

باسمه تعالی



فاز اول پروژه‌ی درس مقدمه‌ای بر یادگیری ماشین

## مدل‌های گرافی و ماشین بولترمن

استاد درس

دکتر محمدحسین یاسائی میبدی

دانشکده‌ی مهندسی برق  
دانشگاه صنعتی شریف

پاییز ۱۴۰۳

آخرین مهلت تحویل:  
۱۷ دی ۱۴۰۳

## فهرست مطالب

۲	۱	مدل‌های گرافی
۲	۱.۱	چند تعریف
۲	۱.۱.۱	گراف بدون جهت
۲	۲.۱.۱	گروه خاص (clique)
۲	۳.۱.۱	گروه خاص ماکسیمال (Maximal clique)
۲	۴.۱.۱	مسیر
۳	۲.۱	مدل گرافی
۳	۱.۲.۱	متغیرهای مستقل شرطی (conditionally independent variables)
۳	۲.۲.۱	فضای حالت
۳	۳.۲.۱	میدان تصادفی مارکف
۳	۴.۲.۱	پوشش مارکف
۳	۵.۲.۱	توزیع به صورت سخت مثبت
۴	۶.۲.۱	Hammersley-Clifford Theorem
۴	۳.۱	Unsupervised Learning
۴	۱.۳.۱	مسئله در میدان تصادفی مارکف
۵	۲.۳.۱	ML
۵	۳.۳.۱	Latent Variables
۵	۴.۳.۱	محاسبه گرادیان لگاریتم درست نمایی میدان تصادفی کارکف با متغیر پنهان
۶	۲	ماشین بولتزمن
۹	۳	تقریب گرادیان، الگوریتم کارا
۱۱	۴	بخش عملی
۱۱	۱.۴	بخش اول: پیش‌پردازش داده‌ها
۱۱	۲.۴	بخش دوم: پیاده‌سازی الگوریتم یادگیری
۱۱	۳.۴	بخش سوم: نمایش روند نمونه‌سازی
۱۱	۴.۴	بخش چهارم: کنترل روی نمونه‌های تولیدی
۱۲	۵	نکات مهم

## ۱ مدل‌های گرافی

مدل‌های گرافی ابزارهای قدرتمندی برای توصیف توزیع‌های احتمالی هستند. این مدل‌ها ارتباط بین متغیرها را به خوبی بیان میکنند و به طور خاص برای بررسی این ارتباطات مناسب هستند. زیرا این امکان را میدهند تا با ابزارهای نظریه گراف این همبستگی‌ها و ارتباطات را تحلیل کنیم. همچنین میتوان محاسبات پیچیده را با الگوریتم‌های گراف با محاسبات کمتری انجام دهیم. انواع متفاوتی از این مدل‌ها وجود دارند برای آشنایی با این مدل‌ها ابتدا با چند تعریف در گراف آشنا میشویم.

### ۱.۱ چند تعریف

#### ۱.۱.۱ گراف بدون جهت

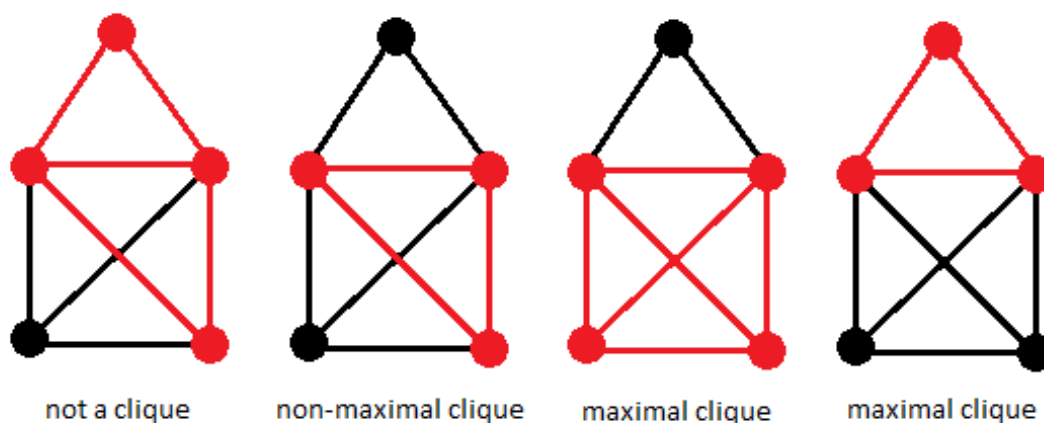
گراف بدون جهت به تاپل دوتایی  $G = (V, E)$  می‌گوییم که در  $V$  نشان دهنده گره‌ها و  $E$  نشان دهنده یال‌هاست به گونه‌ای که اگر یالی بین دو نود  $v, w$  داشته باشیم داریم  $v, w \in E$  و می‌گوییم  $w$  در همسایگی  $v$  است و بالعکس. برخلاف گراف بدون جهت در گراف با جهت ترتیب  $v, w$  اهمیت دارد چون هر یال جهت نیز دارد اما در گراف بدون جهت، جهت نداریم.

#### ۲.۱.۱ گروه خاص (clique)

به زیر مجموعه‌ای از گره‌های یک گراف که که یک گراف کامل را تشکیل میدهند گروه خاص می‌گوییم.

#### ۳.۱.۱ گروه خاص ماکسیمال (Maximal clique)

گروه خاص ماکسیمال به گروه خاصی می‌گوییم که هیچ گره‌ی دیگری در گراف نباشد که با اضافه کردنش گروه همچنان خاص بماند. در شکل ۱ مثال‌هایی را مشاهده می‌کنید.



