

# یک مدل بیزی برای استخراج با مرتب گرامر زبان طبیعی

هشام فیلی<sup>۱</sup>، حمیدرضا قادر<sup>۲</sup> و مرتضی آنالوی<sup>۳</sup>

<sup>۱</sup> آزمایشگاه پردازش زبان طبیعی، دانشکده برق و کامپیوتر، دانشگاه تهران

<sup>۲</sup> گروه هوش مصنوعی و رباتیک، دانشکده کامپیوتر، دانشگاه علم و صنعت

## چکیده

در این مقاله نشان داده‌ایم که مسئله استخراج با مرتب گرامرهای زبان طبیعی، می‌تواند به عنوان ترکیبی پیچیده از تعداد زیادی مسئله انتخاب مدل تعریف شود. مسائل انتخاب مدل به مسائلی گفته می‌شود که در آنها میان مدل‌هایی با پیچیدگی آماری متفاوت تصمیم‌گیری می‌شود. برای ارائه مدل مورد نظر، ابتدا یک مدل بیزی از فرایند شناختی انتخاب مدل را معرفی کردہ‌ایم. این مدل مسئله انتخاب مدل را با تطابق بیشتر با رفتار انسان حل می‌کند. سپس با استفاده از تعیین‌یافته مدل بیزی، که فرایند دیریکله است، به همراه مفهوم گرامرهای مبتنی بر سابقه، یک مدل بیزی مبتنی بر سابقه برای استخراج با مرتب گرامر ارائه کردہ‌ایم. نتایج آزمایشی که با این مدل روی یک پیکره استاندارد برای زبان انگلیسی صورت گرفته است، در مقایسه با مدل مرجع طبق معیار F1 به مقدار ۹/۱٪ پیشرفت نشان می‌دهد.

واژگان کلیدی: مدل بیزی، گرامر جایگزینی درخت، فرایند دیریکله، فرایند رستوران چینی، گرامرهای مبتنی بر سابقه.

صورت گرفته است تا گرامرهایی استخراج شود که نتایج خوبی در تجزیه نحوی جملات زبان طبیعی تولید کنند (Cohn, et al., 2009). این تلاش‌ها شامل ایجاد و به کارگیری روش‌هایی برای تقویت صورت‌گرایی<sup>۴</sup> توصیف کننده گرامر (Johnson, 1998) (Collins, 1999) (Charniak, 2000) روش‌های مسطوح‌سازی پیشرفته (Charniak, et al., 2000) و روش‌های انتخاب ویژگی (Charniak, et al., 2005) است. تمام این تلاش‌ها که در کار ارائه شده توسط گهن و همکاران (Cohn, et al., 2009) تحت عنوان مهندسی گرامر از آنها باد شده، برای پوشش ضعف‌های گرامرهای مستقل از متن در توصیف جملات زبان طبیعی است. بیشتر ضعف‌های گرامرهای مستقل از متن در توصیف زبان‌های طبیعی ناشی از فرض‌های استقلالی است که در این گرامرها می‌شود. به عنوان مثال، برخی خواص زبان‌های طبیعی نظیر خاصیت حذف ضمیر در جایگاه فاعل

## ۱- مقدمه و کارهای پیشین

مجموعه تحقیقات انجام شده روی استخراج گرامر، در دو دهه اخیر، حاکی از اهمیت گرامر محاسباتی در پردازش زبان طبیعی است (Cohn, et al., 2009; Johnson, 1998; Collins, 1999; Charniak, 2000; Charniak, et al., 2005). در استخراج گرامر به صورت بامری<sup>۱</sup>، گرامر محاسباتی از روی درخت‌های تجزیه جملات که از قبل توسط زبان‌شناسان تهیه شده است، استخراج می‌شود. بدین معنا که جملات یک پیکره توسط انسان، تجزیه نحوی شده و درخت‌های تجزیه<sup>۲</sup> این جملات تهیه می‌شود و مجموعه این درخت‌ها با نام بانک درختی<sup>۳</sup> برای استخراج قواعد پایه گرامر مورد استفاده قرار می‌گیرد. تلاش‌های زیادی در گذشته برای استخراج گرامرهای مستقل از متن احتمالاتی<sup>۴</sup>

<sup>1</sup> Supervised

<sup>2</sup> Parse Trees

<sup>3</sup> Tree-bank

<sup>4</sup> Probabilistic Context Free Grammars

<sup>۵</sup> Formalism

مستقل از متن، مبتنی بر سابقه<sup>۶</sup> ارائه شد. در (Collins, 1999) کار دیگری ارائه شد که با ایجاد وابستگی واژگانی<sup>۷</sup> توانست صورت‌گرا را تقویت کند. این کار تجزیه‌گر مناسب برای این نوع گرامر را نیز ارائه کرد و توانست برای تجزیه جملات WSJ100 مقدار ۸۸/۲٪ را برای معیار *F1* به دست آورد. چارنیاک<sup>۸</sup> توانست در (Charniak, 2000) با استفاده از یک تجزیه‌گر و روش مسطح‌سازی<sup>۹</sup> جدید، که هر دو بر پایه آنتروپی بیشینه<sup>۱۰</sup> بودند، مقدار ۸۹/۵٪ را برای جملات موجود در WSJ100 به دست آورد. این کار در (Charniak, et al., 2005) با استفاده از یک الگوریتم رتبه‌بندی<sup>۱۱</sup> تقویت شد و مقدار *F1* را برای جملات WSJ100 به ۹۱/۰٪ رساند. در (Petrov, et al., 2007) روشی برای استخراج گرامرهای مستقل از متن احتمالاتی ارائه شد که با وجود عدم وابستگی واژگانی در گرامر استخراجی، توانست دقیق قابل رقابت با بهترین گرامرهای دارای وابستگی واژگانی به دست آورد. در مدل ارائه شده در این کار، در یک فرایند تکراری، هر غیرپایانه به دو غیرپایانه تقسیم می‌شد تا بدین ترتیب بتواند تخمین بهتری از پارامترهای آزاد مدل ایجاد کند (Cahill, 2008). مقدار *F1* گزارش شده برای جملات موجود در WSJ100 در این کار برابر ۹۰٪ است. براساس دانش نگارندگان، بهترین نتایج تجزیه نحوی مربوط به کار ارائه شده در (McClosky, et al., 2006) است. در این کار با آموزش تجزیه‌گر نحوی روی نتایج تجزیه خود دقت عملکرد تجزیه‌گر افزایش داده شده است. در این کار مقدار گزارش شده برای معیار *F1* در تجزیه جملات موجود در WSJ100 برابر ۹۲/۱٪ است.

مدل *DOP* (Bod, 1992) نیز مدل موفقی است که گرامر استخراجی آن در دسته گرامرهای جایگزینی درخت احتمالاتی قرار می‌گیرد. البته این مدل به طور واضح گرامر استخراج نمی‌کند، بلکه با استفاده از داده‌های آموزشی و زیردرخت‌های پرتکرار در این داده‌ها، بهترین تجزیه نحوی برای جمله ورودی را محاسبه می‌کند. این مدل در (Bod, 2003) به مقدار ۹۰/۷٪ برای معیار *F1* در تجزیه جملات WSJ100 دست یافت.

**مدل‌های بیزی غیرپارامتری در بسیاری از**

در زبان فارسی (فیلی، ۱۳۸۵) و خاصیت واستنگی-Serial در زبان‌های آلمانی، سوئیسی و هلندی (Mohri, et al., 2006) خواص حساس به متن هستند و گرامرهای مستقل از متن توانایی توصیف آنها را ندارند. در مدل‌هایی که مورد اشاره شد، با تغییر بانک درخت با زبان تمامی این تلاش‌ها برای بالا بردن دقت مدل باید دوباره انجام گیرد (Cohn, et al., 2009). در این مقاله، ما مدلی بامribi را ارائه می‌کیم که گرامر جایگزینی درخت احتمالاتی<sup>۱</sup>، که صورت‌گرایی قوی‌تر از گرامرهای مستقل از متن در توصیف زبان طبیعی است، استخراج می‌کند. علاوه‌بر این، مدل ما از دیدگاه آماری در قالب روش‌های بیزی غیرپارامتری<sup>۲</sup> دسته‌بندی می‌شود. در این روش‌ها تعداد پارامترهای مدل محدود فرض نمی‌شود، از این‌رو پیچیدگی مدل با پیچیدگی داده‌های آموزشی به طور خودکار تنظیم شده و مشکلات برازش کم یا زیاد<sup>۳</sup> پیش نمی‌آید (Teh, 2010). بدین ترتیب مدل‌های بیزی غیرپارامتری پیچیدگی‌های داده‌ی آموزشی را به خوبی فرمای گیرند. این خاصیت به همراه قوی‌تر بودن صورت‌گرای استخراجی در مدل ما سبب می‌شود، بسیاری از پیچیدگی‌های زبانی که در مدل‌های پارامتری با روش‌های مهندسی گرامر در مدل کد می‌شوند، به طور خودکار از داده‌های آموزشی فراگرفته شود.

یکی از اولین کارهای شاخص در زمینه استخراج گرامر، که به استخراج گرامر مستقل از متن احتمالاتی پرداخته است، کار ارائه شده توسط چارنیاک (Charniak, 1996) است. در این کار با استفاده از یک مدل مبتنی بر درستنمایی<sup>۴</sup> بیشینه گرامری با ۱۰۶۰۵ قانون استخراج شد. یک برتری شاخص این کار نسبت به مدل‌های موجود در آن زمان، تجزیه‌گر ارائه شده در این کار بود که توانست با استفاده از گرامر ایجاد شده، بالاترین دقت تجزیه جملات در آن زمان را به دست آورد. در این کار برای تجزیه جملات با *WSJ100* طول کمتر از یک‌صد کلمه<sup>۵</sup> موجود در پیکرۀ<sup>۶</sup> *WSJ100* (۷۹/۵٪) برای معیار *F1* گزارش شده است. پس از آن تلاش‌هایی برای توسعه مدل‌های گرامرهای مستقل از متن آغاز شد تا ضعف‌های این گرامرهای را در تجزیه جملات پوشش دهند. در کار ارائه شده توسط جانسون (Johnson, 1998) (روشی برای استخراج گرامرهای

<sup>6</sup> History-based Grammars

<sup>7</sup> Lexicalization

<sup>8</sup> Charniak

<sup>9</sup> Smoothing

<sup>10</sup> Maximum Entropy

<sup>11</sup> Reranking

<sup>1</sup> Probabilistic Tree Substitution Grammar

<sup>2</sup> Non-parametric Bayesian Models

<sup>3</sup> Under-fit or over-fit

<sup>4</sup> Maximum Likelihood

<sup>5</sup> Wall Street Journal



به کار گرفته و نشان می‌دهیم که چگونه این مدل مبنایی، منجر به تعداد زیادی فرایند دیریکله در مسئله استخراج گرامر می‌شود. سپس مفهوم سابقه را نیز در مدل خود وارد کرده و تغییرات این مفهوم بر فرایند دیریکله را نشان می‌دهیم.

مدل ارائه شده در (Johnson, et al., 2007) یک نمونه دیگر برای استخراج گرامر با مدل‌های بیزی غیر پارامتری است. در این کار یک قالب برای استخراج گرامر با مدل‌های بیزی غیرپارامتری ارائه می‌شود که آن را گرامر تطبیقی<sup>۵</sup> می‌نامند. این مدل برای یادگیری لغات از داده گویی در (Johnson, et al., 2009) به کار گرفته شده است. اما این گرامر برای توصیف جملات زبان طبیعی مناسب نیست؛ چرا که بر پایه گرامرهای مستقل از متن بنا شده است.

