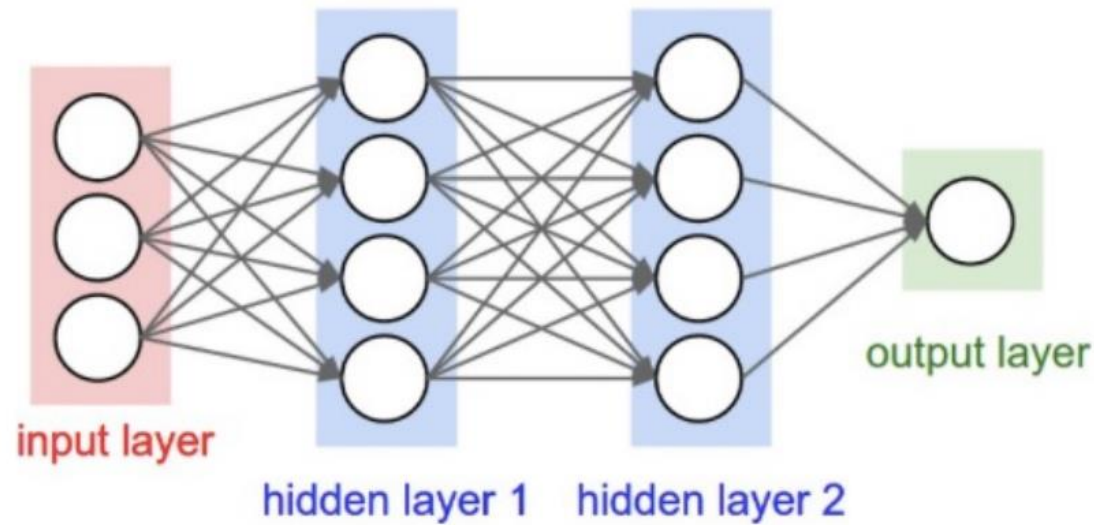


<https://www.akc.org/expert-advice/training/how-to-become-a-dog-trainer/>

Objective Function: Learn Model Parameters that Achieve a Specified Goal



Objective Function: Learn Model Parameters that Achieve a Specified Goal



e.g., make as small as possible the squared error (aka, L2 loss, quadratic loss)

Mean taken over n instances

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

True value

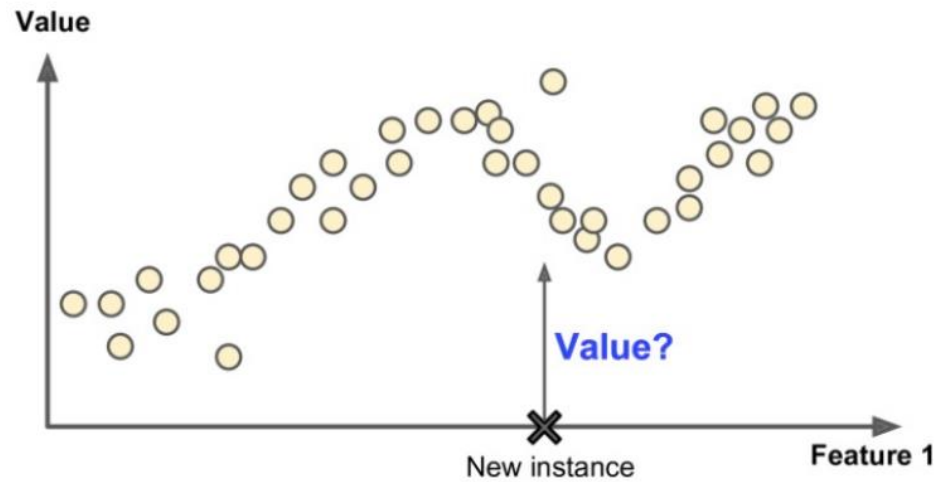
Predicted value

What is the range of possible values?

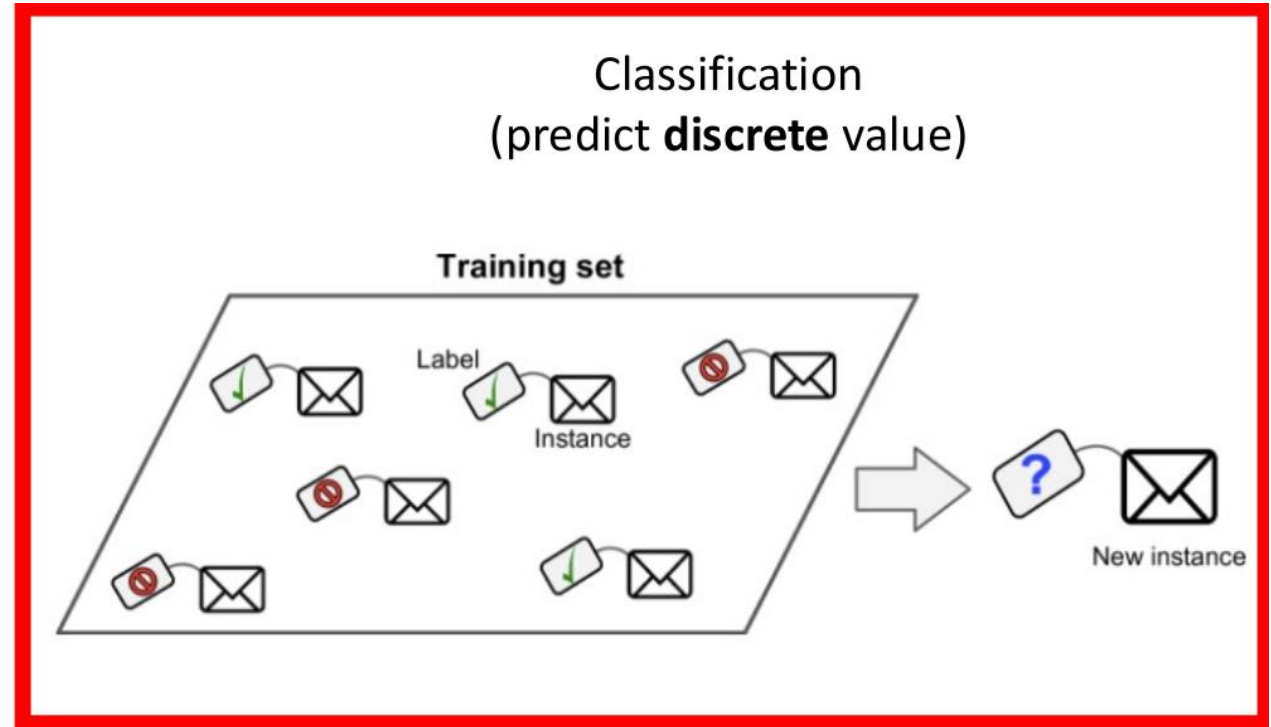
- Minimum: 0
 - i.e., all correct predictions
- Maximum: Infinity
 - i.e., incorrect predictions

