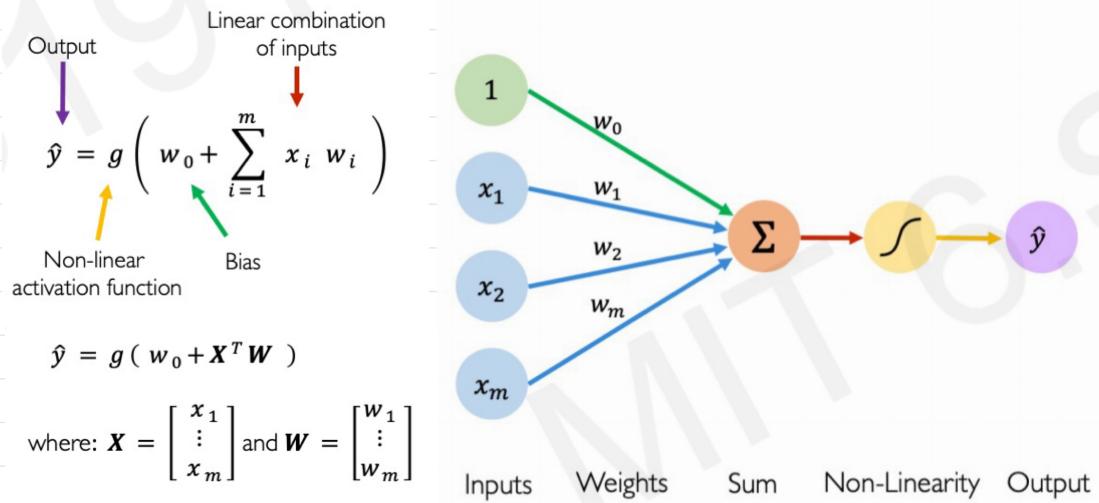



L1-INTRO To DL

BUILDING BLOCKS

PERCEPTRON

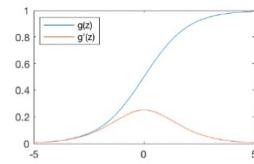
FORWARD PROPAG.



ACTIVATION FUNCTION

- Funzione (non lineare) che mappa gli output in un dato range

Sigmoid Function

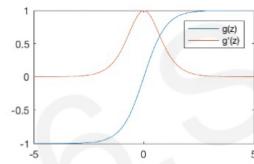


$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$

`tf.math.sigmoid(z)`

Hyperbolic Tangent

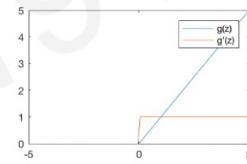


$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$

`tf.math.tanh(z)`

Rectified Linear Unit (ReLU)



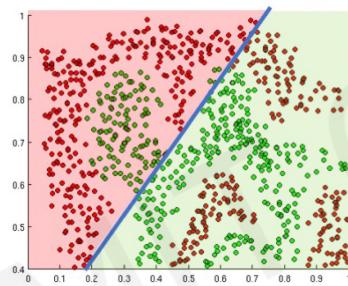
$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

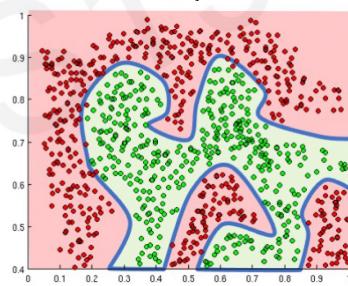
`tf.nn.relu(z)`

- La **NON-LINEARITY** ha lo scopo pratico di approssimare funzioni complesse

↳ Ad esempio per separare i punti colorati

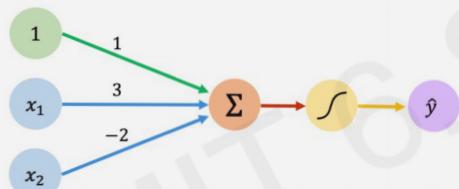


Linear activation functions produce linear decisions no matter the network size

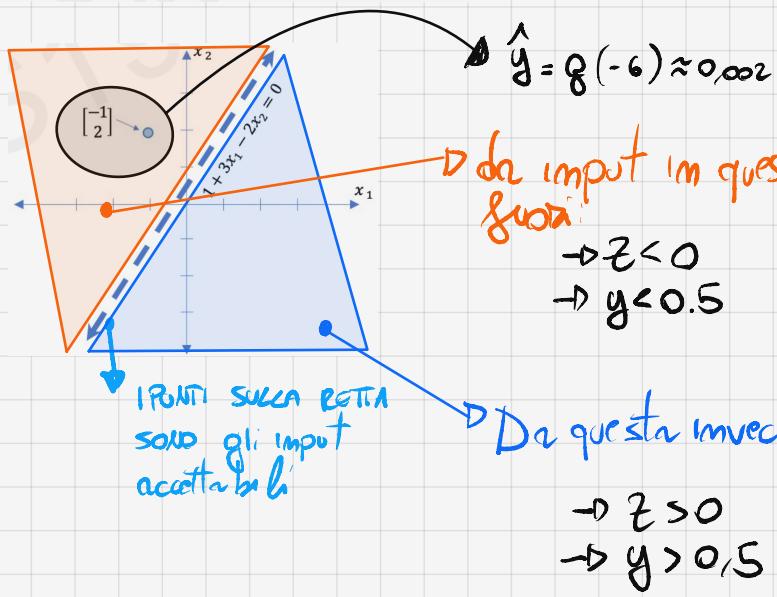


Non-linearities allow us to approximate arbitrarily complex functions

ESEMPIO



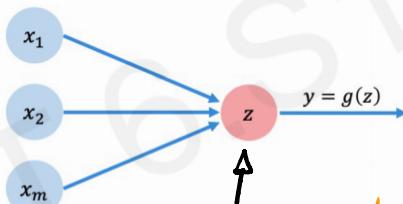
$$\hat{y} = g(1 + 3x_1 - 2x_2) \quad \text{2D}$$



From perceptron to Neural Net

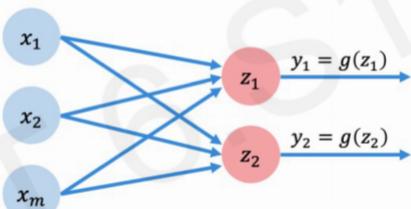
DENSE LAYER

• PERCEP. SEMPLIFICATE



Defn forward prop.

