

Hidden Markov Model

隐马尔可夫模型



杨猛, 郑伟诗

<https://cse.sysu.edu.cn/content/2970>

中山大学

机器智能与先进计算
教育部重点实验室

声明：该PPT只供非商业使用，也不可视为任何出版物。由于历史原因，许多图片尚没有标注出处，如果你知道图片的出处，欢迎告诉我们 at wszheng@ieee.org.

模拟真实世界的现象

- 可以用来检测一个新建的水坝溢流的频率（取决于连续下雨的天数）。为建立这个模型，可以从下面的雨天（R）和晴天（S）开始：

R S S S R R R R R R R R S S S S S R R R R R R R R R R R R R R R

R S R S S S R S R R S R S S S R S S S R R S R S S S S R R S S

时间序列 Times Series

- 随机过程 $\{X_1, X_2, \dots\}$, $X_i \in \mathcal{X}$
 - \mathcal{X} 称为状态空间, 我们假设 $\mathcal{X} = \{1, 2, \dots, N\}$
 - 假设对所有的 i , \mathcal{X} 都相同
 - 假设只处理时间序列, 即 i 代表时间
 - 随机性
- 目的是希望“过去”对“现在”有帮助
 - 即如果有对 X_1, \dots, X_{t-1} 的了解, 能帮助确定 X_t
 - Formally, $P(X_t | X_{1:t-1})$ vs. $P(X_t)$

Markov Property

□ Curse of dimensionality

- $P(X_2|X_1)$ 需要多少存储空间才能指定?
- $P(X_3|X_2, X_1)$ 需要多少存储空间才能指定?
- $P(X_t|X_{1:t-1})$ 需要多少存储空间才能指定?
❖ $N^t!$

□ Markov Property 马尔科夫性质

- 限定: $P(X_t|X_{1:t-1}) = P(X_t|X_{t-1})$, 含义是?
- 无记忆性 memoryless
- 这个假设有效吗?
- 好处是什么?

人物介绍

□ Andrey Markov



□ https://en.wikipedia.org/wiki/Markov_chain

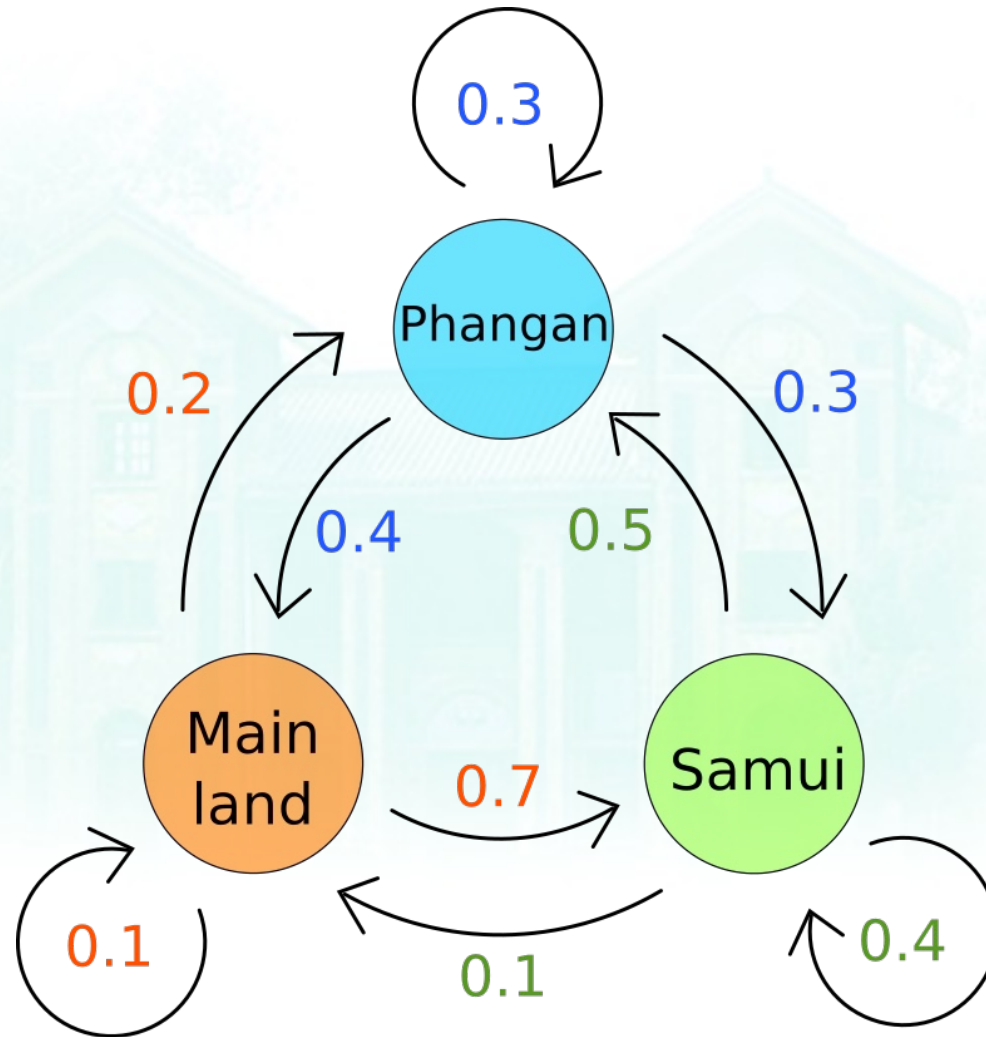
Computational finance
Speech synthesis
Cryptanalysis

Speech recognition
Part-of-speech tagging
Handwriting recognition

Speech synthesis
Single-molecule kinetic analysis
Machine translation

Markov Chain 马尔科夫链

- Markov chain (discrete-time Markov chain or DTMC)



Markov Chain 马尔科夫链

□ Markov chain (discrete-time Markov chain or DTMC)



	A	B
A	$P(A A):0.50$	$P(B A):0.50$
B	$P(A B):0.50$	$P(B B):0.50$

状态空间中有A和B两种状态。共4种可能的转换。

1. 在A时，可以过渡到B或留在A。
2. 在B时，可以过渡到A或者留在B。

在图中，从任意状态到任意状态的转移概率是0.5。

人们会通过使用“转移矩阵”来计算转移概率。状态空间的每个状态都会出现在表格中的一列或者一行中。矩阵的每个单元格指明了从行状态转换到列状态的概率。

状态空间新增一个状态，矩阵将对应增加一行和一系列，向现有的列和行中添加一个单元格。这意味着当我们向马尔可夫链添加状态时，单元格的数目会呈二次方增长。因此，转换矩阵就起到了很大的作用。

Markov Chain 马尔科夫链

□ Markov chain (discrete-time Markov chain or DTMC)

马尔科夫链的一个作用是用计算机模拟现实世界中的现象。

例如，可以用来检测一个新建的水坝溢流的频率（取决于连续下雨的天数）。

为建立这个模型，可以从下面的雨天（R）和晴天（S）开始：

R S S S R R R R R R R R S S S S S R R R R R R R R R R R R R R R R

表述这种模拟天气的方法就是：“有一半的天数是下雨天。所以模拟中的每一天都有50%的概率是下雨的。”这个规则在模拟中所产生的序列如下：

R S R S S S R S R R S R S S S R S S S R R S R S S S S R R S S S

你注意到上面的序列和原来的不太一样了吗？第二个序列似乎具有跳跃性，而第一个（真实数据）似乎具有“粘性”。在真实的数据中，如果某一天是晴天，那么第二天也很可能是晴天。

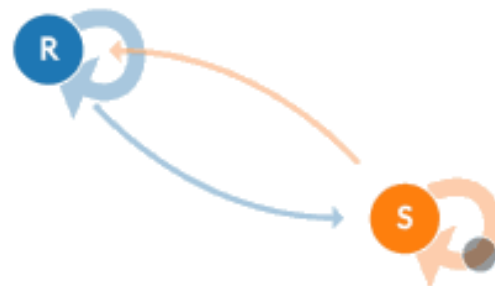
Markov Chain 马尔科夫链

□ Markov chain (discrete-time Markov chain or DTMC)

通过两个状态的马尔可夫链来消除这种“粘性”。

当马尔科夫链处于状态“R”时，它保持在该状态的概率是0.9，状态改变的概率是0.1。同样，“S”状态保持不变的概率是0.9，过渡到“R”状态的概率是0.1。

S S S R R R S R R R R R R S S S



在许多需要对大规模的现象做研究的工作人员手中，马尔科夫链的作用可以变得非常强大。例如，谷歌用于确定搜索结果顺序的算法，称为PageRank，就是一种马尔可夫链。

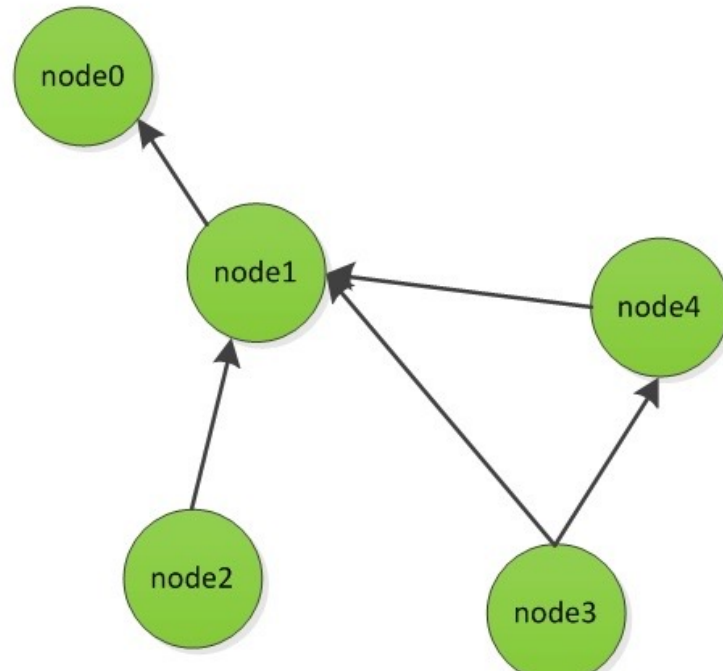
PageRank

- ❑ 随着时代的发展，网页变得越来越多，用人工识别的方式对网页进行识别变得越来越不现实。
- ❑ Google的两名创始人Larry Page与Sergey Brin开始对网页排序问题进行研究，依据学术界评判论文的重要程度的方法，查看论文引用次数，将这种方法用到了网页排序中，PageRank算法就产生了。

问题?

