



**Vilniaus universitetas  
Duomenų mokslo ir skaitmeninių  
technologijų institutas  
L I E T U V A**



---

INFORMATIKA (N009)

---

**KONTEKSTINIS DUOMENŲ, APRAŠOMŲ  
ERDVINIAIS APIBENDRINTAIS  
TIESINIAIS MODELIAIS,  
KLASIFIKAVIMAS**

**Eglė Zikarienė**

2019 m. spalį

Mokslinė ataskaita DMSTI-DS-N009-19-4

VU Duomenų mokslo ir skaitmeninių technologijų institutas, Akademijos g. 4,

Vilnius LT-08412

[www.mii.lt](http://www.mii.lt)

## **Santrauka**

Šis darbas apima teorinę disertacijos dalį, kurioje pateikiama kontekstinio klasifikavimo klaidų tikimybių skaičiavimo ir vertinimo metodika. Siūlomi įvairūs aktualių klasifikavimo klaidų (ang. actual error rates) įvertiniai, kurių reikšmės panaudojamos nustatant klasifikatorių efektyvumą (ang. performance of classifiers) ir palyginant erdviųjų duomenų modelius. Pradėti tirti duomenys apie Šakotųjų banguolių padengimą Baltijos jūroje, analizei taikomas erdvinės Beta regresijos modelis. Trumpai aprašomi parašyti straipsniai šia tema ir koks numatomas darbų praplėtimas disertacijos tema.

**Reikšminiai žodžiai:** Erdvinė analizė, Bajeso diskriminantinė funkcija, Tikroji klaidos tikimybė.

# Turinys

---

1	Įvadas .....	4
2	Klasifikavimo procedūra.....	4
3	Klasifikavimo klaidos vertinimo kriterijai .....	6
4	Erdvinės beta regresijos modelių pritaikymas dumblių padengimo analizei.....	8
5	Išvados .....	8
6	Literatūra.....	8

# 1 Įvadas

Parengta erdvinių duomenų, aprašomų apibendrintais tiesiniais modeliais, kontekstinio klasifikavimo klaidų tikimybių skaičiavimo ir vertinimo metodika. Pagrindinis disertacijos darbo tikslas – atlikti erdvinių duomenų klasifikavimą taikant Bajeso diskriminantinės funkcijos procedūrą. Diskriminantinės funkcijos modeliams paprastai naudojamas tik Gauso skirstinys, todėl šioje disertacijoje planuojamas šių modelių išplėtimas dviem kryptimis: eksponentinių skirstinių šeimos nariais tiek tolydžiu, tiek diskrečiu atveju ir eliptinių skirstinių šeimos nariais.

Orientuojamasi į klasifikavimo metodus, pagrįstus sąlyginiais klasifikuojamo erdvės taško (ang. focal location) skirstiniais, priklausančiais eksponentinių skirstinių šeimai. Siūlomi įvairūs tikrųjų klasifikavimo klaidų (ang. actual error rates) įvertiniai, kurių reikšmės panaudojamos nustatant klasifikatorių efektyvumą (ang. performance of classifiers) ir palyginant erdvinių duomenų modelius. Siūlomi tikros klasifikavimo klaidos vertinimo metodai: R-metodas, aposteriorinės tikimybės metodas, U-metodas.

## 2 Klasifikavimo procedūra

Darbe nagrinėsime atsitiktinio lauko stebinius  $\{Z(s) : s \in D \subset R^p\}$ . Stebėjimo  $Z(s)$  modelis populiacijoje  $\Omega_l$  yra

$$Z(s) = x'(s)\beta_l + \varepsilon(s),$$

čia  $x(s)$  yra  $q \times 1$  neatsitiktinių regresorių vektorius, ir  $\beta_l$  yra  $q \times 1$  parametru vektorius,  $l=1, \dots, L$ . Klaidos narys atsitiktinis laukas  $\{\varepsilon(s) : s \in D\}$  generuojamas stacionarus nulinio vidurkio laukas su kovariacine funkcija, apibrėžiama pagal modelį visiems  $s, u \in D$

$$\text{cov}\{\varepsilon(s), \varepsilon(u)\} = \sigma^2 r(s-u),$$

čia  $\sigma^2$  yra skalės parametras, o  $r(\cdot)$  - erdvinės koreliacijos funkcija. (Dučinskas *et al.* 2015)

Tarkime, nagrinėjamas  $Z_0 = Z(s_0)$  klasifikavimo uždavinys su mokymo imtimi  $T$  į vieną iš aukščiau aprašytų populiacijų.

