

KLASIFIKACIJA

JELENA JOVANOVIĆ

jelena.jovanovic@fon.bg.ac.rs

PREGLED PREDAVANJA

- Šta je klasifikacija?
- Tipovi klasifikacije
- Algoritmi klasifikacije
- Mere uspešnosti klasifikatora

ŠTA JE KLASIFIKACIJA?

- Zadatak određivanja klase kojoj neka instanca pripada
 - instanca je opisana vrednošću atributa
 - skup mogućih klasa je poznat i dat

ŠTA JE KLASIFIKACIJA?

- Zadatak određivanja klase kojoj neka instanca pripada
 - instanca je opisana vrednošću atributa
 - skup mogućih klasa je poznat i dat
- Klase su date kao nominalne vrednosti, npr.
 - klasifikacija email poruka: spam, not-spam
 - klasifikacija novinskih članaka: politika, sport, kultura i sl.

TIPOVI KLASIFIKACIJE

Zavisno od broja klasa, razlikujemo:

- binarnu klasifikaciju (*binary classification*) - postoje dve klase

TIPOVI KLASIFIKACIJE

Zavisno od broja klasa, razlikujemo:

- binarnu klasifikaciju (*binary classification*) - postoje dve klase
- više-klasnu klasifikaciju (*multi-class classification*) - postoji tri ili više klasa i svaka instanca može pripadati samo jednoj klasi

TIPOVI KLASIFIKACIJE

Zavisno od broja klasa, razlikujemo:

- binarnu klasifikaciju (*binary classification*) - postoje dve klase
- više-klasnu klasifikaciju (*multi-class classification*) - postoji tri ili više klasa i svaka instanca može pripadati samo jednoj klasi
- klasifikaciju sa više oznaka (*multi-label classification*) - postoji više klasa i svaka instanca može pripadati proizvoljnom broju klasa ili ni jednoj od njih

ALGORITMI KLASIFIKACIJE

Postoje brojni pristupi/algoritmi za klasifikaciju:

- Logistička regresija
- Naïve Bayes
- Algoritmi iz grupe Stabala odlučivanja
- k-Nearest Neighbor (kNN)
- Algoritmi iz grupe Neuronskih mreža
- Support Vector Machines (SVM)
- ...

MERE USPEŠNOSTI KLASIFIKATORA

Neke od najčešće korišćenih metrika:

- Matrica zabune (Confusion Matrix)
- Tačnost (Accuracy)
- Preciznost (Precision) i Odziv (Recall)
- F mera (F measure)
- Površina ispod ROC krive (Area Under the Curve - AUC)

MATRICA ZABUNE (CONFUSION MATRIX)

Služi kao osnova za računanje mera performansi (uspešnosti) algoritama klasifikacije

		Predicted Class	
		Yes	No
Actual Class	Yes	TP	FN
	No	FP	TN

TP = True Positive

FP = False Positive

TN = True Negative

FN = False Negative

TAČNOST (ACCURACY)

Tačnost (Accuracy) predstavlja procenat slučajeva (instanci) koji su uspešno (korektno) klasifikovani

$$\text{Accuracy} = (TP + TN) / N$$

gde je:

- TP – True Positive; TN – True Negative
- N – ukupan broj uzoraka (instanci) u skupu podataka

		Predicted Class	
		Yes	No
Actual Class	Yes	TP	FN
	No	FP	TN

TAČNOST (ACCURACY)

U slučaju vrlo neravnomerne raspodele instanci između klasa (tzv. skewed classes), ova mera je nepouzdana

Npr. u slučaju klasifikacije poruka na spam vs. not-spam, možemo imati skup za trening sa 0.5% spam poruka

Ako primenimo “klasifikator” koji svaku poruku svrstava u not-spam klasu, dobijamo tačnost od 99.5%

Očigledno je da ova metrika nije pouzdana i da su u slučaju skewed classes potrebne druge metrike

COHEN'S KAPPA (K)

Kappa govori koliko je naš model bolji u odnosu na nasumični klasifikator koji predviđa na osnovu “veličine” klasa (tj broja instanci svake od klasa)

Koristi se umesto Tačnosti u slučaju nebalansiranih skupova podataka

Vrednosti se kreću u opsegu $[-1, 1]$, gde 0 označava performance ekvivalentne nasumičnoj klasifikaciji

COHEN'S KAPPA (K)

Za izračunavanje su potrebne dve vrednosti:

- Opažena saglasnost (*observed agreement*, p_o), što je izračunata tačnost našeg klasifikatora
- Očekivana saglasnost (*expected agreement*, p_e), što je tačnost nasumičnog klasifikatora tj klasifikatora koji vrši predviđanje na osnovu zastupljenosti klasa u trening setu

$$k = \frac{p_o - p_e}{1 - p_e}$$

PRECIZNOST (PRECISION) I ODZIV (RECALL)

Precision = TP / no. predicted positive = TP / (TP + FP)

Npr. od svih poruka koje su *označene kao spam* poruke, koji procenat čine poruke koje su stvarno spam

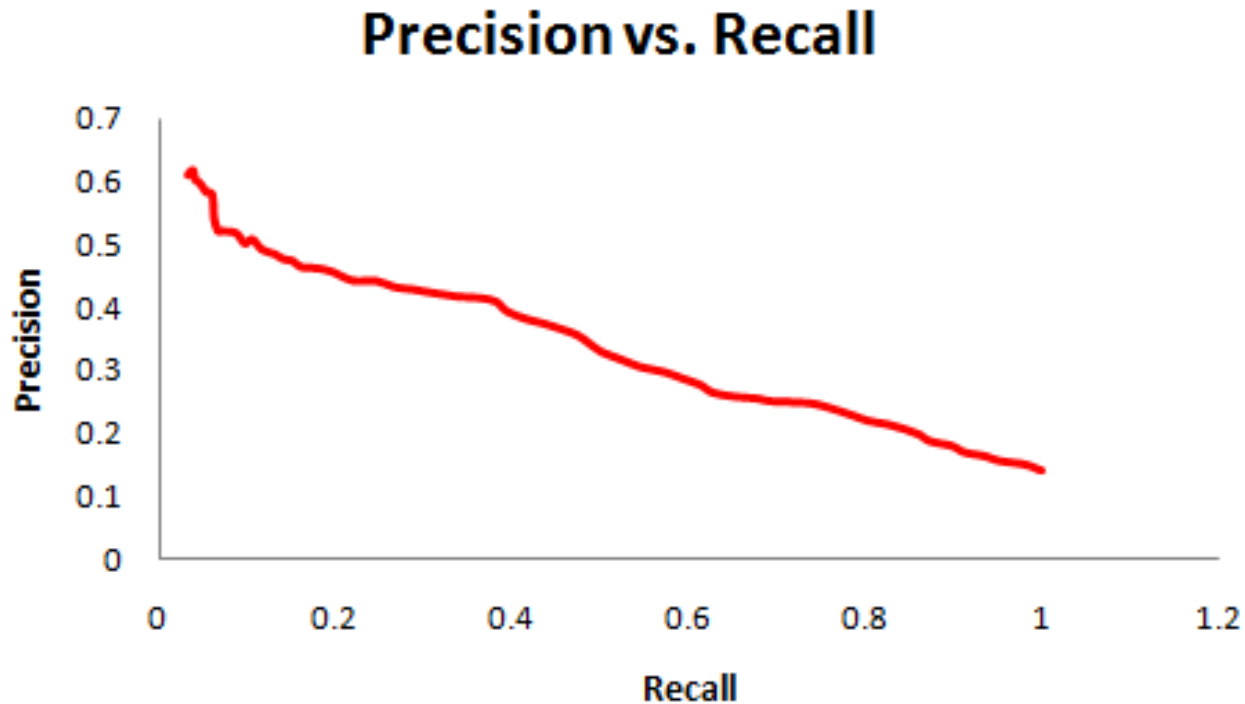
Recall = TP / no. actual positive = TP / (TP + FN)

Npr. od svih poruka koje su *stvarno spam* poruke, koji procenat poruka je detektovan/klasifikovan kao spam

		Predicted Class	
		Yes	No
Actual Class	Yes	TP	FN
	No	FP	TN

PRECIZNOST I ODZIV

U praksi je nužno praviti kompromis između ove dve mere: ako želimo da povećamo Odziv, smanjićemo Preciznost, i obrnuto.



F MERA (F MEASURE)

F mera kombinuje Preciznost i Odziv i omogućuje jednostavnije poređenje dva ili više modela

$$F = (1 + \beta^2) * \text{Precision} * \text{Recall} / (\beta^2 * \text{Precision} + \text{Recall})$$

Parametar β kontroliše koliko više značaja će se pridavati Odzivu u odnosu na Preciznost

U praksi se najčešće koristi tzv. F1 mera („balansirana“ F mera) koja daje podjednak značaj i Preciznosti i Odzivu:

$$F1 = 2 * \text{Precision} * \text{Recall} / (\text{Precision} + \text{Recall})$$

POVRŠINA ISPOD ROC KRIVE

Površina ispod ROC* krive – Area Under the Curve (AUC):

