

# Appendice – Grafici sui risultati ottenuti

## Prestazioni sul Dataset Human-generated E-mails

- **Caso 1:** I modelli ML classici (RF, LR, XGBoost) hanno ottenuto ottime prestazioni; KNN è risultato meno efficace; le feature numeriche iniziali sono risultate già altamente predittive;
- **Caso 2:** L’aggiunta di feature semantiche ha migliorato tutte le metriche; XGBoost ha raggiunto performance perfette comportandosi come un modello ideale; CNN, RNN, LSTM mostrano miglioramenti notevoli;
- **Caso 3:** I modelli LLM di tipo Transformer (RoBERTa, DeBERTa, BERT, DistilBERT, XLNet) hanno superato le reti neurali classiche grazie alla capacità di discriminazione semantica avanzata; i modelli RoBERTa, DeBERTa, DistilBERT, BERT, XLNet si distinguono con F1 e AUC prossimi a 1, confermando l'efficacia dell’approccio LLM nella comprensione del contenuto testuale.

Caso 1 – Feature numeriche iniziali – dataset Human Generated E-mail:

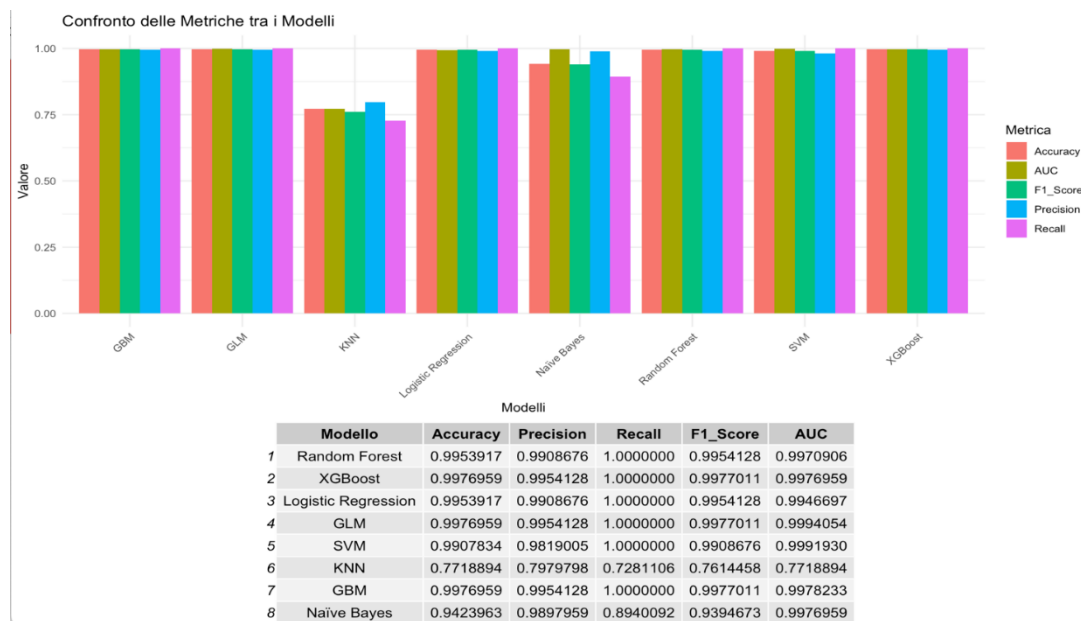


Figura 1 - Confronto metriche modelli ML (Accuracy, Precision, Recall, F1, AUC) – Human generated - Caso 1