شکل ۱: Maximal clique

مجموعه تمام گروه‌های خاص ماکسیمال یک گراف را با  $C$  نشان می‌دهیم.

#### ۴.۱.۱ مسیر

به دنباله‌ای از یال‌های  $v_1, v_2, v_3, \dots, v_m \in V$  که برای  $i = 1, 2, 3, \dots, m-1$  داشته باشیم  $\{v_i, v_{i+1}\} \in E$  یک مسیر از  $v_1$  به  $v_m$  می‌گوییم. به مجموعه گره‌های  $V \subset \mathcal{V}$  جداکننده دو گره  $v \notin \mathcal{V}$

$v, w \notin \mathcal{V}$  می‌گوییم اگر هر مسیر از  $v$  به  $w$  شامل گره‌ای از  $\mathcal{V}$  باشد.

## ۲.۱ مدل گرافی

مدل های گرافی توزیع های احتمالاتی و روابط وابستگی بین آنها را با یک ساختار گراف بیان میکنند، به گونه ای که هر گره گراف نماینده یک متغیر تصادفی است و یال ها بیانگر وابستگی بین متغیرها هستند. انواع متفاوتی از مدل های گرافی وجود دارند که در این جا با میدان تصادفی مارکف آشنا میشوید.

### ۱.۲.۱ متغیرهای مستقل شرطی (conditionally independent variables)

می‌گوییم دو متغیر  $X_1, X_2$  به شرط  $X_3$  مستقل هستند هرگاه داشته باشیم.

$$p(X_1, X_2 | X_3) = p(X_1 | X_3) \cdot p(X_2 | X_3)$$

### ۲.۲.۱ فضای حالت

فضای حالت در یک گراف مارکوف به مجموعه‌ای از تمام حالات ممکن که یک سیستم می‌تواند در آنها باشد اشاره دارد. در یک گراف، متغیرهای تصادفی به هر گره اختصاص داده می‌شوند و فضای حالت شامل تمامی ترکیب های ممکن است که برای این متغیرهای تصادفی قابل تصور است. توجه شود فضای حالت بسیار به فضای نمونه ای نزدیک است فضای حالت یک مفهوم کلی تر از فضای نمونه ای است و بیانگر تمام مقادیری است که برای یک متغیر و یا حالت سیستم قابل تصور است در فضای نمونه ای علاوه بر حالت ها یک تابع باید تعریف شود تا به هر حالت ممکن یک احتمال نسبت دهد.

### ۳.۲.۱ میدان تصادفی مارکف

فرض کنید به هر گره از گراف  $G = (V, E)$  یک متغیر تصادفی نسبت دهیم که مقادیری از فضای حالت  $\Lambda_v$  را اختیار میکند. برای سادگی فرض میکنیم تمام متغیرها فضای حالت معادل  $\Lambda$  دارند. متغیرهای تصادفی  $\mathbf{X} = (X_v)_{v \in V}$  را میدان احتمالاتی مارکف مینامیم اگر توزیع احتمالاتی مشترک  $p$  خاصیت جهانی مارکف را نسبت به گراف برآورده کند. این یعنی برای هر زیر مجموعه  $\mathcal{A}, \mathcal{B}, \mathcal{S} \subset V$  که در آن تمام گره های  $\mathcal{A}$  و  $\mathcal{B}$  توسط  $\mathcal{S}$  از هم جدا شده اند داشته باشیم:  $\mathbf{X} = (X_a)_{a \in \mathcal{A}}$  و  $\mathbf{X} = (X_b)_{b \in \mathcal{B}}$  با داشتن  $\mathbf{X} = (X_s)_{s \in \mathcal{S}}$  به صورت شرطی مستقل اند. یا به عبارتی داریم:

$$p((x_a)_{a \in \mathcal{A}} | (x_t)_{t \in \mathcal{S} \cup \mathcal{B}}) = p((x_a)_{a \in \mathcal{A}} | (x_t)_{t \in \mathcal{S}})$$

به تعبیری دیگر در یک میدان مارکف توزیع گره  $v$  به شرط همسایه هایش مستقل از بقیه گره هاست.

### ۴.۲.۱ پوشش مارکف

مجموعه‌ای از گره‌ها  $\text{MB}(v)$  را پوشش مارکوف گره  $v$  می‌نامند، اگر برای هر مجموعه گره‌های  $B$  با  $v \notin B$  داشته باشیم

$$p(v | \text{MB}(v), B) = p(v | \text{MB}(v)).$$

در یک میدان مارکف همسایه های گره  $v$  پوشش مارکف آن هستند.

### ۵.۲.۱ توزیع به صورت سخت مثبت

یک توزیع احتمال  $p$  به صورت سخت مثبت است اگر و تنها اگر برای هر  $x$  در فضای حالت  $\Lambda$  داشته باشیم:

$$p(x) > 0$$

به عبارت دیگر، این بدان معناست که هیچ مقدار از  $x$  وجود ندارد که احتمال آن صفر باشد و تمامی حالت‌های ممکن دارای احتمال غیرصفر هستند.

