

Computer Vision

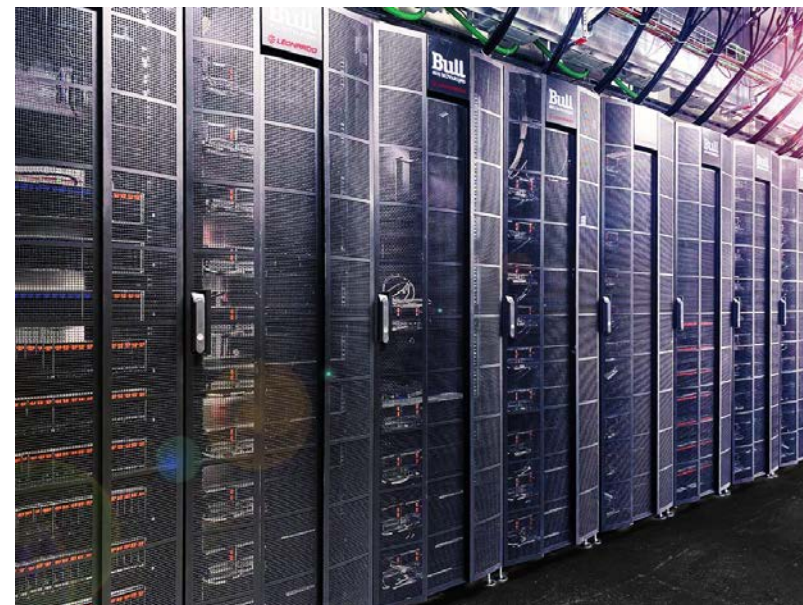
IA per le industrie: corso avanzato

Rosario Di Carlo

PRESENTAZIONE

Rosario Di Carlo

AI Research Fellow @LeonardoLabs, AI-BigData-HPC Lab



Cosa è la Computer Vision?

- **Computer Vision** è il campo di ricerca che include metodologie per l'acquisizione, il processamento, l'analisi e la comprensione delle immagini.
- Strettamente legata alla teoria per costruire sistemi artificiali che ottengono informazioni dalle immagini.
- I dati delle immagini possono assumere molte forme, come una sequenza video, immagini di profondità, viste da più telecamere o dati multidimensionali da uno scanner medico

Introduzione alla Computer Vision

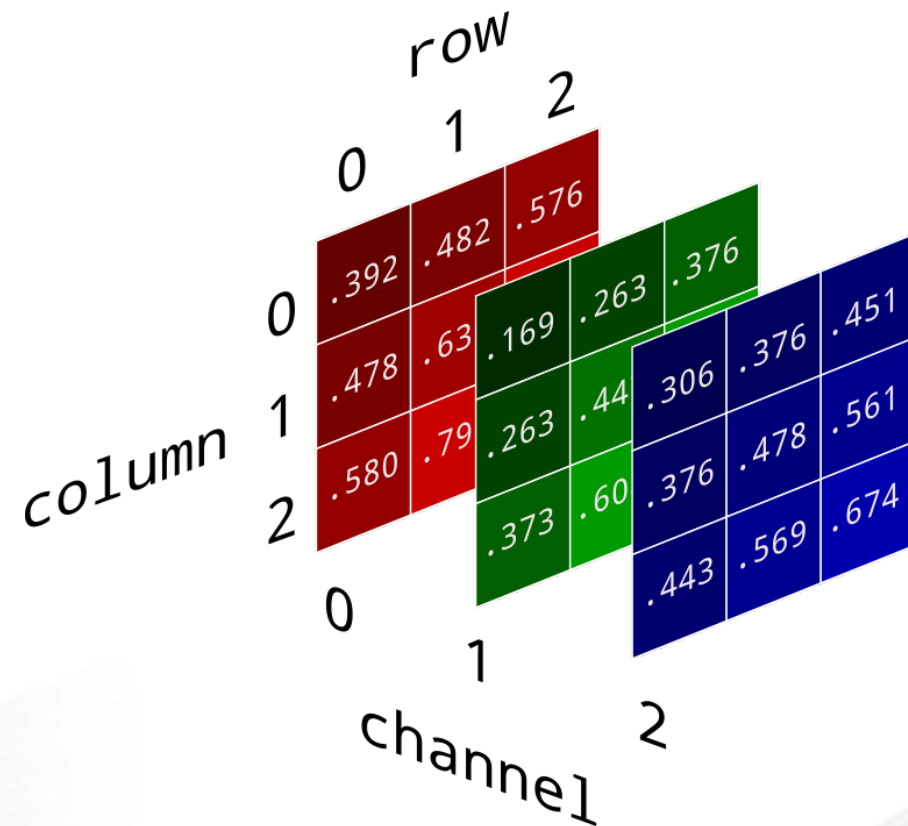
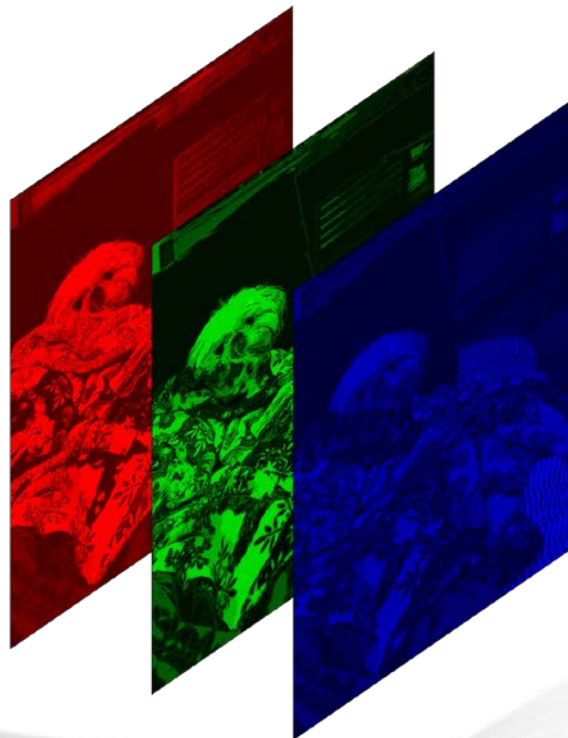
- La prima cosa importante da capire è ciò che un computer vede quando guarda un'immagine: non un contenuto visivo reale, ma un array di valori numerici di pixel.



62	62	63	64	65	66	67	67	69	70	71	72	72	73	73	73	72	71	70	69	67	66	66	65	63	62	61	60	6				
61	62	63	64	66	66	67	68	68	69	70	71	71	72	73	72	72	71	71	70	69	68	66	65	65	63	62	61	60	6			
61	62	63	64	66	66	68	68	69	70	70	71	72	73	73	72	72	71	71	69	68	67	66	65	65	64	63	62	61	6			
61	63	64	64	66	67	68	68	68	69	70	71	71	73	74	73	73	73	71	70	69	68	66	66	65	64	63	62	61	61	6		
61	63	64	65	67	68	69	69	70	70	71	71	72	55	53	69	72	71	71	70	69	68	67	66	65	64	63	62	60	60	6		
63	64	65	66	67	68	69	69	70	70	71	72	42	4	5	11	48	72	71	71	69	68	67	66	65	64	62	62	60	59	5		
63	65	66	66	68	68	69	70	71	71	71	72	18	4	4	7	8	66	71	70	69	68	68	67	66	65	64	63	61	59	59	5	
63	65	67	67	68	69	69	70	71	71	72	64	4	27	24	54	33	29	52	64	68	68	67	66	65	64	63	62	61	59	58	5	
64	65	66	66	68	69	70	71	41	24	24	12	17	24	48	60	37	43	38	52	66	68	67	66	65	64	63	61	60	59	58	5	
65	66	67	67	68	69	71	40	6	6	6	5	34	36	12	47	34	17	29	54	43	63	67	66	65	64	63	62	60	59	58	5	
64	65	66	66	68	69	38	6	6	5	5	7	16	19	4	47	44	27	24	40	67	66	66	65	65	64	63	61	60	59	58	5	
63	64	65	65	67	30	6	6	5	5	5	6	8	9	20	27	51	78	41	44	66	65	65	65	65	64	63	62	60	59	58	5	
63	64	65	65	34	5	5	5	5	5	5	5	4	19	6	7	54	64	20	59	65	65	64	64	64	63	62	61	60	59	57	5	
63	64	64	65	14	5	6	5	5	4	5	4	18	7	5	4	19	10	11	65	64	64	64	63	61	66	62	61	60	59	58	5	
63	64	64	65	53	7	4	5	6	7	10	6	5	5	4	21	24	18	64	64	64	63	62	64	65	62	62	60	59	58	5		
64	64	64	64	65	50	4	4	4	5	11	16	6	6	4	6	35	16	26	66	64	64	63	61	72	67	63	62	61	59	58	5	
64	64	64	64	65	46	4	4	4	5	6	9	8	5	20	10	43	56	20	57	64	64	63	61	70	67	62	64	65	59	59	5	
64	64	64	65	66	27	5	4	4	5	6	6	6	18	66	20	57	60	46	36	75	70	62	61	70	67	62	61	60	59	58	5	
49	50	62	65	57	5	5	6	5	6	6	6	6	41	59	28	60	58	44	22	63	71	72	60	69	68	61	60	58	59	59	5	
42	52	57	52	28	5	5	5	5	5	5	5	5	5	70	50	43	61	62	64	39	42	64	60	62	56	63	65	65	67	61	53	5
32	32	32	33	6	5	5	5	5	5	6	11	39	21	33	51	50	45	46	18	32	36	33	23	44	70	71	51	42	27	3		
50	50	51	39	5	5	5	5	6	5	6	6	42	69	28	34	42	39	43	37	26	29	40	26	29	26	35	42	35	33	18	1	
52	53	51	22	5	5	5	5	6	5	6	5	44	56	17	51	54	53	54	56	51	22	54	54	55	55	54	53	53	53	52	5	
54	54	53	8	5	5	5	5	6	5	6	13	52	42	21	51	54	51	49	49	50	22	41	45	42	42	41	40	41	44	43	4	
52	52	54	38	8	5	5	6	6	5	6	28	55	32	32	54	53	51	51	51	44	25	51	51	49	49	50	49	48	46	4		
54	54	52	53	30	7	5	6	6	5	6	40	54	29	52	51	53	56	55	52	52	51	38	52	52	50	49	46	46	45	46	4	
51	52	51	53	27	14	5	4	5	4	7	47	51	21	39	49	47	49	52	52	49	35	31	48	46	47	47	47	46	46	4		
48	50	51	53	25	14	17	8	4	4	17	46	40	18	43	47	46	49	52	54	53	53	54	18	50	49	46	47	47	47	4		
49	49	49	49	22	12	20	24	6	14	35	51	39	48	48	50	51	51	49	51	51	52	50	41	58	48	47	47	47	45	4		
51	49	50	50	22	13	19	36	13	12	42	50	40	73	50	50	50	49	48	49	49	48	49	45	51	46	44	44	44	42	4		
47	49	49	47	20	16	26	39	21	15	36	48	42	61	47	48	51	47	50	51	51	51	49	47	47	52	47	47	44	43	4		