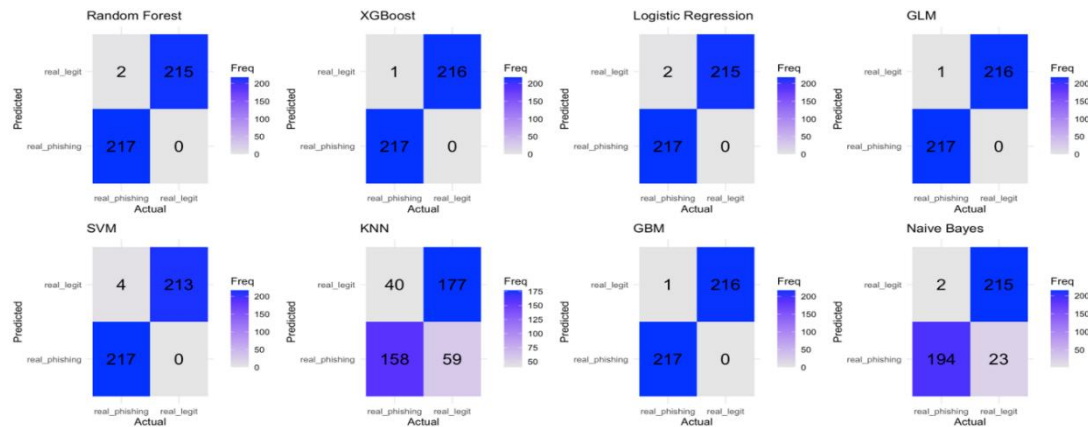


Figura 2 - Matrici di confusione - Human generated - Caso 1

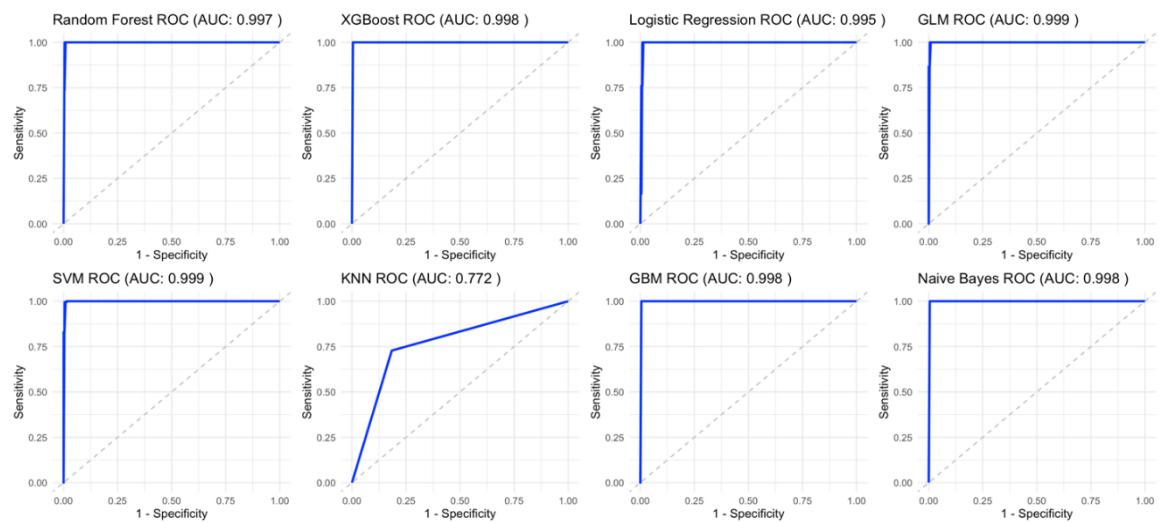


Figura 3 - Curve ROC-AUC - Human generated - Caso 1

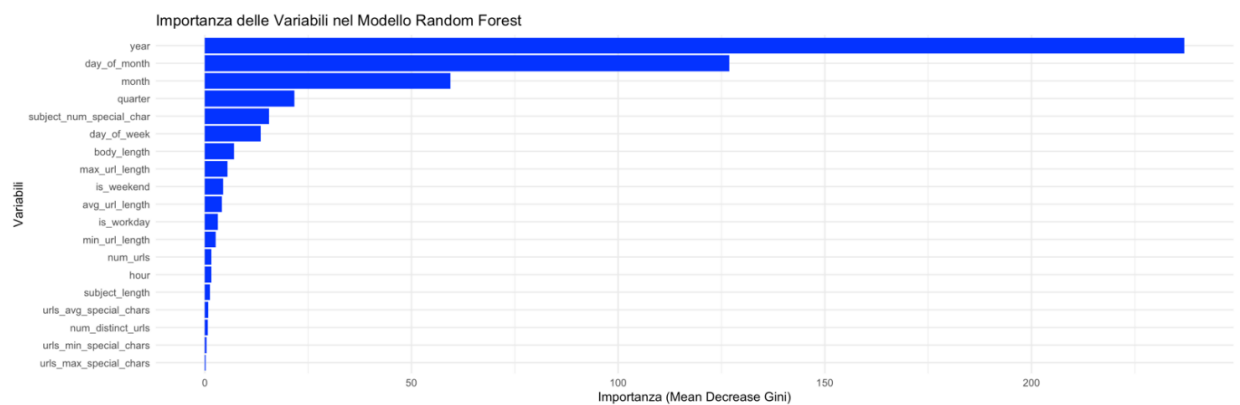


Figura 4 - Importanza variabili Random Forest - Human generated - Caso 1

Caso 2 – Tutte le feature numeriche (con topic, analisi del sentimento e embeddings):

