

Friston等人的最小自由能原理的失误和改进

——两种变分自由能与信息效率

关键词：变分贝叶斯，变分自由能，最小自由能原理，香农互信息，语义互信息，信息效率

本PPT的PDF版下载地址：<http://www.survivor99.com/lcg/fep-error.pdf>

鲁晨光

2026.5

个人网站：<http://www.survivor99.com>

背景1: 变分贝叶斯方法(VB)

- 贝叶斯方法——贝叶斯主义推断方法，强调使用先验分布 $P(z)$
- VB就是变分贝叶斯推断方法。
- https://en.wikipedia.org/wiki/Variational_Bayesian_methods
- Hinton等人: 提出VB用于无监督学习; 优化准则: 最小变分自由能(VFE)准则
- Hinton开始提出赫尔姆霍茨机, 后来又提出有限波尔茨曼机(RBM), 都用这一准则。因为自由能是物理学概念, Hinton获得物理学诺贝尔奖。
- 代表文章, 著作:
- Hinton, G.E.; van Camp, D. [Keeping the neural networks simple by minimizing the description length of the weights.](#) 通过最小化权重编码长度保持神经网络简单性
- Hinton, G.E.; Zemel, R.S. [Autoencoders, minimum description length and Helmholtz free energy.](#) 变分编码器, 最小编码长度和Helmholtz自由能
- Neal, R.; Hinton, G. [A view of the EM algorithm that justifies incremental, sparse, and other variants.](#) 关于EM算法的一个观点...
- 变分自由能(VFE)是什么? 权重长度呈概率分布是 $P(x)$, 如果按预测的概率分布 $P_\theta(x)$ 编码, x_i 的码长是 $-\log P_\theta(x_i)$, 实际编码长度接近交叉熵, 大于或等于香农熵:

$$H_\theta(X) = -\sum_i P(x_i) \log P_\theta(x) \geq H(X)$$

- VB任务: 怎样让交叉熵接近香农熵?
- 反映香农无记忆离散信源无失真编码定理。这是信息论方法, 本来和物理学没啥关系。

背景2：从VB到最小自由能原理（FEP）

- **Friston**研究脑神经科学，把**VB**应用到脑神经和行为科学
- 发展出自小自由能原理（**FEP**）——对标**Jaynes**的最大熵原理，
- 核心思想：人脑和环境相互适应，实现自己的目标；而不仅仅是被动服从和适应约束。
- 最近用于解释生命现象和自组织问题。
- Friston等人代表文章和专著：
 - [Friston, K. The free-energy principle: A unified brain theory? *Nat. Rev. Neurosci.* 2010](#)
 - [Parr, T.; Pezzulo, G.; Friston, K.J. *Active Inference: The Free Energy Principle in Mind, Brain, and Behavior*; MIT Press: Cambridge, MA, USA, 2022.](#)
- 物理学界反对的比较多（主要因为质疑：1，自由能概念；2.作为一般原理）
- 我认为他们的方向是对的——因为最大熵理论要发展，但是数学上有问题
- 最近Friston把FEP用于主动推断Active Inference
- 现在强化主动推断，弱化FEP
- VB、FEP，变分自动编码器VAE，和主动推断有广泛影响
- 核心问题：最小化VFE包含怎样的优化思想？
- 到目前为止还不清楚，VB发明者也没讲清，看网上介绍很难理解

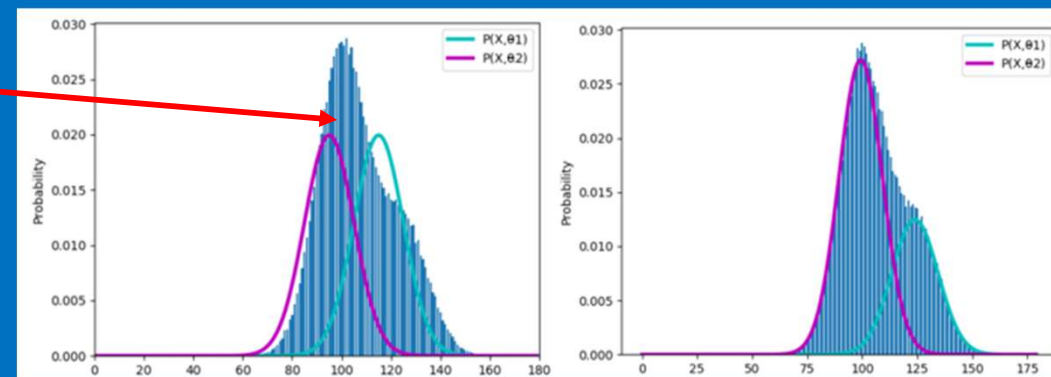
理解变分自由能 (VFE) —— EM算法用于混合模型

- Neal和Hinton, Beal...很多人都用VB解决混合模型问题——无监督学习。
- 样本分布来自真模型:

- $P(x) = P^*(z_1)P(x|z_1^*) + P^*(z_2)P(x|z_2^*)$.
- 混合模型来自模型预测:
- $P_\theta(x) = P(z_1)P(x|z_1, \theta) + P(z_2)P(x|z_2, \theta)$.

开始猜不准

收敛后



- 怎样让 $H_\theta(x) \rightarrow H(x)$?
- 已知 $P(x)$ 和模型类别 (比如高斯), 求解 $P(z)$, 和 $P(x|\theta_j) = P(x|z_j, \theta)$

EM算法是最小化联合交叉熵 $H_\theta(X, Z)$, 采用变分方法

它先假设一个分布 $P(z)$, 和似然函数 $P(x|z_j, \theta)$

用E步 (算期望) 得到 $q^{+1}(z_j | x_i) \leq q(z_j)P(x_i | \theta_j) / Z_i, Z_i = \sum_j q(z_j)[P(x_i | \theta_j)],$

再用M1步: $q^{+1}(z_j) \leq \sum_i P(x_i)q^{+1}(z_j | x_i),$

M2步: $P(x|z_j, \theta) \leq P(x|z_j^*),$ 对于高斯混合, 求期望和标准偏差

重复就能让 $H_\theta(x) \rightarrow H(x)$ 。

任务类似, VB解法有些不同。通常认为VB是比EM算法更一般的方法。

物理学自由能和Boltzmann分布

- 物理学中，赫尔姆霍茨自由能： $F = E(e_i, p_i) - TS(p_i)$

- Hinton等人文章中：

$$F = \sum_i p_i E_i - H$$

能量 熵

$$H = - \sum_i p_i \log p_i$$

- $H = kT \ln Z(x)$
- p_i 是变分，令 $\partial F / \partial S(p_i) = 0$ ，可以求出波尔茨曼分布：

$$p_i = \frac{e^{-E_i}}{\sum_j e^{-E_j}}$$

为什么叫变分自由能 (VFE)

- 物理学中，赫尔姆霍茨自由能： $F = E(e_i, p_i) - TS(p_i)$

- Hinton本意：

$$P(x, z) = P_\theta(x, z)$$

$$= P(z)P(x | z)$$

$$VFE = E_{q(z|x)} \log \frac{q(z|x)}{P(z|x)} + H_\theta(x)$$

$$- E_{q(z|x)} \log \frac{q(z|x)}{P(x, z)} = E_{q(z|x)} \log \frac{1}{P(x, z)} - H(z|x)$$

能量

熵

- 为什么不直接最小化 $H_\theta(x)$ ，因为其中没有 $q(z|x)$
- 为什么常见的变分自由能公式是这样的 →

Wiki Pedia Variational Bayesian Method中，VFE:

为了更像物理学自由能公式，

所以用 $q(z)$ 代替 $q(z|x)$ ，也带来问题。

$$F = E_{q(z)} \log \frac{q(z)}{P(x, z)}$$

- 类似地，令 $\partial VFE / \partial q(x) = 0$ ，求出 $q(z)$ 。

VB使用的 $q(z)$ 是后验分布 $q(z|x)$ ，证据

- 1. 你可以问大模型，它们都说是；
- 2. 在Wikipedia, 词条Variational Bayesian Mehtods中有：

In **variational** inference, the **posterior distribution** over a set of unobserved variables $\mathbf{Z} = \{Z_1 \dots Z_n\}$ given \mathbf{X} is approximated by a so-called **variational distribution**, $Q(\mathbf{Z})$:

$$P(\mathbf{Z} | \mathbf{X}) \approx Q(\mathbf{Z}).$$

$$Q(\mathbf{Z}) = \prod_{i=1}^M q_i(\mathbf{Z}_i | \mathbf{X})$$

$$q_j^*(\mathbf{Z}_j | \mathbf{X}) = \frac{e^{\mathbb{E}_{q_{-j}^*}[\ln p(\mathbf{Z}, \mathbf{X})]}}{\int e^{\mathbb{E}_{q_{-j}^*}[\ln p(\mathbf{Z}, \mathbf{X})]} d\mathbf{Z}_j}$$