در (O'Donnell, et al., 2009) نیز فرایند پیتمن-یور برای استخراج گرامر جایگزینی درخت به کار رفته است. در این کار، استخراج گرامر با مدل‌های بیزی غیرپارامتری از دیدگاه شناختی مدل شده و مورد بررسی و آزمایش قرار گرفته است. کار مذکور یک قالب بیزی غیرپارامتری برای استخراج گرامرهایی که توانایی خوبی در توصیف جملات زبان طبیعی دارند، ارائه می‌کند. در حقیقت قالب ارائه شده در این کار تعمیم‌یافته گرامر تطبیقی است. نتیجه مدل ارائه شده در این کار نیز، همانند مدل ارائه شده در (Cohn, et al., 2009) و (Cohn, et al., 2010)، تعداد زیادی فرایند دیریکله برای یادگیری واحدهای زبانی به شرط گرۀ ریشه است. همچنین، صورت‌گرایی استخراجی در این کار گرامر جایگزینی درخت احتمالاتی است. علاوه‌بر این، در این کار نیز همانند کارهای ارائه شده در (Cohn, et al., 2010) و (Cohn, et al., 2009) مفهوم سابقه در استخراج گرامر دخیل نمی‌شود.

بنابر آنچه گفته شد، در این مقاله یک مدل بامبی، بر پایه روش‌های بیزی غیرپارامتری، برای استخراج گرامر زبان طبیعی ارائه شده است. صورت‌گرایی گرامر استخراج شده در این مقاله یک مدل قوی‌تر از گرامر مستقل از متن است که گرامر جایگزینی درخت احتمالاتی نامیده می‌شود. کار این مقاله به کار ارائه شده در (Cohn, et al., 2009) شباهت دارد و درواقع توسعه‌یافته آن بهشمار می‌رود. نوآوری‌های این کار نسبت به کارهایی که پیش از این مدل‌های مشابهی برای استخراج گرامر ارائه کرده‌اند، دو

کاربردهای پردازش زبان طبیعی به کار گرفته شده‌اند. در (Goldwater, et al., 2007) مسئله بدون ناظر برچسب‌گذاری جزء کلام با استفاده از این مدل‌ها حل شده است. در (Johnson, et al., 2007) قالبی<sup>۱</sup> برای یادگیری لغات از داده گویی ارائه شده است، که بر پایه مدل‌های بیزی غیرپارامتری بنا شده است. در (Brody, et al., 2009) حل مسئله ابهام‌زدایی معنای لغت<sup>۲</sup> با استفاده از این مدل‌ها، مطرح شده است. برچسب‌گذاری نقش معنایی نیز در (Titov, et al., 2011) با این مدل‌ها حل شده است.

مدل‌های بیزی غیرپارامتری پیش از این برای استخراج بامبی گرامر جایگزینی درخت نیز در (Cohn, et al., 2009) و (O'Donnell, et al., 2010) (Cohn, et al., 2010) به کار گرفته شده‌اند. در (Cohn, et al., 2009) و (Cohn, et al., 2009) فرایند دیریکله<sup>۳</sup> و فرایند پیتمن-یور<sup>۴</sup>، که شکل تعمیم‌یافته فرایند دیریکله است، برای استخراج گرامر جایگزینی درخت احتمالاتی به کار رفته‌اند. این دو فرایند نمونه‌هایی از مدل‌های بیزی غیرپارامتری هستند. در این دو کار، هر واحد پایه در گرامر به شرط گرۀ ریشه آن واحد، با استفاده از یکی از دو فرایند مذکور، استخراج می‌شود. بدین ترتیب مفهوم سابقه در این دو کار مد نظر قرار نگرفته است. مفهوم سابقه موجب وابستگی یادگیری هر واحد گرامر به بخشی از بافت جمله که آن واحد در آن رخ می‌دهد، می‌شود. این مفهوم در (Johnson, 1998) برای کاهش اثر سوءفرض‌های استقلال در استخراج گرامرهای مستقل از متن احتمالاتی مطرح شد. در (Johnson, 1998) برای استخراج گرامرهای مبتنی بر سابقه، روشی ارائه شد که براساس آن غیرپایانه‌های درخت‌های تجزیه موجود در داده‌های آموزشی، با برچسب واحد زبانی خود برچسب‌گذاری می‌شوند. بدین ترتیب واحد زبانی رخدادیافته در زیر یک غیرپایانه، به بخشی از بافت جمله که زیر غیرپایانه‌ی پدر ریشه‌اش قرار دارد، وابسته می‌شود و این وابستگی در فرایند یادگیری، تأثیر می‌گذارد. علاوه‌بر این، در این دو کار مبنای شناختی فرایندهای دیریکله و پیتمن-یور مورد اشاره قرار نگرفته است و علت به کارگیری تعداد زیادی فرایند دیریکله و پیتمن-یور محدودیتی بیان شده است که در گرامر استخراجی از یک فرایند، بروز می‌کند. اما در مدل ارائه شده در این مقاله، ما مبنای شناختی فرایند دیریکله را

<sup>1</sup> Framework

<sup>2</sup> Word Sense Disambiguation

<sup>3</sup> Dirichlet Process

<sup>4</sup> Pitman-Yor Process

<sup>5</sup> Adaptor Grammar



آن روش نمونه برداری گیبس<sup>4</sup> را، که یک روش مونت کارلوی زنجیره مارکوف است، توضیح می دهیم. در این مقاله، از نمونه برداری گیبس برای تخمین توزیع احتمال حاصل از فرایندهای دیریکله استفاده می کنیم. در پایان نیز اشاره ای به طور تقریبی به گرامر جایگزینی درخت می کنیم.

## ۱-۲ مدل بیزی از فرایندهای شناختی انتخاب مدل

مدل بیزی ای را که مبنای مدل ارائه شده در این تحقیق تشکیل می دهد یک مدل بیزی از فرایندهای شناختی انتخاب مدل است. بسیاری از مسائل در دنیای واقعی هستند که در آنها باید میان دو یا چند فرضیه که از نظر پیچیدگی با هم متفاوت هستند و فرضیه های غیر همگن به سمار می آیند، تصمیم گیری کرد. این گونه مسائل که در آنها دانش حاصل از داده های مشاهده شده را برای انتخاب میان دو مدل آماری که از نظر پیچیدگی با هم متفاوت هستند به کار می گیریم، انتخاب مدل نامیده می شوند (Griffiths, et al., 2008).

برای درک عملکرد مدل بیزی شناختی از فرایندهای انتخاب مدل یک مثال استاندارد را ارائه می کنیم. این مثال از (Griffiths, et al., 2008) برگرفته شده است. فرض کنید که یک سکه داریم و  $\theta$  نشان دهنده احتمال رخداد شیر در پرتاب آن است. می خواهیم تصمیم بگیریم که این سکه، یک سکه معمولی با  $\frac{1}{2} = \theta$  است یا سکه های است که در آن  $\theta$  مقدار خود را از یک توزیع یکنواخت بین  $0$  و  $1$  گرفته است. این فرضیه به مفهوم آن است که سکه معمولی نیست و احتمال شیر و خط آمدن، لزوماً مساوی نیست؛ ولی اینکه مقدار احتمال چقدر باشد در این فرضیه مهم نیست. مهم آن است که بدانیم سکه اریب است. بنابراین دو فرضیه وجود دارد:

$$1. \quad \theta = \frac{1}{2} \text{ که آنرا } h_0 \text{ می نامیم.}$$

۲.  $\theta$  مقدار خود را از یک توزیع یکنواخت بین  $0$  و  $1$  گرفته است. این فرضیه را  $h_1$  می نامیم.

فرض می کنیم که هیچ دانشی مبنی بر برتری داشتن یکی از فرضیات بر دیگری وجود ندارد. بنابراین می توان فرض کرد که  $\theta = \frac{1}{2} = P(h_0) = P(h_1)$ . حال فرض کنید که ما یک ترتیب ده تایی را از پرتاب این سکه مشاهده می کنیم،

<sup>4</sup> Gibbs Sampling

بخش دارد. اول اینکه ما نشان می دهیم چگونه مدل ارائه شده برای استخراج گرامر از یک مدل بیزی ساده تر، که با عنوان مدل بیزی از فرایندهای شناختی انتخاب مدل<sup>1</sup> معرفی می کنیم، ناشی می شود. در این میان، دلیل اینکه تعداد زیادی فرایندهای دیریکله در مدل ارائه شده برای استخراج گرامر به وجود می آید، روش می شود. نتیجه این بخش از (Cohn, et al., 2009; Cohn, et al., 2010; O'Donnell, et al., 2009) می گردد. دوم اینکه، ما مدل ارائه شده خود را با دخیل کردن مفهوم سابقه طراحی می کنیم و نشان می دهیم که این توسعه منجر به نتایجی بهتر از نتایج ارائه شده در (Cohn, et al., 2009) و (Cohn, et al., 2010) می شود. روشی که در این مقاله معرفی شده روی زبان انگلیسی مورد آزمایش قرار گرفته است. نتایج ارزیابی روی بخش WSJ پیکره استاندارد Penn Treebank، که به طور تقریبی تمامی کارهای معتبر در زبان انگلیسی نیز روی آن تعریف شده اند، نشان دهنده برتری روش ارائه شده در این مقاله نسبت به کارهای مشابه است. در ادامه، در بخش دو پیش زمینه تئوری را که برای درک مدل ارائه شده در این مقاله لازم است معرفی، سپس در بخش سه مدل خود را ارائه و در بخش چهار چگونگی استنتاج در مدل ارائه شده را بیان می کنیم. پس از آن در بخش پنج آزمایش های انجام شده را توضیح داده و نتایج حاصل را مورد تحلیل قرار می دهیم. در بخش پایانی نتیجه گیری خود را از مدل ارائه شده و آزمایش ها بیان می کنیم.

## ۲- پیش زمینه

روش های بیزی غیر پارامتری و استنتاج بیزی<sup>2</sup> با استفاده از روش های مونت کارلوی زنجیره مارکوف<sup>3</sup>، روش هایی هستند که در نیمه دوم دهه گذشته میلادی در حل مسائل حوزه پردازش زبان طبیعی به کار گرفته شده اند (Goldwater, et al., 2007; Goldwater, et al., 2009; Cohn, et al., 2009; Cohn, et al., 2010; Johnson, et al., 2007) بخش، ابتدا به معرفی مدل بیزی از فرایندهای شناختی انتخاب مدل، که مبنای تئوری مدل ارائه شده در این مقاله را تشکیل می دهد می پردازیم. سپس فرایندهای دیریکله را که یک روش بیزی غیر پارامتری است و تعمیم بی نهایت مدل بیزی شناختی از فرایندهای انتخاب مدل است، معرفی کرده و پس از

<sup>1</sup> Model Selection

<sup>2</sup> Gibbs Sampling

<sup>3</sup> Gibbs Sampling



نگرش مدل بیزی در حل این گونه مسائل، این موضوع را دخالت می‌دهد که هر چند درجه آزادی بیشتر امکان تطبیق بیشتر مدل با داده‌های مشاهده شده را فراهم می‌کند، اما این درجه آزادی امکان یک تطبیق بدتر را نیز ایجاد می‌کند. نگرش مدل بیزی در حل این گونه مسائل با انتگرال گیری روی درجات آزادی مدل پیچیده‌تر، میان حالتی که مدل پیچیده‌تر انطباق بهتری در تولید داده ایجاد می‌کند و حالتی که این مدل انطباق بدتری نسبت به مدل ساده، در تولید داده‌ها ایجاد می‌کند، میانگین می‌گیرد (*Griffiths, et al.*, 2008). بنابراین نگرش مدل بیزی، مدل پیچیده‌تر را بر اساس انطباق بیشتر مدل، زمانی که پارامترهای مدل بهترین مقادیر را اختیار کرده‌اند، انتخاب نمی‌کند. بلکه زمانی مدل پیچیده‌تر را بر مدل ساده‌تر برتری می‌دهد که مدل پیچیده‌تر با انتخاب تصادفی مقادیر پارامترها، در بیشتر مواقع، انطباق بهتری بر روی داده‌های مشاهده شده ایجاد کند. این موضوع با قضاوت انسان در حل این مسائل مطابقت دارد، چرا که با فرض در نظر گرفتن مثال پرتاب سکه، انسان نیز، همواره و با مشاهده هر ترتیبی از شیر و خط، مدل پیچیده‌تر را به مدل ساده‌تر ترجیح نمی‌دهد؛ بلکه با فرض اینکه پارامترهای مدل پیچیده‌تر مقدار خود را به صورت تصادفی بگیرند، قضاوت می‌کند که آیا در بیشتر موارد مدل پیچیده‌تر انطباق بهتری روی داده‌ها تولید می‌کند یا مدل ساده‌تر.

