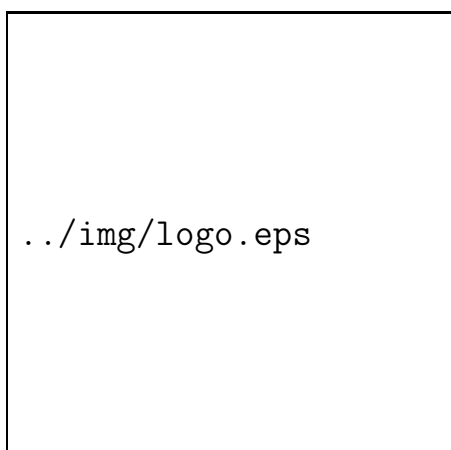


Univerzita Karlova v Praze
Matematicko-fyzikální fakulta

DIPLOMOVÁ PRÁCE



David Marek

Implementace aproximativních bayesovských metod pro odhad stavu v dialogových systémech

Ústav formální a aplikované lingvistiky

Vedoucí diplomové práce: Ing. Mgr. Filip Jurčíček, Ph.D.

Studijní program: program

Studijní obor: obor

Praha 2013

Poděkování.

Prohlašuji, že jsem tuto diplomovou práci vypracoval(a) samostatně a výhradně s použitím citovaných pramenů, literatury a dalších odborných zdrojů.

Beru na vědomí, že se na moji práci vztahují práva a povinnosti vyplývající ze zákona č. 121/2000 Sb., autorského zákona v platném znění, zejména skutečnost, že Univerzita Karlova v Praze má právo na uzavření licenční smlouvy o užití této práce jako školního díla podle §60 odst. 1 autorského zákona.

V dne

Podpis autora

Název práce: Implementace aproximativních bayesovských metod pro odhad stavu
v dialogových systémech

Autor: David Marek

Katedra: Ústav formální a aplikované lingvistiky

Vedoucí diplomové práce: Ing. Mgr. Filip Jurčiček, Ph.D., Ústav formální a aplikované lingvistiky

Abstrakt:

Klíčová slova:

Title:

Author: David Marek

Department: Institute of Formal and Applied Linguistics

Supervisor: Ing. Mgr. Filip Jurčiček, Ph.D., Institute of Formal and Applied Linguistics

Abstract:

Keywords:

Obsah

Úvod	2
0.1 Cíle	2
1 Teorie dialogových systémů	4
1.1 Dialogový systém	4
1.2 Součásti dialogového systému	4
1.2.1 Systém rozpoznávání řeči (ASR)	5
1.2.2 Porozumění mluvené řeči (SLU)	5
1.2.3 Dialogový manager (DM)	6
1.2.4 Generování přirozené řeči (NLG a TTS)	7
1.3 Dialogový stav	7
1.3.1 Aktualizace dialogového stavu	8
2 Bayesovské sítě a inference	10
2.1 Bayesovské sítě	10
2.2 Inference v Bayesovských sítích	11
2.2.1 Exaktní inference	12
2.2.2 Posílání zpráv ve faktor grafu	16
2.2.3 Belief Propagation	18
2.2.4 Loopy Belief Propagation	19
2.2.5 Strategie výběru vrcholu v LBP algoritmu	22
2.3 Propagace s aproximovanými zprávami	23
2.3.1 Rodina exponenciálních rozdělení	24
2.3.2 Expectation Propagation	25
3 Učení parametrů	28
3.1 Grafický model	28
3.2 Výpočet marginálních pravděpodobností	28
3.2.1 Marginální pravděpodobnost proměnných	28
3.2.2 Marginální pravděpodobnost parametrů	30
3.3 Aproximace marginálních pravděpodobností	31
3.4 Algoritmus	35
4 Implementace	36
Závěr	37
Seznam použité literatury	38
Seznam tabulek	40
Seznam použitých zkratk	41
Přílohy	42

Úvod

Dialog je přirozený způsob dorozumívání a sdělování informací mezi lidmi. Počítač, který by dokázal vést dialog s uživatelem, byl vždy snem nejen příznivců vědecko-fantastické literatury. Už pro první počítače vnikaly programy, které se snažily využívat přirozenou řeč pro interakci s uživatelem. Jedním z takových programů byl například Eliza, program, který předstíral, že jej zajímá, co mu uživatel říká. Fungoval na principu rozpoznání textu pomocí gramatiky a následné transformace textu do promluv dle pravidel. Avšak gramatiky a pravidlové systémy se ukázaly nedostačné pro praktické aplikace a tak se vývoj přesunul do statistických metod. S využitím statistických metod a metod strojového učení bylo možné začít s porozumíváním mluveného slova. Přijetí bylo zpočátku chladné a veřejnost byla

0.1 Cíle

1. Dialogové systémy často využívají pouze nejlepší hypotézu ze systému porozumění přirozené řeči. Většina ovšem umí vytvořit seznam n nejlepších hypotéz. Tato práce si klade za cíl představit metody pro inferenci dialogového stavu v dialogovém systému s využitím více hypotéz. Představené metody budou založeny na reprezentaci dialogového stavu pomocí dynamických bayesovských sítí.
2. Bude představen algoritmus Loopy Belief Propagation, který bude implementován pro využití v reálných systémech pro odhad dialogového stavu.
3. Algoritmus bude otestován na datech z reálného dialogového systému Let's Go a porovnán s dalšími systémy, které se účastnily soutěže Dialog State Tracking Challenge 2013.
4. Důležitou částí algoritmů pro inferenci je určení pořadí, v jakém má inference probíhat. Práce bude obsahovat implementaci několika strategií pro inferenci, které budou umožňovat efektivní inferenci pro různé druhy bayesovských sítí (stromy, dynamické sítě, obecné grafy).
5. Nakonec se práce bude zabývat učením parametrů sítě a představí algoritmus Expectation Propagation, který umožňuje inferenci v bayesovských sítích se spojitými náhodnými proměnnými.
6. Při většině reálných použití dochází k aproximacím a úpravám modelu v závislosti na problému tak, aby bylo možné inferenci provádět v reál-

ném čase. Práce bude obsahovat implementaci frameworku, do kterého je možné jednoduše zasadit vlastní moduly pro aproximaci pravděpodobnostních rozdělání, které bude algoritmus Expectation Propagation používat.

7. Jako příklad bude ukázán systém pro učení parametrů pravděpodobnostního rozložení pro pozorování v generickém modelu pro reprezentaci dialogového stavu.

1. Teorie dialogových systémů

1.1 Dialogový systém

Dialogový systém je počítačový systém, který umožňuje uživatelům komunikovat s počítačem ve formě, která je přirozená a efektivní pro použití. Dialogové systémy stále mají spoustu problémů k překonání a pro praktické použití je třeba se uchýlit k několika předpokladům a zjednodušením. Prvním zjednodušením je doménová specializace, v současné době není možné vytvořit dialogový systém, který by se dokázal s uživatelem bavit o libovolném tématu. Vždy je potřeba při vývoji dialogového systému mít ontologii určující, jaké informace má systém poskytovat a o čem se může chtít uživatel bavit.

Další zjednodušení se týká přímo dialogu. Předpokládá se, že dialog probíhá vždy mezi systémem a jedním uživatelem. Navíc se pravidelně střídají v obrátkách. Jedna obrátka dialogu je složená z jedné promluvy systému a jedné promluvy uživatele.

Příkladem dialogového systému může být systém pro nalezení spojení pomocí městské dopravy. Příkazu od uživatele může vypadat např. takto: „Chci jet z Malostranského náměstí na Anděl“. Dialogový systém se nyní může rozhodnout, zda-li mu zadané informace stačí pro nalezení spojení. V tomto případě systém stále neví, kdy chce uživatel jet. Může předpokládat, že uživatel už na zastávce stojí a tak tedy uživateli nalezne nejbližší spojení.

Důležitou vlastností dialogového systému je robustnost. Pokud budeme používat dialogový systém v přirozeném prostředí, musíme se vyrovnat s tím, že často nebude uživateli rozumět. Systém může informaci přeslechnout, anebo si nemusí být jistý tím, co slyšel. Dialogový systém se proto musí umět uživatele doptat na chybějící informace

1.2 Součásti dialogového systému

Dialogový systém se skládá z několika částí, které spolu komunikují. Na vstupu je zvukový záznam uživatele, o jeho převedení do textu se stará systém rozpoznávání řeči (ASR). Z textu je potřeba získat sémantické informace pomocí systému porozumění mluvené řeči (SLU). Nad sémanticky anotovanými informacemi už může pracovat dialogový manager (DM), který zvolí patřičnou odpověď. Výstupem dialogového manageru jsou informace, které se mají předat uživateli. O jejich převedení do textu se stará systém generování přirozené řeči (NLG). Do zvukového záznamu převede text syntetizér řeči (TTS).

1.2.1 Systém rozpoznávání řeči (ASR)

Systém rozpoznávání řeči slouží k převedení mluveného projevu do textové podoby. Po získání textové podoby je teprve možné se zabývat významem textu. Aktuálně nejlepší systémy jsou založené na pravděpodobnostním modelu a využívají skryté Markovské modely (HMM) k určení nejpravděpodobnější sekvence slov pro daný zvukový záznam. Pro tuto část dialogového systému existuje řada dostupných otevřených toolkitů, např. systém HTK [21], Kaldi [14] nebo SPHINX [19]. Existuje i celá řada komerčního software od firem jako IBM nebo Nuance.

Úspěšnost systému rozpoznání řeči je závislá na obtížnosti úlohy a na počtu trénovacích dat, pocházejících ze stejné domény. Pro obecnou doménu se problém stává mnohem těžší a je třeba velké množství dat. Word error rate (WER) je častá metrika pro počítání výkonu ASR. Pro spočítání WER je třeba nejprve provést zarovnání rozpoznávaného a originálního textu. WER je pak počet slov, která jsou změněná, smazaná nebo přidaná, vydělený počtem slov v originálním textu. Systém Let's Go! [15] dosahuje průměrné WER 64.3%.

Systém rozpoznávání řeči může produkovat více hypotéz pro jeden vstup. Často existuje pro jeden zvukový záznam více možných slovních sekvencí, z kterých by mohl pocházet. Reprezentace možných hypotéz může být seznam slovních sekvencí s jejich věrohodností. Věrohodnosti jsou skóre přiřazené hypotézám, které určují jakou důvěru má systém rozpoznávání řeči ve správnost dané slovní sekvence. V ideálním případě je věrohodnost ekvivalentní aposteriorní pravděpodobnosti sekvence slov, dáno vstupní zvuk. Ovšem ne všechny rozpoznávače pracují na pravděpodobnostním principu a pak není možné od nich požadovat skutečné pravděpodobnosti.

Další možností jak reprezentovat výstup je použití konfúzní sítě [3]. Konfúzní síť je vážený orientovaný graf, obsahující startovní a konečný vrchol a hrany označené slovy. Každá cesta ze startovního do konečného vrcholu vede přes všechny ostatní vrcholy. Váhy hran jsou pravděpodobnosti slova přiřazeného dané hraně. Hrany mohou obsahovat i prázdné slovo ϵ . Pravděpodobnost sekvence slov je součinem vah po cestě ze startovního do konečného uzlu. Výhodou konfúzní sítě je, že umožňuje v komprimované podobě uložit mnohem více hypotéz.

1.2.2 Porozumění mluvené řeči (SLU)

Jakmile má systém seznam možných hypotéz toho, co uživatel řekl, musí se pokusit porozumět, co tím uživatel myslel. Dialogový systém nepotřebuje vědět, co přesně uživatel řekl, důležité je pouze zjistit, co se uživatel snaží sdělit. Pokud například uživatel řekne "Chtěl bych nalézt spojení z Malostranského náměstí na Anděl", anebo "Jak se dostanu na Anděl ze zastávky Malostranské náměstí?", tak

Přidat
přík-
lad
sez-
na-
mu
hy-
potéz

Přidat
přík-
lad
kon-
fúzní
sítě

výsledek je stejný, uživatel požaduje informace o spojení mezi dvěma zastávkami, i když v jednom případě jde o větu oznamovací a v druhém případě o otázku.

Semantická reprezentace sdělení uživatele se nazývá dialogový akt (DA), skládá se z jedné nebo více položek dialogového aktu (DAI), které jsou spojené v konjunkci. Každá DAI se skládá z typu, názvu slotu a jeho hodnoty. Typy jsou doménově nezávislé, sloty a jejich hodnoty reprezentují koncepty ontologie. Příklad dialogového aktu z dialogového systému pro hledání restaurací:

`hello()&inform(food="chinese").`

Zde se dialogový akt skládá ze dvou položek, první položka má pouze typ *hello*, značící pozdrav. Druhá položka má typ *inform*, tzn. uživatel nás informuje o svém požadavku. Název slotu je *food* a hodnota je „chinese“, tedy uživatel nám říká, že hledá restauraci, kde servírují čínské jídlo.