• MULTI-OUTPUT



↳ Detto DENSE LAYER perché ogni z_i è comune a ogni x_j

In TensorFlow

CLASS

```
class MyDenseLayer(tf.keras.layers.Layer):
    def __init__(self, input_dim, output_dim):
        super(MyDenseLayer, self).__init__()

        # Initialize weights and bias
        self.W = self.add_weight([input_dim, output_dim])
        self.b = self.add_weight([1, output_dim])

    def call(self, inputs):
        # Forward propagate the inputs
        z = tf.matmul(inputs, self.W) + self.b

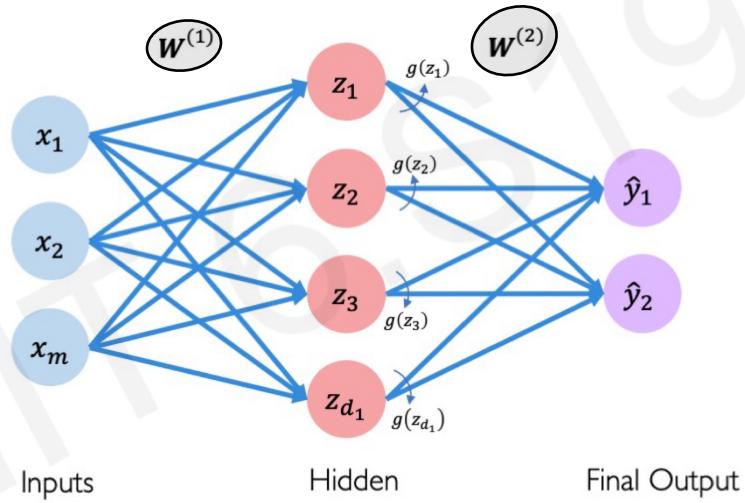
        # Feed through a non-linear activation
        output = tf.math.sigmoid(z)

        return output
```

INSTANCE

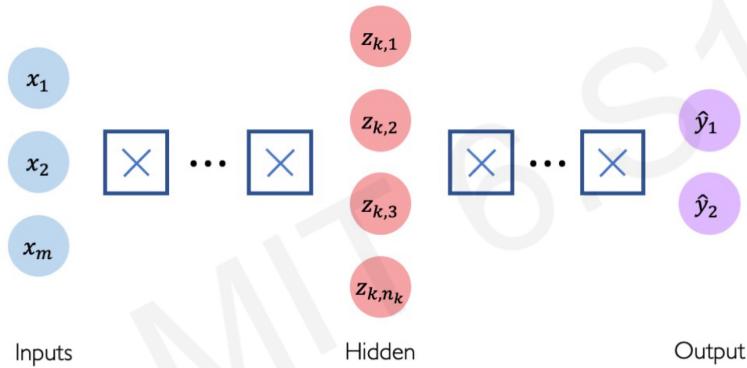
```
TF import tensorflow as tf
layer = tf.keras.layers.Dense(
    units=2)
```

SINGLE LAYER NN



- Ogni layer ha una propria WEIGHT MATRIX

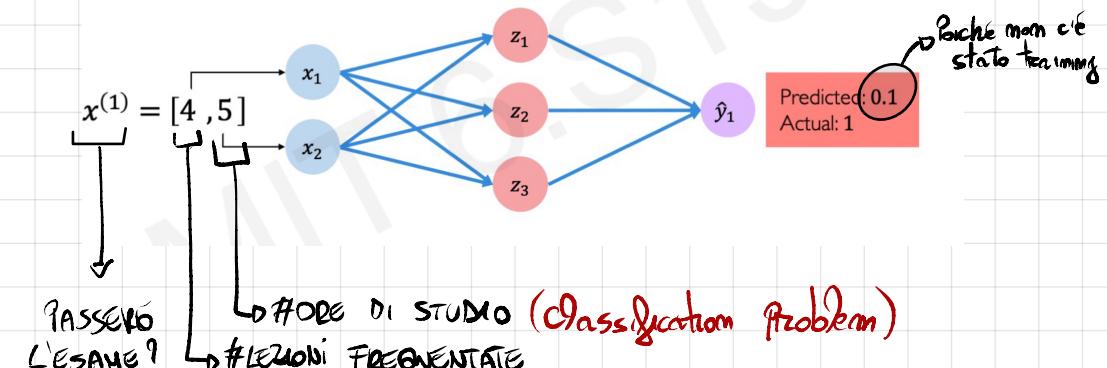
DEEP NN (multi-layer NN)



```
import tensorflow as tf

model = tf.keras.Sequential([
    tf.keras.layers.Dense(n1),
    tf.keras.layers.Dense(n2),
    ...
    tf.keras.layers.Dense(2)
])
```

ESEMPIO



Come diciamo al modello quanto sbaglia?

