

Tehnička Dokumentacija Projekta: Klasifikacija Tumora Dojke

Sadržaj

1. Uvod.....	2
2. Cilj Projekta.....	3
3. Arhitektura Podataka i Karakteristike.....	4
Srednje Vrednosti (Mean Values):.....	4
Standardne Greške (Standard Error - SE):.....	4
Najgore (Worst) Vrednosti:.....	4
4. Metodologija Klasifikacije: Detaljan Postupak Obuke Modela.....	5
4.1. Faze Projekta:.....	5
5. Očekivani Izlazi i Rezultati.....	7
5.1. Metrike Evaluacije Modela (za svaki model).....	7
5.2. Vizuelna Analiza Važnosti Karakteristika (za svaki model, ako je primenjivo).....	7
Logistička Regresija.....	7
Slučajna šuma.....	8
Model K-Najbližih komšija.....	9
Model potpornih vektora(SVM):.....	10
5.3. Interpretacija Rezultata i Preporuke (za svaki model):.....	11
Logistička regresija.....	11
Slučajna šuma.....	11
K-Najbližih komšija.....	12
Model potpornih vektora(SVM).....	12
6. Zaključak.....	13

1. Uvod

Ovaj dokument pruža sveobuhvatan tehnički pregled projekta klasifikacije tumora dojke. Cilj projekta je razvoj prediktivnog modela mašinskog učenja za određivanje malignosti (benigni/maligni) tumora, koristeći numeričke karakteristike ćelijskih jezgara dobijenih iz digitalizovanih slika biopsije. Dokumentacija je strukturirana tako da jasno opisuje svaki korak u procesu obuke modela, naglašavajući metodologiju i pristup, a ne samo krajnje performanse, u duhu principa tehničke dokumentacije usmerene na proizvod.

2. Cilj Projekta

Primarni cilj ovog projekta je uspostavljanje robustnog procesa za razvoj i implementaciju modela mašinskog učenja za preciznu klasifikaciju tumora dojke. Sekundarni ciljevi obuhvataju:

- **Standardizacija Procesnih Koraka:** Jasno definisanje faza od prikupljanja podataka do evaluacije modela.
- **Analiza Karakteristika:** Identifikacija i razumevanje uticaja pojedinačnih karakteristika na klasifikaciju.
- **Vizualizacija Podataka:** Kreiranje vizuelnih prikaza za dublje razumevanje strukture podataka i odnosa između atributa i ishoda.

3. Arhitektura Podataka i Karakteristike

Podaci se sastoje od 30 numeričkih karakteristika izvedenih iz digitalizovanih slika biopsije, koje opisuju ćelijska jezgra. Svaki uzorak tumora obuhvata tri grupe merenja za svaku od deset osnovnih karakteristika:

Srednje Vrednosti (Mean Values):

- **Radius (radijus):** Srednja udaljenost od centra do tačaka na obodu.
- **Texture (tekstura):** Standardna devijacija vrednosti sive skale.
- **Perimeter (obim):** Ukupna dužina konture jezgra.
- **Area (površina):** Površina unutar konture jezgra.
- **Smoothness (glatkoća):** Lokalna varijacija dužina radijusa.
- **Compactness (kompaktnost):** Mera kompaktnosti oblika, izračunata kao $(\text{obim}^2 / \text{površina}) - 1.0$.
- **Concavity (konkavnost):** Ozbiljnost konkavnih delova konture.
- **Concave Points (konkavne tačke):** Broj konkavnih delova konture.
- **Symmetry (simetrija):** Mera simetrije oblika jezgra.
- **Fractal Dimension (fraktalna dimenzija):** Aproksimacija "obalne linije", koja opisuje složenost konture.

Standardne Greške (Standard Error - SE):

Standardna greška svake od gore navedenih karakteristika, odražavajući varijabilnost merenja.

Najgore (Worst) Vrednosti:

"Najgore" ili najveće (ili srednje za najveća tri) vrednosti svake od gore navedenih karakteristika, često ukazuju na ekstremne promene u ćelijama.