Objective Function: Learn Model Parameters that Achieve a Specified Goal

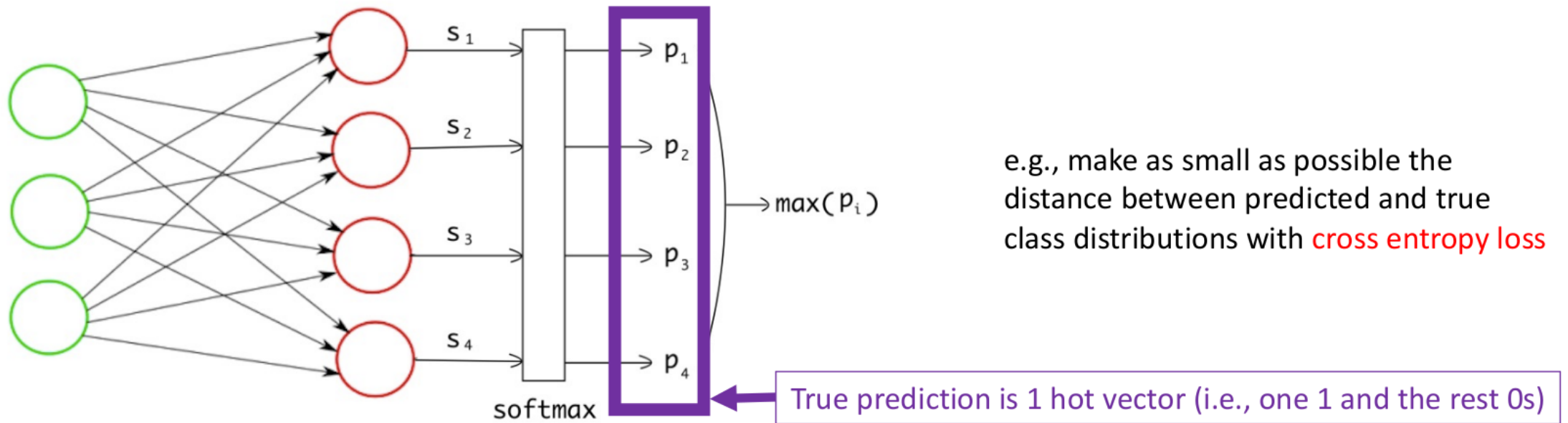
Regression
(predict **continuous** value)



Classification
(predict **discrete** value)



Objective Function: Learn Model Parameters that Achieve a Specified Goal



<https://towardsdatascience.com/multi-label-image-classification-with-neural-network-keras-ddc1ab1afede>

