

پروژه فاز 1

پرهام گیلانی – 400101859

صدرا خنجری – 400101107

بخش تئوری

1. توابع پتانسیل (ψ) در مدل میدان مارکوف (MRF) نشان‌دهنده تأثیر متقابل بین گره‌های یک گروه خاص (clique) در گراف است. این توابع نقش مهمی در توزیع احتمال مشترک متغیرها دارند و توزیع کلی را به صورت زیر فاکتوریزه می‌کنند:

$$p(x) = \psi_C(x_C) \prod_{C \in \text{Cliques}} \frac{1}{Z}$$

که در آن:

- Z یک ثابت نرمال‌سازی است.
- $\psi_C(x_C)$ توابع پتانسیل برای کلیک C است.
- x_C مقادیر متغیرهای گره‌های کلیک C است.

این توابع معمولاً مثبت و غیرمنفی هستند و امکان محاسبه احتمال‌های شرطی را با استفاده از قانون مارکوف فراهم می‌کنند.

مثال ساده

فرض کنید گراف شامل سه گره A, B, C است که به صورت زنجیره‌ای به هم متصل‌اند: A با B و B با C متصل است.

فرض می‌کنیم توابع پتانسیل به شکل زیر تعریف شده‌اند:

$$\psi_{AB}(A, B) = e^{-\frac{AB}{2}}, \psi_{BC}(B, C) = e^{-0.3BC}$$

مقادیر ممکن برای A, B, C باینری هستند، یعنی هر کدام می‌توانند 0 یا 1 باشند.

توزیع احتمال کلی به صورت زیر محاسبه می‌شود:

$$p(A, B, C) = \frac{1}{Z} \psi_{AB}(A, B) \psi_{BC}(B, C)$$

که در آن Z ثابت نرمال سازی است و از مجموع همه مقادیر ممکن محاسبه می‌شود.

محاسبات با مقادیر عددی

تمام ترکیبات ممکن برای A, B, C

A	B	C	$\psi_{AB}(A, B)$	$\psi_{BC}(B, C)$	$\psi_{AB}(A, B) \psi_{BC}(B, C)$
0	0	0	$e^0 = 1$	$e^0 = 1$	$1 \times 1 = 1$
0	0	1	$e^0 = 1$	$e^0 = 1$	$1 \times 1 = 1$
0	1	0	$e^0 = 1$	$e^0 = 1$	$1 \times 1 = 1$
0	1	1	$e^0 = 1$	$e^{-0.3} \approx 0.74$	$1 \times 0.74 = 0.74$
1	0	0	$e^0 = 1$	$e^0 = 1$	$1 \times 1 = 1$
1	0	1	$e^0 = 1$	$e^0 = 1$	$1 \times 1 = 1$
1	1	0	$e^{-0.5} \approx 0.61$	$e^0 = 1$	$1 \times 0.64 = 0.64$
1	1	1	$e^{-0.5} \approx 0.61$	$e^{-0.3} \approx 0.74$	$0.64 \times 0.74 \approx 0.45$

$$Z = \sum_{A,B,C} \psi_{AB}(A, B) \psi_{BC}(B, C) = 1 + 1 + 1 + 0.74 + 1 + 1 + 0.64 + 0.45 = 7.8$$

برای مثال:

$$p(1,1,1) = \frac{1}{Z} \psi_{AB}(1,1) \psi_{BC}(1,1) = \frac{0.64 \times 0.74}{7.8} = 0.058$$

نتیجه‌گیری

این روش نشان می‌دهد چگونه توابع پتانسیل به کمک مقادیر عددی به محاسبه احتمال‌های میدان مارکوف کمک می‌کنند. اگر نیاز به محاسبات دقیق‌تر باشد، می‌توان تمامی ترکیبات را مشابه روش بالا محاسبه کرد.