یک توزیع نسبت به یک گراف بدون جهت  $G$  با گروه های خاص ماکسیمال  $C$  فاکتوریزه می شود اگر مجموعه ای از توابع غیرمنفی  $\{\psi_C\}_{C \in C}$  که توابع پتانسیل نامیده می شوند وجود داشته باشند به طوری که

$$\forall \mathbf{x}, \hat{\mathbf{x}} \in \Lambda^{|V|} : (\mathbf{x}_c)_{c \in C} = (\hat{\mathbf{x}}_c)_{c \in C} \Rightarrow \psi_C(\mathbf{x}) = \psi_C(\hat{\mathbf{x}})$$

و

$$p(\mathbf{x}) = \frac{1}{Z} \prod_{c \in C} \psi_c(\mathbf{x}) = \frac{1}{Z} \prod_{c \in C} \psi_c(\mathbf{x}_c)$$

در اینجا:  $\mathbf{x}_c$  گره هایی از  $\mathbf{x}$  هستند که جزو گروه خاص  $c$  هستند -  $\psi_C(\mathbf{x})$  توابع پتانسیل هستند. -  $Z$  ثابت نرمال سازی است. -  $\Lambda^{|V|}$  فضای حالت است. و  $p(\mathbf{x})$  توزیع مشترک تمام متغیر های تصادفی مربوط به گره های  $G$  هستند. به طور کلی، فاکتوریزه کردن یک توزیع احتمال به معنای تجزیه آن به ضرب تعدادی توابع کوچکتر است که به بخش های مختلف یک گراف مربوط می شوند. این فرآیند باعث ساده تر شدن محاسبات و فهم بهتر ساختار توزیع می شود.

ثابت نرمال سازی  $Z$  به صورت

$$Z = \sum_x \prod_{c \in C} \psi_c(x_c)$$

است و تابع پارتیشن نامیده می شود. اگر  $p$  به صورت سخت مثبت باشد، این موضوع برای توابع پتانسیل نیز صادق است. بنابراین می توانیم بنویسیم:

$$p(x) = \frac{1}{Z} \prod_{c \in C} \psi_c(x_c) = \frac{1}{Z} \exp \left( \sum_{c \in C} \ln \psi_c(x_c) \right) = \frac{1}{Z} \exp(-E(x))$$

که در اینجا  $E := \sum_{c \in C} \ln \psi_c(x_c)$  به عنوان تابع انرژی نامیده می شود. بنابراین، توزیع احتمال هر میدان تصادفی مارکوف (MRF) می تواند به شکل ارائه شده در بالا بیان شود که به آن توزیع گیبس نیز می گویند.

پرسش تئوری ۱. سعی کنید درباره مفهوم توابع پتانسیل در یک میدان مارکف کمی توضیح دهید و یک مثال از یک میدان مارکف به همراه توابع پتانسیل اش بزنید. توجه کنید مثال میتواند بسیار ساده باشد اما حتما باید مقادیر عددی تعیین گردد و توضیحات هم تا حد ممکن کوتاه و هدفمند باشند.

## ۳.۱ Unsupervised Learning

به یادگیری (ویژگی های مهمی) از یک توزیع نامعلوم بر اساس سمپل های آن Unsupervised learning میگوئیم. یادگیری هرگونه ارتباط در بین داده ها به طور مثال کاهش بعد، دسته بندی داده های نزدیک به هم و کشف الگوها در دیتا نمونه هایی از Unsupervised learning هستند.

### ۱.۳.۱ مسئله در میدان تصادفی مارکف

فرض کنید تعدادی متغیر تصادفی دارید که فضای حالت یکسانی دارند و یک میدان تصادفی مارکف تشکیل می دهند. چند نمونه از حالت های سیستم دارید. یک نمونه می تواند به صورت زیر باشد:

$$x_{v_1}, x_{v_2}, \dots, x_{v_n}$$

همچنین فرض کنید ساختار کلی گراف را هم داریم. حال اگر فرض کنیم که توابع انرژی گراف می توانند از خانواده خاصی باشند که این خانواده از توابع با پارامتر  $\Theta$  پارامتریزه شده اند. می توان با استفاده از نمونه ها، پارامترها را یاد گرفت.

---

پرسش تئوری ۳.۲ نوع از توابع پتانسیل را نام برده و پارامترهای آنها را مشخص کنید.

---

### ۲.۳.۱ ML

برای پیدا کردن پارامترها از قاعده درست نمایی بیشینه استفاده میکنیم. با فرض داده های مستقل از هم به صورت زیر:

$$S = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$$

$$\mathbf{x}_1 = \{x_{v_1}, \dots, x_{v_n}\}$$

داریم:

$$L(\theta | S) = \prod_{i=1}^m p(\mathbf{x}_i | \theta)$$

عبارت بالا احتمال مشاهده دیتاست به شرط پارامترهاست. توجه کنید که تنها  $\mathbf{x}_i$  ها از هم مستقل اند و در یک سمپل مقادیر گره های مختلف باهم همبستگی دارند و هدف ما یافتن همین همبستگی ها میباشد.

---

پرسش تئوری ۳.۳ به علت اینکه پیدا کردن آنالیزی پارامترها برای بیشینه سازی عبارت بالا به طور عمومی ممکن نیست در خیلی از مواقع از روش *Gradient Ascent* استفاده میشود. در کمتر از سه خط این روش را توضیح بدهید.

---

### ۳.۳.۱ Latent Variables

پیشنهاد میشود قبل از شروع این بخش این ویدیو را ببینید. همانطور که مشاهده کردید برای مدل کردن یک توزیع مجهول با  $m$  متغیر میتوان گرافی با تعداد بیشتری گره در نظر گرفت که  $m$  گره در آنها مشاهده پذیر اند. اگر  $\mathbf{X} = (X_v)_{v \in V}$  متغیرهای مربوط به یک میدان تصادفی مارکف باشند که یک زیر مجموعه  $m$  تایی از آنها یعنی  $\mathbf{V} = (V_1, \dots, V_m)$  متغیرهای مشاهده پذیر اند، به متغیرهای باقی مانده  $\mathbf{H} = (H_1, \dots, H_m)$  متغیرهای Latent یا پنهان میگوییم. در این شرایط توزیع مشترک میدان مجهول مارکف به صورت زیر است:

$$p(v) = \sum_h p(v, h) = \frac{1}{Z} \sum_h e^{-E(v, h)}$$
$$Z = \sum_{v, h} e^{-E(v, h)}$$

همانطور که در ویدیو مشاهده کردید علت استفاده از متغیرهای Latent مدل کردن و کشف ارتباطات پیچیده ترین متغیرهای مشاهده پذیر است.