Objective Function: Learn Model Parameters that Achieve a Specified Goal

Probability distribution of predicted class

Probability distribution of true class

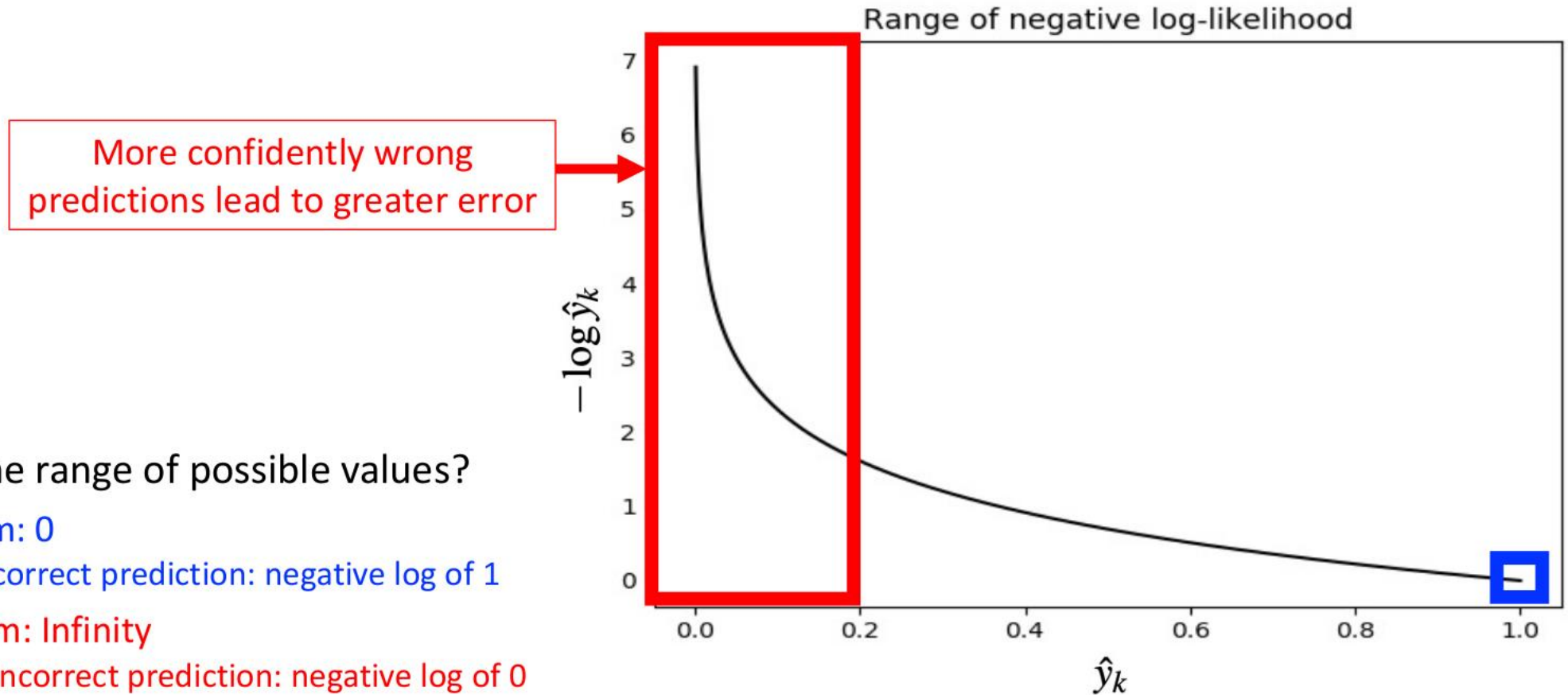
Number of classes?

$$L_{\text{CE}}(\hat{y}, y) = - \sum_{k=1}^K y_k \log \hat{y}_k$$

What is the range of possible values?

- Minimum: 0
 - i.e., correct prediction: negative log of 1
- Maximum: Infinity
 - i.e., incorrect prediction: negative log of 0

Objective Function: Learn Model Parameters that Achieve a Specified Goal

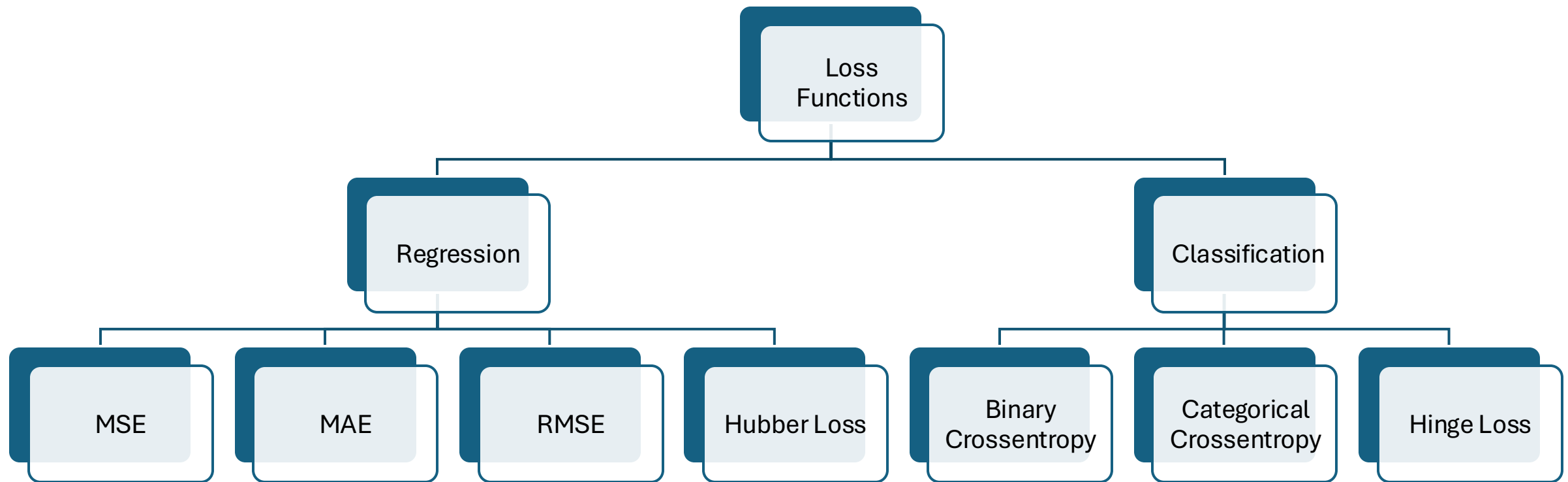


What is the range of possible values?

- **Minimum: 0**
 - i.e., correct prediction: negative log of 1
- **Maximum: Infinity**
 - i.e., incorrect prediction: negative log of 0

Objective Function

- Many objective functions exist, and we will explore popular ones in this course!