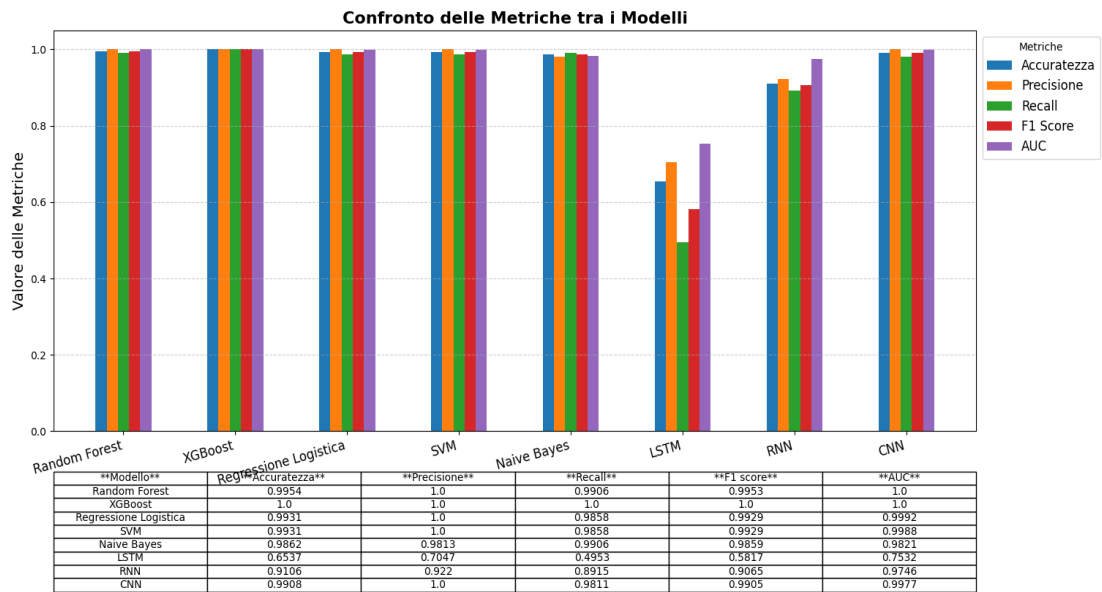


Figura 5 - Confronto metriche - Human generated - Caso 2

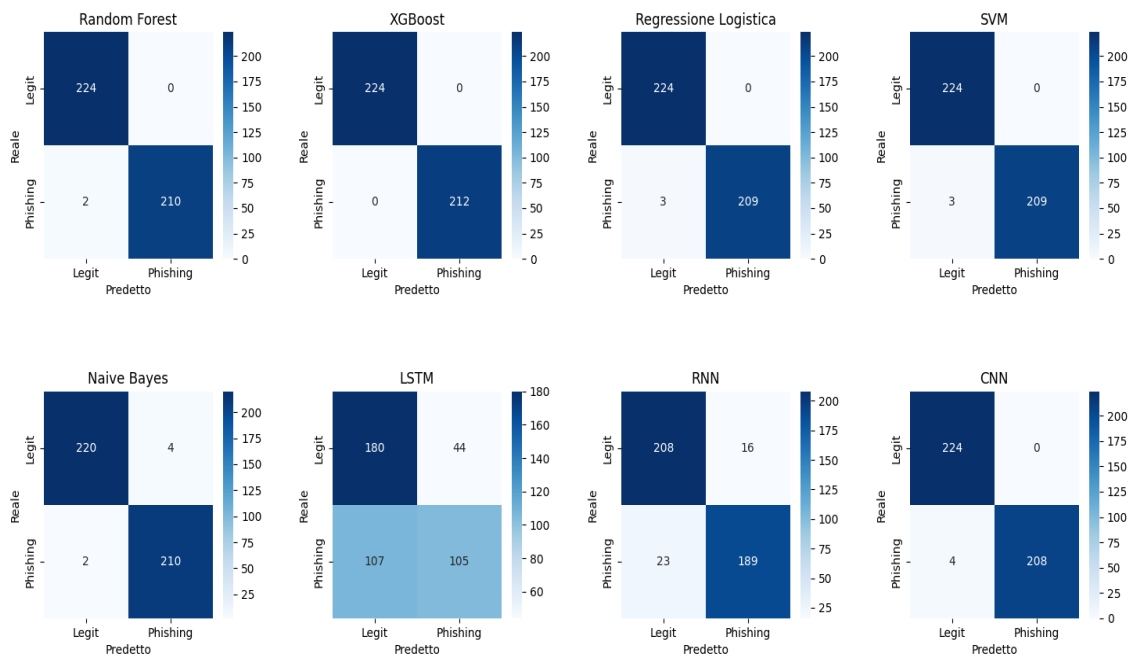


Figura 6 - Matrici di confusione Human generated - Caso 2

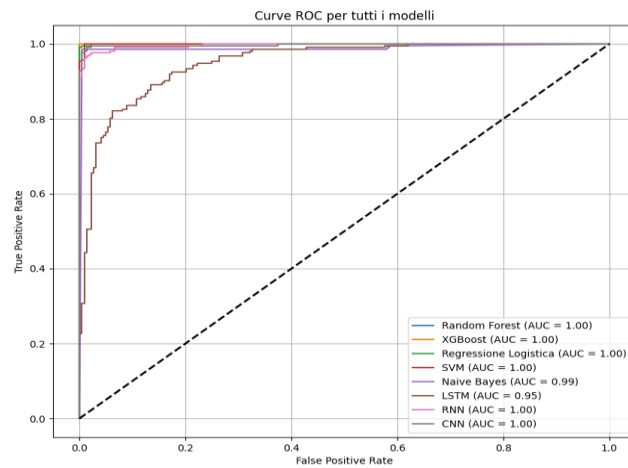


Figura 7 - Curve ROC - Human generated - Caso 2

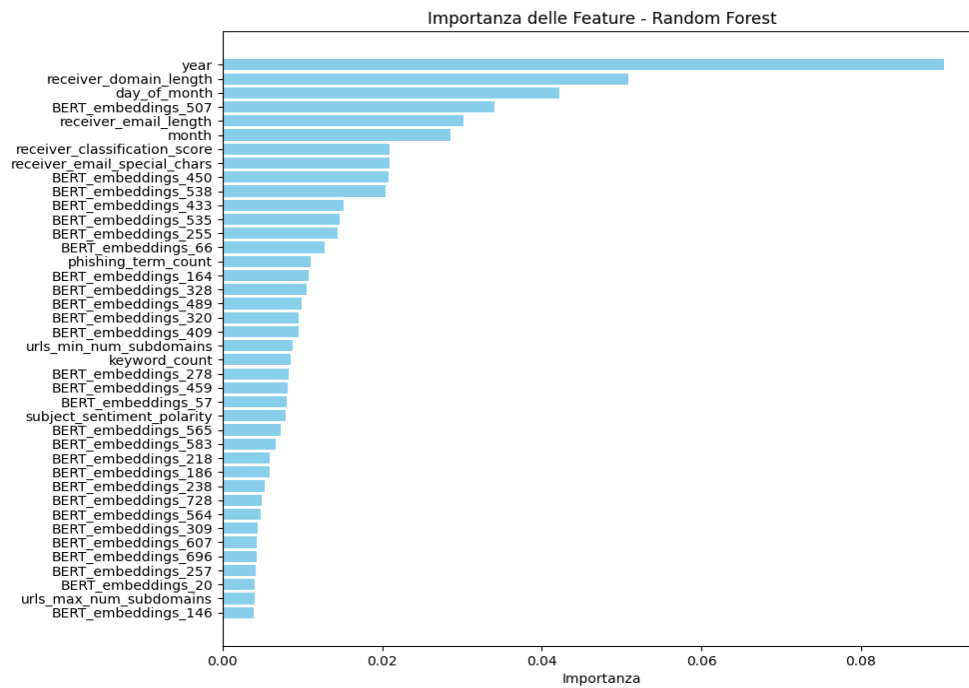


Figura 8 - Importanza feature per Random Forest - Human generated - Caso 2

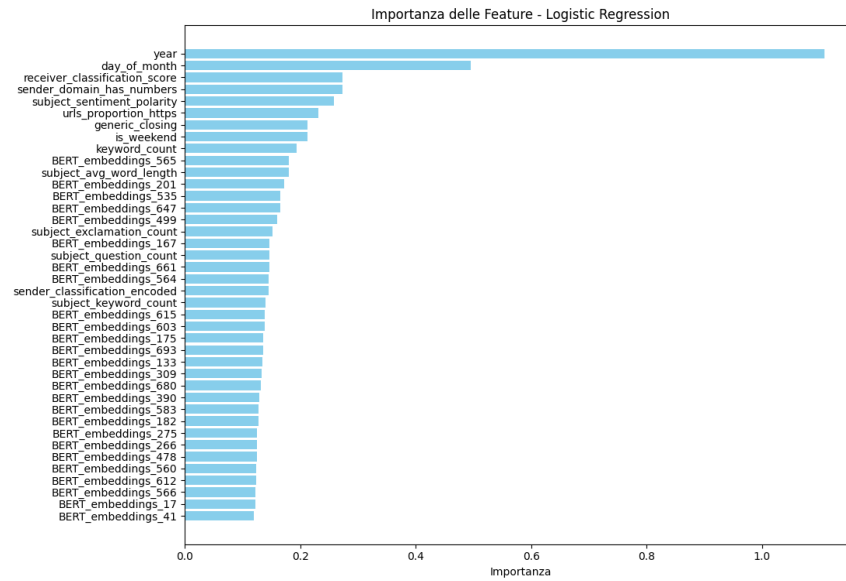


Figura 9 - Importanza feature per Regressione Logistica - Human generated - Caso 2

