

MACHINE LEARNING BIDEZKO MEZU TOXIKOEN DETEKZIOA

Ane Acha Gonzalez

Gradu Amaierako Lana
Matematikako Gradua

Mikel Peñagarikano
irakasleak zuzendutako lana
Leioan, 2024ko uztailaren 11a

SARRERA

HELBURUA

Sare sozial zehatz bateko mezu toxikoak automatikoki iragaztea da, Machine Learning-ean oinarritutako klasifikazio ereduak erabiliz.

AURKIBIDEA

1. Problemaren eta datu basearen azalpena.
2. Naive Bayes.
3. Bektoreen bidezko eredua.
4. Erregresio Logistikoa.
5. Support Vector Machine (Euskarri bektoreko makinak).
6. Decision Tree (Erabaki zuhaitzak) eta Random Forest (Ausazko basoak).
7. Ondorioak.

1. PROBLEMAREN ETA DATU BASEAREN AZALPENA.

- Kaggle → *Quora Insincere Questions Classification*

	qid	question_text	target
0	00002165364db923c7e6	How did Quebec nationalists see their province...	0
1	000032939017120e6e44	Do you have an adopted dog, how would you enco...	0
2	0000412ca6e4628ce2cf	Why does velocity affect time? Does velocity a...	0
3	000042bf85aa498cd78e	How did Otto von Guericke used the Magdeburg h...	0
4	0000455dfa3e01eae3af	Can I convert montra helicon D to a mountain b...	0
...
1306117	ffffcc4e2331aaf1e41e	What other technical skills do you need as a c...	0
1306118	ffffd431801e5a2f4861	Does MS in ECE have good job prospects in USA ...	0
1306119	ffffd48fb36b63db010c	Is foam insulation toxic?	0
1306120	fffec519fa37cf60c78	How can one start a research project based on ...	0
1306121	fffed09fedb5088744a	Who wins in a battle between a Wolverine and a...	0

1. PROBLEMAREN ETA DATU BASEAREN AZALPENA.

- Datuen analisia:
 - Aurreprozesatze teknikak
 - Erantzun aldagaia etiketak: 0 klasea (ez - toxikoa) \rightarrow %94
1 klasea (toxikoa) \rightarrow %6
- Datuen banaketa:
Entrenamendua (% 80), balidazioa (% 10) eta testa (% 10)
- Ebaluazioa:
 - Zehaztasuna (Accuracy)
 - Prezisioa (Precision)
 - Sentikortasuna (Recall)
 - Espezifikotasuna (Specificity)
 - ROC AUC
 - **Precision-Recall AUC**

2. NAIVE BAYES.

2. 1. OINARRI TEORIKOA

Helburua: $P(Y = y_j \mid \mathbf{X} = \mathbf{x}_i)$

$$P(y_j \mid \mathbf{x}_i) = \frac{P(y_j) \prod_{k=1}^n P(x_k \mid y_j)}{P(\mathbf{x}_i)} \quad (1)$$

$$\hat{y} = \operatorname{argmax}_{y_j \in Y} P(y_j \mid \mathbf{x}_i) = \operatorname{argmax}_{y_j \in Y} P(y_j) \prod_{k=1}^n P(x_k \mid y_j) \quad (2)$$

$$\hat{y} = \operatorname{argmax}_{y_j \in Y} \log P(y_j \mid \mathbf{x}_i) = \operatorname{argmax}_{y_j \in Y} \log P(y_j) + \sum_{k=1}^n \log P(x_k \mid y_j) \quad (3)$$

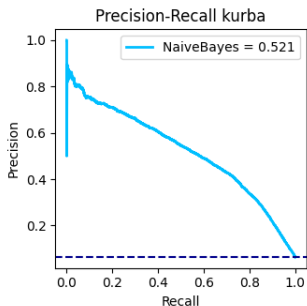
2. NAIVE BAYES.

