

Predikcija vakcinacije ispitanika protiv sezonskog gripa i virusa H1N1

Mašinsko učenje (13M051MU)

student: Mihailo Grbić

profesor: doc. dr Predrag Tadić

Zadatak

- Mišljenje ispitanika o vakcinaciji protiv:
 - Sezonskog gripa (binarno)
 - H1N1 virusa (binarno)
- Anketa sa 37 pitanja —→ podaci
- Postavka
 - Čišćenje i priprema podataka
 - Klasifikacija
 - Predikcija dve labele —→ dva modela
 - Izostavljena labela se ne koristi kao prediktor
 - Evaluacija

Čišćenje i priprema podataka

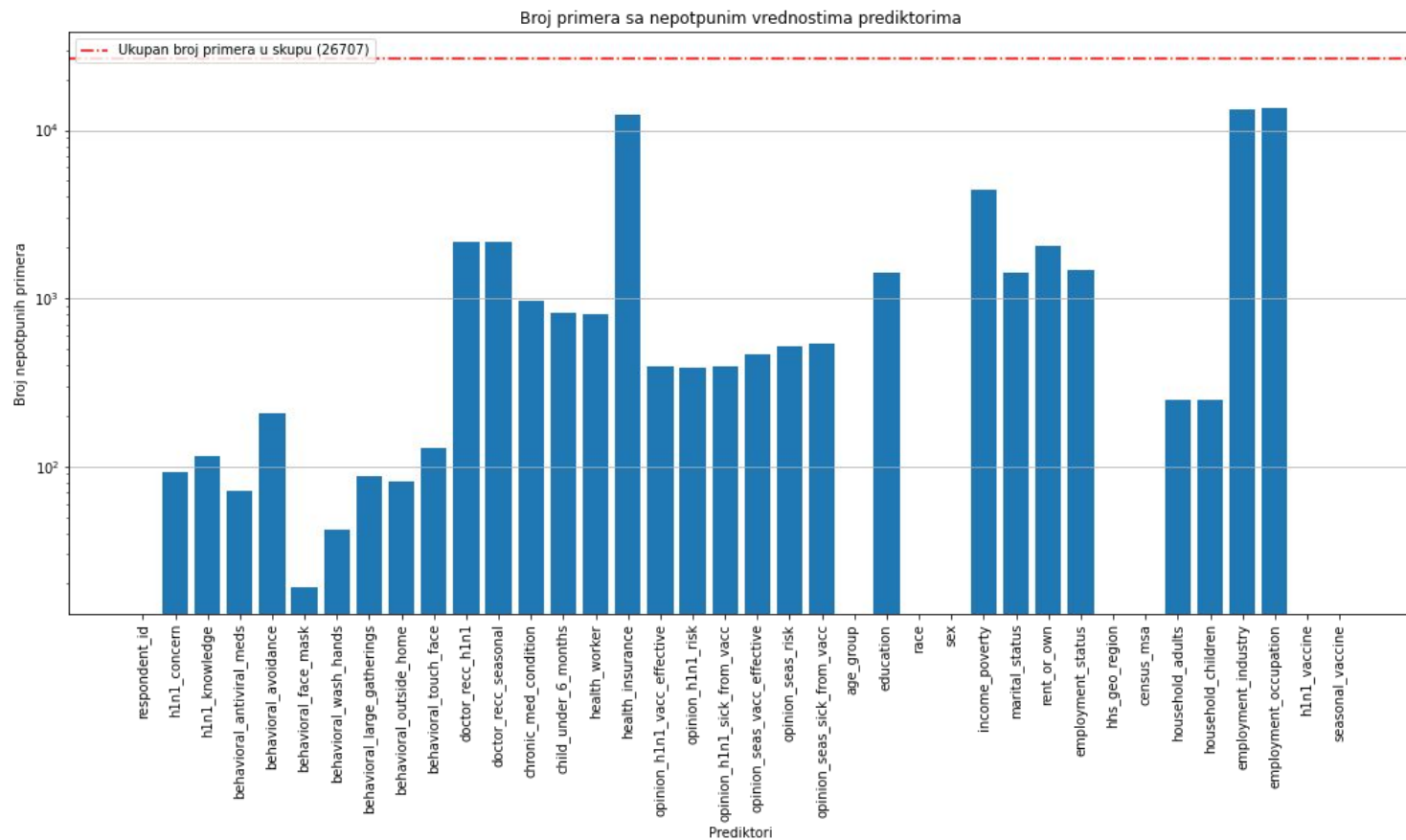
Skup podataka

- Osnovne informacije:
 - *training_set_features.csv* : 35 prediktora
 - numeričkih: 23
 - kategoričkih: 12
 - bez indeks ispitanika
 - *training_set_labels.csv* : 2 izlaza
 - *seasonal_vaccine*
 - *h1n1_vaccine*
 - Broj primera: 26707
 - Problemi:
 - Nepotpuni primeri
 - Kategoričke vrednosti

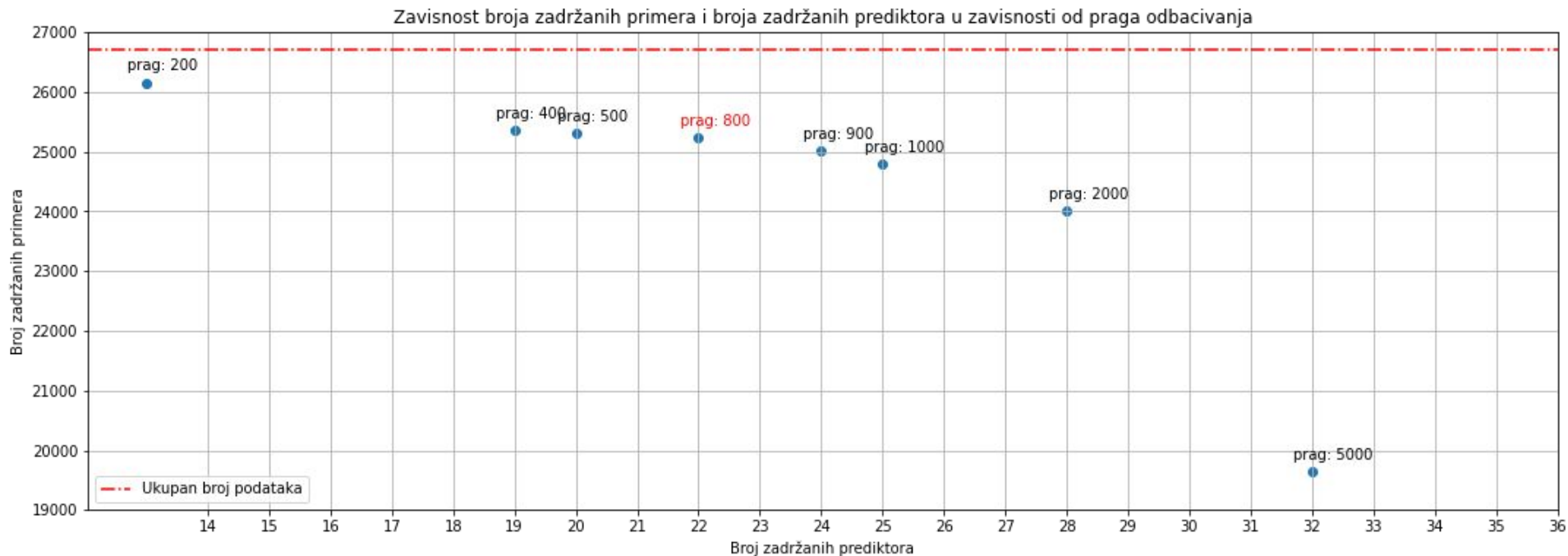
#	Column	Non-Null Count	Dtype
----	-----	-----	-----
0	respondent_id	26707 non-null	int64
1	h1n1_concern	26615 non-null	float64
2	h1n1_knowledge	26591 non-null	float64
3	behavioral_antiviral_meds	26636 non-null	float64
4	behavioral_avoidance	26499 non-null	float64
5	behavioral_face_mask	26688 non-null	float64
6	behavioral_wash_hands	26665 non-null	float64
7	behavioral_large_gatherings	26620 non-null	float64
8	behavioral_outside_home	26625 non-null	float64
9	behavioral_touch_face	26579 non-null	float64
10	doctor_recc_h1n1	24547 non-null	float64
11	doctor_recc_seasonal	24547 non-null	float64
12	chronic_med_condition	25736 non-null	float64
13	child_under_6_months	25887 non-null	float64
14	health_worker	25903 non-null	float64
15	health_insurance	14433 non-null	float64
16	opinion_h1n1_vacc_effective	26316 non-null	float64
17	opinion_h1n1_risk	26319 non-null	float64
18	opinion_h1n1_sick_from_vacc	26312 non-null	float64
19	opinion_seas_vacc_effective	26245 non-null	float64
20	opinion_seas_risk	26193 non-null	float64
21	opinion_seas_sick_from_vacc	26170 non-null	float64
22	age_group	26707 non-null	object
23	education	25300 non-null	object
24	race	26707 non-null	object
25	sex	26707 non-null	object
26	income_poverty	22284 non-null	object
27	marital_status	25299 non-null	object
28	rent_or_own	24665 non-null	object
29	employment_status	25244 non-null	object
30	hhs_geo_region	26707 non-null	object
31	census_msa	26707 non-null	object
32	household_adults	26458 non-null	float64
33	household_children	26458 non-null	float64
34	employment_industry	13377 non-null	object
35	employment_occupation	13237 non-null	object
36	h1n1_vaccine	26707 non-null	int64
37	seasonal_vaccine	26707 non-null	int64

