

Progetto Machine Learning

Davide Pietrasanta, 844824

Fabio D'Elia, 829937

Descrizione del dominio di riferimento e obiettivi dell'elaborato

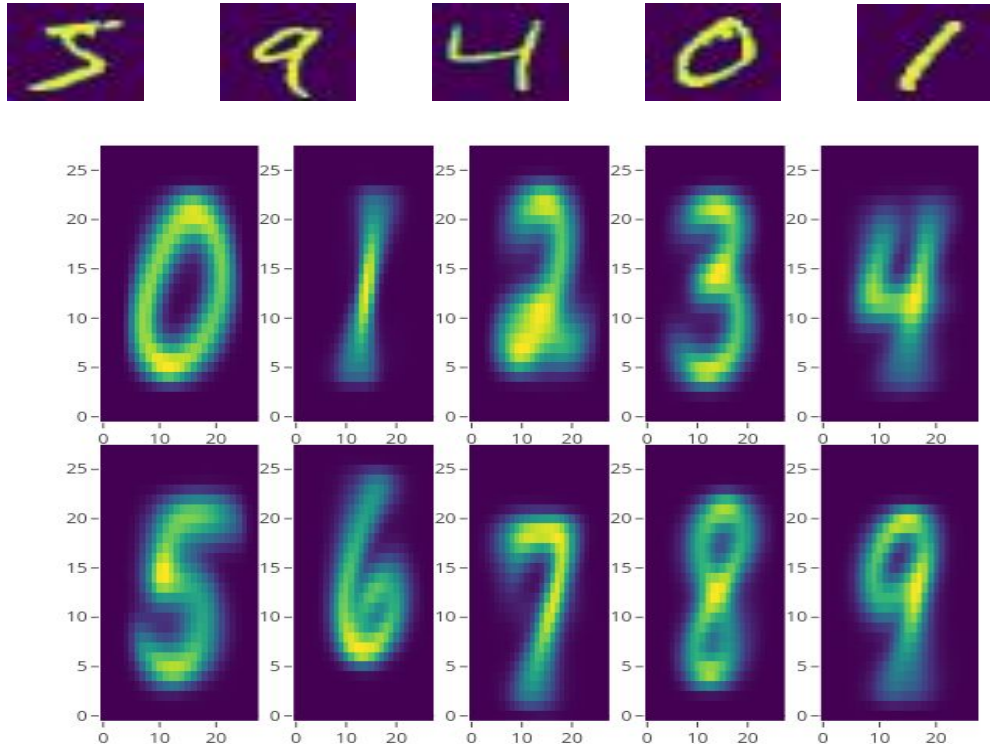
Implementazione di due algoritmi di apprendimento supervisionato per la classificazione di numeri scritti a mano.

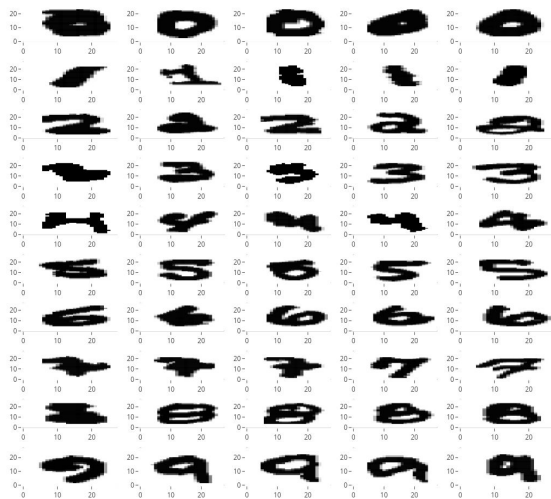
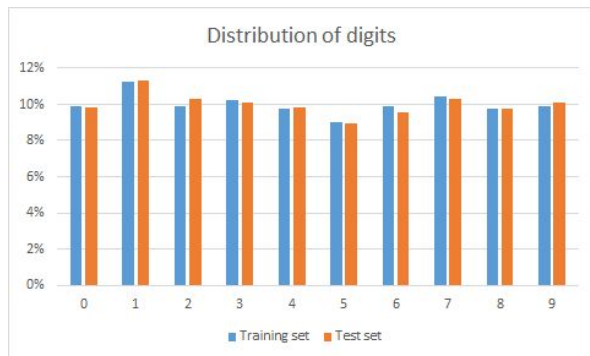
Dataset

Dataset MNIST

Il database MNIST è composto da 60,000 esempi di numeri scritti a mano nel train-set e 10,000 nel test-set.