LOSS

$$\mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

Predicted Actual

$$x = \begin{bmatrix} 4, 5 \\ 2, 1 \\ 5, 8 \\ \vdots \end{bmatrix}$$

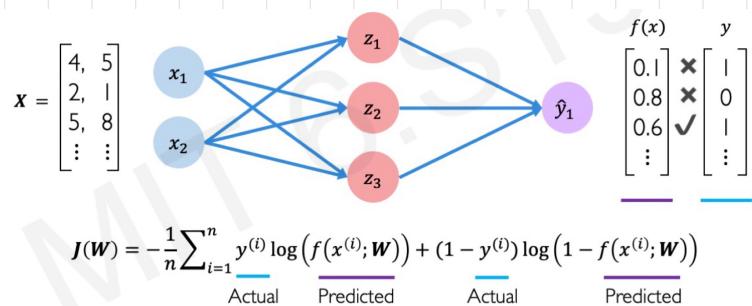
Also known as:
• Objective function
• Cost function
• Empirical Risk

$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

Predicted Actual

EMPIRICAL LOSS

BINARY CROSS EMPTY LOSS

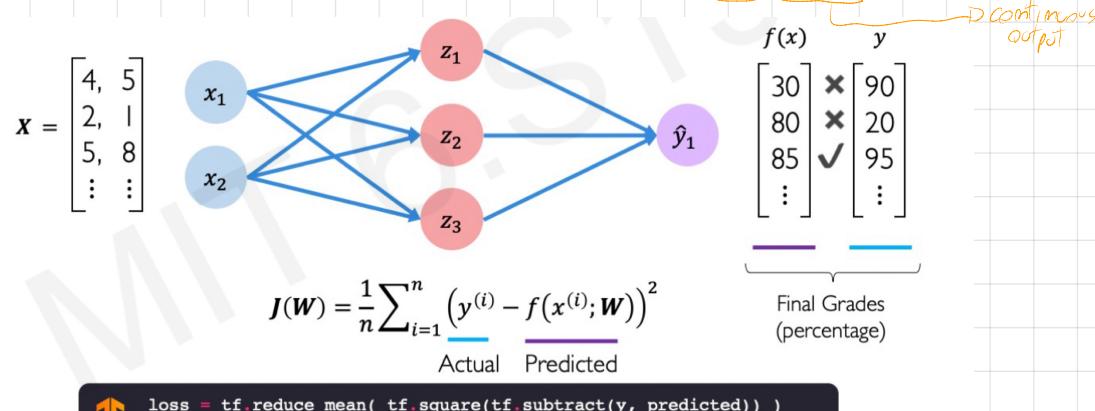


```
TensorFlow: loss = tf.reduce_mean(tf.nn.softmax_cross_entropy_with_logits(y, predicted))
```

→ confronta la distribuzione di probabilità (ATTESA vs OTTENUTA)

MSE

- Se invece di un BINARY OUTPUT volessimo predire il VOTO, ci serve un'altra metrica per l'errore (regression problem)



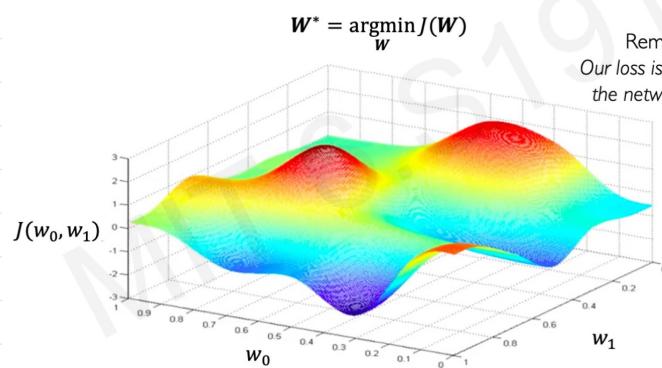
```
TensorFlow: loss = tf.reduce_mean(tf.square(tf.subtract(y, predicted)))  
loss = tf.keras.losses.MSE(y, predicted)
```

TRAINING

- Loss optimization: i pesi che minimizzano la LOSS

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

- Con un problema a due pesi, possiamo plottarne \mathbf{W}^*



es. GRADIENTE

- Scegliamo un punto nel grafico, calcoliamo i GRAD. e ci spostiamo nella sua direzione fino al LOCAL OPT.

GRADIENT DESCENT (ALGO)

NOTA se il learning rate è troppo piccolo restiamo bloccati in LOCAL OPT. se troppo grande si diffida la convergenza

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
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

LEARNING RATE
(Quanto vogliamo imporre
in diretta)

```
import tensorflow as tf

weights = tf.Variable([tf.random.normal()])

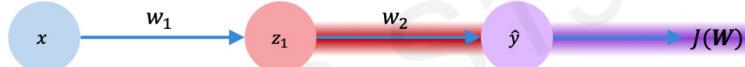
while True: # loop forever
    with tf.GradientTape() as g:
        loss = compute_loss(weights)
    gradient = g.gradient(loss, weights)

    weights = weights - lr * gradient
```

• Come calcoliamo il gradiente?

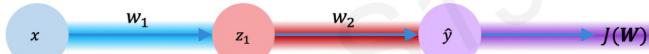
BACKPROPAGATION

- Ci interessa quanto un **PICCOLO CAMBIAMENTO** in un pes w_i influenzi $J(W)$
- Può essere espresso come rapporto tra derivate $\frac{\partial J(W)}{\partial w_2}$ e quale applichiamo la **CHAIN RULE**



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

• Possiamo espanderci a $\frac{\partial J(W)}{\partial w_1}$



$$\frac{\partial J(W)}{\partial w_1} = \frac{\partial J(W)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1} * \frac{\partial z_1}{\partial w_1}$$

- La chain rule si applica a tutti i modi ricorsivamente in NN con più layer

IN Practice

Come sceglio il LEARNING RATE?