dtypes: float64(23), int64(3), object(12)

Čišćenje podataka



Čišćenje podataka



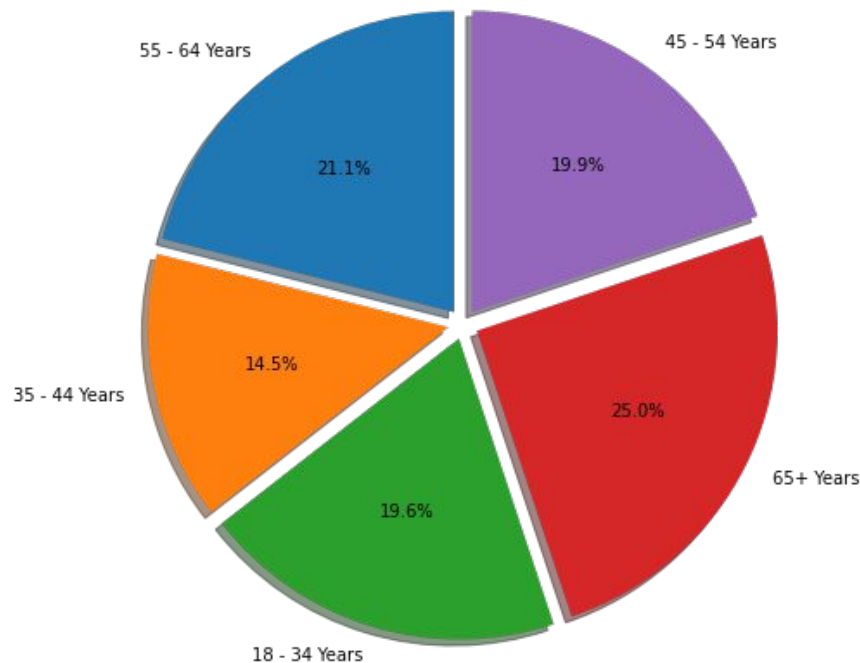
Čišćenje podataka

- Analiza
 - 5/35 “čistih” prediktora
 - 3/35 prediktora kojima nedostaje $\sim \frac{1}{3}$ podataka
- Filtriranje
 - Odbacivanje prediktora koji sadrže više od 800 nepotpunih primera
 - Preostali nepotpuni primeri se odbacuju
- Rezultat
 - Broj zadržanih prediktora: 22/35
 - Broj primera zadržanih primera: 25238
 - ~ 1500 odbačenih primera (5.61% skupa) \longrightarrow Prihvatljivo!

Kategorički prediktori | Starosna grupa

- Punoletni ispitanici
- Pet starosnih grupa
 - Slična zastupljenost
- Gradacioni poredak
 - Kodiranje celobrojnim vrednostima

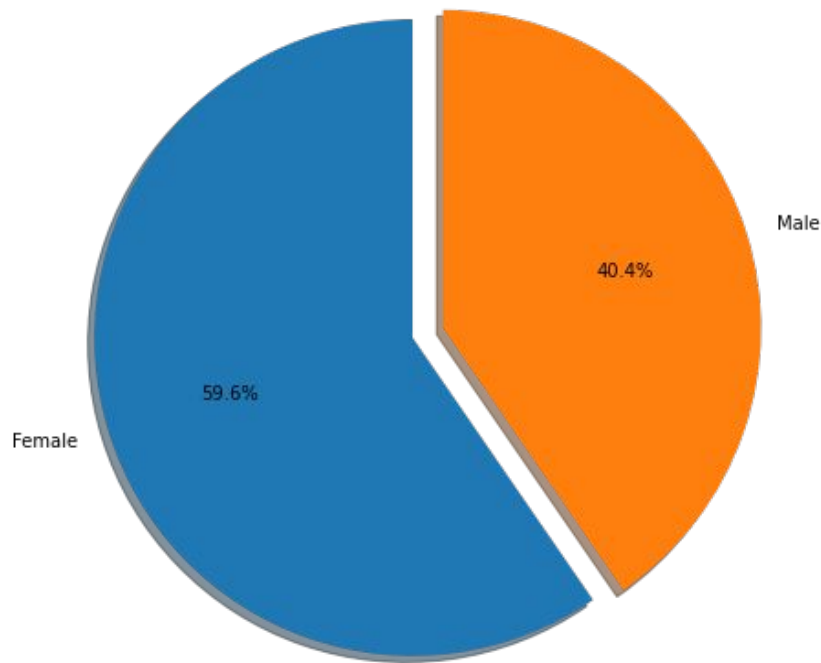
Zastupljenost starosnog doba među ispitanicima



Kategorički prediktori | Pol

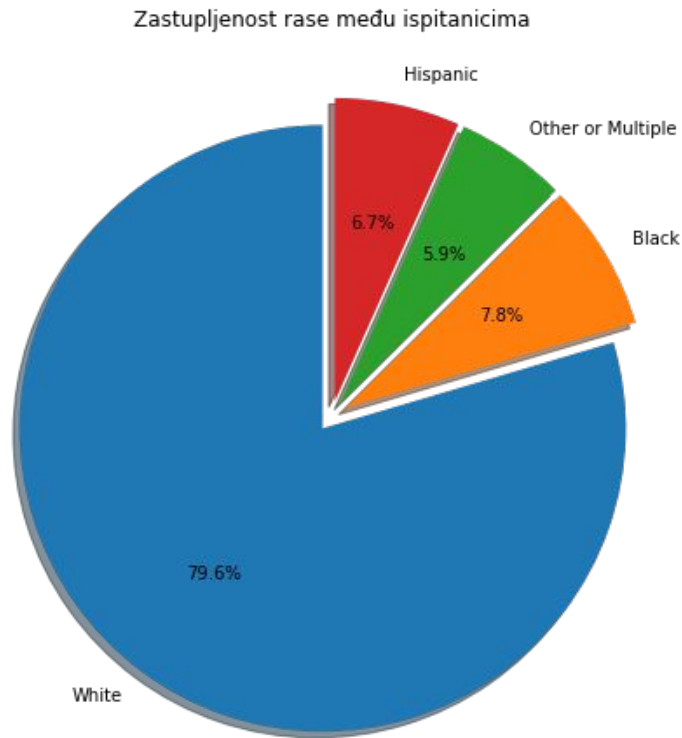
- Dva pola (muški/ženski)
 - Nešto veća zastupljenost žena
- Binarno kodiranje
 - 0: Ženski pol
 - 1: Muški pol