Loss Functions



10 Most Common Loss Functions in Machine Learning

 [blog.DailyDoseofDS.com](https://blog.dailydoseofds.com)

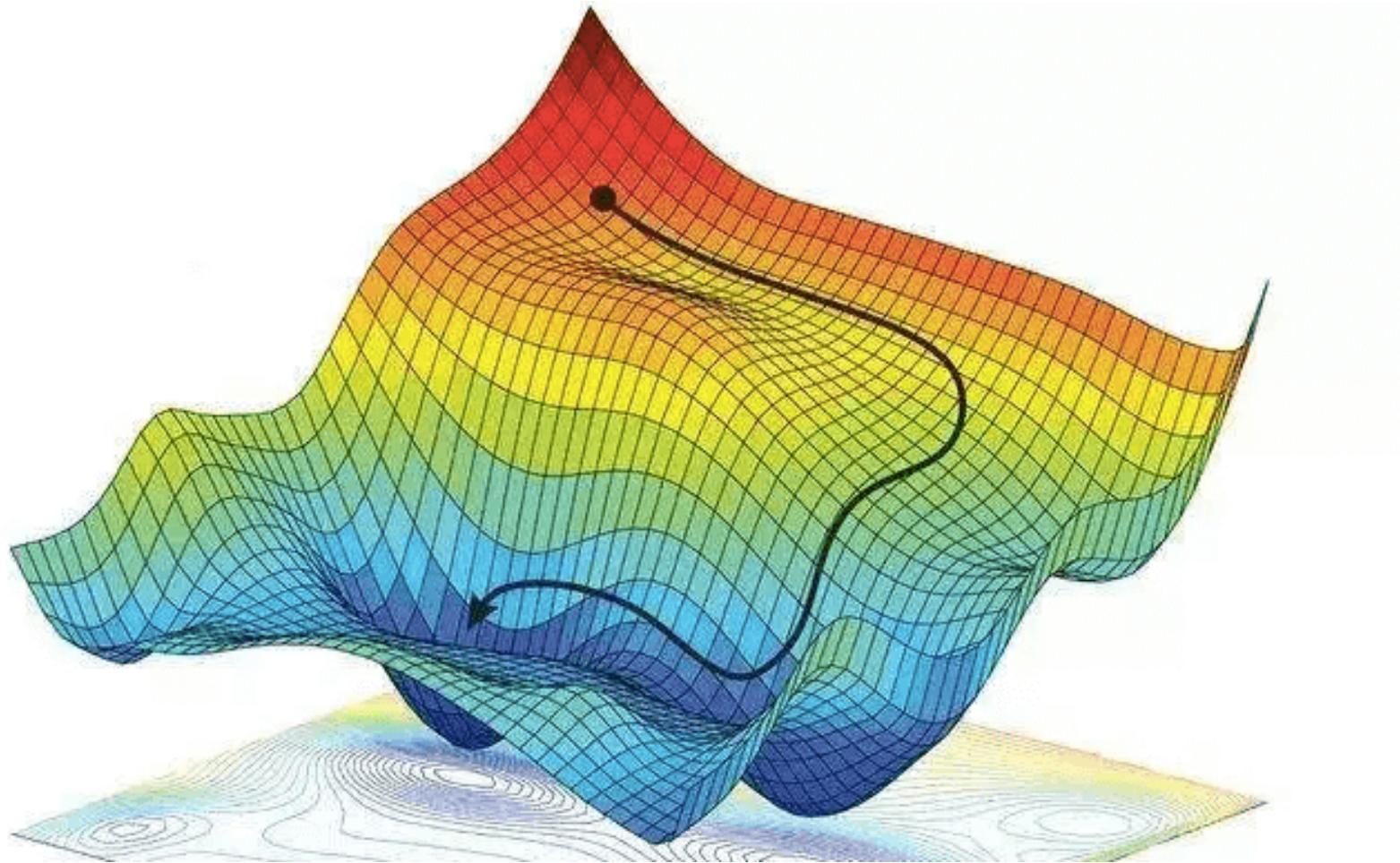
Loss Function Name	Description	Function
Regression Loss Functions		
Mean Bias Error	Captures average bias in prediction. But is rarely used for training.	$\mathcal{L}_{MBE} = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))$
Mean Absolute Error	Measures absolute average bias in prediction. Also called L1 Loss.	$\mathcal{L}_{MAE} = \frac{1}{N} \sum_{i=1}^N y_i - f(x_i) $
Mean Squared Error	Average squared distance between actual and predicted. Also called L2 Loss.	$\mathcal{L}_{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))^2$
Root Mean Squared Error	Square root of MSE. Loss and dependent variable have same units.	$\mathcal{L}_{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))^2}$
Huber Loss	A combination of MSE and MAE. It is parametric loss function.	$\mathcal{L}_{Huberloss} = \begin{cases} \frac{1}{2}(y_i - f(x_i))^2 & : y_i - f(x_i) \leq \delta \\ \delta(y_i - f(x_i) - \frac{1}{2}\delta) & : otherwise \end{cases}$
Log Cosh Loss	Similar to Huber Loss + non-parametric. But computationally expensive.	$\mathcal{L}_{LogCosh} = \frac{1}{N} \sum_{i=1}^N \log(\cosh(f(x_i) - y_i))$
Classification Loss Functions (Binary + Multi-class)		
Binary Cross Entropy (BCE)	Loss function for binary classification tasks.	$\mathcal{L}_{BCE} = \frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(x_i)) + (1 - y_i) \cdot \log(1 - p(x_i))$
Hinge Loss	Penalizes wrong and right (but less confident) predictions. Commonly used in SVMs.	$\mathcal{L}_{Hinge} = \max(0, 1 - (f(x) \cdot y))$
Cross Entropy Loss	Extension of BCE loss to multi-class classification.	$\mathcal{L}_{CE} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} \cdot \log(f(x_{ij}))$ <i>N : samples; M : classes</i>
KL Divergence	Minimizes the divergence between predicted and true probability distribution	$\mathcal{L}_{KL} = \sum_{i=1}^N y_i \cdot \log\left(\frac{y_i}{f(x_i)}\right)$

<https://blog.dailydoseofds.com/p/10-regression-and-classification>

Gradient Descent: How to Learn?