Typů může být libovolné množství, ale existuje několik základních, jejichž použití je ustálené.

- *inform* — sdělujeme informaci, doplňujeme hodnotu do slotu,
- *request* — požadujeme od protějšku doplnění hodnoty pro dotazovaný slot,
- *confirm* — chceme potvrdit hodnotu slotu, potvrzení může být implicitní, anebo explicitní. Při explicitním potvrzení očekáváme odpověď buď „Ano“ nebo „Ne“, U implicitního, pokud se nám nedostane odpovědi předpokládáme, že protějšek souhlasí,
- *select* — žádáme protějšek, aby zvolil z nabízených možností.

Existuje široké množství technik, které lze použít pro porozumění mluvené řeči. Unifikace pomocí šablon anebo gramatiky jsou příklady ručně psaných metod. Metody založené na datech jsou například Hidden Vector State model [7], techniky strojového překladu [20], Combinatory Categorical Grammars [24] nebo Support Vector Machines [10].

1.2.3 Dialogový manager (DM)

Pokud už jsou pravděpodobné dialogové akty dekodovány, je třeba rozhodnout, co bude systém dělat. Komponenta tvořící tato rozhodnutí je dialogový manager. Odpověď systému je zakódována do formy dialogových aktů a nazývá se systémová akce.

Zvolená systémová akce je vybrána z množiny možných akcí $a \in \mathcal{A}$ a závisí na vstupu, který systém obdržel z SLU. Tento vstup se nazývá pozorování $o \in \mathcal{O}$, protože obsahuje vše, co systém pozoroval o uživateli.

Zvolení správné akce potřebuje více znalostí než jen poslední pozorování. Celá historie dialogu a také kontext hrají důležitou roli. Dialogový manager bere na vše ohled pomocí udržování interní reprezentace celého pozorovaného dialogu. Tato reprezentace se nazývá dialogový stav, nebo také belief stav, značí se $b \in \mathcal{B}$. Aktuální dialogový stav závisí na přechodové funkci, která dialogový stav aktualizuje pro každé nové pozorování a systémovou akci. Přechodová funkce je tedy mapování $\mathcal{T} : \mathcal{B} \times \mathcal{A} \times \mathcal{O} \rightarrow \mathcal{B}$. V této práci se budeme věnovat právě metodám aktualizace dialogového stavu.

Chování dialogového stavu definuje dialogová strategie π . Strategie určuje co má systém provést v závislosti na aktuálním dialogovém stavu. Obecně strategie vytvoří pravděpodobnostní rozložení přes možné akce. Pokud $\prod(\mathcal{A})$ značí množinu těchto distribucí, pak dialogová strategie bude zobrazení z dialogového stavu do této množiny, $\pi : \mathcal{B} \rightarrow \prod(\mathcal{A})$.

Pozorování, dialogový stav a akce jsou číslovány podle obrátky. Pokud je časový okamžik důležitý, jsou pozorování, dialogový stav a akce z obrátky číslo t označeny o_t , b_t a a_t .

1.2.4 Generování přirozené řeči (NLG a TTS)

Posledním krokem v tahu dialogového systému je vytvoření odpovědi pro uživatele. Nejprve systém generování přirozené řeči převede dialogové akty na text. Následně je text převeden na zvuk pomocí textového syntetizéru řeči.

Nejjednodušším přístupem ke generování přirozeného jazyka z dialogových aktů je použití šablon. Například pro dialogový akt `inform(type="x")` bude vytvořena šablona „Restaurace servíruje x jídlo“, kde „x“ bude nahrazeno například za „čínské“, „indické“, atd. Šablony se při generování osvědčily, protože počet možných dialogových je většinou zvládnutelný.

Při syntéze řeči existuje mnoho alternativ. Je možné použít segmenty řeči z databáze pro vygenerování zvuků tvořících dohromady celou sekvenci slov. Příkladem těchto systémů je Festival [5] nebo FLite [6].

Alternativní metodou syntézy je použití skrytých Markovských modelů pro generování zvuku, příkladem je HTS systém [23].

1.3 Dialogový stav

Nejistota je základním problémem, s kterým se dialogový systém musí vypořádat. Systémy pro rozpoznávání i porozumění řeči často chybují a tuto možnost musí brát dialogový manager v potaz. Vypořádat se s nejistotou lze pomocí jejího zakomponování do modelu pro odhad dialogového stavu.

Cíle uživatele a další vlastnosti prostředí lze považovat za náhodné částečně pozorovatelné proměnné a je možné je odvodit z pozorování. Pravděpodobnostní rozložení těchto náhodných proměnných dává dobře definovanou reprezentaci nejistoty, navíc je možné je reprezentovat pomocí Bayesovské sítě.

Jedním z možných modelů dialogového stavu je generativní model. Definujeme množinu stavů prostředí, $s \in \mathcal{S}$. Předpokládáme, že pozorování závisí podmíněně pouze na stavu prostředí a definujeme pravděpodobnostní rozdělení pro pozorování, $p(o \mid s)$. Dále předpokládáme, že cíle uživatele se nemění v čase a stav prostředí je tedy závislý pouze na stavu v předchozí obrátce a na poslední akci systému. Tato závislost je zachycena v přechodové pravděpodobnosti $p(s_{t+1} \mid s_t, a_t)$. Předpoklad, že stav prostředí závisí pouze na minulé hodnotě se nazývá Markovská vlastnost.

Pokud vezmeme předchozí předpoklady, tak lze využít bayesovský přístup pro počítání s nejistotou. Stav v čase t označíme s_t . Podle Bayesova vzorce můžeme spočítat pravděpodobnost stavu v čase $t+1$ po přijetí nového pozorování $o_{t+1} = o'$.

$$p(s_{t+1} = s') \propto \sum_{s \in \mathcal{S}} p(s_t = s) p(s_{t+1} = s' \mid s_t = s, a_t = a) p(o_{t+1} = o' \mid s_{t+1} = s') \quad (1.1) \quad \text{\texttt{\{eq:belief}}}$$

Nyní můžeme definovat dialogový stav v čase t , b_t , jako pravděpodobnost přes stavy dáno všechna pozorování až do času t . Množina všech možných dialogových stavů je pravděpodobnost přes všechny možné stavy prostředí $\mathcal{B} = \prod(\mathcal{S})$.

Můžeme přepsat rovnici (1.1) s pomocí dialogových stavů.

$$b(s_{t+1}) \propto \sum_{s \in \mathcal{S}} b(s_t) p(s_{t+1} \mid s_t, a_t) p(o_{t+1} \mid s_{t+1}) \quad (1.2) \quad \text{\texttt{\{eq:belief}}}$$

Rovnice (1.2) nám dává předpis pro přechodovou funkci \mathcal{T} . V praxi ovšem bude množina možných hodnot pro slot s_t příliš velká, protože stav musí obsahovat všechny informace potřebné pro rozhodování, to znamená celou historii dialogu a cíle uživatele. Pokud systém obsahuje sloty, tak každá kombinace hodnot slot je jedním možným cílem uživatele. Tedy velikost stavového prostoru roste exponenciálně.

1.3.1 Aktualizace dialogového stavu

Efektivní metodou pro aktualizaci dialogového stavu je použití dynamických bayesovských sítí [17]. Bayesovské sítě umožňují efektivní výpočet využitím podmíněných nezávislostí mezi sloty. Stále ovšem zůstává problém s výpočtem, pokud i jednotlivé sloty obsahují příliš mnoho hodnot. Lze použít aproximace a počítat

jen s k nejpravděpodobnějšími hodnotami [18].

Alternativním zjednodušením je rozdělit stav prostředí do skupin. Tento přístup se nazývá Hidden Information State (HIS) [22]. Základním předpokladem zde musí být, že uživatel nezmění svůj cíl v průběhu dialogu. Pak lze efektivně provádět aktualizaci, protože rovnice pro aktualizace pravděpodobnosti se nemění mezi jednotlivými skupinami.

V této práci se budeme zabývat prvním přístupem, tedy použitím Bayesovských sítí. Pro inferenci použijeme Loopy Belief Propagation (LBP) algoritmus, který je aproximativní metodou pro sítě s diskrétními náhodnými proměnnými. Pro učení parametrů představíme Expectation Propagation (EP) algoritmus. EP je zobecněním LBP na libovolné pravděpodobnostní rozložení.

2. Bayesovské sítě a inference

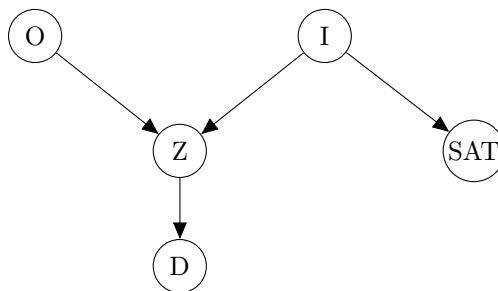
V této kapitole představíme Bayesovské sítě, grafický model pro efektivní reprezentaci pravděpodobnostních rozdělání a nezávislostí mezi náhodnými proměnnými. Bayesovská síť zároveň vytváří koncept pro inferenci, tedy zodpovídání dotazů nad proměnnými v síti. Ukážeme si několik přístupů k inferenci, nejprve naivní výpočet vycházející přímo z definice. Následně využijeme vlastností sítě a konceptů dynamického programování pro jeho zlepšení. Analýzou složitosti exaktní inference dojdeme k závěru, že pro větší a komplexnější modely bude třeba se uchýlit k aproximacím.

Nejprve aproximujeme sdruženou pravděpodobnostní distribuci součinem marginálních distribucí a představíme Loopy Belief Propagation algoritmus. Stále můžeme použít pouze diskrétní pravděpodobnostní distribuce se známými parametry. Pro učení parametrů lze použít Expectation Maximization metodu, pro dialogové systémy je ovšem těžké získat dostatek učících dat. Nakonec se dostaneme k metodě Expectation Propagation, která je zobecněním LBP a je tedy možné ji použít pro libovolné rozdělení. Díky ní budeme schopni vytvořit generativní model pro aktualizaci dialogového stavu, který bude pracovat stejně jako LBP, ale bude schopen adaptace.

2.1 Bayesovské sítě

Bayesovské sítě jsou pravděpodobnostní grafický model, který využívá podmíněných nezávislostí pro úspornou reprezentaci sdružené pravděpodobnosti. Bayesovská síť je orientovaný acyklický graf, jeho vrcholy jsou náhodné proměnné a hrany odpovídají přímé závislosti jednoho uzlu na druhý. Pro každou náhodnou proměnnou v síti platí, že její pravděpodobnost je jednoznačně určena jejími rodiči v grafu. Podmíněná pravděpodobnostní distribuce (CPD) proměnné X popisuje pravděpodobnost proměnné X dáno její rodiče, $P(X \mid \text{parents}(X))$. Pokud proměnná nemá žádná rodiče, pak její podmíněná pravděpodobnostní distribuce je ekvivalentní marginální pravděpodobnostní distribuci.

Příklad Student [8]: firma zvažuje, zda-li přijme studenta. Firma chce přijímat chytré studenty, ale nesmí je testovat na inteligenci (I) přímo. Má však výsledek studentových SAT testů, které ale nemusí stačit pro správné zhodnocení inteligence. Požadují tak tedy i doporučení (D) od jednoho z učitelů. Učitel studentovi napíše doporučující dopis na základě známky (Z), kterou student získal v jeho předmětu. Předměty se ovšem liší v obtížnosti (O) a tak je studentova známka v předmětu závislá nejen na jeho inteligenci, ale také na obtížnosti předmětu. Grafický model reprezentující tento problém je vyobrazen na obrázku 2.1.



{fig:student}

Obrázek 2.1: Bayesovská síť pro příklad se studentem.

V tomto modelu je několik nezávislostí. Obtížnost předmětu a inteligence studenta jsou zjevně nezávislé. Studentova známka z předmětu je závislá na obtížnosti předmětu a inteligenci studenta, ale je podmíněně nezávislá na jeho výsledku ze SAT, dáno studentova inteligence. Konečně doporučení, které student obdrží, je podmíněně nezávislé na všech ostatních proměnných, dáno studentova známka.