### ۴.۳.۱ محاسبه گرادیان لگاریتم درست نمایی میدان تصادفی کارکف با متغیر پنهان

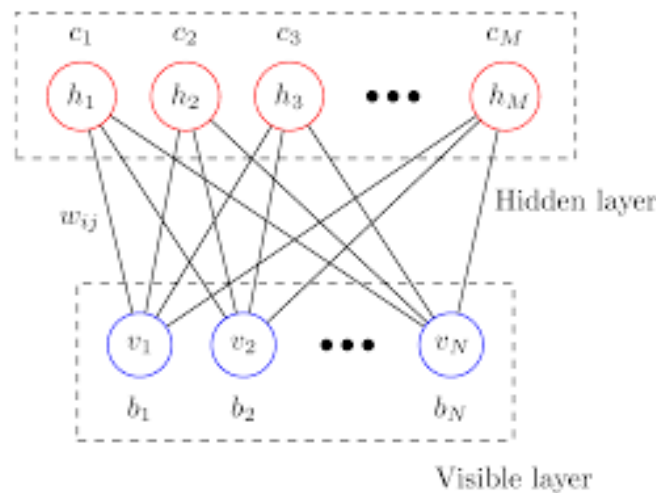
لگاریتم درست نمایی برای میدان تصادفی مارکف با متغیر پنهان به صورت زیر است:

$$\ln L(\theta | v) = \ln p(v | \theta) = \ln \frac{1}{Z} \sum_h e^{-E(v, h)} = \ln \sum_h e^{-E(v, h)} - \ln \sum_{v, h} e^{-E(v, h)}$$

$$\frac{\partial \ln L(\theta | v)}{\partial \theta} = - \sum_h p(h | v) \frac{\partial E(v, h)}{\partial \theta} + \sum_{v, h} p(v, h) \frac{\partial E(v, h)}{\partial \theta}$$

## ۲ ماشین بولتزمن

در این پروژه قصد داریم که با ساختار ماشین بولتزمن Boltzmann Machines آشنا بشویم. ماشین های بولتزمن سعی بر این امر دارد که بتواند توزیع نمونه های ورودی را تخمین بزند و اینگونه می تواند در بخش های مختلف یادگیری ماشین استفاده بشود. ابتدا با ساختار آن آشنا می شویم تا درک بهتری از نوع عملکرد آن داشته باشیم. ماشین بولتزمن شامل دو لایه نورون هست که شامل یک لایه نمایان بوده و یک لایه پنهان. (مطابق شکل زیر)



همانطور که می بینید تمام نورون های لایه نمایان به تمام نورون های لایه پنهان متصل هستند و بالعکس و هیچ دو نورون در یک لایه به یکدیگر متصل نیستند که این باعث مستقل شدن نورون های پنهان به شرط یک نورون نمایان می شود. این لایه پنهان سعی بر این دارد که توزیع ورودی را تخمین بزند. بدین معنی که با استخراج ویژگی های مهم از ورودی می تواند توزیع ورودی را بسازد.

پرسش تئوری ۰۵. با توجه به توضیحات داده شده در مورد نحوه عملکرد ماشین بولتزمن، این ساختار در چه بخش هایی از یادگیری ماشین می تواند استفاده شود؟

حال بایستی یک توزیع مشترک بین تمام حالاتی که نورون های نمایان و پنهان نسبت بدهیم. (توجه کنید که مقدار تمام نورون ها ۰ یا ۱ می باشد) به همین علت از توزیع گیبس Gibbs distribution استفاده می کنیم که فرمول آن به صورت زیر می باشد.

$$p(v, h) = \frac{1}{Z} e^{-E(v, h)}$$

$v$  متناظر با نورون های نمایان و  $h$  متناظر با نورون های پنهان می باشد.

$$E(v, h) = - \sum_{i=1}^n \sum_{j=1}^m w_{ij} h_i v_j - \sum_{j=1}^m b_j v_j - \sum_{i=1}^n c_i h_i.$$

که  $w_{ij}$  وزن تخصیص داده شده به یال بین نورون  $v_j$  و  $h_i$  می باشد و  $b_j$  و  $c_i$  بایاس های متناظر با نورون نمایان  $i$  و نورون پنهان  $j$  می باشد.

همانطور که می دانید برای آموزش دادن مدل بایستی بتوانیم یک سری گرادیان و عبارات را حساب کنیم که در این بخش به آن می پردازیم.

پرسش تئوری ۰۶ ثابت کنید:

$$p(H_i = 1 | \mathbf{v}) = \sigma \left( \sum_{j=1}^m w_{ij} v_j + c_i \right)$$

$$p(V_j = 1 | \mathbf{h}) = \sigma \left( \sum_{i=1}^n w_{ij} h_i + b_j \right).$$

در این قسمت می خواهیم بحث کنیم که توزیع هایی که در پرسش قبل بدست آوردید به چه درد می خورد. در بخش بعدی به طور مفصل نحوه به روزرسانی پارامترها را فراخواهید گرفت اما در این قسمت مفاهیم اولیه آن را با یکدیگر مرور خواهیم کرد.

