

KONTEKSTNO-NEODVISNE GRAMATIKE ZA KODIRANJE IN STISKANJE PODATKOV

JANEZ PODLOGAR

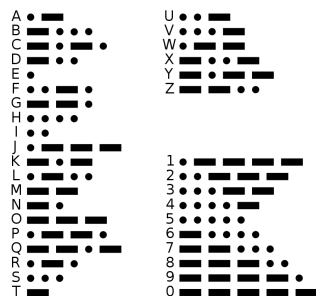
POVZETEK. V delu podamo definicije kodiranja, dekodiranja in stiskanja podatkov ter predstavimo primere, ki motivirajo stiskanje podatkov s kontekstno-neodvisnimi gramatikami.

1. KODIRANJE PODATKOV

Zapis informacije v neki obliki ni primeren za vsakršno rabo. Besedilo, zapisano z pismenkami, je neberljivo za slepe osebe, saj je komunikacijski kanal v tem primeru vid. Prav tako pisanega besedila v pravotni obliki ni mogoče poslati s telegrafom. V tem primeru je komunikacijski kanal žica in pismenke se po njej ne morejo sprehoditi. V obeh primerih je informacija, ki bi jo radi prenesli v neprimerni obliki. V prvem primeru je potrebno besedilo zapisati z Braillovo pisavo. V drugem primeru pa je potrebno besedilo pretvoriti v električni signal. Spreminjanje zapisa sporočila imenujemo *kodiranje*, sistemu pravil, po katerem se kodiranje opravi, pa *kod*.

Primer 1.1. *Morsejeva abeceda* je kodiranje črk, števil in ločil s pomočjo zaporedja kratkih in dolgih signalov:

- Dolžina kratkega signala je ena enota.
- Dolgi signal je trikrat daljši od kratkega signala.
- Razmik med signali znotraj črke je tišina dolžine kratkega signala.
- Razmik med črkami je tišina dolga tri kratke signale oz. en dolgi signal.
- Presledek med besedami je tišina dolga sedmih kratkih signalov.



SLIKA 1. Mednarodna Morsejeva abeceda

Prvotni namen Morsejeve abecede je komunikacija preko telegrama, saj komunikacijski kanal dovoljuje le električne signale in tišino med njimi. Kodiranje črk je takšno, da imajo črke z višjo frekvenco (v angleškem jeziku) krajši zapis. Tako se dolžina kodiranega sporočila skrajša in posledično tudi čas prenosa.

◇

Definicija 1.2. *Abeceda* je končna neprazna množica Σ . Elementom abecede pravimo *črke*. *Množica vseh končnih nizov abecede* Σ je

$$\Sigma^* = \{a_1 a_2 a_3 \cdots a_n \mid n \in \mathbb{N}_0 \wedge \forall i : a_i \in \Sigma\},$$

kjer za $n = 0$ dobimo prazen niz, ki ga označimo z ε . *Dolžino niza* w označimo z $|w|$ in je enaka številu črk v nizu $w \in \Sigma^*$. *Jezik na abecedi* Σ je poljubna podmnožica množice Σ^* .

Opomba 1.3. *Kleenejeva zvezdica* ali *Kleenejevo zaprtje* je operacija, ki abecedi Σ priredi najmanjšo nadmnožico Σ^* , ki vsebuje *prazen niz* ε in je zaprta za konkatenacijo oziroma veriženje. Z drugimi besedami, Σ^* je množica vseh končnih nizov, ki jih lahko generiramo z veriženjem črk abecede Σ .

Za abecedo Σ definirajmo

$$\Sigma^0 = \{\varepsilon\}, \quad \Sigma^1 = \Sigma$$

ter za vsak $i > 0$ rekurzivno

$$\Sigma^{i+1} = \{wa \mid w \in \Sigma^i \text{ in } a \in \Sigma\}.$$

Potem je Kleenejeva zvezdica na Σ enaka

$$\Sigma^* = \bigcup_{i \geq 0} \Sigma^i.$$

Primer 1.4. Naj bo $\Sigma = \{a, b, c\}$ abeceda, potem je

$$ab \in \Sigma^*, \quad ccc \in \Sigma^*, \quad cababcccababcccab \in \Sigma^*.$$

◇

Definicija 1.5. *Kodiranje nizov abecede* Σ je injektivna funkcija $\kappa: \Sigma^* \rightarrow \Sigma_c^*$, kjer je Σ_c *kodirna abeceda* in $\kappa(w)$ imenujemo *koda niza* w . *Dokodiranje kodiranja* κ je funkcija $\kappa^{-1}: C \subseteq \Sigma_c^* \rightarrow \Sigma^*$, da velja

$$\forall w \in \Sigma^*: \kappa^{-1}(\kappa(w)) = w.$$

Opomba 1.6. Zožitev kodomene kodirne funkcije κ na $C \subseteq \Sigma_c^*$ je bijektivna funkcija.

