

Progetto 3 – Computer vision

Lo scopo di questo progetto è quello di realizzare una Convolutional Neural Networks (CNN) in modo da classificare correttamente la categoria di appartenenza di una foto.

Il dataset usato per il training della rete è quello utilizzato nella ricerca [1] composta da 15 categorie.

Il linguaggio scelto per la realizzazione del progetto è stato MATLAB.

Parte 1

La prima parte del progetto è incentrata nella realizzazione della seguente CNN (figura 1) e nel suo allenamento a riconoscere la categoria di appartenenza di una foto.

#	type	size
1	Image Input	64×64×1 images
2	Convolution	8 3×3 convolutions with stride 1
3	ReLU	
4	Max Pooling	2×2 max pooling with stride 2
5	Convolution	16 3×3 convolutions with stride 1
6	ReLU	
7	Max Pooling	2×2 max pooling with stride 2
8	Convolution	32 3×3 convolutions with stride 1
9	ReLU	
10	Fully Connected	15
11	Softmax	softmax
12	Classification Output	crossentropyex

Per la realizzazione della struttura della rete è stato utilizzata l'app deep network design presente in MATLAB, che permette di esportare in codice la struttura della CNN usando un interfaccia grafica.

Figura 1

Manipolazione dei dati

Per il training e la verifica, i dati a disposizione sono stati divisi in un rapporto 85% per il training ed il restante per il testing.

I dati per il training sono stati ulteriormente divisi, utilizzando la funzione di MATLAB, *splitEachLabel()*, in questo modo:

un 85 % per il training e il restante 15 per la validazione.

Adesso le foto, per essere utilizzate dalla CNN, devono essere dimensionate per essere adatte all'ingresso del primo layer (64x64 pixel x1 livello colore) e per far questo, utilizzando la funzione *augmentedImageDatastore()*, possiamo ottenere una dimensione della foto che rispetta il vincolo ed inoltre permette di modificare le foto in modo di avere casuali variazioni in traslazione e rotazione delle stesse.

Impostazione learning

Di seguito si preparano le impostazioni per il training utilizzando la funzione *trainingOptions()* ponendo nel nostro caso:

- 'sgdm', per utilizzare il stochastic gradient descent with momentum;
- 'MiniBatchSize', 32, per avere una minibatch size di 32.

Infine, note le caratteristiche imposte per il progetto, per avere dei pesi iniziali derivati da una Gaussiana con media 0 e deviazione standard 0,01 si modificano i layer *convolution2dLayer()* e *fullyConnectedLayer()* inserendo la seguente componente: 'WeightsInitializer','narrow-normal'. Mentre per avere un bias iniziale uguale a 0 non è richiesta alcuna variabile perché è un'impostazione di default in MATLAB.

Risultati

Adesso si riportano i dati ottenuti durante il training del risultato migliore.

I metodi scelti per terminare il training sono stati i seguenti:

- fissare un limite max di epoche, cioè quante volte l'insieme di foto, usate per il learning (epoch), viene ripetuto prima di fermarsi;
- fissare un valore massimo di quante volte la funzione Loss ha superato il suo valore minimo sui dati di validazione, nel caso migliore il valore era posto a 5.

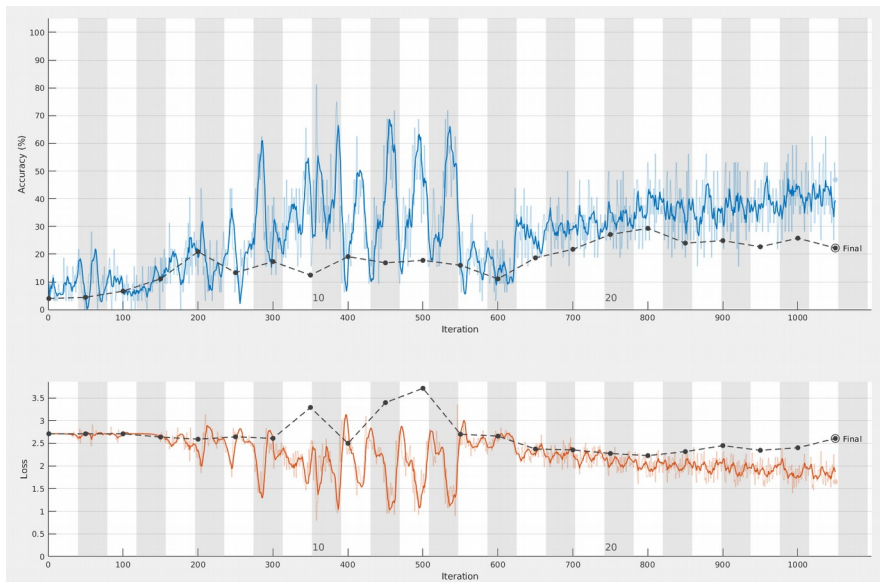


Figura 2

Con una confusion matrix qui riportata (Figura 2):

