

# **Eksploracja danych**

Agnieszka Kulesz, Hania Gibus, Igor Józefowicz





# Opis zbioru

## *Gun violence data*

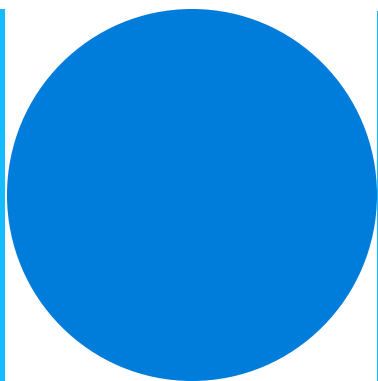
Szczegółowe informacje o incydentach z użyciem broni palnej na terenie Stanów Zjednoczonych w latach 2013–2018

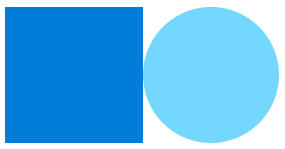


# Cel eksploracji

---

- Predykcja czy incydent z użyciem broni zakończy się ofiarami śmiertelnymi ( $n_{\text{killed}} > 0$ )
- Czułość  $\geq 80\%$
- Swoistość  $\geq 60\%$
- Klasyfikacja binarna



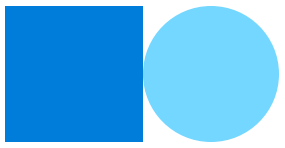


# Dobór algorytmu eksploracji

Wybrany algorytm - drzewo decyzyjne

- Wysoka interpretowalność
- Jasne reguły decyzyjne
- Identyfikacja progów wartości





# Dobór metody testowania wyników

Krosvalidacja 10-krotna

Korzyści:

- Eliminacja wpływu przypadkowego podziału
- Uśrednienie wyników
- Wyższa wiarygodność oceny



# Przygotowanie danych

Braki, transformacja i uzupełnienie, podzbiór danych

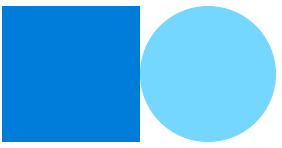




**Braki danych:** Brakujące wartości w kolumnach tekstowych (np. `gun_stolen`, `participant_age_group`, `participant_gender`) zostały ujednolicone jako *"Unknown"*

**Transformacja danych:** Kolumny ze złożonymi łańcuchami tekstowymi (np. `gun_type`) zostały rozbite na unikalne kategorie i zakodowane jako zmienne nominalne.





# Podzbiór danych

Wybrane cechy

- `n_injured` – liczba rannych
- `n_guns_involved` – liczba użytej broni
- `state` – stan USA
- `avg_age` – średni wiek uczestników
- `male_count`, `female_count` – liczba uczestników według płci







# Uzupełnienie danych

**Cel analizy:** Przewidzenie, czy incydent zakończy się ofiarą śmiertelną ( $mortality = 1$  jeśli  $n\_killed > 0$ ).

## Problematyka klas niezbalansowanych

- Incydenty bez ofiar śmiertelnych są liczniejsze, ale odzwierciedlają rzeczywistość

## Metody radzenia sobie z nierównowagą klas

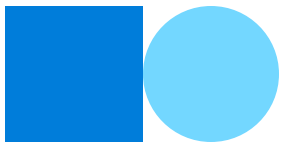
- Wagi klas proporcjonalne do ich częstości
- Metoda **SMOTE** – syntetyczne zwiększenie próbek klasy mniejszościowej



# Utworzenie modelu: Drzewo decyzyjne

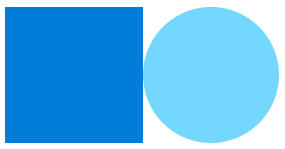
---





# Wykorzystane biblioteki





# Najważniejsze importy

```
from sklearn.model_selection import train_test_split,  
GridSearchCV, StratifiedKFold
```

```
from sklearn.preprocessing import OneHotEncoder
```

```
from sklearn.tree import DecisionTreeClassifier, plot_tree
```

```
from sklearn.ensemble import RandomForestClassifier
```