Sdruženou nezávislost tohoto modelu lze zapsat ve formě podmíněných pravděpodobnostních distribucí s pomocí řetízkového pravidla.

$$P(O, I, Z, S, D) = P(D | Z)P(Z | O, I)P(SAT | I)P(O)P(I) \quad (2.1)$$

Předpokládejme, že obtížnost předmětu, inteligence studenta, doporučující dopis a výsledek SAT jsou binární proměnné. Známka z předmětu pak je ternární proměnná. Pokud bychom zapsali sdruženou pravděpodobnost ve formě tabulky, tak se dostaneme k 48 položkám. Díky rozdělení do podmíněných pravděpodobnostních rozložení, které nám bayesovská síť poskytuje, se dostáváme k $2 + 2 + 12 + 4 + 6 = 26$ položkám. Tedy i v tomto jednoduchém modelu dochází k značné úspoře.

2.2 Inference v Bayesovských sítích

Bayesovské sítě reprezentují pravděpodobnostní model a umožňují nám nad ním provádět dotazy. Můžeme se například ptát na marginální pravděpodobnost jednotlivých proměnných. Tu získáme marginalizací sdružené pravděpodobnosti, pokud vezmeme příklad se studentem a budeme chtít znát marginální pravděpodobnost známek, musíme vysčítat všechny ostatní proměnné.

$$P(Z) = \sum_{O, I, S, D} P(D | Z)P(Z | O, I)P(SAT | I)P(O)P(I) \quad (2.2)$$

Další a asi nejčastější dotaz nastává, pokud některé náhodné proměnné pozorujeme. Pak chceme znát pravděpodobnost jiných proměnných dáno naše po-

zorování, $P(\mathbf{X} \mid \mathbf{E} = \mathbf{e})$, kde \mathbf{X} jsou dotazované proměnné, \mathbf{E} jsou pozorované proměnné a \mathbf{e} jsou pozorované hodnoty. Z definice podmíněné pravděpodobnosti dostáváme

$$P(\mathbf{X} \mid \mathbf{E} = \mathbf{e}) = \frac{P(\mathbf{X}, \mathbf{E} = \mathbf{e})}{P(\mathbf{E} = \mathbf{e})} \quad (2.3)$$

Každou instanci jmenovatele $P(\mathbf{X}, \mathbf{E} = \mathbf{e})$ jde spočítat sumou sdružených pravděpodobností s ohodnoceními proměnných, které jsou kompatibilní s pozorováním a aktuální instancí. Pokud počítáme instanci $P(\mathbf{X} = \mathbf{x}, \mathbf{E} = \mathbf{e})$, pak získáme výsledek marginalizací všech proměnných, kromě \mathbf{X} a \mathbf{E} , které jsou fixovány. Pokud tedy množinu všech proměnných bez \mathbf{X} a \mathbf{E} označíme $\mathcal{W} = \mathcal{X} - \mathbf{X} - \mathbf{E}$, pak pravděpodobnost dané instance je

$$P(\mathbf{X} = \mathbf{x}, \mathbf{E} = \mathbf{e}) = \sum_{\mathbf{w}} P(\mathbf{x}, \mathbf{e}, \mathbf{w}) \quad (2.4)$$

Pro výpočet normalizační konstanty $P(\mathbf{E})$ musíme opět marginalizovat sdruženou pravděpodobnost, anebo si můžeme povšimnout, že platí

$$P(\mathbf{E} = \mathbf{e}) = \sum_{\mathbf{x}} P(\mathbf{x}, \mathbf{e}) \quad (2.5)$$

a tedy můžeme použít už vypočítané hodnoty.

2.2.1 Exaktní inference

V předchozí části jsme viděli, že pomocí definice podmíněné pravděpodobnosti a marginalizace sdružené pravděpodobnosti lze najít odpověď na libovolný dotaz. Nyní si ukážeme algoritmus, který využívá struktury Bayesovské sítě pro inferenci a navíc díky metodám dynamického programování umožňuje samotný výpočet urychlit. Nakonec ovšem zjistíme, že pro velké sítě, které nás většinou zajímají nejvíce, nám přesná inference nebude stačit a musíme se uchýlit k aproximacím.

Začneme s inferencí v jednoduchém modelu $A \rightarrow B \rightarrow C \rightarrow D$. Sdružená pravděpodobnost $P(A, B, C, D)$ je součinem jednotlivých podmíněných pravděpodobnostních distribucí

$$P(A, B, C, D) = P(D \mid C)P(C \mid B)P(B \mid A)P(A) \quad (2.6)$$

Pokud nyní budeme chtít spočítat marginální distribuci D , tak musíme marginalizovat všechny ostatní proměnné

$$P(D) = \sum_{A, B, C} P(D \mid C)P(C \mid B)P(B \mid A)P(A) \quad (2.7)$$

Můžeme si povšimnout, že spousta členů se bude počítat vícekrát. Využitím metod dynamického programování a přeuspořádáním sum si můžeme mezivýsledky uložit a použít vícekrát.

$$P(D) = \sum_C P(D | C) \sum_B P(C | B) \sum_A P(B | A) P(A) \quad (2.8)$$

Při výpočtu pak nejprve spočítáme $\psi_1(A, B) = P(B | A)P(A)$, pak vysčítáme proměnnou A a získáme $\tau_1(B) = \sum_A \psi_1(A, B)$. Pokračujeme obdobně

$$\psi_2(B, C) = P(C | B)\tau_1(B) \quad (2.9)$$

$$\tau_2(C) = \sum_B \psi_2(B, C) \quad (2.10)$$

A nakonec spočítáme finální marginální pravděpodobnost

$$\psi_3(D, C) = P(D | C)\tau_2(C) \quad (2.11)$$

$$P(D) = \sum_C \psi_3(D, C) \quad (2.12)$$

Definice 1. *Nechť \mathcal{X} je množina náhodných proměnných. Potom definujeme faktor ϕ jako zobrazení z $Val(\mathcal{X})$ do \mathbb{R} . Faktor je nezáporný, pokud všechny jeho obrazy jsou nezáporné. Množina proměnných \mathcal{X} je doménou faktoru a značíme ji jako $Dom(\phi)$.*

Faktor, jehož doménu tvoří diskrétní proměnné, si můžeme představit jako tabulku, která obsahuje jednu hodnotu pro každé možné ohodnocení proměnných z domény.

Definice 2. *Nechť $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$ jsou tři disjunktní množiny náhodných proměnných. Nechť $\phi_1(\mathbf{X}, \mathbf{Y})$ a $\phi_2(\mathbf{Y}, \mathbf{Z})$ jsou faktory. Definujeme součin faktorů $\phi_1 \times \phi_2$ jako faktor $\psi : Val(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \rightarrow \mathbb{R}$ následovně:*

$$\psi(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = \phi_1(\mathbf{X}, \mathbf{Y}) \cdot \phi_2(\mathbf{Y}, \mathbf{Z})$$

Násobíme prvky, které mají stejné ohodnocení společných proměnných \mathbf{Y} . Stejný princip použijeme pro všechny matematické operace.

Definice 3. *Nechť \mathbf{X} je množina náhodných proměnných a $Y \notin \mathbf{X}$ náhodná proměnná. Nechť $\phi(\mathbf{X}, Y)$ je faktor. Definujeme marginalizaci Y v ϕ , značenou $\sum_Y \phi$, jako faktor ψ s doménou \mathbf{X} takový, že*

$$\psi(\mathbf{X}) = \sum_Y \phi(\mathbf{X}, Y)$$

Této operaci také říkáme *vysčítání* Y ve ϕ .

Faktory se při počítání chovají jako čísla, všechny operace probíhají po prvcích, pro které je ohodnocení náhodných proměnných z průniku domén faktorů stejné. Proto platí komutativita $\phi_1 \cdot \phi_2 = \phi_2 \cdot \phi_1$ a $\sum_X \sum_Y \phi = \sum_Y \sum_X \phi$. Dále platí asociativita součinu $(\phi_1 \cdot \phi_2) \cdot \phi_3 = \phi_1 \cdot (\phi_2 \cdot \phi_3)$. Nakonec můžeme vyměnit sumu a součin, pokud $X \notin \text{Dom}(\phi_1)$, potom $\sum_X (\phi_1 \cdot \phi_2) = \phi_1 \sum_X \phi_2$.

Sdruženou pravděpodobnost z minulého příkladu tedy můžeme přepsat do formy faktorů.

$$P(A, B, C, D) = \phi_A \cdot \phi_B \cdot \phi_C \cdot \phi_D \quad (2.13)$$

Opět se pokusíme spočítat marginální pravděpodobnost proměnné D .

$$P(D) = \sum_C \sum_B \sum_A \phi_A \phi_B \phi_C \quad (2.14)$$

$$= \sum_C \phi_D \cdot \left(\sum_B \phi_C \cdot \left(\sum_A \phi_A \cdot \phi_B \right) \right) \quad (2.15)$$

Přesuny sum můžeme provést díky doméně jednotlivých faktorů. Faktory ϕ_C a ϕ_D neobsahují proměnnou A a tedy je můžeme vytknout před sumu přes A . Stejně tak faktor ϕ_D neobsahuje proměnnou B a opět jej můžeme vytknout před sumu přes B . Tyto úpravy můžeme provádět v libovolném pořadí, pokud vždy platí, že vysčítáme proměnnou X až poté, co spolu vynásobíme všechny faktory, které ji obsahují.

V obecnosti vždy počítáme výraz, který je ve tvaru

$$\sum_X \prod_{\phi \in \Phi} \phi.$$

Z tohoto také vychází název pro tuto metodu: sum-product. Jednoduchý algoritmus pro exaktní inferenci využívající tuto metodu se nazývá eliminace proměnných. Základní myšlenka je, že máme dán seznam náhodných proměnných v pořadí, v jakém se mají eliminovat. Pro eliminaci proměnné je třeba nejprve vynásobit všechny faktory, které ji obsahují a následně ji vysčítat. Tak získáme faktor, který už tuto proměnnou neobsahuje, tedy jsme ji eliminovali. Eliminace proměnných je popsána v algoritmu 1.

Algoritmus 1 Eliminace proměnných

{alg:ve}

function SUM-PRODUCT-VE(Φ, \mathbf{X}, \prec) Φ množina všech faktorů. \mathbf{X} množina náhodných proměnných, které mají být eliminovány. \prec pořadí proměnných, v jakém mají být eliminovány.Nechť X_1, \dots, X_k je seřazení proměnných z \mathbf{X} , t.ž. $X_i \prec X_j \Leftrightarrow i < j$.**for** $i = 1 \dots k$ **do** $\Phi \leftarrow \text{Sum-Product-Eliminate-Var}(\Phi, X_i)$ **end for** $\phi^* \leftarrow \prod_{\phi \in \Phi} \phi$ **return** ϕ^* **end function****function** SUM-PRODUCT-ELIMINATE-VAR(Φ, X) Φ množina všech faktorů. X , proměnná, která má být eliminována. $\Phi' \leftarrow \{\phi \in \Phi : X \in \text{Dom}(\phi)\}$ $\Phi'' \leftarrow \Phi - \Phi'$ $\psi \leftarrow \prod_{\phi \in \Phi'} \phi$ $\tau \leftarrow \sum_X \psi$ **return** $\Phi'' \cup \{\tau\}$ **end function**

Věta 1. *Nechť \mathbf{X} je množina náhodných proměnných, nechť Φ je množina faktorů, t.ž. pro každé $\phi \in \Phi$, $\text{Dom}(\phi) \subseteq \mathbf{X}$. Nechť $\mathbf{Y} \subset \mathbf{X}$ je množina dotazovaných náhodných proměnných a nechť $\mathbf{Z} = \mathbf{X} - \mathbf{Y}$. Pak pro každé seřazení \prec nad \mathbf{Z} , $\text{Sum-Product-VE}(\Phi, \mathbf{Z}, \prec)$ vrátí faktor $\phi^*(\mathbf{Y})$ takový, že*

$$\phi^*(\mathbf{Y}) = \sum_{\mathbf{Z}} \prod_{\phi \in \Phi} \phi$$

Nyní provedeme analýzu algoritmu eliminace proměnných. Předpokládejme, že na vstupu je n proměnných. Bayesovská síť obsahuje pro každou proměnnou jeden faktor. Pro jednoduchost budeme předpokládat, že algoritmus bude eliminovat všechny proměnné. Běh algoritmu se skládá z jednotlivých eliminačních kroků, při kterých je vždy eliminována jedna proměnná.

Při jednom eliminačním kroku je vybrána proměnná X_i , všechny faktory, které ji obsahují jsou pronásobeny a vytvoří jeden velký faktor ψ_i , proměnná X_i je pak vysčítána z tohoto faktoru. Počet operací pro jeden eliminační krok tedy závisí na velikosti faktoru ψ_i , označme ji N_i . Maximum z velikostí faktorů označme $N_{\max} = \max_i N_i$.