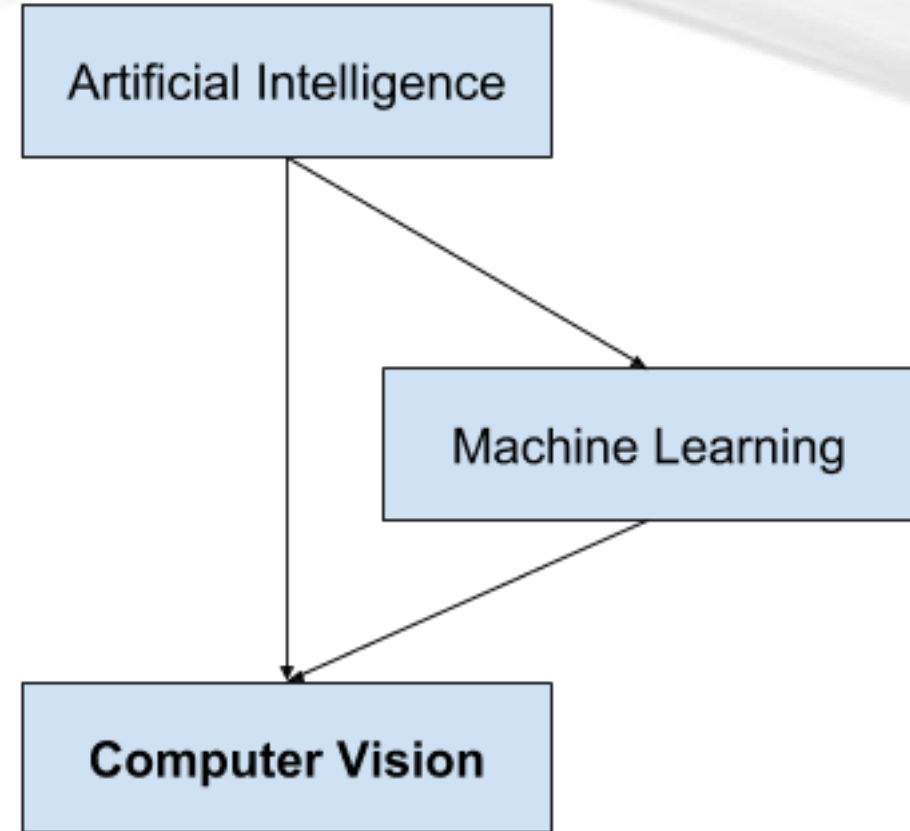
Immagini RGB

- Le immagini a colori sono rappresentate come array tridimensionali, una collezione di tre matrici, uno per ciascuno dei canali rosso, verde e blu. Ognuno di essi, ha un valore per pixel e i loro intervalli sono identici $[0, 255]$



Computer vision e AI

- La Computer Vision è un campo multidisciplinare che potrebbe essere definito in generale un sottocampo dell'intelligenza artificiale e del machine learning, che può comportare l'uso di metodi specifici e l'uso di algoritmi di apprendimento generali.

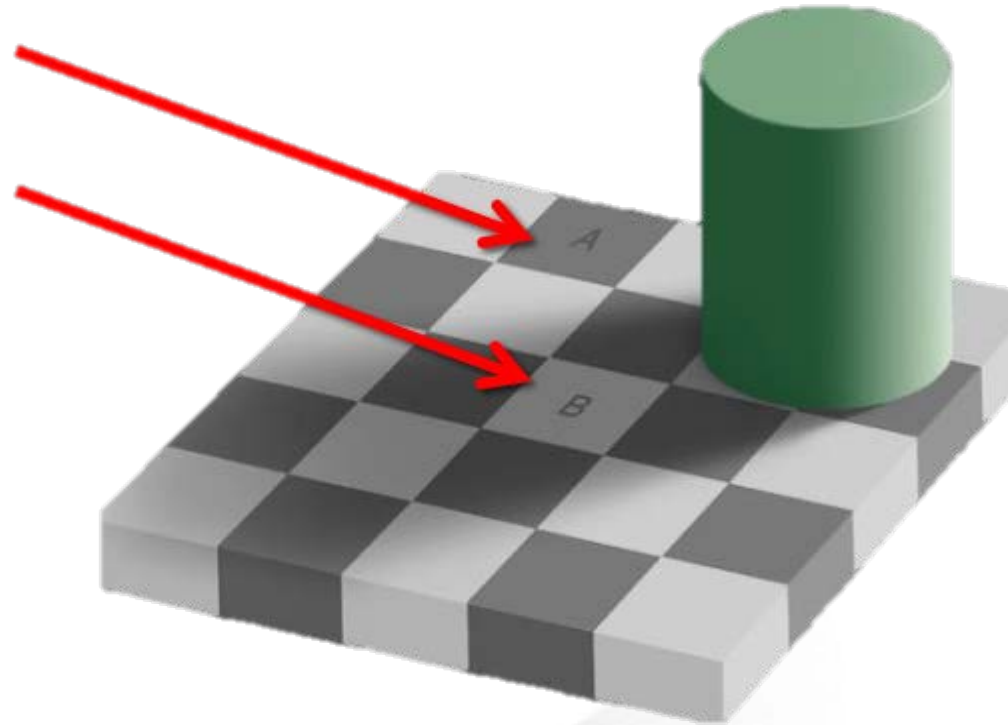


Perché è così difficile?

- La visione non è soltanto image processing

Casella scura

Casella chiara

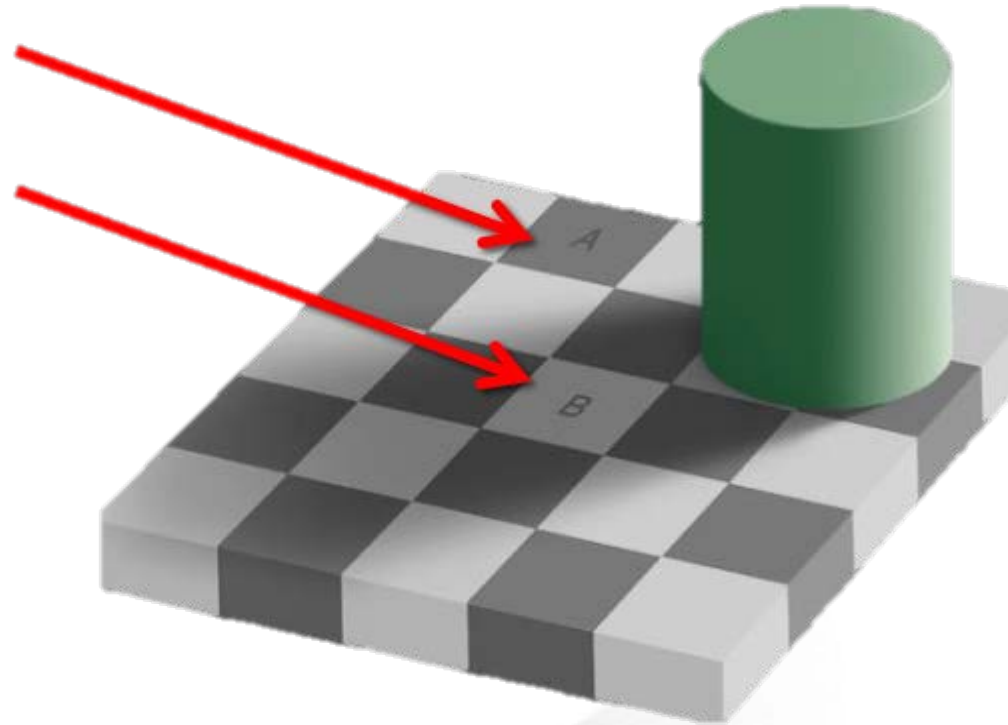


Perché è così difficile?

- La visione non è soltanto image processing

Casella scura

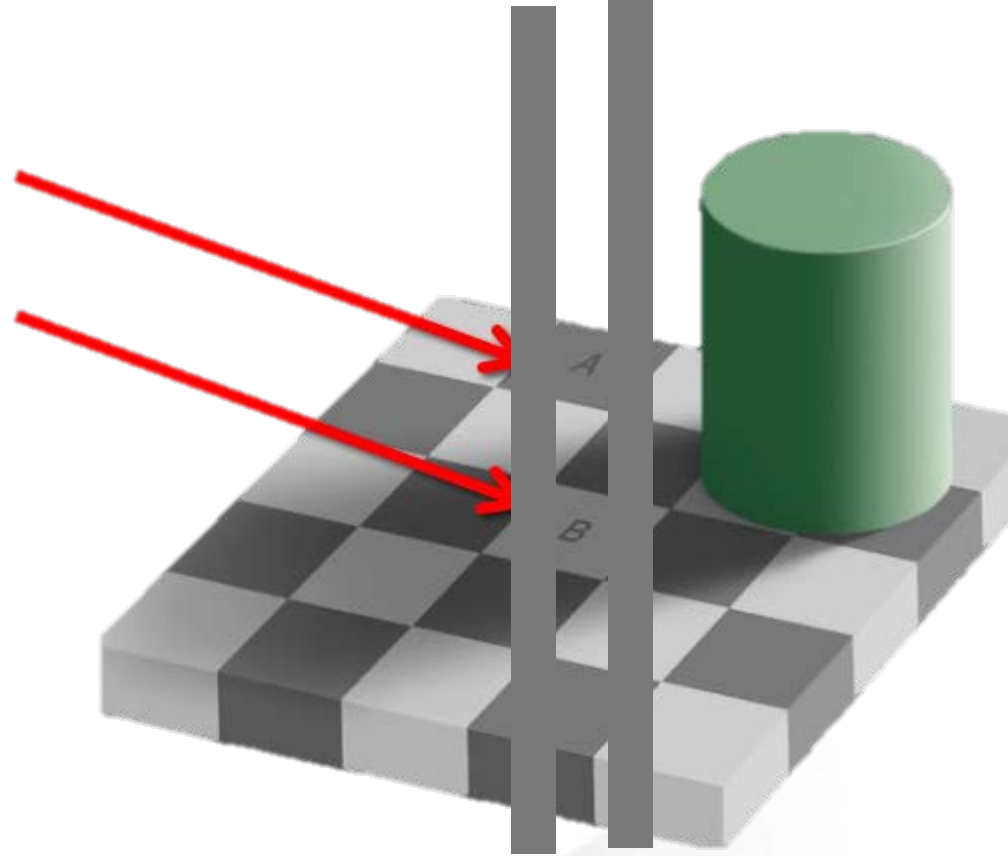
Casella chiara



Perché è così difficile?