4. Metodologija Klasifikacije: Detaljan Postupak Obuke Modela

Proces obuke modela klasifikacije tumora dojke prati standardni tok rada u mašinskom učenju, sa naglaskom na transparentnost i ponovljivost koraka.

4.1. Faze Projekta:

1. Razumevanje Podataka i Problem:

- Detaljna analiza skupa podataka, uključujući tipove karakteristika i njihovu distribuciju.
- Definisanje problema klasifikacije (binarna klasifikacija: benigni vs. maligni).

2. Prikupljanje i Predobrada Podataka:

- **Učitavanje Podataka:** Učitavanje skupa podataka iz izvora (u našem slučaju CSV fajl).
- **Ispitivanje Podataka (EDA - Exploratory Data Analysis):**
 - Provera nedostajućih vrednosti.
 - Analiza distribucije svake karakteristike (histogrami, box plotovi).
 - Identifikacija korelacija između karakteristika.
 - Vizualizacija odnosa između karakteristika i ciljne varijable (tip tumora).
- **Čišćenje Podataka:** Rešavanje eventualnih anomalija ili grešaka u podacima.
- **Skaliranje Podataka:** Normalizacija ili standardizacija numeričkih karakteristika kako bi se sprečila dominacija karakteristika sa većim opsegom vrednosti (npr. StandardScaler).

3. Podela Skupa Podataka:

- Podela celokupnog skupa podataka na trening i test setove (npr. 70% za trening, 30% za testiranje). Ovo osigurava da se model evaluiira na podacima koje ranije nije video, pružajući realniju procenu njegovih performansi.

4. Izbor i Obuka Modela:

- **Izbor Algoritama:** Odabir jednog ili više algoritama mašinskog učenja pogodnih za binarnu klasifikaciju. Primeri algoritama koji se mogu koristiti uključuju:
 - **Logistička Regresija (Logistic Regression):** Jednostavan, ali efikasan linearni model.
 - **Podržavajuće Vektorske Mašine (Support Vector Machines - SVM):** Moćan algoritam za klasifikaciju, posebno efikasan u prostorima visoke dimenzionalnosti.
 - **Drvo Odluke (Decision Tree):** Intuitivan model koji donosi odluke na osnovu serije pravila.
 - **Nasumična Šuma (Random Forest):** Ansambl metod koji kombinuje više stabala odluke za poboljšanje tačnosti i smanjenje prekomernog prilagođavanja.

- **Obuka Modela:** Svaki odabrani model se obučava nezavisno na trening setu podataka. Proces obuke uključuje prilagođavanje internih parametara modela kako bi se minimizovala greška predviđanja.

5. **Evaluacija Modela:**

- Procena performansi obučenih modela na test setu. Metrike evaluacije su ključne za razumevanje koliko dobro model generalizuje na neviđene podatke.