Zastupljenost pola među ispitanicima



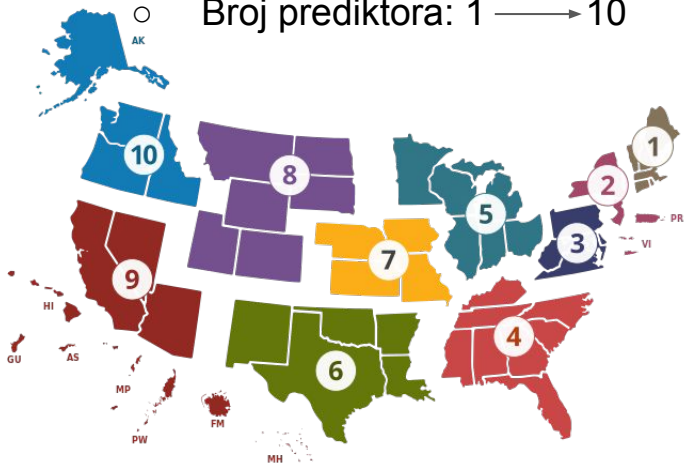
Kategorički prediktori | Rasa

- Podela po boji kože
 - Četiri kategorije
 - Veliku većinu čine ljudi sa belom bojom kože
- Nema smislenog gradacionog poretka
 - *One-Hot Encoding*
 - Broj prediktora: 1 \longrightarrow 4



Kategorički prediktori | Geografski region

- Istraživanje je rađeno na području Sjedinjenih Američkih Država
 - Podela na 10 grupa država
 - Ujednačena zastupljenost
- Nema smislenog gradacionog poretka
 - *One-Hot Encoding*
 - Broj prediktora: 1 → 10



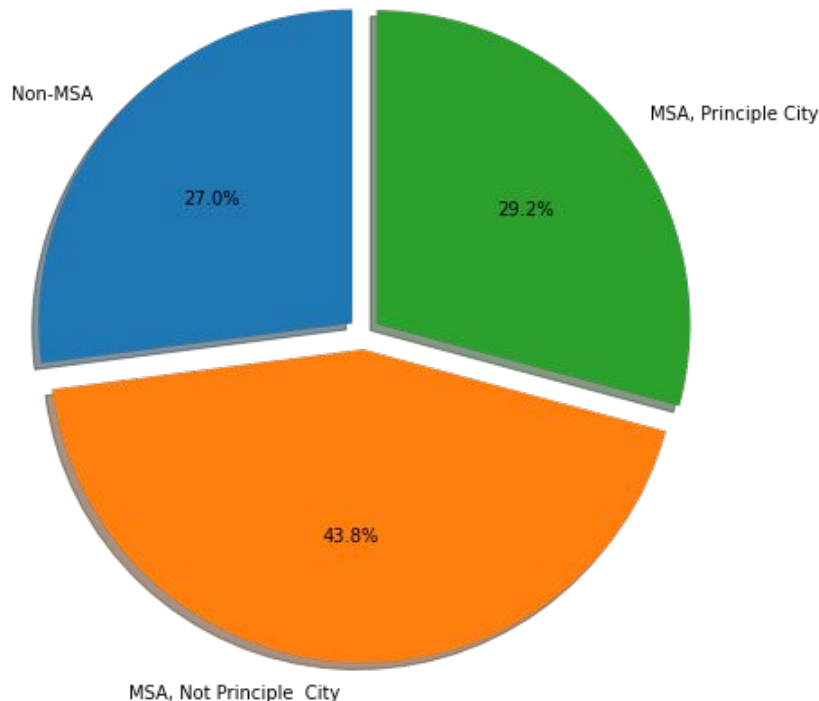
<https://www.hhs.gov/about/agencies/iea/regional-offices/index.html>



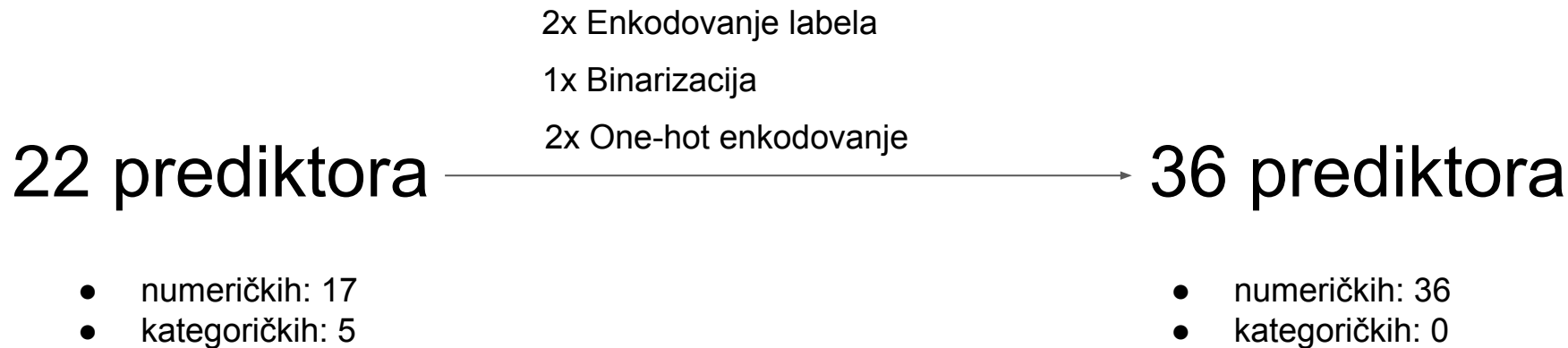
Kategorički prediktori | Urbanost sredine

- Broj stanovnika mesta prebivališta
 - Non-MSA: Populacija < 50 000
 - MSA: Populacija > 50 000
- Najveći udeo ispitanika iz urbanih sredina van glavnog grada
- $\sim\frac{3}{4}$ stanovnika su iz urbanih naselja
- Ima smisla gradaciono porediti kategorije
 - *Non-MSA*: 0
 - *MSA, Not Principle City*: 1
 - *MSA, Principle City*: 2

Zastupljenost urbanih sredina među ispitanicima



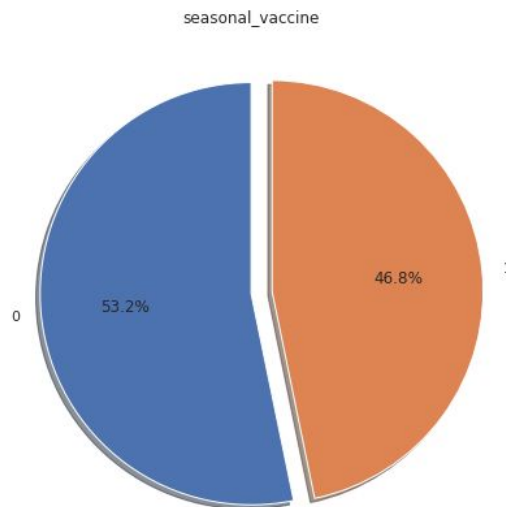
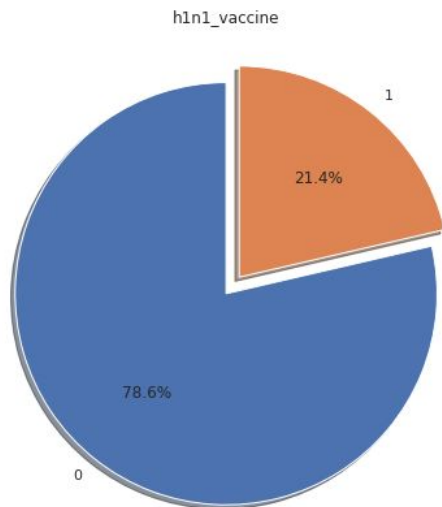
Kategorički prediktori | Rezime