Pristatant klasifikavimo procedūrą paprastumo dėlei pereisime prie modelio dviejų klasių atveju, tai yra  $l=1, 2$ . Tuomet Bajeso diskriminantinė funkcija (BDF), minimizuojanti klaidingo klasifikavimo tikimybę, yra suformuota pagal sąlyginio tankio santykio logaritmą (Fukunaga, 1990). Jos bendrąjį pavidalą galima užrašyti taip:

$$W(Z_0|T) = \ln \left( \frac{p(Z_0|T, \Omega_1)P(\Omega_1)}{p(Z_0|T, \Omega_2)P(\Omega_2)} \right) = \ln \left( \frac{p(Z_0|T, \Omega_1)}{p(Z_0|T, \Omega_2)} \right) + \ln \left( \frac{P(\Omega_1)}{P(\Omega_2)} \right)$$

čia  $Z_0$  - klasifikuojamas stebinis,  $T$  - mokymo imtis,  $\Omega_l$  - objekto klasė,  $l=1, 2$  - klasės numeris,  $p(Z_0|T, \Omega_l)$  - sąlyginis tankis,  $P(\Omega_l)$  - apriorinė klasės tikimybė. Sudarant BDF naudojamos sąlyginio tankio funkcijos. Praktikoje mokymo imtis ir klasifikuojamas taškas yra tame pačiame erdviniam koreliuotame lauke, todėl į aptartas diskriminantines funkcijas įvedama klasifikuojamo taško erdvinė priklausomybė su mokymo imtimi (Dučinskas 2009), (Stabingienė *et al.* 2010).

Eksponentinių skirstinių šeimai bendrąją sąlyginio tankio išraišką galima užrašyti taip:

$$p(Z_0|T; \theta) = \exp\{A_0(T; \theta)T_0(Z_0) - B_0(T; \theta) + C_0(Z_0)\}$$

čia  $T_0(Z_0)$  - minimali pakankama statistika,  $C_0(\cdot)$  yra funkcija nuo  $Z_0$ ,  $A_0(\cdot)$  ir  $B_0(\cdot)$  yra funkcijos nuo  $T$  ir nežinomų parametrų  $\theta$  (Furukawa 2004). Tinkamas  $A_0(\cdot)$  pasirinkimas apibrėžia priklausomybės tipą tarp kaimyninių reikšmių, o  $B_0(\cdot)$  tuomet yra tinkama normalizavimo funkcija.

Šiame darbe nagrinėjant eksponentinių skirstinių šeimą laikomasi prielaidos, kad tarp erdvinių duomenų galioja porinis (*ang. pairwise*) ryšys, kai tankio funkciją galima užrašyti tokiu pavidalu:

$$p(z) = f\left(\alpha + \sum_{1 \leq i \leq n} z_i \beta_i + \sum_{1 \leq i < j \leq n} z_i z_j \beta_{ij}\right)$$

$\alpha, \beta_i, \beta_{ij}$  - funkcijos koeficientai.

Tokiu atveju eksponentinių šeimoms skirstinių modeliams, aprašomiems aukščiau pateikta sąlyginio tankio funkcija su galiojančia prielaida, įvedamas terminas *auto-modeliai* (Besag 1974). Gautos BDF išraiškos atsitiktinių laukų auto modeliams.