- 3. Hinton等人文章中：

$$p_i = \frac{e^{-E_i}}{\sum_j e^{-E_j}} \quad (5)$$

This is exactly the posterior probability distribution obtained when fitting a mixture of Gaussians to an input vector.

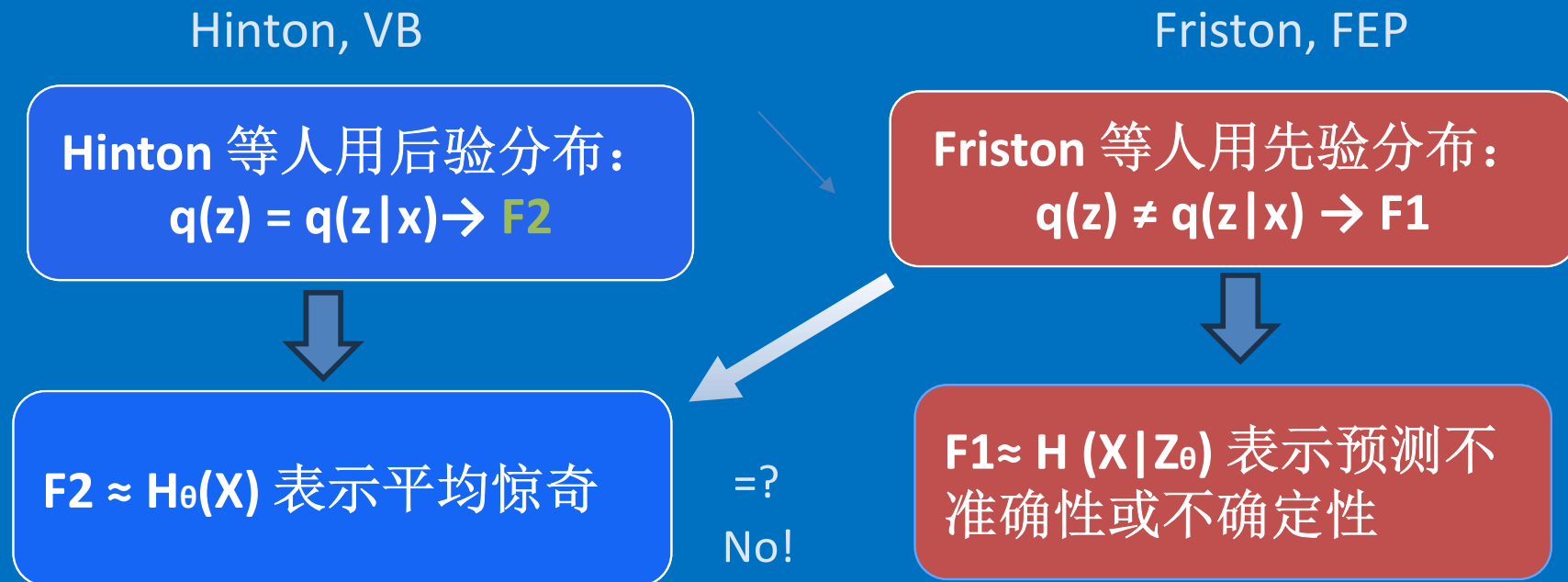
the single Gaussian **posterior**, Q , for a noisy weight is