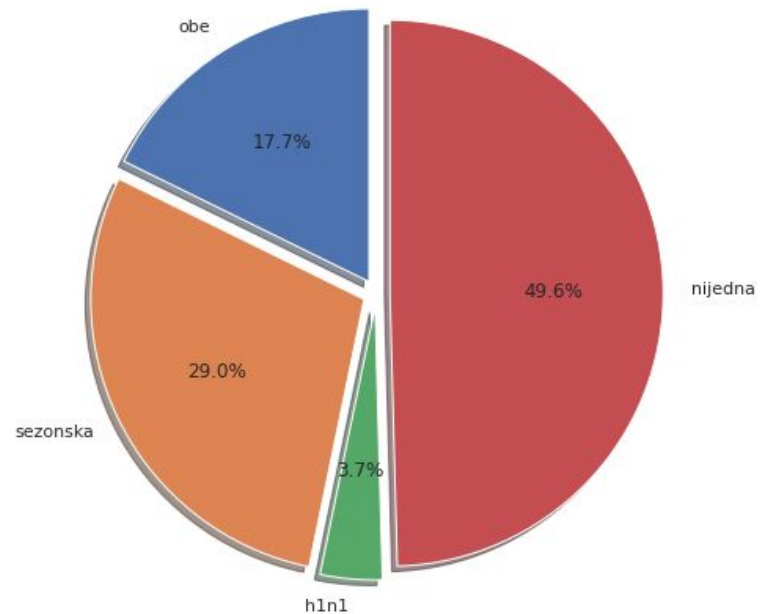
Izlazni podaci

- (Ne)balansiranost klasa
 - 0: “protiv” vakcinacije
 - 1: “za” vakcinaciju

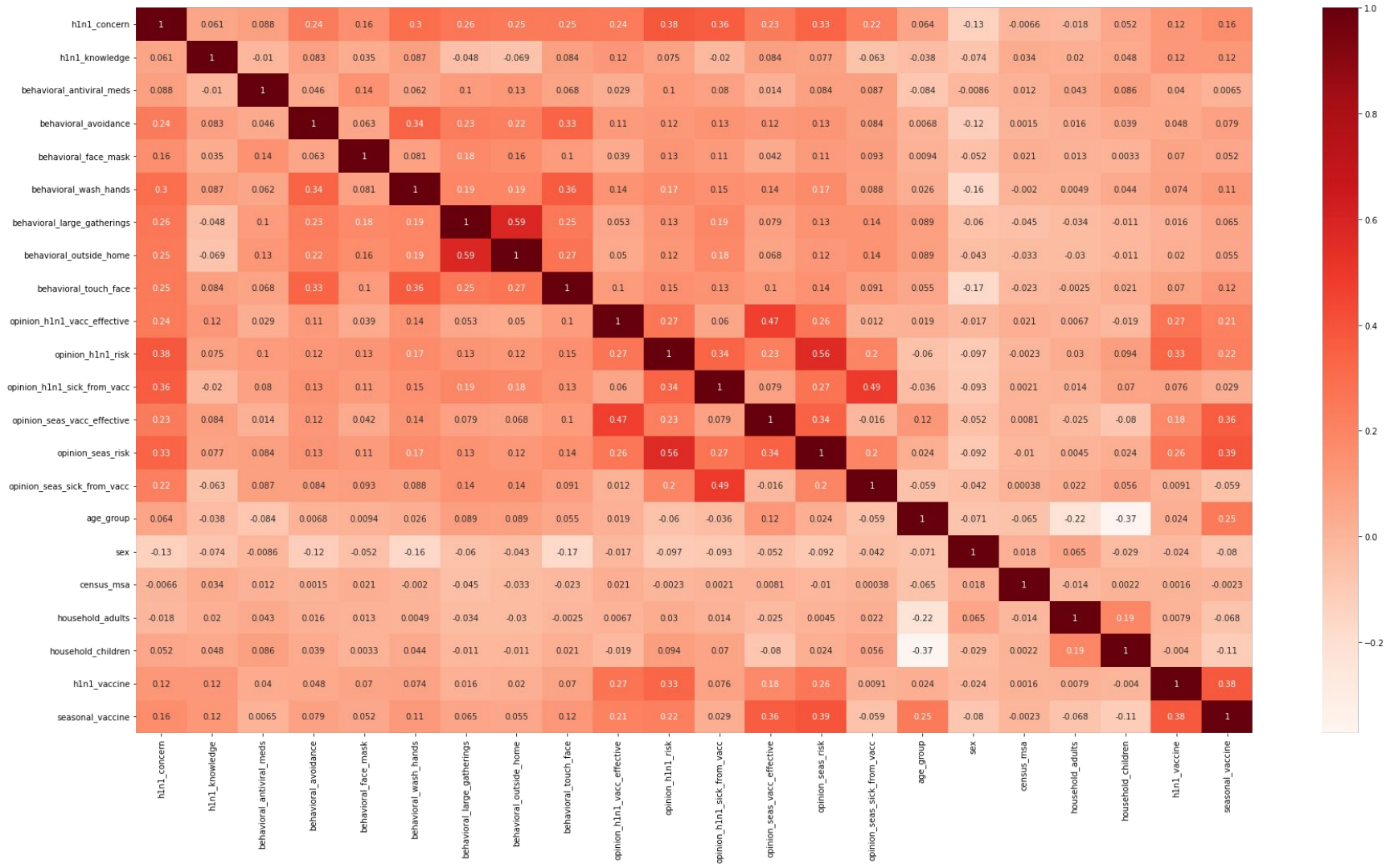
Zastupljenost labela u podacima



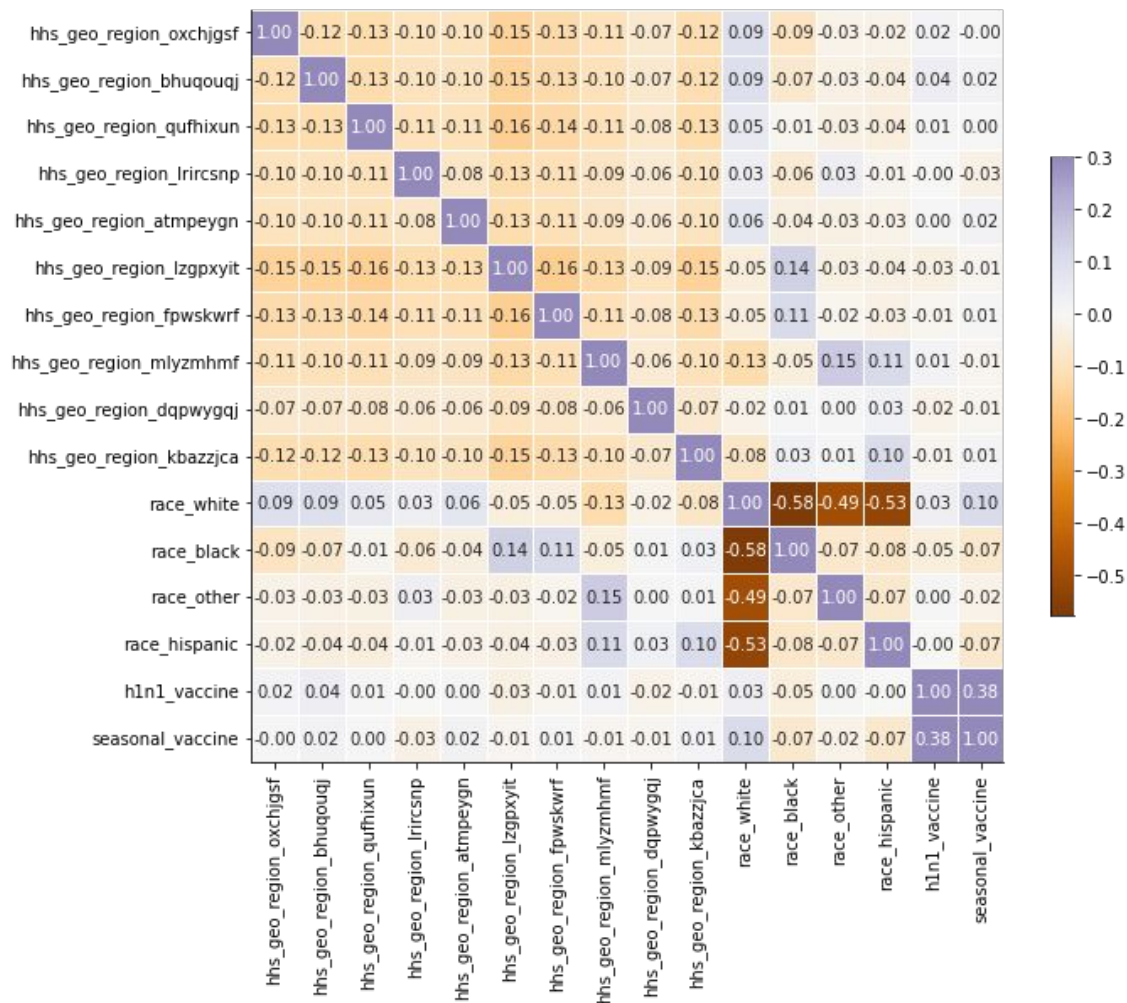
Tip vakcine za koju bi se ispitanici odlučili