در بخش بعد خواهیم دید که ما نیاز داریم با وجود دانستن  $h$  مقدار نورو  $v$  و بالعکس را به صورت تصادفی با توزیعی که بدست آوردید مقدار دهی کنیم. طبیعتاً بایستی به گونه ای نمونه برداریم که توزیع نمونه های به وجود آمده نزدیک به توزیع بدست آمده باشد و برای این مهم از مفاهیم پایه فرآیند های مارکوف و نمونه برداری استفاده می کنیم که شرح زیر می باشد.

## فرایندهای مارکوف و نمونه برداری گیبس

### فرایندهای مارکوف

فرایند مارکوف یک فرآیند تصادفی است که در آن حالت فعلی سیستم تنها به حالت قبلی وابسته است و مستقل از توالی حالت های گذشته است. به طور دقیق، اگر  $\{X(k) | k \in \mathbb{N}_0\}$  یک زنجیره مارکوف باشد، آنگاه برای همه  $k \geq 0$  و حالت های  $i, j$  داریم:

$$P(X(k+1) = j | X(k) = i, X(k-1), \dots, X(0)) = P(X(k+1) = j | X(k) = i).$$

این ویژگی که به عنوان خاصیت مارکوف شناخته می شود، به این معناست که آینده تنها به وضعیت فعلی بستگی دارد و از تاریخچه گذشته مستقل است.

اگر احتمال انتقال از یک حالت به حالت دیگر ثابت باشد، زنجیره مارکوف همگن نامیده می شود. در این حالت، ماتریس انتقال  $P$  تعریف می شود که در آن عنصر  $p_{ij}$  احتمال انتقال از حالت  $i$  به حالت  $j$  است. زنجیره های مارکوف ویژگی های مهمی دارند که برای الگوریتم های نمونه برداری بسیار مفید هستند:

- ایستایی (Stationarity): یک توزیع  $\pi$  به عنوان توزیع پایا تعریف می شود اگر:

$$\pi^T = \pi^T P,$$

به این معنا که اعمال ماتریس انتقال بر روی  $\pi$  تغییری در آن ایجاد نمی کند. اگر زنجیره مارکوف برای زمان کافی اجرا شود، به این توزیع پایا همگرا خواهد شد.

- همگرایی: اگر زنجیره مارکوف غیرقابل کاهش (یعنی از هر حالت بتوان به هر حالت دیگری رسید) و غیرتناوبی باشد، توزیع حالت های زنجیره در نهایت به توزیع پایا نزدیک می شود، مستقل از توزیع اولیه.

### نمونه برداری گیبس

نمونه برداری گیبس یکی از روش های نمونه برداری مونت کارلو مبتنی بر زنجیره مارکوف (MCMC) است که برای تولید نمونه هایی از توزیع های احتمال پیچیده مانند توزیع گیبس در ماشین بولتزمن استفاده می شود. ایده اصلی این روش به شرح زیر است:

- فرض کنید توزیع مشترک متغیرهای تصادفی  $\{X_1, X_2, \dots, X_N\}$  به صورت  $\pi(x) = \frac{1}{Z} e^{-E(x)}$  تعریف شده است، که در آن انرژی سیستم و  $Z$  یک مقدار نرمال سازی است.

- به جای نمونه برداری مستقیم از  $\pi(x)$ ، که ممکن است پیچیده باشد، هر متغیر  $X_i$  به نوبت و بر اساس توزیع شرطی  $\pi(X_i | X_{-i})$  به روزرسانی می شود، که در آن  $X_{-i}$  مجموعه سایر متغیرها است.



در هر گام از الگوریتم:

- یک متغیر  $X_i$  انتخاب می‌شود.
- مقدار جدید  $X_i$  از توزیع شرطی  $P(X_i | X_{-i})$  نمونه‌برداری می‌شود، که به طور مؤثر احتمال تغییر این متغیر را به حالت‌های مختلف در نظر می‌گیرد.

چرا از نمونه‌برداری گیبس استفاده کنیم؟

در مدل‌های پیچیده‌ای مانند ماشین بولترمن محدود (RBM) توزیع احتمال نهایی سیستم (توزیع گیبس) معمولاً بسیار پیچیده است و نمی‌توان مستقیماً از آن نمونه‌برداری کرد. برای مثال، توزیع گیبس به صورت زیر تعریف می‌شود:

$$P(v, h) = \frac{1}{Z} e^{-E(v, h)},$$

که در آن  $E(v, h)$  تابع انرژی سیستم است. برای آموزش مدل و یا ارزیابی آن، نیاز داریم نمونه‌هایی از این توزیع به دست آوریم.

- ساخت یک زنجیره مارکوف: ابتدا زنجیره‌ای طراحی می‌کنیم که توزیع پایای آن همان توزیع گیبس باشد.
- تضمین همگرایی: اگر زنجیره مارکوف غیرقابل کاهش و غیرتناوبی باشد، تضمین می‌شود که با گذشت زمان به توزیع گیبس همگرا شود.
- بازنمایی داده‌ها: نمونه‌هایی که از زنجیره مارکوف پس از همگرایی به دست می‌آیند، به طور مؤثری نماینده توزیع گیبس هستند و می‌توانند برای آموزش یا تحلیل مدل استفاده شوند.

به بیان ساده، نمونه‌برداری گیبس به ما این امکان را می‌دهد که از پیچیدگی توزیع اصلی عبور کرده و نمونه‌هایی تولید کنیم که ساختار مدل را به خوبی بازنمایی می‌کنند. این روش نه تنها محاسبات را ساده می‌کند، بلکه امکان رسیدن به حالت پایدار را نیز فراهم می‌آورد.