#### Auto-Gamma modelis

Įveskime pažymėjimus:  $S^0 = (s_0, S)$  - bendra klasifikuojamo taško ir mokymo imties vietų aibė, kai  $S = (s_1, \dots, s_n)$  - mokymo vietų aibė.  $T'_{0l} = (z'_0, T')$  - klasifikuojamo taško reikšmės ir mokymo imties aibė, čia  $T' = (z_1, \dots, z_n)$  - mokymo imties aibė. Sąlyginį Gamma skirstinį galima užrašyti taip:

$$Z'_0|T \sim \Gamma(\alpha^l, \gamma_0(T))$$

kai  $\alpha^l$  - mastelio parametras kiekvienai klasei  $l$ ,  $\gamma_0(T)$  - formos parametras, kurio išraiška  $\gamma_0(T) = \theta_0 + \sum_{t=1}^n \theta_{0t} z_t$ . Tuomet sąlyginį Gamma skirstinio tankį galima užrašyti taip:

$$p(Z'_0 = z_0|T) = \exp\left\{-\frac{1}{\alpha^l} z_0 + \left(\theta_0 + \sum_{t=1}^n \theta_{0t} z_t\right) \ln z_0 + \left(\theta_0 + \sum_{t=1}^n \theta_{0t} z_t\right) \ln\left(\frac{1}{\alpha^l}\right) - \ln\left(\Gamma\left(\theta_0 + \sum_{t=1}^n \theta_{0t} z_t\right)\right) - \ln z_0\right\}.$$

Sąlyginio Gamma skirstinio atveju gaunama BDF išraiška:

$$W(Z_0|T) = \ln\left(\frac{p_1(Z_0|T, \Omega_1)}{p_2(Z_0|T, \Omega_2)}\right) + \ln\left(\frac{P(\Omega_1)}{P(\Omega_2)}\right) = \left(\frac{\alpha^1 - \alpha^2}{\alpha^1 \alpha^2}\right) z_0 + \left(\theta_0 + \sum_{t=1}^n \theta_{0t} z_t\right) \ln\left(\frac{\alpha^2}{\alpha^1}\right) + \gamma$$

$$\text{kai } \gamma = \ln\left(\frac{P(\Omega_1)}{P(\Omega_2)}\right).$$

#### Auto-Beta modelis

Sąlyginį Beta skirstinį galima užrašyti taip:

$$Z'_0|T \sim \text{Beta}(p^l, q),$$

čia  $p, q$ - mastelio parametrai,  $p$  - skirtingi kiekvienai klasei  $l$ ,  
 $p^l(T) = \theta_0^l + \sum_{t=1}^n \theta_{0t} z_t$ . Tuomet Beta skirstinio tankį galima užrašyti taip:

$$p(Z_0^l = z_0 | T) = \exp \left\{ \left[ \left( \theta_0^l + \sum_{t=1}^n \theta_{0t} z_t \right) - 1 \right] \ln z_0 + (q-1) \ln(1-z_0) + \right. \\ \left. + \ln \left( \frac{\Gamma \left( \theta_0^l + \sum_{t=1}^n \theta_{0t} z_t + q \right)}{\Gamma \left( \theta_0^l + \sum_{t=1}^n \theta_{0t} z_t \right) \Gamma(q)} \right) \right\}$$

Sąlyginio Beta skirstinio atveju gaunama BDF išraiška:

$$W(Z_0 | T) = \ln \left( \frac{p_1(Z_0 | T, \Omega_1)}{p_2(Z_0 | T, \Omega_2)} \right) + \ln \left( \frac{P(\Omega_1)}{P(\Omega_2)} \right) = (p^1 - p^2) \ln z_0 + \left( \frac{\Gamma(p^1 + q) \Gamma(p^2)}{\Gamma(p^1) \Gamma(p^2 + q)} \right) + \ln \left( \frac{P(\Omega_1)}{P(\Omega_2)} \right) = \\ = (\theta_0^1 - \theta_0^2) \ln z_0 + \ln \left( \frac{\Gamma \left( \theta_0^1 + \sum_{t=1}^n \theta_{0t} z_t + q \right) \Gamma \left( \theta_0^2 + \sum_{t=1}^n \theta_{0t} z_t \right)}{\Gamma \left( \theta_0^1 + \sum_{t=1}^n \theta_{0t} z_t \right) \Gamma \left( \theta_0^2 + \sum_{t=1}^n \theta_{0t} z_t + q \right)} \right) + \gamma$$

$$\text{kai } \gamma = \ln \left( \frac{P(\Omega_1)}{P(\Omega_2)} \right).$$

Paprastai klasifikuojant objektus labai dažnai populiacijos parametrai yra nežinomi. Tokiu atveju yra naudojami jų įvertiniai, o pačios funkcijos yra vadinamos įterptomomis Bajeso diskriminantinėmis funkcijomis (*ang. Plug-in Bayesian discriminant function (PBDF)*) (Dučinskas 2009). Įterptoji Bajeso diskriminantinė funkcija, užrašoma taip:

$$W(Z_0 | T, \hat{\theta}) = \ln \left( \frac{p(Z_0 | T, \Omega_1, \hat{\theta}_1) P(\Omega_1)}{p(Z_0 | T, \Omega_2, \hat{\theta}_2) P(\Omega_2)} \right) = \ln \left( \frac{p(Z_0 | T, \Omega_1, \hat{\theta}_1)}{p(Z_0 | T, \Omega_2, \hat{\theta}_2)} \right) + \ln \left( \frac{P(\Omega_1)}{P(\Omega_2)} \right),$$

čia  $\hat{\theta}_l$  - nežinomų parametų vertinimai klasėje  $l$ . Nežinomi parametų įverčiai pasirenkami: ML, REML, Bajeso parametų įvertinimai.

Klasifikuojant stebinius svarbu įvertinti, kaip tiksliai tai atliekama. Dažnai taikomos tokios charakteristikos, kaip bendras klasifikavimo tikslumas bei klaidingo klasifikavimo tikimybių įverčiai.

### 3 Klasifikavimo klaidos vertinimo kriterijai

Klasifikavimo kokybę nusako klaidingo klasifikavimo tikimybių įverčiai, kurie parodo, kokia yra tikimybė suklysti klasifikavimo metu kiekvienai iš klasių (Čekanavičius ir Murauskas 2008).

Tikroji (*angl. or actual error rate or conditional error rate*) – tai klasifikatoriaus, klaidingai klasifikuojančio atsitiktinai parinktą objektą, tikėtina tikimybė. Tai klaidos tikimybė su be galo didele bandymo aibe, sudaryta iš to paties skirstinio, kaip mokymo duomenų.

Tikroji klaidos tikimybė (*ang. actual error rate (AER)*) Bajeso diskriminantinei funkcijai  $W(Z_0 | T, \hat{\theta})$ , yra apibrėžiama taip (Dučinskas 2009):

$$P(\hat{\theta}) = \sum_l \pi_l \hat{P}_{0l}$$

Čia, kai  $l = 1, 2$ ,

$$\hat{P}_{0l}(t) = P_{0l} \left( (-1)^l W(Z_0 | T, \hat{\theta}) > 0 | \Omega_l \right)$$

$P_{0l}$  yra tikimybė stebinio reikšmę  $Z_0$  priskirti ne tai klasei (kai diskriminantinės funkcijos  $W(Z_0 | T, \hat{\theta})$  reikšmė yra didesnė už 0, priskiriame stebinį  $Z_0$  vienai klasei, kai  $W(Z_0 | T, \hat{\theta})$  - kitai klasei). Toliau pateikiami tikrųjų klasifikavimo klaidų (ang. actual error rates) vertinimo metodai, kurie bus naudojami nustatant klasifikatorių efektyvumą (ang. performance of classifiers).

#### R-metodas

R-metodas (ang. Resubstitution estimator), tai stebinių santykis mokymo imtyje iš klasės  $\Omega_l$ , tai yra klaidingo klasifikavimo santykis naudojant pasirinktą klasifikavimo taisyklę. Taikant šį metodą mokymo imtis taip pat panaudojama įvertinti diskriminantinę funkciją. Tai yra, jei  $n_1$  ir  $n_2$  yra imtys atitinkamai iš klasės  $\Omega_1$  ir  $\Omega_2$ , tuomet  $n_1$  ir  $n_2$  naudojamos apskaičiuoti diskriminantinę funkciją. Jei klaidingai suklasifikuotų imties taškų skaičius klasėse  $\Omega_1$  ir  $\Omega_2$  yra atitinkamai  $m_1$  ir  $m_2$ , tai klaidų tikimybių  $P_1$  ir  $P_2$  įvertiniai atitinkamai yra  $\frac{m_1}{n_1}$  ir  $\frac{m_2}{n_2}$ , tuomet R-metodo klaidos tikimybės įvertinimas apibrėžiamas taip (Ikechukwu 2016):