Korelacione mape



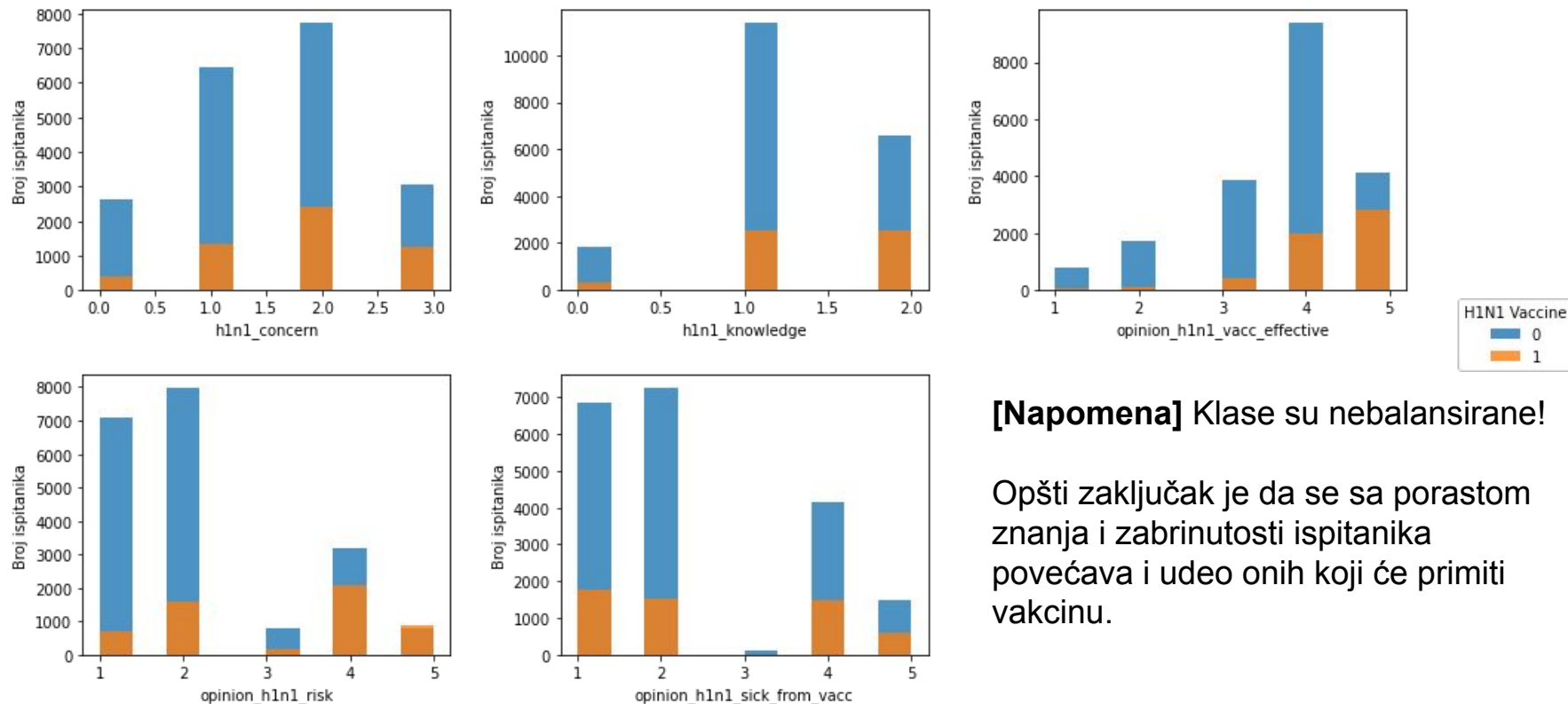
Prediktori transformisani one-hot enkodovanjem



Analiza veze pojedinačnih prediktora i izlaza

Uticaj znanja i mišljenja ispitanika | H1N1 virus

Znanje i mišljenje ispitanika o vakcini protiv H1N1 virusa



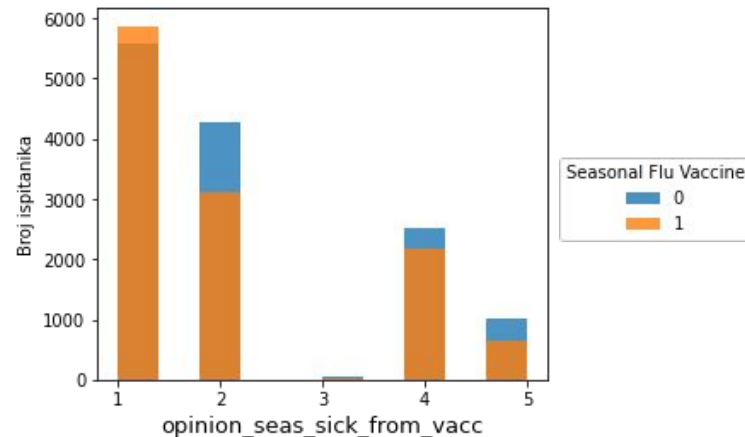
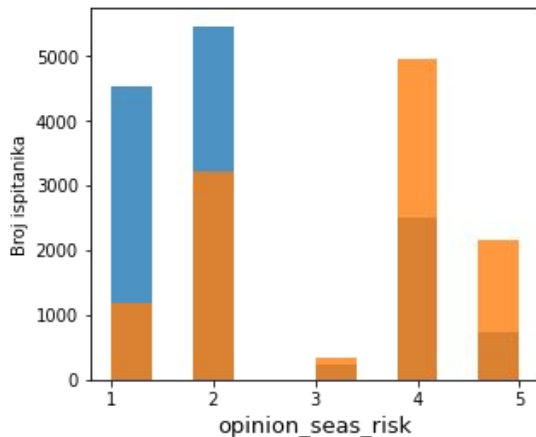
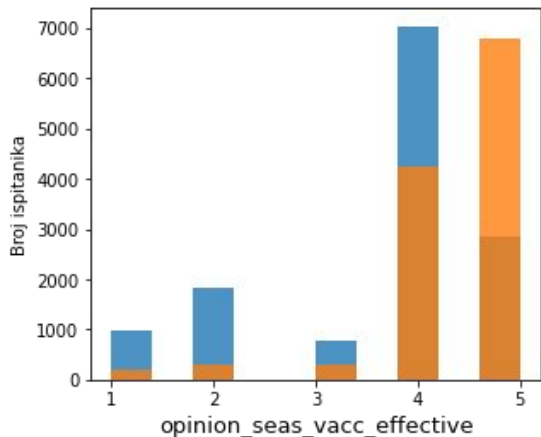
[Napomena] Klase su nebalansirane!