2. 2. ENTRENAMENDUA

- $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p)$ mezuen multzoa eta $\mathbf{x}_i = (x_{i1}, \dots, x_{in})$
- $Y = \{0, 1\}$ klaseen multzo bitarra.
- $\hat{y} = \underset{y_j \in \{0,1\}}{\operatorname{argmax}} \log P(y_j | \mathbf{x}_i) = \underset{y_j \in \{0,1\}}{\operatorname{argmax}} \log P(y_j) + \sum_{k=1}^n \log P(x_k | y_j)$

2.3. EBALUAZIOA

ALGORITMOA	NAIVE BAYES
PR_AUC	0.52
ROC_AUC	0.92
Zehaztasuna	0.92
Prezisia	0.41
Sentikortasuna	0.73
Espezifikotasuna	0.93



3. BEKTOREEN BIDEZKO EREDUA.

3. 1. BEKTORIZAZIOA

$$M = \begin{pmatrix} m_{11} & \dots & m_{1q} \\ \vdots & \ddots & \vdots \\ m_{p1} & \dots & m_{pq} \end{pmatrix} \in \mathbb{N}^{p \times q} \quad (4)$$

3. 2. EREDU BEKTORIALAREN OINARRIAK

$$M \cdot V + J_{p \times 1} \cdot L = \begin{pmatrix} \log(P(\mathbf{x}_1 | 0)P(0)) & \log P((\mathbf{x}_1 | 1)P(1)) \\ \log(P(\mathbf{x}_2 | 0)P(0)) & \log P((\mathbf{x}_2 | 1)P(1)) \\ \vdots & \vdots \\ \log P((\mathbf{x}_p | 0)P(0)) & \log P((\mathbf{x}_p | 1)P(1)) \end{pmatrix} \quad (5)$$

non

$$V = \begin{pmatrix} \log P(x_1 | 0) & \log P(x_1 | 1) \\ \log P(x_2 | 0) & \log P(x_2 | 1) \\ \vdots & \vdots \\ \log P(x_q | 0) & \log P(x_q | 1) \end{pmatrix}, \quad J_{p \times 1} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \text{ eta } L = (\log P(0) \quad \log P(1))$$

4. ERREGRESIO LOGISTIKOA

4. 1. OINARRI TEORIKOA.

Izan bitez Y banaketa binomialari darraion erantzun aldagaia eta $\mathbf{X} = (X_1, X_2, \dots, X_p)^T$ non X_1, X_2, \dots, X_p p aldagai aske diren.

$$P(Y = 1 \mid \mathbf{X}) = p(\mathbf{X}) = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}} \in (0, 1) \quad (6)$$

Logit transformazioa:

$$g(p(\mathbf{X})) = \text{logit}(p(\mathbf{X})) = \ln \left(\frac{p(\mathbf{X})}{1 - p(\mathbf{X})} \right) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p \quad (7)$$

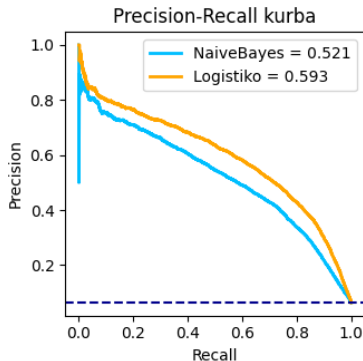
Parametroen estimazioa eta interpretazioa

4. ERREGRESIO LOGISTIKOA

4. 2. ENTRENAMENDUA ETA EBALUAZIOA.

- *LogisticRegression*

ALGORITMOA	NAIVE BAYES	LOGISTIKOA
PR_AUC	0.52	0.59
ROC_AUC	0.92	0.94
Zehaztasuna	0.92	0.95
Prezisia	0.41	0.68
Sentikortasuna	0.73	0.40
Espezifikotasuna	0.93	0.99



5. SUPPORT VECTOR MACHINE.

5. 1. OINARRI TEORIKOA.