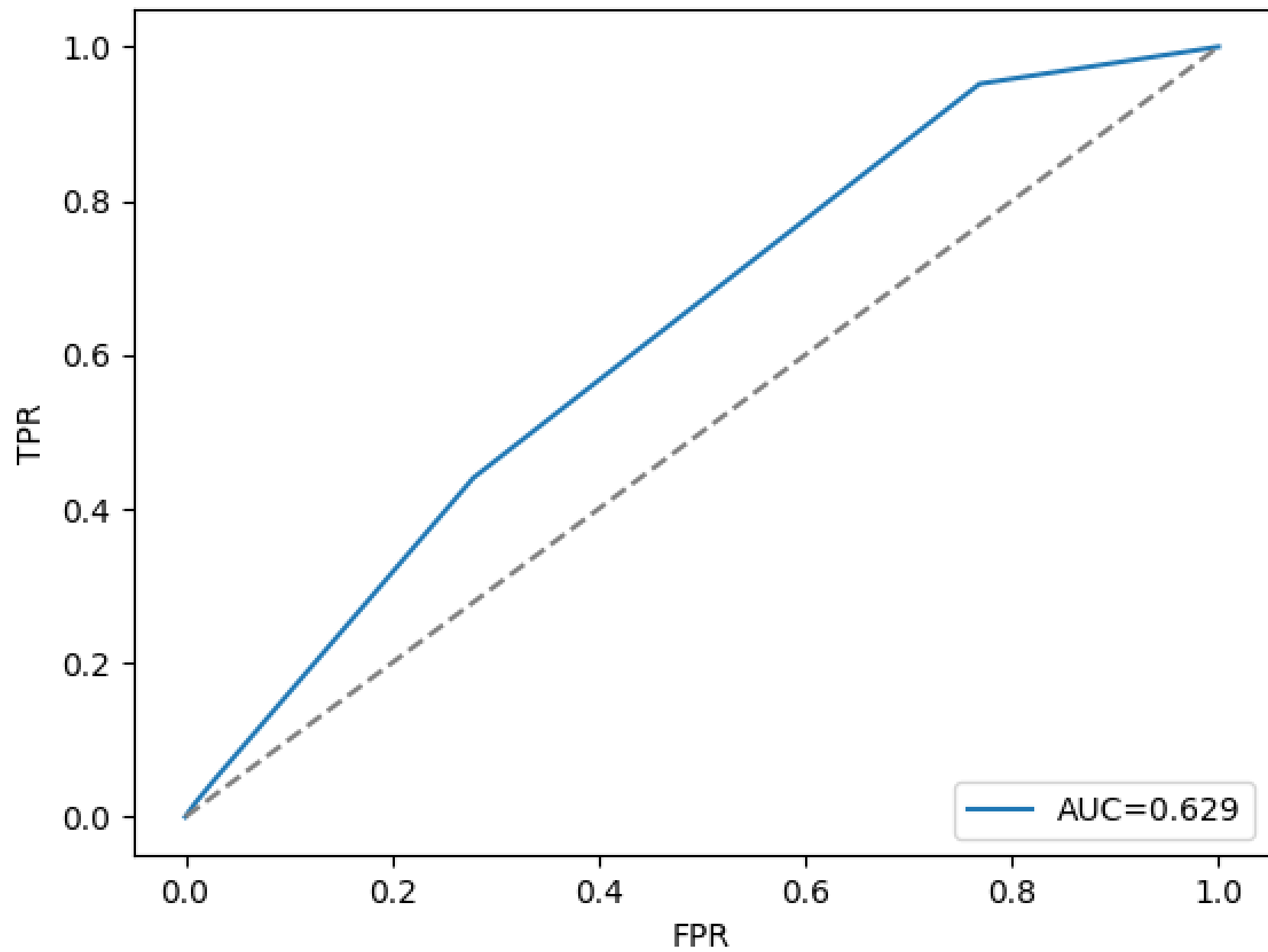




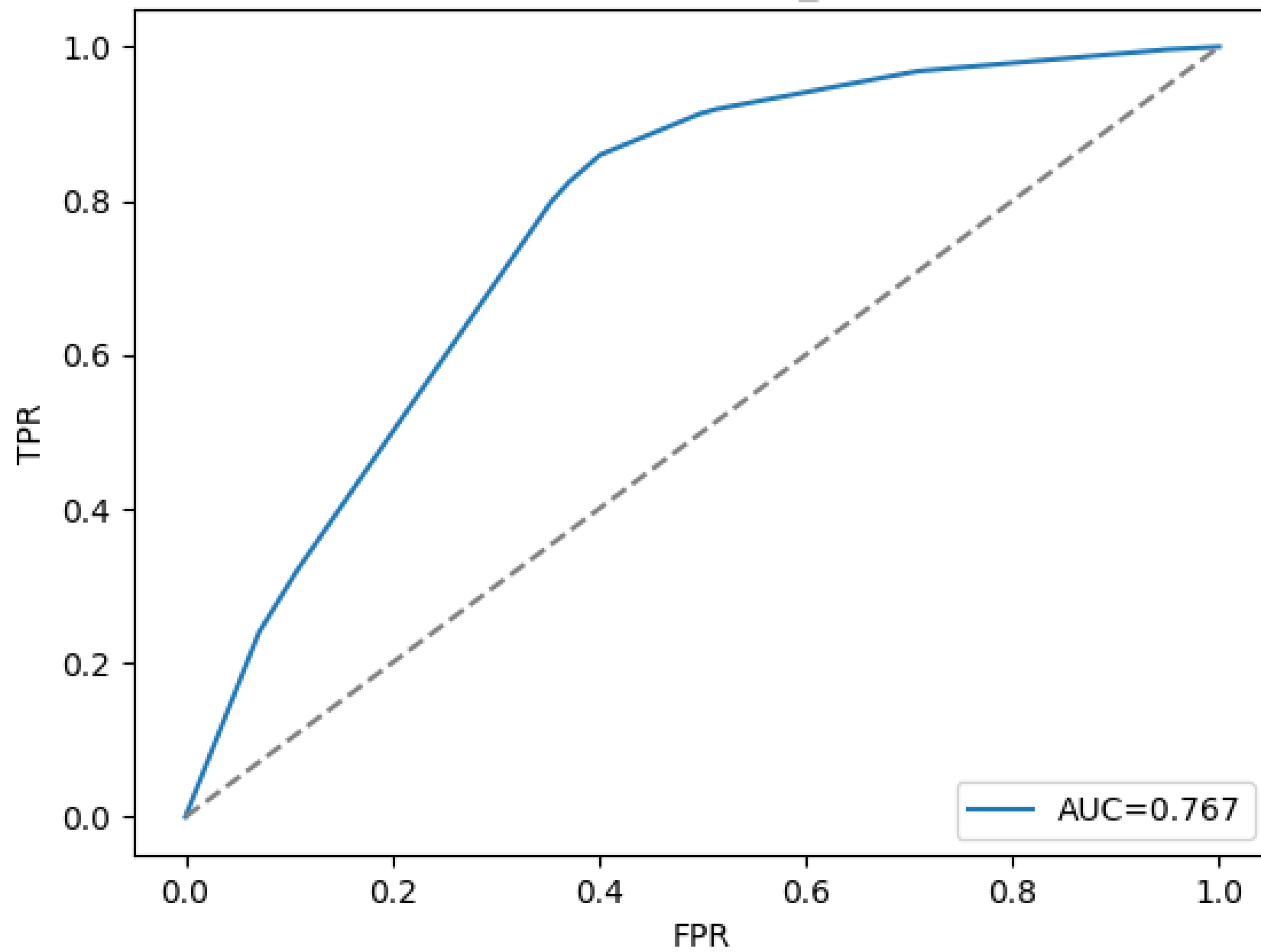
# Wyniki i eksperymenty z modelem i zbiorem danych