Nyní se zaměříme na počet násobení. Celkem vznikne $n + m$ faktorů, kde

m je počet faktorů, které vznikly vysčítáním proměnné. Každý z těchto faktorů je zahrnut do součinu pouze jednou, pokud je eliminována nějaká proměnná, kterou obsahuje. Cena násobení faktorů pro vznik ψ_i je nejvýše N_i . Celkový počet násobení tedy bude nejvýše $(n + m)N_{max}$ což je $\mathcal{O}(nN_{max})$.

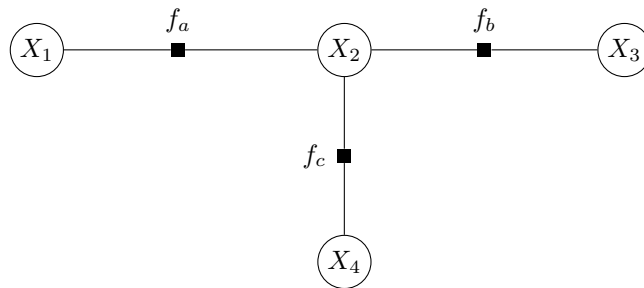
Pokud k je maximum z velikostí domén proměnných, pak velikost faktoru obsahujícího n proměnných může být až k^n . Složitost eliminace proměnných je tedy dominována velikostí faktorů, které vznikají při výpočtu a je exponenciální. Navíc bylo dokázáno, že výběr nejlepšího pořadí proměnných pro eliminaci je NP těžký [1].

2.2.2 Posílání zpráv ve faktor grafu

Algoritmy exaktní inference naráží při reálném použití na příliš velkou složitost násobení faktorů. Hlavním problémem je velikost sdružené pravděpodobnosti, která roste exponenciálně s počtem náhodných proměnných. Většinou nás ovšem zajímá marginální pravděpodobnost jedné nebo jen mála proměnných. Je tedy zbytečné počítat celou sdruženou pravděpodobnost, abychom z ní pak vysčítali většinu proměnných. Řešení se nabízí ve formě faktorizace sdružené pravděpodobnosti.

$$P(X_1, \dots, X_n) = \prod_i P(X_i)$$

Pro výpočet s faktorizovanou distribucí si zavedeme novou datovou strukturu, tzv. faktor graf. Faktor graf je bipartitní graf, kde jednu partitu tvoří faktory a druhou partitu tvoří náhodné proměnné. Hrany ve faktor grafu vedou vždy mezi proměnnou a faktorem, který ji obsahuje.



Obrázek 2.2: Příklad faktor grafu se třemi náhodnými proměnnými X_1, X_2, X_3 a třemi faktory f_a, f_b, f_c .

Nejprve začneme s inferencí na stromech. Pokud zafixujeme jednu náhodnou proměnnou X ve faktor grafu, pak sdruženou pravděpodobnost můžeme spočítat

jako

$$P(\mathbf{X}) = \prod_{s \in ne(X)} F_s(X, \mathbf{X}_s), \quad (2.16) \quad \{\text{eq:fs}\}$$

kde $ne(X)$ jsou faktory obsahující proměnnou X (tedy sousedi ve faktor grafu), F_s je součin všech faktorů v podstromu určeném faktorem f_s a \mathbf{X}_s je množina všech proměnných v daném podstromu.

Pro výpočet marginální pravděpodobnosti X substitujeme (2.16) do výpočtu marginální pravděpodobnosti ze sdružené a po výměně sumy a produktu dostaneme

$$P(X) = \prod_{s \in ne(X)} \sum_{\mathbf{X}_s} F_s(X, \mathbf{X}_s) \quad (2.17)$$

$$= \prod_{s \in ne(X)} \mu_{f_s \rightarrow X}(X) \quad (2.18) \quad \{\text{eq:margx}\}$$

Zavedli jsme funkce

$$\mu_{f_s \rightarrow X} \equiv \sum_{\mathbf{X}_s} F_s(X, \mathbf{X}_s), \quad (2.19) \quad \{\text{eq:dfmsg}\}$$

které můžeme nazývat zprávami z faktoru f_s do proměnné X .

Každé $F_s(X, \mathbf{X}_s)$ je popsáno faktor grafem a tedy může být znova faktori-zováno. Namísto proměnné X nyní vezmeme faktor f_s . Náhodné proměnné sousedící s faktorem f_s bez X si označíme X_1, \dots, X_M . Součin faktorů v podstromech určených těmito proměnnými označíme $G_i(X_i, \mathbf{X}_{si})$. Faktor $F_s(X, \mathbf{X}_s)$ tedy můžeme přepsat jako

$$F_s(X, \mathbf{X}_s) = f_s(X, X_1, \dots, X_M) G_1(X_1, \mathbf{X}_{s1}), \dots, G_M(X_M, \mathbf{X}_{sM}) \quad (2.20)$$

Pokud substitujeme přepsaný faktor $F_s(X, \mathbf{X}_s)$ do definice zprávy z faktoru, dostaneme

$$\mu_{f_s \rightarrow X}(X) = \sum_{X_1} \cdots \sum_{X_M} f_s(X, X_1, \dots, X_M) \prod_{m \in ne(f_s) \setminus X} \sum_{\mathbf{X}_{sm}} G_m(X_m, \mathbf{X}_{sm}) \quad (2.21)$$

$$= \sum_{X_1} \cdots \sum_{X_M} f_s(X, X_1, \dots, X_M) \prod_{m \in ne(f_s) \setminus X} \mu_{X_m \rightarrow f_s}(X_m) \quad (2.22) \quad \{\text{eq:mfsx}\}$$

Zavedli jsme další funkce

$$\mu_{X_m \rightarrow f_s}(X_m) \equiv \sum_{\mathbf{X}_{sm}} G_m(X_m, \mathbf{X}_{sm}), \quad (2.23) \quad \{\text{eq:dfmsg}\}$$

které budeme nazývat zprávami z proměnné X_m do faktoru f_s .

Z rovnice 2.18 vidíme, že marginální pravděpodobnost proměnné vypočítá-

ma jako součin zpráv ze všech okolních faktorů. Každou z těchto zpráv můžeme spočítat jako součin faktoru a zpráv přicházejících z proměnných, které s tímto faktorem sousedí, kromě proměnné, které chceme zprávu posílat. Zbývá nám tedy zjistit v jakém tvaru jsou zprávy z proměnných do faktoru.

Stejně jako u $F_s(X, \mathbf{X}_s)$ nám tady $G_m(X_m, \mathbf{X}_{sm})$ definuje podgraf faktor grafu. V kořeni tohoto podgrafu leží proměnná X_m a tím se dostáváme na už známý případ. Budeme ignorovat faktor f_s , protože neleží v podgrafu určeném $G_m(X_m, \mathbf{X}_{sm})$, a pak

$$G_m(X_m, \mathbf{X}_{sm}) = \prod_{l \in ne(X_m) \setminus f_s} F_l(X_m, \mathbf{X}_{ml}) \quad (2.24)$$

Přepsáno ve formě zpráv

$$\mu_{X_m \rightarrow f_s}(X_m) = \prod_{l \in ne(X_m) \setminus f_s} \sum_{\mathbf{X}_{ml}} F_l(X_m, \mathbf{X}_{ml}) \quad (2.25)$$

$$= \prod_{l \in ne(X_m) \setminus f_s} \mu_{f_l \rightarrow X_m}(X_m) \quad (2.26)$$

Zpráva z náhodné proměnné je tedy součinem zpráv ze všech ostatních faktorů.

Zprávy tedy dokážeme počítat rekurzivně, chybí nám ovšem pravidla pro zprávy z uzlů, které jsou listy. V případě faktoru je odchozí zpráva ekvivalentní faktoru.

$$\mu_{f \rightarrow X}(X) = f(X) \quad (2.27)$$

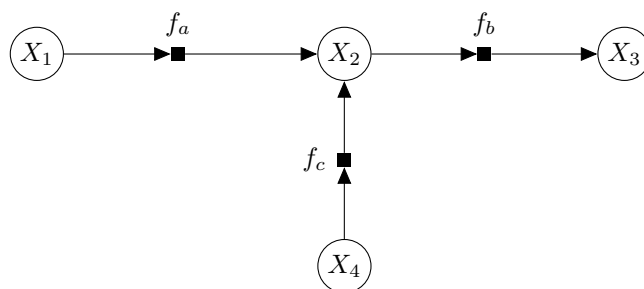
V případě náhodné proměnné je odchozí zpráva

$$\mu_{X \rightarrow f}(X) = 1 \quad (2.28)$$

Pro pozorovanou proměnnou je zpráva vždy rovna pozorované hodnotě, bez ohledu na sousední faktory.

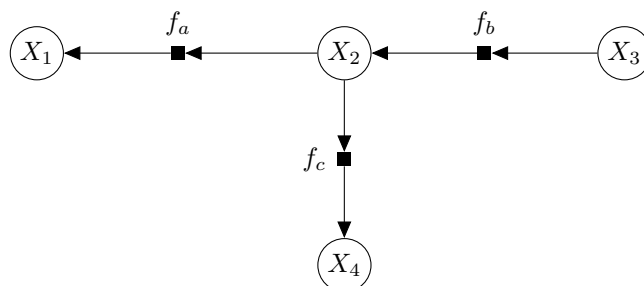
2.2.3 Belief Propagation

Algoritmus exaktní inference na stromech s pomocí posílání zpráv se nazývá belief propagation [13]. Zpráva z vrcholu faktor grafu může být poslána do jiného vrcholu, až když byly obdrženy zprávy ze všech ostatních vrcholů. Jakmile vrchol obdržel zprávy ze všech vrcholů, tak lze spočítat jeho marginální pravděpodobnost. Strom vždy můžeme zakořenit v nějaké náhodné proměnné a posílat zprávy z listů do kořene. Kořen pak obdrží všechny zprávy ze všech sousedních faktorů a tedy je možné spočítat jeho marginální pravděpodobnost. Tomuto se říká dopředný krok propagace.



Obrázek 2.3: Směr posílání zpráv pro strom zakořeněný ve vrcholu X_3 v dopředném kroce.

Belief propagation algoritmus tedy umožňuje na stromech spočítat marginální pravděpodobnost jedné proměnné s lineárním počtem poslaných zpráv. Při rozšíření na spočítání marginální pravděpodobnosti všech proměnných není třeba algoritmus pouštět n -krát. Stačí si uvědomit, že všem vrcholům chybí pouze zpráva od faktoru na cestě ke kořenu. A my můžeme z kořene poslat zprávy zpět k listům hned, jak obdržíme všechny příchozí zprávy. Této části se říká zpětný krok propagace.



Obrázek 2.4: Směr posílání zpráv pro strom zakořeněný ve vrcholu X_3 ve zpětném kroce.

Po propagaci všech zpráv až k listům už každá proměnná získala zprávy od všech sousedních faktorů a tedy je možné spočítat marginální pravděpodobnost všech proměnných.

2.2.4 Loopy Belief Propagation

Algoritmus Belief Propagation funguje pro inferenci na stromech. Hlavní důvod proč nelze použít Belief Propagation na obecných grafech je, že můžeme narazit na cyklus v grafu, což znamená, že žádný z vrcholů v tomto cyklu nebude nikdy mít dostatek příchozích zpráv, aby mohl nějakou zprávu odeslat.

Algoritmus Loopy Belief Propagation řeší tento problém relaxací podmínky na příchozí zprávy. Pro odeslání není třeba znát příchozí zprávy ze všech sousedních vrcholů, kromě toho, kterému je zpráva určena. Chybějící zprávy jsou nastaveny na jedničku. Loopy Belief Propagation je iterativní algoritmus. V grafu nemusí existovat vrchol, který může posílat zprávu k nějakému ze svých sousedů. Je tedy třeba vybrat vrchol podle nějaké strategie a z něj poslat zprávy do všech sousedních vrcholů. V další iteraci je pak vybrán zase jiný vrchol, který bude posílat zprávy. Iterace končí ve chvíli, kdy se už nemění marginální pravděpodobnosti proměnných.

Vzhledem k tomu, že posílané zprávy už neodpovídají skutečným faktorům, je důležitá otázka, zda-li Loopy Belief Propagation vůbec nalezne správné marginální pravděpodobnosti. Existují podmínky za kterých algoritmus bude konvergovat [16]. V obecném případě ovšem může algoritmus konvergovat ke špatným pravděpodobnostem, anebo nemusí konvergovat vůbec a pak dochází k oscilacím. Pokud algoritmus nekonverguje ke správným pravděpodobnostem, tak není ani možné se k nim přiblížit nebo je odvodit z oscilací [12]. V praxi ovšem k problémům většinou nedochází a inference konverguje v rozumném čase.