## ۲-۲- فرایند دیریکله

فرایند دیریکله یک فرایند تصادفی است که در حقیقت به گونه‌ای نتیجه مدل بیزی از فرایند شناختی انتخاب مدل برای مسائلی است که مدل‌های مورد تصمیم‌گیری دارای پیچیدگی نامعلوم هستند. برای به دست آورن درک بهتر از چنین مسائلی یک مثال می‌زنیم: با فرض در نظر داشتن مثال سکه‌ای که در بخش قبل مورد بحث قرار گرفت، فرض کنید ما ترتیبی از پرتاب یک جسم چندوجهی را مشاهده کردیم ولی تعداد وجوده جسم مورد نظر را نمی‌دانیم و می‌خواهیم همانند مثال سکه در مورد توزیع احتمال جسم مورد نظر تصمیم‌گیری کنیم. این یک مسئله انتخاب مدل بسیار پیچیده است.

فرایند دیریکله یک فرایند تصادفی است که در مدل‌های بیزی غیرپارامتری به کار گرفته می‌شود. این فرایند تصادفی در حقیقت یک توزیع احتمال روی توزیع

به عنوان مثال {ش، ش، خ، ش، خ، ش، ش، خ}. اگر ما از روش درستنمایی بیشینه<sup>۱</sup> (MLE) استفاده کنیم برای مشاهده ترتیب فوق، خواهیم داشت  $N = \hat{\theta}$  که در آن  $N$  نشان‌دهنده تعداد رخداد شیر است. در نتیجه کاربرد این روش،  $P(\mathbf{h}|\mathbf{w}, \hat{\theta})$  برای هر ترتیب مشاهده شده‌ای که تعداد نامساوی از رخداد شیر و خط دارد، همواره به  $\mathbf{h}_1$  احتمال بیشتری اختصاص می‌دهد. همچنین برای هر  $\mathbf{h}_0$  و  $\mathbf{h}_1$  برابر می‌شود. بدین ترتیب چنین راه حلی همواره فرض  $\mathbf{h}_1$  را برای جواب این مسئله انتخاب می‌کند که این موضوع به طور کامل مخالف رفتاری است که انسان در حل چنین مسئله‌ای از خود نشان می‌دهد. حال اگر از مدل بیزی مورد نظر استفاده کنیم، حل مسئله به گونه‌ای دیگر خواهد بود. در این حالت ما توزیع احتمال روی پاسخ مسئله را با انتگرال گیری روی مقادیر ممکن برای  $\theta$  پیدا می‌کنیم. یعنی:

$$(1) P(\mathbf{h}|\mathbf{w}) = \int P(\mathbf{h}|\mathbf{w}, \theta)P(\theta|\mathbf{w})d\theta$$

رابطه فوق منجر به مقدار احتمال زیر خواهد شد.

$$(2) P(\mathbf{h}_1|\mathbf{w}) = 1 / (1 + \frac{11!}{N_+ N_-^{210}})$$

این رابطه زمانی که تعداد رخداد شیر و خط برابر باشد بسیار کوچک‌تر از ۵.۰ خواهد بود و فقط زمانی  $\mathbf{h}_1$  را بر  $\mathbf{h}_0$  برتری می‌دهد که  $8 \geq N_+ \geq 2$  باشد. این نتیجه با رفتاری که از انسان در حل چنین مسئله‌ای مشاهده می‌شود مطابقت بیشتری دارد و در حقیقت مدل مناسب‌تری از فرایند شناختی انتخاب مدل در انسان است.

آنچه در حل این مسئله به روش بیز پنهان است بدین ترتیب بیان می‌شود:

مدل‌ها و فرضیات آماری پیچیده‌تر دارای درجه آزادی بیشتری هستند. به همین دلیل همواره می‌توان آنها را بهتر از مدل‌ها و فرضیات ساده بر داده‌های مشاهده شده منطبق کرد؛ به گونه‌ای که داده‌های مشاهده شده تحت مدل پیچیده‌تر احتمال تولید بیشتری داشته باشد (*Griffiths, et al.*, 2008). به عنوان مثال، در مثال پرتاب سکه، برای هر ترتیبی از شیر و خط می‌توان مقادیر برای  $\theta$  یافت که احتمال بیشتری به تولید آن ترتیب خاص نسبت به  $\frac{1}{2} = \theta$  بدهد. بدین ترتیب به نظر می‌رسد که همواره باید در این گونه مسائل مدل پیچیده‌تر را بر مدل ساده برتری داد. اما

<sup>۱</sup> Maximum likelihood

برای توزیع احتمال  $(\theta_{n+1} | \theta_1, \dots, \theta_n)$  با انتگرال گیری از  $P(\theta_{n+1} | \theta_1, \dots, \theta_n)$  روی  $G$  نتیجه می‌شود:

$$P(\theta_{n+1} \in A | \theta_1, \dots, \theta_n) = \frac{1}{\alpha+n} (\alpha H(A) + \sum_{i=1}^n \delta_{\theta_i}(A)) \quad (5)$$

در اینجا  $G$  حکم متغیر نهانی را دارد که در بخش قبل نیز روی آن انتگرال گرفتیم. در رابطه فوق  $\Theta \subset A$  است و  $\delta_{\theta_i}$  تابع دلتای کرونکر<sup>۴</sup> است که تمام احتمال خود را در نقطه  $\theta_i$  متمرکز کرده است (*Teh, 2010*). بنابراین داریم:

$$\theta_{n+1} | \theta_1, \dots, \theta_n \sim \frac{1}{\alpha+n} (\alpha H + \sum_{i=1}^n \delta_{\theta_i}) \quad (6)$$

برای درک این رابطه فرض کنید  $\{\theta_1^*, \dots, \theta_{|c|}^*\}$  نشان‌دهنده تمامی مقادیر مختلفی است که  $\theta_1, \dots, \theta_n$  گرفته‌اند. حال فرض کنید  $\{c_1, \dots, c_n\} = \mathbf{c}$  برداری از نمایه‌ها برای  $\{\theta_1^*, \dots, \theta_{|c|}^*\}$  باشد، به‌طوری‌که  $\theta_i^* = \theta_{c_i}^*$  باشد. در این صورت توزیع  $\theta_{n+1}$  به‌صورت زیر خواهد بود:

$$\theta_n = \begin{cases} \theta_i^* & \text{with prob. } \frac{|\{j: c_j = i\}|}{n-1+\alpha} \\ \theta, \quad \theta \sim H & \text{with prob. } \frac{\alpha}{n-1+\alpha} \end{cases} \quad (7)$$

در رابطه فوق  $| \{j: c_j = i\}|$  بیان‌گر تعداد دفعاتی است که  $\theta_i^*$  در  $\theta_1, \dots, \theta_n$  رخداده است. رابطه فوق نشان می‌دهد که فرایند دیریکله اثر خوشبندی<sup>۵</sup> دارد. قسمت اول رابطه فوق نشان می‌دهد که یک داده مشاهده شده جدید با احتمال  $\frac{|\{j: c_j = i\}|}{n-1+\alpha}$  در یکی از خوشبها موجود قرار می‌گیرد یا با توجه به قسمت دوم رابطه به احتمال  $\frac{\alpha}{n-1+\alpha}$  منجر به ایجاد یک خوشة جدید می‌شود. این خاصیت مبنای یک خاصیت شناختی در فرایند دیریکله نیز است. این خاصیت پیاده‌سازی حافظه به‌وسیله توزیع احتمال است.

انتگرال گیری از پارامتر آزاد فرایند دیریکله که موجب ایجاد رابطه (5) می‌شود، یک نگاه جدید را به فرایند دیریکله مطرح می‌کند، که به آن فرایند رستوران چینی<sup>۶</sup> می‌گویند. در این نگاه، رابطه (5) نشان‌دهنده توزیع احتمال انتخاب میزهای رستوران وارد می‌شوند. این توزیع احتمال با ورود هر مشتری، به‌گونه‌ای تغییر می‌کند که احتمال انتخاب میزهای شلوغ‌تر توسط مشتریان بعدی افزایش پیدا می‌کند. این خاصیت را خاصیت "ثرومند، ثروتمندتر می‌شود"<sup>۷</sup> می‌نامند.

<sup>4</sup> Kronecker Delta Function

<sup>5</sup> Clustering

<sup>6</sup> Chinese Restaurant Process

<sup>7</sup> Rich-get-richer

احتمال‌هاست، یعنی یک رخداد نمونه‌گیری شده از این توزیع احتمال، خود یک توزیع احتمال است. توزیع احتمال‌های نمونه‌گیری شده از فرایند دیریکله توزیع احتمال‌های نایپیوسته (توزیع احتمال جسم  $n$  وجهی)، که  $n$  می‌تواند بی‌نهایت باشد) هستند. این توزیع‌های احتمال را نمی‌توان با تعداد محدودی از پارامترها توصیف کرد و از این‌رو فرایند *Teh* (2010). در حقیقت فرایند دیریکله توزیع احتمالی است که دامنه آن دنیای توزیع‌های احتمال روی یک متغیر تصادفی است.

گفتیم دامنه یک فرایند دیریکله توزیع‌های احتمال روی یک متغیر تصادفی است، بنابراین می‌توان یک توزیع احتمال نمونه‌گیری شده را از یک فرایند دیریکله یک توزیع احتمال تصادفی به حساب آورد.

فرض کنید که  $G$  یک توزیع احتمال تصادفی است؛ چیزی همانند یک متغیر تصادفی، که دارای توزیع فرایند دیریکله است. برای اینکه این چنین باشد، باید توزیع‌های حاشیه‌ای<sup>۸</sup> روی این توزیع احتمال تصادفی دارای توزیع دیریکله<sup>۹</sup> باشند (*Teh, 2010*). برای توضیح این موضوع، فرض کنید که  $H$  یک توزیع روی متغیر تصادفی پیوسته  $\Theta$  است و  $\alpha$  یک عدد حقیقی مثبت است. برای هر افزار متناهی  $R$  که به‌صورت  $A_1, \dots, A_R$  نشان می‌دهیم، بردار  $G(A_1), \dots, G(A_R)$  یک بردار تصادفی است، چرا که  $G$  یک توزیع احتمال تصادفی توزیع شده بر حسب فرایند دیریکله با توزیع مبنای  $H$  است، اگر داشته باشیم:

$$[G(A_1), \dots, G(A_R)] \sim \text{Dirichlet}(\alpha H(A_1), \dots, \alpha H(A_R)) \quad (3)$$

که برای هر افزار متناهی  $(A_1, \dots, A_R)$  از  $\Theta$  برقرار باشد (*Teh, 2010*)

این یعنی توزیع‌های حاشیه‌ای روی توزیع  $G$  دارای توزیع دیریکله باشند. در این صورت می‌گوییم  $G$  دارای توزیع فرایند دیریکله است و نشان می‌دهیم:

$$G \sim DP(\alpha, H) \quad (4)$$

حال فرض کنید که  $\theta_1, \dots, \theta_n$  یک ترتیب نمونه‌گیری شده مستقل با توزیع یکسان<sup>۱۰</sup> (i.i.d.) از  $G$  است.