Primer 1.7. Formalizirajmo Morsejevo abecedo iz Primera 1.1. Abecedi sta

$$\Sigma = \{A, B, \dots, Z\} \cup \{0, 1, \dots, 9\} \cup \{_ \}, \quad \Sigma_c = \{., -, _ \}.$$

Definirajmo kodirno funkcijo črk abecede $\kappa: \Sigma \rightarrow \Sigma_c^*$, ki vsaki črki iz abecede Σ priredi niz črk kodirne abecede Σ_c . Vrednosti funkcije κ so določene s tabelo iz Slike 1, dodatno presledek med besedami $_$ kodiramo v šest kratkih enot tišine

$$\kappa(_) = ______.$$

Za niz $w = a_1 a_2 \dots a_n \in \Sigma^*$ definiramo kodirno funkcijo K po črkah

$$K(w) = \kappa(a_1) _ \kappa(a_2) _ \cdots \kappa(a_n).$$

Poglejmo si dva primera kodiranja in v Morsejevi abecedi

$$K(\text{SOS}) = \cdots \quad - - - \quad \cdots,$$

$$K(\text{AD HOC}) = . - \quad - \cdots \quad \cdots \cdots \quad - - - \quad - \cdots \cdots.$$

Recimo, da smo dobili sporočilo, a se je pošiljatelj zmotil in je namesto kode $- - - - \cdot - \cdots$, ki bi se dekodirala v

$$K^{-1}(- - - - \cdot - \cdots) = \text{QED},$$

poslali kodo

$$- - - - - \cdot - \cdots.$$

Sporočila ne znamo dekodirati, saj se ne nahaja v domeni C dekodirne funkcije K^{-1} . Natančneje ne obstaja dekodiranje signala $\kappa^{-1}(- - \cdot - -)$.

◇

2. STISKANJE PODATKOV

Eden izmed namen kodiranja je tudi doseči čim večjo ekonomičnost zapisa. Želimo, da bi bilo naše sporočilo čim krajše. Kodiranje, ki skrajša zapis podatkov, imenujemo *stiskanje podatkov*.

Definicija 2.1. *Stiskanje* je kodiranje K za katerega velja

$$\exists n \in \mathbb{N} \forall w \in \Sigma^* : |w| \geq n \implies |K(w)| \ll |w|.$$

Opomba 2.2. Ločimo *kodiranje brez izgube*, kjer velja

$$\forall w \in \Sigma^* : \kappa^{-1}(\kappa(w)) = w$$

in *kodiranje z izgubo*, kjer kodiranje ni reverzibilen proces in v grobem velja

$$\forall w \in \Sigma^* : \kappa^{-1}(\kappa(w)) \approx w.$$

Primer 2.3. Za abecedo vzemimo $\Sigma = \{a, b, c\}$ in pogledjmo niz

$$w = cababcccababcccab.$$

Opazimo, da se nam v nizu w večkrat ponovita vzorca ab in ccc . Zato uvedemo novi spremenljivki $A = ab$ in $B = ccc$. Sedaj lahko zapišemo w kot

$$w = cAABAABA.$$

Ponovno se nam pojavi vzorec, tokrat AAB . Uvedemo novo spremenljivko $C = AAB$ in zapišemo w kot

$$w = cCCA.$$

Prvotni niz smo z novimi spremenljivkami skrajšali. Kot bomo videli, smo pretvorili niz w v kontekstno neodvisno gramatiko G_w s produkcijskimi pravili

$$S \rightarrow cCCA,$$

$$A \rightarrow ab,$$

$$B \rightarrow ccc,$$

$$C \rightarrow AAB.$$

◇

3. KONTEKSTNO-NEODVISNE GRAMATIKE

Definicija 3.1. *Formalna gramatika* G so pravila, ki nam iz abecede Σ tvorijo jezik, označimo ga z $L(G)$.