Popis LBP

Popis Loopy Belief Propagation je v algoritmu (2). Při inferenci ve faktor grafu je nejprve třeba inicializovat všechny zprávy. Posílání zpráv budeme iterovat dokud zprávy nedokonvergují, anebo můžeme nastavit pevný počet iterací. Pak vybíráme vrcholy a z každého pošleme zprávu do všech okolních. Způsobem výběru vrcholů se budeme zabývat později. I kdybychom vybírali vrcholy v náhodném pořadí, tak se můžeme dobrat k výsledku, pouze to bude nejspíš trvat déle.

Poslání zprávy se liší podle toho, zda posíláme z proměnné do faktoru, anebo naopak. V obou případech vynásobíme příchozí zprávy ze všech sousedních vrcholů, kromě toho, do kterého zprávu posíláme. Při posílání zprávy z faktoru zprávy vynásobíme faktorem a marginalizujeme všechny proměnné, kromě té, které zprávu posíláme.

Ve faktor grafu se tedy zprávy šíří pouze přes pravděpodobnostní rozdělení jedné proměnné. To může právě vést k oscilacím anebo konvergenci ke špatným hodnotám. Příkladem může být skrytý markovský model se dvěma skrytými binárními proměnnými a pozorováním pro každou skrytou proměnnou. Pokud je mezi proměnnými závislost XOR, pak nikdy nebudeme pozorovat $(1, 1)$, ale přesto bude každá proměnná konvergovat k rovnoměrnému rozdělení. A tedy sdružená pravděpodobnost pro pozorování $(1, 1)$ dostane nenulovou pravděpodobnost.

Při počítání zpráv často musíme pronásobit spoustu zpráv z okolních vrcholů. Přitom se zprávy zas tak často nemění. Častou optimalizací je pro každý vrchol

Algoritmus 2 Loopy Belief Propagation

{alg:lbp}

```
function LBP( $F, S$ )  
   $F$  – faktor graf  
   $S$  – strategie  
  
  INIT-FACTOR-GRAPH( $F$ )  
  repeat  
    for vrchol  $v \in F$  vybraný podle strategie  $S$  do  
      for soused  $n$  vrcholu  $v$  do  
        if  $v$  je faktor then  
           $\mu_{v \rightarrow n}(n) \leftarrow \text{MESSAGE-TO-VAR}(v, n)$   
        else  
           $\mu_{v \rightarrow n}(v) \leftarrow \text{MESSAGE-TO-FACTOR}(v, n)$   
        end if  
      end for  
    end for  
  until konvergence  
  for každá proměnná  $X$  do  
     $P(X) = \prod_{f \in ne(X)} \mu_{f \rightarrow X}(X)$   
  end for  
end function  
  
function INIT-FACTOR-GRAPH( $F$ )  
   $F$  – faktor graf  
  
  for každou proměnnou  $X \in F$  do  
    for sousední faktor  $f$  do  
       $\mu_{X \rightarrow f}(X) \leftarrow 1$   
       $\mu_{f \rightarrow X}(X) \leftarrow 1$   
    end for  
  end for  
end function  
  
function MESSAGE-TO-VAR( $f, v$ )  
   $f$  – zdrojový faktor  
   $v$  – cílová proměnná  
  
   $\mathbf{X} \leftarrow \text{Dom}(f)$   
   $\mu_{f \rightarrow v}(v) \leftarrow \sum_{\mathbf{X} \setminus v} f(\mathbf{X}) \prod_{u \in ne(f) \setminus v} \mu_{u \rightarrow f}(u)$   
end function  
  
function MESSAGE-TO-FACTOR( $v, f$ )  
   $v$  – zdrojová proměnná  
   $f$  – cílový faktor  
  
   $\mu_{v \rightarrow f}(v) \leftarrow \prod_{g \in ne(v) \setminus f} \mu_{g \rightarrow v}(v)$   
end function
```

si pamatovat součin všech jeho zpráv

$$b(v) = \begin{cases} \prod_{u \in ne(v)} \mu_{u \rightarrow v}(v) & v \text{ je proměnná} \\ f_v(\mathbf{X}) \prod_{u \in ne(v)} \mu_{u \rightarrow v}(u) & v \text{ je faktor} \end{cases} \quad (2.29)$$

Pak můžeme při odesílání zprávy z vrcholu v do vrcholu u pouze $b(v)$ vydělit zprávou $\mu_{u \rightarrow v}$, pro faktory marginalizovat, a zpráva je připravena k odeslání. $b(v)$ můžeme aktualizovat po výběru vrcholu v pro posílání zpráv, anebo po každém přijetí zprávy.

2.2.5 Strategie výběru vrcholu v LBP algoritmu

Inference v grafu se může lišit podle typu faktor grafu, ale také podle nároků, které na výsledek máme. Nejjednodušší metodou pro výběr je nechat pořadí na uživateli algoritmu. Pro stromy máme speciální strategii, která zaručí konvergenci po jednom dopředném a jednom zpětném kroku propagace. V dialogových systémech je často používána dynamická bayesovská síť a nás zajímá pravděpodobnost proměnných v poslední vrstvě sítě. V takovém případě můžeme některé zprávy zanedbat, protože už příliš neovlivní proměnné, které nás zajímají.

Inference na stromě

Pro efektivní inferenci je třeba si pro každý vrchol pamatovat, kolik mu chybí zpráv, aby mohl jednu sám odeslat. Pro každý vrchol v s k sousedy je na začátku počet chybějících zpráv $k - 1$. Strom na alespoň dvou vrcholech obsahuje alespoň dva listy. Budeme tedy postupně odebírat vrcholy, jejichž počet chybějících zpráv je nulový. Rozešleme z nich zprávu do souseda, z kterého ještě zpráva nepřišla a snížíme jeho počet chybějících zpráv. Odebráním listu ze stromu vždy dostaneme opět strom. Postupně tak zmenšujeme strom, až dostaneme právě jeden vrchol, který získal všechny zprávy.

Po vypočtení marginální pravděpodobnosti ve stromě o jednom vrcholu můžeme zase přidávat vrcholy v obráceném pořadí než v jakém jsme je odebírali. Do každého přidaného vrcholu pak můžeme poslat zprávu a spočítat jeho marginální pravděpodobnost.

Inference v dynamické bayesovské síti

Dynamická bayesovská síť je nejčastější reprezentace dialogového stavu. V jedné vrstvě sítě je popsána jedna obrátka dialogu. Vrcholy v jedné vrstvě většinou závisí pouze na vrcholech ve stejné nebo předchozí vrstvě. Po každé obrátce nás

pro účely dialogového manageru zajímají hlavně pravděpodobnosti proměnných v poslední vrstvě, tedy po aktuální obrátce.

Inferenci provádíme způsobem podobným indukci. Pro první vrstvu provedeme inferenci libovolným způsobem, může dokonce platit, že v rámci jedné vrstvy se jedná o strom. Po přidání k -té vrstvy předpokládáme, že byla provedena inference na předchozích $k - 1$ vrstvách a tedy zprávy v této části grafu jsou správné. Zprávy ve směru k nové vrstvě jsou stále správné, neznáme ovšem hodnotu zpráv z předposlední vrstvy do nové vrstvy a hodnotu zpráv v nové vrstvě. Pošleme tedy zprávy z předposlední vrstvy a provedeme inferenci v nové vrstvě, opět můžeme použít libovolnou heuristiku.

Stále ještě zbývají zprávy z nové vrstvy zpět v síti. Čím dál do historie ovšem jdeme, tím menší vliv naše nová pozorování budou mít na vrcholy v dané vrstvě. Proto se většinou omezíme jen na posledních několik vrstev (1 až 3).

Implementace Loopy Belief Algoritmu i se strategiemi pro výběr vrcholů bude popsána v kapitole ??.

2.3 Propagace s aproximovanými zprávami

V předchozí sekci byla faktorizována sdružená pravděpodobnost a tak došlo k zjednodušení inference. Mezi jednotlivými vrcholy byly stále posílány posílány exaktní zprávy. V této sekci se zaměříme na inferenci v modelech, kde není možné spočítat zprávy exaktně a je třeba je aproximovat. Příkladem jsou například modely se spojitými náhodnými proměnnými.

Předpokládáme, že máme pravděpodobnostní grafický model, reprezentující sdruženou pravděpodobnost dat a pozorování pomocí součinu faktorů.

$$P(\mathbf{X}, \mathbf{E}) = \prod_i f_i(\mathbf{X}) \quad (2.30)$$

Nás zajímá aposteriorní distribuce $P(\mathbf{X} \mid \mathbf{E})$ pro zjištění stavu, stejně tak jako pravděpodobnost pozorování $P(\mathbf{E})$ pro normalizaci. Aposteriorní distribuci vyjádříme ze sdružené.

$$P(\mathbf{X} \mid \mathbf{E}) = \frac{1}{P(\mathbf{E})} \prod_i f_i(\mathbf{X}) \quad (2.31)$$

a pravděpodobnost pozorování je dána

$$P(\mathbf{E}) = \int \prod_i f_i(\mathbf{X}) d\mathbf{X} \quad (2.32)$$

Pro aproximativní inferenci ukážeme algoritmus Expectation propagation (EP) [11]. Vytvoříme aproximaci aposteriorní distribuce, která je také dána součinem fak-

torů

$$q(\mathbf{X}) = \frac{1}{Z} \prod_i \tilde{f}_i(\mathbf{X}), \quad (2.33)$$

kde každý faktor \tilde{f}_i je aproximace odpovídající skutečnému faktoru f_i a faktor $\frac{1}{Z}$ je normalizační konstanta. Aproximované faktory musíme zkonstruovat tak, abychom byli schopni provádět inferenci.

Pro měření vzdálenosti aproximovaného rozdělení od skutečného používáme Kullback-Leiblerovu divergenci (KL) [9]. KL divergence, také známá jako relativní entropie, mezi dvěma pravděpodobnostními rozděleními $p(x)$ a $q(x)$ je

$$KL(p\|q) = \int p(x) \frac{p(x)}{q(x)} dx. \quad (2.34)$$

Divergence splňuje tři vlastnosti:

1. $KL(p\|p) = 0$,
2. $KL(p\|q) = 0$ právě tehdy když $p = q$,
3. $KL(p\|q) > 0$ pro všechna p, q .

2.3.1 Rodina exponenciálních rozdělení

Rodina exponenciálních rozdělení [2] přes \mathbf{X} , dáno parametry $\boldsymbol{\eta}$ je definována jako množina distribucí ve tvaru

$$P(\mathbf{X} | \boldsymbol{\eta}) = h(\mathbf{X})g(\boldsymbol{\eta}) \exp(\boldsymbol{\eta}^T u(\mathbf{X})), \quad (2.35) \quad \{\text{eq:gj}\}$$

kde \mathbf{X} jsou náhodné proměnné, které mohou být diskrétní nebo spojitě, $\boldsymbol{\eta}$ jsou nazývány přirozené parametry rozdělení, $u(\mathbf{X})$ je funkce \mathbf{X} . Funkce $g(\boldsymbol{\eta})$ může být interpretována jako koeficient, který zajišťuje, že je distribuce normalizována a tedy splňuje

$$g(\boldsymbol{\eta}) \int h(\mathbf{X}) \exp(\boldsymbol{\eta}^T u(\mathbf{X})) d\mathbf{X} = 1 \quad (2.36) \quad \{\text{eq:g1}\}$$

Výhody distribucí z exponenciální rodiny si ukážeme na KL divergenci. Necht' počítáme divergenci $KL(p\|q)$, kde $p(\mathbf{X})$ je zafixovaná distribuce a $q(\mathbf{X})$ je distribuce z exponenciální rodiny. Pokud zapíšeme KL divergenci jako funkci $\boldsymbol{\eta}$, tak dostaneme

$$KL(p\|q) = -\log g(\boldsymbol{\eta}) - \boldsymbol{\eta}^T \mathbb{E}_{p(\mathbf{X})}[u(\mathbf{X})] + \text{const}. \quad (2.37)$$

KL divergenci můžeme minimalizovat tak, že nastavíme první derivaci podle $\boldsymbol{\eta}$ rovnou nule, z čehož dostaneme

$$-\nabla \log g(\boldsymbol{\eta}) = \mathbb{E}_{p(\mathbf{X})}[u(\mathbf{X})]. \quad (2.38) \quad \{\text{eq:k1}\}$$

Nyní se musíme podívat na $-\nabla \log g(\boldsymbol{\eta})$. Vezmeme derivaci obou stran (2.36) podle $\boldsymbol{\eta}$. Získáme

$$\begin{aligned} \nabla g(\boldsymbol{\eta}) \int h(\mathbf{X}) \exp(\boldsymbol{\eta}^T u(\mathbf{X})) d\mathbf{X} \\ + g(\boldsymbol{\eta}) \int h(\mathbf{X}) \exp(\boldsymbol{\eta}^T u(\mathbf{X})) u(\mathbf{X}) d\mathbf{X} = 0 \end{aligned} \quad (2.39)$$

Přeuspořádáním a znova použitím (2.35) a (2.36) získáme

$$-\frac{1}{g(\boldsymbol{\eta})} \nabla g(\boldsymbol{\eta}) = \int h(\mathbf{X}) \exp(\boldsymbol{\eta}^T u(\mathbf{X})) u(\mathbf{X}) d\mathbf{X} = \mathbb{E}[u(\mathbf{X})]. \quad (2.40)$$

Ve výsledku získáváme

$$-\nabla \log g(\boldsymbol{\eta}) = \mathbb{E}[u(\mathbf{X})] \quad (2.41) \quad \{\text{eq:ge}\}$$

Nový poznatek o funkci $g(\boldsymbol{\eta})$ z (2.41) můžeme substituovat do (2.38) a získáme

$$\mathbb{E}_{q(\mathbf{X})}[u(\mathbf{X})] = \mathbb{E}_{p(\mathbf{X})}[u(\mathbf{X})], \quad (2.42) \quad \{\text{eq:ee}\}$$

z čehož vidíme, že pro minimalizaci KL divergence nám stačí najít parametry rozdělení $q(\mathbf{X})$ tak, aby mělo stejné očekávané statistiky jako $p(\mathbf{X})$.