UNIVERSIDADE
CATOLICA
PORTUGUESA
BRAGA



Loss Optimization

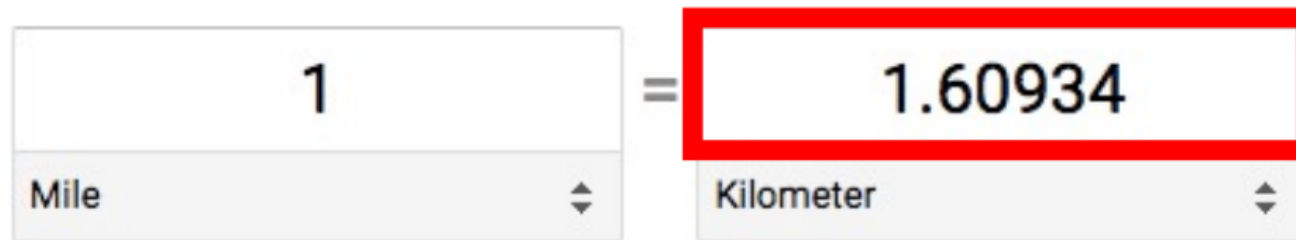
- We want to find the network weights that achieve the lowest loss!

$$W^* = \operatorname{argmin}_W \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; W), y^{(i)})$$

$$W^* = \operatorname{argmin}_W J(W)$$

Gradient Descent: Intuition

- Repeat:
 1. Guess
 2. Calculate error
- e.g., learn linear model for converting kilometers to miles when only observing the input “miles” and output “kilometers”



Gradient Descent: Intuition

- Repeat:
 1. **Guess**
 2. Calculate error



A digital conversion tool interface. On the left, a box contains the number '1' above a dropdown menu labeled 'Mile'. To the right of this is an equals sign. Further right is another box containing the number '1.60934' above a dropdown menu labeled 'Kilometer'. The box containing '1.60934' is highlighted with a thick red border.

1 → Kilometers = miles x **constant**

Gradient Descent: Intuition

- Repeat:
 1. Guess
 2. Calculate error

1	=	1.60934
Mile		Kilometer



Gradient Descent: Intuition

- Repeat:
 1. **Guess**
 2. Calculate error

1	=	1.60934
Mile		Kilometer

1.5 → Kilometers = miles x constant

Gradient Descent: Intuition

- Repeat:
 1. Guess
 2. Calculate error

1	=	1.60934
Mile		Kilometer

1.5 → $\text{Kilometers} = \text{miles} \times \text{constant}$ → $\text{Error} = \text{Guess} - \text{Correct}$

- Idea: iteratively adjust the **constant (i.e., model parameter)** to try to reduce the error

Gradient Descent: Intuition

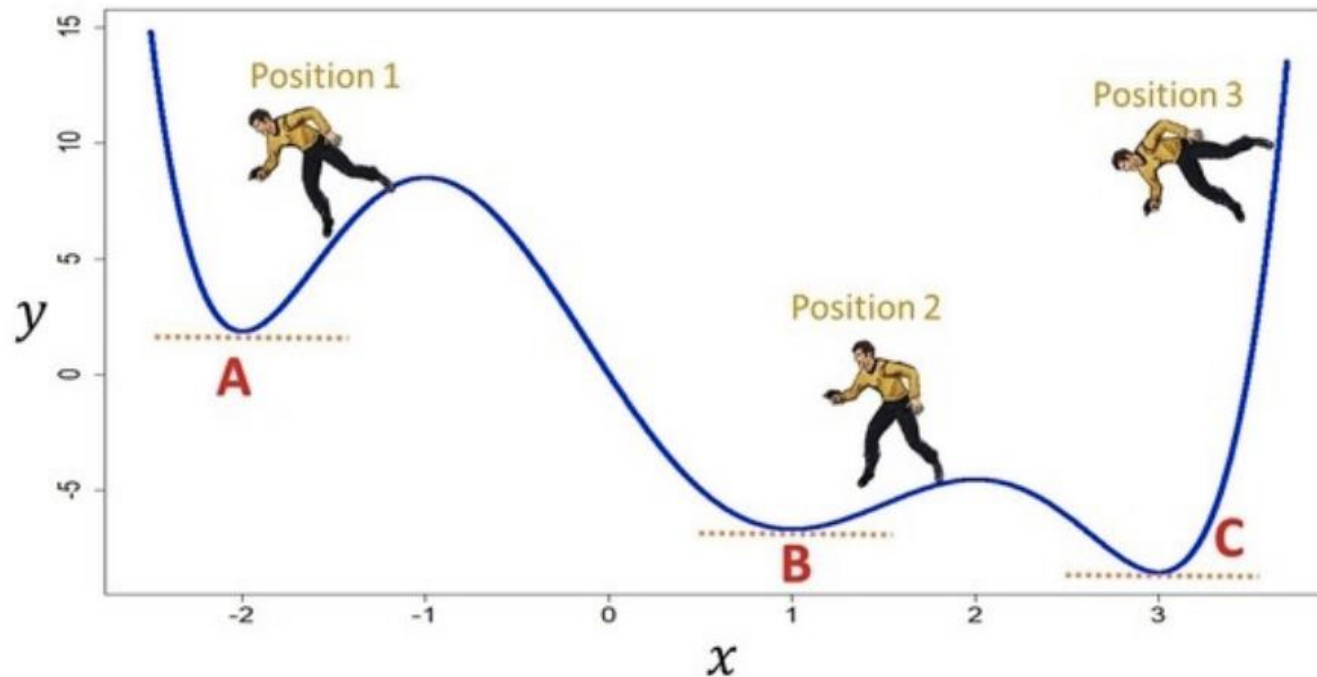
- Iteratively search for values that solve optimization problem (i.e., minimize or maximize an objective function)
- When **minimizing** the objective function, it also is often called interchangeably the **cost function**, **loss function**, or **error function**.