ADAPTIVE LEARN. RATE

- Il GRADIENT DESCENT non funge con molti local opt.

- Il LEARNING RATE va settato per bene

↗ GRAD. DESC. ALGO

- Il LEARN. RATE non è fisso, ma è ottimizzato durante

- IL LEARNING PROCESS

Algorithm	TF Implementation
• SGD	tf.keras.optimizers.SGD
• Adam	tf.keras.optimizers.Adam
• Adadelta	tf.keras.optimizers.Adadelta
• Adagrad	tf.keras.optimizers.Adagrad
• RMSProp	tf.keras.optimizers.RMSProp

```

import tensorflow as tf
model = tf.keras.Sequential([...])
# pick your favorite optimizer
optimizer = tf.keras.optimizer.SGD()
while True: # loop forever
    # forward pass through the network
    prediction = model(x)
    with tf.GradientTape() as tape:
        # compute the loss
        loss = compute_loss(y, prediction)
    # update the weights using the gradient
    grads = tape.gradient(loss, model.trainable_variables)
    optimizer.apply_gradients(zip(grads, model.trainable_variables))

```



Can replace with any TensorFlow optimizer

LA BACKPROP è
possibile?

IDEA

- Impossibile com la BACKPROPAGATION è troppo da calcolare, soprattutto se lo facciamo per tutto il dataset.
- Lo computiamo per un sottoinsieme di punti

STOCHASTIC GRADIENT DESCENT (SGD)

- Ci permette di aumentare la learning rate

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick single data point i
4. Compute gradient $\frac{\partial J_i(\mathbf{W})}{\partial \mathbf{W}}$
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights

com un singolo punto è troppo noisy

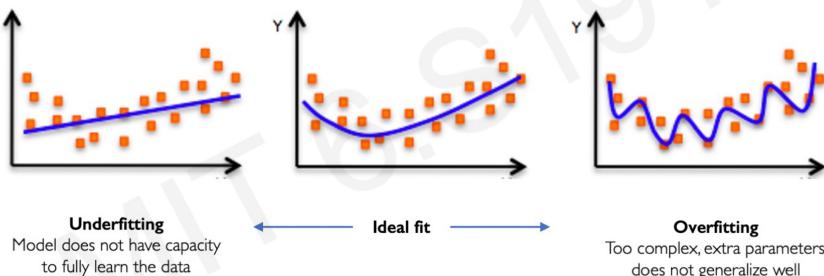
Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick batch of B data points
4. Compute gradient $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(\mathbf{W})}{\partial \mathbf{W}}$
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights

così va meglio

FITTING MODEL

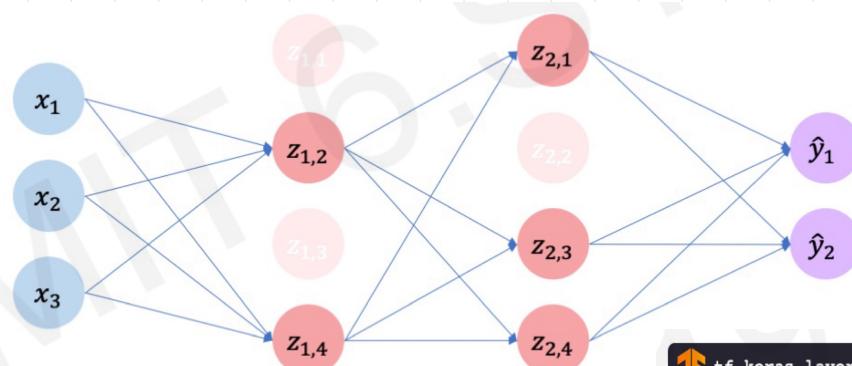
- Come rappresenta i dati il modello



- La linea blu è il modello usato per rappresentare i dati

REGULARIZATION DROPOUT

- Tecnica che impedisce al modello di andare in **OVERTFITTING**
- Durante il training settiamo alcune activation a \emptyset .

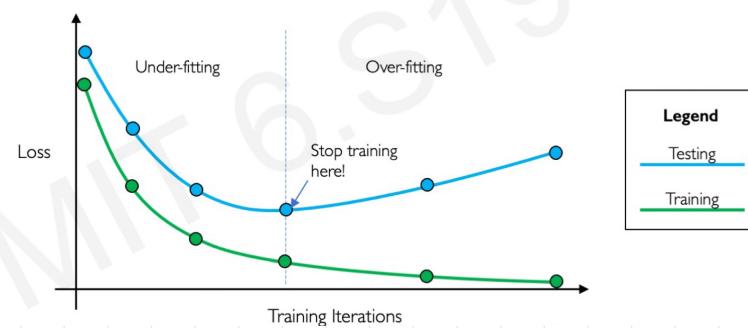


Lo indotte aumenta la velocità di train, poiché deve lavorare su meno modi.

EARLY STOPPING

- Smette il training quando risulta troppo loss

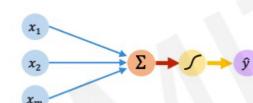
Stop training before we have a chance to overfit



COSA ABBIAMO FATTO?