<sup>1</sup> Marginal Distributions

<sup>2</sup> Dirichlet Distribution

<sup>3</sup> Independent and Identically distributed



فرض کنید یکتابع داریم که با  $f(z)$  آن را نشان می‌دهیم. در صورتی که  $z$  یک متغیر پیوسته باشد رابطه امید ریاضی تابع  $(z)$  به صورت زیر بیان می‌شود:

$$E[f(z)] = \int f(z)P(z)dz \quad (8)$$

در این رابطه،  $(z)$   $P$  توزیع احتمال پیشین روی  $z$  است. این انتگرال را می‌توان با استفاده از رابطه:

$$E_{P(z)}[f(z)] \approx \frac{1}{N} \sum_{i=1}^N f(z_i) \quad (9)$$

به شرط اینکه  $z_i$  ها از توزیع احتمال  $P(z)$  تولید کرده باشیم و تعداد آنها به اندازه کافی زیاد باشد، تخمین زد. این نمونه‌ها را می‌توان از طریق زنجیره مارکوف تولید کرد.

تخمین فوق که در آن تعداد زیادی از نمونه‌های تولید شده از یک تابع برای تخمین آن به کار می‌رود، بیان گر بخش مونت کارلوی این روش‌ها است و تولید نمونه‌ها براساس توزیع احتمال  $P(z)$  به سیله یک زنجیره مارکوف انجام می‌شود. از این‌رو، این روش‌ها را مونت کارلوی زنجیره مارکوف می‌نامند.

حال دقّت کنید که توزیع احتمال یک متغیر تصادفی را می‌توان به صورت توزیع حاشیه‌ای اشتراک آن متغیر با یک متغیر تصادفی دیگر نوشت. به عنوان مثال رابطه زیر را در نظر بگیرید.

$$P(x) = \int P(x|\theta)P(\theta)d\theta \quad (10)$$

رابطه فوق می‌تواند همانند رابطه (۹) یک رابطه امید ریاضی محسوب شود. بنابراین می‌توان توزیع احتمال  $(x)$  را همانند امید ریاضی  $E_{P(z)}[f(z)]$  تخمین زد. بنابراین از روش‌های مونت کارلوی زنجیره مارکوف می‌توان برای تخمین توزیع‌های احتمال بهره گرفت.

#### ۴-۲- نمونه‌برداری گیبس

روش نمونه‌برداری گیبس یکی از روش‌های مونت کارلوی زنجیره مارکوف است که در موقعي که توزیع احتمال مورد نظر، توزیع توأم تعداد زیادی متغیرهای تصادفی است، به کار گرفته می‌شود (Resnik, et al., 2010). ایده مبنایی این روش آن است که هر یک از متغیرهای تصادفی را جداگانه و البته به شرط بقیه متغیرهای تصادفی نمونه‌برداری کند احتمالی نظیر آنچه در زیر می‌بینید، نمونه‌برداری کند.

$$P(rv_i^t | rv_1^t, \dots, rv_{i-1}^t, rv_{i+1}^{t-1}, \dots, rv_k^{t-1}) \quad (11)$$

#### ۳-۲- مونت کارلوی زنجیره مارکوف

استفاده از روش استنتاج بیزی برای یافتن توزیع احتمال پسین<sup>۱</sup> در یک مسئله یادگیری، ممکن است منجر به یافتن توزیع احتمالی شود که به صورت تحلیلی قابل محاسبه نباشد. در این‌گونه موارد یک راه حل پرکاربرد استفاده از روش‌های مونت کارلوی زنجیره مارکوف برای تخمین تابع توزیع مورد نظر است. به عنوان مثال (Cohn, et al., 2010) و (Goldwater, et al., 2007 b) (Johnson, et al., 2007 b) نمونه‌هایی از کاربرد این روش برای تخمین تابع توزیع پسین را در خود دارند. برای تخمین یک توزیع احتمال، روش‌های مونت کارلوی زنجیره مارکوف برای تخدمین تابع تعداد زیادی نمونه از توزیع احتمال مورد نظر تولید کنند. این روش‌ها برای تولید نمونه از یک توزیع احتمال، یک زنجیره مارکوف می‌سازند که توزیع احتمال تولید شده توسط این زنجیره، معادل توزیع احتمال مورد نظر شود (Johnson, et al., 2007 b). به بیان دیگر، این روش‌ها احتمال جایه‌جایی در زنجیره مارکوف را به گونه‌ای تعیین می‌کنند که احتمال مشاهده حالت<sup>۲</sup>  $St_i$  در زنجیره، برابر با احتمال این حالت با توجه به توزیع احتمال مورد نظر باشد. پس از تولید زنجیره مورد نظر، این روش‌ها با جایه‌جایی روی حالت‌های این زنجیره، تعداد زیادی نمونه از توزیع مورد نظر تولید می‌کنند (Johnson, et al., 2007 b). البته زنجیره ایجاد شده به گونه‌ای است که توزیع احتمال آن به سمت توزیع احتمال مورد نظر میل می‌کند؛ بنابراین ابتدا لازم است تعداد زیادی نمونه از زنجیره مورد نظر تولید شده و دور ریخته شود. بدین ترتیب پس از اینکه توزیع احتمال زنجیره به توزیع احتمال هدف به اندازه کافی نزدیک شد، می‌توان نمونه‌های تولید شده را نمونه‌هایی از توزیع هدف تلقی کرد. با تولید نمونه‌های زیاد از توزیع احتمال هدف می‌توان این نمونه‌ها را به طرق مختلف برای محاسبه مقادیر مورد نظر به کار گرفت. به عنوان مثال می‌توان این نمونه‌ها را برای تخمین<sup>۳</sup> MAP توزیع احتمال هدف به کار گرفت (Goldwater, et al., 2007). یک روش دیگر آن است که این نمونه‌ها را برای محاسبه امید ریاضی روی تابع توزیع هدف به کار گرفت (Johnson, et al., 2007 b). حال سعی می‌کنیم با روابط ریاضی ایده مبنایی این روش‌ها را شرح دهیم.

<sup>1</sup> Posterior Distribution

<sup>2</sup> State

<sup>3</sup> Maximum a Posteriori

که نمونه‌های آخر تولید شده از توزیع احتمال هدف حول MAP این توزیع احتمال متمرکز شده‌اند.

در بخش بعد به تعریف رسمی صورت‌گرای استخراجی در این کار خواهیم پرداخت و به طور مختصر نکاتی را در مورد تئوری این صورت‌گرا مورد اشاره قرار خواهیم داد.

## ۵-۲- گرامر جایگزینی درخت احتمالاتی

همان‌طور که گفته شد، صورت‌گرای استخراجی در این کار گرامر جایگزینی درخت احتمالاتی است. این صورت‌گرا قوی‌تر از گرامر مستقل از متن احتمالاتی در توصیف زبان‌های طبیعی است. در این بخش، تعریف این صورت‌گرا و نحوه اشتقاق در آن را توضیح می‌دهیم تا زمینه را برای توصیف مدل مورد نظر برای استخراج این صورت‌گرا فراهم کنیم.

یک گرامر جایگزینی درخت احتمالاتی به صورت یک پنج تایی تعریف می‌شود (*Bod, 1995*):

$$G = \langle V_N, V_T, R, P, \sigma \rangle$$

که در آن  $V_N$  مجموعه نمادهای غیر پایانی،  $V_T$  مجموعه نمادهای پایانی،  $\sigma \in V_N$  غیرپایانه شروع و  $R$  مجموعه‌ای از قواعد تولید هستند که هر یک از این قوانین به صورت یک درخت هستند و به آنها درخت اولیه گفته می‌شود.  $P$  تابعی است که به هر درخت اولیه  $R \in \mathcal{R}$  یک احتمال  $P(t)$  می‌دهد. برای یک درخت  $t$  با ریشه  $\alpha$ ، از  $P(t)$  به عنوان احتمال جانشینی درخت  $t$  در غیرپایانه  $\alpha$  تعبیر می‌شود.

یک درخت اولیه، قطعه‌درختی به عمق  $\leq 2$  است که هر گره داخلی آن یک غیرپایانه بوده و گره‌های برگ در آن پایانه یا غیرپایانه هستند. برگ‌های غیرپایانه نقش مولد را در فرآیند تولید درختان در این گرامر ایفا می‌کنند. هر اشتقاق، درختی را با شروع از ریشه تولید کرده و هر برگ غیرپایانه آن را با یک درخت اولیه جایگزین می‌کند و این کار را تا زمانی ادامه می‌دهد که هیچ برگ غیرپایانه‌ای باقی نماند. برخلاف گرامرهای مستقل از متن، در این گرامر درخت‌های اشتقاق مختلف می‌توانند منجر به تولید یک درخت تجزیه شوند (*Bod, 1995*). این موضوع در شکل زیر نشان داده شده است. غیرپایانه‌هایی که با @ نشان داده شده‌اند محل جایگزینی درخت اولیه هستند.

در این رابطه،  $rv^t$  نشان‌دهنده متغیر تصادفی نام در مرحله آم است. همان‌گونه که از رابطه (11) مشخص است، هر متغیر تصادفی در مرحله آم به آخرین مقادیر نسبت داده شده به دیگر متغیرهای تصادفی وابسته است.

در این روش پس از اینکه از هر متغیر تصادفی یک نمونه ایجاد شد، مجموعه این نمونه‌ها در کنار هم یک نمونه از توزیع احتمال اشتراک متغیرها را تشکیل می‌دهد. کل فرایند به این صورت است که به ازای هر متغیر تصادفی یک توزیع چندجمله‌ای محاسبه می‌شود و سپس به روش Roulette Wheel یک مقدار برای متغیر تصادفی مورد نظر، از توزیع محاسبه شده نمونه‌برداری می‌شود. این کار به ازای تمامی متغیرهای تصادفی تکرار می‌شود تا به ازای هر متغیر تصادفی یک نمونه ایجاد شود. با تکرار این عملیات به اندازه کافی، می‌توان نمونه‌های زیادی از توزیع احتمال هدف تولید کرد که برای محاسبه مقادیر مختلف مانند امید ریاضی و

MAP توزیع احتمال هدف، به کار می‌رond.