2.3.2 Expectation Propagation

Expectation propagation je algoritmus založený na aproximaci aposteriorní distribuce součinem aproximovaných faktorů

$$q(\mathbf{X}) = \frac{1}{Z} \prod_i \tilde{f}_i(\mathbf{X}). \quad (2.43)$$

Z praktických důvodů uvedených v předchozí sekci jsou aproximované faktory z exponenciální rodiny. Díky tomu bude jejich součin také z exponenciální rodiny.

Aproximující distribuci bychom chtěli nalézt pomocí minimalizace KL divergence mezi skutečnou a aproximovanou distribucí.

$$KL(p||q) = KL\left(\frac{1}{p(\mathbf{E})} \prod_i f_i(\mathbf{X}) \left\| \frac{1}{Z} \prod_i \tilde{f}_i(\mathbf{X})\right.\right) \quad (2.44)$$

Pro počítání KL divergence ovšem potřebujeme umět počítat se skutečnou pravděpodobností a když už s ní umíme efektivně počítat, tak proč ztrácet čas s aproximací.

Další možností je aproximovat jednotlivé faktory, díky tomu bychom v jednom kroku našli aproximace pro všechny faktory a měli bychom hotovo. Naším cílem je ale nalézt nejlepší aproximaci celé aposteriorní pravděpodobnosti a to není v

tomto případě zaručeno.

Expectation propagation [4] sice aproximuje jednotlivé faktory, ale vždy v kontextu všech ostatních. Nejprve jsou všechny inicializovány a pak jsou procházeny jeden po druhém a každý je aktualizován. Tento přístup je podobný metodám, které byly ukázány v předešlých sekcích. Předpokládejme, že chceme aktualizovat faktor $\tilde{f}_j(\mathbf{X})$. Nejprve odstraníme faktor z produktu $\prod_{i \neq j} \tilde{f}_i(\mathbf{X})$. Následně nalezneme novou hodnotu pro faktor \tilde{f}_j tak, aby pravděpodobnostní rozložení

$$q^{new}(\mathbf{X}) \propto \tilde{f}_j(\mathbf{X}) \prod_{i \neq j} \tilde{f}_i(\mathbf{X}) \quad (2.45) \quad \{\text{eq:qnew}\}$$

bylo co nejbližší

$$f_j(\mathbf{X}) \prod_{i \neq j} \tilde{f}_i(\mathbf{X}). \quad (2.46)$$

Všechny faktory $i \neq j$ necháváme zafixované. Díky tomu zajistíme, že nová aproximace faktoru je nejpřesnější v oblastech, které mají největší aposteriorní pravděpodobnost definovanou zbývajících faktory.

Odstranění faktoru $\tilde{f}_j(\mathbf{X})$ z aktuální aproximace aposteriorní distribuce provedeme vytvořením nenormalizované distribuce

$$q^{\setminus j}(\mathbf{X}) = \frac{q(\mathbf{X})}{\tilde{f}_j(\mathbf{X})} \quad (2.47)$$

Mohli bychom počítat $q^{\setminus j}(\mathbf{X})$ jako součin všech faktorů kromě j , ale v praxi je dělení rychlejší. Této nenormalizované distribuci se říká cavity distribuce. Její kombinací se skutečným faktorem $f_j(\mathbf{X})$ dostaneme distribuci

$$\frac{1}{Z_j} f_j(\mathbf{X}) q^{\setminus j}(\mathbf{X}), \quad (2.48) \quad \{\text{eq:aprox}\}$$

kde Z_j je normalizační konstanta, dána

$$Z_j = \int f_j(\mathbf{X}) q^{\setminus j}(\mathbf{X}) d\mathbf{X}. \quad (2.49) \quad \{\text{eq:zj}\}$$

Nyní můžeme nalézt novou hodnotu pro faktor \tilde{f}_j minimalizací KL divergence

$$KL\left(\frac{f_j(\mathbf{X}) q^{\setminus j}(\mathbf{X})}{Z_j} \parallel q^{new}(\mathbf{X})\right). \quad (2.50)$$

To uděláme jednoduše použitím poznatku z (2.42), který říká, že nám stačí najít parametry q^{new} tak, aby jeho postačující statistiky odpovídali momentům aproximovaného rozdělení (2.48).

Z nalezené aproximace q^{new} pak získáme novou hodnotu faktoru \tilde{f}_j z (2.45)

vydělením zbývajících faktorů

$$\tilde{f}_j(\mathbf{X}) = K \frac{q^{new}(\mathbf{X})}{q^{j}(\mathbf{X})} \quad (2.51) \quad \{\text{eq:newfj}\}$$

Koeficient K získáme tak, že obě strany (2.51) vynásobíme $q^{j}(\mathbf{X})$ a zintegrováním, čímž získáme

$$K = \int \tilde{f}_j(\mathbf{X}) q^{j}(\mathbf{X}) d\mathbf{X}, \quad (2.52)$$

použili jsme navíc toho, že $q^{new}(\mathbf{X})$ je normalizovaná distribuce. Hodnota K pak může být nalezena srovnáním

$$\int \tilde{f}_j(\mathbf{X}) q^{j}(\mathbf{X}) d\mathbf{X} = \int f_j(\mathbf{X}) q^{j}(\mathbf{X}) d\mathbf{X} \quad (2.53)$$

Z čehož zjistíme, že $K = Z_j$ a tedy může být nalezeno přímo z (2.49).

V praxi je třeba provést několik iterací, a v každé aktualizovat všechny faktory.

Algoritmus 3 Expectation propagation

3. Učení parametrů

3.1 Grafický model

Máme vybraný faktor f , tento faktor je spojený s několika proměnnými $\mathbf{x} = (x_0, x_1, \dots, x_{N_x})$ a množinami parametrů $\Theta = (\theta_1, \dots, \theta_{N_\theta})$. Tento faktor reprezentuje podmíněnou pravděpodobnost:

$$f(\mathbf{x}, \Theta) = p(x_0 | x_1, \dots, x_{N_x}; \Theta)$$

Rodičovské proměnné x_1, \dots, x_{N_x} označujeme jako \mathbf{x}' . Vektor \mathbf{x}' určuje, která množina parametrů bude použita. Protože množiny parametrů jsou číslovány $1, \dots, N_\theta$ a rodičovské proměnné $1, \dots, N_x$, musí být pro vybrání správné množiny parametrů použito mapování $\rho(\mathbf{x}')$. Faktor pak může být zapsán zkráceně:

$$f(\mathbf{x}, \Theta) = p(x_0 | x_1, \dots, x_{N_x}; \Theta) = \theta_{\rho(\mathbf{x}'), x_0}$$

3.2 Výpočet marginálních pravděpodobností

Pro výpočet sdružené pravděpodobnosti používáme plně faktorizovanou distribuci. Pro každou proměnnou anebo množinu parametrů je její marginální pravděpodobnost rovna součinu zpráv přicházejících z faktorů, které jsou s danou proměnnou nebo množinou parametrů propojeny. Pro daný faktor je cavity distribuce $q^\setminus(x_i)$, popř. $q^\setminus(\theta_i)$ rovna součinu zpráv ze všech ostatních faktorů do x_i , popř. θ_i . Aproximovaná marginální pravděpodobnost proměnné je pak součinem cavity distribuce a zprávy z faktoru:

$$q(x_i) = q^\setminus(x_i) m_{f \rightarrow x_i}(x_i)$$

$$q(\theta_i) = q^\setminus(\theta_i) m_{f \rightarrow \theta_i}(\theta_i)$$

Cavity distribuce je právě zpráva z proměnné, popř. množiny parametrů do faktoru.

$$m_{x_i \rightarrow f} = q^\setminus(x_i)$$

$$m_{\theta_i \rightarrow f} = q^\setminus(\theta_i)$$

3.2.1 Marginální pravděpodobnost proměnných

Pokud chceme aktualizovat hodnotu naší aproximace marginální pravděpodobnosti, tak je třeba minimalizovat její vzdálenost od skutečné marginální pravděpo-

dobnosti:

$$p^*(\tilde{x}_j) = \sum_{\mathbf{x}:x_j=\tilde{x}_j} \int_{\Theta} \prod_i q_{\setminus}(x_i) \prod_l q_{\setminus}(\theta_l) f(\mathbf{x}; \Theta) \quad (3.1) \quad \{\text{eq:1}\}$$

$$= \sum_{\mathbf{x}:x_j=\tilde{x}_j} \prod_i q_{\setminus}(x_i) \int_{\theta_{\rho(\mathbf{x}'),x_0}} q_{\setminus}(\theta_{\rho(\mathbf{x}'),x_0}) \theta_{\rho(\mathbf{x}'),x_0} \quad (3.2) \quad \{\text{eq:2}\}$$

$$= \sum_{\mathbf{x}:x_j=\tilde{x}_j} \prod_i q_{\setminus}(x_i) \mathbb{E}_{q_{\setminus}}(\theta_{\rho(\mathbf{x}'),x_0}) \quad (3.3) \quad \{\text{eq:3}\}$$

$$= \sum_{\mathbf{x}:x_j=\tilde{x}_j} \prod_i m_{x_i \rightarrow f}(x_i) \mathbb{E}_{q_{\setminus}}(\theta_{\rho(\mathbf{x}'),x_0}) \quad (3.4)$$

Rovnost (3.1) vychází z definice výpočtu marginální pravděpodobnosti ze sdružené pravděpodobnosti. V (3.2) byla použita definice faktoru, z integrálu byly vytaženy členy, které neobsahují Θ a nakonec bylo využito toho, že pro množiny parametrů, které nejsou spojeny s faktorem f , je jejich jejich cavity distribuce rovná marginální distribuci a tedy $\int_{\theta_i} q(\theta_i) = 1$. V (3.3) byla použita definice očekávané hodnoty.

Marginální pravděpodobnost proměnné x_i tedy je

$$p^*(\tilde{x}) = \sum_{\mathbf{x}:x_j=\tilde{x}_j} \prod_i m_{x_i \rightarrow f}(x_i) \mathbb{E}_{q_{\setminus}}(\theta_{\rho(\mathbf{x}'),x_0}) \quad (3.5)$$

Tady docházíme k výsledku, který je velmi podobný výpočtu marginální pravděpodobnosti v Loopy Belief Propagation algoritmu.