Opšti zaključak je da se sa porastom znanja i zabrinutosti ispitanika povećava i udeo onih koji će primiti vakcinu.

Uticaj znanja i mišljenja ispitanika | Sezonski grip

- Rizična grupa i/ili mišljenje da je vakcina efikasna → veći udeo vakcinisanih
- Učestalo razboljevanje ispitanika nije pokazalo uticaj na opredeljenje za vakcinu

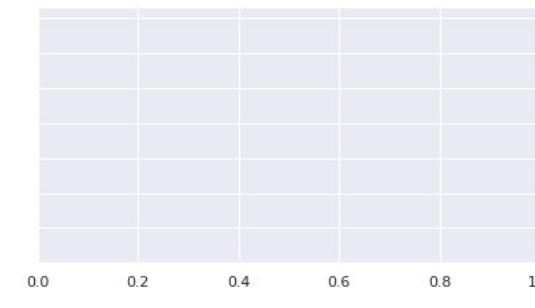
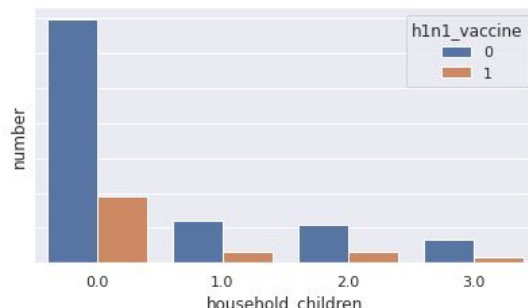
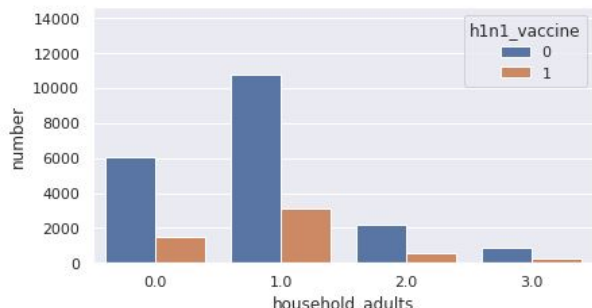
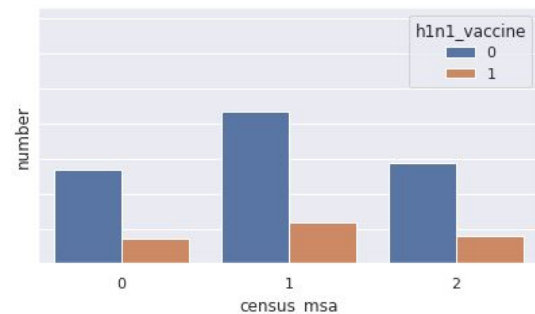
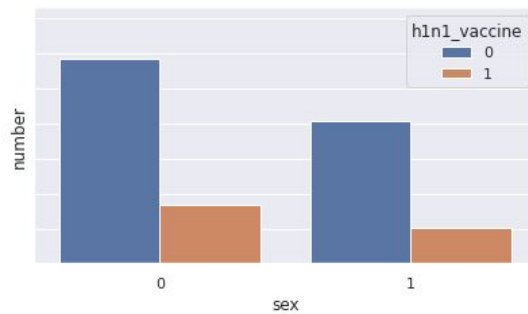
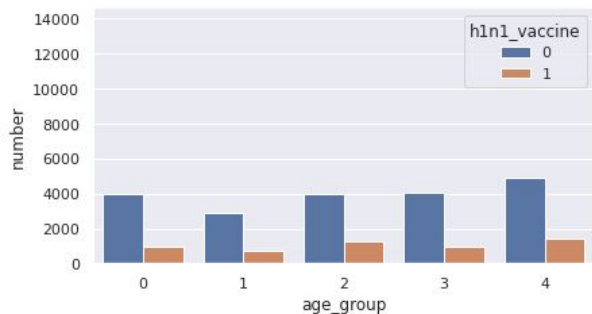
Znanje i mišljenje ispitanika o vakcini protiv sezonskog gripa



Uticaj ostalih prediktora | H1N1

Teško je uočiti trend u podacima koji bi sugerisao na opredeljenje za vakcinu

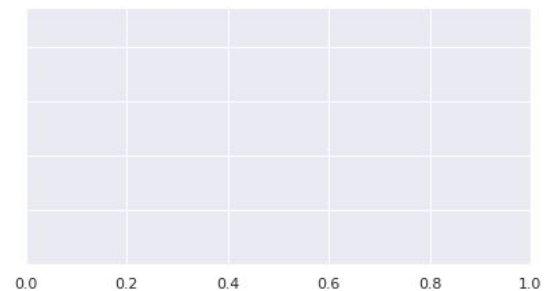
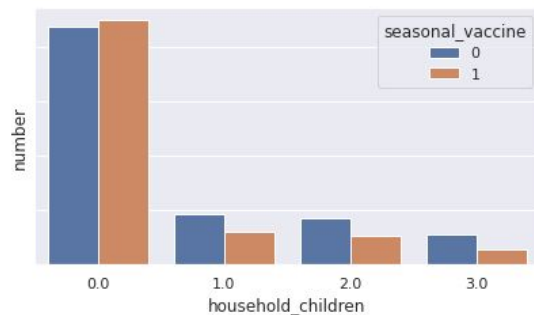
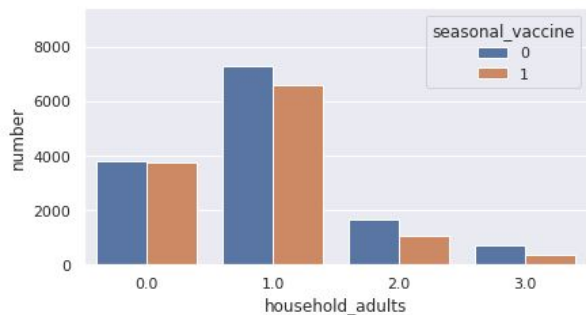
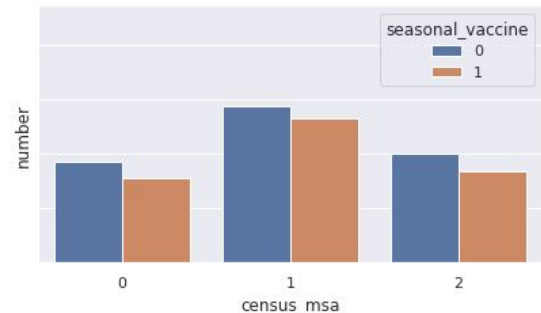
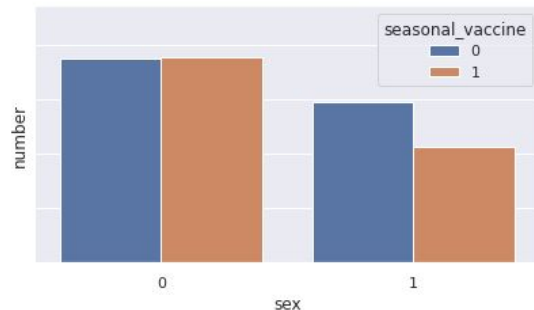
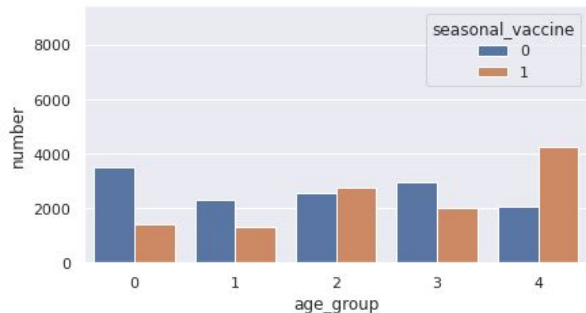
Ostale informacije o ispitanicima u zavisnosti od vakcinacije protiv H1N1 virusa