Gradient Descent: Intuition

- Idea: use derivatives!
 - Derivatives tell us how a small change in the input x affects the output $f(x)$.
 - Specifically, the derivative $f'(x)$ represents the rate of change of the function at a given point, indicating how much $f(x)$ will change in response to a small change in x .
 - Gradient is a vector that indicates how $f(x)$ changes as each function variable changes (i.e., partial derivatives).

Gradient Descent: Intuition

- Gradient descent:
 - Iteratively take steps in the opposite direction of the gradient to minimize the function



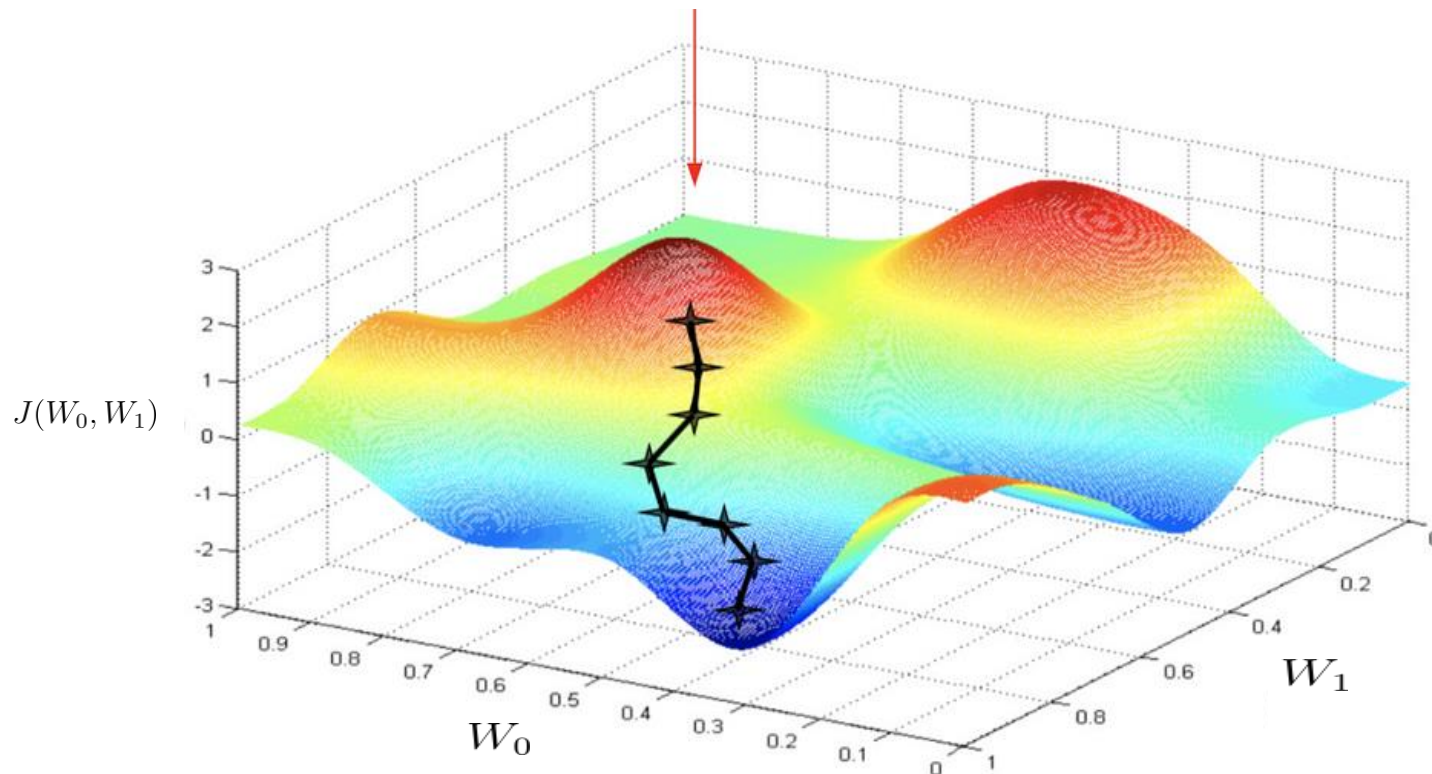
Which letter(s) are the global minima?

Which letter(s) are local minima?

Gradient Descent

$$W^* = \underset{W}{\operatorname{argmin}} J(W)$$

1. Randomly pick an initial (W_0, W_1)



2. Compute the gradient: $\frac{\partial J(W)}{\partial W}$

3. Take a small step in the opposite direction of the gradient.

4. Repeat until convergence.

Gradient Descent

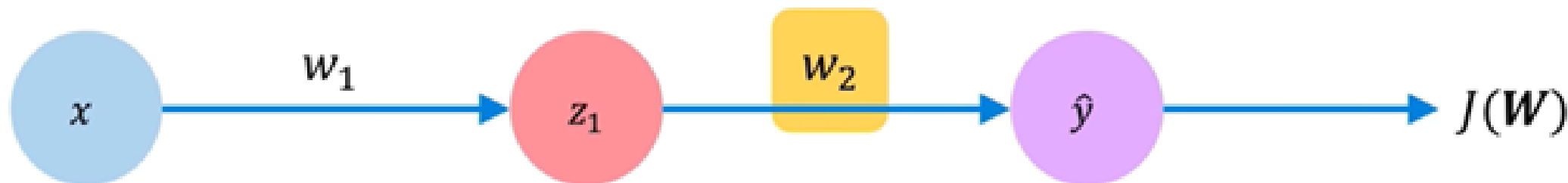
- Algorithm:
 1. Initialize weights randomly $\sim N(0, \sigma^2)$
 2. Loop until convergence
 3. Compute gradient $\frac{\partial J(W)}{\partial W}$
 4. Update weights $W \leftarrow W - \eta \frac{\partial J(w)}{\partial w}$
 5. Return weights

Computing Gradients: Backpropagation



- How does a small change in one weight (e.g., w_2) affect the final loss $J(W)$?