---

پرسش تئوری ۷. برای به روزرسانی پارامترها، نیاز داریم مشتق‌های جزئی تابع هدف را نسبت به پارامترها حساب کنیم. مشتق‌های جزئی زیر را بدست آورده و تا حد امکان آنها را ساده کنید.

$$\frac{\partial \ln \mathcal{L}(\theta | v)}{\partial w_{ij}}$$

$$\frac{\partial \ln \mathcal{L}(\theta | v)}{\partial b_j}$$

$$\frac{\partial \ln \mathcal{L}(\theta | v)}{\partial c_i}$$

---

پرسش تئوری ۸. نشان دهید پاسخ شما به صورت زیر می‌تواند نوشته شود:

$$\frac{\partial \ln \mathcal{L}(\theta | v)}{\partial w_{ij}} = \frac{1}{l} \sum_{v \in \mathcal{S}} [\mathbb{E}_{p(h|v)}[v_i h_j] - \mathbb{E}_{p(h,v)}[v_i h_j]] = \langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{model}} \quad (1)$$

### ۳. تقریب گرادیان، حل مشکل

در قسمت های قبل گرادیان را نسبت به پارامترهای مدل محاسبه کردید، به طور مثال برای وزن  $w_{ij}$  در ماشین بولتزمن، خواهیم داشت.

$$\frac{\partial \ln \mathcal{L}(\theta|v)}{\partial w_{ij}} = \frac{1}{l} \sum_{v \in \mathcal{S}} [\mathbb{E}_{p(h|v)}[v_i h_j] - \mathbb{E}_{p(h,v)}[v_i h_j]] = \langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{model}} \quad (2)$$

با تمامی زیبایی ریاضیاتی خود، معادله ۲ غیر قابل محاسبه است جمله دوم نیاز به محاسبه  $\mathbb{E}_{\text{model}}$  دارد، و جمع روی تعداد حالت های  $v, h$  نیاز به محاسبه  $2^m, 2^n$  حالت مختلف است. یک ایده محاسباتی ساده، این است که همانطور که جمله اول با استفاده از Monte-Carlo estimation از سمپل های داده محاسبه میشود، جمله دوم را نیز همینطور تخمین بزنیم.

$$\mathcal{S} = \{\tilde{v}_s\}_{s=1}^S \rightarrow \mathbb{E}_{p(h,v)}[v_i h_j] = \mathbb{E}_{p(v)} \mathbb{E}_{p(h|v)}[v_i h_j] \approx \frac{1}{S} \sum_{s=1}^S \mathbb{E}_{p(h|v)}[\tilde{v}_i h_j]$$

برای تخمین زدن جمله دوم با استفاده از روش های مونته-کارلو، نیاز به سمپل های زیادی از مدل  $p(v, h)$  دارند، که همانطور که صحبت شد خود نیاز به الگوریتم MCMC دارد. در نتیجه برای دقیق شدن تقریب خود، نیاز به تعداد زیادی نمونه تولید شده توسط یک الگوریتم MCMC هستیم، که خود به اندازه کافی زنجیر طولانی ای بوده که توزیع  $p(v)$  را تخمین بزند باشد. در نتیجه این الگوریتم، خیلی در عمل مناسب نیست. ولی الگوریتم هایی وجود دارند که به این فرایند سرعت بسیاری میدهند، و تنها به  $k$  مرحله در الگوریتم MCMC (۱) در بسیاری از کاربردها نیاز خواهیم داشت. و اینگونه contrastive-divergence معرفی میشود. حال جمله دوم رو تنها با استفاده از یک نمونه تقریبی از مدل با استفاده از  $k$  مرحله MCMC با مقدار اولیه  $v$  انجام میشود، به صورتی که  $v$  یکی از نمونه های داده باشد.

$$\text{CD}_k(\theta, v_*) = - \sum_h p(h|v_*) \frac{\partial E(v_*, h)}{\partial \theta} + \sum_h p(h|v_k) \frac{\partial E(v_k, h)}{\partial \theta} \quad (3)$$

که در اینجا  $v_k$  نمونه تولید شده بعد از  $k$  مرحله MCMC بر روی نمونه اولیه داده  $v_*$  است. توزیع نمونه ها پس از  $k$  مرحله، را با  $p_{\theta}^k$  مشخص میکنیم.

$$v_k \sim p_{\theta}^k \iff v_k \sim \text{MCMC}(p_{\theta}, k, \text{initial}: v_*)$$

مشخص است که  $p_{\theta}^0 = p_{\text{data}}, v_* \sim p_{\text{data}}$  و در حالت نهایی،  $p_{\theta}^{\infty} = p_{\theta}, v^{\infty} \sim p_{\theta}$  پایداری الگوریتم فوق به صورت عملی، و تئوری اثبات شده است.