5. Očekivani Izlazi i Rezultati

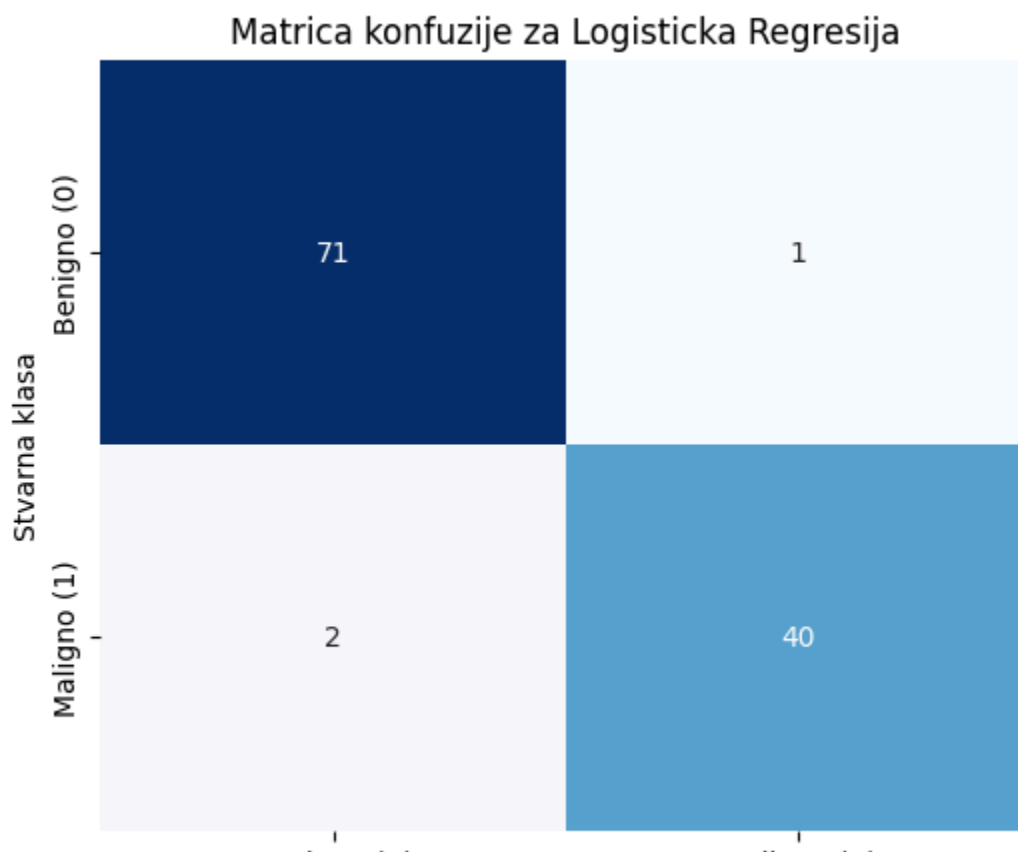
Kao rezultat ovog projekta, očekuju se sledeći artefakti i analize:

5.1. Metrike Evaluacije Modela (za svaki model)

- **Accuracy (tačnost):** Ukupni procenat tačno klasifikovanih uzoraka.
- **Precision (preciznost):** Odnos tačno pozitivnih predviđanja prema ukupnom broju predviđenih pozitivnih (relevantnost pozitivnih predviđanja).
- **Recall (odziv):** Odnos tačno pozitivnih predviđanja prema ukupnom broju stvarnih pozitivnih slučajeva (sposobnost modela da pronađe sve pozitivne slučajeve).
- **F1-Score (F1-mera):** Harmonijska sredina preciznosti i odziva, korisna kada je distribucija klasa neuravnotežena.
- **Matrica Konfuzije:** Tabela koja sumira performanse klasifikacionog modela, prikazujući broj tačnih i netačnih predviđanja za svaku klasu.