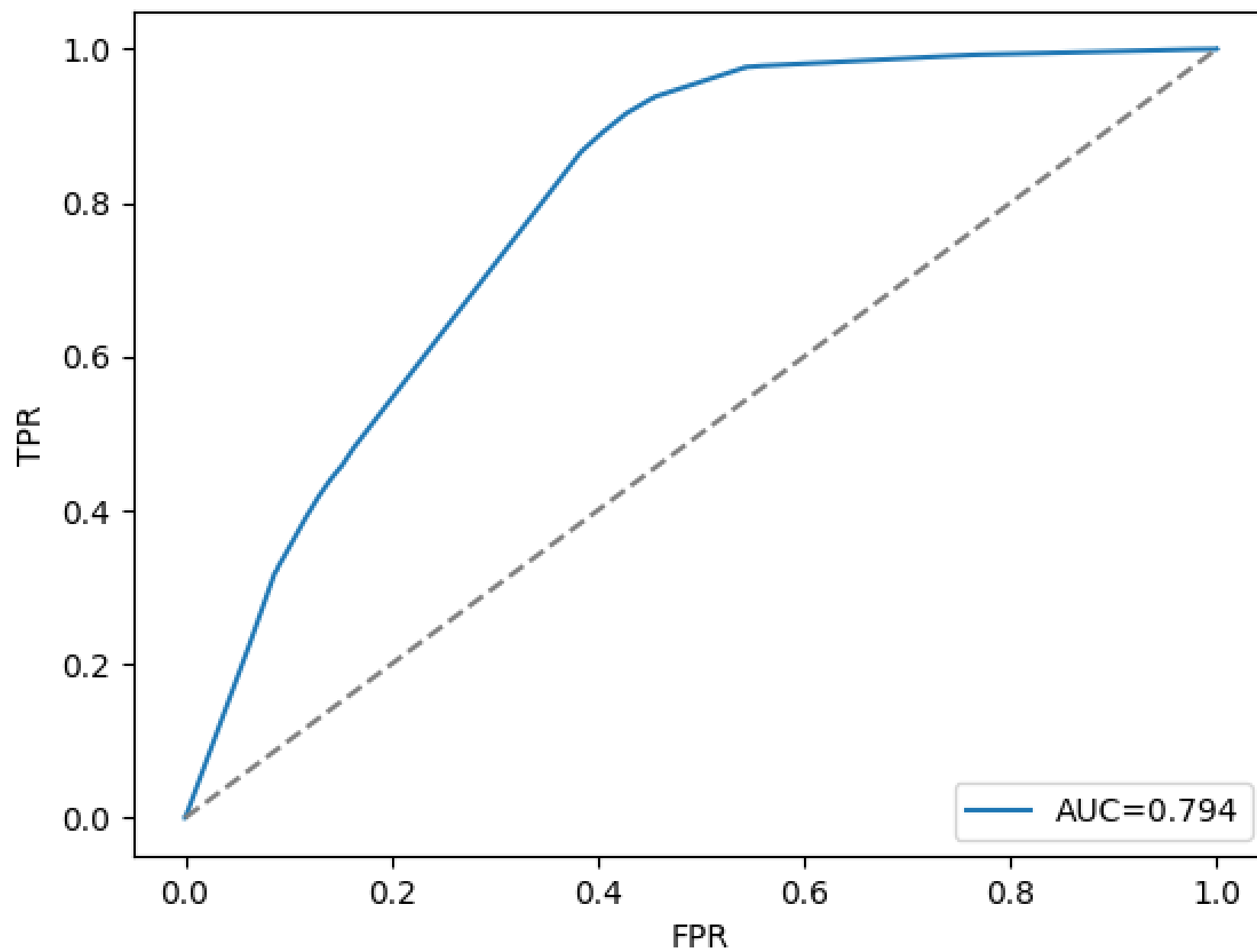
DT ROC - male



DT ROC - male\_ninj

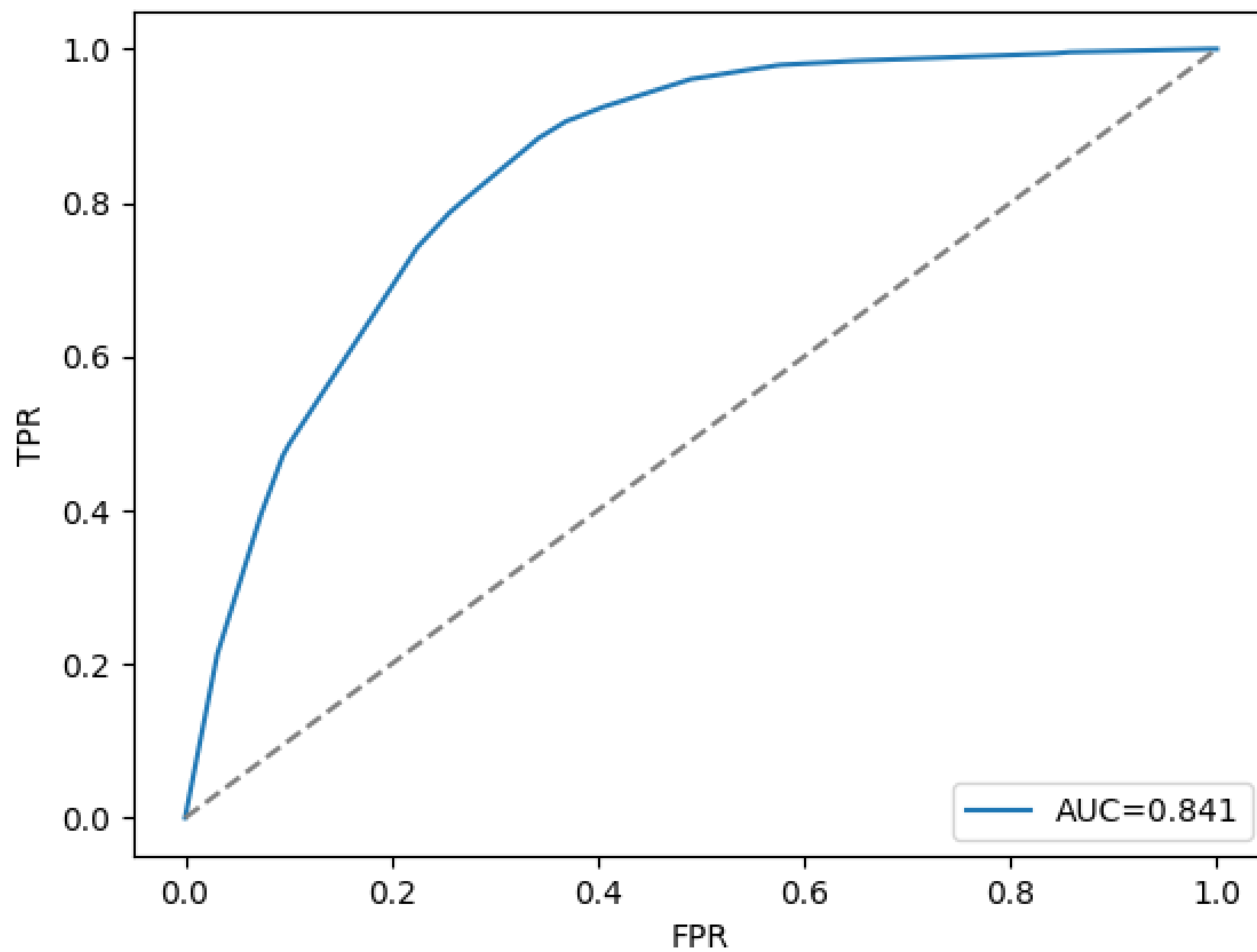


DT ROC - top3

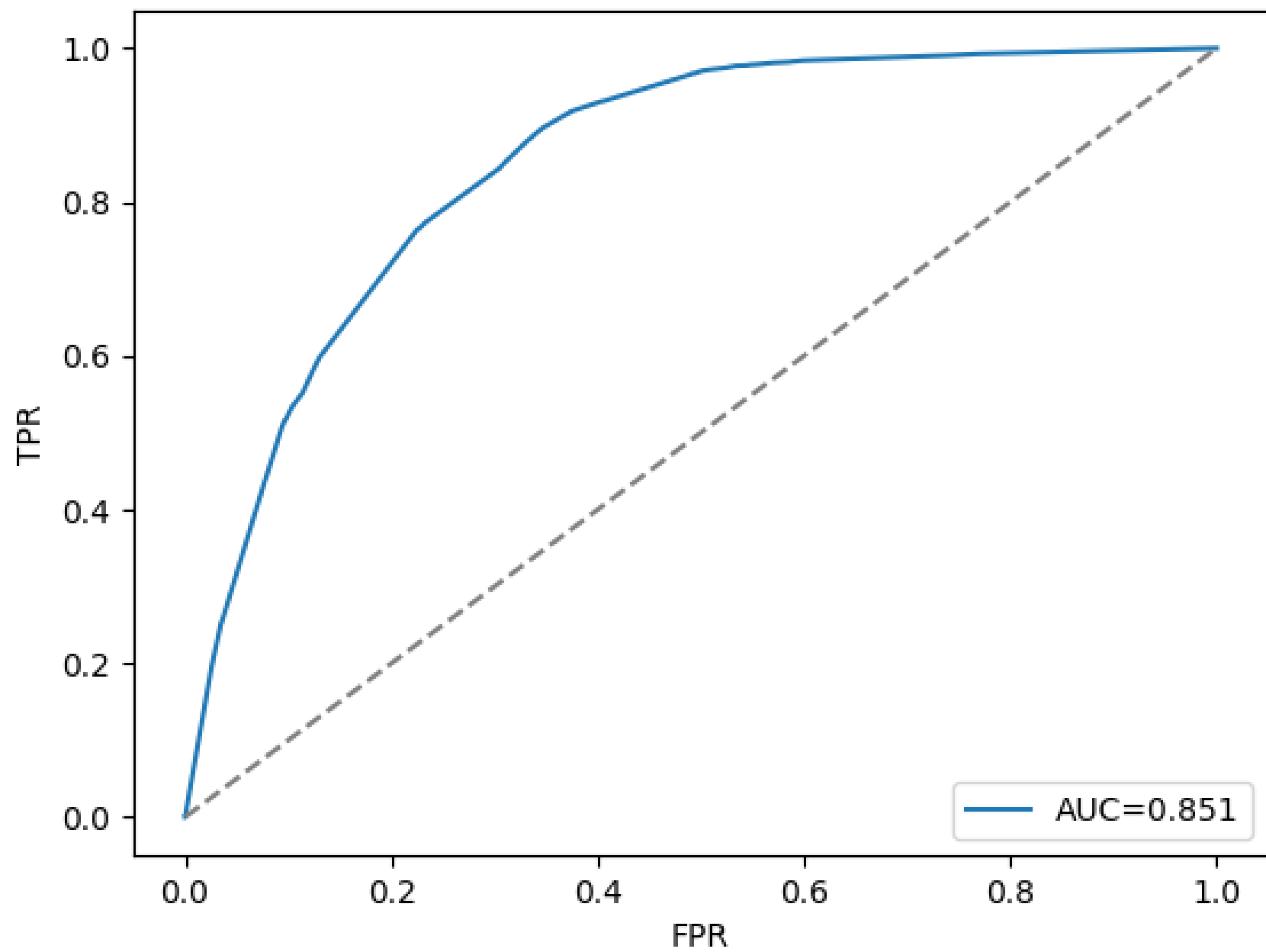


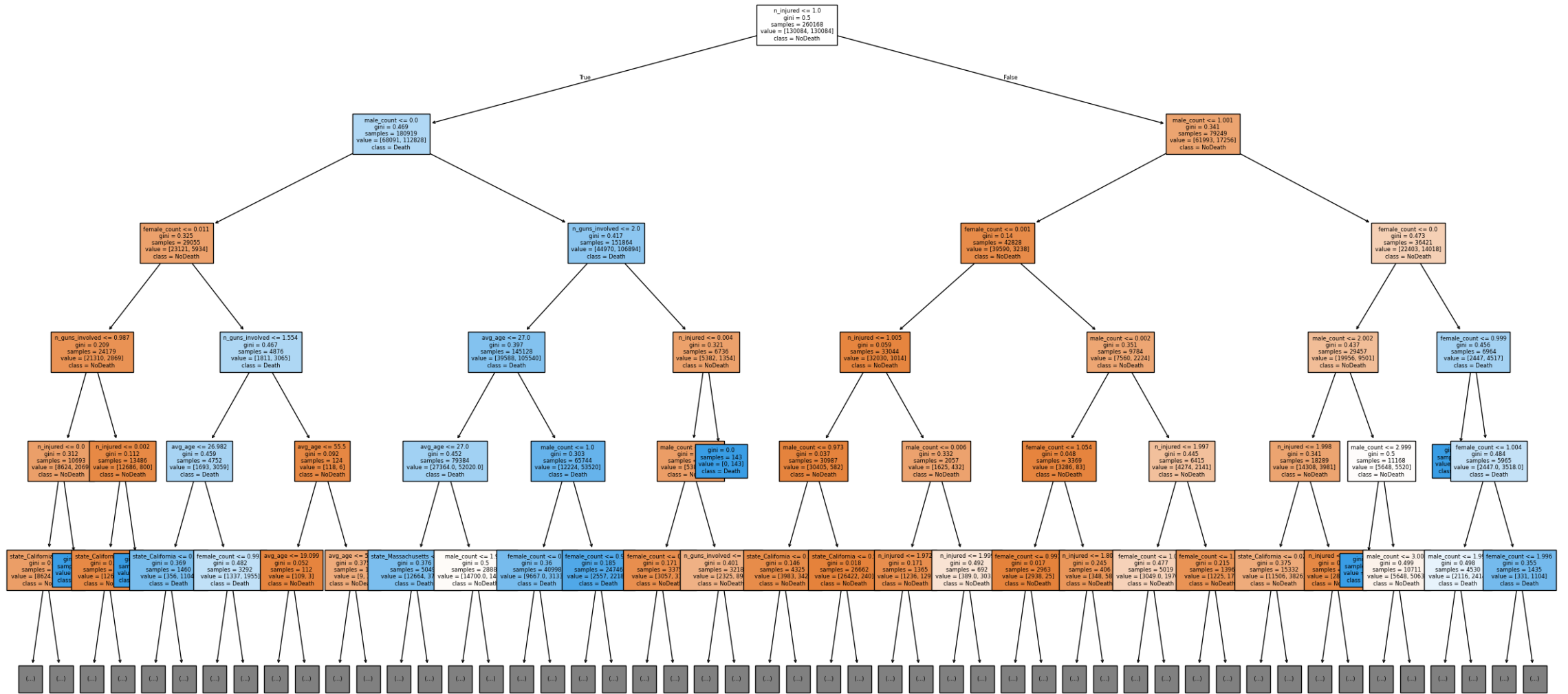


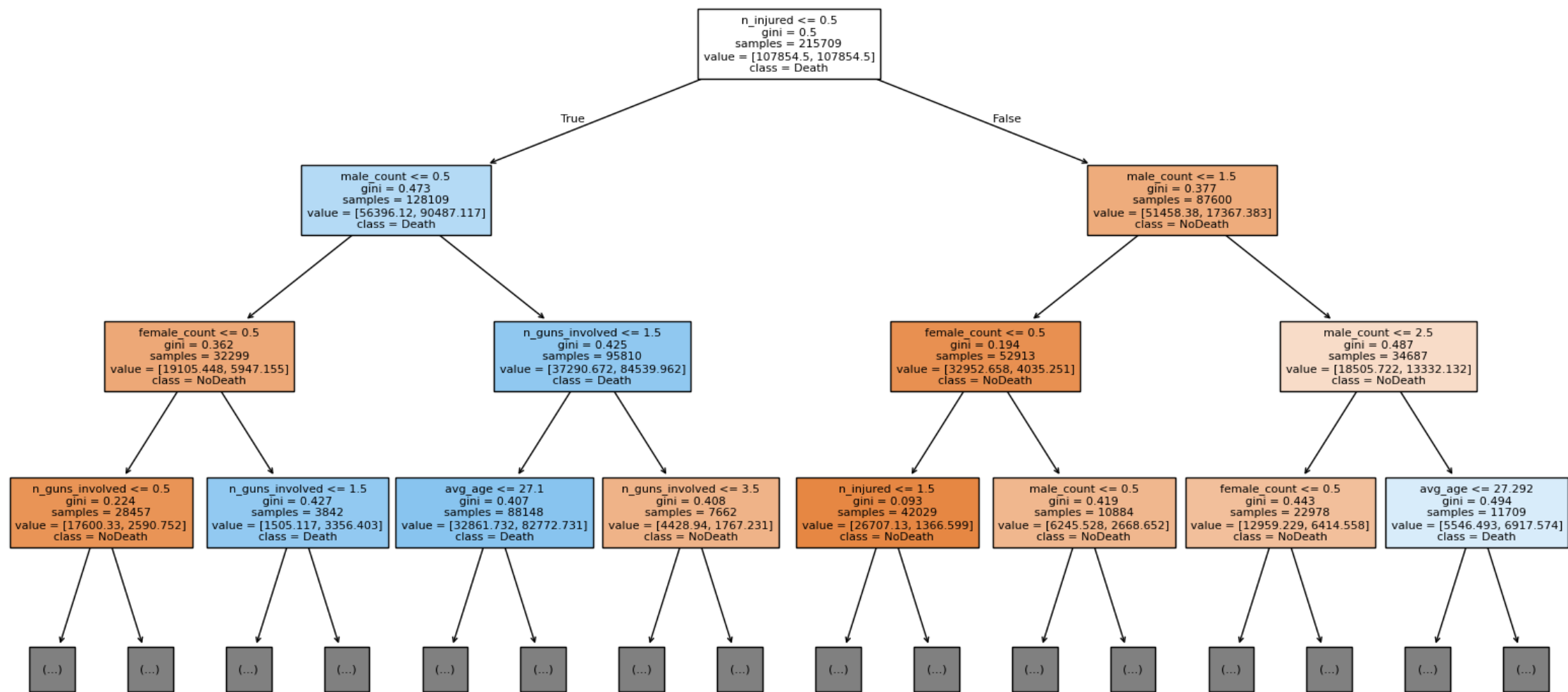
DT ROC - top5

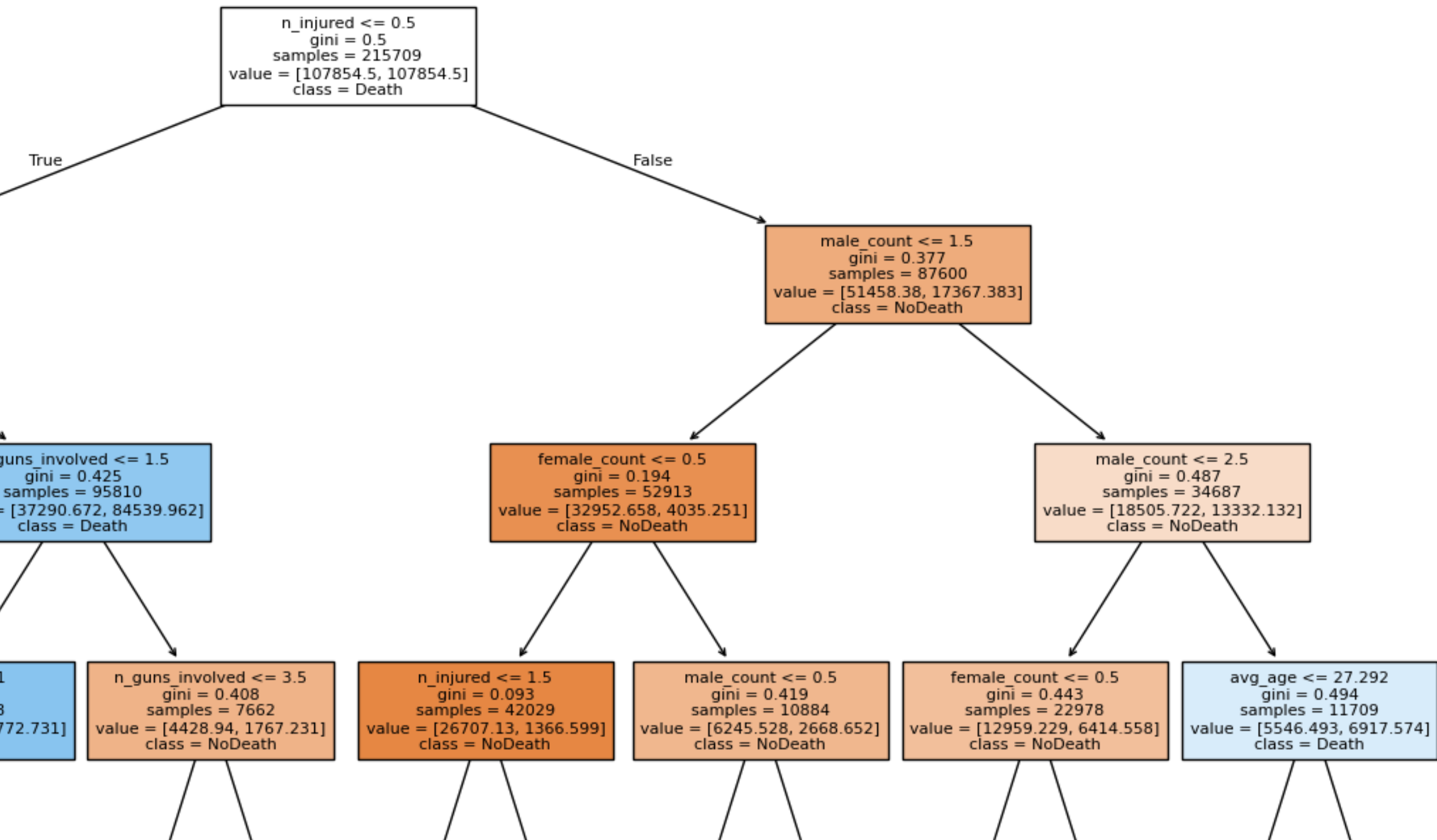