The Perceptron

- Structural building blocks
- Nonlinear activation functions



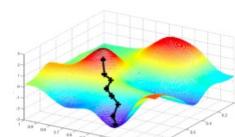
Neural Networks

- Stacking Perceptrons to form neural networks
- Optimization through backpropagation



Training in Practice

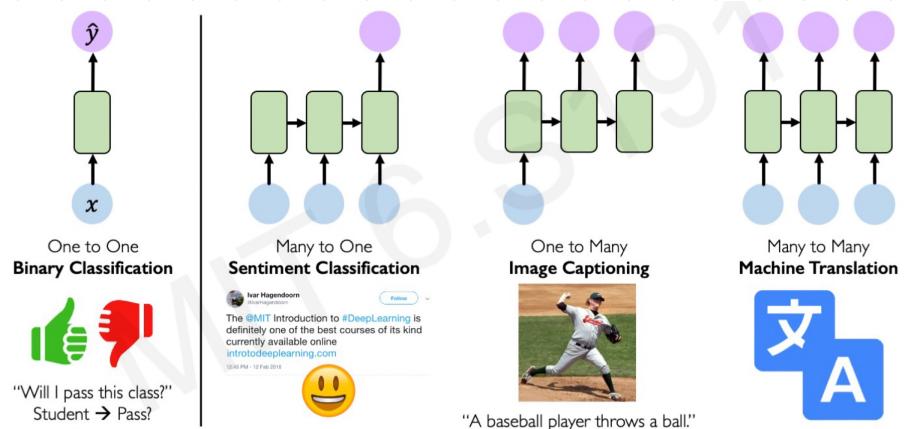
- Adaptive learning
- Batching
- Regularization



DEEP SEQUENCE MODELLING

- Lavoreremo su dati sequenziali nel tempo

Sequence modeling
Applications



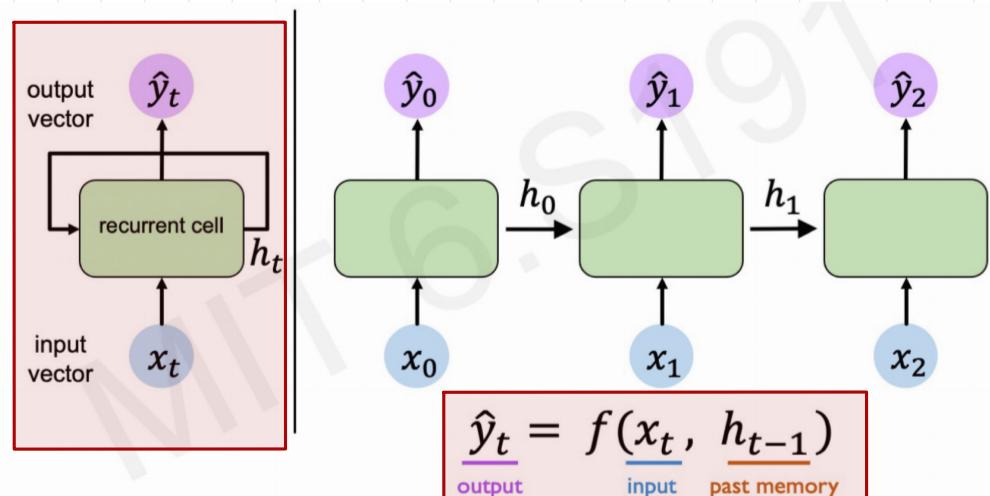
NEURONS WITH RECURRENCE

IDEA:

- i dati sequenziali sono **TIME STEPS** individuali, ci serve un modello che sia in grado di gestirli individualmente, ma anche prendendo in considerazione anche gli stati passati

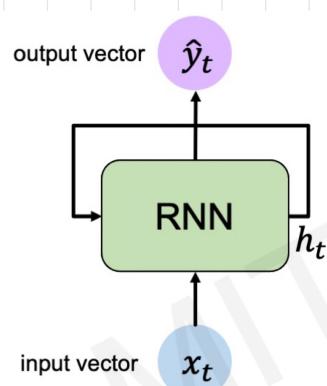
NEURONI CON
RICORRENZA

h_t sta per HIDDEN STATE



RECURRENT NEURAL NETWORKS (RNN)

RECURRENCE
RELATION



- Ad ogni time step calcoliamo l'hidden state h_t associato

$$h_t = f_w(x_t, h_{t-1})$$

cell state function with weights w
input old state

```

my_rnn = RNN()
hidden_state = [0, 0, 0, 0]

sentence = ["I", "love", "recurrent", "neural"]

for word in sentence:
    prediction, hidden_state = my_rnn(word, hidden_state)

next_word_prediction = prediction
# >>> "networks!"

```

Passi del processo

NOTA: a somo 3
WEIGHT MATRICES

(1) INPUT VECTOR x_t

(2) UPDATE di h_t

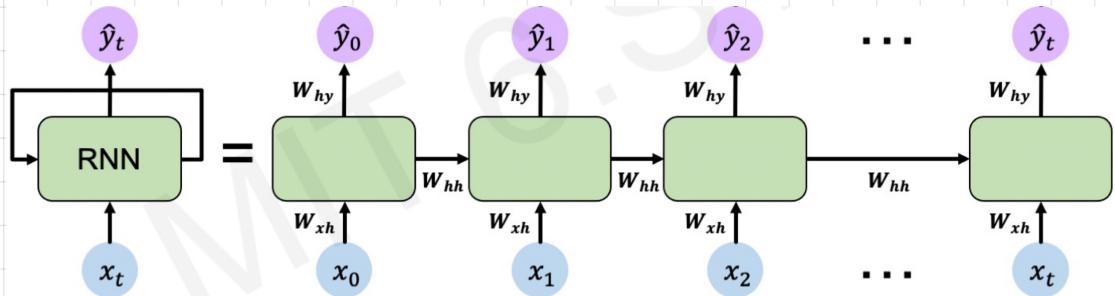
- moltiplichiamo h_{t-1} e x_t per le rispettive WEIGHT MATRICES

$$h_t = \tanh(W_{hh}^T h_{t-1} + W_{xh}^T x_t)$$

(3) CALCOLIAMO \hat{y}_t moltiplicando
il hidden state per un'altra
WEIGHT MATRIX

$$\hat{y}_t = W_{hy}^T h_t$$

... graficamente



OBSs

- Ad ogni output è associata una Loss, i cui insieme ci permette di allenare il modello.
- È importante da notare che le weight matrices sono sempre le stesse a ogni step.