PageRank 搜索算法——马氏链视角

- 下是一个有向图，包含了5个节点，以及5条边。边的起点是一个网页、论文或者人，终点指向的是起点所引用的网页、论文或者人。Node1节点引用node0节点，表示前者从后者获取信息、知识、权力或者财富。



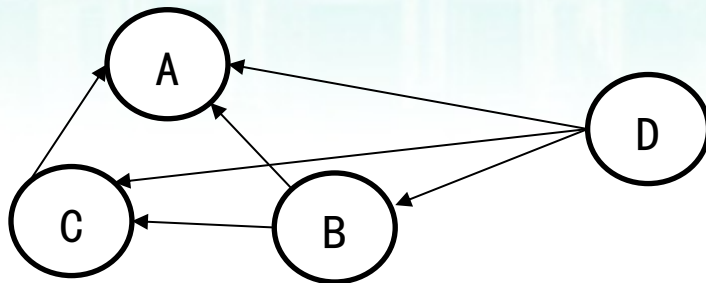
PageRank——马氏链视角

- 假设一个由4个网页组成的群体：A，B，C和D。如果所有页面都只链接至A，那么A的PR（PageRank）值将是B，C及D的Pagerank总和。

$$PR(A) = PR(B) + PR(C) + PR(D)$$

- 假设B链接到A和C，C只链接到A，并且D链接到全部其他的3个页面。一个页面总共只有一票。所以B给A和C每个页面半票。以同样的逻辑，D投出的票只有三分之一算到了A的PageRank上。

$$PR(A) = \frac{PR(B)}{2} + \frac{PR(C)}{1} + \frac{PR(D)}{3}$$



完整的PageRank——马氏链视角

问题：由于存在一些出链为0，也就是那些不链接任何其他网页的网，也称为孤立网页，使得很多网页能被访问到

假设有N个页面，对于一个页面A，那么它的PR值为：

$$PR(A) = (1 - d)/N + d \sum_{i=1}^n \frac{PR(T_i)}{C(T_i)}$$

- d 为阻尼系数，其意义是，在任意时刻，用户到达某页面后并继续向后浏览的概率。



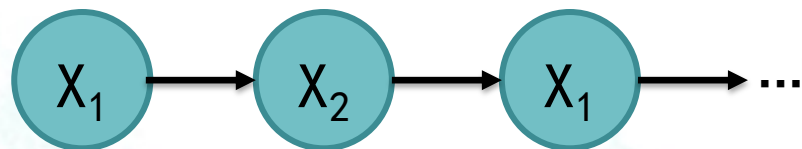
问题？

PageRank——马氏链视角

- 经过迭代计算至收敛，最终页面得到的PR值为该页面的分数。
- 网站被按照他们的PageRank算法分数从大到小排序，那些位于序列顶端的网站被认为是值得信赖的。在这些网站汇总标签为可信赖的网站被选作种子。
- PageRank算法高排名的页面都有很高的入度，说明指向一个网页的重要网页越多这个页面就越重要。

可视化和形式化

□ 可视化:

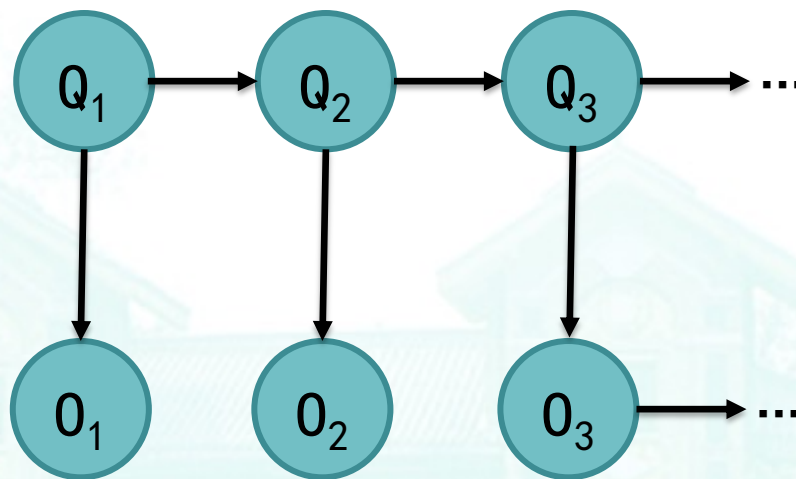


- 注意填充的变量表示观察值（即随机变量值已知）
- 那么，如何形式化定义DTMC？需要哪些量？
 - 系统初始化Initialization: $P(X_1)$ 或者 $X_1 = x_1$
 - Transition probability: $A = P(X_{t+1}|X_t)$
 - 还需要别的吗？
 - 两次运行结果会一样吗？

转移概率矩阵

- Transition probability matrix 转移概率矩阵
 - A 是一个 $N \times N$ 的矩阵
 - $A_{ij} = P(X_t = j | X_{t-1} = i)$
 - 行和为1!
- 如果运行足够久 ($t \rightarrow \infty$)，那么 X_t 的分布在很多情况下将稳定下来，叫 Stationary distribution，记为 π
 - $\pi = A\pi$

隐马尔可夫模型HMM



形式化

- Q : 隐变量(hidden variable), 不可观测的状态
- N : number of states 状态数, N 个可能的状态为 $\{S_1, \dots, S_N\}$
- $O(o)$: 观察值(observations), M 个可能的观察值 $\{V_1, V_2, \dots, V_M\}$
- T : 时间序列的长度
- π : 初始化, $\pi_j = P(Q_1 = S_j)$
- A : transition probability matrix, $A_{ij} = P(Q_{t+1} = S_j | Q_t = S_i)$
- B : emission probability 发出观察值的概率
 - $b_j(k) = \Pr(O_t = V_k | Q_t = S_j)$
 - 假设 B 不随时间变化, 当未知状态为 j 时观察到为 k 的概率
 - 那么, j, k 的取值范围是? B 的行和是?

HMM中要解决的问题

- 怎样设计状态？ -- 自动学习？
- 怎样设计观察值？ -- 根据问题的特点和实践反复设计
- 与具体问题无关
 - 指定一个HMM需要的所有参数： $\lambda = (\boldsymbol{\pi}, A, B)$
 - 问题1：Evaluation估值
 - 问题2：Decoding解码
 - 问题3：Learning学习

Problem 1. Evaluation

□ 输入

- 一个完全指定的HMM模型，即 $\lambda = (\pi, A, B)$ 已知
- 一个完全观测的输出序列 $O_1 O_2 \cdots O_T$ ，或 $\mathbf{O} = O_{1:T}$

□ 输出

- $P(\mathbf{O}|\lambda)$ - 含义是？
- 在这个模型 λ 中观察到特定输出 \mathbf{O} 的概率

□ 作用是？

- 可以看出此模型对该观察序列的成绩score
- 可以用来从多个模型中选择最适合的模型

Problem 1. Evaluation

假设状态已知

- 已知 $\lambda, o_{1:T}$, 求 $P(o_{1:T}|\lambda)$
- 若假设oracle已告知所有的隐变量的值 $q_{1:T}$
 - $\Pr(o_{1:T}|\lambda, q_{1:T}) = \prod_{i=1}^T \Pr(o_t|q_t, \lambda) = \prod_{i=1}^T b_{q_i}(o_i)$
 - 证明? 含义?
- λ 的存在只是表明概率的大小是基于该模型参数计算的, 可以去除而不影响计算

Problem 1. Evaluation

一种naïve的计算方法

- 那么隐变量序列 $q_{1:T}$ 的可能性多大呢？
 - $\Pr(q_{1:T}|\lambda) = \pi_{q_1} A_{q_1 q_2} A_{q_2 q_3} \cdots A_{q_{T-1} q_T}$
 - 含义？
- 用全概率公式对**所有可能的** $q_{1:T}$ 求和可以得到 $\Pr(o_{1:T}|\lambda)$
 - $\Pr(o_{1:T}|\lambda) = \sum_{all\ Q} \Pr(o_{1:T}|\lambda, q_{1:T})\Pr(q_{1:T}|\lambda)$, 复杂度？
 - $O(T \times N^T)$

Problem 1. Evaluation

那么，如何快速计算？

动态规划！

只看最后一步 ($t = T$)，该如何计算？

1. 最后一步 ($t = T$) 时一共可能有 N 种状态： $q_T = S_1, \dots, S_N$ ，其概率 $\Pr(o_{1:T-1}, Q_T = S_i | \lambda) = ?$
2. 若最后一步状态为 S_i ，那么观察到输出 o_T 的概率是多少？
3. 所求的值是多少？（全概率公式）

$$\Pr(o_{1:T} | \lambda) = \sum_{i=1}^N \Pr(o_{1:T-1}, Q_T = S_i | \lambda) b_i(o_T)$$

只限于最后一步吗？

Problem 1. Evaluation

如何计算 $\Pr(o_{1:T-1}, Q_T = S_i | \lambda)$?

