

Kontekstno-neodvisne gramatike za kodiranje in stiskanje podatkov

Janez Podlogar

Univerza v Ljubljani, Fakulteta za matematiko in fiziko

16. 9. 2024

2024-09-15

Gramatike za stiskanje podatkov

Kontekstno-neodvisne gramatike
za kodiranje in stiskanje podatkov

Janez Podlogar

Univerza v Ljubljani, Fakulteta za matematiko in fiziko

16. 9. 2024

- Kodiranje je spreminjanje zapisa sporočila.
- Stiskanje podatkov je zapis sporočila v zgoščeni obliki.
- Gramatike je nabor pravil za tvorjenje nizov.

- Ideja je: Stisniti niz z “majhno” gramatiko.

Primer

Naj bo $V = \{S\}$, $\Sigma = \{a, b, c\}$, $P = \{S \rightarrow aSb, S \rightarrow \epsilon\}$ in $S = S$.
Izpeljemo nize

$$S \Rightarrow \epsilon, S \Rightarrow aSb \Rightarrow ab, S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow aabb, \dots$$

- Elementom P pravimo *prepisovalna pravila*.
- Prazen niz ϵ je enota za operacijo stik.
- Niz $\alpha A \gamma$ se *prepiše s pravilom* $A \rightarrow \beta$ v $\alpha \beta \gamma$, pišemo $\alpha A \gamma \Rightarrow \alpha \beta \gamma$.
- Niz α *izpelje* niz β , če lahko α prepišemo v β z uporabo končno mnogo *prepisovalnih pravil*.

Kontekstno-neodvisna gramatika

Definicija

Kontekstno-neodvisna gramatika je četverica $G = (V, \Sigma, P, S)$, kjer je

- V abeceda *nekončnih simbolov*;
- Σ abeceda *končnih simbolov* taka, da $\Sigma \cap V = \emptyset$;
- $P \subseteq V \times (V \cup \Sigma)^*$ celovita relacija;
- $S \in V$ *začetni simbol*.

Gramatike za stiskanje podatkov

└ Kontekstno-neodvisna gramatika

└ Kontekstno-neodvisna gramatika

2024-09-15

Definicija

Kontekstno-neodvisna gramatika je četverica $G = (V, \Sigma, P, S)$, kjer je

- V abeceda *nekončnih simbolov*;
- Σ abeceda *končnih simbolov* taka, da $\Sigma \cap V = \emptyset$;
- $P \subseteq V \times (V \cup \Sigma)^*$ celovita relacija;
- $S \in V$ *začetni simbol*.

- Abeceda je končna neprazna množica.
- $(V \cup \Sigma)^*$ je množica vseh nizov končne dolžine končnih in nekončnih simbolov. Vsebuje tudi prazen niz ε , saj je to niz dolžine 0.
- Okrajšamo z KNG.

Primer

Jezik je $\{a^n b^n \mid n \geq 0\}$.

Kontekstno-neodvisna gramatika

Definicija

Kontekstno-neodvisna gramatika je četverica $G = (V, \Sigma, P, S)$, kjer je

- V abeceda *nekončnih simbolov*;
- Σ abeceda *končnih simbolov* taka, da $\Sigma \cap V = \emptyset$;
- $P \subseteq V \times (V \cup \Sigma)^*$ celovita relacija;
- $S \in V$ *začetni simbol*.

Definicija

Jezik KNG je množica nizov, ki jih lahko izpeljemo iz začetnega simbola in ne vsebujejo nekončnih simbolov.

2024-09-15

Gramatike za stiskanje podatkov

└ Kontekstno-neodvisna gramatika

└ Kontekstno-neodvisna gramatika

- Abeceda je končna neprazna množica.
- $(V \cup \Sigma)^*$ je množica vseh nizov končne dolžine končnih in nekončnih simbolov. Vsebuje tudi prazen niz ε , saj je to niz dolžine 0.
- Okrajšamo z KNG.

Primer

Jezik je $\{a^n b^n \mid n \geq 0\}$.