2. توابع پتانسیل در میدان مارکوف به انواع مختلفی تقسیم می‌شوند که می‌توان به موارد زیر اشاره کرد:

توابع پتانسیل گره‌ای (Node Potential)

این توابع به متغیرهای منفرد در گره‌ها اعمال می‌شوند. به عنوان مثال، برای یک گره A ، می‌توان تعریف کرد:

$$\psi_A(A) = e^{-\alpha A}, \alpha = \text{constant}$$

اگر A فقط مقادیر 0 یا 1 بگیرد و $\alpha = 0.7$:

$$\psi_A(0) = e^{-0 \times 0.7} = 1, \psi_A(1) = e^{-1 \times 0.7} \approx 0.5$$

توابع پتانسیل یالی (Edge Potential)

این توابع تأثیر متقابل بین دو گره متصل را نشان می‌دهند. برای مثال، گره‌های متصل A و B :

$$\psi_{AB}(A, B) = e^{-\beta AB}, \beta = \text{constant}$$

اگر A, B مقادیر 0 یا 1 بگیرند و $\beta = 0.5$:

$$\begin{aligned} \psi_{AB}(0,0) &= e^{-\frac{0 \times 0}{2}} = 1, \psi_{AB}(0,1) = e^{-\frac{0 \times 1}{2}} = 1 \\ \psi_{AB}(1,0) &= e^{-\frac{1 \times 0}{2}} = 1, \psi_{AB}(1,1) = e^{-\frac{1 \times 1}{2}} \approx 0.61 \end{aligned}$$

توابع پتانسیل گروه خاص (Clique Potential)

این توابع برای گروه‌های خاص (cliques) در گراف تعریف می‌شوند و تأثیر متقابل بیش از دو گره را نشان می‌دهند. به عنوان مثال، برای یک گروه خاص شامل سه گره: A, B, C

$$\psi_{ABC}(A, B, C) = e^{-\gamma(A+B+C)}, \gamma = \text{constant}$$

اگر A, B, C مقادیر 0 یا 1 بگیرند و $\gamma = 0.2$:

$$\begin{aligned} \psi_{ABC}(0,0,0) &= e^{-0.2 \times (0+0+0)} = 1, \psi_{ABC}(0,0,1) = e^{-0.2 \times (0+0+1)} \approx 0.82 \\ \psi_{ABC}(0,1,1) &= e^{-0.2 \times (0+1+1)} \approx 0.67, \psi_{ABC}(1,1,1) = e^{-0.2 \times (1+1+1)} \approx 0.55 \end{aligned}$$

توابع پتانسیل می‌توانند برای گره‌ها، یال‌ها، یا گروه‌های خاص تعریف شوند.

در هر مورد، مقدار آن‌ها به متغیرهای گره‌ها و ثابت‌های مشخص وابسته است.

محاسبات عددی این توابع به ما کمک می‌کند تا ارتباط متغیرها در میدان مارکوف را بهتر درک

کنیم.

3. روش **MLE** پارامترهای مدل مارکوف را با بیشینه‌کردن احتمال داده‌های مشاهده شده تخمین می‌زند. تابع درست‌نمایی $\ln L(\theta|S)$ شامل توابع پتانسیل و ثابت نرمال‌سازی $Z(\theta)$ است. گرادیان برای بیشینه‌سازی محاسبه و از روش‌هایی مثل **Gradient Ascent** استفاده می‌شود. در مثال عددی، با داده‌های $S = \{(1,1), (1,0), (0,0)\}$ ، ثابت $Z(w) = 3 + e^{-w}$ و $\ln L(\theta|S)$ محاسبه شده و w بهینه‌سازی می‌شود.

4. گرادیان لگاریتم درست‌نمایی برای یک مدل مارکوف با متغیر پنهان به صورت زیر تعریف می‌شود:

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

که v متغیرهای مشاهده‌شده، h متغیرهای پنهان، $E(v, h)$ انرژی سیستم، Z ثابت نرمال‌سازی است.