Caso 3 – Solo cleaned\_body testuale:

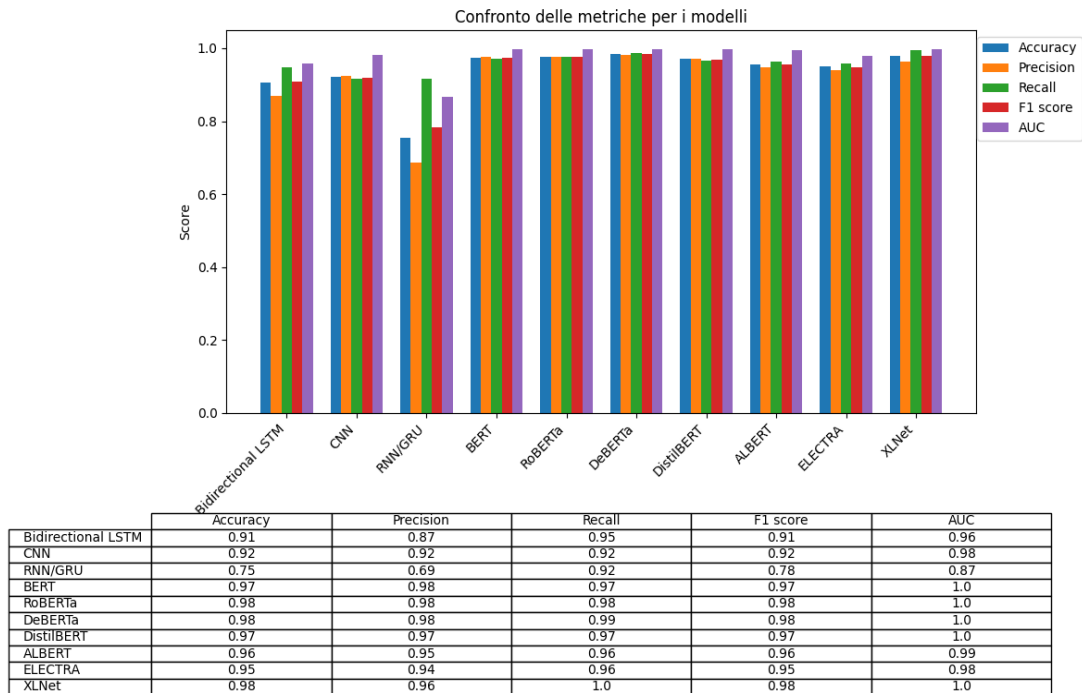


Figure 10 - Confronto metriche - Human generated - caso 3

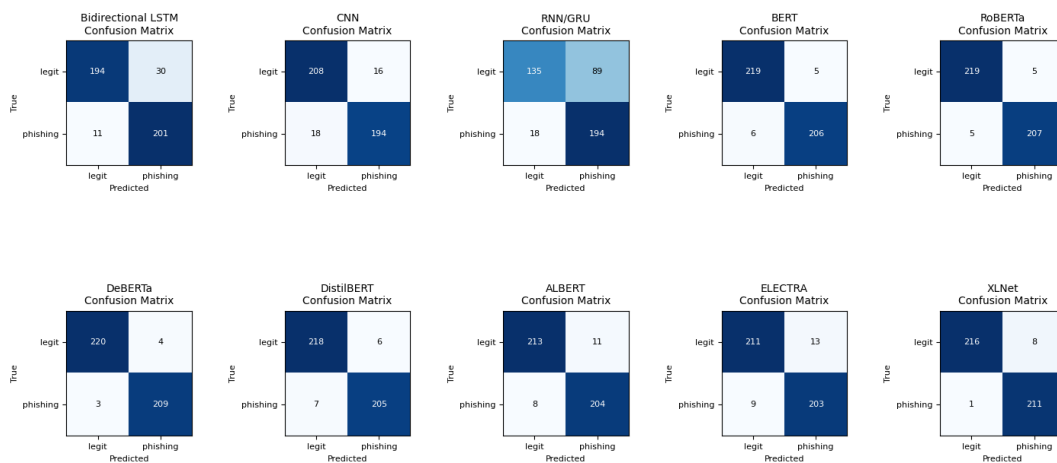


Figura 1 - Matrici di confusione - Human generated - caso 3

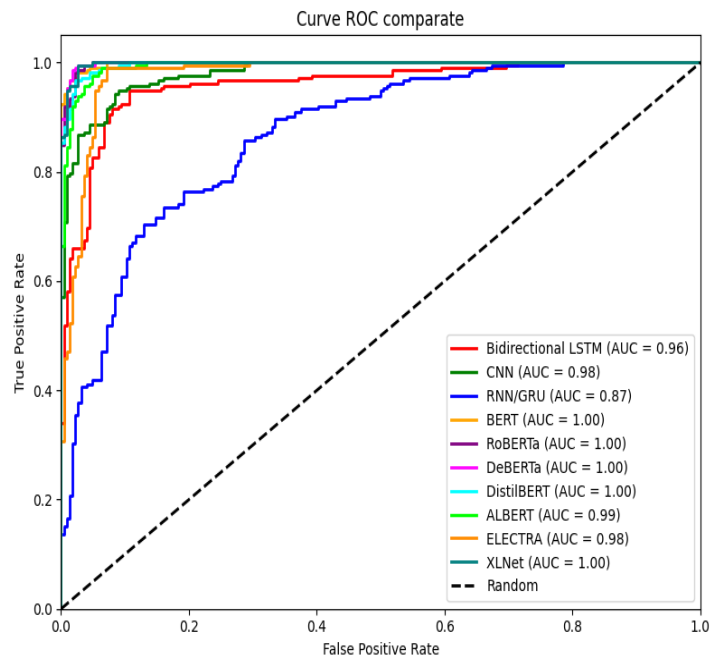


Figura 3 - Curve ROC - Human generated - caso 3

## Prestazioni sul Dataset LLM-generated E-mails

I risultati mostrano una diminuzione delle performance complessive, evidenziando come le e-mail generate da LLM siano meno prevedibili e più difficili da rilevare rispetto a quelle generate da esseri umani:

- **Caso 1:** Tutti i modelli hanno registrato un calo significativo delle prestazioni; le feature numeriche iniziali non sono sufficienti a rilevare i testi generati da LLM;
- **Caso 2:** L'integrazione di feature semantiche migliora nettamente le performance; la Regressione Logistica risulta essere il miglior modello; CNN, SVM e XGBoost mostrano performance molto elevate, mentre LSTM risulta debole;

- **Caso 3:** I modelli LLM (DistilBERT, RoBERTa, XLNet) mostrano le migliori performance, confermando la loro efficacia nel riconoscere testo artificiale; DistilBERT e XLNet si confermano come i modelli più performanti.

Caso 1 – Feature numeriche iniziali:

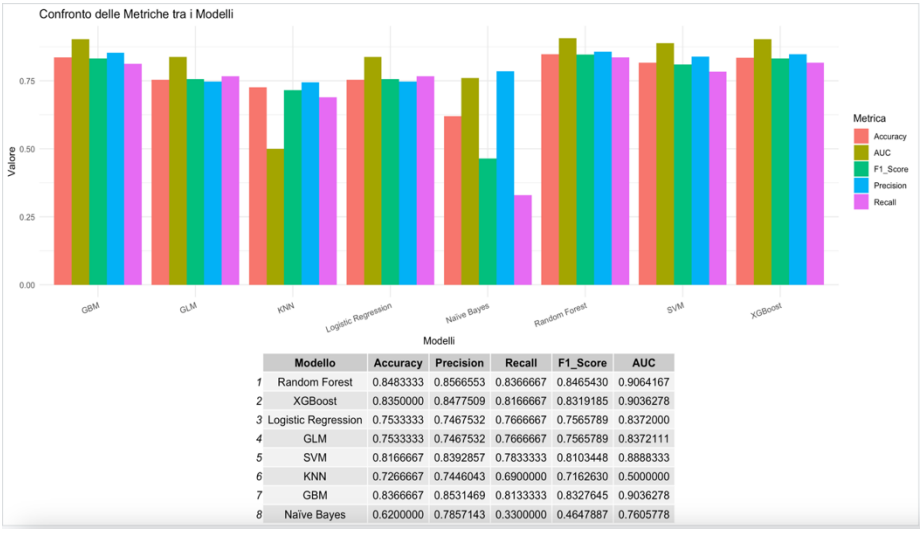


Figura 13 - Confronto metriche - LLM generated - caso 1

Caso 2 – Tutte le feature numeriche:

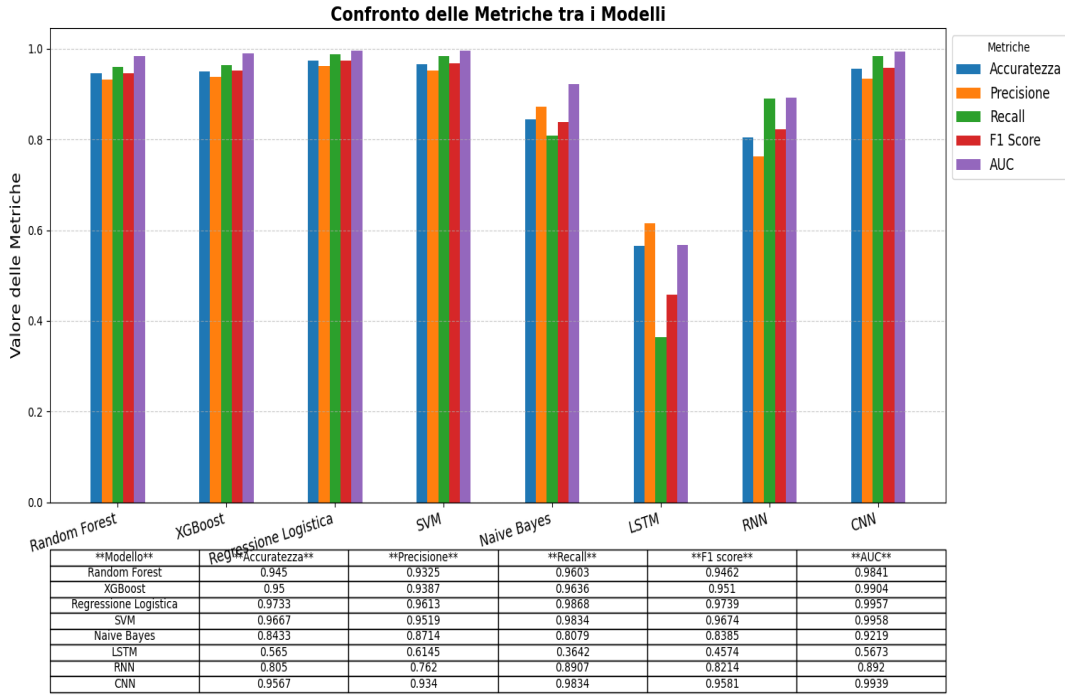


Figura 14 - Confronto metriche - LLM generated - Caso 2

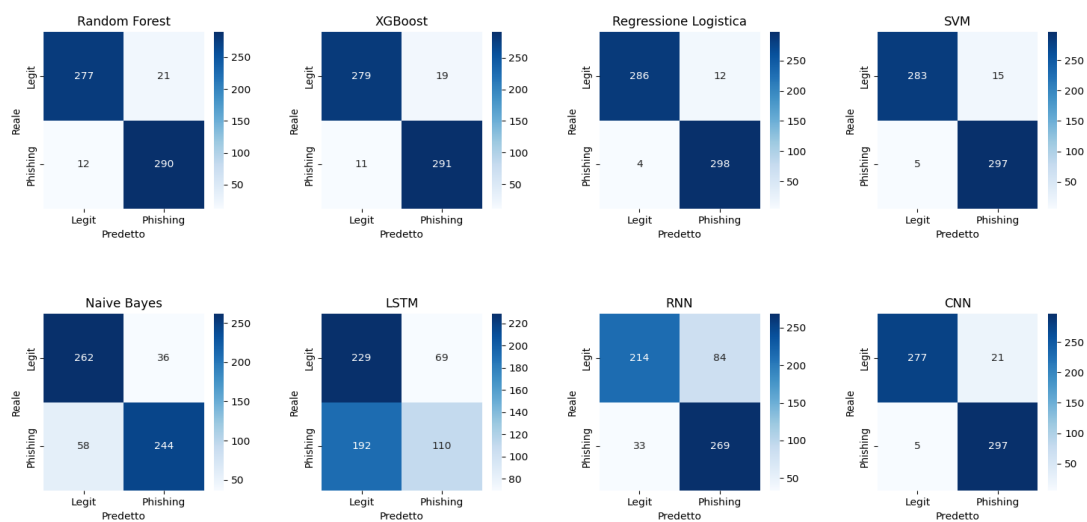


Figura 15 - Matrici di confusione - LLM generated - Caso 2

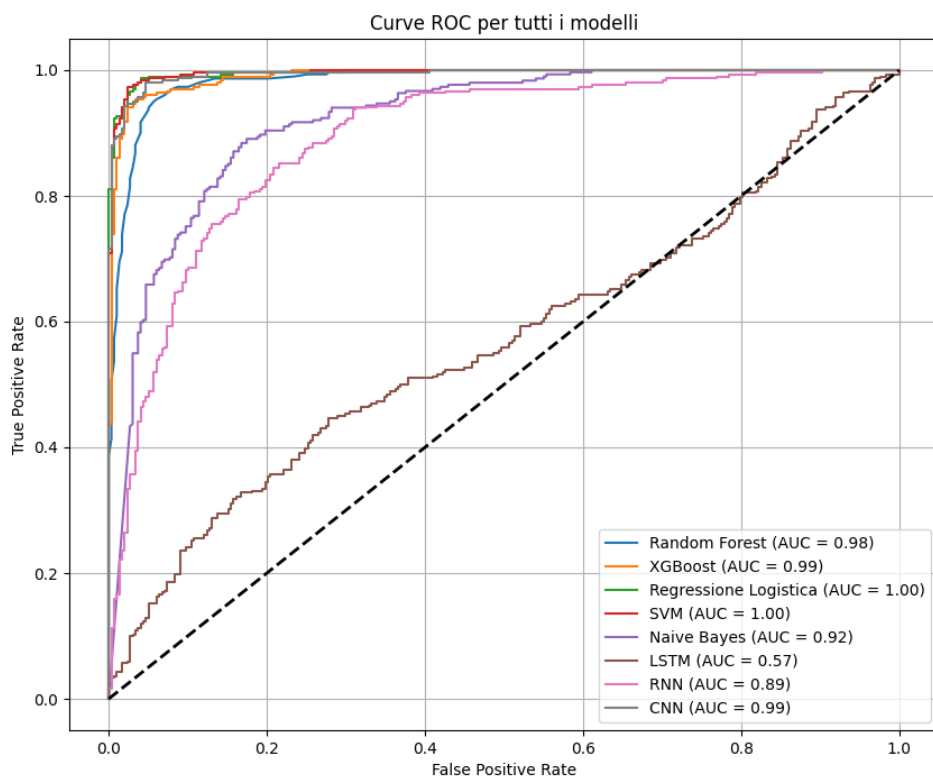


Figura 16 - Curve ROC - LLM generated - Caso 2

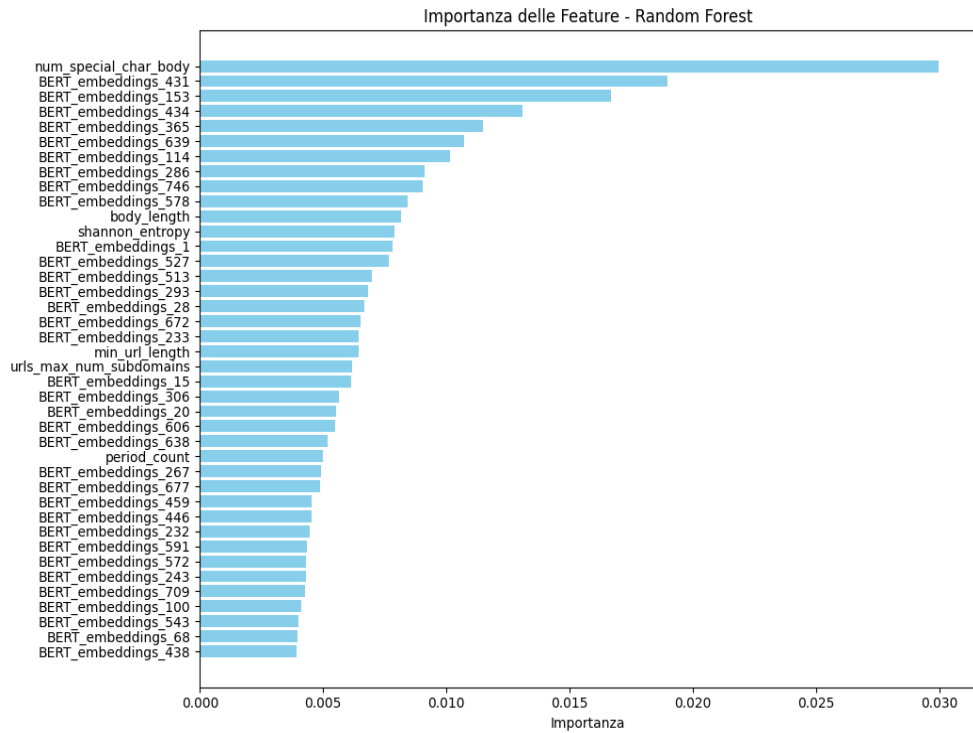


Figura 17 - Importanza delle variabili per Random Forest - LLM generated - Caso 2