Definicija

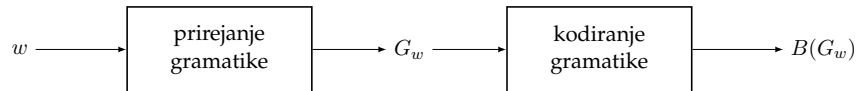
Kontekstno-neodvisna gramatika je četverica $G = (V, \Sigma, P, S)$, kjer je

- V abeceda *nekončnih simbolov*;
- Σ abeceda *končnih simbolov* taka, da $\Sigma \cap V = \emptyset$;
- $P \subseteq V \times (V \cup \Sigma)^*$ celovita relacija;
- $S \in V$ *začetni simbol*.

Definicija

Jezik KNG je množica nizov, ki jih lahko izpeljemo iz začetnega simbola in ne vsebujejo nekončnih simbolov.

Stiskanje niza

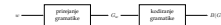


2024-09-15

Gramatike za stiskanje podatkov

- └ Kontekstno-neodvisna gramatika
- └ Stiskanje niza

Stiskanje niza



- Glavna ideja je, da je jezik G_w je $\{w\}$, saj lahko enolično rekonstruiramo niz w .
- Alternativni pristop: KNG G je poznana tako kodirniku kot dekodirniku in jezik KNG G vsebuje vse nize, ki jih želimo stisniti. Posamezni niz stisnemo tako, da stisnemo izpeljavo niza iz začetnega simbola.

Dopustna gramatika

Definicija

KNG je *deterministična*, če vsak nekončen simbol $A \in V$ nastopa natanko enkrat kot leva stran nekega prepisovalnega pravila.

2024-09-15

Gramatike za stiskanje podatkov

- └ Kontekstno-neodvisna gramatika
 - └ Dopustna gramatika
 - └ Dopustna gramatika

- Prejšnji primer ni deterministična gramatika.
- Če odstranimo odstarnimo pravilo $S \rightarrow aSb$, postane deterministična, a je jezik $\{\varepsilon\}$.
- Če odstranimo $S \rightarrow \varepsilon$ postane deterministična, a je njen jezik prazna množica.
- Determinizem sam po sebi ni dovolj močan, da prepreči praznost jezika. Zato bomo od dopustnih gramatik zahtevali, da je njihov jezik neprazen in da ne vsebuje pravila oblike $A \rightarrow \varepsilon$.

Dopustna gramatika

Definicija

KNG je *deterministična*, če vsak nekončen simbol $A \in V$ nastopa natanko enkrat kot leva stran nekega prepisovalnega pravila.

Trditev

Jezik deterministične KNG je enojec ali prazna množica.

2024-09-15

Gramatike za stiskanje podatkov

- └ Kontekstno-neodvisna gramatika
 - └ Dopustna gramatika
 - └ Dopustna gramatika

Definicija

KNG je *deterministična*, če vsak nekončen simbol $A \in V$ nastopa natanko enkrat kot leva stran nekega prepisovalnega pravila.

Trditev

Jezik deterministične KNG je enojec ali prazna množica.

- Prejšnji primer ni deterministična gramatika.
- Če odstranimo odstarnimo pravilo $S \rightarrow aSb$, postane deterministična, a je jezik $\{\varepsilon\}$.
- Če odstranimo $S \rightarrow \varepsilon$ postane deterministična, a je njen jezik prazna množica.
- Determinizem sam po sebi ni dovolj močan, da prepreči praznost jezika. Zato bomo od dopustnih gramatik zahtevali, da je njihov jezik neprazen in da ne vsebuje pravila oblike $A \rightarrow \varepsilon$.

2024-09-15

Gramatike za stiskanje podatkov

└ Kontekstno-neodvisna gramatika

└ Dopustna gramatika

Definicija

KNG ne vsebuje neuporabnih simbolov, če se vsak simbol $y \in V \cup \Sigma$, $y \neq S$ pojavi vsaj enkrat v izpeljavi niza, ki je v jeziku KNG.

Definicija

KNG ne vsebuje neuporabnih simbolov, če se vsak simbol $y \in V \cup \Sigma$, $y \neq S$ pojavi vsaj enkrat v izpeljavi niza, ki je v jeziku KNG.