$$G(P, Q) = \int Q(w) \log \frac{Q(w)}{\sum_i \pi_i P_i(w)} dw \quad (16)$$

- VB难懂，这是主要原因！

Friston等人的一个错误结论

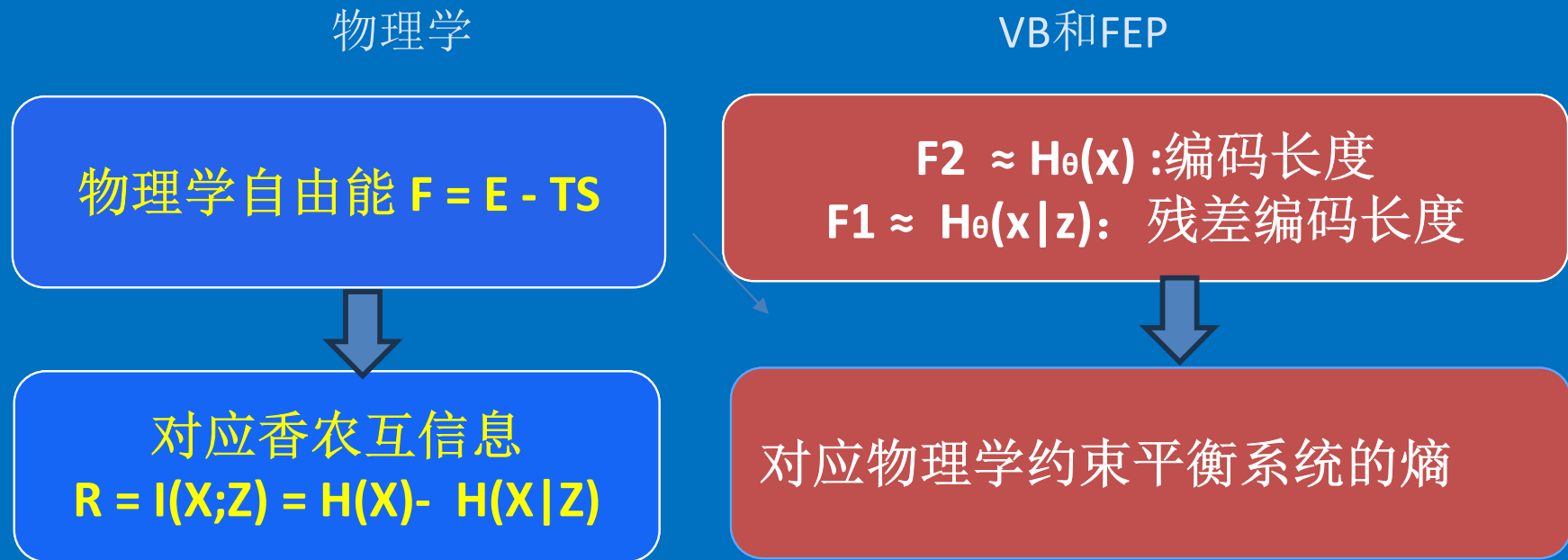
- 期望变分自由能 (VFE) 既表示平均惊奇, 又表示预测不准确性。
- 但平均惊奇 $H_{\theta}(X)$ 与预测不准确性 $H_{\theta}(X|Z) = H(X|Z_{\theta})$ 通常差异很大。
 $H_{\theta}(X) \gg H_{\theta}(X|Z)$
- 两种期望变分自由能来自两种变分 $q(z), q(z|x)$



$$F_2 = D_{KL}(q(x|z) || P(z)) + H_{\theta}(X)(z) \approx H_{\theta}(X) \gg F_1 = D_{KL}(q(z) || P(z)) + H(X|Z_{\theta}) \approx H(X|Z_{\theta})$$

物理学自由能 vs 变分自由能

- 我发表于Entropy 《Improving the Minimum Free Energy Principle to the Maximum Information Efficiency Principle》中讨论过变分自由能概念的问题。



- 通信中类似的是信息，两者都是好东西，通常增大好；但是用起来要节省。
- 物理学中熵增大，自由能减小，这是被动减小；而VB主动减小VFE。

两种自由能之间关系

- F1~条件交叉熵 $H_{\theta}(X|Z)$

- F2~交叉熵 $H_{\theta}(X)$

F1
预测不确定性
负的效用

+ R
香农互信息
成本

= F2
总目标: 追求效率

- 最小变分自由能准则包含追求最大信息效率思想!