Linealki banandutako kasuak

$$\mathbf{w} \cdot \mathbf{x} + b = 0 \quad (8)$$

Hiperplano optimoaren bilaketa \mathbf{w} eta b balioak aurkitzeko optimizazio problema gisa adieraz daiteke:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 \geq 0 \quad & i = 1, \dots, n \end{aligned} \quad (9)$$

Problema duala:

$$\begin{aligned} \max_{\alpha} \mathcal{L}(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j \\ \sum_{i=1}^n \alpha_i y_i = 0 \quad \text{eta} \quad \alpha_i \geq 0 \quad & i = 1, \dots, n \end{aligned} \quad (10)$$

5. SUPPORT VECTOR MACHINE.

Linealki kuasi-bananduak

Optimizazio problema:

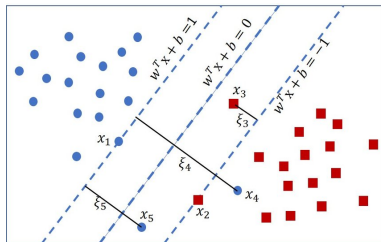
$$\min_{\mathbf{w}, b, \xi} \quad \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \quad (11)$$

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) + \xi_i - 1 \geq 0, \quad \xi_i \geq 0, \quad i = 1, \dots, n$$

Lortzen den problema duala:

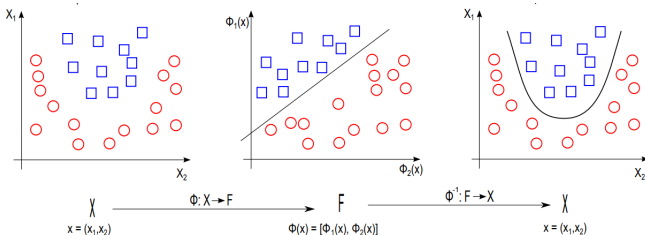
$$\max_{\alpha} \quad \mathcal{L}(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j \quad (12)$$

$$\sum_{i=1}^n \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C, \quad i = 1, \dots, n$$



5. SUPPORT VECTOR MACHINE.

Linealki banangarriak ez diren kasuak



Problema duala:

$$\max_{\alpha} \mathcal{L}(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) \quad (13)$$

$$\sum_{i=1}^n \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C, \quad i = 1, \dots, n$$

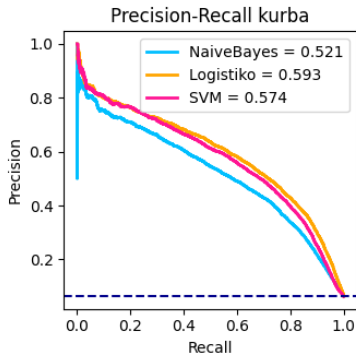
non $K(\mathbf{x}_i, \mathbf{x}_j)$ Kernelen funtzioa den.

5. SUPPORT VECTOR MACHINE.

5. 2. ENTRENAMENDUA ETA EBALUAZIOA

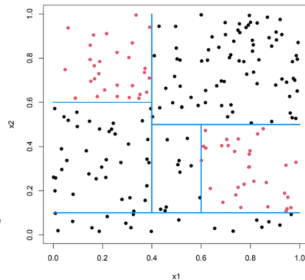
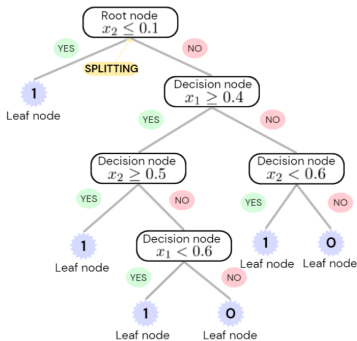
- *LinearSVC*

ALGORITMOA	NAIVE BAYES	LOGISTIKOA	SMV
PR_AUC	0.52	0.59	0.57
ROC_AUC	0.92	0.94	0.93
Zehaztasuna	0.92	0.95	0.95
Prezisia	0.41	0.68	0.67
Sentikortasuna	0.73	0.40	0.39
Espezifikotasuna	0.93	0.99	0.99



6. DECISION TREE ETA RANDOM FOREST.

6. 1. OINARRI TEORIKOA



6. DECISION TREE ETA RANDOM FOREST.

Zatiketa bitarraren ezpurutasunaren murrizketa:

$$\Delta i(s, t) = i(t) - p_L i(t_L) - p_R i(t_R) \quad (14)$$

Zatiketa optimoa:

$$s^* = \arg \max_{s \in Q} \Delta i(s, t) \quad (15)$$

Shannon-en entropia eta *Gini indizea*, hurrenez hurren:

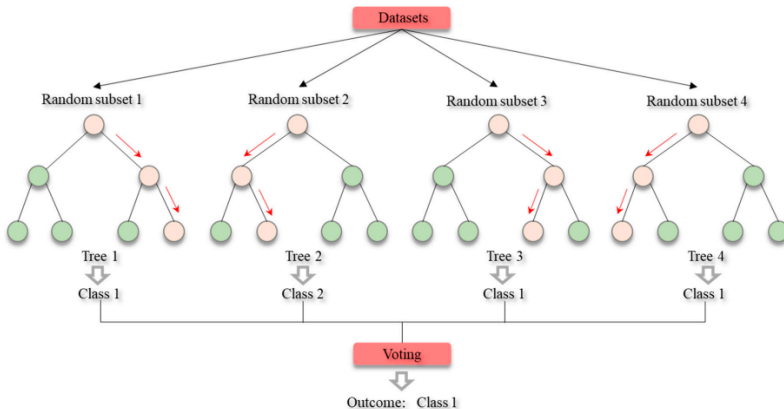
$$i_H(t) = - \sum_{k=1}^J p(c_k|t) \log_2(p(c_k|t)) \quad (16)$$

$$i_G(t) = \sum_{k=1}^J p(c_k|t) (1 - p(c_k|t)) \quad (17)$$

non $p(c_k|t)$ t nodoa c_k klasekoa izateko probabilitatea den.

6. DECISION TREE ETA RANDOM FOREST.

Gehiegizko doikuntza (overfitting) ekiditeko: Random Forest.



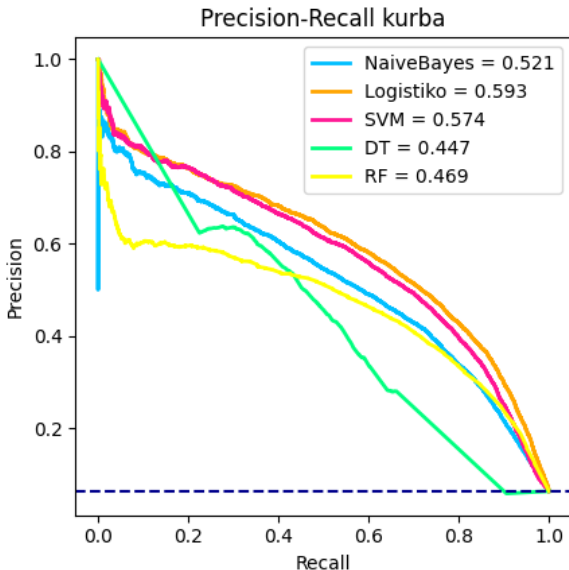
6. DECISION TREE ETA RANDOM FOREST.

6. 2. ENTRENAMENDUA ETA EBALUAZIOA

- *DecisionTreeClassifier* eta *RandomForestClassifier*

ALGORITMOA	NAIVE BAYES	LOGISTIKOA	SMV	DECISION TREE	RANDOM FOREST
PR_AUC	0.52	0.59	0.57	0.45	0.47
ROC_AUC	0.92	0.94	0.93	0.75	0.92
Zehaztasuna	0.92	0.95	0.95	0.94	0.94
Prezisia	0.41	0.68	0.67	0.62	0.60
Sentikortasuna	0.73	0.40	0.39	0.33	0.15
Espezifikotasuna	0.93	0.99	0.99	0.99	0.99

6. DECISION TREE ETA RANDOM FOREST.



7. ONDORIOAK.

- Algoritmo bakoitzak bere abantailak eta desabantailak.
- Hainbat arazori egin behar diete aurre.
- Erregresio logistikoa eta SVM eredurik eraginkorrenak.
- Ezinbestekoa da hauen mugak ezagutzea, matematikoki ulertzea eta egoki ebaluatzea.

Eskerrik asko !