برای یافتن MAP یک توزیع احتمال، می‌توان از روش‌های جستجوی نقاط بهینه توابع استفاده کرد. یک روش معمول برای محاسبه MAP یک توزیع احتمال، همراه‌سازی فرایند نمونه‌گیری با الگوریتم شبیه‌سازی ذوب فلز<sup>۱</sup> است. در حقیقت، این کار منجر به یک جستجوی تپه‌نوردی تصادفی<sup>۳</sup> می‌شود که برای یافتن نقطه بیشینه تابع توزیع احتمال هدف به کار می‌رود. این روش در (*Cohn, et al., 2010*) (*Goldwater, et al., 2007*) به کار گرفته شده است. این کار با الگوریتم نمونه‌برداری گیبس به راحتی قابل انجام است. برای این منظور یک متغیر دما  $T$  تعریف می‌کنند. سپس در هنگام نمونه‌برداری از توزیع احتمال هر متغیر تصادفی به شرط متغیرهای تصادفی دیگر، این توزیع احتمال را به توان  $T/1$  می‌رسانند. متغیر دما در هر دور نمونه‌برداری از تمامی متغیرهای تصادفی، ثابت نگه داشته می‌شود. اما در دور بعدی با یک گام مشخص، کاهش پیدا می‌کند. این کار را تا جایی ادامه می‌دهند که متغیر دما به صفر برسد. اثر این کار آن است که در الگوریتم تپه‌نوردی، در مراحل اولیه، حالت‌های بدتر شانس بیشتری برای تولید می‌یابند و در ادامه این شانس کمتر و کمتر می‌شود. در (*Goldwater, et al., 2007*) اشاره شده است که اگر این تغییر دما به اندازه کافی کند صورت گیرد، تضمین می‌شود

<sup>1</sup> Simulated Annealing

<sup>3</sup> Stochastic hill climbing



شده می‌تواند ترکیب‌های مختلفی از برتتاب تاس‌ها تولید شده باشد. بنابراین در هر تجزیهٔ نحوی ممکن، از داده‌ها تعداد محدودی از تاس‌های موجود شرکت دارند. حال اگر توزیع احتمال تاس‌هایی را که در تولید داده‌ها شرکت داشته‌اند، بیابیم توزیع احتمال  $e$  را یافته‌ایم. تصمیم‌گیری راجع به توزیع احتمال هر یک از تاس‌ها یک مسئلهٔ انتخاب مدل است. اگر بخواهیم برای تصمیم‌گیری راجع به توزیع احتمال هر یک از تاس‌ها از مدل بیزی از فرایند شناخت انتخاب مدل استفاده کنیم، باید از تعمیم یافتهٔ این مدل برای تاس بی‌نهایت وجهی استفاده کنیم. همان‌طور که در قبل هم اشاره کردیم، این تعمیم همان فرایند دیریکله است. بنابراین حل مسئلهٔ استخراج گرامر جایگزینی درخت، به صورت بامتری منجر به تعداد زیادی فرایند دیریکله می‌شود که هر کدام، یک مسئلهٔ انتخاب مدل را برای ما حل می‌کند.

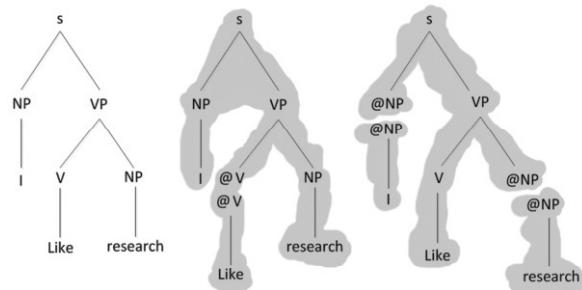
بنابرآنچه در بند قبل توضیح دادیم، به ازای ترکیب هر غیرپایانه و پدرش یک فرایند دیریکله خواهیم داشت. اگر فرض کنیم که  $e'$  نشان‌دهندهٔ تمامی زیردرخت‌هایی باشد که گره ریشه و پدر ریشه مشترک دارند، می‌توانیم مدل را به صورت زیر بیان کنیم:

$$G_{N,p(N)} | \alpha_{N,p(N)}, H_{N,p(N)}, E \sim DP(\alpha_{N,p(N)}, H_{N,p(N)}) \quad (12-1)$$

$$e' | N, p(N), E \sim G_{N,p(N)} \quad (12-2)$$

که در آن  $N$  نشان‌دهندهٔ یک غیرپایانه و  $p(N)$  نشان‌دهندهٔ پدر آن است.  $H_{N,p(N)}$  توزیع پایهٔ فرایند دیریکله  $\alpha_{N,p(N)}$  پارامتر تمرکز این فرایند است.  $G_{N,p(N)}$  نیز نشان‌دهندهٔ متغیر تصادفی از نوع توزیع احتمال است که به ازای ترکیب هر غیرپایانه و پدرش یکی وجود دارد. اگر از دیدگاه مدل استخراج گرامر بخواهیم توصیف کنیم،  $G_{N,p(N)}$  تمامی توزیع‌های ممکن برای قوانینی است که دارای ریشه  $N$  دارند؛ و پدر ریشه آنها  $(N)p$  بوده است. مجموعه  $G_{N,p(N)}$  ها به ازای تمامی  $N$  و  $(N)p$ ‌های موجود، یک متغیر تصادفی است که دامنه آن تمامی گرامرهای ممکن است. حال اگر همانند آنچه در رابطه با مدل بیزی از فرایند شناخت انتخاب مدل در حل مسئلهٔ سکه انجام دادیم از پارامتر آزاد مدل انتگرال بگیریم، انتگرالی شبیه رابطه (۱) حاصل خواهد شد. این انتگرال گیری طبق فرمول (۶) در این مسئله رابطه زیر را نتیجه خواهد داد:

$$P(e'_i | e'_{-i}, N, p(N), \alpha_{N,p(N)}, H_{N,p(N)}, E) = \frac{\alpha_{N,p(N)} H(e'_i | N, p(N)) + \text{cnt}_{-i}(e'_i | N, p(N))}{\alpha_{N,p(N)} + \text{cnt}_{-i}(e'' | N, p(N))} \quad (13)$$



(شکل ۱): دو اشتراق مختلف از یک درخت تجزیه. در سمت راست دو اشتراق مختلف را که منجر به تولید درخت تجزیه سمت چپ می‌شوند، مشاهده می‌کنید. گره‌های جایگزینی با  $@$  علامت‌گذاری شده‌اند.

### ۳- مدل

یک پیکره از جملات تجزیه شده را در نظر بگیرید، آن را  $E$  می‌نامیم. فرض کنید  $e$  نیز نشان‌دهندهٔ ترتیب درخت‌های پایه‌ای است که درخت‌های پیکره  $E$  را تجزیه می‌کنند. ما در این مسئله می‌خواهیم بهترین  $e$  که پیکره مورد نظر را تجزیه می‌کند، پیدا کنیم. برای این منظور می‌توان توزیع احتمال  $e$  به شرط داده‌های مشاهده شده را یافت و از این توزیع احتمال استفاده کرد تا بهترین  $e$  که پیکره  $E$  را تجزیه می‌کند، به دست آوریم. به عنوان مثال، می‌توانیم  $MAP$  این توزیع احتمال را با استفاده از روش‌های موجود بیابیم. در ادامه راجع به این موضوع بیشتر توضیح خواهیم داد.

می‌توان نشان داد که مسئلهٔ استخراج گرامر جایگزینی درخت شبیه به مثال سکه‌ای است که در بخش ۱-۲ برای توصیف مدل بیزی از فرایند شناختی انتخاب مدل، ارائه شد و در حقیقت ترکیب پیچیده‌ای از این مسئله است.

فرض کنید ترکیب هر غیرپایانه و پدرش معادل یک تاس  $n$  وجهی است. پس به تعداد ترکیب همهٔ غیرپایانه‌ها و پدرشان تاس خواهیم داشت. در مسئلهٔ استخراج گرامر جایگزینی درخت، هر غیرپایانه به بی‌نهایت شکل ممکن می‌تواند گسترش یابد. بنابراین ترکیب هر غیرپایانه و پدر آن نیز به بی‌نهایت شکل ممکن می‌تواند گسترش یابد. درنتیجه، هر تاس ما یک تاس بی‌نهایت وجهی است. در حقیقت هر تاس به این شکل است که با ترکیب یک غیرپایانه و پدرش نام‌گذاری شده و روی هر وجه آن یک درخت اولیه ممکن چاپ شده است. با توجه به این نکته که هر درخت تجزیه می‌تواند با درخت‌های اشتراق مختلفی تولید شده باشد، می‌توان نتیجه گرفت که داده‌های مشاهده

مدل ارائه شده، موجب می‌شود توزیع احتمال فوق را بتوان توزیع اشتراک تاس‌ها دانست. یعنی اگر تاس‌ها را متغیرهای تصادفی فرض کنیم، توزیع احتمال فوق توزیع اشتراک این متغیرهای تصادفی است. همان‌طور که در بخش دو اشاره کردیم، این‌گونه توزیع‌ها را که توزیع احتمال اشتراک تعداد زیادی متغیرهای تصادفی است، با استفاده از روش نمونه‌برداری گیبس که یک نسخه از روش‌های مونت کارلوی زنجیره مارکوف است، می‌توان تخمین زد. برای این کار باید توزیع یک متغیر را به شرط متغیرهای دیگر محاسبه کنیم و هر بار از یکی از متغیرها به شرط بقیه متغیرها نمونه تولید کنیم تا پس از تولید کرده باشیم (*Resnik, et al., 2010*)<sup>1</sup>. بدین منظور ما از مدل ارائه شده در (*Cohn, et al., 2009*)<sup>2</sup> تبعیت کردیم. بدین شکل که فرض کنید برای هر گره غیر ریشه در پیکره، یک متغیر باینری در نظر می‌گیریم که مشخص کند این گره محل جایگزینی است یا نه. در حقیقت با درنظر گرفتن تاس‌هایی که گفتیم، مجموعه این متغیرها تعیین می‌کنند که کدام تاس‌ها باید در این مرحله پرتاب شوند.

فرض کنید برای دور اول نمونه‌برداری گیبس را به صورت تصادفی، گره‌هایی که محل جایگزینی هستند، تعیین کرده‌ایم. یعنی به طور تصادفی تاس‌هایی را که در دور اول نمونه‌برداری گیبس پرتاب می‌شوند، تعیین کنیم. حال نمونه‌بردار گیبس ما روی تمامی گره‌هایی که در مرحله قبل برچسب جایگزینی نخورده‌اند، حرکت کرده و با تولید نمونه تعیین می‌کند که گره مورد بررسی، محل جایگزینی بشود یا نه. این طریقه عملکرد سبب می‌شود که عملیات پرتاب هر تاس نیز در دل فرایند نمونه‌برداری گیبس انجام شود. این بدان معناست که فرایند پرتاب یک تاس پس از نمونه‌برداری از چند گره صورت می‌گیرد. به طوری که زمانی فرایند پرتاب یک تاس کامل و یک وجه آن ظاهر می‌شود، که در تمامی مسیرها از یک گره جایگزینی به سمت پایین درخت، یک گره جایگزینی بهوسیله نمونه‌برداری تعیین شود. یعنی مرزهای یک درخت اولیه تعیین گردد.

حال فرض کنید که می‌خواهیم از متغیر تصادفی مرتبط با گره  $N$  نمونه‌گیری کنیم که آیا محل جایگزینی بشود یا نه. همان‌طور که از جمله قبل مشخص است و قبلاً از این هم گفتیم این متغیر تصادفی یک متغیر باینری است. برای نمونه‌گیری ابتدا باید توزیع احتمال آن را در این مرحله به شرط بقیه متغیرها تعیین کرده و سپس از آن نمونه‌گیری

در رابطه فوق  $e^i$  نمونه زیردرخت  $\alpha^i$  است که از فرایند دیریکله تولید می‌شود.  $e^i$  بردار  $\mathbf{1} - \mathbf{i}$  نمونه قبلي است که از فرایند دیریکله نمونه‌گیری شده‌اند.  $(\mathbf{e}^i | N, p(N))$  نیز تابعی است که نشان‌دهنده تعداد رخداد زیردرخت  $e^i$  با گره ریشه  $N$  و پدر ریشه  $p(N)$  است که در  $\mathbf{1} - \mathbf{i}$  نمونه‌گیری قبلی تولید شده است. نشان‌دهنده تمامی زیردرخت‌های تولید شده است. نتیجه انتگرال گیری از پارامتر آزاد مدل در فرایند دیریکله، فرایند رستوران چینی را نتیجه خواهد داد. یعنی رابطه فوق رابطه توزیع احتمال انتخاب میز در فرایند رستوران چینی است. به بیان دیگر این توزیع احتمال همان توزیع احتمالی است که پیاده‌ساز حافظه برای ذخیره‌سازی واحدهای پرکاربرد است. اثر برچسب‌زننده تاس در مدل بالا با گره ریشه و پدر ریشه مبتنی بر سابقه شدن مدل نتیجه شده خواهد بود. اما با این کار تعداد تاس‌های در گیر در مسئله بیشتر خواهد شد. یعنی از دیدگاه فرایند رستوران چینی تعداد رستوران‌ها افزایش پیدا کرده است. در نتیجه این تغییر تعداد میزهایی که در هر رستوران پر می‌شود کمتر خواهد بود. این مسئله باعث می‌شود در رابطه بالا مخرج کسر کوچک‌تر و فاصله بین احتمال‌های کوچک و بزرگ بیشتر شده و در نتیجه یادگیری یا ذخیره‌سازی تصادفی بهتر انجام گیرد. یعنی خاصیت "ثروتمند ثروتمندتر می‌شود" بدین ترتیب تشدید می‌شود.