Computing Gradients: Backpropagation

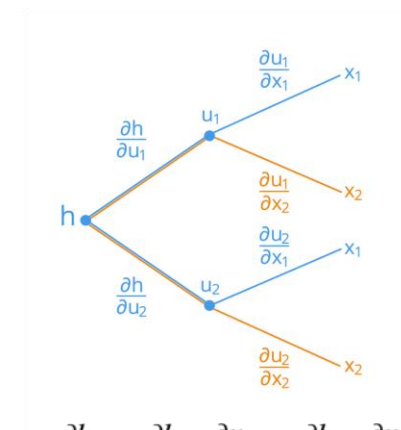


$$\frac{\partial J(W)}{\partial w_2} =$$

Let's use the chain rule!

$$\frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx}$$

$$\frac{dy}{dx} = \left(\begin{array}{c} \text{Differentiate} \\ \text{outer function} \\ \text{Keep the inside} \\ \text{the same} \end{array} \right) \left(\begin{array}{c} \text{Differentiate} \\ \text{inner function} \end{array} \right)$$

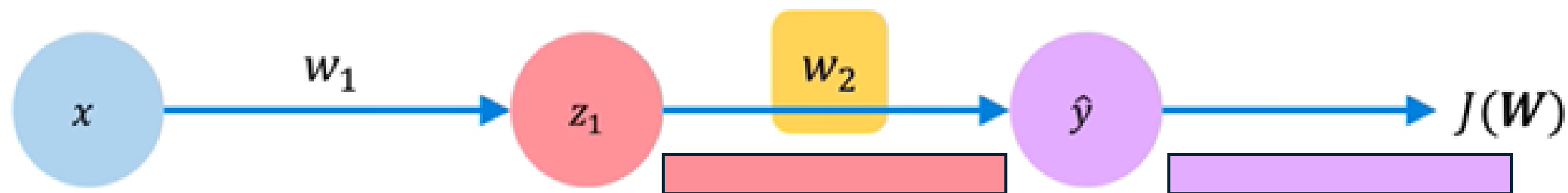


$$\frac{\partial h}{\partial x_1} = \frac{\partial h}{\partial u_1} \cdot \frac{\partial u_1}{\partial x_1} + \frac{\partial h}{\partial u_2} \cdot \frac{\partial u_2}{\partial x_1}$$

$$\frac{\partial h}{\partial x_2} = \frac{\partial h}{\partial u_1} \cdot \frac{\partial u_1}{\partial x_2} + \frac{\partial h}{\partial u_2} \cdot \frac{\partial u_2}{\partial x_2}$$

$$y(x) = (x^2 + 1)^3 \longrightarrow ?$$

Computing Gradients: Backpropagation



$$\frac{\partial J(W)}{\partial w_2} = \underbrace{\frac{\partial J(W)}{\partial \hat{y}}}_{\text{purple bar}} * \underbrace{\frac{\partial \hat{y}}{\partial w_2}}_{\text{red bar}}$$

Computing Gradients: Backpropagation

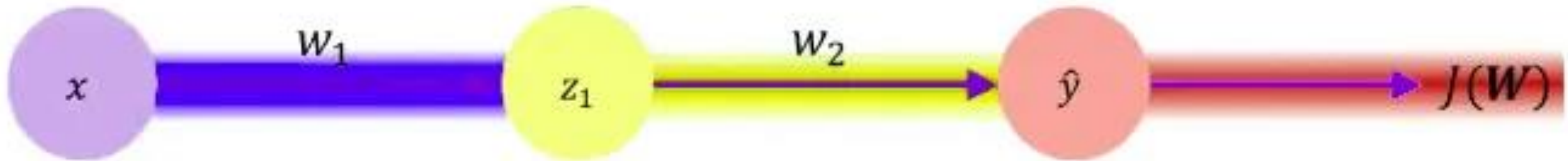


$$\frac{\partial J(W)}{\partial w_1} = \underbrace{\frac{\partial J(W)}{\partial \hat{y}}}_{\text{red box}} * \underbrace{\frac{\partial \hat{y}}{\partial w_1}}_{\text{yellow box}}$$

↑
Apply the chain rule!

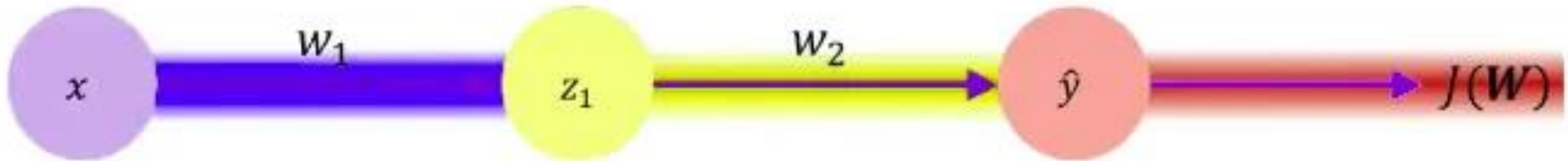
↑
Apply the chain rule!