Uticaj ostalih prediktora | Sezonski grip

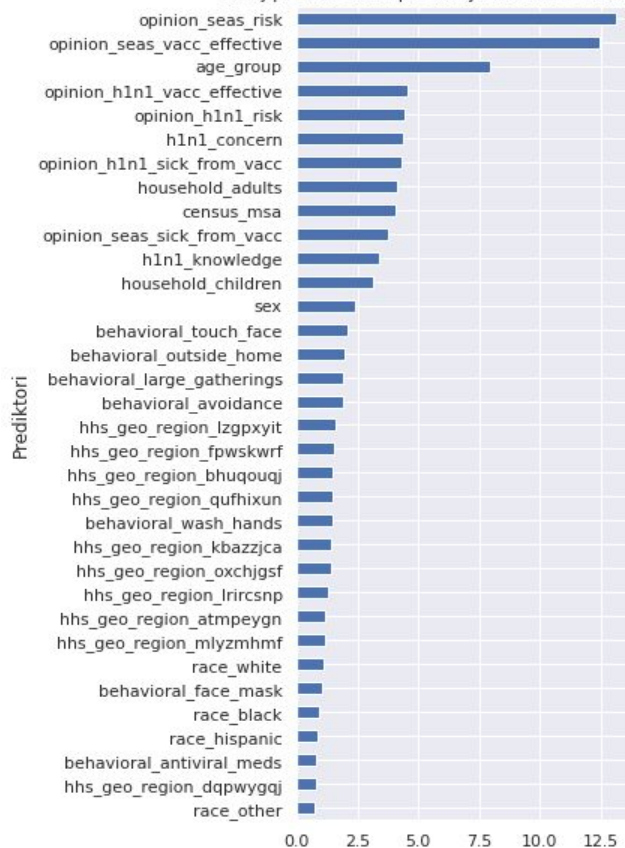
- Ispitanici sa preko 65 godina se najčešće opredeljuju za vakcinu
- Žene se nešto češće opredeljuju za vakcinu od muškaraca, a muškarci češće odbijaju vakcinu
- Ostali prediktori nisu naročito informativni

Ostale informacije o ispitanicima u zavisnosti od vakcinacije protiv sezonskog gripa



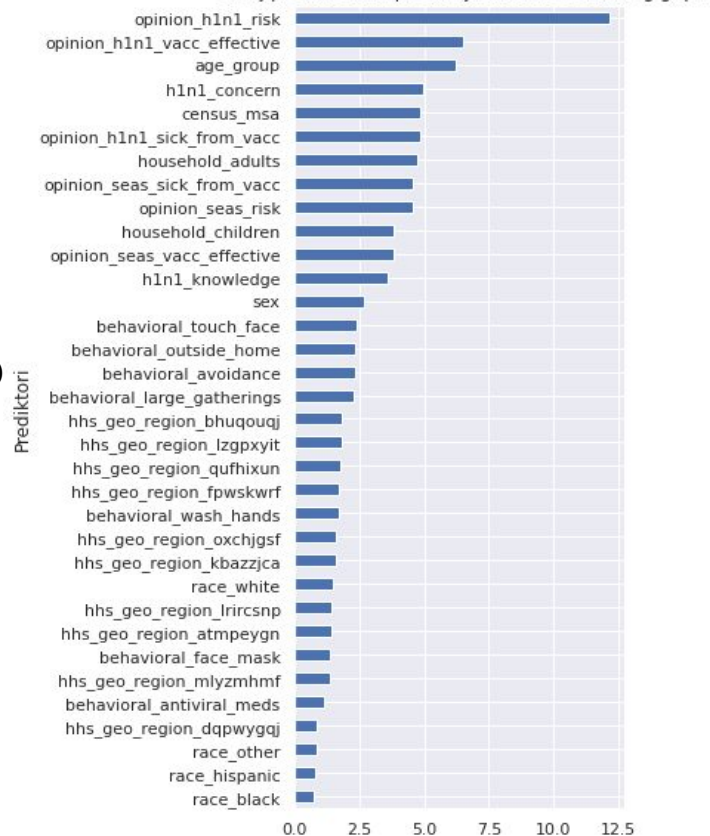
Značaj prediktora

Uticaj prediktora za predikciju vakcine sezonskog gripa



- Random Forest
 - 200 estimatora
- Isti prediktori su u top 3 najinformativnijih

Uticaj prediktora za predikciju vakcine sezonskog gripa



Izbor modela

Metodologija

- Podela inicijalnog skupa
 - *obučavajući skup : testirajući skup = 90% : 10% = 22714 : 2524*
- Standardizacija podataka
 - Učenje statističkih parametara na obučavajućem skupu
 - Transformacija prediktora obučavajućeg i testirajućeg skupa
- Unakrsna validacija
 - Obučavajući skup \longrightarrow 4 struka
 - Odnos zastupljenosti klasa kao u inicijalnim podacima
- Pretraga optimalnih hiper-parametara njihovim iscrpnim kombinovanjem
- Izbor na osnovu srednje balansirane tačnosti na validacionom skupu

[Napomena] Metodologija je nezavisno primenjena na obe labele!

Modeli i mogući hiper-parametri

Vakcina protiv sezononskog gripa

- **Logistička regresija**
 - 'penalty': ['l1', 'l2'],
 - 'C': [0.5, 1, 5, 10, 20],
 - 'solver': ['liblinear'],
 - 'class_weight': ['balanced', None]
- **Metoda nosećih vektora**
 - Linearni / Polinomijalni / Gausovski kernel
 - 'C': [1, 5, 10]
 - Polinomijalni kernel
 - 'degree': [2, 3, 5]
- **Random forest**
 - 'n_estimators': [50, 100, 200]
 - 'criterion': ['gini', 'entropy']
 - 'max_depth': [3, 5, 10]
 - 'max_features': ['sqrt', 'log2']

Vakcina protiv H1N1 virusa

- **Logistička regresija**
 - 'penalty': ['l1', 'l2'],
 - 'C': [0.5, 1, 5, 10, 20],
 - 'solver': ['liblinear'],
 - 'class_weight': ['balanced']
- **Metoda nosećih vektora**
 - Linearni / Polinomijalni / Gausovski kernel
 - 'C': [1, 5, 10]
 - 'class_weight': ['balanced']
 - Polinomijalni kernel
 - 'degree': [2, 3, 5]
- **Random forest**
 - 'n_estimators': [50, 100, 200]
 - 'criterion': ['gini', 'entropy']
 - 'max_depth': [3, 5, 10]
 - 'max_features': ['sqrt', 'log2']
 - 'class_weight': ['balanced']

Modeli i optimalni hiper-parametri

Vakcina protiv sezonskog gripa

- **Logistička regresija**
 - 'penalty': ['**l1**', 'l2'],
 - 'C': [0.5, 1, 5, **10**, 20],
 - 'solver': ['**liblinear**'],
 - 'class_weight': ['balanced', **None**]
- **Metoda nosećih vektora**
 - Linearni / Polinomijalni / **Gausovski kernel**
 - 'C': [**1**, 5, 10]
 - Polinomijalni kernel
 - 'degree': [2, 3, 5]
- **Random forest**
 - 'n_estimators': [50, 100, **200**]
 - 'criterion': ['gini', '**entropy**']
 - 'max_depth': [3, 5, **10**]
 - 'max_features': ['sqrt', '**log2**']

Vakcina protiv H1N1 virusa

- **Logistička regresija**
 - 'penalty': ['**l1**', 'l2'],
 - 'C': [0.5, 1, **5**, 10, 20],
 - 'solver': ['**liblinear**'],
 - 'class_weight': ['**balanced**']
- **Metoda nosećih vektora**
 - **Linearni** / Polinomijalni / Gausovski kernel
 - 'C': [**1**, 5, 10]
 - 'class_weight': ['**balanced**']
 - Polinomijalni kernel
 - 'degree': [2, 3, 5]
- **Random forest**
 - 'n_estimators': [**50**, 100, 200]
 - 'criterion': ['**gini**', 'entropy']
 - 'max_depth': [3, **5**, 10]
 - 'max_features': ['sqrt', '**log2**']
 - 'class_weight': ['**balanced**']

Evaluacija modela sa optimalnim hiper-parametrima

Tumačenje:

precision: Među svim ljudima za koje konačni model bude predvideo da [će se / se neće] vakcinisati, [$\text{precision } 1$ / $\text{precision } 0$]% će se vakcinisati

recall: Među svim ljudima koji [će se / se neće] vakcinisati, model će predvideti [$\text{recall } 1$ / $\text{recall } 0$]%.