- Povedano drugače: simbol je neupraven, če se ne pojavi v nobeni izpeljavi niza, ki je v jeziku KNG.

Definicija

KNG je *dopustna gramatika*, če je:

- deterministična,
- ne vsebuje neuporabnih simbolov,
- ima neprazen jezik,
- prazen niz ne nastopa kot desna stran kateregakoli prepisovalnega pravila.

2024-09-15

Gramatike za stiskanje podatkov

└ Kontekstno-neodvisna gramatika

└ Dopustna gramatika

Definicija

KNG je *dopustna gramatika*, če je:

- deterministična,
- ne vsebuje neuporabnih simbolov,
- ima neprazen jezik,
- prazen niz ne nastopa kot desna stran kateregakoli prepisovalnega pravila.

Definicija

KNG je *dopustna gramatika*, če je:

- deterministična,
- ne vsebuje neuporabnih simbolov,
- ima neprazen jezik,
- prazen niz ne nastopa kot desna stran kateregakoli prepisovalnega pravila.

Posledica

Jezik dopustne gramatike je enojec.

2024-09-15

Gramatike za stiskanje podatkov

└ Kontekstno-neodvisna gramatika

└ Dopustna gramatika

Definicija

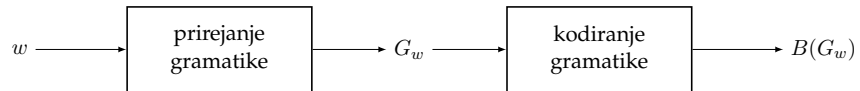
KNG je *dopustna gramatika*, če je:

- deterministična,
- ne vsebuje neuporabnih simbolov,
- ima neprazen jezik,
- prazen niz ne nastopa kot desna stran kateregakoli prepisovalnega pravila.

Posledica

Jezik dopustne gramatike je enojec.

Prirejanje gramatike

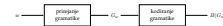


2024-09-15

Gramatike za stiskanje podatkov

└ Prirejanj gramatike

└ Prirejanje gramatike



Asimptotsko kompaktno prirejanje gramatike

Definicija

Prirejanje gramatike nizu abecede \mathcal{A} je *asimptotsko kompaktno*, če za vsak niz $w \in \mathcal{A}^+$ velja

$$\lim_{n \rightarrow \infty} \max_{w \in \mathcal{A}^n} \frac{|G_w|}{|w|} = 0.$$

Gramatike za stiskanje podatkov

└ Prirejanj gramatike

└ Asimptotsko kompaktno prirejanje gramatike

└ Asimptotsko kompaktno prirejanje gramatike

- Z $|G_w|$ označimo vsoto dolžin desnih strani prepisovalnih pravil KNG
- Primer je bisekcijsko prirejanje gramatike.

Primer

Naj bo $w = 000101$. Podčrtamo podnize 000101.

$A_w \rightarrow A_{0001}A_{01}$, $A_{0001} \rightarrow A_{00}A_{01}$, $A_{00} \rightarrow A_00$, $A_{01} \rightarrow A_01$, $A_0 \rightarrow 1$, $A_1 \rightarrow 1$.

$$\lim_{n \rightarrow \infty} \max_{w \in \mathcal{A}^n} \frac{|G_w|}{|w|} = 0.$$

Neskrčljivo prirejanje gramatike

Definicija

Pravimo, da je KNG G *neskrčljiva gramatika*, če:

- 1 Vsak $A \in V$, $A \neq S$ nastopa vsaj dvakrat v desni strani prepisovalnih pravil;
- 2 Ne obstajata $y_1, y_2 \in V \cup \Sigma$, da niz $y_1 y_2$ nastopa kot podniz desne strani kateregali prepisovalnega pravila več kot enkrat na neprekrivajočih se mestih.