- 有 N 种可能, 即 $T - 1$ 时刻状态为 $q_{T-1} = S_j$, $j = 1, 2, \dots, N$, 然后通过概率 A_{ji} 转移
- 全概率公式, again!

$$\begin{aligned} & \Pr(o_{1:T-1}, Q_T = S_i | \lambda) \\ &= \sum_{j=1}^N \Pr(o_{1:T-1}, Q_{T-1} = S_j | \lambda) A_{ji} \end{aligned}$$

Problem 1. Evaluation

快速计算小结

- $\Pr(o_{1:T} | \lambda) = \sum_{i=1}^N \Pr(o_{1:T-1}, Q_T = S_i | \lambda) b_i(o_T) = \sum_{i=1}^N (b_i(o_T) \sum_{j=1}^N \Pr(o_{1:T-1}, Q_{T-1} = S_j | \lambda) A_{ji})$
- 红色部分是什么？
 - 一个规模小一点儿的相同问题 ($T - 1$)
 - 但是需要对所有 j 的可能取值计算
 - 正如DTW中一样，可以通过动态规划解决，但是需要解决比原问题更多数目的小规模子问题
 - 但是，复杂的是，目前牵涉两个数值而不是一个： $\Pr(o_{1:T-1}, Q_T = S_i | \lambda)$ 和 $P(o_{1:T} | \lambda)$
 - 计算的方向应该是什么？

Problem 1. Evaluation

动态规划算法（前向forward算法）

$$\square P(o_{1:T}|\lambda) = \sum_{i=1}^N \Pr(o_{1:T-1}, Q_T = S_i | \lambda) b_i(o_T) = \sum_{i=1}^N (b_i(o_T) \sum_{j=1}^N \Pr(o_{1:T-1}, Q_{T-1} = S_j | \lambda) A_{ji})$$

定义

○ $\alpha_t(i) = P(o_{1:t}, Q_t = S_i | \lambda)$ - 含义是?

○ Initialization: $\alpha_1(i) = \pi_i b_i(o_1), \quad 1 \leq i \leq N$

○ Induction: For $1 \leq t \leq T - 1$

$$\alpha_{t+1}(i) = \left[\sum_{j=1}^N \alpha_t(j) A_{ji} \right] b_i(o_{t+1}), \quad 1 \leq i \leq N$$

○ Termination (output): $\Pr(o_{1:T} | \lambda) = \sum_{i=1}^N \alpha_T(i)$

Problem 1. Evaluation

后向算法backward algorithm

- 定义 $\beta_t(i) = \Pr(o_{t+1:T} | Q_t = S_i, \lambda)$
 - 若在时刻 t 状态为 S_i , 将来观测到 $o_{t+1:T}$ 的概率
- 初始化: $\beta_T(i) = 1, 1 \leq i \leq N$
- 反向更新: $t = T - 1, T - 2, \dots, 2, 1$

$$\beta_t(i) = \sum_{j=1}^N A_{ij} b_j(o_{t+1}) \beta_{t+1}(j), \quad 1 \leq i \leq N$$

- 输出: $\beta_1(i) = \Pr(o_{2:T} | Q_1 = S_i, \lambda)$

$$P(o_{1:T} | \lambda) = \sum_{i=1}^N \pi_i b_i(o_1) \beta_1(i)$$

Problem 2: Decoding

□ 输入

- 一个完全指定的HMM模型，即 $\lambda = (\pi, A, B)$ 已知
- 一个完全观测的输出序列 $O_1 O_2 \cdots O_T$ ，或 $\mathbf{O} = O_{1:T}$
- 某个标准criterion

□ 输出

- 一个完全指定的隐变量序列 $X_{1:T}$ 的值

□ 作用是？

- 如，语音识别中状态可能有实际意义（各音节）
 - ◆ 唯一吗？
- 可以用来观察模型结构，优化模型

Problem 2: Decoding

发现“最好”的隐变量值

- 标准1：对于每个时刻，发现其后验概率最大的状态
 - 定义 $\gamma_t(i) = \Pr(Q_t = S_i | o_{1:T}, \lambda)$ ，当观测到输出为 $o_{1:T}$ 时，时刻 t 时隐变量为第 i 个状态的后验概率
 - 那么，对于一个输出序列 $o_{1:T}$ ，选择
$$q_t = \operatorname{argmax}_{1 \leq i \leq N} \gamma_t(i), \quad t = 1, 2, \dots, T$$
 - 可能出现什么问题？
 - 不存在这样的路径 $q_{1:T}$

Problem 2: Decoding

怎样计算 γ

- $\alpha_t(i)\beta_t(i) = P(o_{1:T}, Q_t = S_i|\lambda)$
 - 为什么?

- 贝叶斯定理

$$\gamma_t(i) = \Pr(Q_t = S_i | o_{1:T}, \lambda) = \frac{\Pr(o_{1:T}, Q_t = S_i | \lambda)}{\Pr(o_{1:T} | \lambda)} = \frac{\alpha_t(i)\beta_t(i)}{\Pr(o_{1:T} | \lambda)}$$

- $\Pr(o_{1:T} | \lambda) = \sum_{i=1}^N \alpha_t(i)\beta_t(i)$ for any t !
- 三种计算方法计算 $P(o_{1:T}|\lambda)$ 了

- 或者 1) $\gamma_i = \alpha_t(i)\beta_t(i)$ 2) L1 normalize: $\gamma_i \leftarrow \frac{\gamma_i}{\sum_i \gamma_i}$

Problem 2: Decoding

寻找最大概率的路径

- 一共有 N^T 种可能的路径，有些的概率可能为0
 - 比如通过准则1得到的路径
 - 那么，如果寻找所有可能路径里面概率最大的那个呢？

$$q_{1:T} = \underset{Q_{1:T}}{\operatorname{argmax}} \Pr(Q_{1:T} | o_{1:T}, \lambda) = \underset{Q_{1:T}}{\operatorname{argmax}} \Pr(Q_{1:T}, o_{1:T} | \lambda)$$

- Naïve的方法复杂性是 N^T ，有没有更好的方法？
 - Viterbi方法

Problem 2: Decoding

Viterbi decoding

- $q_{1:T} = \underset{Q_{1:T}}{\operatorname{argmax}} \Pr(Q_{1:T}, o_{1:T} | \lambda)$

- 定义更多的子问题

$$\delta_t(i) = \max_{Q_{1:t-1}} \Pr(Q_{1:t-1}, Q_t = S_i, o_{1:t} | \lambda)$$

- 含义：当限定两个条件1) 前 t 个时刻的输出为 $o_{1:t}$ ，2) 第 t 个时刻的隐状态为第 i 个状态的时候，最佳路径所能取得的最大概率
- 怎么取得 q_t ?
 - ❖ 用另外一个变量 $\psi_t(i)$ 做记录
- 怎么从 t 进展到 $t + 1$?

Problem 2: Decoding

两个步骤

□ 从 t 进展到 $t + 1$

○ $\delta_{t+1}(i) = \max_j([\delta_t(j)A_{ji}] b_i(o_{t+1}))$

○ $\delta_{t+1}(i)$ 是概率，如果只需要发现概率最大那个状态， $b_i(o_{t+1})$?

□ 所以在时刻 $t + 1$ ，需要用另外一个变量 $\psi_t(i)$ 记录最大概率的路径在时刻 t 是哪一个状态

○ $\psi_{t+1}(i) = \operatorname{argmax}_{1 \leq j \leq N}([\delta_t(j)A_{ji}])$

Problem 2: Decoding

Viterbi算法

- 初始化: $\delta_1(i) = \pi_i b_i(o_1)$, $\psi_1(i) = 0$, $1 \leq i \leq N$
- 递归: $2 \leq t \leq T$, $1 \leq i \leq N$
$$\delta_t(i) = \max_{1 \leq j \leq N} \left([\delta_{t-1}(j) A_{ji}] b_i(o_t) \right)$$
$$\psi_t(i) = \operatorname{argmax}_{1 \leq j \leq N} \left([\delta_{t-1}(j) A_{ji}] \right)$$
- 输出:
 - 最大概率: $P^* = \max_{1 \leq i \leq n} \delta_T(i)$
 - 时刻 T 的最佳路径变量: $q_T^* = \operatorname{argmax}_{1 \leq i \leq N} (\delta_T(i))$
 - 时刻 $T-1, T-2, \dots, 2, 1$ 的最佳路径变量: $q_t^* = \psi_{t+1}(q_{t+1}^*)$

Problem 2: Decoding

分析

- ❑ 问题1的动态规划 $\alpha_{t+1}(i) = \sum_{j=1}^N \alpha_t(j) A_{ji}$
- ❑ 问题2的动态规划 $\delta_t(i) = \max_j \left([\delta_{t-1}(j) A_{ji}] b_i(o_t) \right)$
- ❑ 最重要的操作分别是sum-product和max-product
 - 其复杂性均为 N^2T
 - 和naïve方法的 TN^T 比较，极其巨大的速度提高

Problem 3: Learning

学习系统的参数

- 发现 $\lambda = (A, B, \pi)$ ，使得对于固定的 N ， T ，和观察值 \mathbf{O} ，似然 (likelihood) $P(\mathbf{O}|\lambda)$ 最大
 - 目前没有方法能发现全局最优的解
 - 常用的方法是 Baum-Welch 算法，发现一个局部最优的解

Problem 3: Learning

- 输入
 - 网络结构，状态数、输出数
 - 若干观测序列 $\{\mathbf{O}\}$
- 输出
 - 最优的参数 $\lambda = (\boldsymbol{\pi}, A, B)$ 使得 $P(\{\mathbf{O}\}|\lambda)$ 最大
- 作用
 - 显而易见
 - 最重要的问题
 - 有时候一个足够长的观测序列就够了

$$\xi_t(i, j) = Pr(Q_t = S_i, Q_{t+1} = S_j | o_{1:T}, \lambda) = \frac{\alpha_t(i) A_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{\Pr(o_{1:T} | \lambda)}$$