## ۴- فرایند یادگیری مدل

همان‌طور که در بخش قبل نیز اشاره شد، هدف در مدل ارائه شده، یافتن توزیع احتمال  $P(e|E)$  است. این توزیع احتمال با توجه به تعریف مدل، به صورت زیر تغییر خواهد کرد:

$$P(e|E, \alpha, H)$$

علت این تغییر آن است که پارامترهای فرایندهای دیریکله ( $\alpha_{N,p(N)}$ ,  $H_{N,p(N)}$ ) متغیرهایی هستند که باید مقادیر آنها در مدل تعیین شود.

$\alpha$  و  $H$  برداری از تمامی پارامترهای مرکز و توزیع اولیه برای تمام فرایندهای دیریکله هستند. این توزیع احتمال توزیع پیچیده‌ای است که راه حل ریاضی برای محاسبه تحلیلی آن نداریم. اما می‌توانیم با استفاده از راه حل‌های موجود برای تخمین انتگرال‌ها، نظریه مونت کارلوی زنجیره مارکوف، این توزیع را تخمین بزنیم.



برای محاسبه روابط ارائه شده در این بخش باید پارامترهای فرایندهای دیریکله تعیین شده باشد. ما هر دو پارامتر فرایند دیریکله ( $H_{N,p(N)}$ ,  $\alpha_{N,p(N)}$ ) را به تبعیت از (Cohn, et al., 2009) تعیین کردیم. به این ترتیب که  $H$  را تابع تصادفی در نظر گرفتیم که روی هر گره، حرکت کرده و یک سکه پرتاب می‌کند و براساس نتیجه سکه تصمیم می‌گیرد که گره مذکور گسترش پیدا کند یا نه. هر متغیر  $\alpha$  نیز از توزیع گاما ( $0.001, 1000$ ) نمونه‌گیری می‌شود.

فرایند نمونه‌برداری گیبس باید تعداد زیادی نمونه از توزیع احتمال هدف تولید کند تا بتوان با استفاده از این نمونه‌ها تخمین مناسبی از توزیع احتمال هدف تولید کرد. طبق آنچه در (Goldwater, et al., 2007) بیان شده، اگر فرایند نمونه‌برداری گیبس را با ذوب فلز شبیه‌سازی شده<sup>2</sup>، همراه کنیم و فرایند کاهش دما را به اندازه کافی گند انجام دهیم، پس از تولید نمونه‌های زیادی از توزیع احتمال هدف، نمونه‌های تولید شده حول نقطه MAP توزیع احتمال هدف متمرکز می‌شوند. بدین ترتیب نمونه تولید شده را در مرحله آخر، نمونه MAP به حساب آورده. این نمونه تولید شده نهایی همان تجزیه نهایی پیکره یادگیری است که از قوانین گرامر استخراجی تشکیل شده است.

در صورتی که در مدل ارائه شده در این تحقیق فرایند دیریکله به فرایند پیتمن-یور تغییر یابد، نسخه مبتنی بر سابقه مدل ارائه شده در (Cohn, et al., 2010) به دست خواهد آمد. به همین دلیل ما از پیاده‌سازی انجام شده برای کار (Cohn, et al., 2010) استفاده کردیم تا مدل خود را ایجاد کنیم. برای ایجاد مدل خود ما فرایند پیتمن-یور پیاده‌سازی شده را با تغییر پارامترهای فرایند به فرایند دیریکله تبدیل کردیم. سپس با روش برچسب‌گذاری پدر که در (Johnson, 1998) ارائه شده است پیاده‌سازی مذکور را به پیاده‌سازی مدل خود تبدیل کردیم.

با توجه به اینکه الگوریتم یادگیری، در هر تکرار به ازای هر گره تصمیم می‌گیرد که گره مذکور یک گره داخلی است یا یک گره مرزی، هزینه اجرای الگوریتم یادگیری  $n \times n$  خواهد بود. در اینجا  $n$  نشان‌دهنده تعداد تکرار الگوریتم و  $n$  نمایانگر مجموع تعداد گره‌های درخت‌های تجزیه نهایی موجود در پیکره آموزشی است.

کنیم. بسته به مقدار تعیین شده برای این متغیر ممکن است دو حالت اتفاق بیفتد. یکی اینکه گره مرتبط با این متغیر، محل جایگزینی انتخاب نشود و درنتیجه باید گره داخلی یک درخت  $e_m$  شود. دیگری اینکه این گره محل جایگزینی تعیین شده و مرز دو درخت  $e_a$  و  $e_b$  شود (Cohn, et al., 2009) پس باید احتمال رخداد این دو حالت را حساب کرده و یک توزیع دوجمله‌ای<sup>1</sup> به دست آوریم و از آن نمونه تولید کنیم. این توزیع به صورت زیر خواهد بود:

$$P(e_m | N_{e_m}, p(N_{e_m}), e_{-i}, \alpha_{N_{e_m}, p(N_{e_m})}, H_{N_{e_m}, p(N_{e_m})}, E) \quad (14)$$

$$= \frac{H(e_m | N_{e_m}, p(N_{e_m})) \alpha_{N_{e_m}, p(N_{e_m})} + cnt_{-i}(e_m | N_{e_m}, p(N_{e_m}))}{cnt_{-i}(e'' | N_{e_m}, p(N_{e_m})) + \alpha_{N_{e_m}, p(N_{e_m})}} \quad (15)$$

$$\begin{aligned} & P(e_a, e_b | N_{e_a}, p(N_{e_a}), e_{-i}, \alpha_{N_{e_m}, p(N_{e_m})}, H_{N_{e_m}, p(N_{e_m})}, E) \\ & = \frac{H(e_a | N_{e_a}, p(N_{e_a})) \alpha_{N_{e_a}, p(N_{e_a})} + cnt_{-i}(e_a | N_{e_a}, p(N_{e_a}))}{cnt_{-i}(e'' | N_{e_a}, p(N_{e_a})) + \alpha_{N_{e_a}, p(N_{e_a})}} \\ & \times \frac{H(e_b | N_{e_b}, p(N_{e_b})) \alpha_{N_{e_b}, p(N_{e_b})} + cnt_{-i}(e_b | N_{e_b}, p(N_{e_b})) + \delta(e_a, e_b)}{cnt_{-i}(e'' | N_{e_b}, p(N_{e_b})) + \alpha_{N_{e_b}, p(N_{e_b})} + \delta((N_{e_b}, p(N_{e_b})), (N_{e_a}, p(N_{e_a})))} \end{aligned}$$

تمامی درخت‌های اولیه‌ای هستند که درنهایت  $e_{-i}$

یک گره غیرپایانه با  $e_a$ ،  $e_b$  یا  $e_m$  اشتراک دارند.

رابطه (14) همان رابطه (13) است که متغیرهای

آن با متغیرهای جدید جایگزین شده‌اند. همچنین به راحتی قابل مشاهده است که رابطه (15) نیز با استفاده از رابطه (13) و رابطه زیر که قانون زنجیری است، به دست می‌آید:

$$P(e_a, e_b | N_{e_a}, p(N_{e_a}), e_{-i}, \alpha_{N_{e_m}, p(N_{e_m})}, H_{N_{e_m}, p(N_{e_m})}, E) =$$

$$P(e_a | N_{e_a}, p(N_{e_a}), e_{-i}, \alpha_{N_{e_m}, p(N_{e_m})}, H_{N_{e_m}, p(N_{e_m})}, E) \times$$

$$P(e_b | e_a, N_{e_a}, p(N_{e_a}), e_{-i}, \alpha_{N_{e_m}, p(N_{e_m})}, H_{N_{e_m}, p(N_{e_m})}, E)$$

فرایند استخراج گرامر در ادامه به این شکل خواهد بود که به ازای هر گره، که در گام قبلی به عنوان گره جایگزینی انتخاب شده، روابط بالا محاسبه می‌شود و از آنها نمونه تولید می‌شود. هر نمونه از توزیع‌های فوق یک درخت اوّلیه (قانون گرامر استخراجی) را نتیجه می‌دهد. بدین ترتیب پس از هر دور حرکت روی پیکره و تولید نمونه از توزیع‌های فوق، به ازای هر گره جایگزینی، یک تجزیه نحوی از کل جملات پیکره خواهیم داشت. در مرحله پایانی این فرایند تکراری تجزیه نهایی جملات پیکره تولید خواهد شد. این تجزیه نهایی از درخت‌های اوّلیه‌ای ساخته شده است که در حقیقت قوانین گرامر نهایی هستند.

<sup>1</sup> Binomial

<sup>2</sup> Simulated Annealing

## ۵- نتایج و ارزیابی

می‌گیرد. از آنجایی که ما در نظر داشتیم نتایج خود را با نتایج ارائه شده در (Cohn, et al., 2010) مقایسه کنیم، تصمیم گرفتیم از تجزیه‌گرهای ارائه شده در این کار برای تجزیه جملات بخش آزمایشی پیکره استفاده کنیم تا الگوریتم تجزیه‌گر، به عنوان یک عامل تغییر در نتایج تجزیه (Cohn, et al., 2010) ابتدا گرامر جایگزینی درخت استخراج شده را، با استفاده از روش کاهش به گرامر مستقل از متن احتمالاتی<sup>۶</sup> (Goodman, 2002)، به گرامر مستقل از متن احتمالاتی تبدیل می‌کنند. در این فرایند برخی قوانین که احتمال رخداد کمتری دارند حذف می‌شوند. سپس گرامر مستقل از متن احتمالاتی تولید شده، به وسیله تجزیه‌گر ارائه شده در (Johnson, et al., 2007 b) برای تولید تعداد زیادی نمونه از فضای احتمال درخت‌های تجزیه به کار گرفته می‌شوند. پس از آن، این نمونه تجزیه‌ها با استفاده از یک پذیرنده که بر اساس الگوریتم متروپلیس-هستینگس<sup>۷</sup> (Cohn, et al., 2009) و از گرامر جایگزینی درخت ایجاد شده است، قبول یا رد می‌شوند.

## ۲-۵- آزمایش

برای این آزمایش مدل ارائه شده در این تحقیق و مدل ارائه شده در (Cohn, et al., 2010) را روی بخش دو تا ۲۱ بخش WSJ پیکره Penn آموزش دادیم. برای این کار تعداد تکرار فرایند نمونه‌برداری ۵۰۰۰ تکرار تنظیم کردیم و الگوریتم شبیه‌سازی ذوب فلز را در ۴۰۰۰ تکرار اولیه به کمترین دمای خود رساندیم. در این فرایند دما از بیشینه پنج به مقدار کمینه یک کاهش یافت. سپس گرامر نتیجه شده از مرحله آموزش را به وسیله تجزیه‌گرهای گزارش شده در (Cohn, et al., 2010) برای تجزیه نحوی بخش ۲۲ از WSJ به کار گرفتیم. دلیل ما برای انتخاب بخش ۲۲ به جای بخش ۲۳، برای آزمایش، آن بود که بتوانیم نتایج را با نتایج گزارش شده در (Cohn, et al., 2010) مقایسه کنیم، چرا که آنها بخش ۲۲ را برای آزمایش بکار گرفته‌اند.