- 两者关系:

$$F_2 = F_1 + R$$

- 推导公式

香农互信息R

$$\begin{aligned}
 F_2 &= E_{q(x)} E_{q(z|x)} \log \frac{q(z|x)}{P(x,z)} \\
 &= E_{q(x,z)} \log \frac{q(z|x)}{P(z)} - E_{q(x,z)} \log P_{\theta}(x|z) \\
 &= D_{KL}(q(z|x) \| q(z)) + \underbrace{D_{KL}(q(z) \| P(z)) - E_{q(x,z)} \log P_{\theta}(x|z)}_{F_1} \\
 &= R + F_1
 \end{aligned}$$

- VB问题: 怎么计算信息效率? 如何权衡最大语义信息和最大信息效率?

语义信息G理论

- 1993, 1997:



- 1999: A generalization of Shannon's information theory" , Int. J. of General Systems, 28: (6) 453-490 1999
- 2025: A Semantic Generalization of Shannon's Information Theory and Applications. *Entropy* 2025, 27, 461. <https://doi.org/10.3390/e27050461>
- 基本公式1: 真值函数(反映语义) 和似然函数之间的关系:
- 语义贝叶斯公式

$$P(x | \theta_j) = \frac{T(\theta_j | x)P(x)}{T(\theta_j)}, \quad T(\theta_j) = \sum_i T(\theta_j | x_i)P(x_i).$$

语义互信息公式的性质 及G和F1的关系

- 基本公式2：语义互信息公式：

- G和F1之间的关系：

$$G=I(X;Z_\theta)=\sum_j \sum_i P(x_i|z_j)P(z_j) \log \frac{P(x_i|\theta_j)}{P(x_i)} \approx H(X) - F_1$$

- G和平均失真之间的关系：

$$= \sum_j \sum_i P(x_i)P(z_j|x_i) \log \frac{T(\theta_j|x_i)}{T(\theta_j)} = H_\theta(Z) - \bar{d}$$

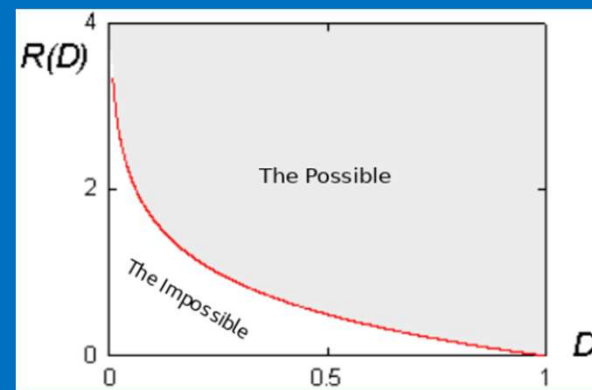
- 公式的性质：

- 1) 当语义信道匹配香农信道，即 $P(x|\theta_j)=P(x|z)$ 或 $T(\theta_j|x)$ 正比于 $q(z|x)$ 时，语义信息量最大达，达到其上界R。
- 2) 当香农信道匹配语义信道时，即 $P(x|z)=P(x|\theta_j)$ 或 $q(z|x)$ 正比于 $T(\theta_j|x)$ 时，信息差R-G最小，或信息效率G/R最大。

语义变分贝叶斯——来自R(D)函数推广

(8)

- 信息率失真函数R(D)推广到信息率逼真函数R(G)
- 目标函数: $R-sG$, $s=1$, 就是追求最大信息效率1
- 用变分方法得到优化香农信道迭代公式:



$$q^{+1}(z_j | x_i) = q(z_j)[P(x_i | \theta_j)]^s / Z_i, \quad Z_i = \sum_j q(z_j)[P(x_i | \theta_j)]^s,$$