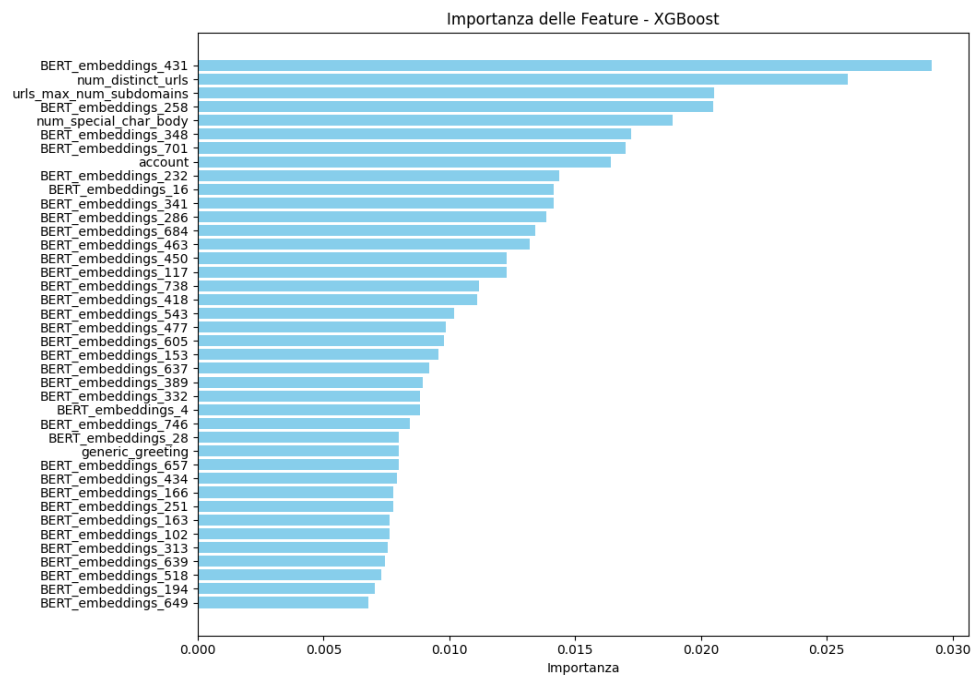


Figura 18 - Importanza variabili per XGBoost - LLM generated - Caso 2

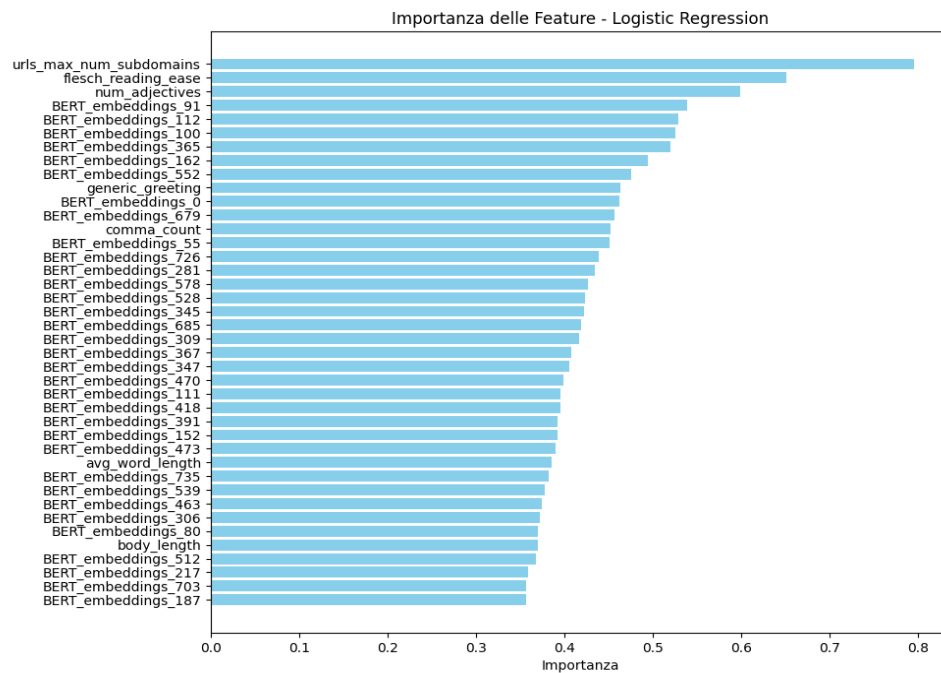


Figura 19 - Importanza variabili per Regressione Logistica - LLM generated - Caso 2

Caso 3 – Solo Body testuale ripulito:

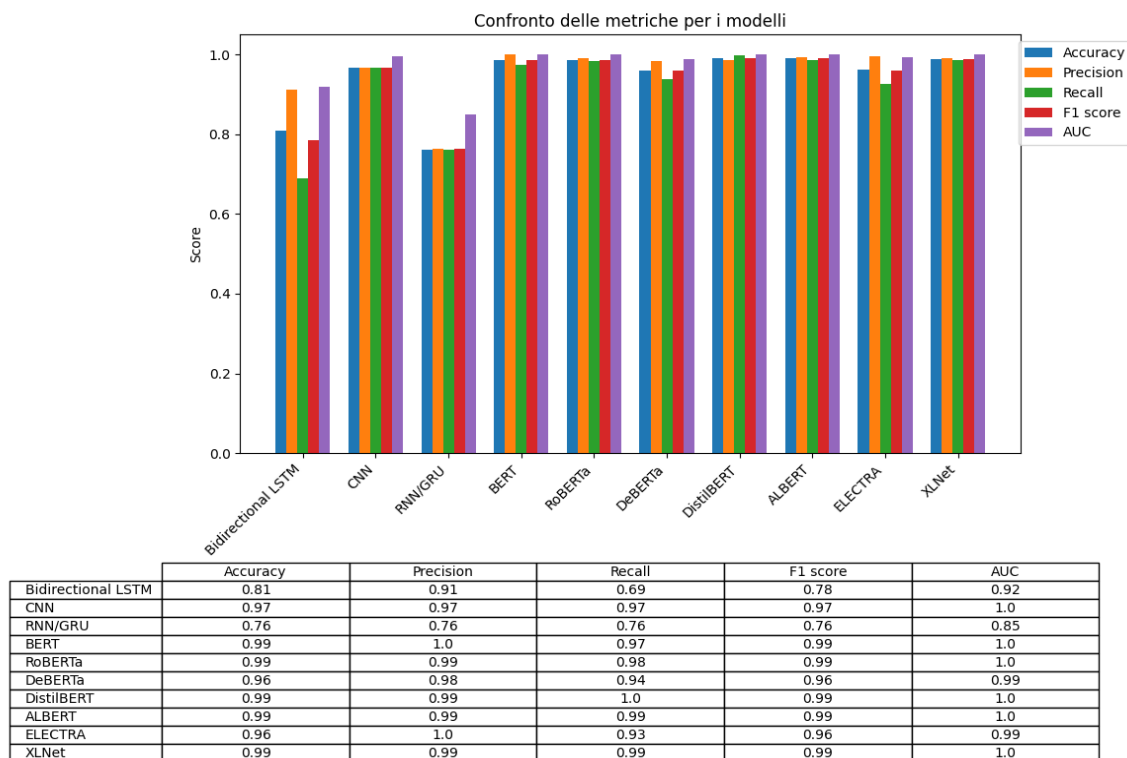


Figura 20 - Confronto metriche - LLM generated - Caso 3

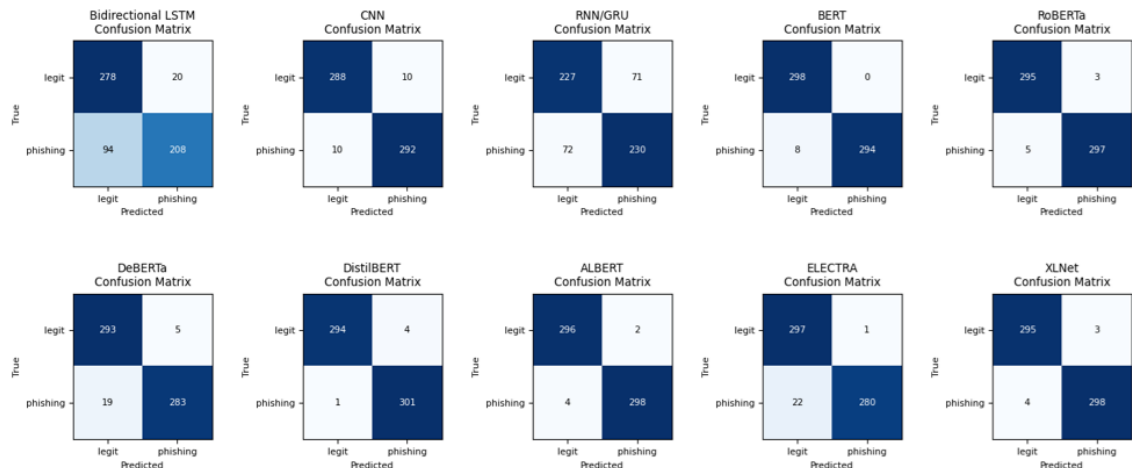


Figura 21 - Matrici di confusione - LLM generated - Caso 3

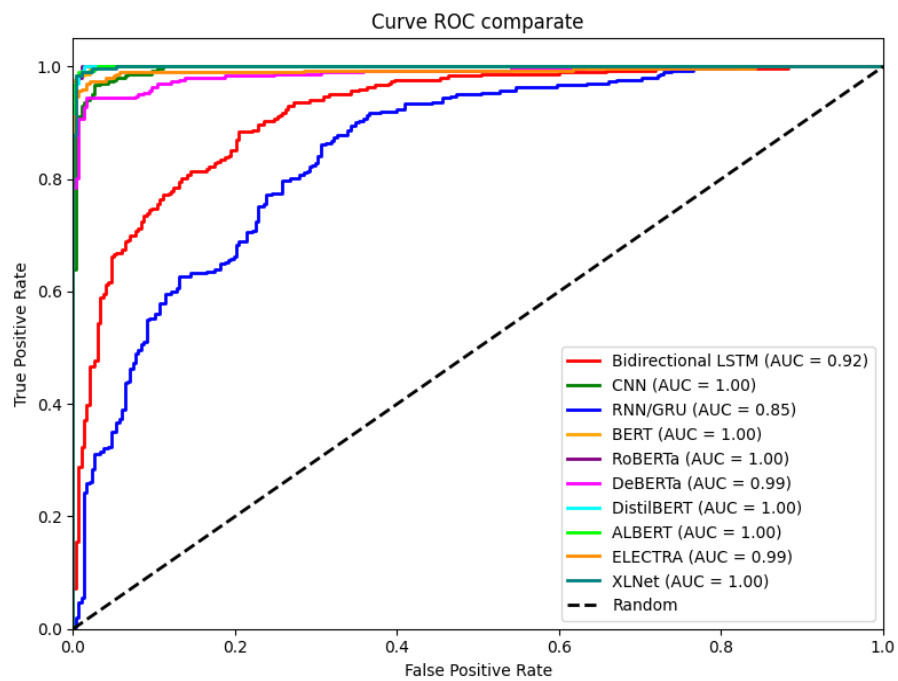


Figura 22 - Curve ROC-AUC - LLM generated - Caso 3