

LATVIJAS UNIVERSITĀTE
EKSAKTO ZINĀTŅU UN TEHNOLOĢIJU FAKULTĀTE
MATEMĀTIKAS NODAĻA

GADĪJUMA MEŽA VISPĀRINĀŠANAS KĻŪDAS AUGŠĒJO ROBEŽU ANALĪZE

REFERĀTS

Autors: Kirills Bobkovs
Studenta apliecības nr.: kb25028

RĪGA 2025

SATURA RĀDĪTĀJS

1. Ievads	2
2. Gadījuma mērs un tā precizitāte	3
3. Praktiskā daļa un rezultāti	9
4. Secinājumi	14
Izmantotā literatūra un avoti	15
Pielikumi	16
A. Programmu kodi	16

1. IEVADS

Gadījuma meža (angl. *Random Forest*) klasifikācijas algoritmi ir plaši izmantoti gan datorzinātnē, gan statistikā, pateicoties to augstajai precizitātei un spējai samazināt pārmācīšanās risku, salīdzinot ar atsevišķiem lēmumu kokiem. Tomēr, lai pilnībā izprastu šo metožu darbību un to kļūdas uzvedību, nepieciešams analizēt un formulēt teorētisko pamatu vispārināšanas kļūdas (angl. *generalization error*) augšējām robežām.

Šajā darbā galvenā uzmanība tiek pievērsta gadījuma meža vispārināšanas kļūdas un tās augšējās robežas izpētei, īpašu uzmanību pievēršot diviem kritiskiem parametriem: individuālo klasifikatoru precizitātei (angl. *strength*) un to savstarpējai korelācijai. Šie parametri tiek definēti, izmantojot jaunus jēdzienus, piemēram, *marginal function* un *raw margin function*, kas ļauj kvantitatīvi novērtēt meža klasifikatoru kvalitāti un to savstarpējo atkarību.

Darba mērķis ir iepazīties ar teorētiskajām pamatformulām un praktiski ilustrēt vispārināšanas kļūdas augšējās robežas gadījuma mežos. Mērķa sasniegšanai tika izvirzīti šādi darba uzdevumi:

1. analizēt un izpētīt svarīgākos jēdzienus vispārināšanas kļūdas un tās augšējās robežas gadījuma mežos;
2. izpētīt, vai modelis pārmācās, palielinoties koku skaitam mežā;
3. izvest vispārināšanas kļūdas augšējo robežu, izmantojot divus parametrus: individuālo klasifikatoru precizitāti un to savstarpējo korelāciju.

Darba struktūra sastāv no piecām nodaļām: Ievads, Gadījuma mežs un tā precizitāte, Praktiskā daļa un rezultāti, Secinājumi un Pielikums (programmas kods).

Šajā darbā izmantotās teorētiskās formulas un izklāsti balstās uz literatūru par gadījuma mežiem ([2]). Simulācijas rezultāti ir ģenerēti ar šī darba autora *R* kodu (skat. Pielikums A).

2. GADĪJUMA MEŽS UN TĀ PRECIZITĀTE

2.1. Definīcija. Gadījuma mežs ir klasifikators, kas sastāv no koka strukturētiem klasifikatoriem $\{h(\mathbf{x}, \Theta_k), k = 1, \dots, K\}$, kur $\{\Theta_k\}$ ir neatkarīgi un vienādi sadalīti (iid) gadījuma lielumu vektori, un katrs koks piešķir vienu balsi vispopulārākajai klasei pie dotās ievades \mathbf{x} .

Tā kā katrs koks tiek izveidots, izmantojot mācīšanās datu kopu un gadījuma vektoru Θ_k , tad kopu $\{\Theta_k\}$ var interpretēt kā vienu no veidiem, kā koka konstruēšanas procesā ieviest jeb *injecēt* (angl. *inject*) nejaušību.

Piemēram, ja tiek izmantota tā sauktā iepakojšana (angl. *bagging*), tad gadījuma lieluma vektors Θ tiek ģenerēts kā skaitu sadalījums N kastēs, kas rodas, nejauši iemetot N šautriņas šajās kastēs, kur N ir piemēru skaits mācīšanās datu kopā. Citiem vārdiem sakot, tiek izmantota bootstrap metode jeb izlase ar atkārtotšanu (*sampling with replacement*), kas katram kokam piešķir indeksu kopu $\{i_1, i_2, \dots, i_N\}$, kur N ir datu punktu skaits mācīšanās kopā. Šī indeksu kopa nosaka, kuri datu punkti no mācīšanās datu kopas tiek izmantoti konkrētā koka apmācībai.

Savukārt gadījumā, kad tiek izmantota nejauša sadalījumu izvēle (*random split selection*), vektors Θ sastāv no neatkarīgiem gadījuma veseliem skaitļiem intervālā no 1 līdz M , kur M ir mainīgo skaits datu punktam, kas nosaka, kuri prediktori jeb mainīgie tiek apsvērti katrā sadalījuma solī *CART* algoritma.

Tādējādi var secināt, ka vektora Θ daba un dimensija ir tieši atkarīga no tā, kādā veidā nejaušība tiek izmantota koka konstruēšanas procesā.

Vairāk par *CART* algoritmu skatīt šajā avotā: [1].

2.2. Definīcija. Lai $h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_K(\mathbf{x})$ būtu klasifikatoru kopa, un mācīšanas izlase būtu paņemta no gadījuma lielumu vektora Y, \mathbf{X} , tad marginālo funkciju definē šādi:

$$\text{mg}(\mathbf{X}, Y) = v_Y(\mathbf{X}) - \max_{j \neq Y} v_j(\mathbf{X}),$$

kur $v_j(\mathbf{X}) = \frac{1}{K} \sum_{k=1}^K \mathbf{I}(h_k(\mathbf{X}) = j)$ apzīmē koku balsu proporciju klasei j pie ievades \mathbf{X}

Margināla funkcija parāda, cik lielā mērā vidējais balsu skaits, pie ievades \mathbf{X} un klases Y , par pareizo klasi pārsniedz vidējo balsu skaitu par jebkuru citu klasi. Jo lielāka marža, jo lielāka pārliecība par klasifikāciju.

2.3. Definīcija. Par $PE^* = P_{\mathbf{X}, Y}(\text{mg}(\mathbf{X}, Y) < 0)$ tiek saukta vispārināšanas kļūda.

Mašīnmācīšanās kontekstā vispārīnāšanas kļūda mēra, cik precīzs ir modelis uz mācīšanās laikā neredzētajiem datiem.

2.1. Teorēma. *Palielinoties koku skaitam, t.i. $K \rightarrow \infty$, PE^* visām secībām $\Theta_1, \Theta_2, \dots$ gandrīz droši (angl. *almost surely*) konverģē uz*

$$P_{\mathbf{X},Y} \left(P_{\Theta}(h(\mathbf{X},\Theta) = Y) - \max_{j \neq Y} P_{\Theta}(h(\mathbf{X},\Theta) = j) < 0 \right).$$

Pierādījums. Vispirms vēlreiz aprakstīsim vajadzīgo.