**پرسش تئوری ۹.** نشان دهید که بهینه کردن  $\text{CD}_k(\theta, v_*)$  معادل با بهینه سازی  $D_{\text{KL}}(p_{\text{data}} \| p_{\theta}) - D_{\text{KL}}(p_{\theta}^k \| p_{\theta})$  است، که در آن:

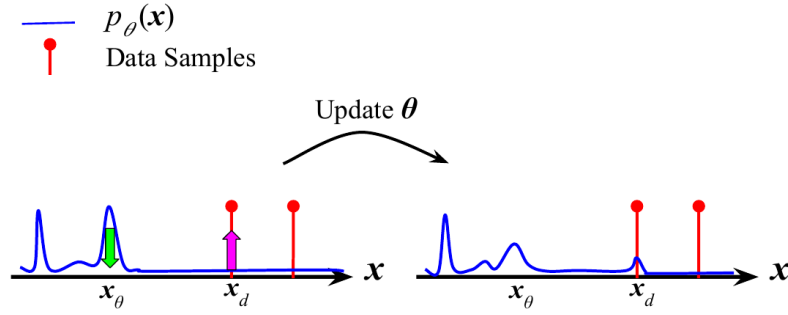
$$D_{\text{KL}}(P \| Q) = \sum_x P(x) \log \left( \frac{P(x)}{Q(x)} \right) = \mathbb{E}_{x \sim P} \left[ \log \frac{P(x)}{Q(x)} \right]$$

نشان دهید که با افزایش  $k$  بهینه کردن CD معادل با تخمینگر بیشینه درست نمایی است.

با بهینه کردن CD دو اتفاق همزمان رخ میدهد.

$$\theta \longleftarrow \theta - \nabla_{\theta} \text{CD} \Rightarrow \begin{cases} p_{\theta}(x) \uparrow & x \sim p_{\text{data}} \\ p_{\theta}^k(x) \downarrow & x \sim p_{\theta} \end{cases}$$

به صورت خلاصه، حال یک الگوریتم کلی برای یادگیری RBM به صورت شهودی، فرایند به شکل زیر است.



شکل ۲: الگوریتم شهودی معمولی RBM

به صورت الگوریتمی، به شکل زیر است.

---

**Algorithm 1.**  $k$ -step contrastive divergence

---

**Input:** RBM  $(V_1, \dots, V_m, H_1, \dots, H_n)$ , training batch  $S$

**Output:** gradient approximation  $\Delta w_{ij}$ ,  $\Delta b_j$  and  $\Delta c_i$  for  $i = 1, \dots, n$ ,  
 $j = 1, \dots, m$

```

1  init  $\Delta w_{ij} = \Delta b_j = \Delta c_i = 0$  for  $i = 1, \dots, n, j = 1, \dots, m$ 
2  forall the  $v \in S$  do
3       $v^{(0)} \leftarrow v$ 
4      for  $t = 0, \dots, k - 1$  do
5          for  $i = 1, \dots, n$  do sample  $h_i^{(t)} \sim p(h_i | v^{(t)})$ 
6          for  $j = 1, \dots, m$  do sample  $v_j^{(t+1)} \sim p(v_j | h^{(t)})$ 
7      for  $i = 1, \dots, n, j = 1, \dots, m$  do
8           $\Delta w_{ij} \leftarrow \Delta w_{ij} + p(H_i = 1 | v^{(0)}) \cdot v_j^{(0)} - p(H_i = 1 | v^{(k)}) \cdot v_j^{(k)}$ 
9           $\Delta b_j \leftarrow \Delta b_j + v_j^{(0)} - v_j^{(k)}$ 
10          $\Delta c_i \leftarrow \Delta c_i + p(H_i = 1 | v^{(0)}) - p(H_i = 1 | v^{(k)})$ 

```

---

شکل ۳: الگوریتم بهینه سازی اولیه برای RBM

---

پرسش تئوری ۱۰. یافته های قسمت های قبل خود را، با الگوریتم  $k$ -step contrastive divergence توجیح کنید، در مورد پیچیدگی این الگوریتم، و انتظار شما از خروجی مدل بعد از اعمال یادگیری با استفاده از این الگوریتم بحث کنید.

---



---

پرسش تئوری ۱۱. همانطور که تا الان مشاهده کردید، ماشین بولتزمن مناسب کار کردن با داده های باینری است، و  $h$  و  $v$  هر دو متغیرهای باینری هستند، آیا ممکن است این مدل ها را برای داده های پیوسته استفاده کرد؟ پیشنهاد شما چیست، برای  $v, h \in \mathbb{R}$  و همچنین  $v, h \in [0, 1]$  آیا با الگوریتم های مشابه با الگوریتم قبل میتوان این مدل ها را آموزش داد؟

---

## ۴ بخش عملی

### ۱.۴ بخش اول: پیش‌پردازش داده‌ها

- داده‌های MNIST را بارگذاری کرده و آن‌ها را به حالت باینری تبدیل کنید (مثلاً با اعمال یک آستانه روی شدت پیکسل‌ها).
- نمونه‌هایی از داده‌های باینری‌شده را نمایش دهید تا با داده‌ها آشنا شوید.