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

محاسبه گرادیان

$$\ln L(\theta|v) = \ln \left(\sum_h e^{-E(v,h)} \right) - \ln \left(\sum_{v,h} e^{-E(v,h)} \right)$$

مشتق نسبت به θ :

$$\frac{\partial \ln L(\theta|v)}{\partial \theta} = \frac{\partial}{\partial \theta} \ln \left(\sum_h e^{-E(v,h)} \right) - \frac{\partial}{\partial \theta} \ln \left(\sum_{v,h} e^{-E(v,h)} \right)$$

گرادیان بخش اول:

$$\frac{\partial}{\partial \theta} \ln \left(\sum_h e^{-E(v,h)} \right) = \frac{1}{\sum_h e^{-E(v,h)}} \sum_h \left(-\frac{\partial E(v,h)}{\partial \theta} e^{-E(v,h)} \right)$$

این عبارت را می‌توان به صورت زیر نوشت:

$$p(h|v) = \frac{e^{-E(v,h)}}{\sum_h e^{-E(v,h)}} \rightarrow \frac{\partial}{\partial \theta} \ln \left(\sum_h e^{-E(v,h)} \right) = - \sum_h p(h|v) \frac{\partial E(v,h)}{\partial \theta}$$

گرادیان بخش دوم:

$$\frac{\partial}{\partial \theta} \ln \left(\sum_{v,h} e^{-E(v,h)} \right) = \sum_{v,h} p(v,h) \frac{\partial E(v,h)}{\partial \theta}, p(v,h) = \frac{e^{-E(v,h)}}{Z}$$

ترکیب دو بخش:

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

5. کاربرد های ماشین بولتزمن:

یادگیری ویژگی های پنهان (Unsupervised Learning): ماشین بولتزمن می تواند برای یادگیری ویژگی های پنهان در داده های بدون برچسب استفاده شود. با استفاده از این مدل، می توان ساختارهای پنهان در داده ها را شناسایی و ویژگی های مهم را استخراج کرد.

شبکه های مولد (Generative Models): ماشین بولتزمن یک مدل مولد است که می تواند برای تولید داده های جدید مشابه داده های آموزشی استفاده شود. به این صورت که می تواند نمونه هایی شبیه به داده های ورودی تولید کند، به ویژه در مسائلی مانند تولید تصاویر، متن یا صدا.

شبکه های عصبی پنهان (Deep Belief Networks - DBN): ماشین بولتزمن یکی از اجزای اصلی شبکه های عصبی عمیق (Deep Networks) به نام شبکه های باور عمیق (DBN) است. در DBN، از چندین لایه ماشین بولتزمن برای یادگیری ویژگی های پیچیده تر در سطوح مختلف داده ها استفاده می شود.

مدل سازی توزیع های پیچیده: ماشین بولتزمن قادر است توزیع های پیچیده و غیرخطی را مدل سازی کند. به عنوان مثال، در مدل سازی توزیع های پنهان یا توزیع های احتمال در داده هایی که دارای روابط پیچیده ای هستند، مانند داده های تصویری یا صوتی، می تواند مفید باشد.

یادگیری تقویتی (Reinforcement Learning): در برخی از مسائل یادگیری تقویتی، ماشین بولتزمن می تواند برای انتخاب عمل ها یا سیاست های بهینه استفاده شود. به ویژه در مدل های دارای محیط های پیچیده و فضای حالت بزرگ، استفاده از ماشین بولتزمن برای مدل سازی توزیع های پیچیده و بهینه سازی انتخاب ها کاربرد دارد.

کاهش ابعاد: (Dimensionality Reduction) از ماشین بولتزمن می‌توان برای کاهش ابعاد داده‌ها و استخراج ویژگی‌های مهم و مختصر از داده‌های پیچیده استفاده کرد. به طور خاص، مدل‌های مولد مانند ماشین بولتزمن می‌توانند به شبیه‌سازی ویژگی‌های پنهان داده‌های پیچیده بپردازند و ابعاد داده را کاهش دهند.

6. توزیع مشترک احتمالات بین نرون‌های مشاهده پذیر v و نرون‌های پنهان h به صورت زیر تعریف می‌شود:

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

$$E(v, h) = - \sum_{i, j} w_{ij} h_i v_j - \sum_j b_j v_j - \sum_i c_i h_i$$

اثبات برای $P(H_i = 1|v)$

برای محاسبه از تعریف توزیع شرطی استفاده می‌کنیم:

$$p(H_i = 1|v) = \frac{p(H_i = 1, v)}{p(v)}$$

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

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

$$E(v, h) = -H_i \left(\sum_j w_{ij} v_j + c_i \right) + \text{terms independent to } H_i$$