Baum-We lch算法

- Baum-We lch算法
- 1: 初始化参数 $\lambda^{(1)}$ (例如随机地)
- 2: $r \leftarrow 1$
- 3: **while** 似然尚未收敛 **do**
- 4: 对所有 $t(1 \leq t \leq T)$ 和所有 $i(1 \leq i \leq N)$, 使用前向过程基于 $\lambda^{(r)}$ 计算 $\alpha_t(i)$
- 5: 对所有 $t(1 \leq t \leq T)$ 和所有 $i(1 \leq i \leq N)$, 使用后向过程基于 $\lambda^{(r)}$ 计算 $\beta_t(i)$
- 6: 对所有 $t(1 \leq t \leq T)$ 和所有 $i(1 \leq i \leq N)$, 根据公式计算 $\gamma_t(i)$
- 7: 对所有 $t(1 \leq t \leq T - 1)$ 和所有 $i, j(1 \leq i, j \leq N)$, 根据表 12.1中的公式计算 $\xi_t(i, j)$
- 8: 更新参数为 $\lambda^{(r+1)}$
 - $$\pi_i^{(r+1)} = \gamma_1(i) \quad 1 \leq i \leq N$$
 - $$A_{ij}^{(r+1)} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)} \quad 1 \leq i, j \leq N$$
 - $$b_j^{(r+1)}(k) = \frac{\sum_{t=1}^T [[o_t=k]] \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)} \quad 1 \leq j \leq N \quad 1 \leq k \leq M$$
- 9: $r \leftarrow r + 1$
- 10: **end while**

怎样在模式识别中发挥更大作用

□ 语音识别

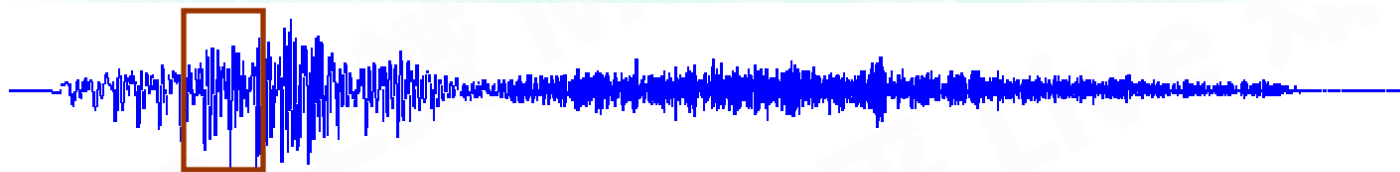
语音识别的目的是将声音信号映射为文字信息，如何实现这种映射？

分帧：声音实际上是一种波，要对声音进行分析，需要对声音分帧，也就是把声音切开成一小段一小段，每小段称为一帧。分帧操作一般不是简单的切开，而是使用移动窗函数来实现，帧与帧之间一般有交叠。

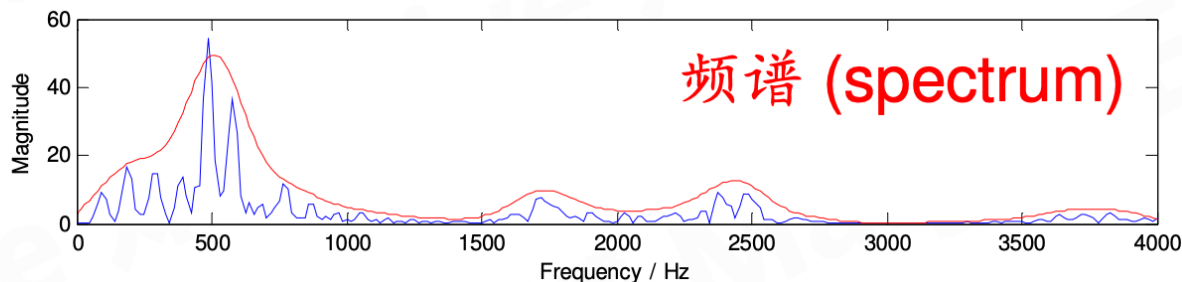
怎样在模式识别中发挥更大作用

声学特征提取

- 分帧后，语音变成了很多小段，根据人耳的生理特性，把每一帧波形变成一个多维向量，可以简单地理解为这个向量包含了这帧语音的内容信息。为什么转化为向量，因为数据驱动模型和算法基本都是从数据向量或者矩阵开始。下图为例，声音信号变成了12行（假设声学特征是12维）、N列的矩阵，每一帧都用一个12维的向量表示，色块的颜色深浅表示向量值的大小。



傅里叶变换



怎样在模式识别中发挥更大作用

HMM建模（隐马尔可夫模型）

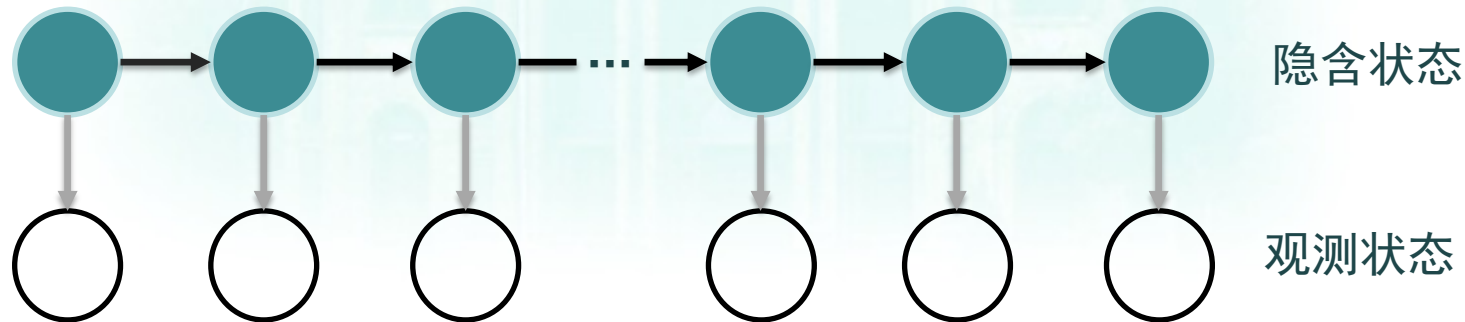
1. 求值：给定模型参数和语音，求 $P(\text{语音} | \text{模型})$
 - 把 $P(\text{语音}, \text{对齐} | \text{模型})$ 对所有对齐方式求和
 - 动态规划算法：Forward algorithm
2. 解码：给定模型参数和语音，求最佳对齐方式
 - 动态规划算法：Viterbi decoding
 - 最佳对齐方式的概率，可以作为总概率的近似
3. 训练：给定模型参数和语音，求参数模型
 - 帧就是观测序列
 - 状态就是状态序列
 - 状态变化存在转移概率，状态与观测之间存在观测概率

这就构成了一个典型的HMM模型，知道每帧语音对应的状态，得到语音识别结果。

怎样在模式识别中发挥更大作用

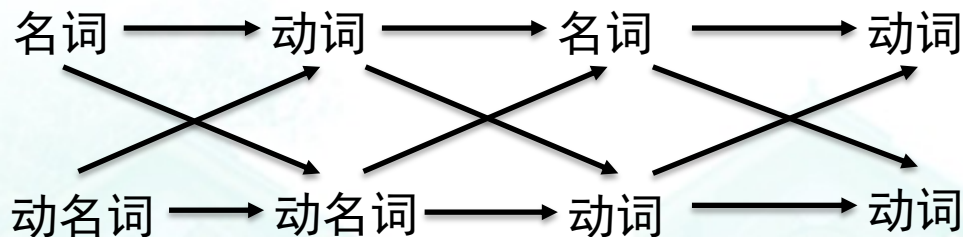
□ HMM用于NLP词性标注

- 对句子“教授喜欢画画”进行词性标注，分词之后的结果可能是“教授/喜欢/画/画”，“教授”词性可以是名词和动名词，“喜欢”词性可以是动词和动名词，“画”词性可以是名词和动词，画成图可以表示为：



怎样在模式识别中发挥更大作用

“教授喜欢画画”



- 隐马是个生成模型，生成的过程是先生成状态节点，根据状态节点再生成观测节点。
- 首先生成“教授”词性是“名词”，然后生成词“教授”；
- 根据“教授”的词性节点“名词”生成“喜欢”的词性节点“动词”，然后生成词“喜欢”；
- 根据“喜欢”的词性“动词”生成“画”的词性“动词”，然后生成词“画”。

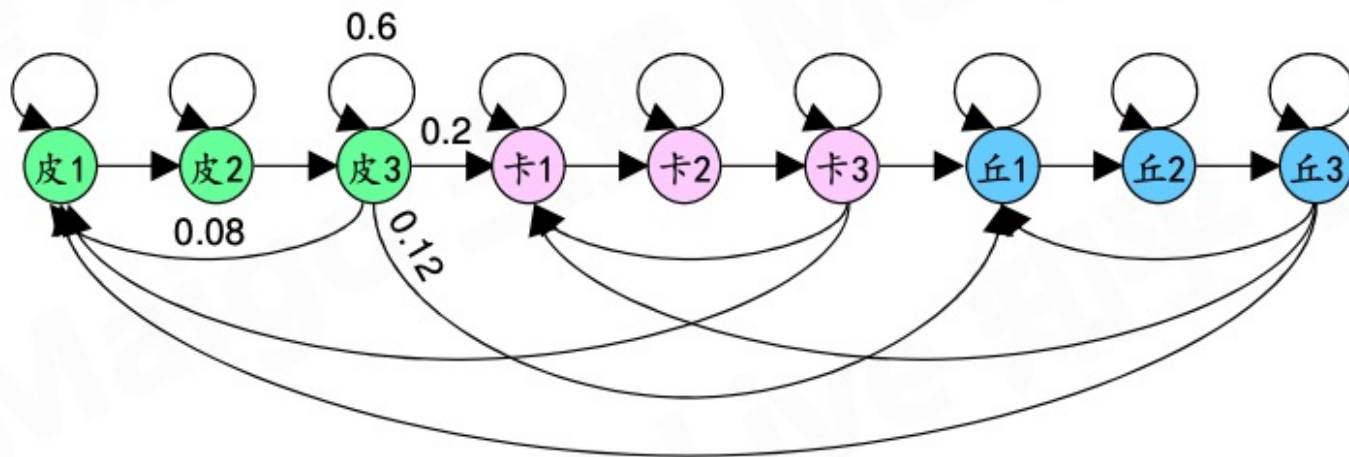
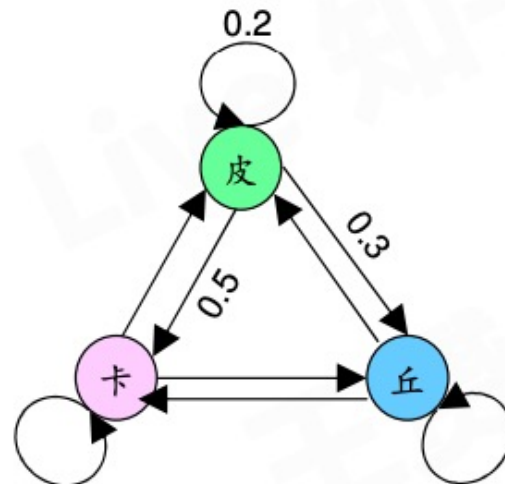
怎样在模式识别中发挥更大作用