### ۲.۴ بخش دوم: پیاده‌سازی الگوریتم و یادگیری

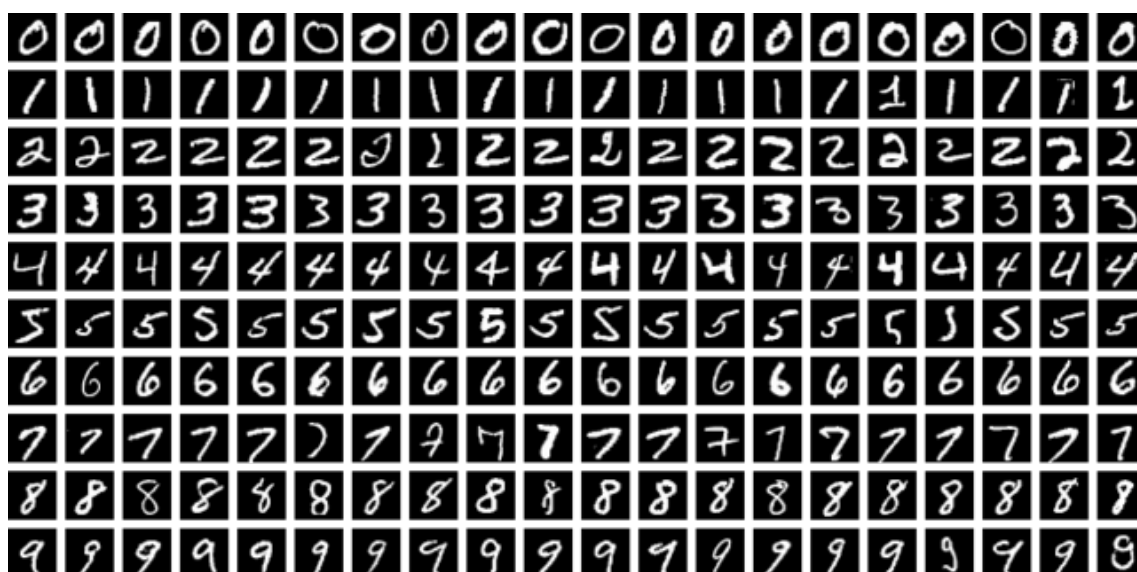
- الگوریتم  $k$ -step contrastive divergence را برای یادگیری مدل خود پیاده‌سازی کنید.
- به ازای مقادیر مختلف  $k$  (مثلاً  $k = 1, 5, 10$ ) عملیات یادگیری را انجام دهید.
- نمونه‌های تولیدی مدل را به ازای  $k$ های مختلف تولید کرده و نمایش دهید.

### ۳.۴ بخش سوم: نمایش روند نمونه‌سازی

پویانمایی‌ای از نمونه‌های تولیدی در طول اجرای الگوریتم MCMC ( $v_0, v_1, \dots, v_k$ ) بسازید و روند آن را نمایش دهید. این پویانمایی باید نشان دهد که چگونه نمونه‌ها در طول تکرارها بهبود می‌یابند.

### ۴.۴ بخش چهارم: کنترل روی نمونه‌های تولیدی

- بررسی کنید که آیا می‌توان از مدل برای تولید نمونه‌هایی از اعداد مشخص (مثلاً ۰ یا ۱) استفاده کرد یا خیر.
- پیشنهاد خود را برای ایجاد کنترل بر خروجی مدل توصیف کنید.
- اگر امکان دارد، راهکار خود را پیاده‌سازی کرده و نمونه‌های کنترل‌شده تولیدی را نمایش دهید.



شکل ۴: نمونه‌هایی از داده‌های MNIST

## ۵ نکات مهم

لطفاً به نکات زیر دقت کنید:

۱. پروژه شامل دو فاز خواهد بود.
۲. پروژه را میتوانید به صورت انفرادی یا به شکل گروه های دو نفره انجام دهید. دقت کنید چه به شکل انفرادی و چه به صورت گروهی باید تمام بخش های پروژه را انجام دهید و انجام انفرادی آن امتیاز اضافه ای برای شما نخواهد داشت.
۳. دو فاز این پروژه در مجموع ۲ تا ۳ نمره از نمره درس را تشکیل می دهند.
۴. پس از پایان پروژه یک روز برای تحویل حضوری پروژه در نظر گرفته می شود و باید کد ها و خروجی های خود را در حضور دستیاران آموزشی ارائه دهید و به پرسش های دستیاران پاسخ دهید. دقت کنید که تمام اعضای گروه باید به تمام بخش های پروژه مسلط باشند. در نهایت برای تمام اعضای گروه یک نمره در نظر گرفته خواهد شد.
۵. برای فاز نخست پروژه میتوانید حداکثر ۳ روز تاخیر مجاز استفاده نمایید اما به دلیل وجود ددلاین ثبت نمرات ددلاین فاز دوم پروژه سخت خواهد بود و امکان استفاده از تاخیر برای آن وجود ندارد.
۶. تمام بخش های هر دو فاز پروژه اجباری هستند و نمره امتیازی برای آن در نظر گرفته نشده است.
۷. تمامی شبیه سازی ها باید با کمک زبان Python انجام شود. همچنین مجاز هستید از تمام کتابخانه هایی که در طول تمرین ها از آنها استفاده کرده اید مانند numpy، scipy و pytorch استفاده نمایید اما دقت کنید پیاده سازی الگوریتم ها باید توسط شما انجام شده باشد و نمیتوانید از کتابخانه هایی که الگوریتم را به صورت آماده پیاده سازی کرده اند استفاده نمایید.
۸. تحویل پروژه به صورت گزارش و کدهای نوشته شده است. گزارش باید شامل پاسخ پرسش ها، تصاویر و نمودارها و نتیجه گیری های لازم باشد. در نهایت یک فایل شامل کد ها و یک گزارش به فرمت pdf را در سامانه CW آپلود نمایید. آپلود کردن پروژه توسط یکی از اعضای گروه کافی میباشد.
۹. اگر برای پاسخ به پرسش ها، از منبعی (کتاب، مقاله، سایت و...) کمک گرفته اید، حتماً به آن ارجاع دهید.
۱۰. در صورت مشاهده ی تقلب، نمره ی هردو فرد صفر منظور خواهد شد.
۱۱. مسئول پروژه آقای سلیمان بیگی میباشد و در صورت داشتن مشکلاتی در گروه بندی، زمان تحویل حضوری و ... به ایشان (@amirr62a) پیام دهید.
۱۲. در صورت داشتن پرسش در بخش ۱ به @amirrezazameni، در بخش ۲ به @Mahdi\_h721 و در بخش های ۳ و ۴ به @BornaKhodabandeh پیام دهید.

موفق باشید!