DT ROC - all

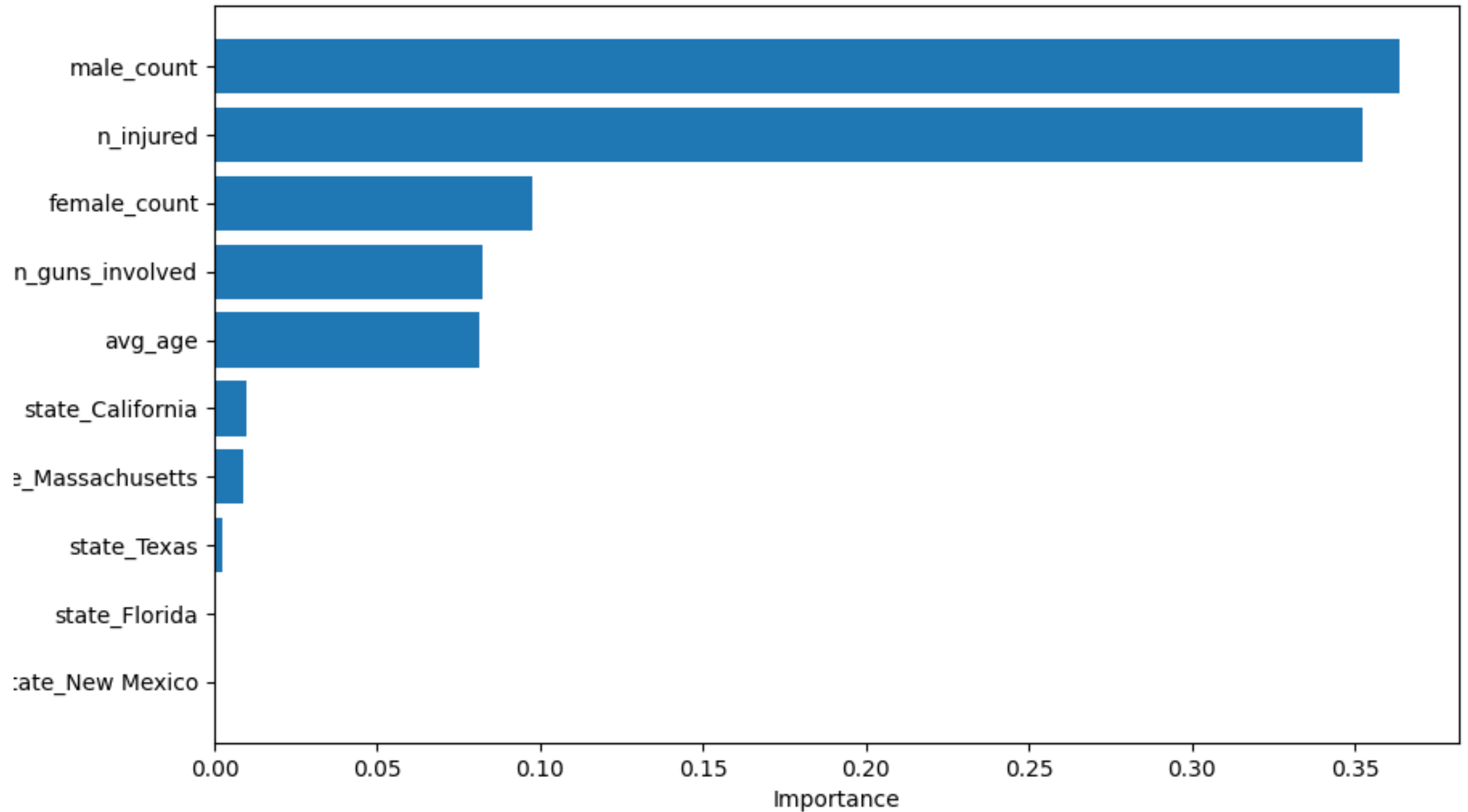








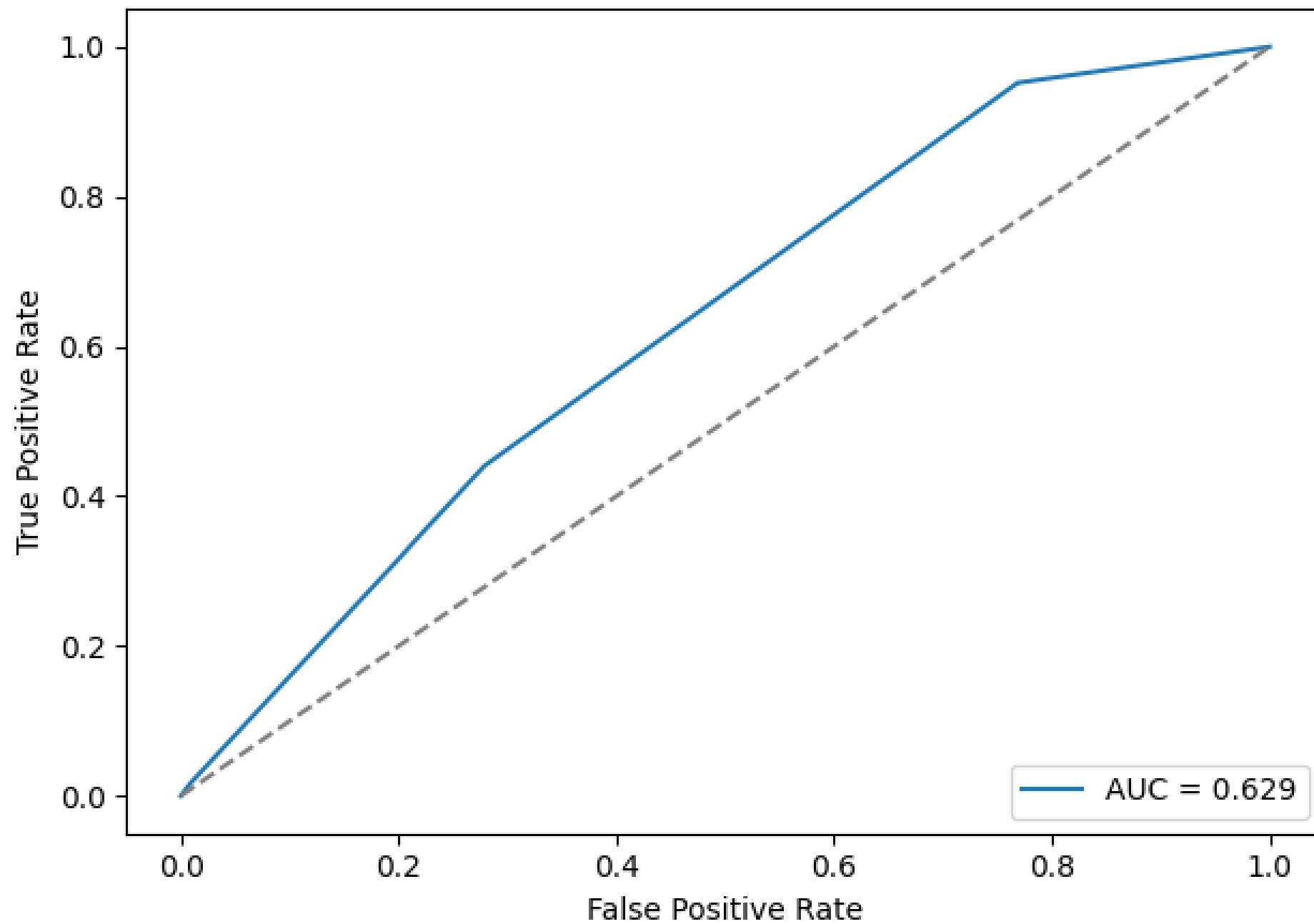
Top 10 Feature Importances





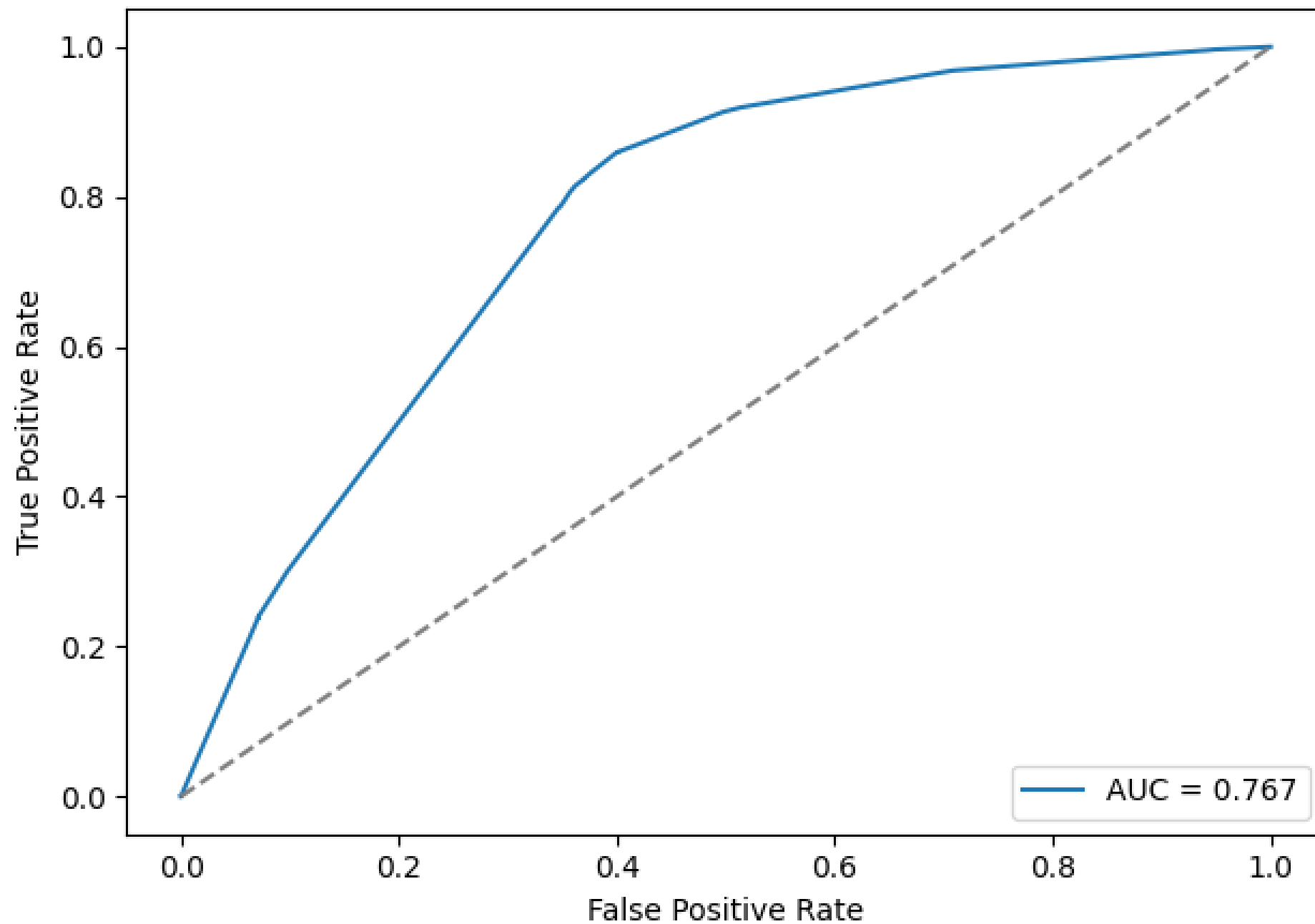
# Las losowy

RandomForest ROC - top1

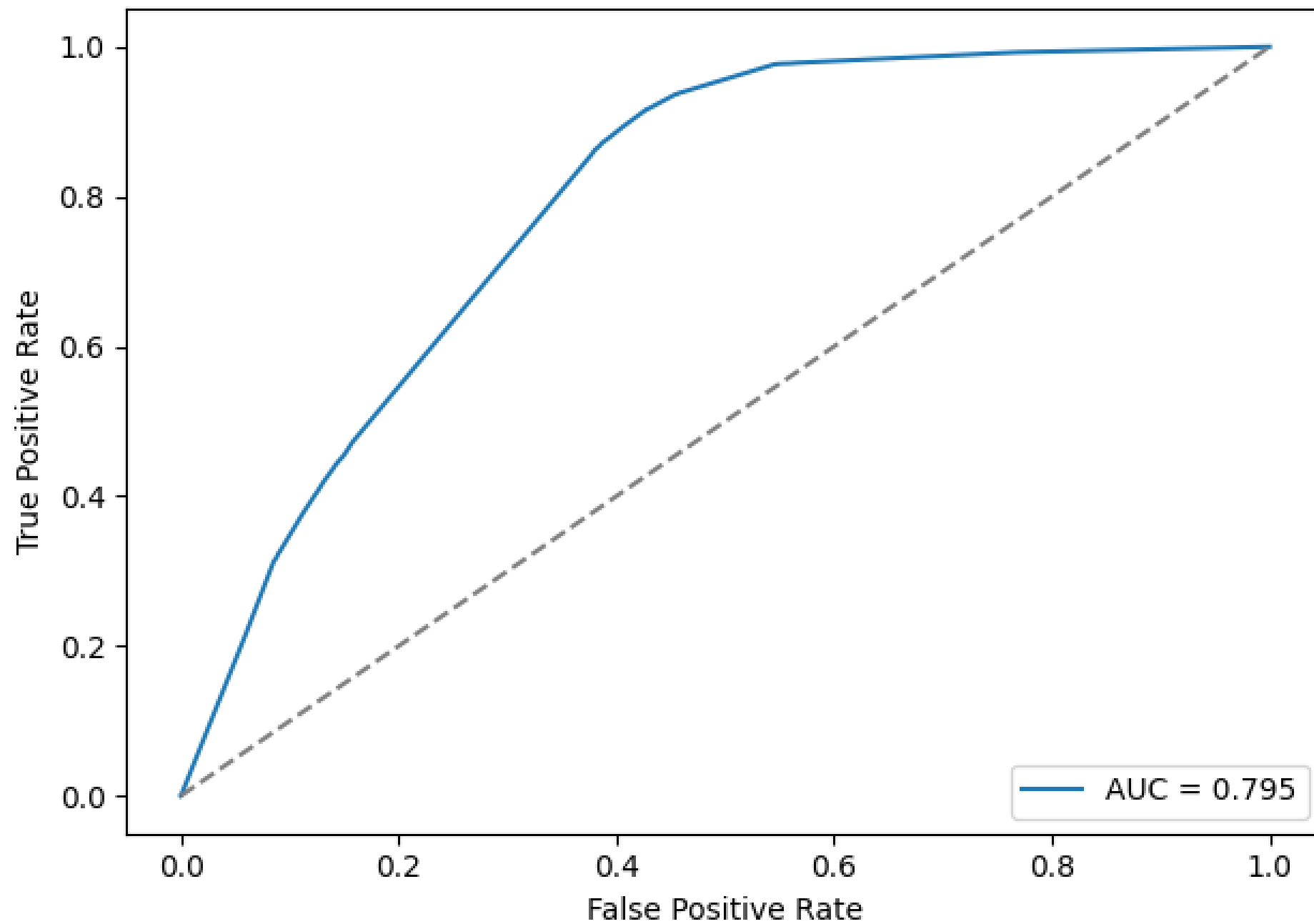




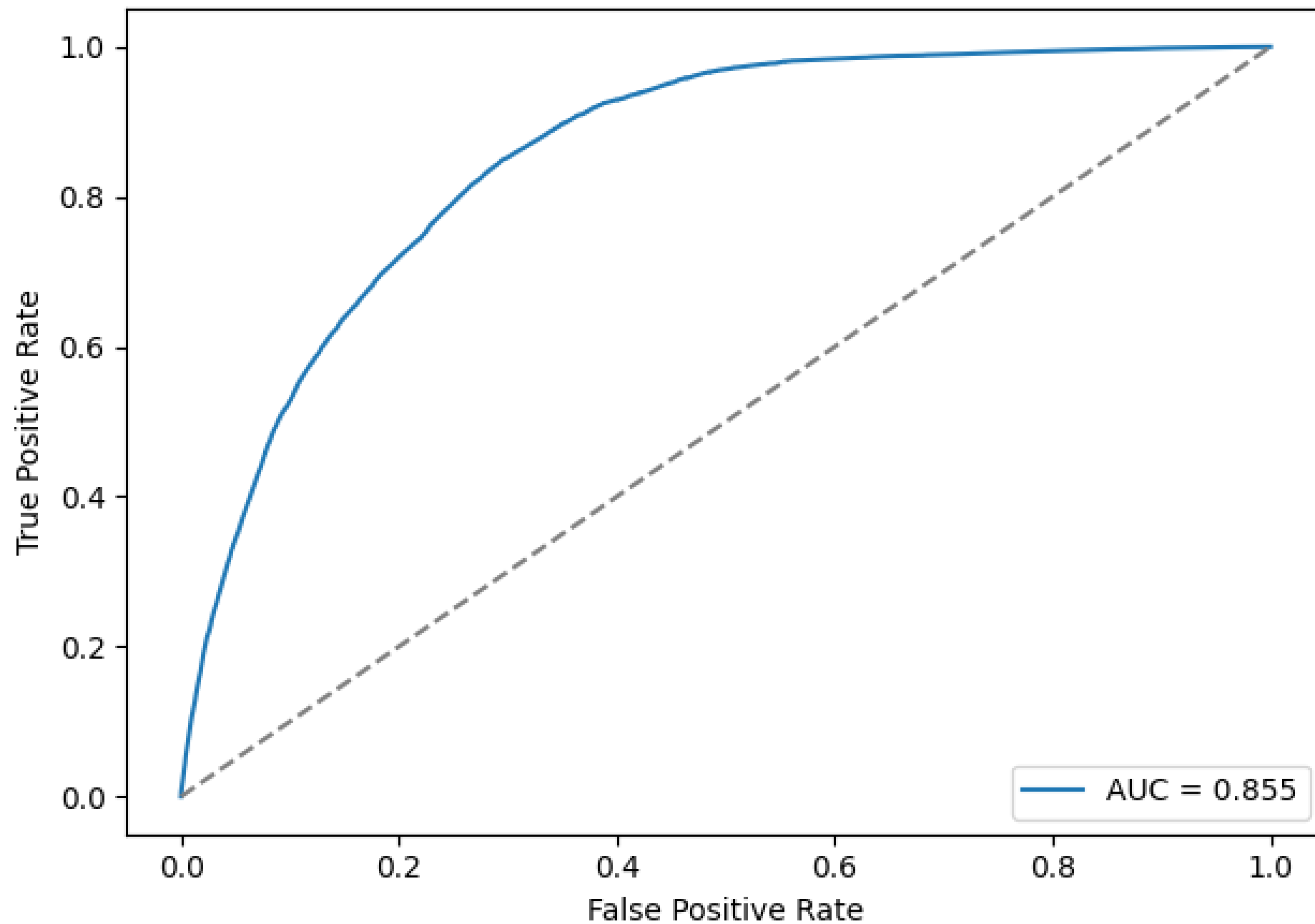