□ Bigram是马尔可夫模型

- 下一个词只与当前词有关
- 模型是遍历的，不是单向的

□ 可与单词的声学模型复合

- 得到一门语言的HMM

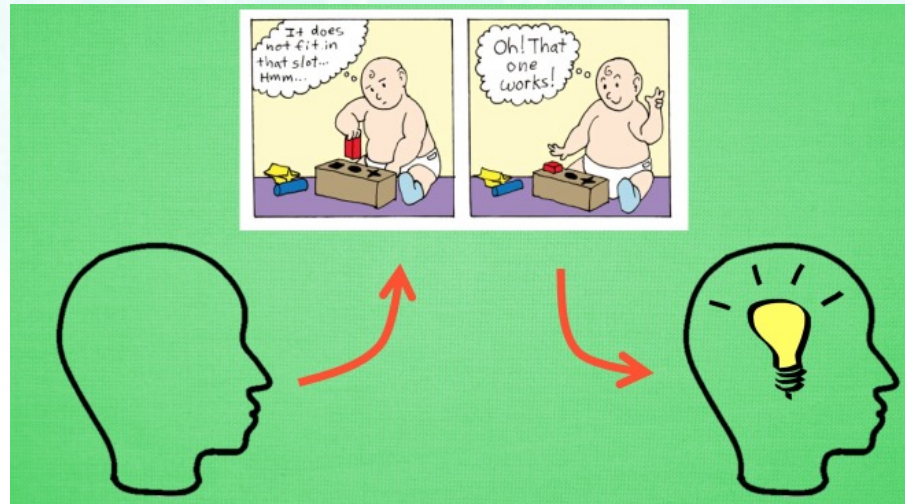




强化学习模型与马氏链

强化学习概念

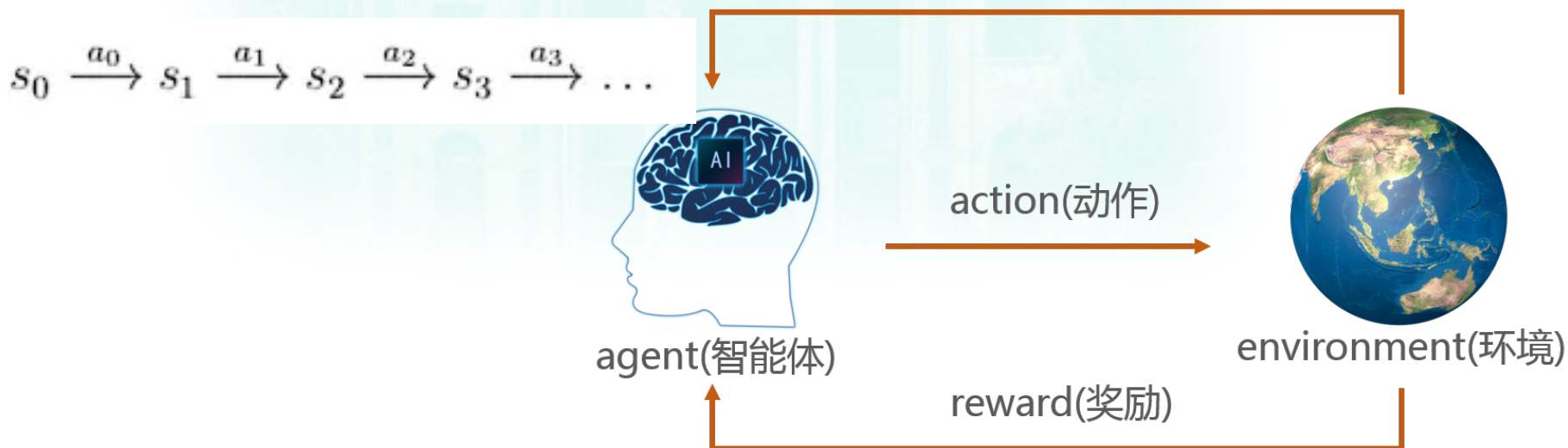
- ❑ **Reinforcement learning (RL)** is an area of machine learning inspired by *behaviorist psychology*, concerned with how software *agents* ought to take *actions* in an *environment* so as to maximize some notion of cumulative *reward*. ——[维基百科](#)



强化学习基础

- 基于马尔可夫决策——MDP $\langle S, A, R, P, \gamma \rangle$
 - S 状态集, A 动作集, R 奖励函数, P 状态转换概率, γ 折扣因子
 - ✓ 目标: Agent 最大化累积奖励 G (Cumulative Reward) 期望 $E[G]$

$$G = R(s_0, a_0) + \gamma R(s_1, a_1) + \gamma^2 R(s_2, a_2) + \dots_{\text{state(状态)}}$$



强化学习基础

- 策略 (policy) 是指任意一个从状态到行动的映射函数 $\pi: S \rightarrow A$
 - 价值函数: $V^\pi(s) = E[G | s_0 = s, \pi]$
 - 贝尔曼方程: $V^\pi(s) = R(s, a) + \gamma \sum_{s' \in S} P_{s\pi(s)}(s') V^\pi(s')$
 - ✓ 价值迭代
 - 最优价值函数: $V^*(s) = \max_{\pi} V^\pi(s)$
 - 策略迭代
 - 最优策略 $\pi^*(s) = \arg \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V^*(s')$
 - 对于所有状态 s 与策略 π , 有: $V^*(s) = V^{\pi^*}(s) \geq V^\pi(s)$.
- 模型 (model)
 - 预测 agent 行动后的状态

Diffusion模型概述

-从Markov链过渡到Diffusion模型

前言-对抗式生成网络 (GAN)

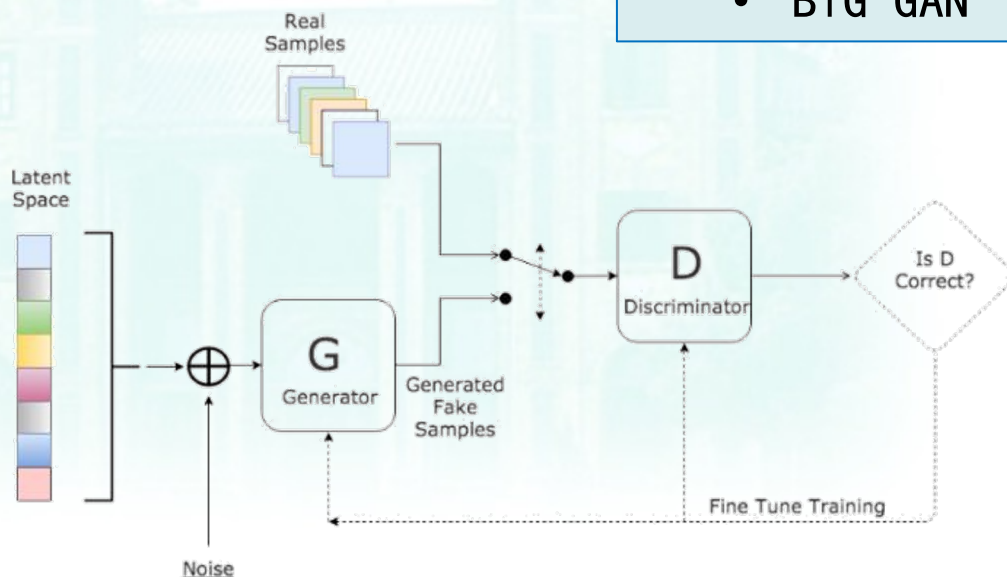
- **起源:** 生成模型缺乏固定的指标进行训练
- **动机:** 零和博弈与纳什均衡
- **方法:** 生成-对抗-再生成
- **结构:** 生成器与判别器
- **缺陷:** 不够稳定

应用:

- 图像生成
- 超分辨率
- 音频

后续改进:

- Conditional GANs
- DC GAN
- BiG GAN

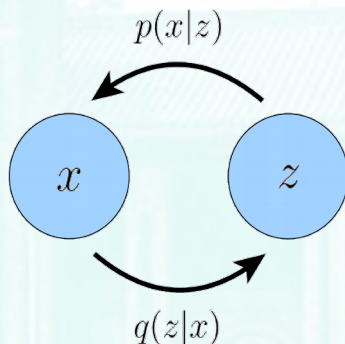
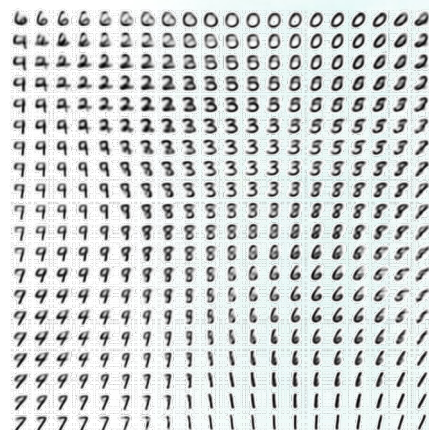