- La visione non è soltanto image processing

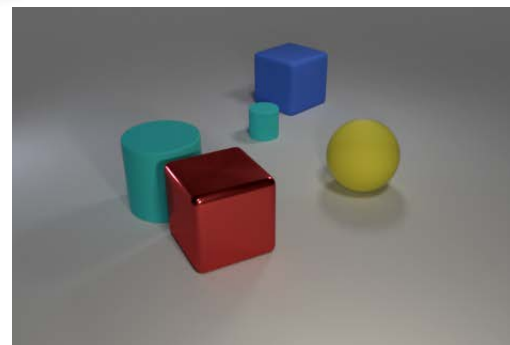
Casella scura
Casella chiara



Perché è così difficile?

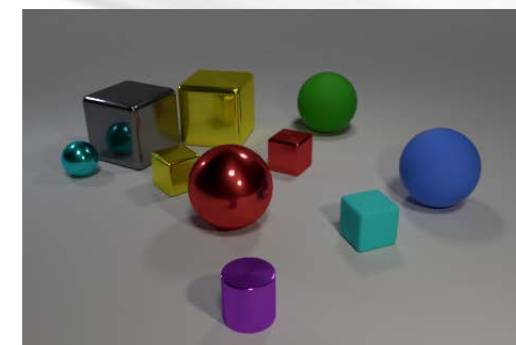
Nell'esempio precedente, i due quadrati hanno esattamente la stessa misura di intensità. Quindi, vedere non equivale a misurare le proprietà dell'immagine.

Piuttosto, "vedere" è costruire una percezione di ciò che è nel mondo sulla base delle misurazioni effettuate da un sensore di immagini.



Q: What **shape** is the **red thing**?

A: Cube



Q: There is a **large blue** ball that is on the **right side** of the **rubber sphere** that is **behind** the **tiny cyan object** on the **left side** of the **large green matte thing**; what is its **color**?

A: Blue

Il task di **Visual Reasoning** prova ad allenare un modello a comprendere il contenuto dell'immagine

Classificazione

- Quando noi umani vediamo e classifichiamo qualcosa, il nostro cervello dà un senso ad esso etichettando, prevedendo e riconoscendo pattern specifici. Allo stesso modo, i classificatori di immagini si basano su strutture che prendono ispirazione dalla biologia del cervello: le cosiddette **reti neurali convoluzionali** (CNN),
- Quello che vogliamo in un task di classificazione di immagini è prendere la matrice di numeri che costituisce un'immagine e produrre la probabilità che l'immagine sia di una certa classe (ad esempio 90% per il cane, 10% per il gatto).

Convoluzione

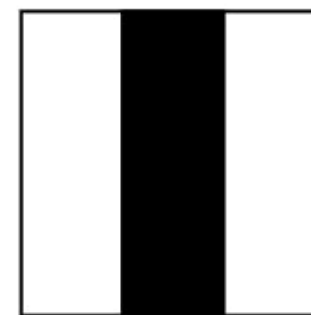
I layer di convoluzione sono i principali elementi costitutivi dei classificatori di immagini. Il termine matematico convoluzione si riferisce alla combinazione di due funzioni (f e g) che producono una terza funzione (z).

Un layer di convoluzione prende un input, applica un filtro (**kernel**) e produce una **feature map**. La feature map (z) è una combinazione di input e filtro (f e g), da cui il nome layer di convoluzione. L'obiettivo della convoluzione è di estrarre delle caratteristiche significative di un'immagine.

Una feature è una caratteristica specifica dell'immagine originale, come punti, bordi o la forma del naso di un cane. Similmente all'immagine che viene elaborata come un array di numeri, una feature si traduce in una casella di valori numerici dei pixel. Questa matrice serve come rilevatore di feature.

Example: Vertical line detection

Feature



Filter matrix

0	1	0
0	1	0
0	1	0

Convoluzione

- Per scansionare l'immagine, questo kernel si muove sull'immagine, blocco di pixel per blocco di pixel. Per ogni sottoregione, viene calcolato un valore basato su quanto è buona la corrispondenza tra il filtro e l'immagine.

Convolution: Image matrix * Filter matrix = Feature map

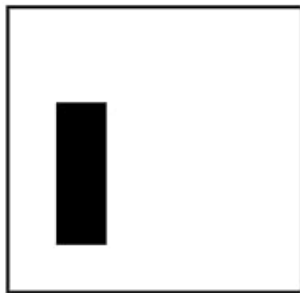


Image (6x6)

0	0	0	0	0	0
0	0	0	0	0	0
0	1	0	0	0	0
0	1	0	0	0	0
0	1	0	0	0	0
0	0	0	0	0	0

Image matrix



0	1	0
0	1	0
0	1	0

Filter matrix



1	0	0	0
2	0	0	0
3	0	0	0
2	0	0	0

Feature map

Example calculation, top left: $0*0 + 0*1 + 0*0 + 0*0 + 0*1 + 0*0 + 0*0 + 1*1 + 0*0 = 1$

Convoluzione

- In computer vision "tradizionale" i Kernel sono definiti manualmente, in modo tale da **estrarre determinate caratteristiche (features)** di interesse dall'immagine di input.

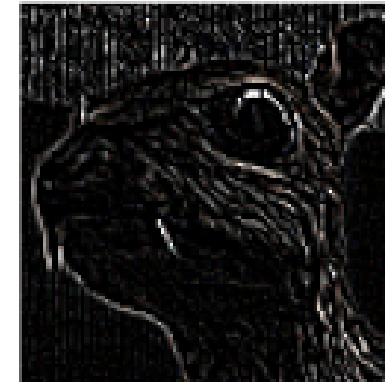
Input image



Convolution
Kernel

$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

Feature map



A kernel (filtri) diversi, corrispondono diverse features estratte

Convoluzione

- Ci sono due parametri che controllano come l'operazione di convoluzione viene effettuata: lo **stride** (orizzontale e verticale) ed il **padding**.
- Questi due parametri altro non sono che il numero di pixel da "saltare" quando si fa scorrere il kernel sull'immagine di input (stride), e da aggiungere come cornice all'input (padding)

0 ₂	0 ₀	0 ₁	0	0	0	0
0 ₁	2 ₀	2 ₀	3	3	3	0
0 ₀	0 ₁	1 ₁	3	0	3	0
0	2	3	0	1	3	0
0	3	3	2	1	2	0
0	3	3	0	2	3	0
0	0	0	0	0	0	0

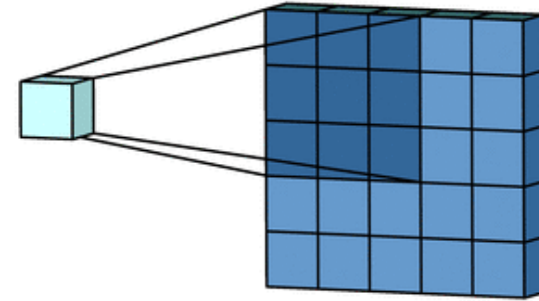
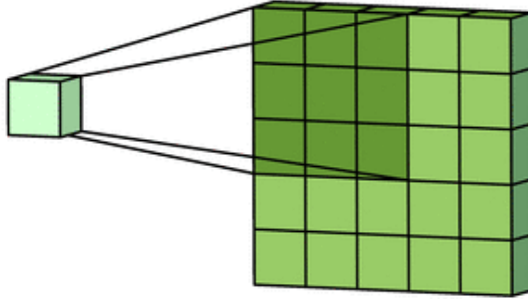
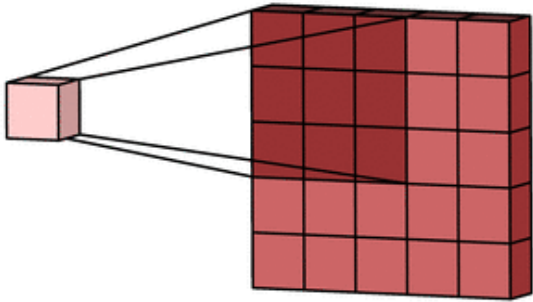
1	6	5
7	10	9
7	10	8

Convoluzione tra un immagine 5×5, con uno **zero padding** 1×1 ed un kernel 3×3, usando stride 2×2

Convoluzione

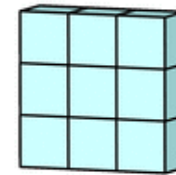
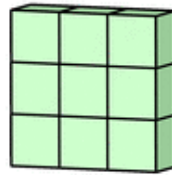
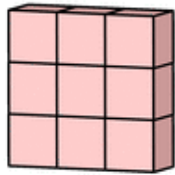
Data un immagine RGB, la convoluzione con un kernel 3×3 viene effettuata facendo:

- La convoluzione di ogni canale per il rispettivo filtro a profondità 1



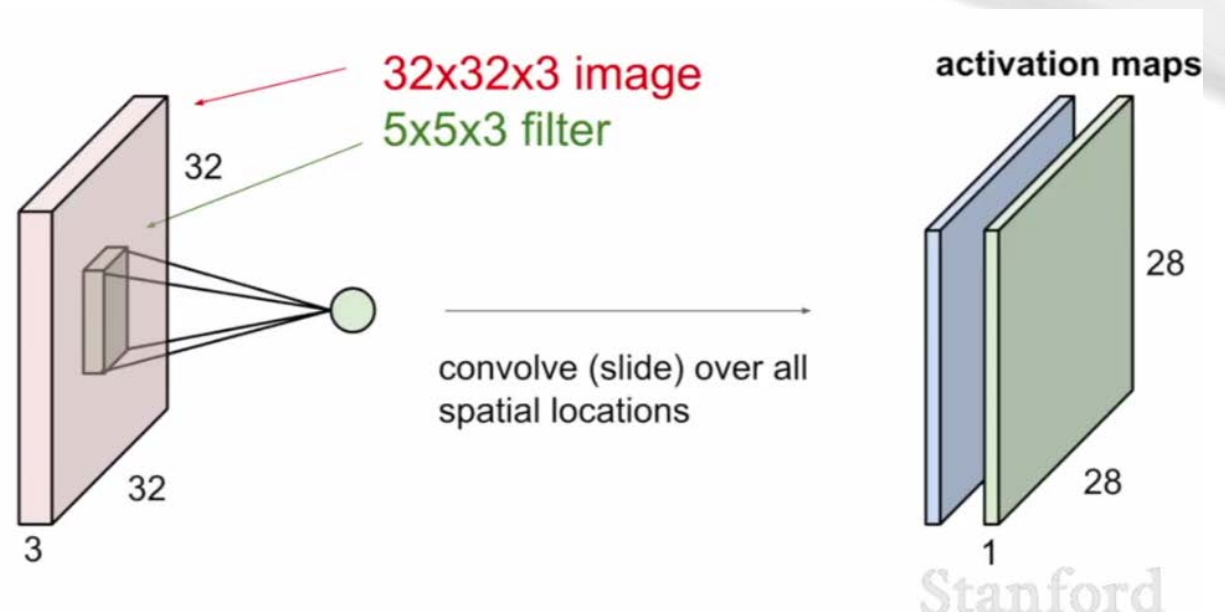
Convoluzione

- La somma delle feature map estratte



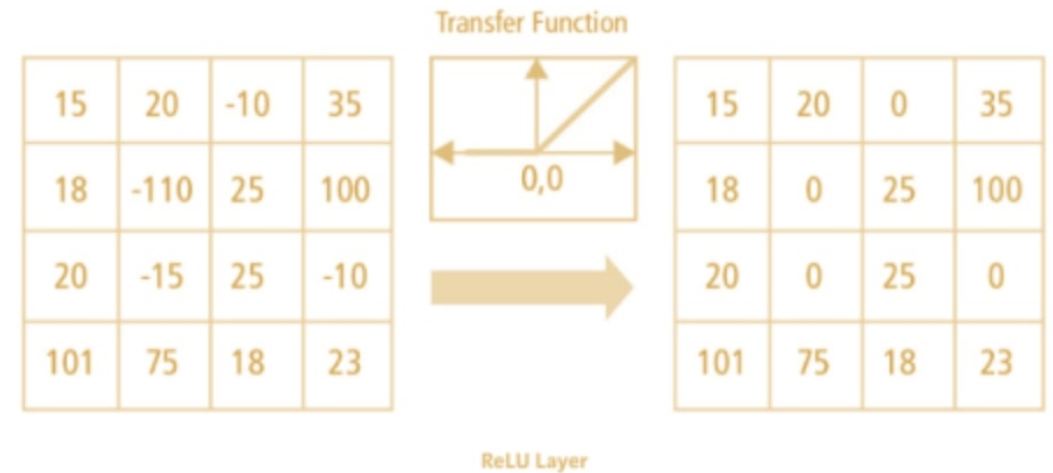
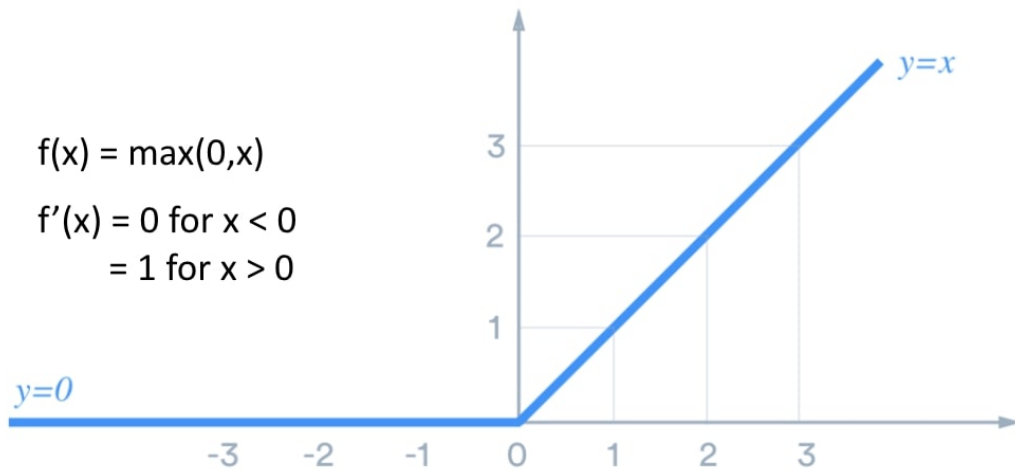
Imparare i kernel

- Apprendere filtri convoluzionali consiste nel **definire un numero arbitrario N** di filtri da apprendere ed eseguire N convoluzioni tra volumi, in maniera indipendente.
- Ogni filtro convoluzionale è un **neurone** che osserva una regione locale dell'immagine (sulla quale scorre).
- Ogni convoluzione produce una feature map, che passa attraverso una funziona di attivazione e diventa una **activation map**.



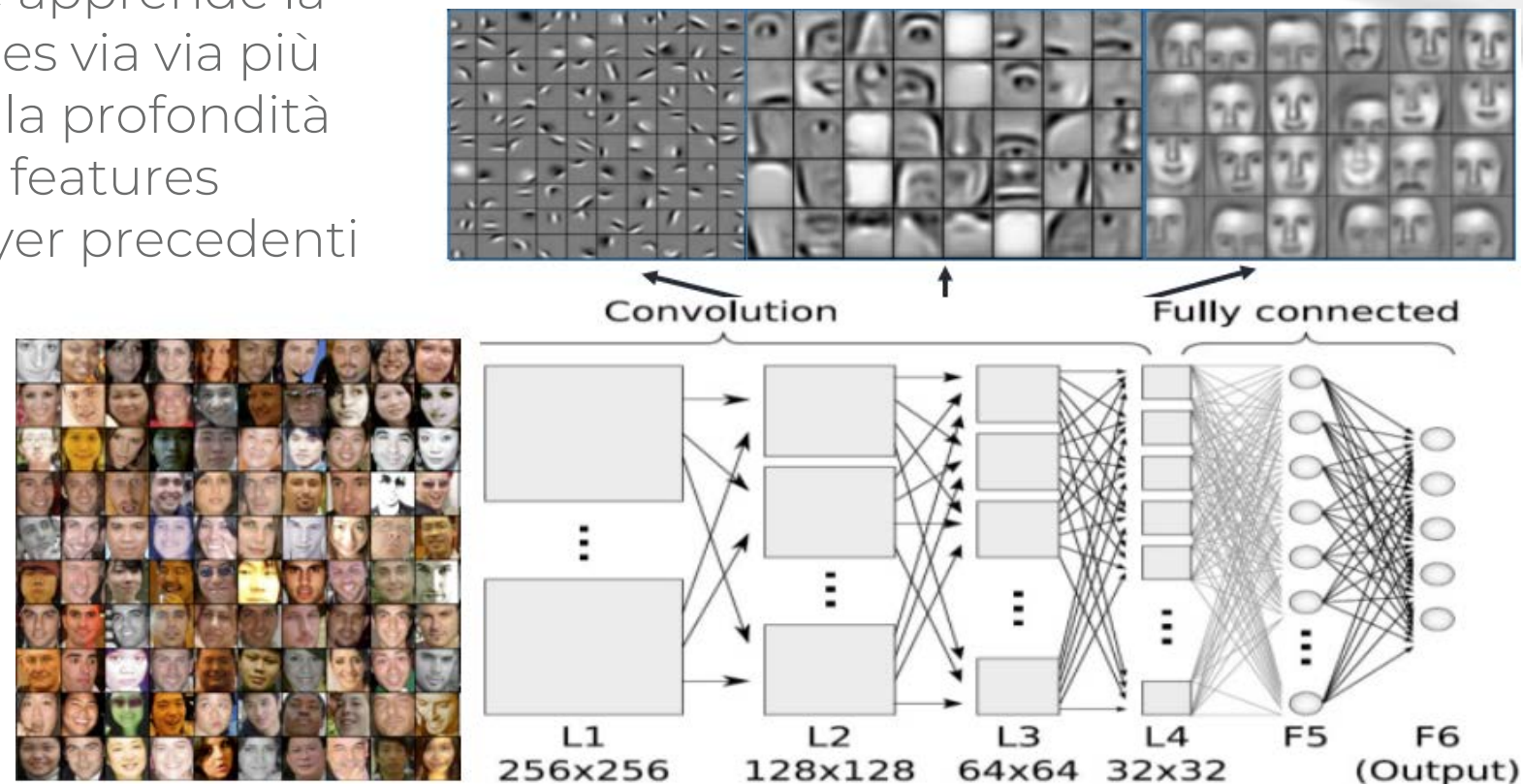
ReLU

- ReLU è l'abbreviazione di unità lineare rettificata ed è una cosiddetta **funzione di attivazione**. L'obiettivo principale dell'uso di una funzione di attivazione è di aggiungere non linearità al calcolo. Uno strato ReLU prende una feature map e rettifica qualsiasi valore negativo a zero. I numeri positivi rimangono invariati.



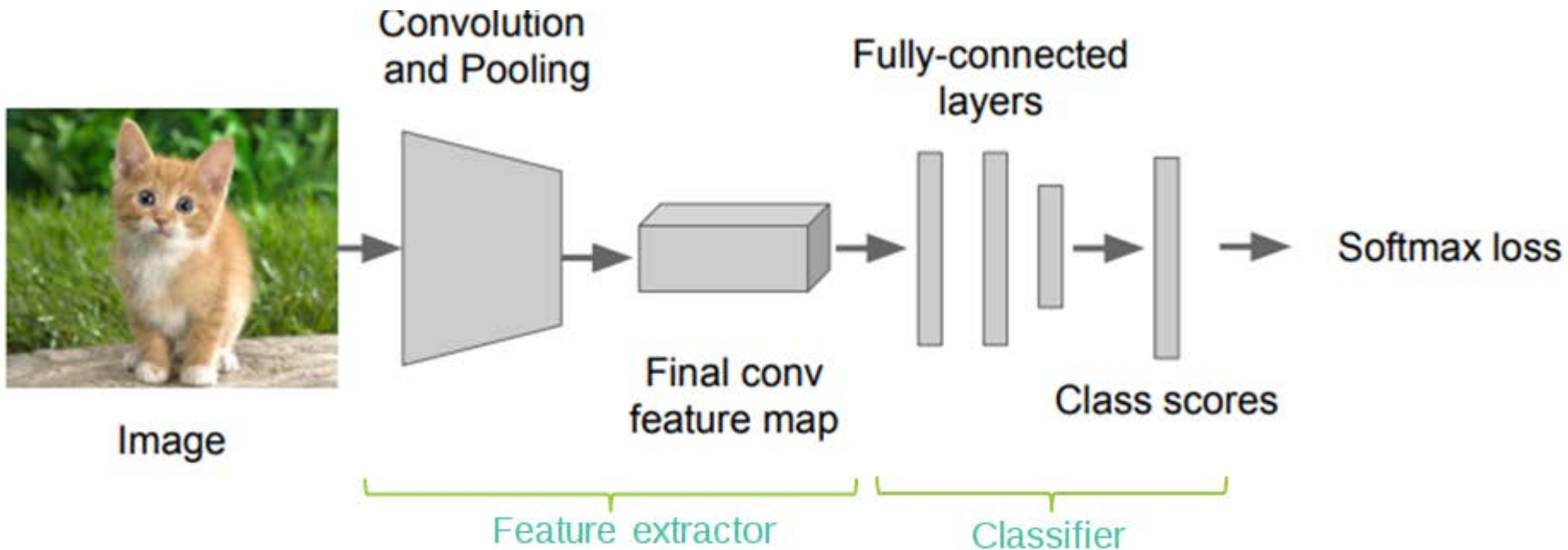
Layers di convoluzione

- Ogni layer convoluzionale apprende la capacità di estrarre features via via più astratte all'aumentare della profondità della rete, combinando le features elementari estratte dai layer precedenti



Classificazione ed estrazione di features

- Essendo a tutti gli effetti una rete neurale, possiamo definire l'architettura in modo tale da estrarre un numero arbitrario di features (a bassa dimensionalità) ed utilizzarlo come input di una classifikatore (rete FC).