Vakcina protiv sezonskog gripa | *train*

Logistička regresija				
	precision	recall	f1-score	support
0	0.76	0.78	0.77	12090
1	0.74	0.72	0.73	10624
accuracy			0.75	22714
macro avg	0.75	0.75	0.75	22714
weighted avg	0.75	0.75	0.75	22714
Metoda nosećih vektora				
	precision	recall	f1-score	support
0	0.80	0.83	0.81	12090
1	0.80	0.76	0.78	10624
accuracy			0.80	22714
macro avg	0.80	0.79	0.80	22714
weighted avg	0.80	0.80	0.80	22714
Random forest				
	precision	recall	f1-score	support
0	0.80	0.83	0.82	12090
1	0.80	0.76	0.78	10624
accuracy			0.80	22714
macro avg	0.80	0.80	0.80	22714
weighted avg	0.80	0.80	0.80	22714

Vakcina protiv H1N1 virusa | *train*

Logistička regresija				
	precision	recall	f1-score	support
0	0.90	0.72	0.80	17858
1	0.41	0.70	0.52	4856
accuracy			0.72	22714
macro avg	0.65	0.71	0.66	22714
weighted avg	0.79	0.72	0.74	22714
Metoda nosećih vektora				
	precision	recall	f1-score	support
0	0.89	0.75	0.82	17858
1	0.42	0.67	0.52	4856
accuracy			0.73	22714
macro avg	0.66	0.71	0.67	22714
weighted avg	0.79	0.73	0.75	22714
Random forest				
	precision	recall	f1-score	support
0	0.90	0.72	0.80	17858
1	0.41	0.71	0.52	4856
accuracy			0.72	22714
macro avg	0.65	0.71	0.66	22714
weighted avg	0.79	0.72	0.74	22714

Evaluacija modela sa optimalnim hiper-parametrima

Tumačenje:

precision: Među svim ljudima za koje konačni model bude predvideo da [će se / se neće] vakcinisati, [$\text{precision } 1 > \text{precision } 0$]% će se vakcinisati

recall: Među svim ljudima koji [će se / se neće] vakcinisati, model će predvideti [$\text{recall } 1 > \text{recall } 0$]%.

Vakcina protiv sezononskog gripa | test

Logistička regresija

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.75	0.77	0.76	1344
1	0.73	0.71	0.72	1180

accuracy			0.74	2524
macro avg	0.74	0.74	0.74	2524
weighted avg	0.74	0.74	0.74	2524

Metoda nosećih vektora

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.75	0.79	0.77	1344
1	0.74	0.71	0.72	1180

accuracy			0.75	2524
macro avg	0.75	0.75	0.75	2524
weighted avg	0.75	0.75	0.75	2524

Random forest

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.76	0.78	0.77	1344
1	0.74	0.72	0.73	1180

accuracy			0.75	2524
macro avg	0.75	0.75	0.75	2524
weighted avg	0.75	0.75	0.75	2524

Vakcina protiv H1N1 virusa | test

Logistička regresija

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.89	0.70	0.78	1984
1	0.38	0.69	0.49	540

accuracy			0.70	2524
macro avg	0.64	0.69	0.64	2524
weighted avg	0.78	0.70	0.72	2524

Metoda nosećih vektora

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.89	0.73	0.80	1984
1	0.40	0.66	0.50	540

accuracy			0.71	2524
macro avg	0.64	0.69	0.65	2524
weighted avg	0.78	0.71	0.73	2524

Random forest

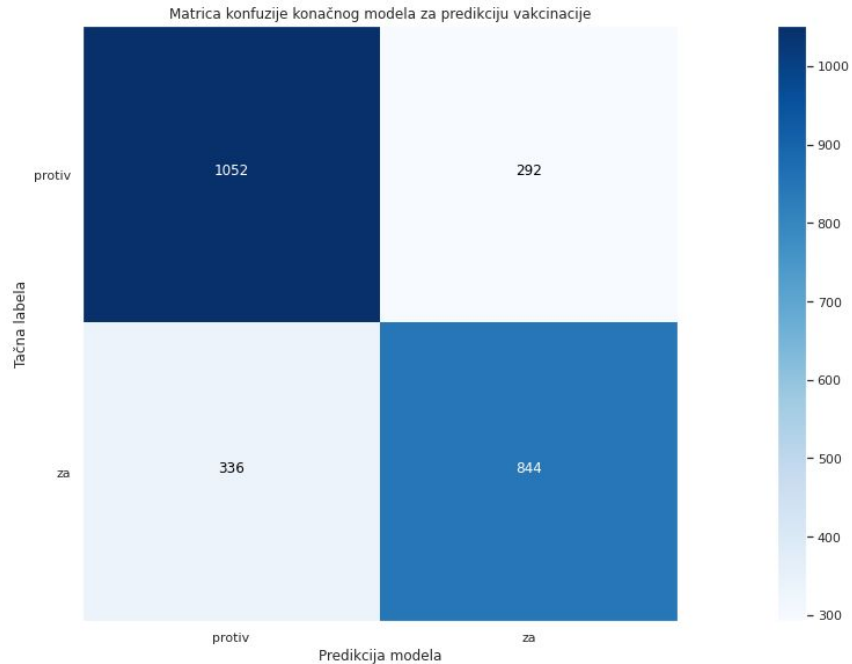
	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.89	0.70	0.79	1984
1	0.39	0.70	0.50	540

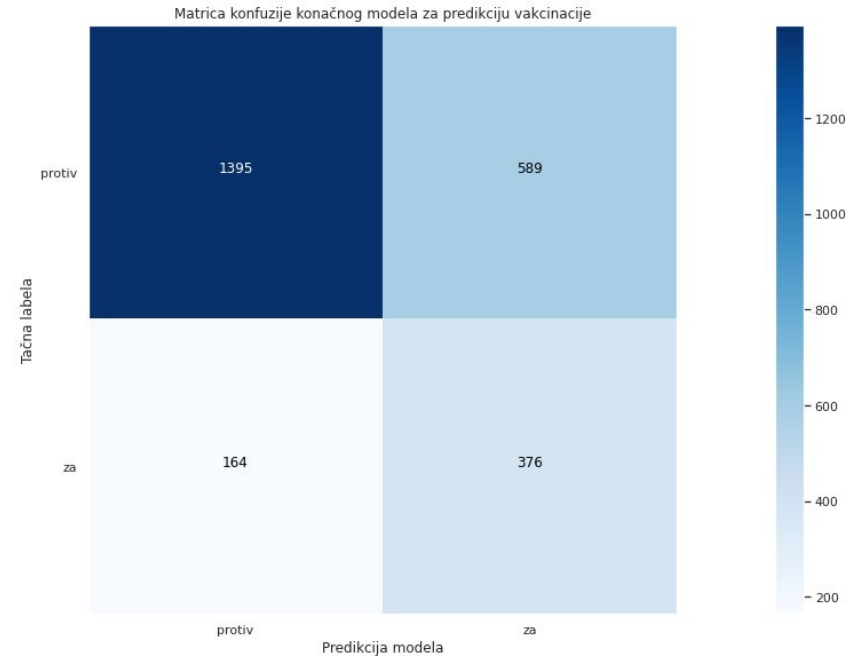
accuracy			0.70	2524
macro avg	0.64	0.70	0.64	2524
weighted avg	0.79	0.70	0.73	2524

Matrica konfuzije na testirajućem skupu

Vakcina protiv sezonskog gripa



Vakcina protiv H1N1 virusa



Hvala na pažnji!