损失函数不是纯粹的最小化

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

前言-变分自编码器VAE

- 起源：数据的真实分布很难计算
- 动机：贝叶斯进行后验近似估计
- 方法：重建图像并优化ELBO
- 结构：编码器与解码器
- 缺陷：过于连续导致图像边界模糊



$$\begin{aligned}\log p(x) &= \log \int p(x, z) dz \\ &= \log \int \frac{p(x, z) q_\phi(z|x)}{q_\phi(z|x)} dz \\ &= \log \mathbb{E}_{q_\phi(z|x)} \left[\frac{p(x, z)}{q_\phi(z|x)} \right] \\ &\geq \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p(x, z)}{q_\phi(z|x)} \right] \\ &= \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p_\theta(x|z)p(z)}{q_\phi(z|x)} \right] \\ &= \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)] + \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p(z)}{q_\phi(z|x)} \right] \\ &= \underbrace{\mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_\phi(z|x) \parallel p(z))}_{\text{prior matching term}}\end{aligned}$$

重建图像与原始图像的差距

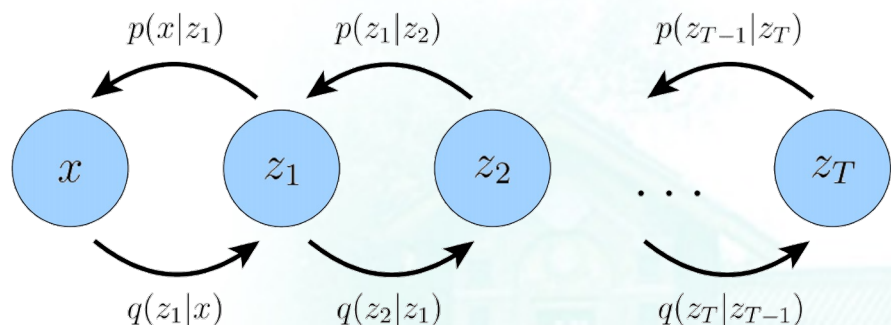
学习得到的分布q与先验假设p的相似

这导致了VAE的Loss中KL散度一项拥有解析解，而重构项只需要MSE就可以计算！

与GAN的区别：GAN中，数据的**真实分布未知**；而在VAE中，假设了数据服从Gaussian分布，**只是方差均值未知**；

前言-Markov假设的HMVAE

VAE难以生成精致图片->用很多个VAE迭代生成可以解决。



x 是原始数据, z_1 是第一层隐空间表征, \dots , z_T 是最终的高斯分布

$$\mathbb{E}_{q_\phi(z_{1:T}|\mathbf{x})} \left[\log \frac{p(z_T)p_\theta(\mathbf{x}|z_1) \prod_{t=2}^T p_\theta(z_{t-1}|z_t)}{q_\phi(z_1|\mathbf{x}) \prod_{t=2}^T q_\phi(z_t|z_{t-1})} \right]$$

为什么是Markov?

我需要这样的链来生成精致图片, 但是我不想要解析解过于复杂

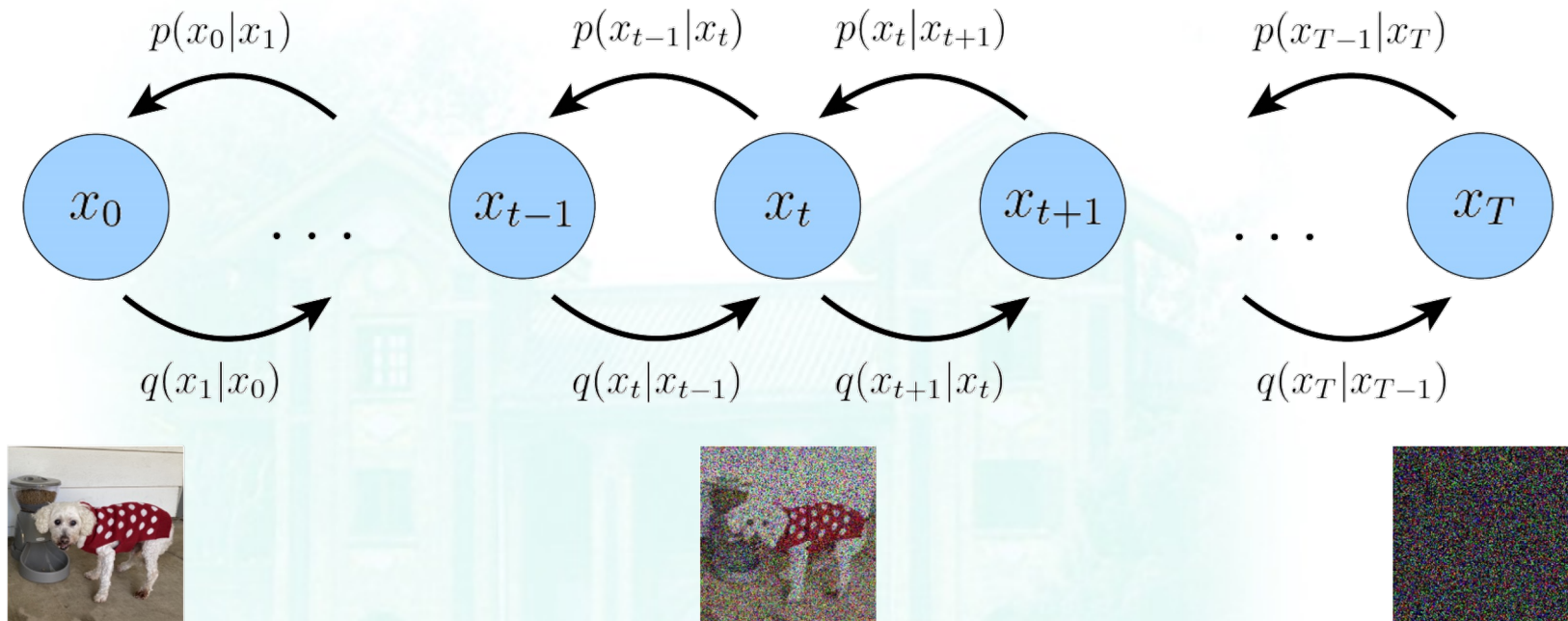
仍然存在的问题?

虽然用Markov链的性质消除了许多过长的条件概率, 但是因为先验扰动 q 仍然是含参的, 进行优化仍然复杂。

➡ 那么能否消除先验扰动 q 的参数呢?

扩散概率模型 Diffusion Probabilistic Model

答案：可以！



将HMVAE中，先验过程分布的方差和均值也设定为已知的，通常是一个高斯分布

GAN（分布未知，均值与方差未知） \Rightarrow VAE（先验分布已知，均值方差未知）
 \Rightarrow Diffusion（可解析的先验分布）

扩散概率模型 Diffusion Probabilistic Model

具体怎么优化？

DPM的对数似然为：

$$\log p(\mathbf{x}) = \log \int p(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T}$$

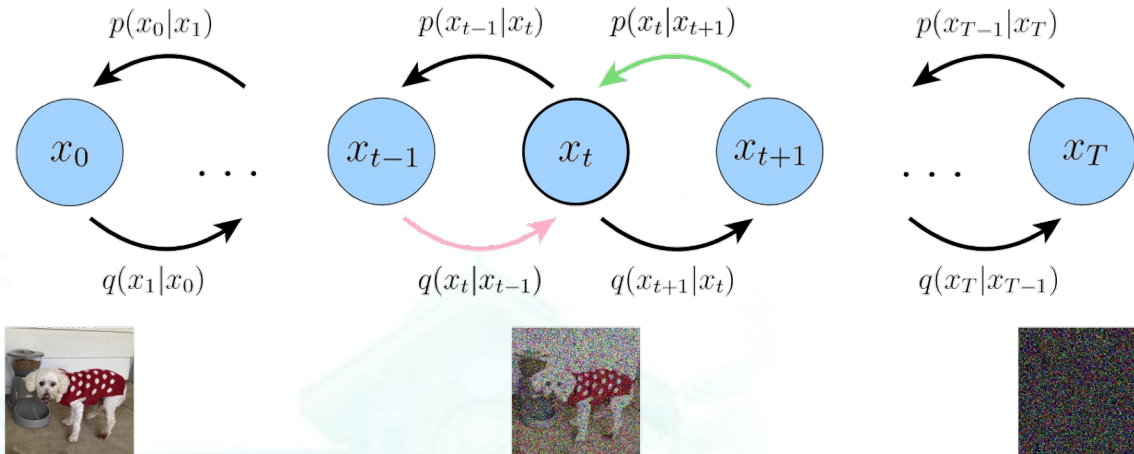
$$\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right]$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right]$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\sum_{t=1}^{T-1} \log \frac{p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right]$$

$$= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{\mathbb{E}_{q(\mathbf{x}_{T-1}|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_{T-1}) \| p(\mathbf{x}_T))]}_{\text{prior matching term}}$$

$$- \sum_{t=1}^{T-1} \underbrace{\mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_{t+1}|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_t|\mathbf{x}_{t-1}) \| p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1}))]}_{\text{consistency term}}$$



正向过程 $q(\mathbf{x}_t|\mathbf{x}_{t-1}) \Rightarrow$ 对前一步的结果加噪声
 逆向过程 $p(\mathbf{x}_t|\mathbf{x}_{t+1}) \Rightarrow$ 根据已加噪声进行恢复