RandomForest ROC - top2



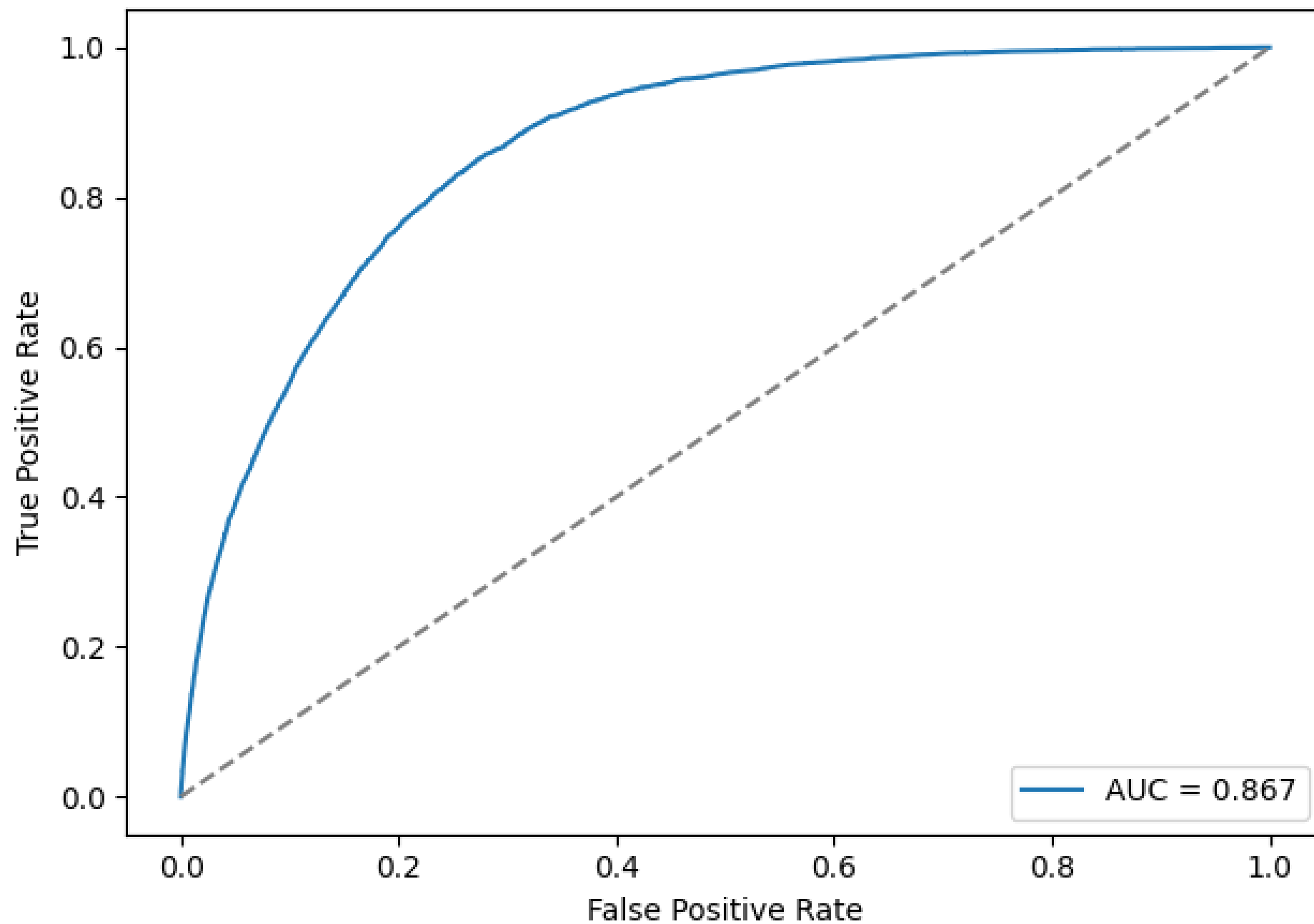
RandomForest ROC - top3



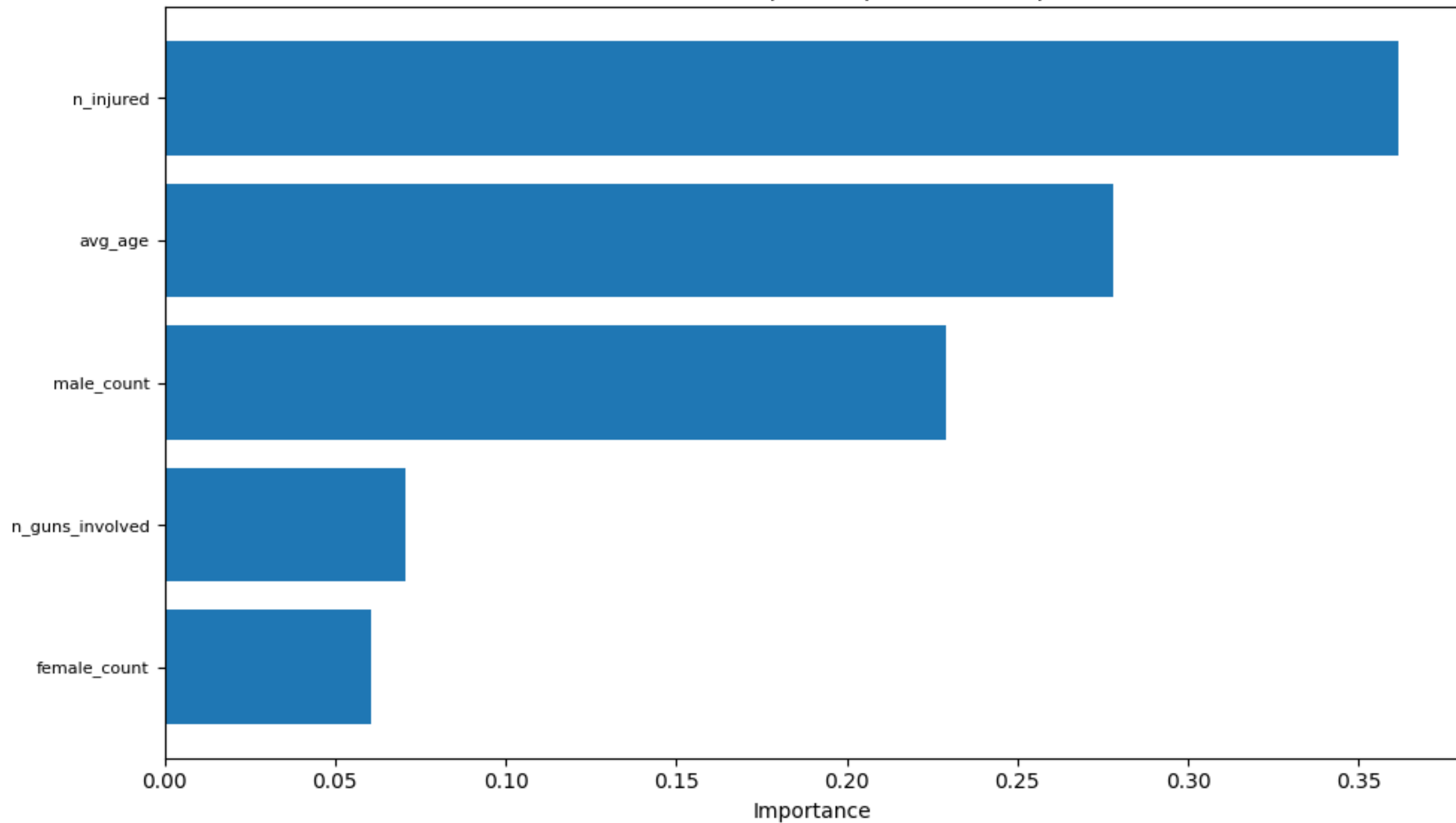
RandomForest ROC - top5



RandomForest ROC - all



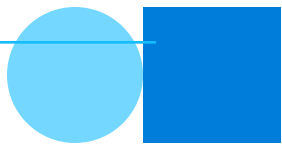
RandomForest Top 10 Importances - top5