Katrs koks tiek ģenerēts ar nejaušu "recepti" Θ . Šī Θ ietver visu nejaušību. Tātad $\Theta_1, \Theta_2, \dots$ ir secība no N neatkarīgiem koku ģenerēšanas mēģinājumiem, un $h(\Theta_n, x)$ ir n -tā koka prognoze punktam x .

Indikatoru funkcija

$$I(h(\Theta_n, x) = j)$$

ir 1, ja n -tais koks prognozē klasi j , un 0 citādi. Šie indikatori ir neatkarīgi un vienādi sadalīti (iid.), jo Θ_n ir iid.

Katrs koks Θ_n dala x telpu uz lapu mezgliem (angl. *leaf nodes*), kas ir hipertaisnstūri (angl. *hyperrectangles*). Tādēļ jebkurai klasei j , kopa $S_n := \{x : h(\Theta_n, x) = j\}$ ir hipertaisnstūra reģionu apvienojums.

Tā kā jebkuram kokam ir galīgs skaits lapu mezglu, eksistē tikai galīgs dažādu hipertaisnstūra reģionu apvienojumu skaits, pie visām iespējamām Θ . Apzīmē S_1, S_2, \dots, S_K , kur $K \leq N$ un N ir koku skaits.

Tiek definēta funkcija

$$\phi(\Theta) = k, \quad \text{ja } \{x : h(\Theta, x) = j\} = S_k,$$

kas katram kokam Θ_n piešķir kādu kopu S_k klasei j .

Tālāk tiek skaitīts cik reizes katrs S_k parādās pirmajos N kokās. $N_k :=$ koku skaits no pirmajiem N , kuriem $\phi(\Theta_n) = k$, jeb

$$N_k = \sum_{n=1}^N I(\phi(\Theta_n) = k)$$

Tad

$$\frac{1}{N} \sum_{n=1}^N I(h(\Theta_n, x) = j) = \frac{1}{N} \sum_{k=1}^K N_k I(x \in S_k)$$

Tas ļauj izmantot stipro lielo skaitļu likumu uz N_k/N .

Kad $N \rightarrow \infty$, $\frac{N_k}{N} \rightarrow P(\phi(\Theta) = k)$ gandrīz droši.

Tad jebkuram x :

$$\frac{1}{N} \sum_{n=1}^N I(h(\Theta_n, x) = j) = \sum_{k=1}^K \frac{N_k}{N} I(x \in S_k) \xrightarrow{\text{g.d.}} \sum_{k=1}^K P(\phi(\Theta) = k) I(x \in S_k) = P_{\Theta}(h(\Theta, x) = j).$$

Tā kā, palielinoties koku skaitam $K \rightarrow \infty$, katra koka balsu proporcija

$$v_j(\mathbf{X}) = \frac{1}{K} \sum_{k=1}^K \mathbf{I}(h_k(\mathbf{X}) = j)$$

gandrīz droši konverģē uz atbilstošo varbūtību

$$v_j(\mathbf{X}) \xrightarrow{\text{g.d.}} P_{\Theta}(h(\mathbf{X}, \Theta) = j),$$

tad arī marginālā funkcija gandrīz droši konverģē:

$$\text{mg}(\mathbf{X}, Y) \xrightarrow{\text{g.d.}} P_{\Theta}(h(\mathbf{X}, \Theta) = Y) - \max_{j \neq Y} P_{\Theta}(h(\mathbf{X}, \Theta) = j),$$

un līdz ar to vispārīnāšanas kļūda PE^* gandrīz droši konverģē uz

$$PE^* \xrightarrow{\text{g.d.}} P_{\mathbf{X}, Y} \left(P_{\Theta}(h(\mathbf{X}, \Theta) = Y) - \max_{j \neq Y} P_{\Theta}(h(\mathbf{X}, \Theta) = j) < 0 \right).$$

□

Ņemot vērā, ka, palielinoties koku skaitam, atsevišķu koku balsu proporcijas gandrīz droši konverģē uz atbilstošajām varbūtībām, tika ieviesta marginālā funkcija, kas raksturo gadījuma meža asimptotisko (bezgalīga koku skaita) klasifikācijas uzvedību.

2.4. Definīcija. Gadījuma meža marginālā funkcija tiek definēta kā

$$\text{mr}(\mathbf{X}, Y) = P_{\Theta}(h(\mathbf{X}, \Theta) = Y) - \max_{j \neq Y} P_{\Theta}(h(\mathbf{X}, \Theta) = j).$$

2.5. Definīcija. Klasifikatoru kopas $\{h(\mathbf{x}, \Theta)\}$ stiprums tiek definēts kā

$$s = \mathbb{E}_{\mathbf{X}, Y} [\text{mr}(\mathbf{X}, Y)].$$

2.1. Apgalvojums. $PE^* \leq \text{var}(mr)/s^2$

Pierādījums. Pieņemsim, ka gadījuma meža stiprums ir $s \geq 0$, kur

$$s = \mathbb{E}_{\mathbf{X}, Y} [\text{mr}(\mathbf{X}, Y)].$$

Atcerēsimies, ka vispārīnāšanas kļūda ir definēta kā

$$PE^* = P_{\mathbf{X}, Y}(\text{mr}(\mathbf{X}, Y) < 0).$$

Ir spēkā: $\{\text{mr}(\mathbf{X}, Y) < 0\} \subseteq \{|\text{mr}(X, Y) - s| \geq s\}$.

Pārejot pie varbūtībām, seko

$$PE^* = P_{\mathbf{X}, Y}(\text{mr}(X, Y) < 0) \leq P_{\mathbf{X}, Y}(|\text{mr}(\mathbf{X}, Y) - s| \geq s).$$

Pielietojot Čebiševa nevienādību, iegūstam

$$PE^* \leq P_{\mathbf{X}, Y}(|\text{mr}(\mathbf{X}, Y) - s| \geq s) \leq \frac{\text{Var}(\text{mr}(\mathbf{X}, Y))}{s^2}. \quad (2.1)$$

□

Marginālās funkcijas dispersijas izteiksmi var uzrakstīt atklātākā formā, bet vispirms jādefinē papildu jēdzieni.

Ar \hat{j} definē:

$$\hat{j}(\mathbf{X}, Y) := \arg \max_{j \neq Y} P_{\Theta}(h(X, \Theta) = j).$$

Tad

$$\begin{aligned} \text{mr}(\mathbf{X}, Y) &= P_{\Theta}(h(\mathbf{X}, \Theta) = Y) - P_{\Theta}(h(\mathbf{X}, \Theta) = \hat{j}(\mathbf{X}, Y)) \\ &= \mathbb{E}_{\Theta}[\mathbf{I}(h(\mathbf{X}, \Theta) = Y) - \mathbf{I}(h(\mathbf{X}, \Theta) = \hat{j}(\mathbf{X}, Y))]. \end{aligned}$$

2.6. Definīcija. Neapstrādātā (angl. *raw*) marginālā funkcija tiek definēta kā

$$\text{rmg}(\Theta, \mathbf{X}, Y) = \mathbf{I}(h(\mathbf{X}, \Theta) = Y) - \mathbf{I}(h(\mathbf{X}, \Theta) = \hat{j}(\mathbf{X}, Y)).$$

Tātad $\text{mr}(\mathbf{X}, Y) = \mathbb{E}_{\Theta}[\text{rmg}(\Theta, \mathbf{X}, Y)]$ ir neapstrādātās marginālās funkcijas $\text{rmg}(\Theta, \mathbf{X}, Y)$ matemātiskā cerība attiecībā pret Θ .