5.2. Vizuelna Analiza Važnosti Karakteristika (za svaki model, ako je primenjivo)

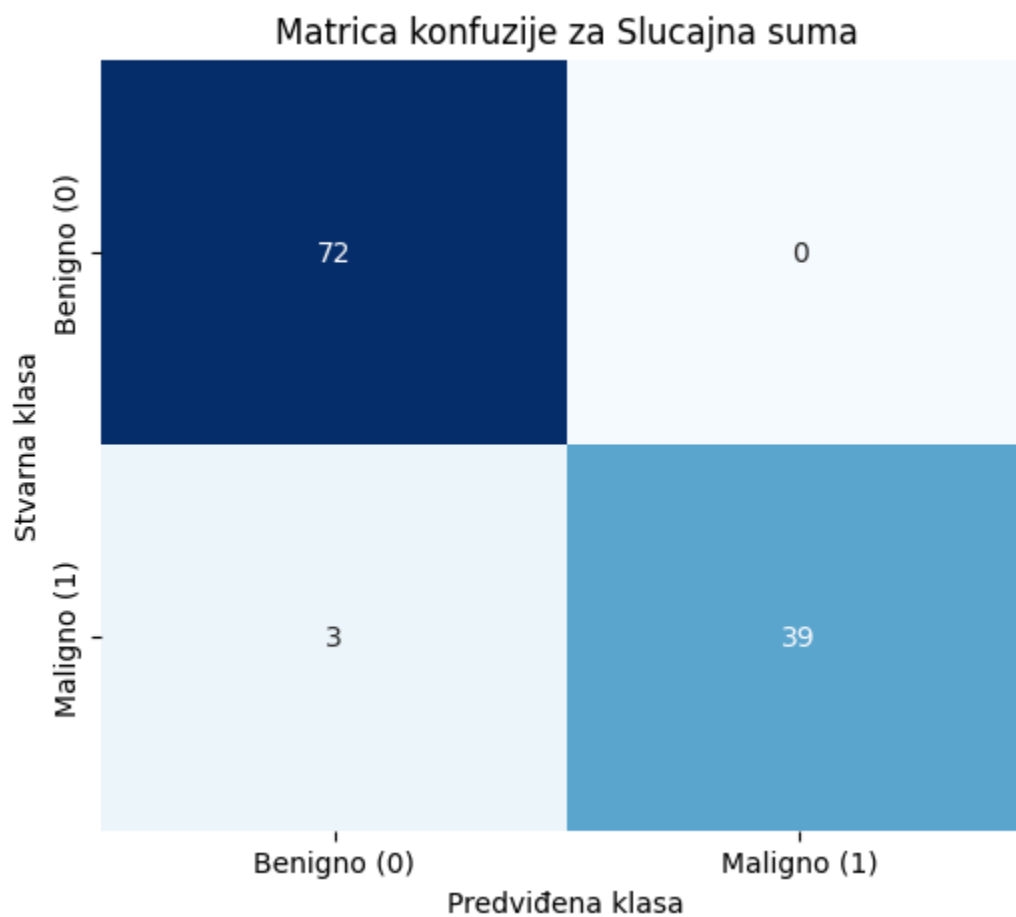
Logistička Regresija



Rezultati:

	preciznost	osetljivost	f1-score	podrška
0	0.97	0.99	0.98	72
1	0.98	0.95	0.96	42
tačnost			0.97	114

Slučajna suma



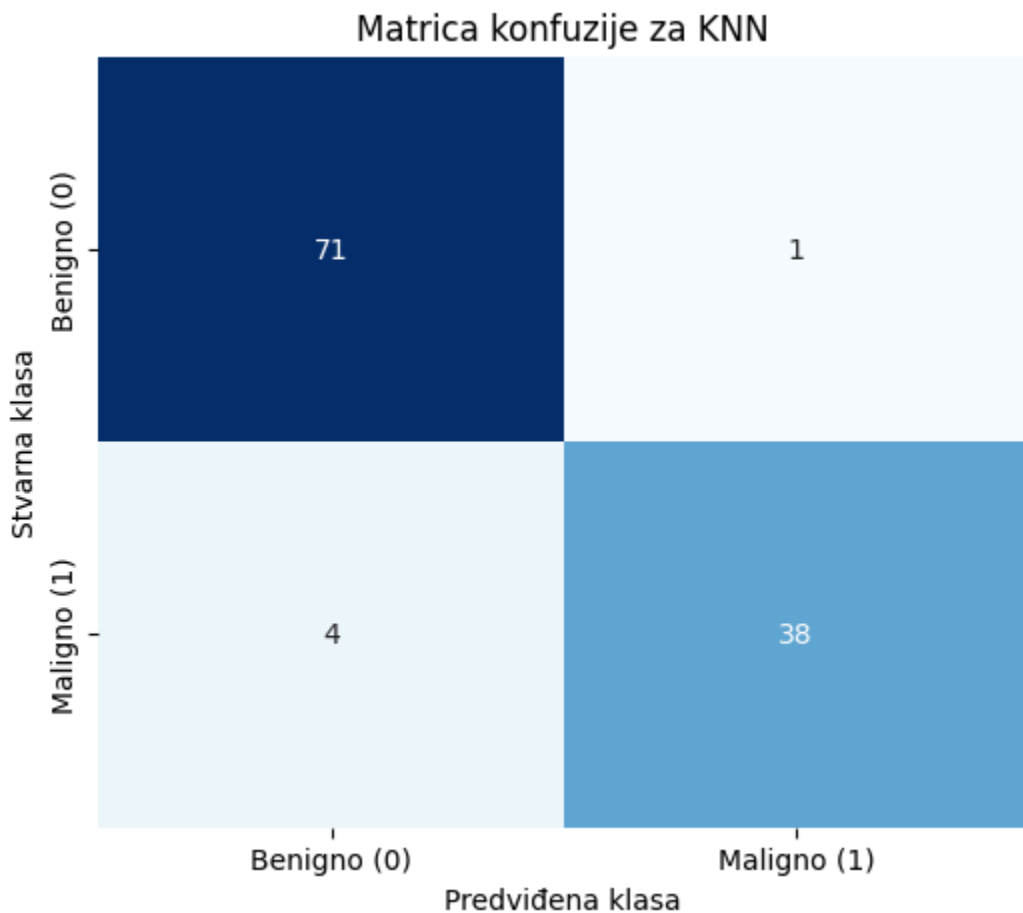
Rezultati:

	preciznost	osetljivost	f1-score	podrška
0	0.96	1	0.98	72

1	1	0.93	0.96	42
tačnost			0.97	114

Model

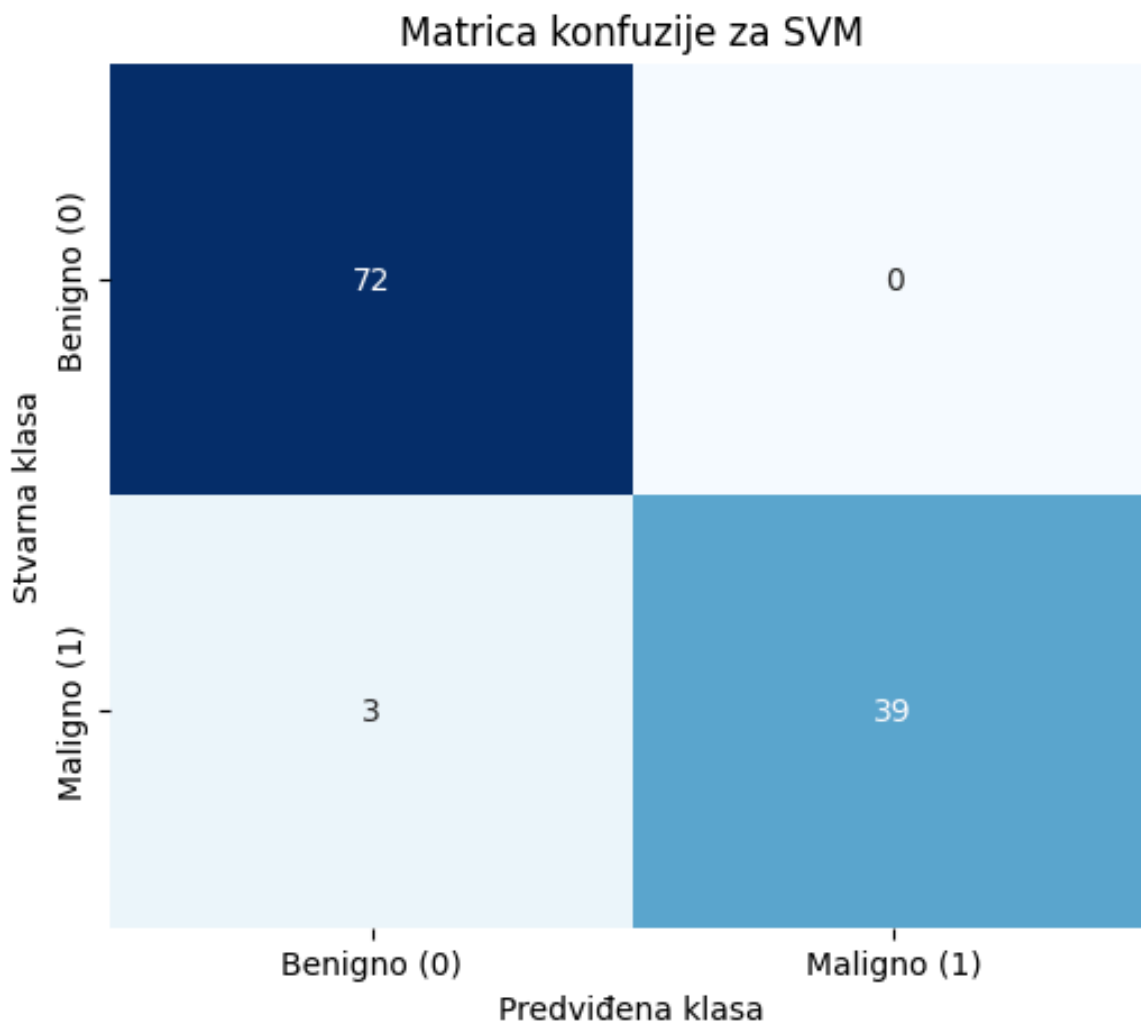
K-Najbližih komšija



Rezultati:

	preciznost	osetljivost	f1-score	podrška
0	0.95	0.99	0.97	72
1	0.97	0.90	0.94	42
tačnost			0.96	114

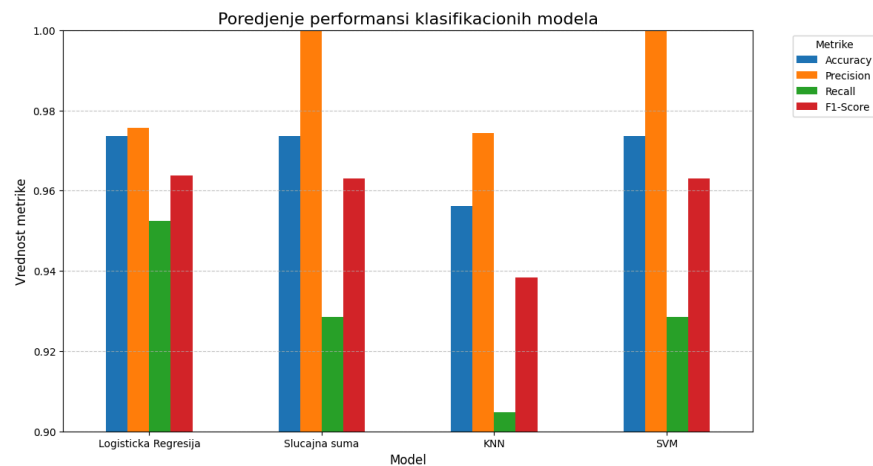
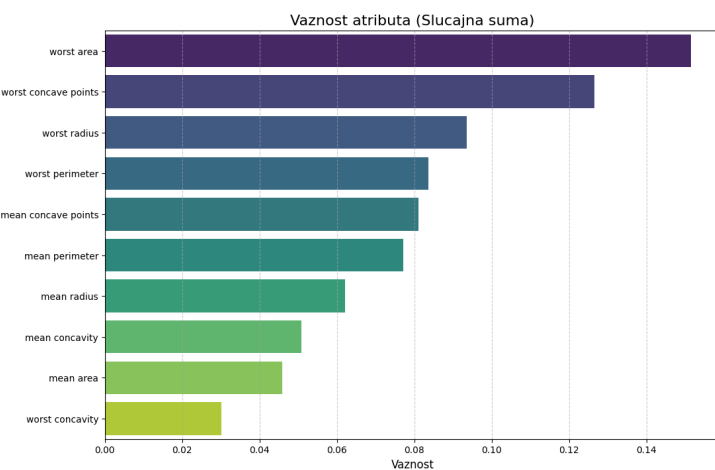
Model potpornih vektora(SVM):



Rezultati:

	preciznost	osetljivost	f1-score	podrška
0	0.96	1	0.98	72
1	1	0.93	0.96	42
tačnost			0.97	114

5.3. Interpretacija Rezultata i Preporuke (za svaki model):



Na osnovu rezultata evaluacije modela, uključujući metrike preciznosti, osetljivosti, F1-mere i tačnosti za Logističku Regresiju, Slučajnu Šumu, Model K-Najbližih Komšija i Model Potpornih Vektora (SVM), u nastavku će biti predstavljen uvod u analizu ovih performansi.