# Las losowy - podsumowanie

Zestaw cech	AUC	Recall	Specificity
Wszystkie	0.867	0.830	0.747
Top 1	0.629	0.440	0.721
Top 2	0.767	0.814	0.638
Top 3	0.795	0.863	0.620
Top 5	0.855	0.834	0.718





# Obie metody - podsumowanie

Zestaw cech	Decision Tree (AUC / Recall / Spec)	Random Forest (AUC / Recall / Spec)
Wszystkie	0.851 / 0.772 / 0.768	0.867 / 0.823 / 0.753
Top 1	0.629 / 0.440 / 0.721	0.629 / 0.440 / 0.721
Top 2	0.767 / 0.825 / 0.629	0.767 / 0.814 / 0.639
Top 3	0.794 / 0.865 / 0.617	0.795 / 0.871 / 0.613
Top 5	0.841 / 0.801 / 0.732	0.855 / 0.835 / 0.718



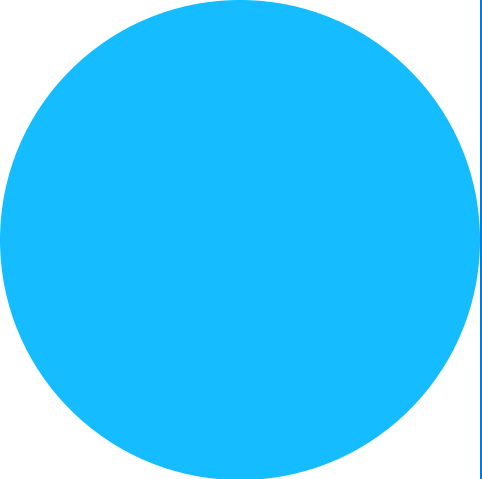


# Podsumowanie projektu ✨ 🎉

- Udało nam się osiągnąć cele eksploracji danych:
  - Czułość (recall)  $\geq 80\%$
  - Swoistość (specificity)  $\geq 60\%$
- Wykonaliśmy 10 eksperymentów z różnymi cechami
  - 5 z drzewami decyzyjnymi
  - 5 z lasami losowymi







# **Dziękujemy za uwagę**

Agnieszka Kulesz, Hania Gibus, Igor Józefowicz