- YOLO è stato introdotto per la prima volta nel 2016 ed è stato una pietra miliare nella ricerca sul rilevamento degli oggetti grazie alla sua capacità di rilevare gli oggetti in tempo reale con una migliore precisione.
- È stato proposto da **Joseph Redmon**, laureato dell'Università di Washington. L'articolo che descrive YOLO ha vinto il premio OpenCV People's Choice Award nella conferenza su Computer Vision e Pattern Recognition (**CVPR**) nel 2016.

Versioni di YOLO di Joseph Redmon:

- Version 1
'You Only Look Once: Unified, Real-Time Object Detection' (2016)
- Version 2
'YOLO9000: Better, Faster, Stronger' (2017)
- Version 3
'YOLOv3: An Incremental Improvement' (2018)

- L'implementazione principale di YOLO di Redmon è basata su **Darknet**, che è un framework di reti neurali open source scritto in C e CUDA. Darknet imposta l'architettura di base della rete e la usa come framework per l'addestramento di YOLO. Questa implementazione è stata introdotta da Redmon stesso ed è veloce, facile da installare e supporta il calcolo su CPU e GPU.
- Tuttavia, nel febbraio 2020, Joseph Redmon, il creatore di YOLO ha annunciato di aver interrotto la sua ricerca nel campo della computer vision! Ha inoltre dichiarato che ciò era dovuto a diverse preoccupazioni riguardanti il potenziale impatto negativo del suo lavoro.



Joseph Redmon
@pjreddie

I stopped doing CV research because I saw the impact my work was having. I loved the work but the military applications and privacy concerns eventually became impossible to ignore.

1. Introduction

Sometimes you just kinda phone it in for a year, you know? I didn't do a whole lot of research this year. Spent a lot of time on Twitter. Played around with GANs a little. I had a little momentum left over from last year [12] [1]; I managed to make some improvements to YOLO. But, honestly, nothing like super interesting, just a bunch of small changes that make it better. I also helped out with other people's research a little.

Yolov3 - introduzione

YOLO v4

Il ritiro di Redmon non è stata la fine di YOLO. In ricordo di molti nella comunità della computer vision, la quarta generazione di YOLO è stata rilasciata nell'aprile 2020. È stato introdotto in un documento intitolato '[*YOLOv4: Optimal Speed and Accuracy of Object Detection*](#)' by Alexey Bochkovskiy et al.

Inoltre, il lavoro di Redmon è stato continuato da Alexey nel [*fork*](#) del repository [*principale*](#). YOLO v4 è stato considerato il modello in tempo reale più veloce e accurato per il rilevamento di oggetti.



Joseph Redmon
@pjreddie



Doesn't matter what I think! At this point [@alexeyab84](#) has the canonical version of darknet and yolo, he's put a ton of work into it and everyone uses it, not mine haha.

- Dopo il rilascio di YOLO v4, in soli due mesi di tempo, è stata rilasciata un'altra versione di YOLO chiamata YOLO v5! E' opera di Glenn Jocher, già noto nella comunità per aver creato la popolare implementazione PyTorch di YOLO v3.
- Il 9 giugno 2020, Jocher ha dichiarato che la sua implementazione di YOLO v5 è rilasciata pubblicamente ed è raccomandata per l'uso in nuovi progetti. L'implementazione di YOLOv5 di Jocher differisce dalle versioni precedenti in alcuni modi notevoli. In primo luogo, Jocher non ha (ancora) pubblicato un documento per accompagnare il suo rilascio. Secondo, Jocher ha implementato YOLOv5 nativamente in PyTorch, mentre tutti i modelli precedenti della famiglia YOLO sfruttano Darknet.
- Il rilascio di YOLO v5 ha attirato molta attenzione e ha causato accese discussioni nelle piattaforme della comunità Machine Learning. Questo era principalmente dovuto a diversi fatti su un articolo pubblicato dal team Roboflow riguardo a YOLO v5.

YOLO architettura

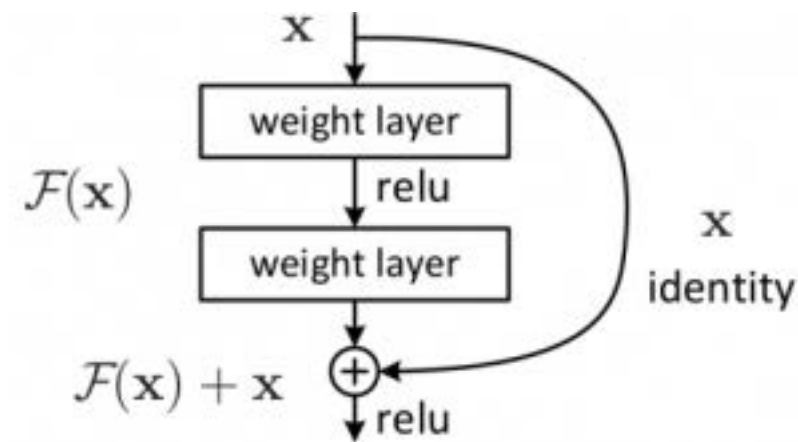
- YOLO fa uso solo di layer convoluzionali, rendendola una rete completamente convoluzionale (FCN). Ha **75 layer di convoluzione**, con **skip connection** e **layer di upsampling**. Non viene utilizzata alcuna forma di pooling, e uno layer di convoluzione con stride 2 viene utilizzato per ricampionare le feature map. Questo aiuta a prevenire la perdita di feature di basso livello.
- Essendo una FCN, YOLO è invariante alla dimensione dell'immagine di ingresso. Tuttavia, in pratica, potremmo voler attenerci a una dimensione di input costante a causa di vari problemi implementativi.
- Uno di questi problemi è che se vogliamo elaborare le nostre immagini in **batch** (le immagini in batch possono essere elaborate in parallelo dalla **GPU**, con conseguente aumento della velocità), abbiamo bisogno di avere tutte le immagini di altezza e larghezza fissa. Questo è necessario per concatenare più immagini in unico batch.
- La rete scala l'immagine di un fattore chiamato **stride della rete**. Per esempio, se lo stride della rete è **32**, allora un'immagine di input di dimensioni **416 x 416** produrrà un output di dimensioni **13 x 13**. In generale, lo stride di qualsiasi strato della rete è uguale al fattore per cui l'uscita dello strato è più piccola dell'immagine in ingresso alla rete.

YOLO – skip connection

Le reti neurali deep possono imparare funzioni complesse in modo più efficiente delle loro controparti non-deep. Tuttavia, durante l'addestramento delle reti neurali deep, le prestazioni del modello diminuiscono con l'aumento della profondità dell'architettura. Questo è noto come problema della **degradazione**.

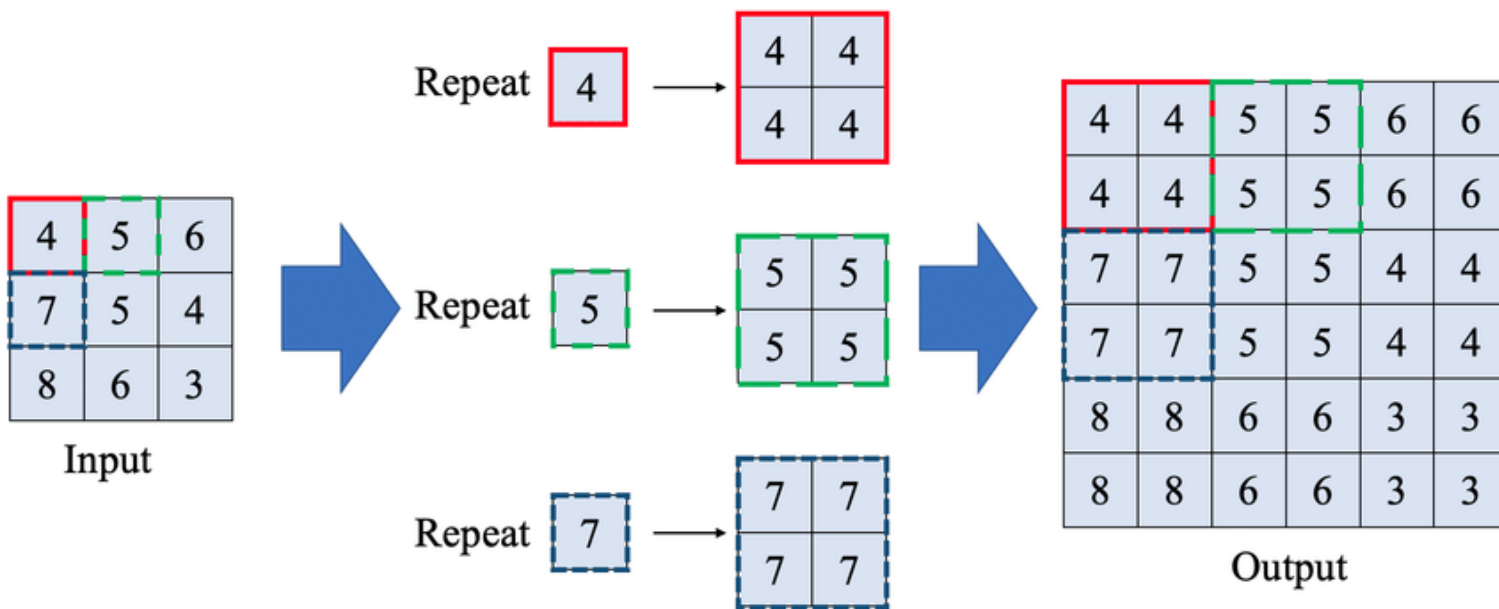
Skip Connections (o Shortcut Connections), come suggerisce il nome, salta alcuni dei layer della rete neurale e alimenta l'uscita di un layer come ingresso dei layer successivi.

Le Skip Connections sono state introdotte per risolvere diversi problemi in diverse architetture. Nel caso di ResNets, skip Connections hanno risolto il problema della degradazione.



YOLO – layer di upsampling

Un layer di upsampling aumenta la larghezza e l'altezza dell'array di input. Il metodo più semplice per fare upsampling è ripetere l'elemento della matrice di input in due direzioni.



YOLO vs Sliding window

Nel metodo **sliding window** quello che veniva fatto era una scansione lineare dell'immagine alla ricerca di oggetti conosciuti. Se all'interno di una finestra viene riconosciuto un oggetto, allora non solo potremmo dire che quell'oggetto è presente nell'immagine, ma anche indicarne la posizione.

Questa tecnica presenta però dei difetti. In particolare ogni oggetto riconoscibile tramite questo metodo deve avere una dimensione adatta per rientrare nella finestra di scansione e non deve trovarsi in una posizione tale da finire sulla linea di demarcazione tra due finestre adiacenti.