oooPython

```
class MyRNNCell(tf.keras.layers.Layer):
    def __init__(self, rnn_units, input_dim, output_dim):
        super(MyRNNCell, self).__init__()

        # Initialize weight matrices
        self.W_xh = self.add_weight([rnn_units, input_dim])
        self.W_hh = self.add_weight([rnn_units, rnn_units])
        self.W_hy = self.add_weight([output_dim, rnn_units])

        # Initialize hidden state to zeros
        self.h = tf.zeros([rnn_units, 1])

    def call(self, x):
        # Update the hidden state
        self.h = tf.math.tanh( self.W_hh * self.h + self.W_xh * x )

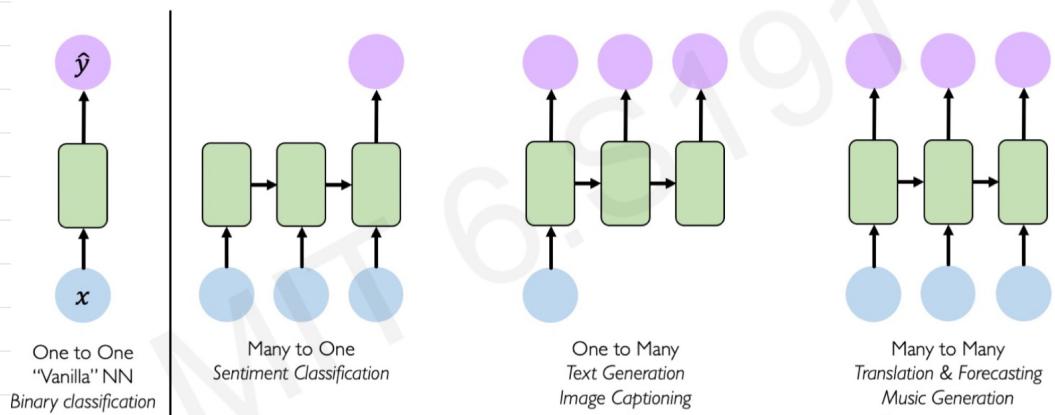
        # Compute the output
        output = self.W_hy * self.h

        # Return the current output and hidden state
        return output, self.h
```

- Ovviamente ci pensa tensorflow, non dobbiamo nulla implementarlo

```
tf.keras.layers.SimpleRNN(rnn_units)
```

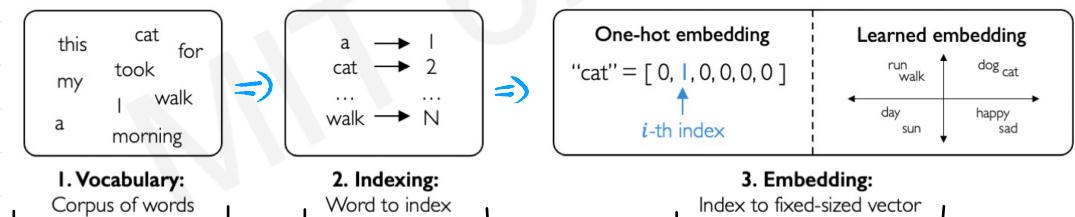
RNNs per Sequential modeling



ENCODING ISSUE



SOL



↳ CREA UN VOCABOLARIO DI TUTTE LE PAROLE

↳ INDIVIDUA IL VOCABOLARIO

↳ ONE-HOT ENCODING

↳ LEARNED ENCODING: crea gruppi di parole semanticamente simili

DESIGN CRITERIA

• 1) Modello deve avere l'abilità di:

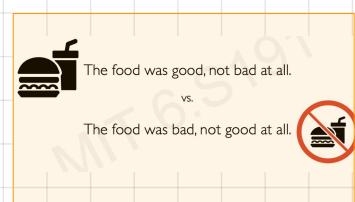
(1) Gestire seq. di **LUNGHEZZA VARIABILE**



(2) Trascurare **LONG-TERM DEPENDENCIES**?

"France is where I grew up, but I now live in Boston. I speak fluent ____."

(3) Mantenere informazioni sull'**ORDINE** → delle sequenze

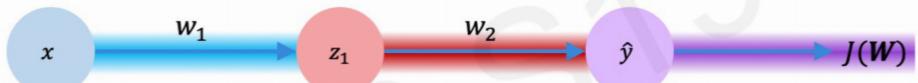


(4) CONDIVIDERE parametri tra le sequenze

BACKPROPAGATION THROUGH TIME (BPTT)

RIC

- Abbiamo visto che in BACKPROPAGATION consiste nel calcolare i GRADIENTI DELLA LOSS TOTAL ($\partial J(W)$) rispetto ogni parametro ($\frac{\partial J(W)}{\partial w_1}$), e poi modificare i parametri per diminuire la loss.