- meri diskriminacionu moć klasifikatora tj. sposobnost da razlikuje instance koje pripadaju različitim klasama
- primenjuje se za merenje performansi binarnih klasifikatora
- vrednost za AUC se kreće u intervalu 0-1
- za metodu slučajnog izbora važi da je $AUC = 0.5$; što je AUC vrednost klasifikatora > 0.5 , to je klasifikator bolji
 - 0.7–0.8 se smatra prihvatljivim; 0.8–0.9 jako dobrim; sve > 0.9 je odlično

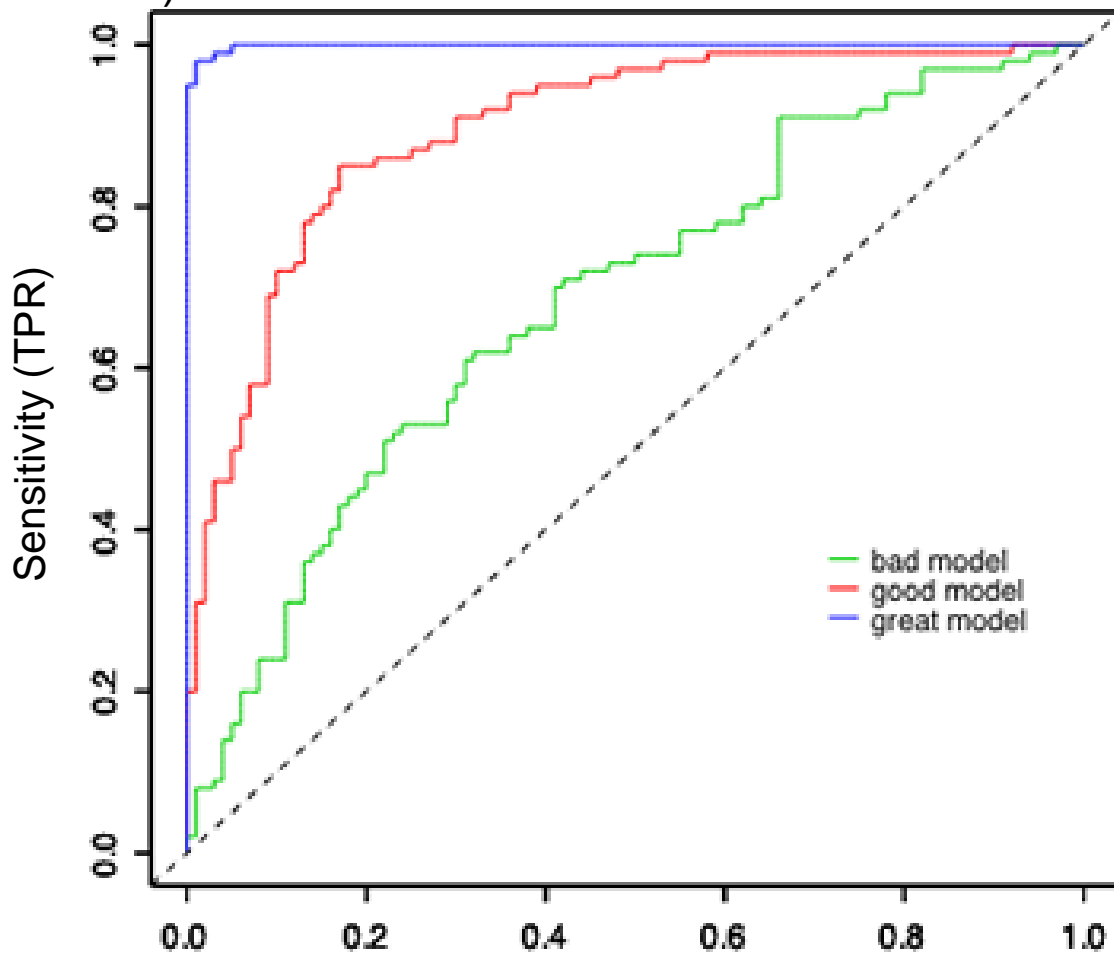
*ROC = Receiver Operating Characteristic;

http://en.wikipedia.org/wiki/Receiver_operating_characteristic

POVRŠINA ISPOD ROC KRIVE

Sensitivity ili True Positive Rate (TPR) ili Recall

$$\text{TPR} = \text{TP}/(\text{TP} + \text{FN})$$



1 - Specificity (TNR)

True Negative Rate (TNR)

$$\text{TNR} = \text{TN}/(\text{FP} + \text{TN})$$

(Anonimni) upitnik za vaše
komentare, predloge, kritike:

<https://forms.gle/7vdwezRAumx72VnXA>