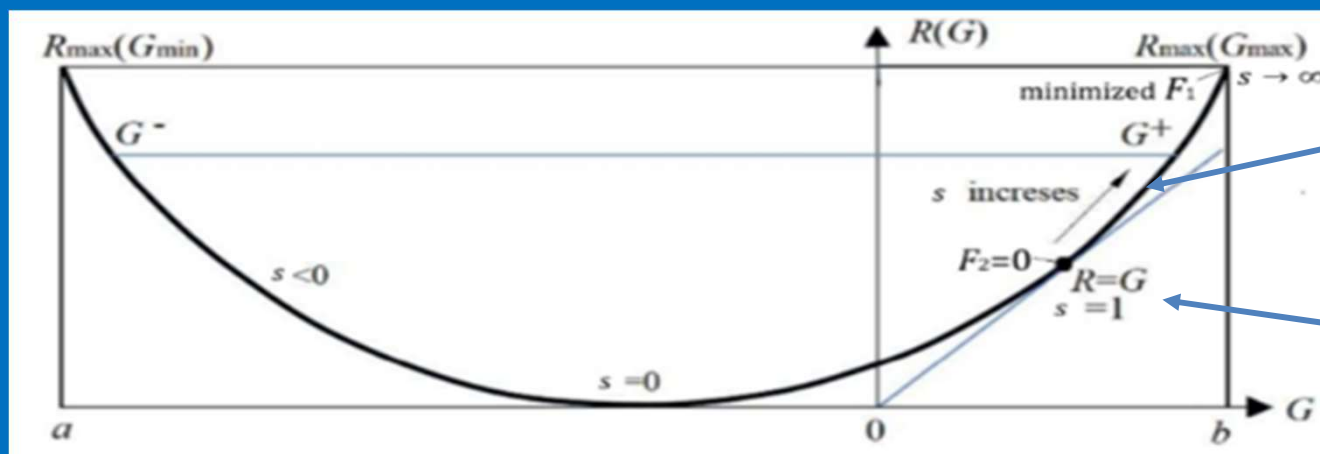
$$q^{+1}(z_j) = \sum_i P(x_i)q(z_j | x_i),$$

- 增大s, 可在最大语义信息和最大信息效率之间权衡

S=1, 最小化F2
最大信息效率

最小化R-sG
权衡

S → ∞, 最小化F1
最大语义信息



s从1增大, 牺牲效率G/R, 增加效用G

最小化F2, 使得R=G, s=1, G/R=1

F2 和香农互信息，语义互信息及信息效率之间的关系

- F1: 预测不确定性
- F2: 预测 不确定性+ 信息成本
- IND: 信息差 R- G
- 关系推导:

$$\begin{aligned} F_2 &= E_{q(x,z)} \left[\log \frac{q(z|x)}{P(z)} - \log P_\theta(x|z) \right] \\ &= E_{q(x,z)} \left[\log \frac{q(z|x)}{q(z)} + \log \frac{q(z)}{P(z)} \right] + H_\theta(x|z) \\ &= \underbrace{R - G}_{\text{信息差}} + \underbrace{D_{\text{KL}}(q(z) \parallel P(z))}_0 + \underbrace{H(x)}_{\text{香农熵=常数}} \approx \text{IND} + H(x). \end{aligned}$$

- 所以最小化F2 → 最小化信息差 → 最大化信息效率
- 这就是最小化VFE包含的优化思想

最小化F1和F2 的不同目标和优缺点

- F1: 预测不确定性, 最小化F1等价于最大化语义互信息和对数似然度
- 机器学习中用于优化模型参数 θ ;
- 分类, 主动推断中用于最小化不确定性或误差
- F2: 用于发现隐含变量, 混合模型 (聚类), 无监督学习。

VB: F2
混合模型/求解隐含变量

FEP: 最小化F1
预测/分类/主动推断

- 缺点:
- Friston等人错在没用F2, 没有区分两种任务, 在数学上漏掉了VB提高通信效率的努力;
- Hinton等人用 $q(z)$ 代替 $q(z|x)$, 容易引起混乱, 也不便于表达香农互信息。
- 改进: 用 $IND = R - sG = R - s [H(x) - F2]$ 代替F2和F1。
- $S=1$ 时, $\min IND$ 等价于 $\min F2$
- $S=\infty$ 时, $\min IND$ 等价于 $\min F1$.

主动推断误用F1的后果——变成误差控制

- 预测模型推断：优化模型参数，使主观预测符合客观事实；
- 主动推断：固定预测模型，改变客观事实，使客观事实符合主观期望——约束控制。

$$G=I(X;Z_\theta) = \sum_j \sum_i P(x_i | z_j) P(z_j) \log \frac{P(x_i | \theta_j)}{P(x_i)}$$

- 如果最小化F2, 两者相等VFE最小；而最小化F1, $P(x|z)$ 集中分布F1最小。
- 例如**主动推断**任务之一：想象公交车一路停靠很多站点。
- 给定初始状态 x_0 和一系列目标约束——用似然函数表示。因为有随机干扰，控制可能不准。优化准则是最小化 $F1 = H(x|z, \theta) = -E \log P(x|z, \theta)$ 。
- 和混合模型不同，**不要求** $q(x)=P(x)$ 。
- 假设似然函数是高斯分布，