Le immagini sono state normalizzate e centrate in una matrice 28x28.





Analisi dei dati

Abbiamo verificato la distribuzione uniforme delle classi nel dataset usando il chi-quadro test.

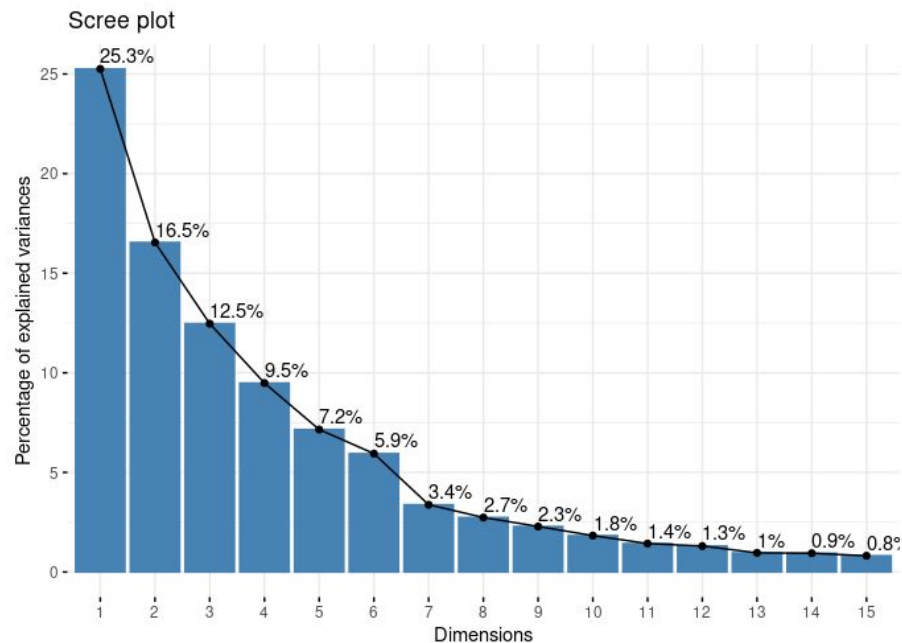
Un attento studio della distanza di ogni campione dal centroide della classe di appartenenza, ci ha fatto notare i possibili problemi di classificazione.

Principal Component Analysis

PCA

Abbiamo deciso di ridurre il numero di dimensioni del dataset in input, utilizzando la PCA.

Abbiamo deciso di mantenere 20 PCs che spiegavano il 95% della varianza cumulata.



Modelli

NAIVE BAYES

Si pensa che l'uso della probabilità possa aiutare nella predizione delle cifre.

Per esempio, una istanza con il centro bianco avrà poche probabilità di essere uno zero.

NEURAL NETWORK

La scelta dell'utilizzo di una rete neurale è dovuta dal fatto che una rete neurale si presenta come un modello flessibile e molto simile alle CNN.