2.2. Teorēma. Augšējā robeža vispārināšanas kļūdai ir

$$PE^* \leq \frac{\bar{\rho}(1 - s^2)}{s^2}.$$

Pierādījums. Vispārīgai funkcijai f ir identitāte

$$(\mathbb{E}_{\Theta} f(\Theta))^2 = \mathbb{E}_{\Theta, \Theta'}[f(\Theta)f(\Theta')], \quad (2.2)$$

kur Θ un Θ' ir i.i.d., jo

$$\mathbb{E}_{\Theta, \Theta'}[f(\Theta)f(\Theta')] = \iint f(\theta)f(\theta') dP_{\Theta}(\theta) dP_{\Theta}(\theta') = \left(\int f(\theta) dP_{\Theta}(\theta)\right)^2 = (\mathbb{E}_{\Theta}[f(\Theta)])^2.$$

Tātad identitāte turas tieši no Θ un Θ' neatkarības.

Kā iepriekš jau bija parādīts, $\text{mr}(\mathbf{X}, Y)$ var uzrakstīt kā:

$$\text{mr}(\mathbf{X}, Y) := \mathbb{E}_{\Theta}[\text{rmg}(\Theta, \mathbf{X}, Y)]. \quad (2.3)$$

Lai aprēķinātu marginālās funkcijas dispersiju attiecībā pret (X, Y) , izmantojam dispersijas definīciju:

$$\begin{aligned}\text{var}(\text{mr}) &:= \text{var}_{\mathbf{X}, Y}(\text{mr}(\mathbf{X}, Y)) \\ &= \mathbb{E}_{\mathbf{X}, Y}[\text{mr}(\mathbf{X}, Y)^2] - \left(\mathbb{E}_{\mathbf{X}, Y}\text{mr}(\mathbf{X}, Y)\right)^2.\end{aligned}\quad (2.4)$$

Sākam ar pirmo locekli no 2.4. Izmantojot 2.3, iegūstam:

$$\text{mr}(\mathbf{X}, Y)^2 = \mathbb{E}_{\Theta, \Theta'}[\text{rmg}(\Theta, \mathbf{X}, Y) \text{rmg}(\Theta', \mathbf{X}, Y)],$$

kur Θ' un Θ ir i.i.d. $((X, Y)$ šeit ir fiksēti).

Pielietojot Fubini teorēmu, samainām matemātisko cerību secību:

$$\mathbb{E}_{\mathbf{X}, Y}[\mathbb{E}_{\Theta, \Theta'}[\cdot]] = \mathbb{E}_{\Theta, \Theta'}[\mathbb{E}_{\mathbf{X}, Y}[\cdot]].$$

Tādējādi

$$\mathbb{E}_{\mathbf{X}, Y}[\text{mr}(\mathbf{X}, Y)^2] = \mathbb{E}_{\Theta, \Theta'}[\mathbb{E}_{\mathbf{X}, Y}[\text{rmg}(\Theta, \mathbf{X}, Y) \text{rmg}(\Theta', \mathbf{X}, Y)]] . \quad (2.5)$$

Tagad aplūkojam otro locekli no (2.4).

$$\begin{aligned}\mathbb{E}_{\mathbf{X}, Y}\text{mr}(\mathbf{X}, Y) &= \mathbb{E}_{\mathbf{X}, Y}\mathbb{E}_{\Theta}[\text{rmg}(\Theta, \mathbf{X}, Y)] \\ &= \mathbb{E}_{\Theta}\mathbb{E}_{\mathbf{X}, Y}[\text{rmg}(\Theta, \mathbf{X}, Y)],\end{aligned}$$

kur pēdējā vienādība atkal izriet no neatkarības un Fubini teorēmas.

Pacelot kvadrātā un ieviešot neatkarīgu kopiju Θ' (pielietojot 2.2), iegūstam

$$\left(\mathbb{E}_{\mathbf{X}, Y}\text{mr}(\mathbf{X}, Y)\right)^2 = \mathbb{E}_{\Theta, \Theta'}[\mathbb{E}_{\mathbf{X}, Y}\text{rmg}(\Theta, \mathbf{X}, Y) \mathbb{E}_{\mathbf{X}, Y}\text{rmg}(\Theta', \mathbf{X}, Y)] . \quad (2.6)$$

Ievietojot izteiksmes (2.5) un (2.6) dispersijas formulā (2.4), iegūstam

$$\begin{aligned}\text{var}(\text{mr}) &= \mathbb{E}_{\Theta, \Theta'}\left[\mathbb{E}_{\mathbf{X}, Y}[\text{rmg}(\Theta, \mathbf{X}, Y) \text{rmg}(\Theta', \mathbf{X}, Y)]\right. \\ &\quad \left.- \mathbb{E}_{\mathbf{X}, Y}\text{rmg}(\Theta, \mathbf{X}, Y) \mathbb{E}_{\mathbf{X}, Y}\text{rmg}(\Theta', \mathbf{X}, Y)\right].\end{aligned}$$

Secinām, ka

$$\text{var}(\text{mr}) = \mathbb{E}_{\Theta, \Theta'}[\text{cov}_{\mathbf{X}, Y}(\text{rmg}(\Theta, \mathbf{X}, Y), \text{rmg}(\Theta', \mathbf{X}, Y))] . \quad (2.7)$$

Izmantojot šo izteiksmi, iegūstam

$$\text{var}(\text{mr}) = \mathbb{E}_{\Theta, \Theta'}[\text{cov}_{\mathbf{X}, Y}(\text{rmg}(\Theta, \mathbf{X}, Y), \text{rmg}(\Theta', \mathbf{X}, Y))] \quad (2.8)$$

$$= \mathbb{E}_{\Theta, \Theta'}[\rho(\Theta, \Theta') \text{sd}(\Theta) \text{sd}(\Theta')], \quad (2.9)$$

kur $\rho(\Theta, \Theta')$ ir neapstrādātās marginālās funkcijas $\text{rmg}(\Theta, \mathbf{X}, Y)$ un $\text{rmg}(\Theta', \mathbf{X}, Y)$ korelācija, kad Θ, Θ' tiek turēti fiksēti, un $\text{sd}(\Theta)$ ir $\text{rmg}(\Theta, \mathbf{X}, Y)$ standarta novirze, turot Θ fiksētu.