Logistička regresija

Veoma snažan model koji se ističe kao **najbolji "detektiv"**.

- **(+) Glavna snaga:** Najbolja osetljivost(recall). Ovaj model je **najuspešniji u pronalaženju malignih slučajeva**, sa samo 2 propuštena uzorka – najmanje od svih modela.
- **(-) Glavna slabost:** Napravio je jednu grešku "lažne uzbune", gde je zdrav uzorak pogrešno klasifikovao kao malignan.
- **Presuda: Odlična alternativa** ako je apsolutni prioritet pronaći što veći broj malignih slučajeva, čak i po cenu povremenih lažnih alarma.

Slučajna šuma

Izuzetno pouzdan i robustan model, čija je glavna odlika **opreznost**.

(+) Glavna snaga: Savršena preciznost (100%) kod maligne klase. Kada model kaže da je nešto maligno, on je apsolutno u pravu i ne pravi greške "lažne uzbune".

(-) Glavna slabost: Zbog svoje opreznosti i fokusa na izbegavanje lažnih alarma, "propustio" je 3 stvarna maligna slučaja.

Presuda: Najbolji izbor ako je prioritet maksimalna pouzdanost i izbegavanje nepotrebnog stresa i procedura kod zdravih pacijenata.

K-Najbližih komšija

Objektivno **najslabiji model** u ovoj grupi.

- **(+) Glavna snaga:** I dalje poseduje visoku ukupnu tačnost (96%), što pokazuje da nije loš model u apsolutnom smislu.
- **(-) Glavna slabost:** Pravi najviše najopasnijih grešaka. Sa **4 propuštena maligna slučaja**, najmanje je pouzdan za kritičnu dijagnostiku.
- **Presuda:** **Ne preporučuje se** za ovaj zadatak u poređenju sa ostala tri modela, jer su superiorniji u svakom važnom pogledu.

Model potportnih vektora(SVM)

Moćan i efikasan model koji, kao i Slučajna šuma, pruža savršeno pouzdanu "pozitivnu" dijagnozu.

- **(+) Glavna snaga:** Apsolutna preciznost (100%) kod predviđanja maligne klase. Nije napravio nijednu grešku "lažne uzbune", što znači da je njegova pozitivna dijagnoza neupitna.
- **(-) Glavna slabost:** Identično kao Slučajna šuma, njegova mana je što je propustio 3 maligna slučaja, svrstavajući ih u benignu kategoriju.
- **Presuda:** Takođe najbolji izbor kada je cilj eliminisati lažno pozitivne rezultate i imati najviši stepen poverenja u dijagnozu maligniteta.

Na osnovu detaljne analize i definisanog prioriteta da je **propuštanje maligne dijagnoze najkritičnija greška**, zaključuje se sledeće:

Najbolji model za ovaj zadatak je Logistička Regresija.

Iako Slučajna šuma i SVM daju pouzdanije pozitivne dijagnoze (bez lažnih uzbuna), Logistička Regresija je superiornija jer je **najuspešnija u primarnom cilju – pronalaženju malignih slučajeva**. Sa samo dva propuštena slučaja, ona nudi najbolji kompromis i najmanji rizik od najopasnije moguće greške

6. Zaključak

Ovaj projekat pruža solidnu osnovu za klasifikaciju tumora dojke, sa fokusom na jasan i ponovljiv metodološki pristup. Dalji razvoj može uključivati istraživanje naprednijih algoritama, primenu tehnika za smanjenje dimenzionalnosti, korišćenje većih i raznovrsnijih skupova podataka, te potencijalnu integraciju sa kliničkim sistemima za podršku dijagnostici. Kontinuirana evaluacija i iterativno poboljšanje su ključni za optimizaciju modela i maksimizovanje njegove korisnosti.