Gramatike za stiskanje podatkov

└ Prirejanj gramatike

└ Neskrčljivo prirejanje gramatike

└ Neskrčljivo prirejanje gramatike

2024-09-15

Definicija

Pravimo, da je KNG G *neskrčljiva gramatika*, če:

- 1 Vsak $A \in V$, $A \neq S$ nastopa vsaj dvakrat v desni strani prepisovalnih pravil;
- 2 Ne obstajata $y_1, y_2 \in V \cup \Sigma$, da niz $y_1 y_2$ nastopa kot podniz desne strani kateregali prepisovalnega pravila več kot enkrat na neprekrivajočih se mestih.

- Primer prekrivajočih mest je 111 in 111.
- Primer neprekriavjočih mest je 11 11.
- Primer je metodo najdaljšega ujemaajočega podniza.

Primer

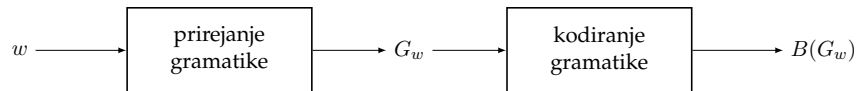
Naj bo $w = 00101010001$. Začnemo z trivialno dopustno gramatiko in podčratno najdaljši podniz, ki se ponovi vsaj dvakrat.

$S \rightarrow \underline{00101010001}$.

$S \rightarrow A\underline{01} \underline{010}A$, $A \rightarrow \underline{001}$.

$S \rightarrow ABB\underline{0}A$, $A \rightarrow \underline{0}B$, $B \rightarrow \underline{01}$.

Binarno kodiranje gramatike



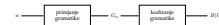
2024-09-15

Gramatike za stiskanje podatkov

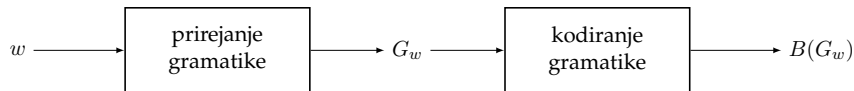
- Binarno kodiranje gramatike

- Binarno kodiranje gramatike

- $H(\mathcal{G})$ je entropija KNG G .



Binarno kodiranje gramatike



Izrek

Obstaja bijektivno brezpredponsko binarno kodiranje gramatike, da

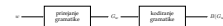
$$\forall G \in \mathcal{G}(\mathcal{A}): |B(G)| \leq |\mathcal{A}| + 4|G| + \lceil H(G) \rceil.$$

2024-09-15

Gramatike za stiskanje podatkov

└ Binarno kodiranje gramatike

└ Binarno kodiranje gramatike



Izrek

Obstaja bijektivno brezpredponsko binarno kodiranje gramatike, da

$$\forall G \in \mathcal{G}(\mathcal{A}): |B(G)| \leq |\mathcal{A}| + 4|G| + \lceil H(G) \rceil.$$

- $H(\mathcal{G})$ je entropija KNG G .

$$\kappa: \mathcal{A}^+ \rightarrow \{0, 1\}^+, \\ w \mapsto B(\pi(w)),$$

Definicija

Stiskanje niza abecede \mathcal{A} z gramatikami je par preslikav kodne in dekodne preslikave $\Phi = (\kappa, \delta)$. Kodna preslikava je

$$\kappa: \mathcal{A}^+ \rightarrow \{0, 1\}^+, \\ w \mapsto B(\pi(w)),$$

kjer je π prirejanje gramatike nizu abecede \mathcal{A} in B binarno kodiranje dopustne gramatike.

- Odvečnost meri količino ponavljajočih se ali predvidljivih podatkov znotraj sporočila, ki jih je mogoče odstraniti, da se prihrani prostor, brez izgube informacije.
- Odvečnost stiskanja z asimptotsko kompaktnim prirejanjem konvergira proti 0 v odvisnosti od izbire kodiranja znotraj razreda.
- Stiskanje z neskrčljivim prirejanjem tudi stiskanje z asimptotsko kompaktnim prirejanjem in odvečnost konvergira enakomerno proti 0 za vsa stiskanja z neskrčljivim prirejanjem, vsaj tako hitro kot $\frac{\log_2 \log_2(n)}{\log_2(n)}$ pomnoženo z neko konstanto.