Tad iegūstam

$$\text{var}(\text{mr}) = \frac{\mathbb{E}_{\Theta, \Theta'} [\rho(\Theta, \Theta') \text{sd}(\Theta) \text{sd}(\Theta')] \mathbb{E}_{\Theta, \Theta'} [\text{sd}(\Theta) \text{sd}(\Theta')]}{\mathbb{E}_{\Theta, \Theta'} [\text{sd}(\Theta) \text{sd}(\Theta')]} = \bar{\rho} (\mathbb{E}_{\Theta} \text{sd}(\Theta))^2 \leq \bar{\rho} \mathbb{E}_{\Theta} \text{var}(\Theta), \quad (2.10)$$

kur $\bar{\rho}$ ir vidējā korelācija, t.i.,

$$\bar{\rho} := \frac{\mathbb{E}_{\Theta, \Theta'} [\rho(\Theta, \Theta') \text{sd}(\Theta) \text{sd}(\Theta')]}{\mathbb{E}_{\Theta, \Theta'} [\text{sd}(\Theta) \text{sd}(\Theta')]}.$$

$$\text{Tā kā } \text{var}(\Theta) = \text{var}_{\mathbf{X}, Y}(\text{rmg}(\Theta, \mathbf{X}, Y)) = \mathbb{E}_{\mathbf{X}, Y}[\text{rmg}(\Theta, \mathbf{X}, Y)^2] - (\mathbb{E}_{\mathbf{X}, Y}[\text{rmg}(\Theta, \mathbf{X}, Y)])^2,$$

un $\text{rmg}(\Theta, \mathbf{X}, Y)$ ir ierobēžota, t.i., $\text{rmg}(\Theta, \mathbf{X}, Y)^2 \leq 1$, tad iegūstam

$$\mathbb{E}_{\mathbf{X}, Y}[\text{rmg}(\Theta, \mathbf{X}, Y)^2] \leq 1.$$

$$\text{Tādējādi } \text{var}(\Theta) \leq 1 - (\mathbb{E}_{\mathbf{X}, Y}[\text{rmg}(\Theta, \mathbf{X}, Y)])^2.$$

Nemot cerību pēc Θ :

$$\mathbb{E}_{\Theta} \text{var}(\Theta) \leq \mathbb{E}_{\Theta} [1 - (\mathbb{E}_{\mathbf{X}, Y}[\text{rmg}(\Theta, \mathbf{X}, Y)])^2] = 1 - \mathbb{E}_{\Theta} (\mathbb{E}_{\mathbf{X}, Y}[\text{rmg}(\Theta, \mathbf{X}, Y)])^2. \quad (2.11)$$

Izmantojot Jensena nevienādību (Jensen's inequality) izliektai funkcijai x^2 , iegūstam

$$\mathbb{E}_{\Theta} (\mathbb{E}_{\mathbf{X}, Y}[\text{rmg}(\Theta, \mathbf{X}, Y)])^2 \geq (\mathbb{E}_{\Theta} \mathbb{E}_{\mathbf{X}, Y}[\text{rmg}(\Theta, \mathbf{X}, Y)])^2 = s^2, \quad (2.12)$$

Apvienojot 2.11 un 2.12, iegūstam galīgo nevienādību:

$$\mathbb{E}_{\Theta} \text{var}(\Theta) \leq 1 - s^2. \quad (2.13)$$

Apvienojot 2.1, 2.10 un 2.13, iegūstam augšējo robežu vispārināšanas kļūdai:

$$PE^* \leq \frac{\bar{\rho}(1 - s^2)}{s^2}. \quad (2.14)$$

□

2.1. Piezīme. Nevienādība 2.10

$$(\mathbb{E}_{\Theta} \text{sd}(\Theta))^2 \leq \mathbb{E}_{\Theta} \text{var}(\Theta)$$

tiek iegūta, izmantojot Koši–Švarca nevienādību:

$$(\mathbb{E}A)^2 \leq \mathbb{E}(A^2),$$

kur $A = \text{sd}(\Theta) = \sqrt{\text{var}(\Theta)}$.

3. PRAKTISKĀ DAĻA UN REZULTĀTI

Šajā sadaļā tiek analizēta gadījuma meža vispārināšanas kļūdas augšējās robežas uzvedība, īpaši tās stabilizēšana, palielinoties koku skaitam. Iegūta robeža 2.14 palīdz izprast divus galvenos faktorus, kas ietekmē nejaušā meža kļūdu: *atsevišķo koku stiprumu* un *to savstarpējo korelāciju* neapstrādātās marginālās funkcijas izteiksmē. Attiecība $\bar{\rho}/s^2$ (vidējā korelācija dalīta ar stipruma kvadrātu) darbojas kā indikators — jo mazāka šī attiecība, jo labāka meža vispārināšanas spēja.

Lai ilustrētu iegūto teorētisko rezultātu, tika veikta simulācija R vidē. Tika izmantota klasiskā *Iris* [3] datu kopa, kas sastāv no četrām pazīmēm un trim klašu kategorijām. Visas teorētiskajās formulās sastopamās matemātiskās cerības tika aproksimētas ar empīrisko vidējo vērtību, izmantojot galīgu datu punktu un gadījuma koku realizāciju skaitu.

Simulācijā tika ģenerēti gadījuma meži ar dažādu koku skaitu, un katram koku skaitam tika veiktas vairākas inicializācijas, izmantojot atšķirīgas `set.seed` vērtības, lai novērtētu vispārināšanas kļūdas augšējās robežas stabilitāti attiecībā pret gadījuma meža konstrukciju.

Simulācijas rezultātā iegūtie grafiki parāda, kā pārsvarā augšējā robeža stabilizējas ar palielinātu koku skaitu. Tomēr nevar teikt, ka tas izpildās vienmēr, jo daži grafiki nedod drošu iespaidu par augšējās robežas konvergenci. Svarīgi atzīmēt, ka teorētiskā konverģence uz kādu vērtību ir pierādīta tikai PE^* vērtībai, nevis tās augšējām robežām. Tāpēc mūsu brīva (angl. *loose*) augšējā robeža var nerādīt PE^* īsto konverģences dinamiku. Rezultāti ir attēloti 1, 2, 3, 4.

Papildus pārbaude teorēmai (2.1) tika veikta, simulējot gadījuma meža veikspēju sirds slimību klasifikācijas uzdevumā ([4]), izmantojot *Python*. Datu kopā ietilpst 14 klīniskie un demogrāfiskie prediktori, kas raksturo pacientus, piemēram, vecumu, dzimumu, krūšu sāpju tipu, asinsspiedienu, holesterīna līmeni, elektrokardiogrammas rādītājus un maksimālo sirdsdarbības frekvenci, kā arī citi sirds veselību raksturojoši parametri. Mērķis bija prognozēt mainīgo `num`, kas norāda sirds slimības esamību vai neesamību.