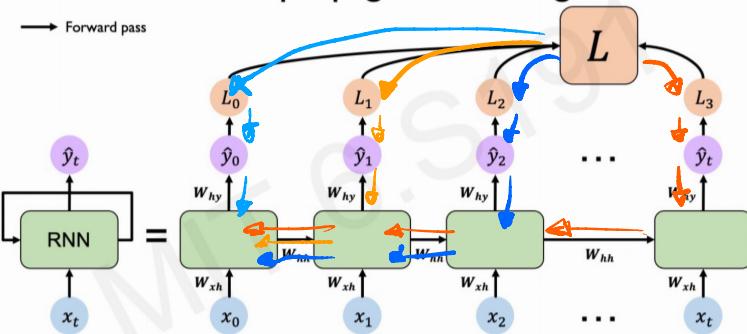


$$\frac{\partial J(W)}{\partial w_1} = \frac{\partial J(W)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1} * \frac{\partial z_1}{\partial w_1}$$

NON GRADIENT FLOW

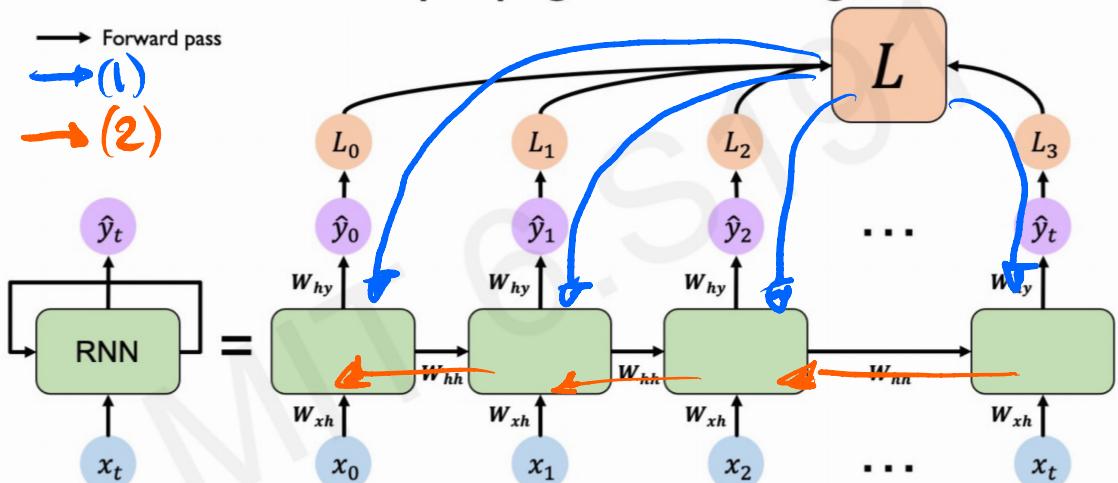
NO !

- Dala total loss invece di computare la Backprop. per ogni TIME STEP



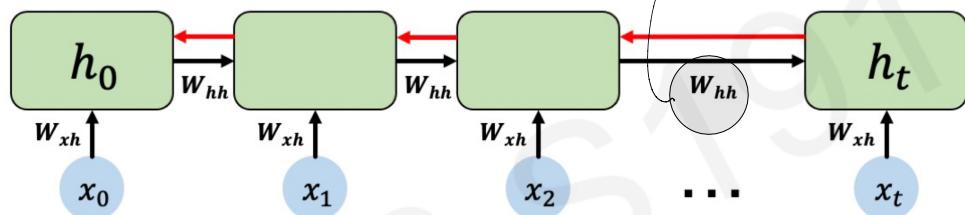
⇒ conviene farze un solo viaggio, dal TIME STEP più recente.
In questo modo l'errore ce lo potremo in una sola direzione.

SI !



- Calcolare il GRADIENTE rispetto h_0 , comporre:

- Molti moltiplicazioni per W_{hh}
- Computare molti gradienti



Exploding Gradient

Vanishing Gradient

- Se abbiamo molti valori > 1 in questo modo, i gradienti crescono ed overze moltiplicazioni per il GRADIENTE
- se abbiamo invece valori < 1 i GRADIENTI diminuiscono troppo

Why are vanishing gradients a problem?

Multiply many small numbers together

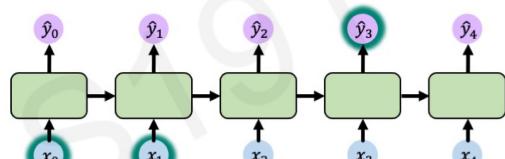


Errors due to further back time steps have smaller and smaller gradients

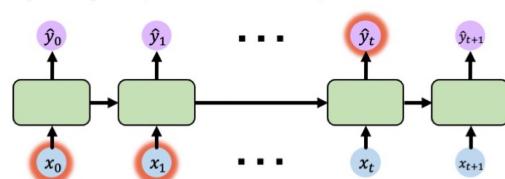


Bias parameters to capture short-term dependencies

"The clouds are in the ___"



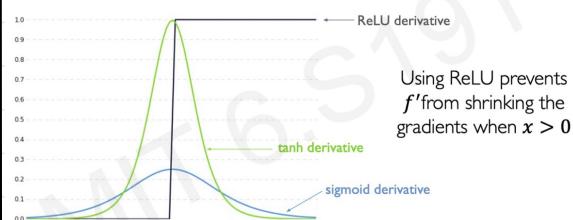
"I grew up in France, ... and I speak fluent ___"



SOL

- Ce sono vari trick per gestire i VANISHING GRADIENTS

Trick #1: Activation Functions



Using ReLU prevents f' from shrinking the gradients when $x > 0$

Trick #2: Parameter Initialization

Initialize **weights** to identity matrix

$$I_n = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{pmatrix}$$

Initialize **biases** to zero

This helps prevent the weights from shrinking to zero.