Zprávu z faktoru f do vrcholu x_j pak získáme vydělením zprávy z x_j z marginální pravděpodobnosti.

$$m_{f \rightarrow x_j}(x_j) = \sum_{\mathbf{x}:x_j=\tilde{x}_j} \prod_{i \neq j} m_{x_i \rightarrow f}(x_i) \mathbb{E}_{q_{\setminus}}(\theta_{\rho(\mathbf{x}'),x_0}) \quad (3.6) \quad \{\text{eq:msgfrom}\}$$

3.2.2 Marginální pravděpodobnost parametrů

Pro množiny parametrů se jejich marginální pravděpodobnost spočítá podobně jako pro proměnné.

$$p^*(\tilde{\boldsymbol{\theta}}_j) = \sum_{\mathbf{x}} \int_{\boldsymbol{\Theta}:\boldsymbol{\theta}_j=\tilde{\boldsymbol{\theta}}_j} \prod_i q^{\setminus}(x_i) \prod_l q^{\setminus}(\boldsymbol{\theta}_l) f(\mathbf{x}; \boldsymbol{\Theta}) \quad (3.7) \quad \{\text{eq:ep:the}\}$$

$$= \sum_{l \neq j} \sum_{\mathbf{x}:\rho(\mathbf{x}')=l} \prod_i q^{\setminus}(x_i) \int_{\boldsymbol{\Theta}:\boldsymbol{\theta}_j=\tilde{\boldsymbol{\theta}}_j} \prod_k q^{\setminus}(\boldsymbol{\theta}_k) \theta_{l,x_0} + \quad (3.8) \quad \{\text{eq:ep:the}\}$$

$$+ \sum_{\mathbf{x}:\rho(\mathbf{x}')=j} \prod_i q^{\setminus}(x_i) \int_{\boldsymbol{\Theta}:\boldsymbol{\theta}_j=\tilde{\boldsymbol{\theta}}_j} \prod_k q^{\setminus}(\boldsymbol{\theta}_k) \tilde{\theta}_{j,x_0} \\ = \left[\sum_{l \neq j} \sum_{\mathbf{x}:\rho(\mathbf{x}')=l} \prod_i q^{\setminus}(x_i) \mathbb{E}_{q^{\setminus}(\boldsymbol{\theta}_l)}(\theta_{l,x_0}) \right] q^{\setminus}(\tilde{\boldsymbol{\theta}}_j) + \quad (3.9) \quad \{\text{eq:ep:the}\}$$

$$+ \sum_{\mathbf{x}:\rho(\mathbf{x}')=j} \prod_i q^{\setminus}(x_i) \tilde{\theta}_{j,x_0} q^{\setminus}(\tilde{\boldsymbol{\theta}}_j) \\ = w_0 q^{\setminus}(\tilde{\boldsymbol{\theta}}_j) + \sum_k w_k \tilde{\theta}_{j,k} q^{\setminus}(\tilde{\boldsymbol{\theta}}_j), \quad (3.10) \quad \{\text{eq:ep:the}\}$$

$$= w_0 m_{\tilde{\boldsymbol{\theta}}_j \rightarrow f}(\tilde{\boldsymbol{\theta}}_j) + \sum_k w_k \tilde{\theta}_{j,k} m_{\tilde{\boldsymbol{\theta}}_j \rightarrow f}(\tilde{\boldsymbol{\theta}}_j), \quad (3.11)$$

kde

$$w_0 = \sum_{l \neq j} \sum_{\mathbf{x}:\rho(\mathbf{x}')=l} \prod_i m_{x_i \rightarrow f}(x_i) \mathbb{E}_{q^{\setminus}(\boldsymbol{\theta}_l)}(\theta_{l,x_0}) \quad (3.12)$$

$$w_k = \sum_{\mathbf{x}:\rho(\mathbf{x}')=j, x_0=k} \prod_i m_{x_i \rightarrow f}(x_i) \quad (3.13)$$

Opět vycházíme z výpočtu marginální pravděpodobnosti ze sdružené pravděpodobnosti. V rovnici (3.8) jsme rozdělili sumu přes \mathbf{x} na ty, pro které se ve faktoru použije množina parametrů $\tilde{\boldsymbol{\theta}}_j$ a na ty ostatní. Také jsme z integrálu vytknuli součin cavity distribucí pro proměnné. V dalším kroku (3.9) jsme opět použili toho, že integrál přes $\boldsymbol{\Theta}$ je ve skutečnosti několik integrálů přes jednotlivé množiny parametrů. A tedy je můžeme vložit mezi jednotlivé členy produktu cavity distribucí pro množiny parametrů. Ve výsledku získáme $q^{\setminus}(\tilde{\boldsymbol{\theta}}_j) \int_{\boldsymbol{\theta}_l} q^{\setminus}(\boldsymbol{\theta}_l) \theta_{l,x_0}$ a pak zbylé členy, které zmizí.

Docházíme k vyjádření skutečné marginální pravděpodobnosti, ve které není třeba integrovat přes všechny množiny parametrů, ale stačí jen očekávaná hodnota těchto parametrů.

3.3 Aproximace marginálních pravděpodobností

Stále tu ovšem zůstává problém, že spočítat aproximující distribuci $q(\theta_j)$ může být příliš složité, protože skutečná marginální distribuce je směs několika distribucí a ta nemusí být v obecném případě vyjádřitelná. Je tedy třeba model dále aproximovat. Pro zjednodušení výpočtu jsou zprávy z faktoru do množiny parametrů, $m_{f \rightarrow \theta_i}(\theta_i)$, ve tvaru Dirichletovského rozdělení s parametry $\alpha_{f \rightarrow \theta_i}$:

$$m_{f \rightarrow \theta_i}(\theta_i) = \text{Dir}(\theta_i; \alpha_{f \rightarrow \theta_i}) = \frac{\Gamma(\sum_j \alpha_{f \rightarrow \theta_i, j})}{\prod_j \Gamma(\alpha_{f \rightarrow \theta_i, j})} \prod_j \theta_{i, j}^{\alpha_{f \rightarrow \theta_i, j} - 1} \quad (3.14)$$

kde Γ je Gamma funkce (zobecnění faktoriálu):

$$\Gamma(z) = \int_0^\infty t^{z-1} \exp(-t) dt \quad (3.15)$$

Dirichletovské rozdělení bylo zvoleno, protože má důležité vlastnosti pro součin, které budou využity dále pro výpočet cavity distribuce a celkové aproximace. Pokud označíme aproximované faktory indexem β a každý bude mít vlastní parametry $\alpha_{f \rightarrow \theta_i}$, tak výsledná aproximace bude tvaru:

$$q(\theta_i) \propto \prod_{\beta} m_{f \rightarrow \theta_i}(\theta_i) \quad (3.16)$$

$$\propto \prod_{\beta} \prod_j \theta_{i, j}^{\alpha_{f \rightarrow \theta_i, j} - 1} \quad (3.17)$$

$$\propto \text{Dir}(\theta_i; \sum_{\beta} \alpha_{f \rightarrow \theta_i} - (|\beta| - 1)\mathbf{1}) \quad (3.18)$$

$$= \text{Dir}(\theta_i; \alpha_i) \quad (3.19) \quad \{\text{eq:ep:aprx}\}$$

kde $\alpha_i = \sum_{\beta} \alpha_{f \rightarrow \theta_i} - (|\beta| - 1)\mathbf{1}$.

Při aktualizaci faktoru $\tilde{\beta}$ tedy cavity distribuce bude:

$$q^{\setminus \tilde{\beta}}(\theta_i) \propto \prod_{\beta \neq \tilde{\beta}} m_{f \rightarrow \theta_i}(\theta_i) \quad (3.20)$$

$$\propto \text{Dir}(\theta_i; \alpha_i - \alpha_{f \rightarrow \theta_i} + \mathbf{1}) \quad (3.21)$$

Nášim cílem je nalézt parametry α^* aproximované marginální pravděpodobnosti (3.19), které minimalizují vzdálenost od skutečné marginální pravděpodobnosti (3.10). Pro měření vzdálenosti mezi dvěma pravděpodobnostními rozloženími se používá Kullback-Leiblerova divergence:

$$KL(p||q) = \int_{-\infty}^{\infty} p(x) \log \left(\frac{p(x)}{q(x)} \right) dx \quad (3.22)$$

Pro nalezení minima použijeme algoritmus Expectation Propagation a budeme tedy minimalizovat $KL(p^*||q)$.

Pokud se podíváme na skutečnou marginální pravděpodobnost $p^*(\theta_i)$, zjistíme, že můžeme některé její členy upravit. Využijeme také vlastnosti gamma funkce $\Gamma(x) = (x-1)\Gamma(x-1)$.

$$w_j \theta_j Dir(\theta; \alpha) \propto w_j \theta_j \frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \prod_i \theta_i^{\alpha_i-1} \quad (3.23) \quad \{\text{eq:4}\}$$

$$\propto w_j \frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \theta_j^{\alpha_j} \prod_{i \neq j} \theta_i^{\alpha_i-1} \quad (3.24)$$

$$\propto w_j \frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \frac{\Gamma(\alpha_j + 1) \prod_{i \neq j} \Gamma(\alpha_i)}{\Gamma(1 + \sum_i \alpha_i)} Dir(\theta; \alpha + \delta_j) \quad (3.25)$$

$$\propto w_j \frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \frac{\alpha_j \Gamma(\alpha_j) \prod_{i \neq j} \Gamma(\alpha_i)}{(\sum_i \alpha_i) \Gamma(\sum_i \alpha_i)} Dir(\theta; \alpha + \delta_j) \quad (3.26)$$

$$\propto w_j \frac{\alpha_j}{\sum_i \alpha_i} Dir(\theta; \alpha + \delta_j) \quad (3.27)$$

$$(3.28)$$

Díky této úpravě lze p^* vyjádřit jako směs Dirichletovských rozdělení.

$$p^*(\theta) = w_0^* Dir(\theta; \alpha) + \sum_j w_j^* Dir(\theta; \alpha + \delta_j) \quad (3.29) \quad \{\text{eq:margin}\}$$

kde

$$w_0^* \propto w_0 \quad (3.30)$$

$$w_j^* \propto w_j \frac{\alpha_j}{\sum_i \alpha_i} \quad (3.31)$$

$$\sum_{i=0}^k w_i^* = 1 \quad (3.32)$$

Pro minimalizaci KL divergence mezi dvěma rozděleními z exponenciální rozdělení stačí, pokud se budou rovnat jejich postačující statistiky. Dokážeme jednoduše spočítat první dva momenty Dirichletovského rozdělení a tedy použijeme aproxi-maci a budeme počítat pouze s nimi a zbylé momenty zanedbáme. Je tedy třeba nalézt střední hodnotu a rozptyl proměnných z $p^*(\theta)$.

$$\mathbb{E}_{p^*}[\boldsymbol{\theta}] = \int \boldsymbol{\theta} p^*(\boldsymbol{\theta}) d\boldsymbol{\theta} \quad (3.33)$$

$$= \int \boldsymbol{\theta} (w_0^* \text{Dir}(\boldsymbol{\theta}; \boldsymbol{\alpha}) + \sum_j w_j^* \text{Dir}(\boldsymbol{\theta}; \boldsymbol{\alpha} + \boldsymbol{\delta}_j)) d\boldsymbol{\theta} \quad (3.34)$$

$$= w_0^* \int \boldsymbol{\theta} \text{Dir}(\boldsymbol{\theta}; \boldsymbol{\alpha}) d\boldsymbol{\theta} + \sum_j w_j^* \int \boldsymbol{\theta} \text{Dir}(\boldsymbol{\theta}; \boldsymbol{\alpha} + \boldsymbol{\delta}_j) d\boldsymbol{\theta} \quad (3.35)$$

$$= w_0^* \mathbb{E}_{\text{Dir}(\boldsymbol{\alpha})}[\boldsymbol{\theta}] + \sum_j w_j^* \mathbb{E}_{\text{Dir}(\boldsymbol{\alpha} + \boldsymbol{\delta}_j)}[\boldsymbol{\theta}] \quad (3.36)$$

Střední hodnotu proměnných $\boldsymbol{\theta}$ podle rozdělení p^* lze tedy spočítat jako vážený součet středních hodnot $\boldsymbol{\theta}$ podle jednotlivých Dirichletovských distribucí, z kterých se p^* skládá. Střední hodnota proměnné X_i podle Dirichletovského rozdělení je

První moment tedy máme spočítaný, pro výpočet rozptylu můžeme využít přímo definici:

$$\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 \quad (3.37)$$

Chybí nám tedy ještě výpočet střední hodnoty druhé mocniny proměnné $\boldsymbol{\theta}$ podle p^* . Můžeme ji vyjádřit z definice střední hodnoty.

$$\mathbb{E}_{p^*}[\boldsymbol{\theta}^2] = \int \boldsymbol{\theta}^2 p^*(\boldsymbol{\theta}) d\boldsymbol{\theta} \quad (3.38)$$

$$= w_0^* \int \boldsymbol{\theta}^2 \text{Dir}(\boldsymbol{\theta}; \boldsymbol{\alpha}) d\boldsymbol{\theta} + \sum_j w_j^* \int \boldsymbol{\theta}^2 \text{Dir}(\boldsymbol{\theta}; \boldsymbol{\alpha} + \boldsymbol{\delta}_j) d\boldsymbol{\theta} \quad (3.39)$$

$$= w_0^* \mathbb{E}_{\text{Dir}(\boldsymbol{\alpha})}[\boldsymbol{\theta}^2] + \sum_j w_j^* \mathbb{E}_{\text{Dir}(\boldsymbol{\alpha} + \boldsymbol{\delta}_j)}[\boldsymbol{\theta}^2] \quad (3.40)$$

Opět získáváme vážený součet středních hodnot podle Dirichletovských rozdělení. Střední hodnotu druhé mocniny proměnné podle Dirichletovského rozdělení lze opět jednoduše odvodit z definice.