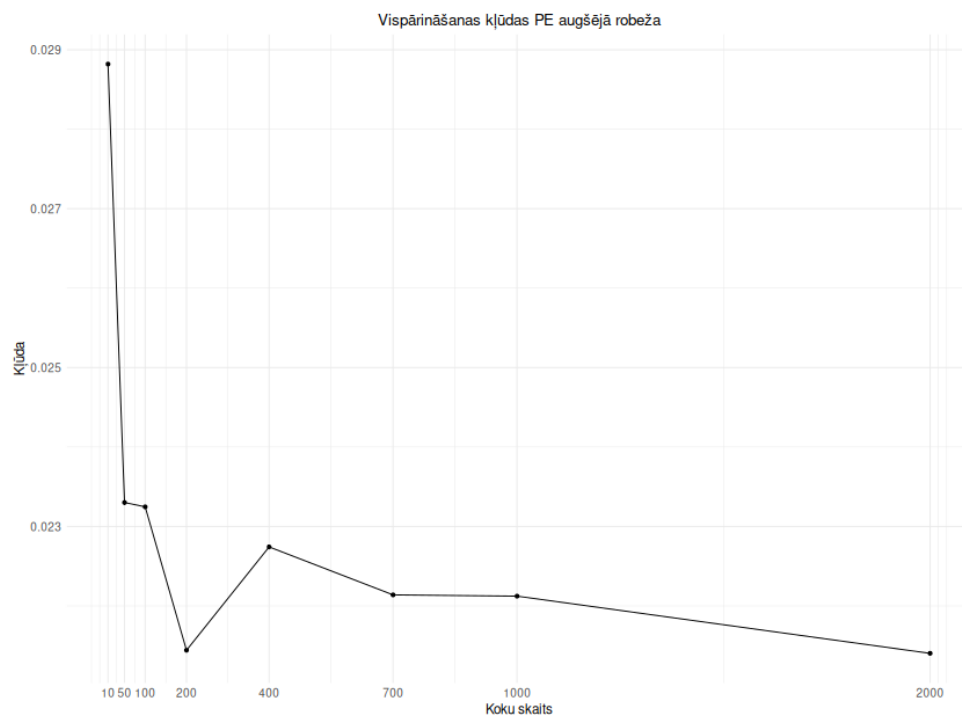
Simulācijā tika ģenerēti gadījuma meži ar dažādu koku skaitu, saglabājot citus hiperparametrus nemainīgus. Katram koku skaita parametram tika veiktas 20 inicializācijas ar dažādiem `seed` iestatījumiem, aprēķinot AUC vidējo vērtību un standartnovirzi. Rezultāti, parādīti tabulā 1, empiriski ilustrē, ka AUC vērtības pieaug, palielinot koku skaitu, līdz sasniedz noteiktu K vērtību, pēc kuras novērojams plateau efekts — modelis vairs neuzlabojas, bet aprēķina laiks pieaug.

Koku skaits	AUC vidējā vērtība	std(AUC)	Vidējais aprēķina laiks, s
1	0.749633	0.029051	0.054637
5	0.869783	0.013673	0.098653
10	0.890504	0.009983	0.139046
25	0.903147	0.004656	0.179914
50	0.907483	0.004686	0.324881
100	0.910726	0.003512	0.456434
200	0.913392	0.002666	1.007332
500	0.914867	0.001961	2.541950
1000	0.914772	0.001178	4.329879
1500	0.915053	0.000732	6.402032