		Confusion Matrix																															
Output Class	Bedroom	8	0	0	1	4	3	3	1	4	0.1%	0.3%	6	8	0.2%	0.8%	1	2	3	1	0	0%	2	1	14.3%								
	Coast	15	211	23	70	20	11	13	12	96	12	4.3%	127	13	8	4	27	31.9%	0.5%	7.1%	0.8%	2.3%	0.7%	0.4%	0.4%	0.9%	68.1%						
	Forest	2	0	37	2	4	20	6	4	2	8	3	19	4	10	10	28.2%	0.1%	0%	1.2%	0.1%	0.1%	0.7%	0.2%	0.1%	0.3%	71.8%						
	Highway	1	2	1	32	4	6	2	3	6	1	11	0	4	6	1	40.0%	0.0%	0%	0%	1.1%	0.1%	0.2%	0.1%	0.2%	0%	60.0%						
	Industrial	18	13	16	11	60	13	19	37	29	9	27	14	11	7	42	18.4%	0.6%	0.4%	0.5%	0.4%	2.0%	0.4%	0.6%	1.2%	1.0%	0.3%	0.9%	81.6%				
	InsideCity	4	1	37	2	21	51	2	19	14	6	14	31	29	16	16	19.4%	0.4%	0.1%	1.2%	0.1%	0.7%	1.7%	0.1%	0.6%	0.5%	0.2%	0.5%	80.6%				
	Kitchen	3	0	8	1	5	7	10	10	6	5	2	11	3	1	15	11.5%	0.1%	0%	0.3%	0%	0%	0.2%	0.3%	0.3%	0.2%	0.2%	0.1%	88.5%				
	LivingRoom	6	1	9	1	19	10	6	19	6	4	6	22	9	6	15	13.7%	0.2%	0%	0.3%	0%	0%	0.2%	0.3%	0.2%	0.6%	0.2%	0.1%	86.3%				
	Mountain	14	16	12	4	8	9	8	14	43	3	26	16	8	6	15	21.3%	0.5%	0.5%	0.4%	0.1%	0.3%	0.3%	0.3%	0.5%	0.3%	0.2%	0.5%	78.7%				
	Office	17	6	17	9	27	17	19	27	20	40	28	24	10	16	41	12.6%	0.6%	0.2%	0.6%	0.4%	0.9%	0.6%	0.9%	0.7%	1.3%	0.9%	0.8%	0.3%	87.4%			
	OpenCountry	10	4	9	8	11	12	3	5	18	0	44	8	2	5	16	28.4%	0.3%	0.1%	0.3%	0.3%	0.4%	0.4%	0.1%	0.2%	0.6%	0.0%	1.5%	0.3%	0.1%	0.2%	71.6%	
	Store	4	0	18	1	4	8	4	6	2	6	5	18	2	1	10	20.2%	0.1%	0%	0%	0.8%	0%	0.1%	0.3%	0.1%	0.2%	0.1%	0.3%	79.8%				
	Street	12	5	12	13	8	22	5	15	15	8	6	17	2	88	34	33.0%	0.4%	0.2%	0.4%	0.4%	0.3%	0.7%	0.2%	0.5%	0.5%	0.3%	0.2%	0.6%	67.0%			
	Suburb	0	0	11	1	5	11	5	6	1	3	5	6	12	27	3	28.1%	0.0%	0%	0%	0%	0%	0.2%	0.4%	0.2%	0.4%	0.9%	0.1%	71.9%				
	TallBuilding	2	1	14	2	12	7	5	4	10	2	4	13	1	1	36	31.6%	0.1%	0.1%	0.5%	0.1%	0.4%	0.2%	0.4%	0.3%	0.1%	0.4%	0.0%	0%	68.4%			
		6.9%	81.2%	16.2%	20.0%	28.4%	24.5%	9.1%	10.1%	15.7%	34.8%	14.2%	8.4%	45.8%	19.1%	14.1%	24.3%	93.1%	18.8%	83.8%	80.0%	71.6%	75.5%	90.9%	89.9%	84.3%	65.2%	85.8%	91.6%	54.2%	80.9%	85.9%	75.7%
		Bedroom	Coast	Forest	Highway	Industrial	InsideCity	Kitchen	LivingRoom	Mountain	Office	OpenCountry	Store	Street	Suburb	TallBuilding																	
		Target Class																															

Figura 3

Si nota che per il caso migliore (Figura 2) il sistema è intervenuto prima della fine del valore massimo di epoche fissate ottenendo il valore massimo di accuratezza del 21,64 % sulle immagini di test.