$$\mathbb{E}_{\text{Dir}(\boldsymbol{\alpha})}[x_i^2] = \int x_i^2 \text{Dir}(\mathbf{x}; \boldsymbol{\alpha}) d\mathbf{x} \quad (3.41)$$

$$= \int x_i^2 \frac{\Gamma(\alpha_0)}{\prod_{j=1}^N \Gamma(\alpha_j)} \prod_{j=1}^N x_j^{\alpha_j-1} d\mathbf{x} \quad (3.42)$$

Nyní jsme v podobné situaci jako v (3.23). Budeme postupovat stejně, vyjádříme nové Dirichletovské rozdělení.

$$\mathbb{E}_{Dir(\boldsymbol{\alpha})}[x_i^2] = \int \frac{\Gamma(\alpha_0 + 2)\alpha_i(\alpha_i + 1)}{\alpha_0(\alpha_0 + 1)\Gamma(\alpha_i + 2) \prod_{j \neq i} \Gamma(\alpha_j)} x_i^{\alpha_i+1} \prod_{j \neq i} x_j^{\alpha_j-1} d\mathbf{x} \quad (3.43)$$

$$= \frac{\alpha_i(\alpha_i + 1)}{\alpha_0(\alpha_0 + 1)} \int \frac{\Gamma(\beta_0)}{\prod_i \Gamma(\beta_i)} \prod_i x_i^{\beta_i-1} d\mathbf{x} \quad (3.44)$$

$$= \frac{\alpha_i(\alpha_i + 1)}{\alpha_0(\alpha_0 + 1)} \int Dir(\mathbf{x}; \boldsymbol{\beta}) d\mathbf{x} \quad (3.45)$$

$$= \frac{\alpha_i(\alpha_i + 1)}{\alpha_0(\alpha_0 + 1)} \quad (3.46)$$

Vyjádřili jsme $\Gamma(\alpha_0)$ a $\Gamma(\alpha_i)$ s pomocí $\Gamma(\alpha_0 + 2)$ a $\Gamma(\alpha_i + 2)$

$$\Gamma(\alpha_0) = \frac{\Gamma(\alpha_0 + 2)}{\alpha_0(\alpha_0 + 1)} \quad (3.47)$$

$$\Gamma(\alpha_i) = \frac{\Gamma(\alpha_i + 2)}{\alpha_i(\alpha_i + 1)} \quad (3.48)$$

Následně jsme vytvořili nové parametry $\boldsymbol{\beta}$:

$$\beta_i = \alpha_i + 2 \quad (3.49)$$

$$\beta_{j \neq i} = \alpha_j \quad (3.50)$$

$$\beta_0 = \sum_i \beta_i \quad (3.51)$$

Nyní tedy dokážeme spočítat $\mathbb{E}_{p^*}[\boldsymbol{\theta}]$ a $\mathbb{E}_{p^*}[\boldsymbol{\theta}^2]$. Parametry aproximovaného rozdělení nalezneme následovně

$$\frac{\mathbb{E}[X_1] - \mathbb{E}[X_1^2]}{\mathbb{E}[X_1^2] - \mathbb{E}[X_1]^2} = \frac{\frac{\alpha_1}{\alpha_0} - \frac{\alpha_1(\alpha_1+1)}{\alpha_0(\alpha_0+1)}}{\frac{\alpha_1(\alpha_1+1)}{\alpha_0(\alpha_0+1)} - \frac{\alpha_1^2}{\alpha_0^2}} \quad (3.52) \quad \{\text{eq:5}\}$$

$$= \frac{\frac{\alpha_1(\alpha_0+1) - \alpha_1(\alpha_1+1)}{\alpha_0(\alpha_0+1)}}{\frac{\alpha_0\alpha_1(\alpha_1+1) - \alpha_1^2(\alpha_0+1)}{\alpha_0^2(\alpha_0+1)}} \quad (3.53)$$

$$= \frac{\alpha_0\alpha_1(\alpha_0 - \alpha_1)}{\alpha_1(\alpha_0\alpha_1 + \alpha_0 - \alpha_0\alpha_1 - \alpha_1)} \quad (3.54)$$

$$= \alpha_0 \quad (3.55)$$

$$\alpha_i = \mathbb{E}[X_i]\alpha_0 \quad (3.56) \quad \{\text{eq:6}\}$$

Z rovnice (3.52) vypočítáme sumu všech parametrů α_0 . Protože střední hodnota proměnné z Dirichletovského rozdělení je právě $\frac{\alpha_i}{\alpha_0}$, tak jednotlivé parametry získáme z rovnice (3.56).

3.4 Algoritmus

Algoritmus 4 Expectation Propagation pro učení parametrů

{alg:ep}

Parametry zpráv z faktoru β do množiny parametrů θ_i označíme $\alpha_{f_\beta \rightarrow \theta_i}$.
 Parametry zpráv z množiny parametrů θ_i do faktoru β označíme $\alpha_{\theta_i \rightarrow f_\beta}$.
 Parametry marginální distribuce množiny parametrů θ_i označíme α_i .

init

Nastav zprávy mezi faktory a proměnnými na 1.

Nastav parametry $\alpha_{f_\beta \rightarrow \theta_i}$ na 1.

Nastav parametry α_i na apriorní hodnotu.

end init

repeat

Vyber faktor $f_{\tilde{\beta}}$, který se bude aktualizovat.

Spočítej všechny zprávy z parametrů:

for každý parametr θ_i spojený s faktorem $f_{\tilde{\beta}}$ **do**

Parametry zprávy z θ_i do $f_{\tilde{\beta}}$: $\alpha_{\theta_i \rightarrow f_{\tilde{\beta}}} = \alpha_i - \alpha_{f_{\tilde{\beta}} \rightarrow \theta_i} + 1$.

end for

Aktualizuj zprávy z faktoru do proměnných:

for každou proměnnou X_i , spojenou s faktorem $f_{\tilde{\beta}}$ **do**

Zpráva z $f_{\tilde{\beta}}$ do X_i podle (3.6):

$$\hat{f}(x_j) = \sum_{\mathbf{x}: x_j = \tilde{x}_j} \mathbb{E}_{q \setminus (\theta_{\rho(\mathbf{x}'), x_0})} \prod_{i \neq j} m_{x_i \rightarrow f_{\tilde{\beta}}}(x_i)$$

end for

Aktualizuj marginální pravděpodobnost parametrů:

for každý parametr θ_i spojený s faktorem $f_{\tilde{\beta}}$ **do**

Spočítej parametry α_i^* pro Dirichletovské rozdělení, které nejlépe aproximuje cílovou marginální distribuci (3.29). Metoda popsána v předchozí sekci.

Parametry zprávy z $f_{\tilde{\beta}}$ do θ_i :

$$\alpha_{f_{\tilde{\beta}} \rightarrow \theta_i} = \alpha_i^* - \alpha_{\theta_i \rightarrow f_{\tilde{\beta}}} + 1$$

Aktualizuj parametry marginální distribuce $q(\theta_i)$

$$\alpha_i = \alpha_i^* = \alpha_{f_{\tilde{\beta}} \rightarrow \theta_i} + \alpha_{\theta_i \rightarrow f_{\tilde{\beta}}}$$

end for

for každou proměnnou X_i , spojenou s faktorem $f_{\tilde{\beta}}$ **do**

Aktualizuj zprávy z proměnných do faktoru:

$$m_{x_i \rightarrow f_{\tilde{\beta}}}(x_i) = \prod_{\beta \neq \tilde{\beta}} m_{f_\beta \rightarrow x_i}(x_i)$$

end for

until konvergence

4. Implementace

V předchozích kapitolách byly popsány teoretické základy nutné pro implementaci inference v bayesovských sítích. V této kapitole kapitole bude popsána

Závěr

Literatura

- [1] Arnborg, S.; Corneil, D. G.; Proskurowski, A.: Complexity of finding embeddings in ak-tree. *SIAM Journal on Algebraic Discrete Methods*, ročník 8, č. 2, 1987: s. 277–284.
- [2] Bernardo, J. M.; Smith, A. F.: *Bayesian theory*, ročník 405. Wiley. com, 2009.
- [3] Bertoldi, N.; Federico, M.: A new decoder for spoken language translation based on confusion networks. In *Automatic Speech Recognition and Understanding, 2005 IEEE Workshop on*, IEEE, 2005, s. 86–91.
- [4] Bishop, C. M.; Nasrabadi, N. M.: *Pattern recognition and machine learning*, ročník 1. springer New York, 2006.
- [5] Black, A.; Taylor, P.; Caley, R.; aj.: The Festival Speech Synthesis System, Version 1.4. 2. *Unpublished document available via <http://www.cstr.ed.ac.uk/projects/festival.html>*, 2001.
- [6] Black, A. W.; Lenzo, K. A.: Flite: a small fast run-time synthesis engine. In *4th ISCA Tutorial and Research Workshop (ITRW) on Speech Synthesis*, 2001.
- [7] He, Y.; Young, S.: Semantic processing using the hidden vector state model. *Computer speech & language*, ročník 19, č. 1, 2005: s. 85–106.
- [8] Koller, D.; Friedman, N.: *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- [9] Kullback, S.: *Information theory and statistics*. Courier Dover Publications, 1997.
- [10] Mairesse, F.; Gasic, M.; Jurcicek, F.; aj.: Spoken language understanding from unaligned data using discriminative classification models. In *Acoustics, Speech and Signal Processing, 2009. ICASSP 2009. IEEE International Conference on*, IEEE, 2009, s. 4749–4752.
- [11] Minka, T. P.: Expectation propagation for approximate Bayesian inference. In *Proceedings of the Seventeenth conference on Uncertainty in artificial intelligence*, Morgan Kaufmann Publishers Inc., 2001, s. 362–369.
- [12] Murphy, K. P.; Weiss, Y.; Jordan, M. I.: Loopy belief propagation for approximate inference: An empirical study. In *Proceedings of the Fifteenth conference on Uncertainty in artificial intelligence*, Morgan Kaufmann Publishers Inc., 1999, s. 467–475.
- [13] Pearl, J.: *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann Pub, 1988.

- [14] Povey, D.; Ghoshal, A.; Boulianne, G.; aj.: The Kaldi Speech Recognition Toolkit. In *IEEE 2011 Workshop on Automatic Speech Recognition and Understanding*, IEEE Signal Processing Society, Proseminar 2011, IEEE Catalog No.: CFP11SRW-USB.
- [15] Raux, A.; Bohus, D.; Langner, B.; aj.: Doing research on a deployed spoken dialogue system: One year of Let's Go! experience. In *Proc. Interspeech*, 2006, s. 65–68.
- [16] Tatikonda, S. C.; Jordan, M. I.: Loopy belief propagation and Gibbs measures. In *Proceedings of the Eighteenth conference on Uncertainty in artificial intelligence*, Morgan Kaufmann Publishers Inc., 2002, s. 493–500.
- [17] Thomson, B.; Schatzmann, J.; Young, S.: Bayesian update of dialogue state for robust dialogue systems. In *Acoustics, Speech and Signal Processing, 2008. ICASSP 2008. IEEE International Conference on*, IEEE, 2008, s. 4937–4940.
- [18] Thomson, B.; Young, S.: Bayesian update of dialogue state: A POMDP framework for spoken dialogue systems. *Computer Speech & Language*, ročník 24, č. 4, 2010: s. 562–588.
- [19] Walker, W.; Lamere, P.; Kwok, P.; aj.: Sphinx-4: A flexible open source framework for speech recognition. 2004.
- [20] Wong, Y. W.; Mooney, R.: Learning synchronous grammars for semantic parsing with lambda calculus. In *ANNUAL MEETING-ASSOCIATION FOR COMPUTATIONAL LINGUISTICS*, ročník 45, 2007, str. 960.
- [21] Young, S.; Evermann, G.; Gales, M.; aj.: The HTK book. *Cambridge University Engineering Department*, ročník 3, 2002.
- [22] Young, S.; Gašić, M.; Keizer, S.; aj.: The Hidden Information State Model: a practical framework for POMDP-based spoken dialogue management. *Computer Speech & Language*, ročník 24, č. 2, 2010: s. 150–174.
- [23] Zen, H.; Nose, T.; Yamagishi, J.; aj.: The HMM-based speech synthesis system (HTS) version 2.0. In *Proc. of Sixth ISCA Workshop on Speech Synthesis*, 2007, s. 294–299.
- [24] Zettlemoyer, L. S.; Collins, M.: Online learning of relaxed CCG grammars for parsing to logical form. In *In Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL-2007)*, Citeseer, 2007.

Seznam tabulek

Seznam použitých zkratek

Přílohy