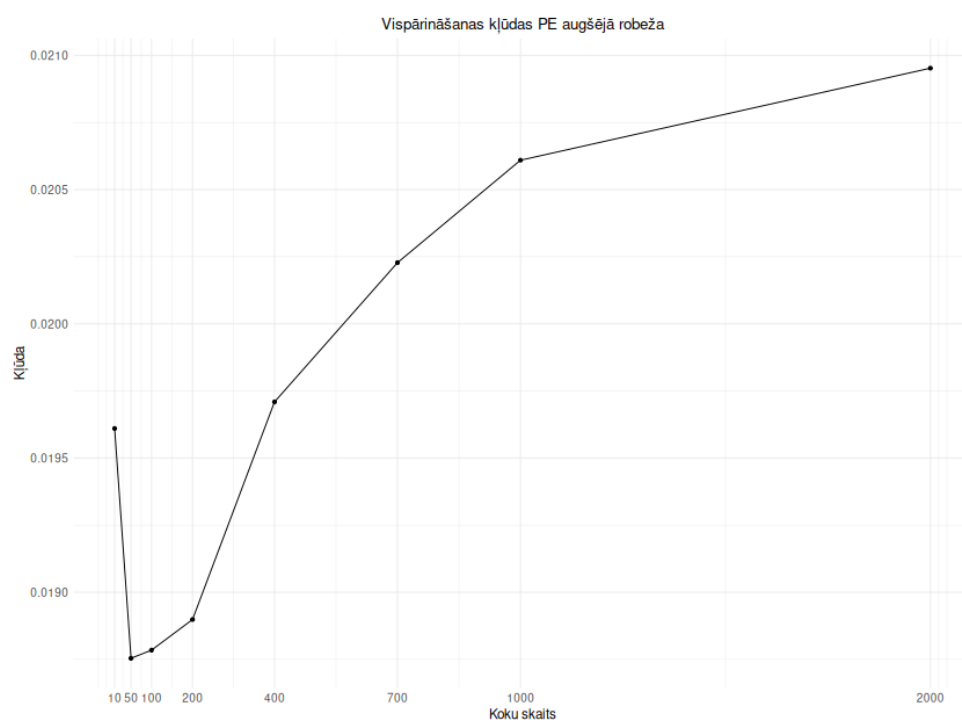
1. tabula

**Gadījuma meža AUC veikspēja un aprēķina laiks dažādam koku skaitam
(atkārtojumu skaits: 20).**

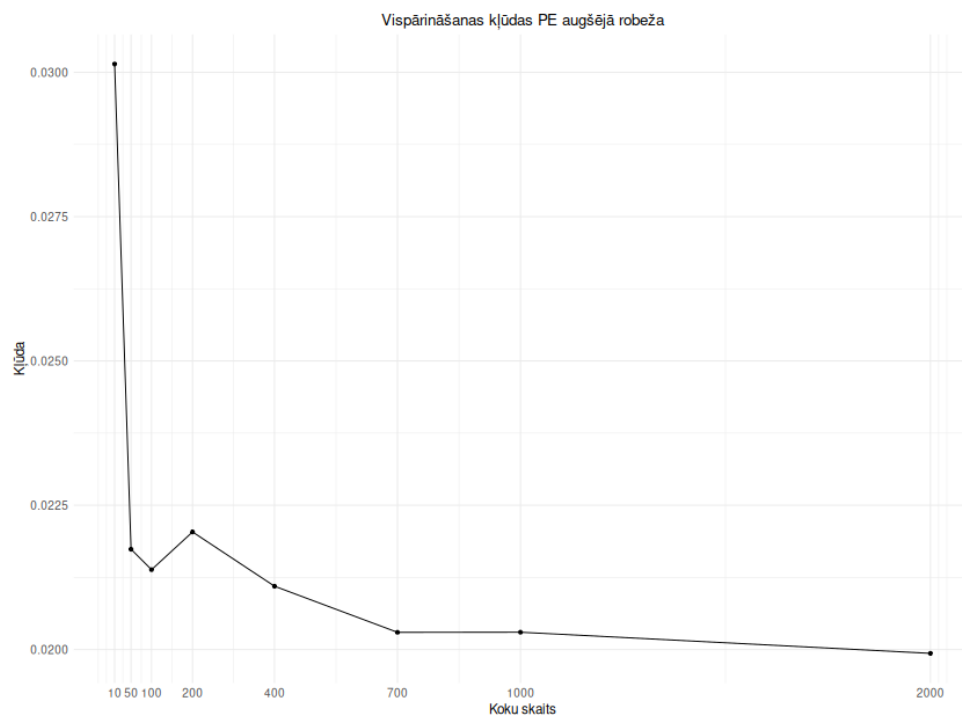
Rezultāti tiek vizualizēti divos grafikos: *AUC konverģence un standartnovirzes izmaiņas* (5) parāda, kā AUC vidējās vērtības un standartnovirze stabilizējas pie palielināta koku skaita, savukārt *Gadījuma meža veikspēja un resursu patēriņš atkarībā no koku skaita* (6) ilustrē, ka palielinoties koku skaitam, izpildes laiks pieaug, saglabājot citus parametrus nemainīgus. Šie rezultāti empīriski apstiprina teorētiskās prognozes par vispārīnāšanas kļūdas konverģenci.



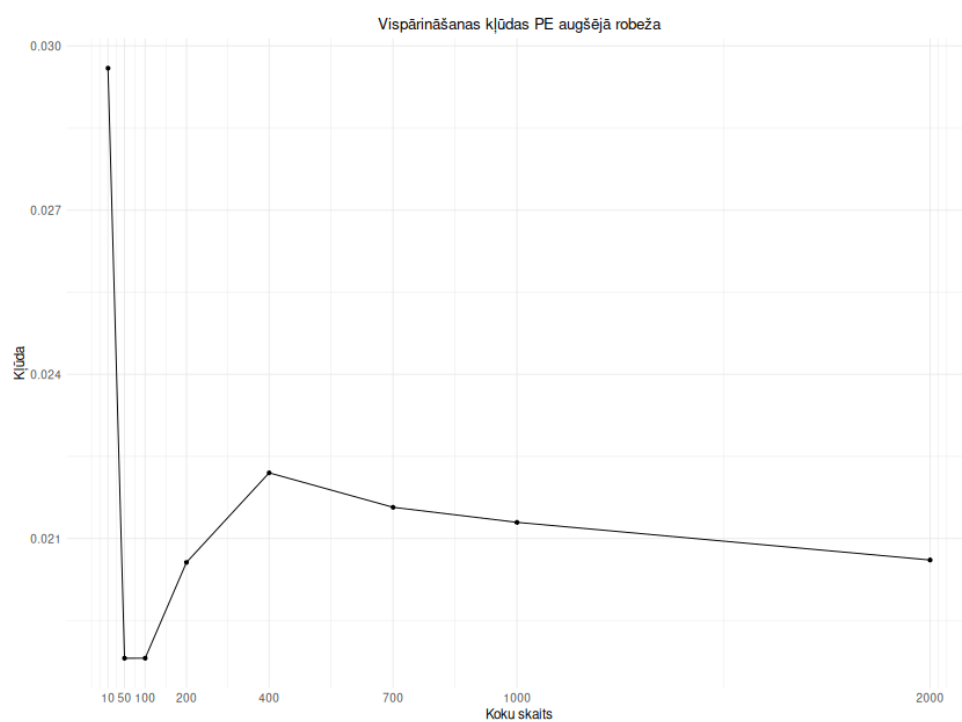
1. att. Simulācijas rezultāts: Kļūdas PE^* augšējās robežas atkarība no koku skaita (nejaušā sēkla: 42).



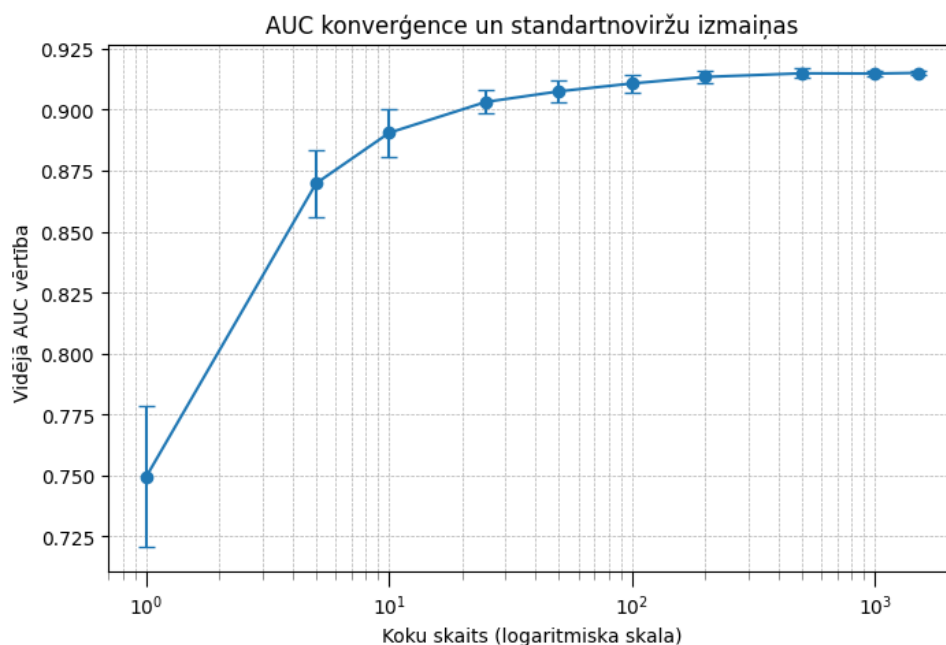
2. att. Simulācijas rezultāts: Kļūdas PE^* augšējās robežas atkarība no koku skaita (nejaušā sēkla: 1).