怎么解决？

问题：虽然ELBO是可解析的形式，然而迭代项存在 $\mathbf{x}_{t-1}\mathbf{x}_t\mathbf{x}_{t+1}$ ，
 这会导致次优的解，这不是一个好的做法

扩散概率模型 Diffusion Probabilistic Model

贝叶斯！ 代入原始的对数似然中

可得：

$$\begin{aligned}
 \log p(\mathbf{x}) &\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p_\theta(\mathbf{x}_T)p_\theta(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)p_\theta(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)p_\theta(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} \right] + \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t, \mathbf{x}_{t-1}|\mathbf{x}_0)} \left[\log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \\
 &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}
 \end{aligned}$$

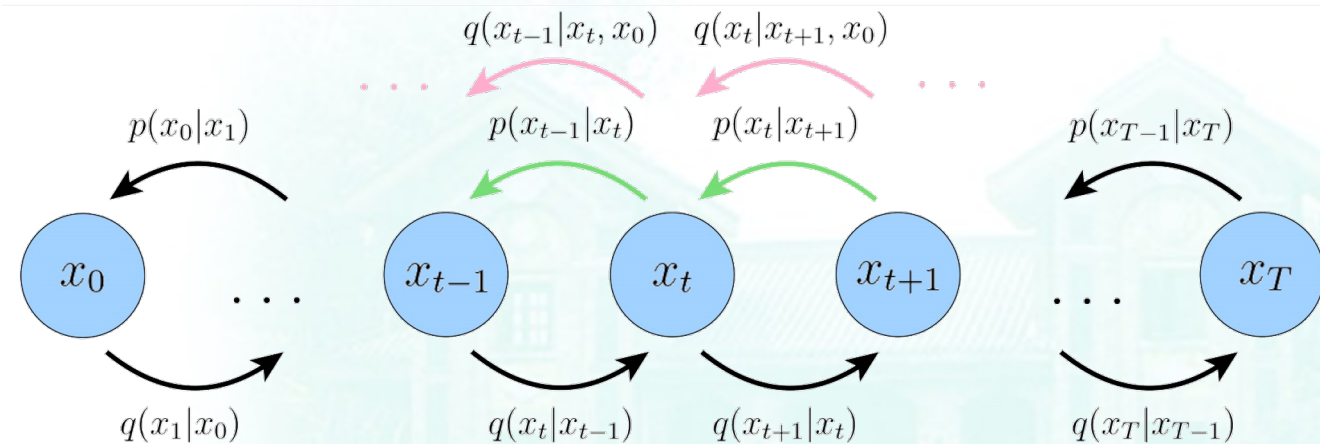
$$q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) = \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}$$

在原始ELBO的优化推导中引入贝叶斯后，消除了迭代项中的 $\mathbf{x}_{t-1}\mathbf{x}_t\mathbf{x}_{t+1}$ 之间的计算

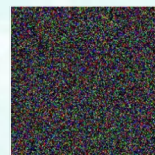
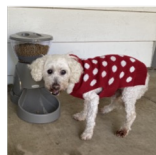
$$\begin{aligned}
 &\log \prod_{t=2}^T \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}} = \log \prod_{t=2}^T \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} * \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} \\
 &= \log \prod_{t=2}^T \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} + \log \prod_{t=2}^T \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} \\
 &= \log \prod_{t=2}^T \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} + \log \frac{\prod_{t=2}^T q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{\prod_{t=2}^T q(\mathbf{x}_t|\mathbf{x}_0)} \\
 &= \log \prod_{t=2}^T \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \\
 &+ \log \frac{\cancel{q(\mathbf{x}_{T-1}|\mathbf{x}_0)} * \cancel{q(\mathbf{x}_{T-2}|\mathbf{x}_0)} * \cancel{q(\mathbf{x}_{T-3}|\mathbf{x}_0)} * \dots * \cancel{q(\mathbf{x}_3|\mathbf{x}_0)} * q(\mathbf{x}_2|\mathbf{x}_0)}{\cancel{q(\mathbf{x}_T|\mathbf{x}_0)} * \cancel{q(\mathbf{x}_{T-1}|\mathbf{x}_0)} * \cancel{q(\mathbf{x}_{T-2}|\mathbf{x}_0)} * \dots * \cancel{q(\mathbf{x}_3|\mathbf{x}_0)} * \cancel{q(\mathbf{x}_2|\mathbf{x}_0)}} \\
 &= \log \prod_{t=2}^T \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{q(\mathbf{x}_T|\mathbf{x}_0)}
 \end{aligned}$$

扩散概率模型 Diffusion Probabilistic Model

直观上的改变 \Rightarrow Loss可以被计算，并有**最优解**



注意：目前，仍然仅需要一个Markov假设就能使如此复杂的过程成立。



$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

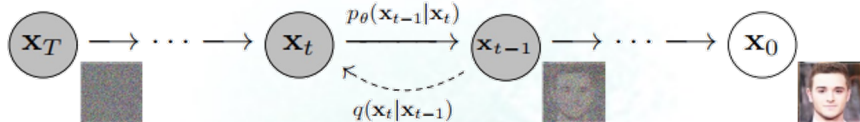
\Rightarrow 我们需要求出来，得到精确的Loss解析形式

扩散概率模型 Diffusion Probabilistic Model

$$\begin{aligned}
 q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) &= \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} \\
 &= \frac{\mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1-\alpha_t)\mathbf{I})\mathcal{N}(\mathbf{x}_{t-1}; \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0, (1-\bar{\alpha}_{t-1})\mathbf{I})}{\mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1-\bar{\alpha}_t)\mathbf{I})} \\
 &\propto \exp \left\{ - \left[\frac{(\mathbf{x}_t - \sqrt{\alpha_t}\mathbf{x}_{t-1})^2}{2(1-\alpha_t)} + \frac{(\mathbf{x}_{t-1} - \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0)^2}{2(1-\bar{\alpha}_{t-1})} - \frac{(\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0)^2}{2(1-\bar{\alpha}_t)} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left[\frac{(\mathbf{x}_t - \sqrt{\alpha_t}\mathbf{x}_{t-1})^2}{1-\alpha_t} + \frac{(\mathbf{x}_{t-1} - \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0)^2}{1-\bar{\alpha}_{t-1}} - \frac{(\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0)^2}{1-\bar{\alpha}_t} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left[\frac{(-2\sqrt{\alpha_t}\mathbf{x}_t\mathbf{x}_{t-1} + \alpha_t\mathbf{x}_{t-1}^2)}{1-\alpha_t} + \frac{(\mathbf{x}_{t-1}^2 - 2\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_{t-1}\mathbf{x}_0)}{1-\bar{\alpha}_{t-1}} + C(\mathbf{x}_t, \mathbf{x}_0) \right] \right\} \\
 &\propto \exp \left\{ - \frac{1}{2} \left[- \frac{2\sqrt{\alpha_t}\mathbf{x}_t\mathbf{x}_{t-1}}{1-\alpha_t} + \frac{\alpha_t\mathbf{x}_{t-1}^2}{1-\alpha_t} + \frac{\mathbf{x}_{t-1}^2}{1-\bar{\alpha}_{t-1}} - \frac{2\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_{t-1}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left[\left(\frac{\alpha_t}{1-\alpha_t} + \frac{1}{1-\bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}\mathbf{x}_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left[\frac{\alpha_t(1-\bar{\alpha}_{t-1}) + 1-\alpha_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \mathbf{x}_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}\mathbf{x}_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left[\frac{\alpha_t - \bar{\alpha}_t + 1 - \alpha_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \mathbf{x}_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}\mathbf{x}_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left[\frac{1-\bar{\alpha}_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \mathbf{x}_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}\mathbf{x}_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left(\frac{1-\bar{\alpha}_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \right) \left[\mathbf{x}_{t-1}^2 - 2 \frac{\left(\frac{\sqrt{\alpha_t}\mathbf{x}_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right)}{\frac{1-\bar{\alpha}_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}} \mathbf{x}_{t-1} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left(\frac{1-\bar{\alpha}_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \right) \left[\mathbf{x}_{t-1}^2 - 2 \frac{\left(\frac{\sqrt{\alpha_t}\mathbf{x}_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right) (1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} \mathbf{x}_{t-1} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left(\frac{1}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \right) \left[\mathbf{x}_{t-1}^2 - 2 \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)\mathbf{x}_0}{1-\bar{\alpha}_t} \mathbf{x}_{t-1} \right] \right\} \\
 &\propto \mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)\mathbf{x}_0}{1-\bar{\alpha}_t}}_{\mu_q(\mathbf{x}_t, \mathbf{x}_0)}, \underbrace{\frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}}_{\Sigma_q(t)} \mathbf{I})
 \end{aligned}$$

扩散概率模型 Diffusion Probabilistic Model

模型训练 尽管DPM的原理已经在2015年就给出了完整的数学证明，但直到2020年，才广泛进入到研究者的视野，即：Denoising Diffusion Probabilistic Models