NAIVE BAYES

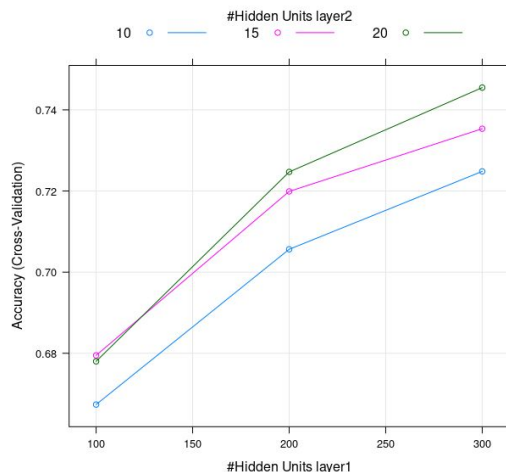
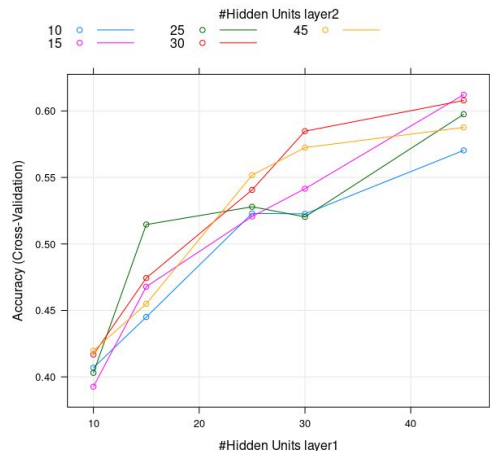
È stata applicata la 10-fold cross validation, in modo da valutare le performance del modello.

- **Accuracy:** 0.8520
- **Precision:** 0.8556952
- **Recall:** 0.8546793
- **F-measure:** 0.8546174
- **AUC:** 0.92027810
- **Accuracy of test:** 0.8621

Confusion Matrix Naive Bayes

		Target									
		9	8	7	6	5	4	3	2	1	0
Prediction	9	4744 79.7%	155 2.6%	274 4.4%	3 0%	40 0.7%	562 9.6%	137 2.2%	75 1.3%	11 0.2%	15 0.3%
	8	78 1.3%	4641 79.3%	31 0.5%	36 0.6%	101 1.8%	40 0.7%	275 4.5%	199 3.3%	53 0.8%	30 0.5%
	7	233 3.9%	43 0.7%	5439 86.8%	49 0.9%	21 0.4%	62 1%	118 2%	21 0.3%	7 0.1%	
	6	40 0.7%	62 1.1%	19 0.3%	529 89.4%	110 2%	92 1.6%	46 0.8%	169 2.8%	20 0.3%	125 2.1%
	5	147 2.5%	310 5.3%	42 0.7%	346 5.8%	4520 83.5%	101 1.7%	387 6.2%	94 1.6%	64 1%	189 3.2%
	4	502 8.4%	65 1.1%	157 2.5%	55 0.9%	93 1.7%	4923 84.3%	10 0.2%	93 1.6%	3 0%	11 0.2%
	3	80 1.3%	307 5.2%	11 0.2%	18 0.3%	357 6.6%	1 0%	4983 84.3%	136 2.3%	37 0.6%	24 0.4%
	2	47 0.8%	82 1.4%	112 2%	74 1.3%	60 1.1%	54 0.9%	151 2.5%	4948 83%	157 2.7%	29 0.5%
	1	45 0.8%	150 2.6%	151 2.4%	32 0.5%	15 0.3%	40 0.7%	45 0.7%	64 1.1%	6376 84.6%	
	0	33 0.6%	36 0.6%	29 0.5%	63 1.1%	70 1.3%	8 0.1%	35 0.6%	62 1%		5493 92.7%