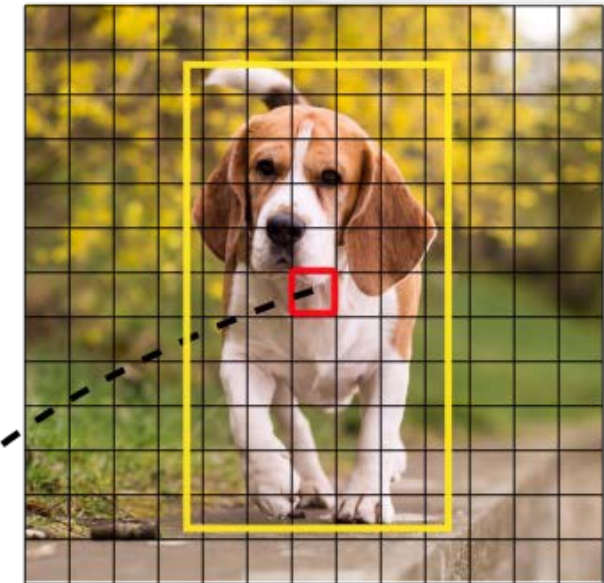
YOLO (You Only Look Once) permette di riconoscere più oggetti e, allo stesso tempo, di identificarne la posizione e lo spazio occupato effettuando una sola "lettura" dell'immagine di input.



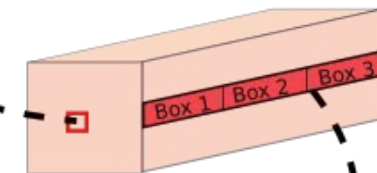
YOLO – feature map

- Consideriamo un esempio dove l'immagine di input è 416×416 , e lo stride della rete è 32. Come indicato in precedenza, le dimensioni feature map saranno 13×13 . Dividiamo quindi l'immagine di input in 13×13 celle.
- La cella (sull'immagine di input) che contiene il valore di Ground Truth è scelta per essere quella responsabile della predizione dell'oggetto. Nell'immagine, è la cella segnata in rosso.
- La cella rossa è la 7a cella nella 7a riga della griglia. Ora assegniamo la 7a cella nella 7a riga sulla feature map come responsabile del rilevamento del cane.

Image Grid. The Red Grid is responsible for detecting the dog



Prediction Feature Map



YOLO - output

- Ognuna delle bounding box ha 5 + C attributi, dove C è il numero delle classi di oggetti che il modello può classificare.

t_x : Coordinata x del centro della bounding box

t_y : Coordinata y del centro della bounding box

t_w : Larghezza della bounding box

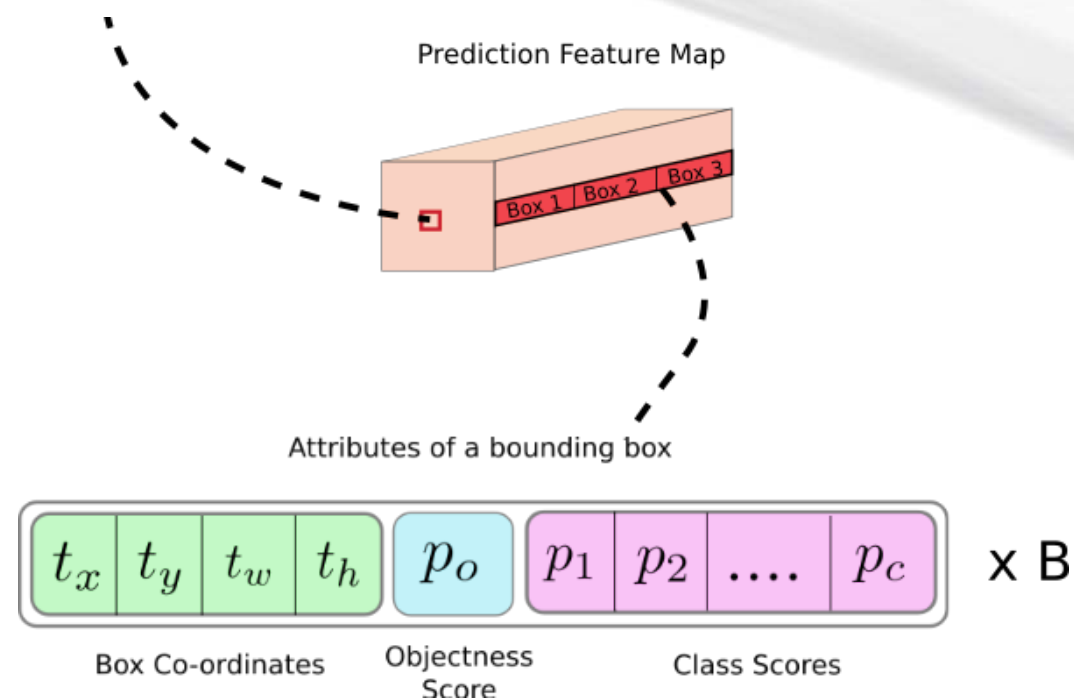
t_h : Altezza della bounding box

p_o : Score

p_1 : Probabilità che l'oggetto appartenga alla classe 1

...

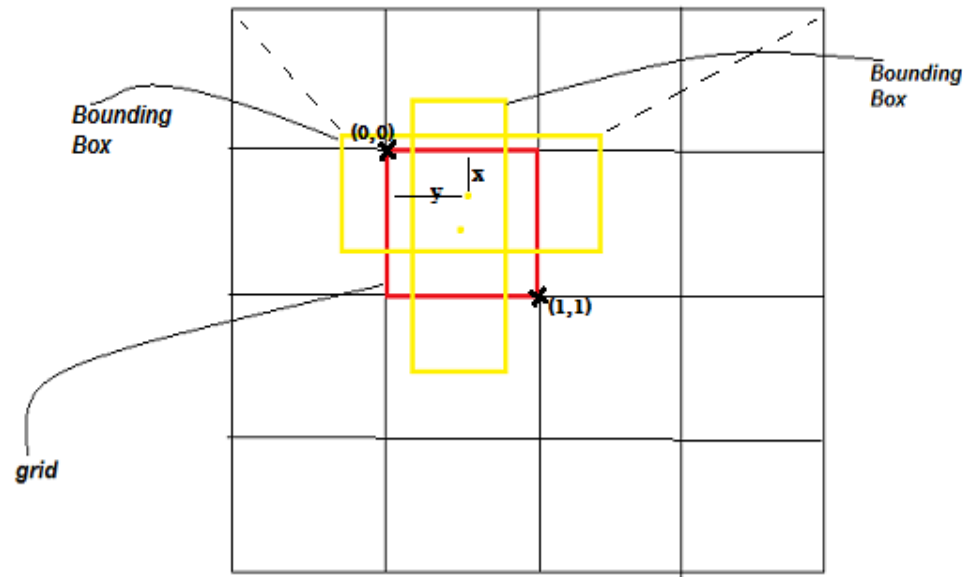
p_c : Probabilità che l'oggetto appartenga alla classe c



YOLO - output

- Inoltre in YOLO ogni cella può predire **tre** bounding box, quindi il vettore di output finale avrà dimensione:

$$3 \times (5 + N_{\text{classi}})$$

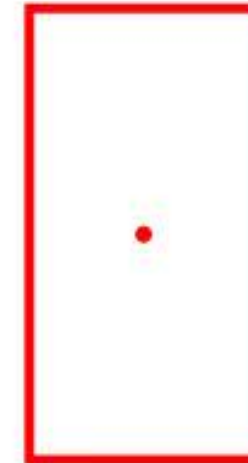


YOLO – anchor boxes

Anchor Boxes

- Potrebbe avere senso prevedere la larghezza e l'altezza della bounding box, ma in pratica, questo porta a gradienti instabili durante l'allenamento. Invece, la maggior parte dei moderni classificatori predice le trasformazioni log-space, o semplicemente gli offset delle bounding box predefinite chiamate **ancore**.
- Poi, queste trasformazioni sono applicate alle anchor boxes per ottenere la predizione. YOLO utilizza tre ancore, che si traducono nella predizione di tre bounding box per cella.
- La bounding box responsabile del rilevamento del cane sarà quella la cui ancora ha il più alto **IoU** con il riquadro del GT.

Anchor box 1



Anchor box 2



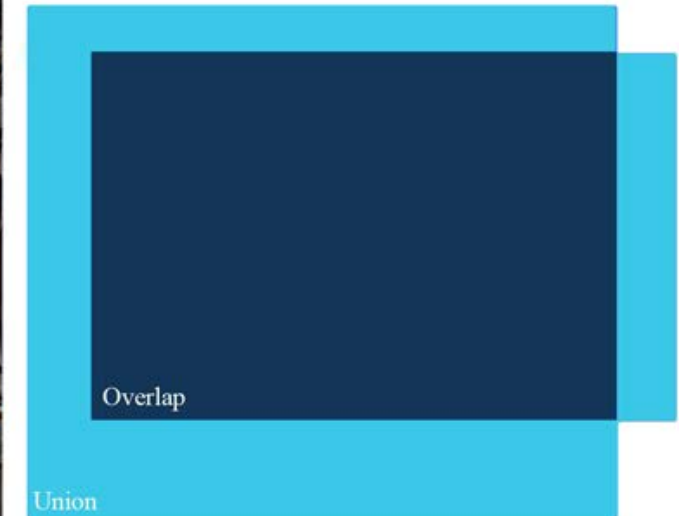
YOLO - IoU

IoU

Per decidere se una bounding box predetta è a tutti gli effetti una "detection" è possibile utilizzare l'intersection over union (delle aree):

$$\text{IoU} = \frac{A \cap B}{A \cup B}$$

Il suo valore è tra 0 (nessuna detection) ed 1 (overlap completo).

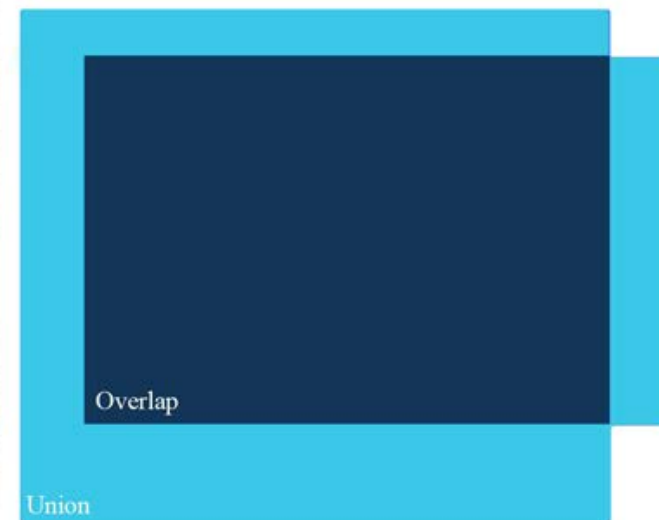


YOLO - IoU

Per *decidere* se una bounding box predetta è da considerare una **detection**, si fissa una soglia di IoU e vengono considerate come detection solo quelle predizioni di bounding box che superano la soglia.