برای آزمایش مدل ارائه شده باید گرامر حاصل از آموزش مدل را بر روی یک پایگاه درختی، در عملیات تجزیه نحوی جملات به کار گرفت. از آنجایی که استخراج گرامر، سابقه طولانی در حوزه پردازش زبان طبیعی دارد، منبع استاندارد برای آزمایش تجزیه‌گرهای نحوی در زبان انگلیسی وجود دارد.

علاوه بر این، براساس معیارهای استاندارد ارزیابی تجزیه نحوی<sup>۸</sup> (Harrison, et al., 1991) (Collins, 1996) یک ابزار استاندارد به نام<sup>۹</sup> EVALB<sup>۱۰</sup> برای ارزیابی تجزیه نحوی جملات ایجاد شده است، که امروزه تمامی تجزیه‌گرهای نحوی از این ابزار برای ارزیابی نتایج خود بهره می‌گیرند. معیارهای استاندارد ارزیابی تجزیه نحوی شامل دقت پرانتزگذاری با برچسب<sup>۱۱</sup>، فراخوانی پرانتزگذاری با برچسب<sup>۱۲</sup>، معیار F و تقاطع پرانتزگذاری<sup>۱۳</sup> است.

منبع استاندارد برای آموزش و آزمایش تجزیه‌گرهای نحوی در زبان انگلیسی پیکره (Marcus, et al., Penn Charniak, et al., 1993) است. همانند کارهای ارائه شده در (Charniak, et al., 2005; Charniak, 2000; Cohn, et al., 2009; Cohn, et al., 2010) بخش WSJ در این پیکره برای آموزش و آزمایش تجزیه‌گرهای نحوی به شکلی که در جدول زیر مشاهده می‌کنید، به صورت استاندارد تقسیم‌بندی شده است.

(Cohn, et al., 2010): تقسیم‌بندی بخش WSJ از پیکره Penn

Partition	Sections	Sentences
Training	2-21	33180
Development	22	1700
Testing	23	2416

برای ارزیابی عملکرد مدل، باید آن را روی بخش‌های آموزشی پیکره آموزش داده و روی بخش ارزیابی، آزمایش کرد. برای آزمایش، گرامر استخراج شده توسط مدل، به وسیله یک تجزیه‌گر نحوی به کار گرفته می‌شود تا جملات بخش ارزیابی پیکره، تجزیه نحوی شوند. سپس نتیجه این تجزیه نحوی براساس معیارهای استاندارد ارزیابی EVALB<sup>۱۴</sup> و با استفاده از ابزار PARSEVAL<sup>۱۵</sup> که پیاده‌سازی استاندارد این معیارهای است، مورد ارزیابی قرار

<sup>1</sup> PARSEVAL Measures

<sup>2</sup> <http://nlp.cs.nyu.edu/evalb/>

<sup>3</sup> Labeled Precision

<sup>4</sup> Labeled Recall

<sup>5</sup> Cross Bracketing



کاهش قوانین با ریشه مشترک با توجه به فرمول (۱۳) موجب اختلاف بیشتر احتمال‌های کوچک و بزرگ برای میزهای (درخت‌های اولیه) هر رستوران می‌شود. همین مسئله، آنطور که در قبل هم گفته شد، موجب تشدید خاصیت "ثروتمند، ثروتمندتر می‌شود" می‌گردد. از دیدگاه شناختی ارائه شده در (O'Donnell, et al., 2009)

ذخیره‌سازی تصادفی در مدل ما نسبت به مدل گهن و همکاران تقویت شده است و قطعات پر تکرار بهتر تشخیص داده شده و به حافظه سپرده شده‌اند.

به منظور اطمینان از عملکرد بهتر مدل ارائه شده و اعتباربخشی بیشتر به ارزیابی مدل، گرامر استخراجی برای تجزیه نحوی بخش استاندارد پیکره Penn (بخش ۲۳) نیز به کار گرفته شد. بدین ترتیب می‌توان نتایج حاصل از این مدل را با مدل‌های بیشتری نیز مقایسه کرد. البته بسیاری از مدل‌های دیگر استخراج گرامر و تجزیه نحوی سرعت بسیار بالاتری از مدل ارائه شده در این مقاله در هنگام تجزیه نحوی جملات فراهم می‌کنند، و سرعت بسیار پایین این مدل در هنگام تجزیه نحوی جملات ضعفی است که حاصل هزینه اجرایی بالای الگوریتم تجزیه نحوی است. هزینه اجرایی این الگوریتم با استفاده از گرامر استخراجی توسط مدل ارائه شده، در اوخر همین بخش محاسبه و گزارش شده است. در (جدول ۴) نتایج آزمایش گرامر حاصل از مدل ارائه شده در تجزیه نحوی بخش استاندارد پیکره Penn (بخش ۲۳) ارائه شده است.

(جدول ۴): درصد به دست آمده برای معیار F1 برای مدل بیزی شناختی که در این کار معرفی شده و مدل ارائه شده در مقاله (Cohn, et al., 2010) برای تجزیه بخش ۲۳ پیکره Penn. سطر آخر اعداد گزارش شده برای گرامر مستقل از متن احتمالاتی است که با روش مبتنی بر درست‌نمایی بیشینه استخراج شده است.

مطابقت کامل	فراخوانی	دقت	F1	روش
50.75	91.87	91.76	91.82	مدل بیزی شناختی [این مقاله]
34.72	86.59	86.44	86.52	مدل گهن و همکاران آزمایش شده توسط نگارنده

(جدول ۲): درصد به دست آمده برای معیار F1 برای مدل بیزی شناختی که در این کار معرفی شده و مدل ارائه شده در مقاله (Cohn, et al., 2010) برای تجزیه بخش ۲۲ پیکره Penn. سطر آخر اعداد گزارش شده برای گرامر مستقل از متن احتمالاتی است که با روش مبتنی بر درست‌نمایی بیشینه استخراج شده است.

مطابقت	فراخوانی	دقت	F1	روش
کامل				
52.92	92.99	92.94	92.9	مدل بیزی شناختی [این مقاله]
26.2	-	-	83.8	گهن و همکاران (Cohn, et al., 2010)
6.7	-	-	63.1	(Cohn, MLE PCFG et al., 2010)

ما مدل خود را با تجزیه گر<sup>۱</sup> (Cohn, et al., 2010) مورد آزمایش قرار دادیم. مطابق با آنچه در رابطه با تئوری این تجزیه گر در (Cohn, et al., 2010) بیان شده است، این تجزیه گر با استفاده از روش‌های مونت کارلوی زنجیره مارکوف، احتمال درخت‌های تجزیه ممکن برای جمله را تخمین می‌زند و محتمل‌ترین درخت تجزیه را برمی‌گزیند. همان‌طور که در ۰ مشاهده می‌شود، مقدار معیار F1 حاصل از مدل ما، با اختلاف بنسپه قابل توجهی بیشتر از مقدار حاصل شده از مدل گهن و همکاران، با استفاده از تجزیه گر MPT است. این موضوع را این‌گونه تحلیل می‌کنیم که با افزایش تعداد رستوران‌ها در مدل ما پیچیدگی توزیع احتمال هر رستوران کاهش پیدا کرده است یعنی تعداد میزهای کمتری در هر رستوران پر شده است. بررسی قوانین گرامر استخراج شده نیز این مسئله را تأیید می‌کند. چرا که قوانینی (درخت‌های اولیه) که دارای ریشه مشترک هستند، در گرامر استخراجی، ما به طور متوسط بسیار کمتر از قوانین با ریشه مشترک در گرامر استخراجی مدل گهن و همکاران است. (۰ نشان‌دهنده میانگین تعداد قوانین با ریشه مشترک در گرامرهای استخراجی توسط مدل ما و مدل گهن و همکاران است.

(جدول ۳): میانگین تعداد قوانین با ریشه مشترک

مدل	میانگین قوانین با ریشه مشترک
مدل بیزی شناختی [این مقاله]	6.35
گهن و همکاران (Cohn, et al., 2010)	24.46

<sup>۱</sup> Maximum Probability Tree

تنکی داده یادگیری نیز بشود و بدین ترتیب میزان یادگیری کاهش یابد.

آزمایش‌های انجام شده در این تحقیق، از دیدگاه زمان، آزمایش‌هایی بسیار پر هزینه هستند. دلیل این امر هزینه بالای زمانی الگوریتم تجزیه نحوی MPT، بهویژه پس از دخیل کردن اطلاعات سابقه، است. هزینه زمانی اجرای این الگوریتم، به صورت زیر بیان می‌شود:

$$O(N^3 n^3 + iN^2 n^2)$$

که در آن  $N$  تعداد غیرپایانه‌های گرامر،  $n$  طول جمله و  $i$  تعداد نمونه‌های است که در فرایند مونت کارلوی زنجیره مارکوف از درخت تجزیه یک جمله تولید می‌شود. پس از دخیل کردن اطلاعات سابقه در فرایند یادگیری، هزینه تجزیه، به صورت زیر تغییر می‌کند:

$$O(N^4 n^3 + iN^2 n^2)$$

همان‌طور که مشاهده می‌شود، ترتیب زمانی تجزیه نحوی جملات در این حالت  $N$  برابر شده است. در آزمایش‌های انجام شده در این تحقیق مقدار  $N$  برابر با ۴۵ است. این موضوع سبب شد در آزمایش‌های انجام شده، زمان تجزیه نحوی جملات پس از دخیل کردن اطلاعات سابقه، به شدت افزایش یابد. بنابراین، برای تجزیه نحوی جملات بخش ۲۲ بانک درختی Penn در زمانی منطقی، مجبور به اجرای موازی الگوریتم روی حدود ۳۸۰ پردازش‌گر موازی شدیم.

## ۶- نتیجه‌گیری

در این مقاله، ما یک مدل بیزی غیرپارامتری برای استخراج گرامر مبتنی بر سابقه ارائه کردیم. مدل ما علاوه بر مبتنی بر سابقه بودن، تفاوت دیگری نیز با مدل‌های مشابه خود که در (Cohn, et al., 2009) (Cohn, et al., 2010) (O'Donnell, et al., 2009) مسئله استخراج گرامر جایگزینی درخت را به تعداد زیادی مسئله انتخاب مدل تقسیم می‌کند و هر مسئله انتخاب مدل را با روشی که بر مبنای یک مدل بیزی از فرایند شناختی انتخاب مدل، موجب ساده‌تر شدن هر مسئله انتخاب مدل می‌شود و در عوض تعداد این مسائل را در مسئله اصلی افزایش می‌دهد. ساده‌تر شدن مسائل انتخاب مدل موجب عملکرد بهتر فرایند دیریکله یا فرایند رستوران چینی در فرآگیری واحدهای پایه پر کاربرد می‌شود. نتایج حاصل از مدل ما در مقایسه با مدل‌های (Cohn, et al., 2010) و (Cohn, et al., 2009) در آزمایش‌های مشابه، این موضوع را در عمل نشان می‌دهد. علاوه‌بر این، مدل ما مدل مبتنی بر

(جدول ۵) نیز نتایج حاصل از مدل‌های پیشرو پیشین استخراج گرامر و تجزیه‌ی نحوی را بر روی بخش استاندارد پیکره Penn (بخش ۲۳) گزارش کرده است که می‌تواند با نتایج حاصل از مدل ارائه شده در این مقاله مورد مقایسه قرار گیرد.

(جدول ۵): نتایج حاصل از مدل‌های پیشین استخراج گرامر و

تجزیه نحوی روی بخش ۲۳ پیکره Penn.