1. 说明任意时刻的数据 x_t 服从于高斯分布
2. 说明Loss中的迭代项是高斯分布

$$\begin{aligned}
 x_t &= \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\epsilon_{t-1}^* \\
 &= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}}x_{t-2} + \sqrt{1 - \alpha_{t-1}}\epsilon_{t-2}^* \right) + \sqrt{1 - \alpha_t}\epsilon_{t-1}^* \\
 &= \sqrt{\alpha_t\alpha_{t-1}}x_{t-2} + \sqrt{\alpha_t - \alpha_t\alpha_{t-1}}\epsilon_{t-2}^* + \sqrt{1 - \alpha_t}\epsilon_{t-1}^* \\
 &= \sqrt{\alpha_t\alpha_{t-1}}x_{t-2} + \sqrt{\alpha_t - \alpha_t\alpha_{t-1} + \sqrt{1 - \alpha_t}^2}\epsilon_{t-2}^* \\
 &= \sqrt{\alpha_t\alpha_{t-1}}x_{t-2} + \sqrt{\alpha_t - \alpha_t\alpha_{t-1} + 1 - \alpha_t}\epsilon_{t-2}^* \\
 &= \sqrt{\alpha_t\alpha_{t-1}}x_{t-2} + \sqrt{1 - \alpha_t\alpha_{t-1}}\epsilon_{t-2}^* \\
 &= \dots \\
 &= \sqrt{\prod_{i=1}^t \alpha_i}x_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i}\epsilon_0 \\
 &= \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_0 \\
 &\sim \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I})
 \end{aligned}$$

$$\begin{aligned}
 q(x_{t-1}|x_t, x_0) &= \frac{q(x_t|x_{t-1}, x_0)q(x_{t-1}|x_0)}{q(x_t|x_0)} \\
 &= \frac{\mathcal{N}(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)\mathbf{I})\mathcal{N}(x_{t-1}; \sqrt{\bar{\alpha}_{t-1}}x_0, (1 - \bar{\alpha}_{t-1})\mathbf{I})}{\mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I})} \\
 &\propto \exp \left\{ - \left[\frac{(x_t - \sqrt{\alpha_t}x_{t-1})^2}{2(1 - \alpha_t)} + \frac{(x_{t-1} - \sqrt{\bar{\alpha}_{t-1}}x_0)^2}{2(1 - \bar{\alpha}_{t-1})} - \frac{(x_t - \sqrt{\bar{\alpha}_t}x_0)^2}{2(1 - \bar{\alpha}_t)} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left[\frac{(x_t - \sqrt{\alpha_t}x_{t-1})^2}{1 - \alpha_t} + \frac{(x_{t-1} - \sqrt{\bar{\alpha}_{t-1}}x_0)^2}{1 - \bar{\alpha}_{t-1}} - \frac{(x_t - \sqrt{\bar{\alpha}_t}x_0)^2}{1 - \bar{\alpha}_t} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left[\frac{(-2\sqrt{\alpha_t}x_t x_{t-1} + \alpha_t x_{t-1}^2)}{1 - \alpha_t} + \frac{(x_{t-1}^2 - 2\sqrt{\bar{\alpha}_{t-1}}x_{t-1}x_0)}{1 - \bar{\alpha}_{t-1}} + C(x_t, x_0) \right] \right\} \\
 &\propto \exp \left\{ - \frac{1}{2} \left[\frac{2\sqrt{\alpha_t}x_t x_{t-1} + \alpha_t x_{t-1}^2}{1 - \alpha_t} + \frac{x_{t-1}^2}{1 - \bar{\alpha}_{t-1}} - \frac{2\sqrt{\bar{\alpha}_{t-1}}x_{t-1}x_0}{1 - \bar{\alpha}_{t-1}} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left[\left(\frac{\alpha_t}{1 - \alpha_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) x_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}x_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}x_0}{1 - \bar{\alpha}_{t-1}} \right) x_{t-1} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left[\frac{\alpha_t(1 - \bar{\alpha}_{t-1}) + 1 - \alpha_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} x_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}x_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}x_0}{1 - \bar{\alpha}_{t-1}} \right) x_{t-1} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left[\frac{\alpha_t - \bar{\alpha}_t + 1 - \alpha_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} x_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}x_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}x_0}{1 - \bar{\alpha}_{t-1}} \right) x_{t-1} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left[\frac{1 - \bar{\alpha}_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} x_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}x_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}x_0}{1 - \bar{\alpha}_{t-1}} \right) x_{t-1} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left(\frac{1 - \bar{\alpha}_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \right) \left[x_{t-1}^2 - 2 \frac{\left(\frac{\sqrt{\alpha_t}x_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}x_0}{1 - \bar{\alpha}_{t-1}} \right)}{\frac{1 - \bar{\alpha}_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}} x_{t-1} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left(\frac{1 - \bar{\alpha}_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \right) \left[x_{t-1}^2 - 2 \frac{\left(\frac{\sqrt{\alpha_t}x_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}x_0}{1 - \bar{\alpha}_{t-1}} \right) (1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_{t-1} \right] \right\} \\
 &= \exp \left\{ - \frac{1}{2} \left(\frac{1}{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}} \right) \left[x_{t-1}^2 - 2 \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t} x_{t-1} \right] \right\} \\
 &\propto \mathcal{N}(x_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t}}_{\mu_q(x_t, x_0)}, \underbrace{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})\mathbf{I}}_{\Sigma_q(t)})
 \end{aligned}$$

扩散概率模型 Diffusion Probabilistic Model

3. 根据两个高斯分布的KL可解，求出具体的解析形式

$$L_{t-1} = \mathbb{E}_q \left[\frac{1}{2\sigma_t^2} \|\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) - \mu_\theta(\mathbf{x}_t, t)\|^2 \right] + C$$

$$D_{\text{KL}}(\mathcal{N}(\mathbf{x}; \mu_x, \Sigma_x) \parallel \mathcal{N}(\mathbf{y}; \mu_y, \Sigma_y)) = \frac{1}{2} \left[\log \frac{|\Sigma_y|}{|\Sigma_x|} - d + \text{tr}(\Sigma_y^{-1} \Sigma_x) + (\mu_y - \mu_x)^T \Sigma_y^{-1} (\mu_y - \mu_x) \right]$$

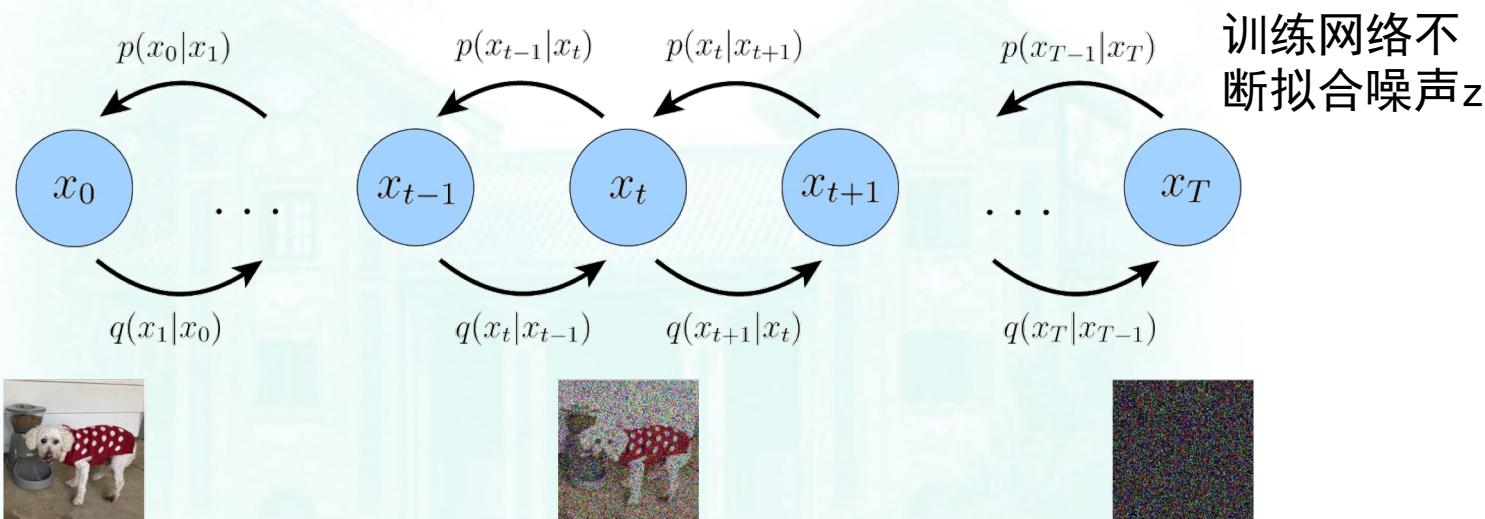
$$\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \mathbf{x}_0$$

配完化简后的结果，但是我哪有X0啊现在

之前咱们说Xt可以由X0计算得到，现在逆一下 $\hat{\mathbf{x}}_0 = \frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \mathbf{z}_t)$

$$\tilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \mathbf{z}_t \right)$$

最终结果



4. 求和，简化后=>

$$L_{\text{simple}}(\theta) := \mathbb{E}_{t, \mathbf{x}_0, \epsilon} \left[\|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2 \right]$$

扩散概率模型 Diffusion Probabilistic Model

Algorithm 1 Training

- 1: **repeat**
 - 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
 - 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
 - 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 5: Take gradient descent step on
$$\nabla_{\theta} \left\| \epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t) \right\|^2$$
 - 6: **until** converged
-

Algorithm 2 Sampling

- 1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 2: **for** $t = T, \dots, 1$ **do**
 - 3: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$
 - 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
 - 5: **end for**
 - 6: **return** \mathbf{x}_0
-

$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \epsilon_{t-1}^*$$

Jonathan Ho et al. Denoising Diffusion Probabilistic Models. arXiv:2006.11239

扩散概率模型Diffusion Probabilistic Model

Stable Diffusion应用示例

Prompt: A silver mech horse running in a dark valley, in the night, Beppe, Kaino University, high-definition picture, unreal engine, cyberpunk

应用场景

- 图像生成
- 多模态
- 语音、文本等领域
- 例如：DALI2、Imagen、GLIDE



缺陷与其他改进

- 需要迭代采样，速度相对来说较慢；目前已经有了不少对应的改进方法，例如DDIM和ANALYTIC-DPM等。
- 需要整个过程中，输入输出的大小保持相同维度；目前还没有比较好的方法。
- 加噪过程的最终分布是假设的标准高斯噪声，在风格转换之类的条件生成中受到挑战；已有了不少改进，例如Cold Diffusion已经取消了标准高斯的假设。
- 因为前两者的限制导致计算量过大，例如Stable Diffusion据说使用了4000台 A100 显卡集群，用了一个月时间。