

**Tim\_5****Branislav Stojković SV64/2020****Anja Petković SV22/2020****Predikcija stadijuma raka pluća**

**Tema projekta:** Predikcija stadijuma raka pluća kod pacijenata korišćenjem mašinskog učenja.

**Opis problema:** Rak pluća je jedna od najsmrtonosnijih bolesti širom sveta. Precizna predikcija stadijuma raka za pacijente sa rakom pluća može značajno pomoći lekarima u ranom otkrivanju bolesti i donošenju odluka o tretmanu. Ovaj projekat ima za cilj da razvije model mašinskog učenja koji može predvideti stadijum raka pluća koristeći različite demografske, medicinske i tretmanske podatke.

**Podaci:** Za potrebe ovog projekta koristićemo sintetički dataset koji pruža detaljne informacije o pacijentima sa rakom pluća. Dataset obuhvata sledeće atribute:

- **Patient\_ID:** Identifikacioni broj pacijenta
- **Age:** Godine starosti pacijenta
- **Gender:** Pol pacijenta
- **Smoking\_History:** Istorija pušenja pacijenta (npr. Current Smoker, Never Smoked, Former Smoker)
- **Tumor\_Size\_mm:** Veličina tumora u milimetrima
- **Tumor\_Location:** Lokacija tumora (npr. Lower Lobe)
- **Stage:** Stadijum raka (npr. Stage I, Stage II, Stage III)
- **Treatment:** Tretman koji pacijent prima (npr. Surgery, Radiation Therapy, Chemotherapy)
- **Survival\_Months:** Broj meseci preživljavanja nakon dijagnoze
- **Ethnicity:** Etnička pripadnost pacijenta
- **Insurance\_Type:** Tip osiguranja (npr. Medicare, Private, Other)
- **Family\_History:** Porodična istorija raka (Yes/No)
- **Comorbidity\_Diabetes:** Prisustvo dijabetesa (Yes/No)
- **Comorbidity\_Hypertension:** Prisustvo hipertenzije (Yes/No)
- **Comorbidity\_Heart\_Disease:** Prisustvo srčanih oboljenja (Yes/No)
- **Comorbidity\_Chronic\_Lung\_Disease:** Prisustvo hroničnih plućnih bolesti (Yes/No)
- **Comorbidity\_Kidney\_Disease:** Prisustvo bolesti bubrega (Yes/No)
- **Comorbidity\_Autoimmune\_Disease:** Prisustvo autoimunih bolesti (Yes/No)
- **Comorbidity\_Other:** Prisustvo drugih komorbiditeta (Yes/No)
- **Performance\_Status:** Performans status pacijenta (skala od 0 do 4)
- **Blood\_Pressure\_Systolic:** Sistolički krvni pritisak
- **Blood\_Pressure\_Diastolic:** Dijastolički krvni pritisak
- **Blood\_Pressure\_Pulse:** Puls
- **Hemoglobin\_Level:** Nivo hemoglobina
- **White\_Blood\_Cell\_Count:** Broj belih krvnih zrnaca
- **Platelet\_Count:** Broj trombocita
- **Albumin\_Level:** Nivo albumina

- **Alkaline\_Phosphatase\_Level:** Nivo alkalne fosfataze
- **Alanine\_Aminotransferase\_Level:** Nivo alanin aminotransferaze
- **Aspartate\_Aminotransferase\_Level:** Nivo aspartat aminotransferaze
- **Creatinine\_Level:** Nivo kreatinina
- **LDH\_Level:** Nivo laktat dehidrogenaze
- **Calcium\_Level:** Nivo kalcijuma
- **Phosphorus\_Level:** Nivo fosfora
- **Glucose\_Level:** Nivo glukoze
- **Potassium\_Level:** Nivo kalijuma
- **Sodium\_Level:** Nivo natrijuma
- **Smoking\_Pack\_Years:** Broj godina pušenja (izraženo u pakovanjima godišnje)

Link do dataset-a: <https://www.kaggle.com/datasets/rashadrmammadov/lung-cancer-prediction>

### Metodologija:

#### 1. Prikupljanje i analiza podataka:

- Preuzimanje dataset-a i pregled strukture podataka.
- Obrada podataka: čišćenje podataka, rukovanje nedostajućim vrednostima, kodiranje kategorijalnih podataka i normalizacija numeričkih atributa.
- Moguća redukcija dimenzionalnosti.
- Podela inicijalnog skupa podataka na skup za trening i skup za test u razmeri 70:30.

#### 2. Model:

- Korišćenje ansambl modela (Bagging, Stacking, Voting) za izradu prediktivnog modela.
- Optimizacija hiperparametara korišćenjem kros validacije.

#### 3. Evaluacija modela:

- Evaluacija performansi modela korišćenjem mikro F1 skora.
- Evaluacija performansi modela korišćenjem makro F1 skora.
- Evaluacija ponašanja modela (Precision, Recall i F1 score) po klasama.

#### 4. Implementacija i testiranje:

- Testiranje modela na nepoznatim podacima i procena njegove generalizacije.

**Evaluacija:** Evaluacija rešenja će se vršiti korišćenjem sledećih metrika:

- **Mikro F1 skor:** Harmonijska sredina preciznosti i odziva, uzimajući u obzir sve instance. Ova metrika pruža balans između preciznosti i odziva na mikro nivou, tj. na nivou svih instanci zajedno.
- **Makro F1 skor:** Harmonijska sredina preciznosti i odziva za svaku klasu posebno, pružajući balans između preciznosti i odziva na nivou svake klase.