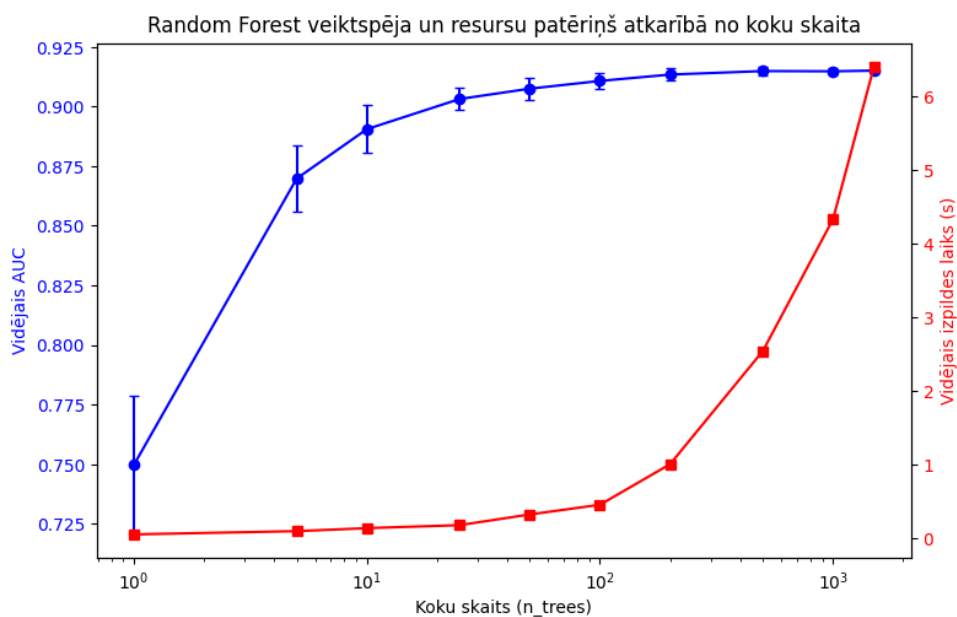
3. att. Simulācijas rezultāts: Kļūdas PE^* augšējās robežas atkarība no koku skaita (nejaušā sēkla: 666).



4. att. Simulācijas rezultāts: Kļūdas PE^* augšējās robežas atkarība no koku skaita (nejaušā sēkla: 69).



5. att. Simulācijas rezultāts: Vidējā AUC vērtība un standartnovirze atkarībā no koku skaita (nejaušā sēkla: 42).



6. att. Simulācijas rezultāts: Vidējā AUC vērtība, standartnovirze un vidējais izpildes laiks (s) atkarībā no koku skaita (nejaušā sēkla: 42).

4. SECINĀJUMI

1. Tika analizēti un izpētīti svarīgākie jēdzieni, kas saistīti ar vispārināšanas kļūdu un tās augšējo robežu;
2. Tika parādīts, ka kopējais gadījuma meža algoritms ir salīdzinoši stabils un neuzrāda pārmācīšanās pazīmes, palielinoties koku skaitam;
3. Tika izvestas formulas vispārināšanas kļūdas augšējai robežai, kurās iekļauti divi kritiski parametri: individuālo klasifikatoru precizitāti un to savstarpējo korelāciju.

Lai gan iegūtā augšējā robeža, visticamāk, ir konservatīva (*loose*), tā sniedz informatīvu funkciju gadījuma mežiem.

Praktiskajā daļā tika izveidoti simulācijas grafiki, kas ilustrē augšējās robežu uzvedību un modeļa veikspēju, mainoties koku skaitam mežā, ļaujot vizuāli novērtēt teorētisko prognožu atbilstību empīriskajiem datiem. Darba rezultāti liecina, ka aprēķinātā augšējā robeža ne vienmēr stabilizējas, liecinot, ka robeža būtu jāprecizē un jāaspiež tuvāk patiesajai PE^* vērtībai. Korelācijas starp klasifikatoriem un individuālais stiprums ievērojami ietekmē augšējās robežas lielumu. Savukārt rezultāti, kas iegūti sirds slimību klasifikācijas uzdevumā, parāda, ka ROC, kalpojot kā PE^* *proksi* jeb aproksimācija, konverģē uz noteikto vērtību, apstiprinot teorētiskās prognozes.

IZMANTOTĀ LITERATŪRA UN AVOTI

- [1] Timothy C. Au. Random forests, decision trees, and categorical predictors: The “absent levels” problem. *Journal of Machine Learning Research*, 19:1–30, 2018. Pieejams: <https://jmlr.org/papers/v19/16-474.html>. Skatīts: 2026-01-05.
- [2] Leo Breiman. Random forests. *Machine Learning*, 45(1), 2001. Pieejams: <https://doi.org/10.1023/A:1010933404324>. Skatīts: 2026-01-05.
- [3] R. A. Fisher. Iris dataset. UCI Machine Learning Repository, 1936. Pieejams: <https://archive.ics.uci.edu/ml/datasets/iris>. Skatīts: 2026-01-05.
- [4] Redwan Karim and M. S. Ony. Heart disease dataset. <https://www.kaggle.com/datasets/redwankarimsony/heart-disease-data>, 2021. Skatīts: 2026-01-05.

PIELIKUMI

A. PROGRAMMU KODI

```
1 library(randomForest)
2 library(dplyr)
3
4 data(iris)
5
6 X <- iris[, 1:4]
7 Y <- iris$Species
8
9 compute_rmg <- function(rf, X, Y) {
10   pred_trees <- predict(rf, X, predict.all = TRUE)$individual
11   n_obs <- nrow(pred_trees)
12   n_trees <- ncol(pred_trees)
13   rmg <- matrix(0, nrow=n_obs, ncol=n_trees)
14   for(i in 1:n_trees) {
15     rmg[,i] <- ifelse(pred_trees[,i] == Y, 1, -1)
16   }
17   return(rmg)
18 }
19
20 compute_PE_margin <- function(rmg) {
21   margins <- rowMeans(rmg)
22   mean(margins < 0)
23 }
24
25 compute_strength_rho <- function(rmg) {
26   # Strength: mean margin across all data points
27   margins <- rowMeans(rmg)
28   s <- mean(margins)
29
30   # Standard deviation per tree
31   sd_trees <- apply(rmg, 2, sd)
32
33   # Weighted numerator and denominator for average correlation
34   n_trees <- ncol(rmg)
35   num <- 0
36   den <- 0
```

```

37   for(i in 1:n_trees) {
38     for(j in 1:n_trees) {
39       rho_ij <- cor(rmg[,i], rmg[,j])
40       if(is.na(rho_ij)) rho_ij <- 0 # Treat undefined correlation as 0
41       num <- num + rho_ij * sd_trees[i] * sd_trees[j]
42       den <- den + sd_trees[i] * sd_trees[j]
43     }
44   }
45   rho_bar <- num / den
46   return(list(s = s, rho_bar = rho_bar))
47 }
48
49 compute_PE_bound <- function(rf, X, Y) {
50   rmg <- compute_rmg(rf, X, Y)
51   stats <- compute_strength_rho(rmg)
52   s <- stats$s
53   rho_bar <- stats$rho_bar
54   #PE_bound <- ifelse(s>0,(rho_bar * (1 - s^2)) / (s^2),NA)
55   eps <- 1e-10
56   PE_bound <- (rho_bar * (1 - s^2)) / (s^2 + eps)
57   return(list(s = s, rho_bar = rho_bar, PE_bound = PE_bound))
58 }
59
60
61 # Test different numbers of trees
62 tree_list <- c(10, 50, 100, 200, 400,700,1000,2000)
63 results <- data.frame(ntree=integer(), s=numeric(), rho_bar=numeric(),
64                       PE_bound=numeric(), stringsAsFactors = FALSE)
65 seed=1
66 for(nt in tree_list) {
67   set.seed(seed)
68   rf <- randomForest(x=X, y=Y, ntree=nt)
69   res <- compute_PE_bound(rf, X, Y)
70   results <- rbind(results, data.frame(ntree=nt, s=res$s, rho_bar=res$rho_
71   bar, PE_bound=res$PE_bound))
72 }
73 print(results)
74
75 library(ggplot2)