پس $p(H_i = 1|v)$ به صورت یک توزیع برنولی نتیجه میشود.

$$p(H_i = 1|v) = \sigma \left(\sum_j w_{ij} v_j + c_i \right)$$

اثبات برای $p(v_j = 1|h)$

اثبات $p(v_j = 1|h)$ مشابه است، با این تفاوت که توزیع شرطی بدین گونه است:

$$p(v_j = 1|h) = \frac{p(v_j = 1, h)}{p(h)} \rightarrow p(v_j = 1|h) = \sigma \left(\sum_j w_{ij} v_j + c_i \right)$$

7. ابتدا انرژی کلی سیستم را یادآوری می‌کنیم:

$$E(v, h) = - \sum_{i,j} w_{ij} h_i v_j - \sum_j b_j v_j - \sum_i c_i h_i$$

مشتق نسبت به w_{ij}

فقط ترم اول انرژی شامل w_{ij} است:

$$E(v, h) = -w_{ij} h_i v_j + (\text{terms independent of } w_{ij}) \rightarrow$$

$$\frac{\partial E(v, h)}{\partial w_{ij}} = -h_i v_j$$

جایگذاری در فرمول کلی مشتق:

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

به صورت ساده‌شده:

$$\frac{\partial \ln L(\theta|v)}{\partial w_{ij}} = \langle h_i v_j \rangle_{data} - \langle h_i v_j \rangle_{model}$$

اینجا:

$\langle h_i v_j \rangle_{data}$ امید ریاضی تحت توزیع شرطی $p(h|v)$

$\langle h_i v_j \rangle_{model}$ امید ریاضی تحت توزیع مدل $p(v, h)$

مشتق نسبت به b_j

فقط ترم دوم انرژی شامل b_j است:

$$E(v, h) = -b_j v_j + (\text{terms independent of } b_j) \rightarrow$$

$$\frac{\partial E(v, h)}{\partial b_j} = -v_j$$

جایگذاری در فرمول کلی مشتق:

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

به صورت ساده‌شده:

$$\frac{\partial \ln L(\theta|v)}{\partial w_{ij}} = \langle v_j \rangle_{data} - \langle v_j \rangle_{model}$$

مشتق نسبت به c_i

فقط ترم سوم انرژی شامل c_i است:

$$E(v, h) = -c_i h_i + (\text{terms independent of } c_i) \rightarrow$$

$$\frac{\partial E(v, h)}{\partial c_i} = -h_i$$

جایگذاری در فرمول کلی مشتق:

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

به صورت ساده‌شده:

$$\frac{\partial \ln L(\theta|v)}{\partial c_i} = \langle h_i \rangle_{data} - \langle h_i \rangle_{model}$$

8. تعریف تابع هدف

$$\begin{aligned} \ln L(\theta|v) &= \ln P(v|\theta) \\ p(v|\theta) &= \frac{1}{Z} \sum_h e^{-E(v,h)}, Z = \sum_{v,h} e^{-E(v,h)} \rightarrow \\ \ln L(\theta|v) &= \ln \left(\sum_h e^{-E(v,h)} \right) - \ln Z \end{aligned}$$

مشتق لگاریتم درست‌نمایی نسبت به w_{ij}

$$\begin{aligned} \frac{\partial}{\partial w_{ij}} \ln \left(\sum_h e^{-E(v,h)} \right) &= \frac{1}{\sum_h e^{-E(v,h)}} \frac{\partial}{\partial w_{ij}} \sum_h e^{-E(v,h)} \\ \frac{\partial}{\partial w_{ij}} \sum_h e^{-E(v,h)} &= \sum_h \frac{\partial}{\partial w_{ij}} e^{-E(v,h)} = - \sum_h e^{-E(v,h)} \frac{\partial E(v,h)}{\partial w_{ij}} \\ \frac{\partial E(v,h)}{\partial w_{ij}} &= -v_i h_j \rightarrow \frac{\partial}{\partial w_{ij}} \ln \left(\sum_h e^{-E(v,h)} \right) = \frac{\sum_h e^{-E(v,h)} v_i h_j}{\sum_h e^{-E(v,h)}} \rightarrow \\ \frac{\partial}{\partial w_{ij}} \ln \left(\sum_h e^{-E(v,h)} \right) &= \sum_h p(h|v) v_i h_j \end{aligned}$$