Definicija 3.2. *Kontekstno-neodvisna gramatika* je četverica $G = (V, \Sigma, P, S)$, kjer je V končna množica spremenljivk, abeceda Σ množica končnih simbolov tako, da $\Sigma \cap V = \emptyset$, $P \subseteq V \times (V \cup \Sigma)^*$ relacija, ki ji pravimo *produkcijsko pravilo* in $S \in V$ začetna spremenljivka.

Definicija 3.3. Naj bo $G = (V, \Sigma, P, S)$ kontekstno-neodvisna gramatika. Naj bodo $\alpha, \beta, \gamma \in (V \cup \Sigma)^*$ nizi spremenljivk in končnih simbolov, $A \in V$ spremenljivka ter naj bo $(A, \beta) \in P$ produkcijsko pravilo, označimo ga z $A \rightarrow \beta$. Pravimo, da se $\alpha A \gamma$ *prepiše s pravilom* $A \rightarrow \beta$ v $\alpha \beta \gamma$, pišemo $\alpha A \gamma \Rightarrow \alpha \beta \gamma$. Pravimo, da α *porodi* β , če je $\alpha = \beta$ ali če za $k \geq 0$ obstaja zaporedje $\alpha_1, \alpha_2, \dots, \alpha_n \in (V \cup \Sigma)^*$ tako, da

$$\alpha \Rightarrow \alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n \Rightarrow \beta$$

in pišemo $\alpha \xRightarrow{*} \beta$.

Posledica 3.4. Jezik kontekstno neodvisne gramatike G je

$$L(G) = \{w \in \Sigma^* \mid S \xRightarrow{*} w\}.$$

Opomba 3.5. Ime kontekstno-neodvisna gramatika izvira iz oblike produkcijskih pravil. Na levi strani produkcijskega pravila mora vedno stati samo spremenljivka. Torej vsebuje samo pravila oblike

$$A \rightarrow \alpha,$$

kjer je $A \in V$ in $\alpha \in (V \cup \Sigma)^*$. Ne sme pa vsebovati pravila oblike

$$\alpha A \gamma \rightarrow \alpha \beta \gamma,$$

kjer je $A \in V$ in so $\alpha, \beta, \gamma \in (V \cup \Sigma)^*$, saj je pravilo odvisno od predhodnega konteksta. Odvisno je od tega ali pred njim stoji α in za njim β .

Primer 3.6. Formalizirajmo gramatiko iz Primera 2.3, ki smo jo generirali z nizom $w = cababcccababcccab$. Označimo jo z $G_w = (V, \Sigma, P, S)$, kjer je

$$\begin{aligned} V &= \{S, A, B, C\}, \\ \Sigma &= \{a, b, c\}, \\ P &= \{S \rightarrow cCCA, A \rightarrow ab, B \rightarrow ccc, C \rightarrow AAB\}, \\ S &= S. \end{aligned}$$

Vidimo, da G_w ustreza naši definiciji kontekstno-neodvisne gramatike in res kodira w , saj je

$$L(G_w) = \{w\}.$$

Dolžina gramatike G_w je enaka številu črk abecede $\{S, A, B, C, \rightarrow\}$, ki smo jih porabili za opis gramatike in je enaka $|P| = 20$. Vidimo, da niza w nismo skrajšali, saj je $|w| = 17$. Niz w je bil prekratek, da bi ga lahko zares stisnili. Naj bo sedaj

$$w = cababcccababcccabababccc.$$

Gramatika, ki jo generira novi niz se od prejšnje gramatike razlikuje le v

$$P = \{S \rightarrow cCCAC, A \rightarrow ab, B \rightarrow ccc, C \rightarrow AAB\}.$$

Sedaj smo stisnili w z G_w , saj je

$$|w| = 24 > 21 = |P|.$$

◇

Za razliko od klasičnih metod stiskanja, naša metoda ne izvaja procesa stiskanja direktno na nizu w , temveč poiščemo gramatiko G_w , ki generira enojec $\{w\}$ za svoj jezik. Med njimi poiščemo “najmanjšo”, ki jo nato kodiramo. Ker gramatike G_w generira w in ker je gramatika “majhna” jo bomo kodirali v “kratko” kodo. Tako bomo preko gramatike “dobro” stisnili niz w .