Così facendo ci ritroviamo nel caso di una classificazione binaria, in cui possiamo avere:

- True positives: detection
- False positives: detection errate (non passano la soglia)



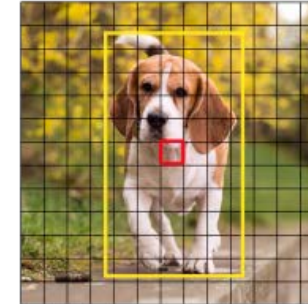
YOLO - Scales

YOLO v3 fa previsioni su **3 scale** diverse. Il layer di detection è utilizzato per predire feature map di tre dimensioni diverse, con **stride 32, 16, 8** rispettivamente. Ciò significa che, con un input di **416 x 416**, si effettuano detection su scale **13 x 13, 26 x 26 e 52 x 52**.

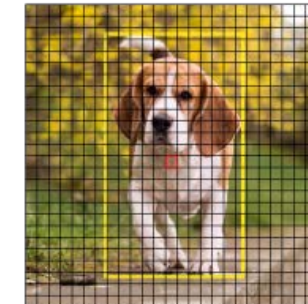
La rete ricampiona l'immagine di input fino al primo layer di detection, dove viene effettuato una prima detection utilizzando le feature map di un layer con stride 32.

Inoltre, gli strati sono sovracampionati di un fattore 2 e concatenati con le feature map di uno strato precedente con dimensioni identiche. Un'altra detection viene quindi fatta al layer con stride 16. La stessa procedura di upsampling è ripetuta, e una detection finale è fatta al layer con stride 8.

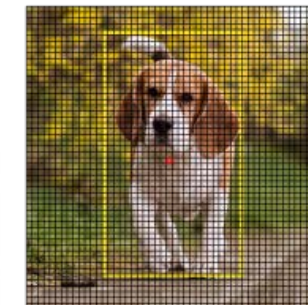
Prediction Feature Maps at different Scales



13 x 13



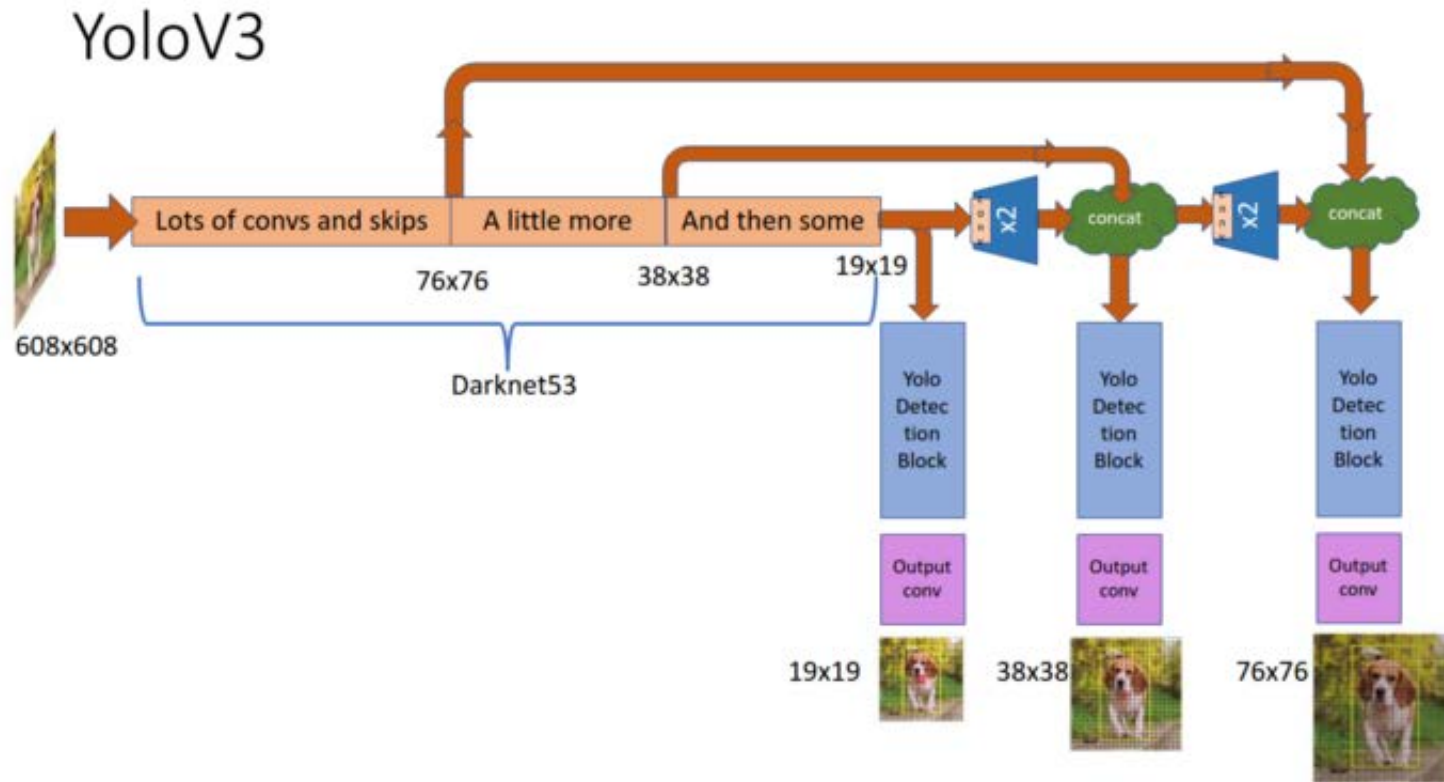
26 x 26



52 x 52

YOLO - Scales

Ad ogni scala, ogni cella predice 3 bounding box usando 3 ancore, rendendo il numero totale di ancore usate 9. (Le ancore sono diverse per scale diverse)



YOLO - NMS

Per un'immagine di dimensioni 416 x 416, YOLO predice $((52 \times 52) + (26 \times 26) + 13 \times 13) \times 3 = 10647$ bounding box.

Algoritmo Non-max suppression

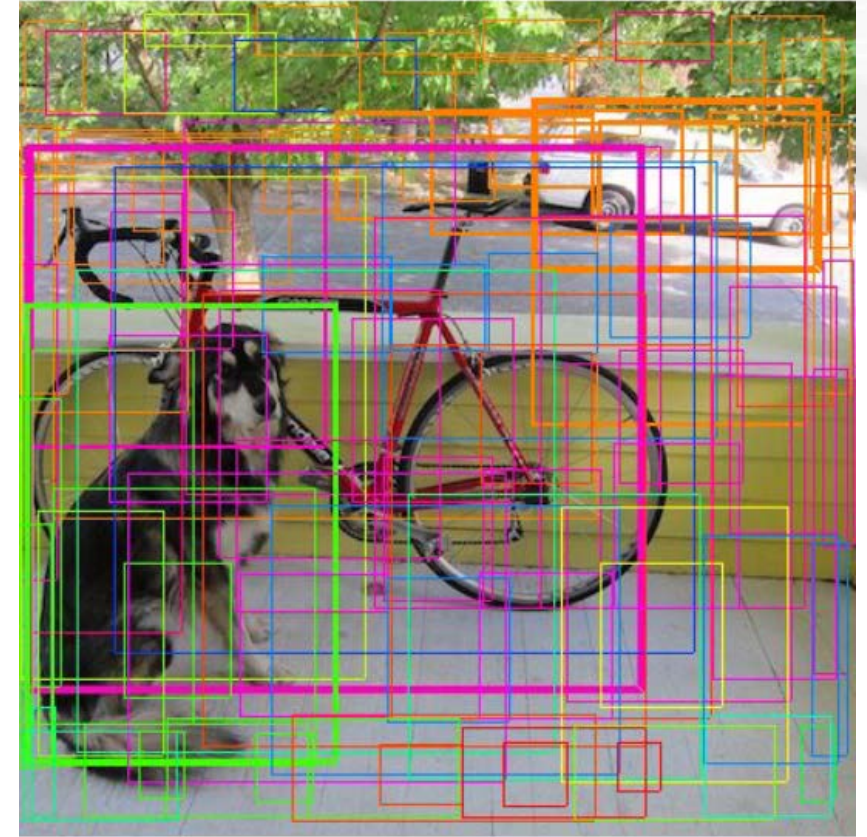
- Rimuovere tutte le bounding box con uno score $p_o < 0.5$

Con le rimanenti detection:

- Selezionare la bounding box con il più alto score p_o e utilizzarla come output della predizione
- Rimuovere tutte le rimanenti bounding box con

$$\text{IoU} > 0.45$$

con la bounding box utilizzata come output allo step precedente



YOLO - NMS

Per un'immagine di dimensioni 416 x 416, YOLO predice $((52 \times 52) + (26 \times 26) + 13 \times 13) \times 3 = 10647$ bounding box.

Algoritmo Non-max suppression

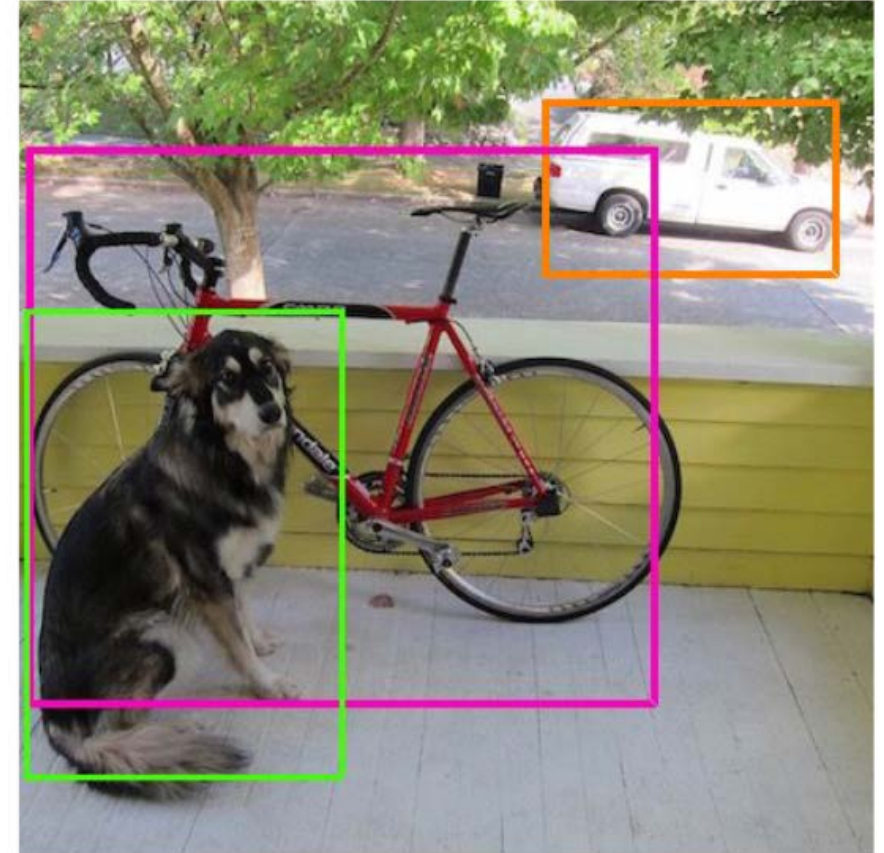
- Rimuovere tutte le bounding box con uno score $p_o < 0.5$

Con le rimanenti detection:

- Selezionare la bounding box con il più alto score p_o e utilizzarla come output della predizione
- Rimuovere tutte le rimanenti bounding box con

$$\text{IoU} > 0.45$$

con la bounding box utilizzata come output allo step precedente



YOLO – Funzione di Loss

$$\begin{aligned} & \lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B 1_{ij}^{obj} [(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2] \\ & + \lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B 1_{ij}^{obj} [(\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2] \\ & + \sum_{i=0}^{S^2} \sum_{j=0}^B 1_{ij}^{obj} (C_i - \hat{C}_i)^2 + \lambda_{noobj} \sum_{i=0}^{S^2} \sum_{j=0}^B 1_{ij}^{noobj} (C_i - \hat{C}_i)^2 \\ & + \sum_{i=0}^{S^2} 1_i^{obj} \sum_{c \in classes} (p_i(c) - \hat{p}_i(c))^2 \end{aligned}$$

x_i, y_i è la posizione del centroide della bounding box

w_i, h_i è la larghezza e l'altezza della bounding box

C_i è l'Objectness, cioè il punteggio di confidenza della presenza o meno di un oggetto

$p_i(c)$ è la loss di classificazione.