$$R = \frac{m_1 + m_2}{n_1 + n_2}.$$

#### U-metodas

Naudojant U-metodą arba „vieno neįtraukimo“ (ang. leave-one-out) vertinimo metodą visi išskyrus vieną stebinį yra naudojami apibrėžiant klasifikavimo taisyklę, o ši taisyklė yra naudojama klasifikuojant neįtrauktąjį stebinį. Ši procedūra kartojama kiekvienam stebiniui, todėl  $N$  dydžio imtyje kiekvienas stebinis yra klasifikuojamas naudojant funkciją su  $N-1$  stebinių. Bendru atveju, galima sakyti, kad tai  $m$  – lypė kryžminė patikra, naudojama klaidos tikimybės įvertinimui,  $R(cv)$ , pagal (Lachenbruch, 1967), dviejų klasių atveju  $R(cv) = \sum_{i=1}^2 \sum_{j=1}^2 n_{ij} / n_i$ . Daugelio autorių tyrimai, įskaitant ir (Efron, 1983), parodė, kad  $m$ -lypė kryžminė patikra turi didelę dispersiją. Todėl, nors  $R(cv)$  gali būti nepaslinktas įvertinimas, pasitikėjimas, su kuriuo galima tikėtis, kad  $R(cv)$  priartės prie  $R(T)$ , nėra didelis. Pagrindinis šio metodo privalumas yra tai, kad jis gauna nepaslinktą klasifikavimo vidutinės tikrosios klaidos tikimybę su mokymo imtimis  $n_1 - 1$  ir  $n_2$  iš klasių  $\Omega_l$ . Tačiau negalima daryti prielaidos, kad U-metodas turi nedidelį poslinkį lyginant su sąlygine tikrosios klaidos tikimybe. Vienas šio įvertinimo trūkumų yra tai, kad jam reikia daugiau skaičiavimų, nei R-metodui. Kitas trūkumas yra didelė dispersija. Pagrindinis dėmesys lyginant šiuos metodus skiriamas klasifikavimo vidutinės tikrosios klaidos tikimybės vertinimo poslinkiui, bet dispersija taip pat labai svarbus faktorius.

Aposteriorinės tikimybės metodas

Šis vertinimo metodas aprašytas (Moore, 1973). Kai daroma prielaida, jog apriorinės klasių tikimybės yra lygios, jei parametrai  $\theta$  yra žinomi ir turima diskriminantinė funkcija  $W$ , klaidingo klasifikavimo aposteriorinė tikimybė yra:

$$\left[ \min \{f(z, \theta_1), f(z, \theta_2)\} \right] / \left[ f(z, \theta_1) + f(z, \theta_2) \right],$$

kai  $\theta$  yra nežinomi, klaidingo klasifikavimo aposteriorinės tikimybės, paremtos diskriminantine funkcija  $W$  išraiška, įvertinamos taip:

$$\left[ \min \{f(z, \hat{\theta}_1), f(z, \hat{\theta}_2)\} \right] / \left[ f(z, \hat{\theta}_1) + f(z, \hat{\theta}_2) \right].$$

## 4 Erdvinės beta regresijos modelių pritaikymas dumblių padengimo analizei

Pagrindinis dėmesys skiriamas apibendrintiems tiesiniams ir adityviems (GAM) modeliams, su tolydžiu priklausomu kintamuoju, pasiskirsčiusiu pagal Beta dėsnį. Pradėti tirti duomenys apie Šakotųjų banguolių padengimą Baltijos jūroje, analizei taikomas erdvinės Beta regresijos modelis. Sprendžiamas modelio, besisiskiriančio erdvinės informacijos įtraukimo lygmeniu, parinkimo ir parametų vertinimo uždavinys. Atliekamas modelių palyginimas pagal AIC kriterijų. Modelių parametų vertinimo algoritmai realizuoti statistinių skaičiavimų programinėje aplinkoje R. Rezultatai pristatyti konferencijose:

Zikarienė E., Dučinskas K. Application of spatial beta regression for modelling of the algae concentration index. Spatial statistics 2019, Sitges, Ispanija 2019 m. liepos 9 - 13 d.

Zikarienė E., Dučinskas K. Implementation of generalized additive models for spatial beta regression. Computer data analysis and modeling: stochastics and data science. Minskas, Baltarusija 2019 m. rugsėjo 18-22 d.