F1	روش
88.2	کالینز (Collins, 1999)
89.5	چارنیاک (Charniak, 2000)
91.0	چارنیاک و جانسون (Charniak, et al., 2005)
90.7	باد (Bod, 2003)
90.0	پتروف و کلین (Petrov, et al., 2007)
92.1	مک‌کلاوسکی و همکاران (McClosky, et al., 2006)
91.82	مدل بیزی شناختی [این مقاله]

همان‌طور که مشاهده می‌شود نتیجه حاصل از مدل ارائه شده در این مقاله فقط از مدل ارائه شده توسط مک‌کلاوسکی و همکاران عملکرد ضعیفتری داشته است.

از آنجایی که فرایند یادگیری و تجزیه نحوی در مدل ارائه شده یک فرایند تصادفی است، لازم است اختلاف مقدار معیار F1 حاصل از مدل ارائه شده و مدل کهن و همکاران [۱۹] از دیدگاه ارزشمندی آماری مورد آزمون قرار گیرند تا بتوان نتیجه گرفت که اختلاف مشاهده شده معنادار است. برای محاسبه ارزشمندی آماری این اختلاف، ما از آزمون t استفاده کردیم. میزان احتمال درستی فرضیه خلف از طریق این آزمون برابر با  $1.62 * 10^{-12}$  شد. فرضیه خلف در این آزمون آن بود که مدل ما در تجزیه جملات، بدتر مساوی مدل کهن و همکاران عمل می‌کند. احتمال محاسبه شده برای فرضیه خلف، نشان‌دهنده نادرستی این فرضیه است. بدین ترتیب می‌توان نتیجه گرفت که اختلاف نتایج حاصل از آزمایش مدل ارائه شده و مدل کهن و همکاران (Cohn, et al., 2010) از دیدگاه آماری ارزشمند است. لازم به ذکر است که افزایش سطح به کارگیری اطلاعات سابقه می‌تواند موجب تولید گرامرهای دقیق‌تری شود (Feili, et al., 2006). اما همانطور که در (Feili, et al., 2006) اشاره شده است، این کار می‌تواند موجب افزایش نمایی در اندازه گرامر تولیدی شود و عملیات تجزیه نحوی را غیر ممکن کند. علاوه‌بر این، افزایش سطح اطلاعات سابقه می‌تواند باعث بروز مشکل



Cohn Trevor, Goldwater Sharon and Blunsom Phil, 2009. Inducing Compact but Accurate Tree-Substitution Grammars. NAACL.

Collins Michael, 1996. A new statistical parser based on bigram lexical dependencies. Proceedings of the 34th Meeting of the Association for Computational Linguistics - .Santa Cruz, CA. - pp. 184-191.

Collins Michael, 1999. Head-Driven Statistical Models for Natural Language Parsing. PhD thesis / University of Pennsylvania.

Feili Heshaam and Ghassem-Sani Gholamreza, 2006. Unsupervised Grammar Induction Using History Based Approach. Computer Speech and Language. - Vol. 20. - pp. 644-658.

Goldwater Sharon and Griffiths Thomas L., 2007, A Fully Bayesian Approach to Unsupervised Part-of-Speech Tagging. Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics. The Association for Computational Linguistics. - pp. 744-751.

Goldwater Sharon, Griffiths Thomas L. and Johnson Mark, 2009. A Bayesian framework for word segmentation: Exploring the effects of context. Cognition. - Vol. 112. - pp. 21-54.

Goodman Joshua, 2002. Efficient Parsing of DOP with PCFG-reductions. Data Oriented Parsing / ed. Bod Rens, Sima'an Khalil and Scha Remko. University of Chicago Press.

Griffiths Thomas, Kemp Charles and Tenenbaum Joshua B, 2008. Bayesian models of cognition. The Cambridge Handbook of Computational Psychology / book auth. Sun Ron. Cambridge University Press.

Harrison Philip [et al.], 1991. Evaluating Syntax Performance of Parser/Grammars of English. Natural Language Processing Systems Evaluation Workshop / ed. Neal Jeannette G.and Walter Sharon M. ACL - pp. 71-78.

Johnson Mark, 1998. The Effect of Alternative Tree Representations on Tree Bank Grammars. Proceedings of the Joint Conferences on New Methods in Language Processing and Computational Natural Language Learning - .Sydney, Australia. - pp. 39-48.

Johnson Mark, Goldwater Sharo, 2009. Improving nonparametric Bayesian inference: Experiments on unsupervised word segmentation with adaptor grammars. Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics. Boulder, Colorado.

Johnson Mark, Griffiths Thomas L. and Goldwater Sharon, 2007. Adaptor grammars: A framework for specifying compositional nonparametric Bayesian models. Advances in Neural Information Processing Systems. - Vol 19.

انتخاب مدل را نسبت به مدل‌های مبتنی بر درستنمایی بیشینه، در حل مسائل انتخاب، مدل نشان می‌دهد.

## ۷- تشرک و قدردانی

در پایان لازم می‌بینیم از پژوهشگاه دانش‌های بنیادی (IPM) برای فراهم کردن امکان اجرای آزمایش‌های این تحقیق بر روی رایانه خوشه‌ای آن پژوهشگاه (www.cluster.hpc.ipm.ac.ir) تشرک کنیم.

## ۸- مراجع

Bod Rens, 1992. A computational model of language performance: data oriented parsing. Proceedings COLING'92. - Nantes, France. - pp. 855-859.

Bod Rens, 2003. An efficient implementation of a new DOP model. EACL '03 Proceedings of the tenth conference on European chapter of the Association for Computational Linguistics. - Budapest, Hungary. - Vol. 1. - pp. 19-26.

Bod Rens, 1995. The problem of computing the most probable tree in data-oriented parsing and stochastic tree grammars. Proceeding EACL '95 Proceedings of the seventh conference on European chapter of the Association for Computational Linguistics - .San Francisco, CA, USA.

Brody Samuel, Lapata Mirella, 2009. Bayesian word sense induction. Proceedings of the 12th Conference of the European Chapter of the Association for Computational Linguistics - .Stroudsburg, PA, USA . -pp. 103-111.

Cahill Aoife, 2008. Treebank-Based Probabilistic Phrase Structure Parsing. Language and Linguistics Compass. - Vol. 2. - pp. 18-40.

Charniak Eugene, 2000. A maximum-entropy-inspired parser. Proceedings of the First Annual Meeting of the North American chapter of the Association for Computational Linguistics (NAACL 2000). - pp. 132-139.

Charniak Eugene, 1996. Treebank grammars. Proceedings of the Thirteenth National Conference on Artificial Intelligence. - Menlo Park. - pp. 1031-1036.

Charniak Eugene, Johnson Mark, 2005. Coarse-to-fine n-best parsing and MaxEnt discriminative reranking. Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics. - pp. 173-180.

Cohn Trevor, Blunsom Phil and Goldwater Sharon, 2010. Inducing Tree-Substitution Grammars. Journal of Machine Learning Research. - pp. 3053-3096.



دانشکده فنی دانشگاه تهران است. زمینه‌های تحقیقاتی مورد علاقه ایشان عبارتند از: "پردازش هوشمند متن و زبان طبیعی"، "متترجم ماشینی"، "داده کاوی"، "بازیابی اطلاعات" و "شبکه‌های اجتماعی".

ایشان به مدت یک دهه تجربه مدیریت یک تیم نرمافزاری در زمینه تولید سامانه مترجم ماشینی را بر عهده دارند. همچنین از ایشان حدود ۵۰ مقاله در نشریات معتبر و کنفرانس‌های بین‌المللی و ملی به چاپ رسیده است.

نشانی رایانامه ایشان عبارت است از:

hfaili@ut.ac.ir



**حمیدرضا قادر** تحصیلات خود را در مقاطع کارشناسی مهندسی کامپیوتر گرایش نرمافزار در دانشکده مهندسی برق و کامپیوتر دانشگاه تهران در سال ۱۳۸۸ به پایان رساند. سپس در سال

۱۳۹۰ در مقاطع کارشناسی ارشد در گرایش هوش مصنوعی از دانشکده مهندسی کامپیوتر دانشگاه علم و صنعت ایران فارغ التحصیل شد. وی هم اکنون عضو تیم تحقیقاتی توسعه متترجم ماشینی در آزمایشگاه پردازش متن و زبان طبیعی دانشگاه تهران است. زمینه‌های تحقیقاتی مورد علاقه وی شامل "پردازش متن و زبان طبیعی"، "متترجم ماشینی"، "یادگیری ماشین" و "ساخت اتوماتیک منابع زبانی" است.

نشانی رایانامه ایشان عبارت است از:

h.ghader@ece.ut.ac.ir



**مرتضی آنالویی** دانشیار دانشکده مهندسی کامپیوتر در دانشگاه علم و صنعت ایران است. تحقیقات دانشگاهی و فعالیت‌های صنعتی ایشان معطوف به فناوری‌های روز در

حوزه‌های "شبکه‌های اجتماعی و اقتصادی"، "هوش مصنوعی و علوم شناختی"، "بهینه‌سازی" و "مهندسی زبان‌های طبیعی" می‌باشد.

نشانی رایانامه ایشان عبارت است از:

analoui@iust.ac.ir

Johnson Mark, Griffiths Thomas L. and Goldwater Sharon, 2007 b. Bayesian Inference for PCFGs via Markov Chain Monte Carlo. Proceedings of the North American Conference on Computational Linguistics (NAACL '07). The Association for Computational Linguistics. - pp. 139-146.

Marcus Mitchell P., Santorini Beatrice and Marcinkiewicz Mary Ann, 1993. Building a large annotated corpus of English: the Penn treebank. Computational Linguistics. - Vol. 19. - pp. 313-330.

McClosky David, Charniak Eugene and Johnson Mark, 2006. Effective Self-Training for Parsing. The Conference on Human Language Technology and North American chapter of the Association for Computational Linguistics (HLT-NAACL 2006). Brooklyn, New York .

Mohri Mehryar, Sproat Richard, 2006. On a Common Fallacy in Computational Linguistics. SKY Journal of Linguistics / ed. Suominen Mickael . - Vol. 19. - pp. 432-439.

O'Donnell Timothy J., Goodman Noah D. and Tenenbaum Joshua B., 2009. Fragment Grammars: Exploring Computation and Reuse in Language. Massachusetts Institute of Technology - .Cambridge: MIT-CSAIL-TR-2009-013.

Petrov Slav and Klein Dan, 2007. Improved Inference for Unlexicalized Parsing. Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Proceedings of the Main Conference. - pp. 404-411.

Resnik Philip, Hardisty Eric, 2010. Gibbs Sampling for the Uninitiated. Technical Report. UMIACS.

Teh Yee Whye, 2010. Dirichlet Processes. Encyclopedia of Machine Learning. Springer.

Titov Ivan, Klementiev Alexandre, 2011. A Bayesian Model for Unsupervised Semantic Parsing. The 49th Annual Meeting of the Association for Computational Linguistics (ACL). Portland, USA.

**فیلی هشام**، ۱۳۸۵. استنتاج استقرائي یک گرامر محاسباتي برای زبان طبیعی با استفاده از روش بی مری. پایان نامه دکتری / دانشگاه صنعتی شریف. - تهران.

**هشام فیلی** تحصیلات خود را در مقاطع کارشناسی نرمافزار در دانشکده مهندسی کامپیوتر دانشگاه صنعتی شریف با رتبه یک در سال ۱۳۷۶ به پایان رساند. سپس مقاطع کارشناسی ارشد نرمافزار و دکتری هوش مصنوعی را به ترتیب در سال‌های ۱۳۷۸ و ۱۳۸۵ در همان دانشکده تکمیل کرد. از سال ۱۳۸۷ تا کنون عضو هیئت علمی دانشکده مهندسی برق و کامپیوتر