```

```

76 ggplot(results , aes(x = ntree)) +
77   geom_line(aes(y = PE_bound)) +
78   geom_point(aes(y = PE_bound)) +
79   scale_x_continuous(breaks = tree_list) +
80   labs(
81     title = "āšVisprinanas ļūkdas PE šēāaugj žrobea",
82     x = "Koku skaits",
83     y = "ļūKda",
84     color = ""
85   ) +
86   theme_minimal(base_size = 6) +
87   theme(
88     plot.title = element_text(size = 8, hjust = 0.5),
89     axis.title = element_text(size = 7),
90     axis.text = element_text(size = 6),
91   )
92 a

```

```

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import time

from sklearn.model_selection import train_test_split
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline

from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import roc_auc_score
from sklearn.preprocessing import OneHotEncoder
from sklearn.impute import SimpleImputer
from matplotlib.ticker import LogLocator
import warnings
warnings.filterwarnings("ignore")

```

```

# Kods ir replicējams
RANDOM_SEED = 42
np.random.seed(RANDOM_SEED)

df = pd.read_csv("heart_disease_uci.csv")

TARGET = "num"
df[TARGET] = df[TARGET].transform(lambda x: 1 if x > 0 else 0)
print(df[TARGET].unique())
X = df.drop(columns=[TARGET, "id"])
y = df[TARGET]

# Identificējam kategoriālos un skaitliskos (nepārtrauktos) mainīgus
numeric_cols = X.select_dtypes(include=["int64", "float64"]).columns
categorical_cols = X.select_dtypes(include=["object"]).columns
categorical_transformer = Pipeline(
    steps=[
        ("imputer", SimpleImputer(strategy="most_frequent")),
        ("onehot", OneHotEncoder(handle_unknown="ignore"))
    ]
)
numeric_transformer = Pipeline(
    steps= [
        ("imputer", SimpleImputer(strategy="median"))
    ]
)
preprocessor = ColumnTransformer(
    transformers=[
        ("num", numeric_transformer, numeric_cols),
        ("cat", categorical_transformer, categorical_cols)
    ]
)

X_train, X_test, y_train, y_test = train_test_split(

```

```

X, y,
test_size=0.3,
random_state=RANDOM_SEED,
stratify=y # svarīgi, lai testa un mācīšanās datos būtu vienāda "1" un "0" pro
)

def train_rf_auc(n_trees:int, X_train, X_test, y_train, y_test, seed=42):
    """
    Trenē gadījumu mežu ar n_trees kokiem
    un atdot testa AUC.
    """
    model = RandomForestClassifier(
        n_estimators=n_trees,
        max_depth=None, # vērtība pēc noklusējuma
        min_samples_split=2, # default vērtība
        min_samples_leaf = 1, # default vērtība
        max_features="sqrt", # default vērtība
        bootstrap=True, # default
        random_state=seed,
        n_jobs=-1
    )

    pipe = Pipeline(
        steps=[
            ("preprocess", preprocessor),
            ("rfc", model)
        ]
    )

    pipe.fit(X_train, y_train)

    y_prob = pipe.predict_proba(X_test)[:, 1] ## paņemam varbūtības
    auc = roc_auc_score(y_test, y_prob)

```

```

    return auc

tree_counts = [1, 5, 10, 25, 50, 100, 200, 500, 1000, 1500]
auc_scores = []
for k in tree_counts:
    auc = train_rf_auc(
        n_trees=k,
        X_train=X_train,
        X_test=X_test,
        y_train=y_train,
        y_test=y_test,
        seed=RANDOM_SEED
    )
    auc_scores.append(auc)

plt.figure(figsize=(8, 5))
plt.plot(tree_counts, auc_scores, marker="o")
#plt.xscale("log")
plt.xlabel("Koku skaits")
plt.ylabel("AUC")
plt.title("Gadījuma meža AUC kā funkcija no koku skaita")
plt.grid(True)
plt.show()

def repeated_auc(n_trees, n_repeats):
    """
    Funkcija, kas simulē AUC vidējo vērtību un standartnovirzi
    katram koku skaita parametram.
    """
    aucs = []
    times = []
    for i in range(n_repeats):
        start_time = time.perf_counter()
        aucs.append(

```

```

        train_rf_auc(
            n_trees,
            X_train, X_test,
            y_train, y_test,
            seed=RANDOM_SEED + i
        )
    )

    total_time = time.perf_counter()-start_time
    times.append(total_time)

return np.mean(aucs), np.std(aucs), np.mean(times)

means, stds, mean_times = [], [], []
n_repeats=20
for k in tree_counts:
    mean_auc, std_auc, mean_time = repeated_auc(k, n_repeats=n_repeats)
    means.append(mean_auc)
    stds.append(std_auc)
    mean_times.append(mean_time)

results= pd.DataFrame({"Koku skaits": tree_counts,
    f"AUC vidēja vērtība (atkartojumu skaits: {n_repeats})":means,
    f"std(AUC) (atk. skaits: {n_repeats})":stds,
    f"Vidējais aprēķina laiks, s (atk. skaits: {n_repeats})":mean_times})

results

plt.figure(figsize=(8, 5))
plt.errorbar(tree_counts, means, yerr=stds, marker="o", capsize=4)
plt.xscale("log")
plt.gca().xaxis.set_minor_locator(LogLocator(base=10.0, subs=range(2,10)))
plt.gca().xaxis.set_major_locator(LogLocator(base=10.0))
plt.grid(True, which="both", linestyle="--", linewidth=0.5)

plt.tick_params(axis='x', which='minor', length=4, color='gray')
plt.tick_params(axis='x', which='major', length=8, color='black')

```

```

plt.xlabel("Koku skaits (logaritmiska skala)")
plt.ylabel("Vidējā AUC vērtība ")
plt.title("AUC konverģence un standartnoviržu izmaiņas")
plt.show()

fig, ax1 = plt.subplots(figsize=(8,5))
color = 'blue'
ax1.set_xlabel('Koku skaits (n_trees)')
ax1.set_ylabel('Vidējais AUC', color=color)
ax1.errorbar(tree_counts, means, yerr=stds, color=color, fmt='-o', capsize=3, label=
ax1.tick_params(axis='y', labelcolor=color)
ax1.set_xscale("log")

ax2 = ax1.twinx() # otra y ass laika vizualizācijai
color = 'red'
ax2.set_ylabel('Vidējais izpildes laiks (s)', color=color)
ax2.plot(tree_counts, mean_times, '-s', color=color, label='Vidējais laiks')
ax2.tick_params(axis='y', labelcolor=color)

fig.tight_layout()
plt.title('Gadījuma meža veikspēja un resursu patēriņš atkarībā no koku skaita')
plt.show()

```