Solution #3: Gated Cells

Idea: use a more complex recurrent unit with gates to control what information is passed through

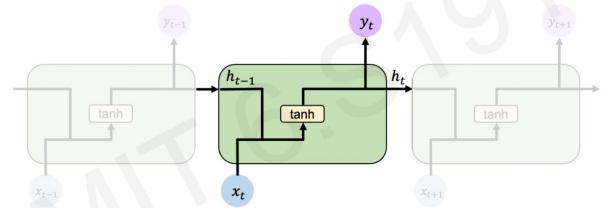
gated cell

LSTM, GRU, etc.

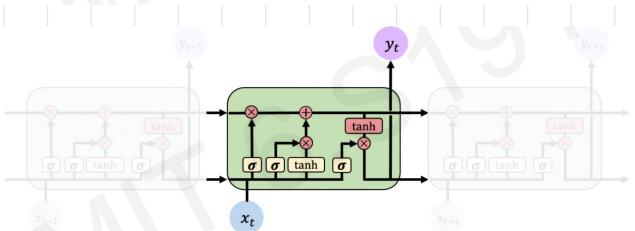
Long Short Term Memory (LSTMs) networks rely on a gated cell to track information throughout many time steps.

LONG SHORT TERM MEMORY

- STANDARD RNN



- LSTM RNN



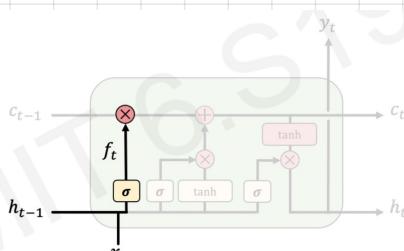
`tf.keras.layers.LSTM(num_units)`

- I GATES decidono quale informazione far passare o meno
- Lo fanno tramite, ad esempio, **SIGMOIDI** (come nelle imgs)

CELL STATE
(TIPPI) \rightarrow

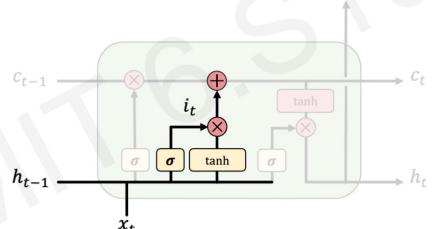
- Il parametro **C** è una variabile contenente informazioni passate
- Si divide in 4 fasi

(1)

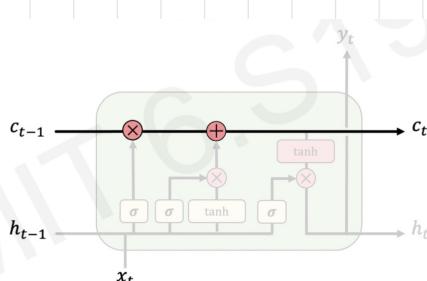


\rightarrow **FORGET**: butta parti irrilevanti dello stato precedente

(2) **STORE**: salva le info irrilevanti sullo stato prossimo

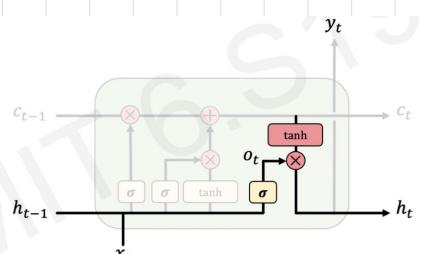


(3)

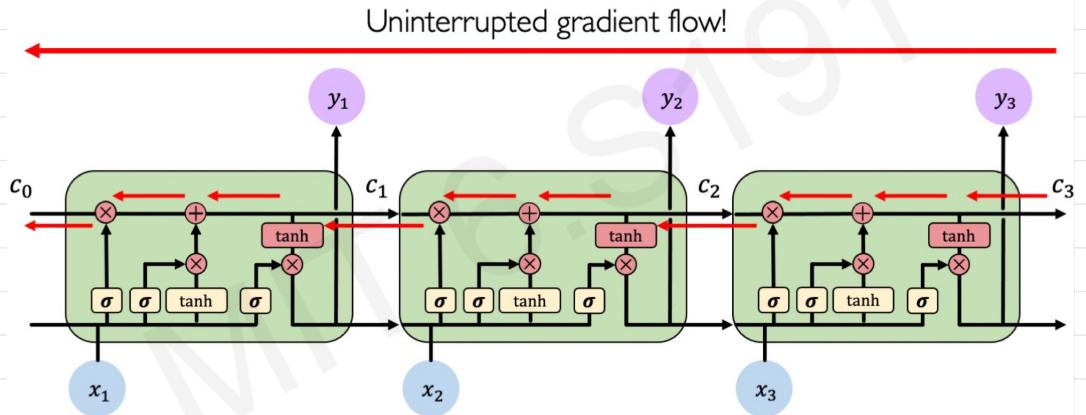


\rightarrow **UPDATE**: aggiorna il valore di C_t

(4) **OUTPUT**: controlla q'info \rightarrow invia alla stato successivo



• i9 GRADIENTE E ora calcolato attraverso i9 **CELL GATES** e
elimina i9 **VANISHING GRADIENT PROBLEM**



Concetti di LSTM

1. Maintain a **separate cell state** from what is outputted
2. Use **gates** to control the **flow of information**
 - **Forget** gate gets rid of irrelevant information
 - **Store** relevant information from current input
 - Selectively **update** cell state
 - **Output** gate returns a filtered version of the cell state
3. Backpropagation through time with **uninterrupted gradient flow**

Cosa abbiamo fatto?

Deep Learning for Sequence Modeling: Summary

1. RNNs are well suited for **sequence modeling** tasks
2. Model sequences via a **recurrence relation**
3. Training RNNs with **backpropagation through time**
4. Gated cells like **LSTMs** let us model **long-term dependencies**
5. Models for **music generation**, classification, machine translation, and more