NEURAL NETWORK

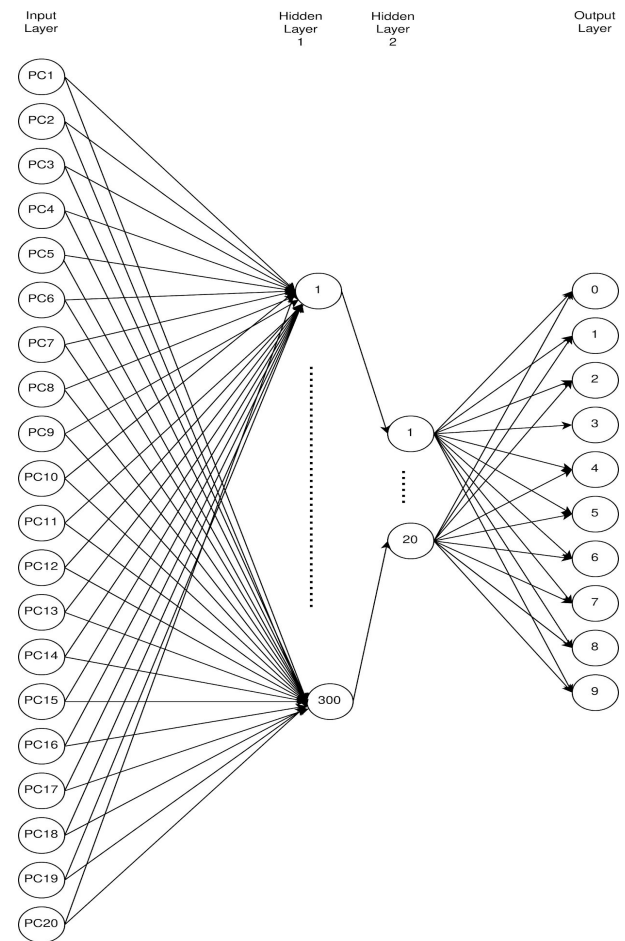


È stata applicata la 10-fold cross validation, in modo da valutare le performance del modello.

Attraverso 2 fasi di tuning differenti si è deciso di utilizzare una NN con un layout composto da 2 strati nascosti, con 300 nodi nel primo e 20 nel secondo.

- **Accuracy:** 0.7449164
- **Precision:** 0.742
- **Recall:** 0.72
- **F-measure:** 0.715
- **AUC:** 0.8453
- **Accuracy of test:** 0.7226