Computing Gradients: Backpropagation



$$\frac{\partial J(W)}{\partial w_1} = \underbrace{\frac{\partial J(W)}{\partial \hat{y}}}_{\text{blue}} * \underbrace{\frac{\partial \hat{y}}{\partial z_1}}_{\text{yellow}} * \underbrace{\frac{\partial z_1}{\partial w_1}}_{\text{red}}$$

Computing Gradients: Backpropagation



$$\frac{\partial J(W)}{\partial w_1} = \underbrace{\frac{\partial J(W)}{\partial \hat{y}}}_{\text{blue}} * \underbrace{\frac{\partial \hat{y}}{\partial z_1}}_{\text{yellow}} * \underbrace{\frac{\partial z_1}{\partial w_1}}_{\text{red}}$$

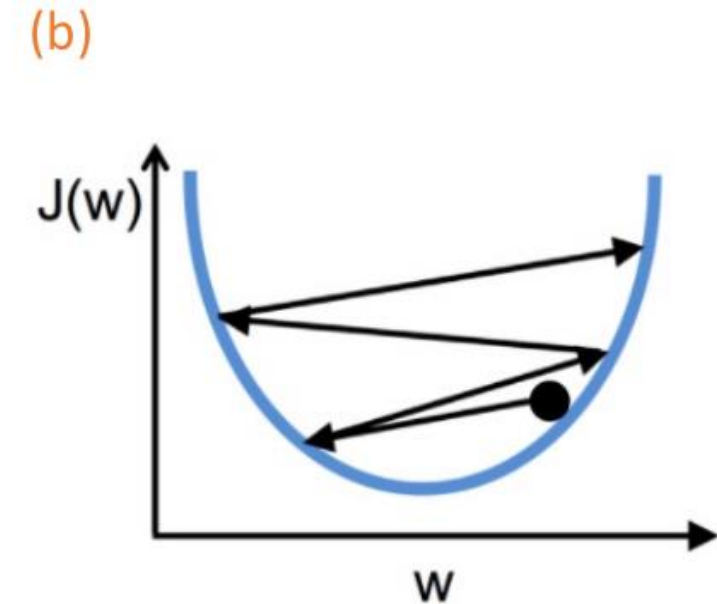
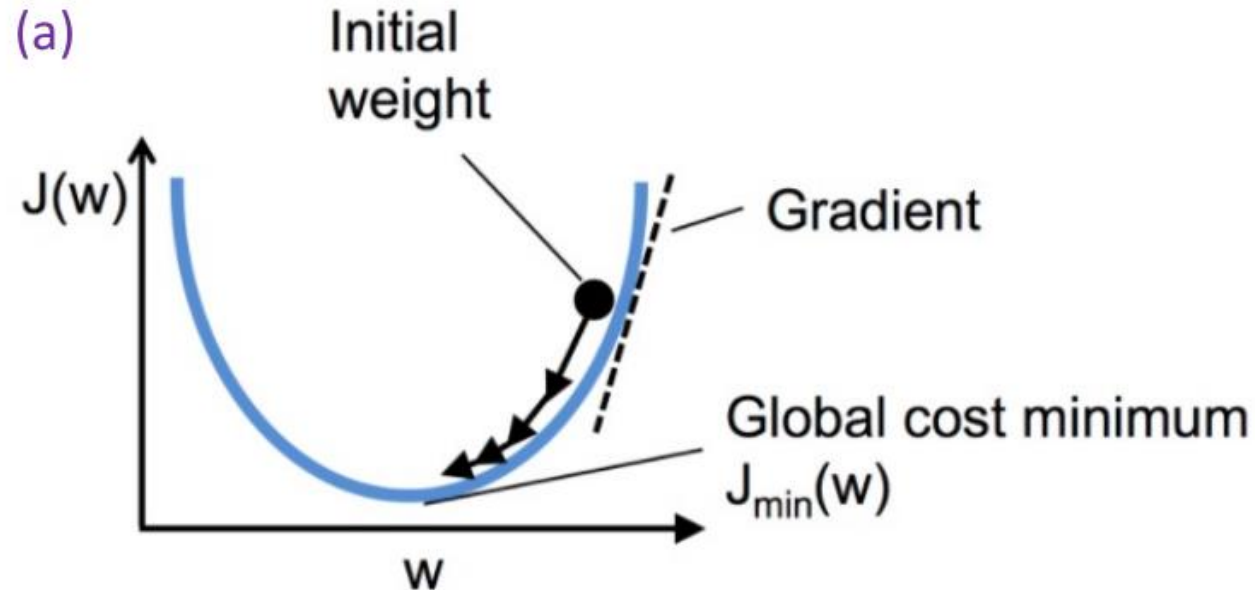
Repeat this for **every weight in the network** using gradients from later layers.

Gradient Descent: How Often to Update?

- **Batch Gradient Descent:**
 - Uses all training examples in each update
 - Less noisy updates but slow or infeasible for large datasets
- **Stochastic Gradient Descent (SGD):**
 - Updates with a single example per iteration
 - Fast and memory-efficient for large datasets, but updates can be noisy
- **Mini-batch Gradient Descent:**
 - Updates with a subset of examples per iteration
 - Reduces noise compared to SGD, scalable to large datasets, but can still be slow with large data

Gradient Descent: How Much to Update?

- Step size = learning rate
 - (a) When learning rate is too small, convergence to good solution will be slow!
 - (b) When learning rate is too large, convergence to a good solution is not possible!



Gradient Descent: How to Choose the Learning Rate?



Idea 1

Try lots of different learning rates and see what works better!

Gradient Descent: How to Choose the Learning Rate?



UNIVERSIDADE
CATOLICA
PORTUGUESA
BRAGA

Idea 3

Do something smarter!

Design an adaptive learning rate that "adapts" to the landscape!



Adaptive Learning Rates

- Learning rates are no longer fixed!
- Can be made larger or smaller depending on:
 - How large the gradient is.
 - How fast learning is happening.
 - The size of particular weights.
 - etc...