YOLO - Funzione di Loss

$$\lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B 1_{ij}^{obj} [(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2]$$

$$+ \lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B 1_{ij}^{obj} [(\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2]$$

$$+ \sum_{i=0}^{S^2} \sum_{j=0}^B 1_{ij}^{obj} (C_i - \hat{C}_i)^2 + \lambda_{noobj} \sum_{i=0}^{S^2} \sum_{j=0}^B 1_{ij}^{noobj} (C_i - \hat{C}_i)^2$$

$$+ \sum_{i=0}^{S^2} 1_i^{obj} \sum_{c \in classes} (p_i(c) - \hat{p}_i(c))^2$$

1. Penalizza la cattiva localizzazione del centro delle celle
2. Penalizza il rettangolo di selezione con altezza e larghezza imprecise.
3. Rende il punteggio di confidenza uguale all'IOU tra l'oggetto e la predizione quando c'è un oggetto
4. Rende il punteggio di confidenza vicino a 0 quando non ci sono oggetti nella cella
5. Loss di classificazione

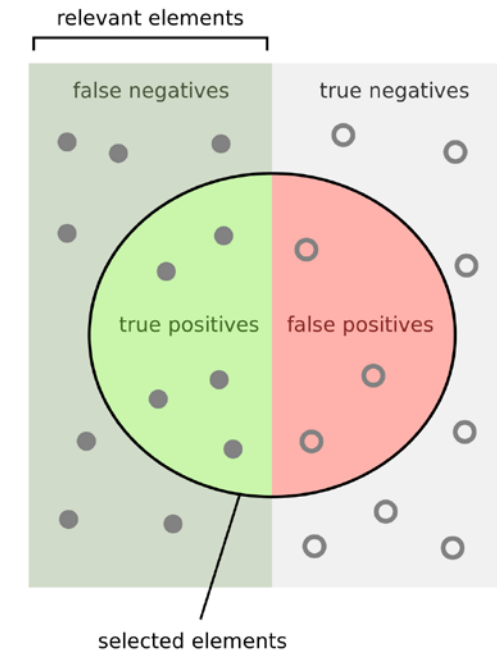
Metriche

La **precision** misura quanto sono accurate le vostre previsioni, cioè la percentuale di previsioni corrette. Misura quante delle previsioni che il modello ha fatto erano effettivamente corrette.

La **recall** misura quanto bene trovi tutti i positivi. Per esempio, possiamo trovare l'80% dei possibili casi positivi nelle nostre prime K previsioni.

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$



How many selected items are relevant?

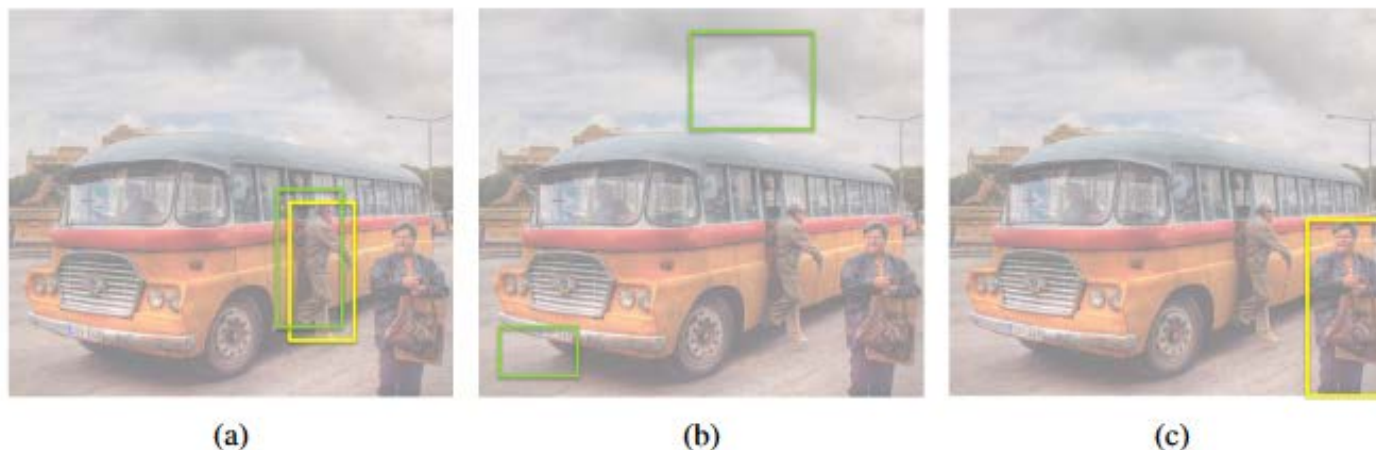
$$\text{Precision} = \frac{\text{Green Circle}}{\text{Green Circle} + \text{Red Circle}}$$

How many relevant items are selected?

$$\text{Recall} = \frac{\text{Green Circle}}{\text{Green Circle} + \text{Green Square}}$$

Metriche

- a) **True positive**, perché la sovrapposizione del GT (box giallo) e della predizione (box verde) è superiore allo 0,5
- b) **False positive**, poiché le bounding box di predizione (verdi) non si sovrappongono a nessuna delle annotazioni di GT.
- c) **False negative**, poiché l'annotazione di GT (giallo) non viene rilevato dal modello



$$\text{class}(IoU) = \begin{cases} \text{Positive} \rightarrow IoU \geq \text{Threshold} \\ \text{Negative} \rightarrow IoU < \text{Threshold} \end{cases}$$

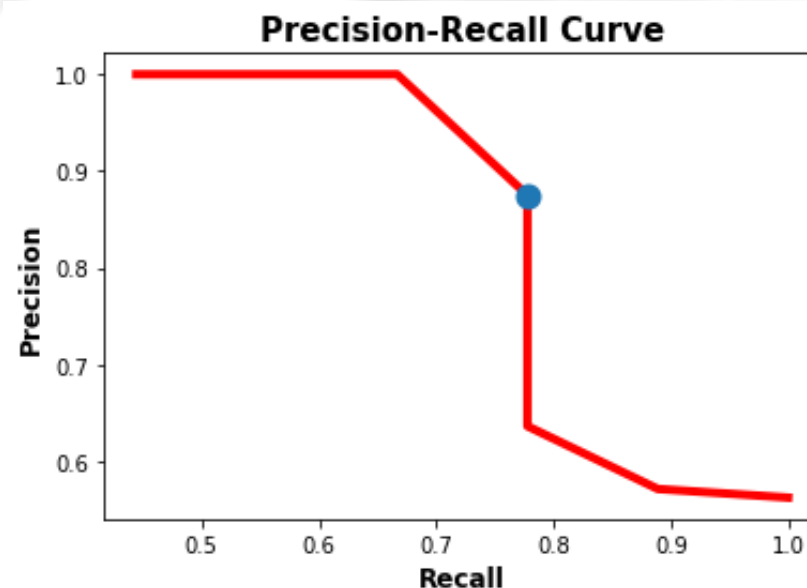
Metriche

La definizione generale per l'**Average Precision (AP)** è trovare l'area sotto la curva precision-recall.

mAP (mean Average Precision) è la media dell'AP su tutte le classi di oggetti.

$$AP = \int_0^1 P(R) dR$$

$$mAP = \frac{1}{|Q_R|} \sum_{q=Q_R} AP(q)$$



Curva precision (asse y) e recall (asse x).
Il modello è considerato un buon modello predittivo se la precisione rimane alta mentre il richiamo aumenta

1. Exercise.py (da completare)
Effettuare la detection di un'immagine utilizzando il modello YOLOv5s pre-allenato



Small
YOLOv5s

14 MB_{FP16}
2.0 ms_{V100}
37.2 mAP_{COCO}



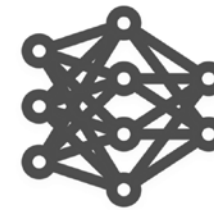
Medium
YOLOv5m

41 MB_{FP16}
2.7 ms_{V100}
44.5 mAP_{COCO}



Large
YOLOv5l

90 MB_{FP16}
3.8 ms_{V100}
48.2 mAP_{COCO}



XLarge
YOLOv5x

168 MB_{FP16}
6.1 ms_{V100}
50.4 mAP_{COCO}

Yolov5 models

2. Annotare un nuovo dataset utilizzando il tool RoboFlow, utilizzare il tool di data Augmentation per ingrandire il dataset

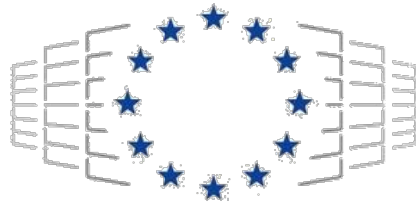
<https://app.roboflow.com/>

3. Re-allenare il modello sul nuovo dataset.

```
$ python train.py --img 640 --batch 16 --epochs 3 --data data.yaml --weights yolov5s.pt
```

4. Visualizzare le metriche del modello allenato.
 - 4.1 Test freeze backbone
 - 4.2 Test cambiare il numero di epoche
 - 4.3 Test cambiare IoU threshold
5. Fare inference con il nuovo modello su un'immagine
6. Fare inference con il nuovo modello su un video

THANK YOU!



EuroHPC
Joint Undertaking

This project has received funding from the European High-Performance Computing Joint Undertaking (JU) under grant agreement No 951732. The JU receives support from the European Union's Horizon 2020 research and innovation programme and Germany, Bulgaria, Austria, Croatia, Cyprus, Czech Republic, Denmark, Estonia, Finland, Greece, Hungary, Ireland, Italy, Lithuania, Latvia, Poland, Portugal, Romania, Slovenia, Spain, Sweden, United Kingdom, France, Netherlands, Belgium, Luxembourg, Slovakia, Norway, Switzerland, Turkey, Republic of North Macedonia, Iceland, Montenegro