مشتق بخش دوم:

$$\begin{aligned}\frac{\partial}{\partial w_{ij}} \ln Z &= \frac{1}{Z} \frac{\partial Z}{\partial w_{ij}} \\ \frac{\partial Z}{\partial w_{ij}} &= \sum_{v,h} \frac{\partial}{\partial w_{ij}} e^{-E(v,h)} = - \sum_{v,h} e^{-E(v,h)} \frac{\partial E(v,h)}{\partial w_{ij}} \\ \frac{\partial E(v,h)}{\partial w_{ij}} &= -v_i h_j \rightarrow \frac{\partial Z}{\partial w_{ij}} = \sum_{v,h} e^{-E(v,h)} v_i h_j \rightarrow \\ \frac{\partial}{\partial w_{ij}} \ln Z &= \frac{\sum_{v,h} e^{-E(v,h)} v_i h_j}{Z} = \sum_{v,h} p(v,h) v_i h_j\end{aligned}$$

ترکیب دو بخش

$$\begin{aligned}\frac{\partial}{\partial w_{ij}} \ln L(\theta|v) &= \sum_h p(h|v) v_i h_j - \sum_{v,h} p(v,h) v_i h_j \\ &= \langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_{model}\end{aligned}$$

اینجا:

$\langle h_i v_j \rangle_{data}$ امید ریاضی تحت توزیع شرطی $p(h|v)$

$\langle h_i v_j \rangle_{model}$ امید ریاضی تحت توزیع مدل $p(v,h)$

9. تابع کلاسبندی متقاطع کلبِرگ-لیبلر: (KL Divergence)

برای دو توزیع احتمال P و Q، $D_{KL}(P||Q)$ به صورت زیر تعریف می‌شود:

$$D_{KL}(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

توزیع‌های مربوطه:

P_{data} توزیع داده‌های واقعی.

P_θ توزیع مدل بولتزمن پارامتری.

P_k^θ توزیع حاصل از k مرحله نمونه‌برداری گیبس. (Gibbs Sampling)

تابع $CD_k(\theta, v_0)$

$$CD_k(\theta, v_0) = - \sum_h P(h, v_0) \frac{\partial E(v_0, h)}{\partial \theta} + \sum_h P(h, v_k) \frac{\partial E(v_k, h)}{\partial \theta}$$

که v_k خروجی حاصل از k مرحله نمونه برداری گیبس است.

اثبات معادل بودن $CD_k(\theta, v_0)$ با $D_{KL}(P_k^\theta || P_\theta) - D_{KL}(P_{data} || P_\theta)$

مرحله 1: مشتق $D_{KL}(P_{data} || P_\theta)$

$$D_{KL}(P_{data} || P_\theta) = \sum_v P(v) \log \frac{P_{data}(v)}{P_\theta(v)} \rightarrow$$

$$\frac{\partial}{\partial \theta} D_{KL}(P_{data} || P_\theta) = - \sum_v P_{data}(v) \frac{\partial}{\partial \theta} \ln P_\theta(v), P_\theta(v) = \frac{1}{Z} \sum_h e^{-E(v, h)}$$

$$\rightarrow \frac{\partial}{\partial \theta} \ln P_\theta(v) = \sum_h P(h|v) \frac{\partial E(v, h)}{\partial \theta} - \sum_{v, h} P(h, v) \frac{\partial E(v, h)}{\partial \theta}$$

مرحله 2: مشتق $D_{KL}(P_k^\theta || P_\theta)$

$$D_{KL}(P_k^\theta || P_\theta) = \sum_v P_k^\theta(v) \log \frac{P_k^\theta(v)}{P_\theta(v)} \rightarrow$$

$$\frac{\partial}{\partial \theta} D_{KL}(P_k^\theta || P_\theta) = - \sum_v P_k^\theta(v) \frac{\partial}{\partial \theta} \ln P_\theta(v) \rightarrow$$

$$\frac{\partial}{\partial \theta} D_{KL}(P_k^\theta || P_\theta) = \sum_v P_k^\theta(v) \left(\sum_{v',h} P(h, v') \frac{\partial E(v', h)}{\partial \theta} - \sum_h P(h|v) \frac{\partial E(v, h)}{\partial \theta} \right)$$