Taip pat parašyta mokslinė publikacija:

Zikarienė E., Dučinskas K. 2019. Implementation of generalized additive models for spatial beta regression. Proceedings of the XII International Conference. Computer data analysis and modeling: stochastics and data science. p. 341-343.

## 5 Išvados

Išvestos Bajeso diskriminantinės funkcijos išraiškos eksponentinių skirtinių šeimos nariams tolydžiu atveju: Gama ir Beta skirstiniams. Išanalizuoti tikslios klasifikavimo klaidos vertinimo kriterijai, kuriuos planuojama panaudoti nustatant klasifikatorių efektyvumą (ang. performance of classifiers) ir palyginant erdviųjų duomenų modelius. Atlikta realiųjų duomenų apie Šakotųjų banguolių padengimą Baltijos jūroje, analizė, kai taikomas erdvinės Beta regresijos modelis. Sudaryti modeliai su skirtingu erdvinės informacijos įtraukimu.

## 6 Literatūra

Besag J. 1974. Spatial Interaction and the Statistical Analysis of Lattice Systems. Journal of the Royal Statistical Society, vol. 36, p. 192-236.



- Čekanavičius V., Murauskas G. 2009. Statistika ir jos taikymai III. Vilnius: TEV.
- Efron B. 1983. Estimating the error rate of the error rate of a prediction rule: improvement on cross validation. *Journal of the American Statistical Association*, vol. 78, p. 316-331.
- Fukunaga K. 1990. *Introduction to Statistical Pattern recognition*. Second edition. New York: Academic press.
- Furman E. 2008. On a multivariate gamma distribution. *Statistics & probability letters*, vol. 78, p. 2353-2360.
- Furukawa K. 2004. Development of Markov random field models based on exponential family conditional distributions. *Retrospective Theses and Dissertations*, vol. 939.
- Dučinskas K. 2009. Approximation of the expected error rate in classification of the Gaussian random field observations. *Statistics and Probability Letters*, vol. 79, p. 138-144.
- Dučinskas K., Dreičienė L. 2011. Supervised classification of the scalar Gaussian random field observations under a deterministic spatial sampling design. *Austrian Journal of Statistics*, vol. 40 (1-2), p. 25-36.
- Dučinskas K., Dreičienė L., Zikarienė E. 2015. Multiclass classification of the scalar Gaussian random field observation with known spatial correlation function. *Statistics & probability letters*, vol. 98, p. 107-114. ISSN: 0167-7152. Prieiga per internetą: <http://www.sciencedirect.com/science/article/pii/S0167715214004118#>.
- Dučinskas K., Zikarienė E. 2015. Actual error rates in classification of the T-distributed random field observation based on plug-in linear discriminant function. *Informatica*, vol. 26, no. 4, p. 557-568. ISSN: 0868-4952. Prieiga per internetą: <http://www.mii.lt/Informatica/htm/INFO1061.htm>.
- Dučinskas K., Zikarienė E., Dreičienė L. 2014. Comparison of Performances of Plug-in Spatial Classification Rules Based on Bayesian and ML Estimators. *Proceedings of the 3rd International Conference on Pattern Recognition Applications and Methods*, p. 161-166.
- Ikechukwu E. 2016. Evaluation of Error Rate Estimators in Discriminant Analysis with Multivariate Binary Variables. *American Journal of Theoretical and Applied Statistics*, vol. 5, p. 173-179.
- Lachenbruch P. A. 1967. An almost unbiased method of obtaining confidence intervals for the probability of misclassification in discriminant analysis. *Biometrics*, vol. 23, p. 639-645.
- Moore D. H. 1973, Evolution of five Discriminant procedures for binary variables. *Journal of American Statistical Association*, vol. 68, p. 399-404.
- Stabingienė L., Stabingis G., Dučinskas K. 2010. Comparison of linear discriminant functions in image classification. *Lietuvos matematikos rinkinys. LMD darbai*, vol. 51, p. 227-231. ISSN 0132-2818.
- Zikarienė E., Dučinskas K. 2019. Implementation of generalized additive models for spatial beta regression. *Proceedings of the XII International Conference*. p. 341-343.