Se il modello non viene precedentemente fermato si ottiene un peggioramento dell'accuratezza del 3% sui test data.

Parte 2

La seconda parte del progetto consiste nel migliorare le performance del modello utilizzando le seguenti variabili:

- aumentare le capacità di learning aggiungendo alle foto delle traslazioni casuali, rotazioni e riflessioni rispetto agli assi principali. Utilizzando la funzione `imageDataAugmenter()` e ponendo le seguenti variabili alla funzione `'RandXReflection',1` ; `'RandYReflection',1` ; `'RandXTranslation',[-3 3]` ; `'RandYTranslation',[-3 3]` e `'RandRotation',[-5,5]` ha portato ad un aumento del 5-10% all'accuratezza (ora 25-30%) sulle immagini per il test, il tutto a discapito di un maggiore tempo di training;
- aggiungere un batch normalization layer prima del layer ReLU, il che ha portato ad ottenere un'accuratezza del 50% circa, con un tempo di learning molto minore;

- modificare la dimensione dei layer di convoluzione usando le seguenti dimensioni 3x3, 5x5 e 7x7 rispettivamente per il primo, secondo e terzo layer, il che ha portato ad un aumento dell'accuratezza media del 5%;
- modificare eventuali altri parametri:
 - modificando la minibatchsize incrementandola, a parità di maxepoch di 50, il che ha portato ad un aumento del 1-2 %, data la poca differenza tra 128 e 256 ho lasciato il valore su 128;
 - passando ad un maxepoch di 100, che porta il sistema in regime di overfitting, ha mantenuto mediamente costante il risultato di accuratezza;
 - scegliendo un altro metodo di risoluzione, il metodo rmsprop non ha modificato di tanto il risultato, invece il metodo adam ne ha aumentato il risultato medio del 5 %
 - variando i parametri per i pesi, il sistema non ha modificato di molto il risultato portando a lasciare i valori di default
- aggiungere un layer di dropout non ha permesso ai grafici di accuratezza e di validazione di distanziarsi troppo, permettendo di non andare in overfitting ed ottenendo una percentuale di accuratezza simile al precedente;
- applicare un ensemble of networks di 5 reti CNN. Tali reti hanno avuto il training set uguale e poi attraverso un sistema di votazione è stata scelta la predizione più comune. Le CNN sono state modificate variando il maxepoch, minibatch size e ripetendo un insieme di layer, modificandone la dimensione e il kernel. Queste operazioni hanno portato ad avere un accuratezza media del 65 % con un picco di precisione del 66,30%

In conclusione, inizialmente il sistema aveva un'accuratezza media del 23% ed alla fine, dopo aver applicato le operazioni sopra descritte, abbiamo raggiunto un'accuratezza del 65%. Alla fine il metodo migliore è stato quello di porre un normalization layer prima del layer ReLU (incremento del 20%).

Parte 3

La terza parte si basa sul transfer learning utilizzando una rete già addestrata (AlexNet), questa parte si divide in due sezioni:

- la prima, predire la categoria di appartenenza modificando solo l'ultimo fully connected layer fc8 tenendo fissi i pesi dei layer precedenti (utilizzando la funzione di MATLAB *freezeWeights()*). Per far questo, dopo aver sostituito gli ultimi layer e modificando il layer fc8 per avere in output le 15 categorie del nostro training set, abbiamo cercato di ottenere il miglior risultato modificando le impostazioni sul learning dei pesi. Alla fine abbiamo ottenuto un accuratezza max del 81%;
- la seconda, addestrare una rete SVM sulle feature estratte dall'ultimo layer di convoluzione di AlexNet, MATLAB dispone della funzione *activation()* che fornisce le feature per un determinato layer (in questo caso 'conv5'). Presi i dettagli delle feature, utilizzando l'error correcting Output code (funzione MATLAB *fitcecoc()*) e ponendo il metodo risolutivo un sistema SVM lineare con metodo di comparazione one vs one abbiamo addestrato il sistema a capire quali determinate feature appartengono a una categoria. Il miglior risultato ottenuto è stato del 73,53%.

Bibliografia:

[1] Lazebnik, S., Schmid, C., and Ponce, J. (2006). Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06), volume 2, pages 2169–2178.