مرحله 3: ترکیب دو مشتق

با ترکیب نتایج، مشتق ترکیبی برای اختلاف دو KL Divergence به صورت زیر است:

$$\begin{aligned} & \frac{\partial}{\partial \theta} (D_{KL}(P_{data} || P_\theta) - D_{KL}(P_k^\theta || P_\theta)) = \\ & - \sum_v P_k^\theta(v) \left(\sum_h P(h|v) \frac{\partial E(v, h)}{\partial \theta} - \sum_{v',h} P(h, v') \frac{\partial E(v', h)}{\partial \theta} \right) \end{aligned}$$

این معادله دقیقاً با تعریف $CD_k(\theta, v_0)$ مطابقت دارد، زیرا:

بخش اول: (P_{data}) توزیع داده.

بخش دوم: (P_k^θ) توزیع تخمینی حاصل از نمونه‌برداری گیبس.

پس تابع $CD_k(\theta, v_0)$ معادل با $D_{KL}(P_{data} || P_\theta) - D_{KL}(P_k^\theta || P_\theta)$ است.

10. شرح الگوریتم k-step Contrastive Divergence

الگوریتم CD برای یادگیری مدل‌های انرژی مانند Restricted Boltzmann Machine (RBM) طراحی شده و بر نمونه‌برداری زنجیره مارکوف (MCMC) مبتنی است:

در CD، به جای استفاده از تعداد زیادی گام مارکوف برای همگرا شدن به توزیع پایدار، از تعداد محدودی گام k استفاده می‌شود.

مراحل اصلی الگوریتم:

- شروع با داده واقعی: (v_0) نمونه‌هایی از داده‌های آموزشی به مدل ورودی داده می‌شوند.
- نمونه‌برداری مارکوف: از حالت اولیه v_0 ، با استفاده از توزیع شرطی مدل $P(h|v)$ و $P(v|h)$ ، گام k نمونه‌برداری می‌شود تا v_k به دست آید.
- محاسبه گرادیان‌ها:

$$\Delta\theta \propto E_{data}[\nabla E(v_0, h) - E_{model}[\nabla E(v_k, h)]]$$

که در آن $E(v, h)$ انرژی مدل است.

- به‌روزرسانی پارامترها: گرادیان محاسبه‌شده برای بهینه‌سازی پارامترها استفاده می‌شود.

پیچیدگی الگوریتم و توجیه آن

پیچیدگی زمانی:

الگوریتم CD برای هر گام k نمونه‌برداری، به محاسبه شرطی $P(h|v)$ و $P(v|h)$ نیاز دارد. این محاسبات به تعداد واحدهای نورونی در لایه‌های پنهان و نمایان بستگی دارد. پیچیدگی کلی برابر است با $O(k \times (m \times n))$ که m تعداد نورون‌های لایه نمایان و n تعداد نورون‌های لایه پنهان است.

مزایا:

برخلاف Gibbs Sampling کامل که به تعداد زیادی گام برای همگرا شدن نیاز دارد، CD با مقدار کوچک (مثلاً $k = 1$) سرعت اجرا را به طور قابل توجهی افزایش می‌دهد و تقریب ساده‌ای ارائه می‌دهد که در عمل برای یادگیری پارامترها مؤثر است.

معایب:

با k کوچک، همگرایی کامل به توزیع مدل تضمین نمی‌شود و انتخاب مقدار k تأثیر مستقیمی بر کیفیت یادگیری دارد.

انتظارات از مدل خروجی

بازسازی داده‌ها: مدل باید بتواند داده‌هایی مشابه داده‌های آموزشی ورودی تولید کند.

نزدیکی توزیع‌ها: توزیع داده‌های تولیدی مدل به توزیع داده‌های آموزشی نزدیک تر می‌شود.

کاهش خطای بازسازی: تفاوت بین داده اصلی و داده بازسازی شده (reconstruction error) باید در طول زمان کاهش یابد.