NEURAL NETWORK

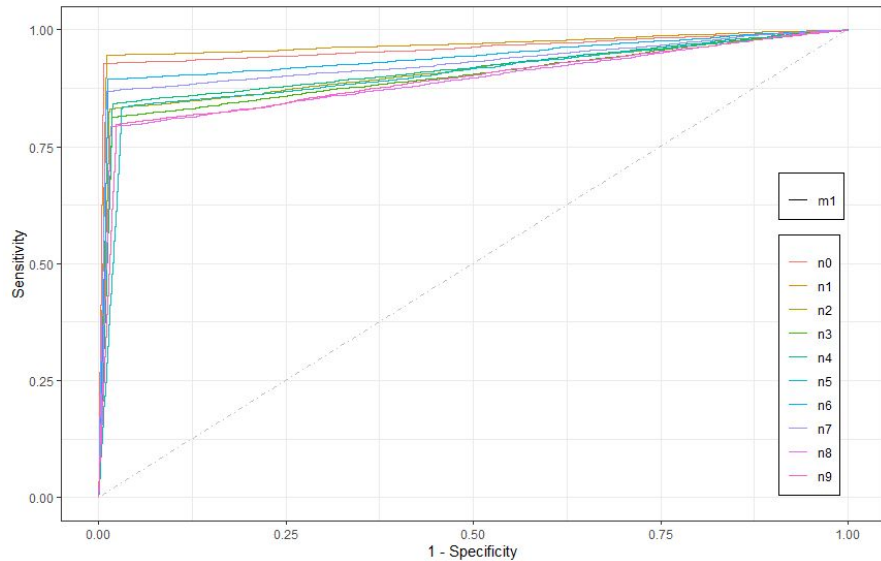


Confusion Matrix Neural Network

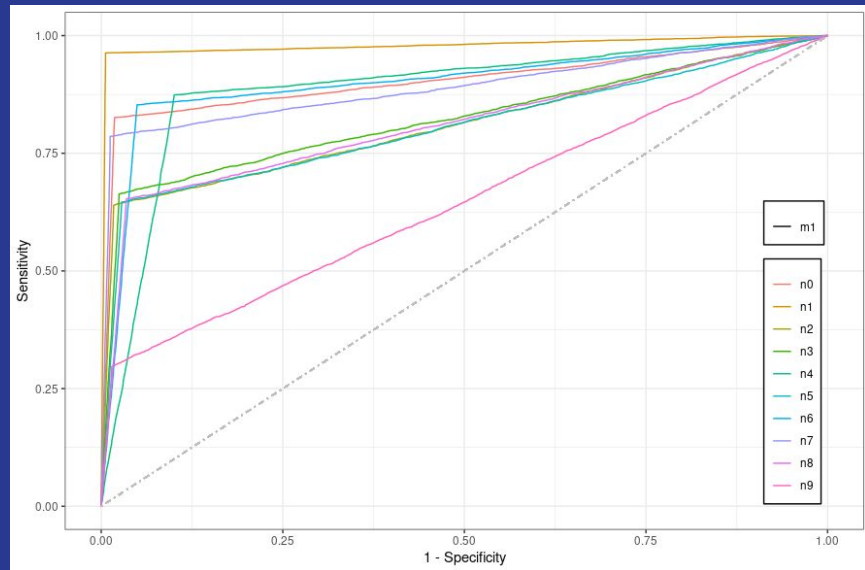
		Target									
		9	8	7	6	5	4	3	2	1	0
Prediction	9	1767 29.7%	53 0.9%	405 6.5%		51 0.9%	241 4.1%	24 0.4%	4 0.2%	10 0.1%	
	8	153 2.6%	3822 65.3%	133 2.1%	67 1.1%	282 5.2%	74 1.3%	748 13.1%	310 5.4%	49 0.7%	68 1.1%
	7	302 5.1%	75 1.3%	4923 78.6%		64 1.1%	102 1.7%	72 1.3%	28 0.5%	14 0.3%	7 0.1%
	6	40 0.7%	376 6.4%	8 0.1%	5048 85.3%	351 6.5%	188 3.2%	191 3.1%	974 16.3%	36 0.5%	508 8.6%
	5	177 3%	307 5.2%	77 1.2%	140 2.4%	3505 64.7%	49 0.8%	584 9.5%	87 1.5%	18 0.3%	112 1.9%
	4	3329 56%	386 6.6%	619 9.9%	320 5.4%	333 6.1%	5106 87.4%	114 1.9%	213 3.6%	28 0.4%	99 1.7%
	3	85 1.4%	479 8.2%	9 0.1%	23 0.4%	466 8.6%	9 0.2%	4069 66.4%	178 3%	57 0.8%	15 0.3%
	2	9 0.2%	178 3%	31 0.5%	139 2.3%	86 1.6%	25 0.4%	203 3.3%	3810 63.9%	38 0.6%	220 3.7%
	1	43 0.7%	92 1.6%	42 0.7%	27 0.5%	16 0.3%	37 0.6%	38 0.6%	19 0.3%	6492 96.3%	2 0%
	0	44 0.7%	83 1.4%	18 0.3%	154 2.6%	267 4.9%	11 0.2%	88 1.4%	335 5.7%		4892 82.6%

ROC e AUC

Naive Bayes AUC: 0.920278



Neural network AUC: 0.8453



Conclusioni

Entrambi i modelli, nonostante non siano modelli ottimali per il problema definito, forniscono delle discrete performance, anche senza usare l'informazione spaziale.

Si può notare che Naive Bayes fornisce predizioni in modo più accurato, ma con tempi maggiori, mentre Neural Network fornisce predizioni più veloci ma in modo meno accurato.

Pensiamo inoltre che con un miglior tuning degli iperparametri della Neural Network si possano raggiungere migliori risultati.

	NAIVE BAYES	NEURAL NETWORK
Train (s)	1.37	272.79
Predict test (s)	42.79	0.482
Predict train (s)	248.4	5.321
Accuracy (10-f CV)	0.852	0.745
Accuracy (test)	0.862	0.723