تولید نمونه‌های جدید: مدل توانایی تولید داده‌هایی که شبیه داده‌های آموزشی باشند را پیدا می‌کند.

11. ماشین‌های بولتزمن (BM) و مدل‌های مشابه مانند ماشین‌های بولتزمن محدود (RBM)، به طور سنتی برای کار با داده‌های باینری طراحی شده‌اند. در این مدل‌ها، متغیرهای نمایان (v) و پنهان (h) به صورت باینری در نظر گرفته می‌شوند و تابع انرژی مدل به این صورت است که به طور طبیعی برای داده‌های باینری بهینه می‌شود. فرمول انرژی در این مدل‌ها معمولاً به شکل زیر است:

$$E(v, h) = - \sum_{i,j} w_{ij} h_i v_j - \sum_j b_j v_j - \sum_i c_i h_i$$

استفاده از ماشین‌های بولتزمن برای داده‌های پیوسته (که مقادیر آن‌ها می‌توانند هر عددی از یک بازه خاص باشند) نیاز به برخی تغییرات در مدل و نحوه محاسبات دارد. در اینجا چند پیشنهاد برای استفاده از این مدل‌ها برای داده‌های پیوسته آورده شده است:

تغییر در توزیع‌ها:

برای داده‌های باینری، توزیع‌های برنولی معمولاً برای مدل سازی استفاده می‌شود. اما برای داده‌های پیوسته، باید از توزیع‌های پیوسته مانند توزیع گوسی استفاده کرد و برای داده‌های پیوسته، متغیرهای v و h می‌توانند از توزیع‌های نرمال (گوسی) به جای توزیع‌های باینری پیروی کنند.

به عنوان مثال، فرض کنید که v_i و h_j به جای مقادیر باینری، از مقادیر پیوسته با توزیع نرمال $v_i \sim N(0, \sigma_v^2)$ پیروی کنند.

تغییر در تابع انرژی:

برای داده‌های پیوسته، می‌توان تابع انرژی را به گونه‌ای تغییر داد که شامل متغیرهای پیوسته باشد. برای مثال، تابع انرژی به صورت زیر در نظر گرفته می‌شود:

$$E(v, h) = - \sum_{i,j} w_{ij} v_i h_j - \sum_i b_i v_i - \sum_j c_j h_j$$

این تابع می‌تواند به گونه‌ای تغییر یابد که متغیرهای v_i و h_j به جای مقادیر باینری، مقادیر پیوسته بگیرند.

استفاده از الگوریتم‌های نمونه‌برداری مناسب:

برای داده‌های باینری، از نمونه برداری گیبز برای تولید نمونه‌های جدید استفاده می‌شود. برای داده‌های پیوسته، به دلیل اینکه توزیع‌های گوسی داریم، باید از نمونه برداری گوسی یا روش‌هایی مانند نمونه برداری از زنجیره‌های گوسی مارکوف استفاده کرد.

(Gibbs sampling for Gaussian distributions)

می‌توان از الگوریتم‌هایی مانند Contrastive Divergence (CD) و Gibbs Sampling برای داده‌های پیوسته نیز استفاده کرد، اما به طور خاص برای داده‌های پیوسته، باید روش‌های نمونه برداری را تغییر داد تا با توزیع‌های پیوسته سازگار شوند.

تغییرات مورد نیاز در الگوریتم‌ها:

در الگوریتم Contrastive Divergence، اگر داده‌های ورودی به صورت پیوسته باشند، به جای استفاده از توزیع برنولی، باید از توزیع‌های گوسی برای نمونه برداری استفاده کرد. در الگوریتم Gibbs Sampling، در هر گام، برای متغیرهای پیوسته، از نمونه برداری گوسی استفاده می‌شود، نه نمونه برداری باینری.

پس ماشین‌های بولتزمن می‌توانند برای داده‌های پیوسته با تغییرات در توزیع‌ها و الگوریتم‌های نمونه برداری استفاده شوند. به طور خاص، استفاده از توزیع‌های گوسی به جای توزیع‌های باینری و تغییر در الگوریتم‌های Gibbs Sampling و Contrastive Divergence برای داده‌های پیوسته ضروری است.