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



- $-\text{Log}f(x)$ =相对误差。
- **最小化F1**就等价于最小化平均相对误差，主动推断就变成误差控制。控制信息成本的努力就漏掉了。好像新手频繁打方向盘（用F1），老手转得少（用F2）
- 如果用F2, 我们可能使用较小的控制量，近似满足一系列约束。

Friston等人真地误用先验分布作为变分？ 证据？

1. 在他们很多文章中， $q(z)$ 和 $q(z|x)$ 同时出现在一个公式中。而在VB中不会有(同时出现好表示香农互信息)

$R-G = D_{KL}(P(o)||P\pi(o)) + D_{KL}(P_\pi(s)||P(s))$

We can express the expected free energy of a policy as a bound on information gain and expected log (model) evidence (a.k.a., Bayesian risk):

$$G(\pi) = \underbrace{\mathbb{E}_Q[D_{KL}[Q(s_\tau, A|o_\tau, \pi)||P(s_\tau, A|o_\tau)]]}_{\text{Expected evidence bound}} - \underbrace{\mathbb{E}_Q[\log P(o_\tau)]}_{\text{Expected log evidence}} \quad \text{F2}$$

$$- \underbrace{\mathbb{E}_Q[D_{KL}[Q(s_\tau, A|o_\tau, \pi)Q(s_\tau, A|\pi)]]}_{\text{Expected information gain}} \quad \text{R=Shannon mutual information}$$

$$\geq - \underbrace{\mathbb{E}_Q[\log P(o_\tau)]}_{\text{Expected log evidence}} - \underbrace{\mathbb{E}_Q[D_{KL}[Q(s_\tau, A|o_\tau, \pi)||Q(s_\tau, A|\pi)]]}_{\text{Expected information gain}}$$

cross entropy Shannon Mutual information R (14)

2. 我和Friston先生本人有过交流，他说他用的期望自由能是F1。
3. 他也说从F1公式能推导出F2公式，但是他推荐的文章中，实际推导的是

$$F1 = H_\theta(x) - I(X;Y) = F2 - R,$$

这等价于我推导的 $F2 = R + F1$ 。

4. 他赞成语义互信息 $G=H(x)-F1$ ，所以F1接近 $H_\theta(x|z)$ 。

语义变分贝叶斯 (SVB) 与 β -VAE比较

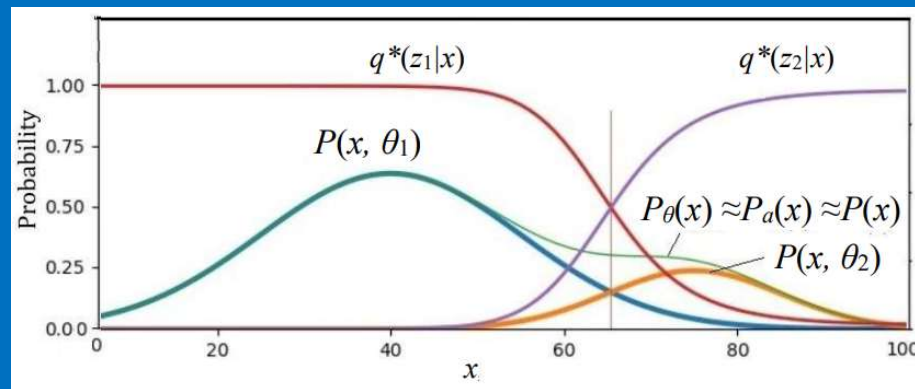
- 新潮流，改进VB到 β -VAE，涉及信息效率。
- Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S., Lerchner, A.: beta-VAE: Learning basic visual concepts with a constrained variational framework. In: Proceedings of the International Conference on Learning Representations
- 目标函数（最大化）：
 - $L = E_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - \beta D_{\text{KL}}(q_{\phi}(z|x) \parallel p(z))$,
- 等价于（最小化）：
 - $F2 = F1 + \beta R'$
- 比较SVB中目标函数：
 - $\text{IND} = R - sG$
- 但其信息项不是真正香农互信息。
- 我的SVB可直接计算R、G与G/R，允许约束函数是真值，隶属，相似，或失真函数。

实验1：三种任务优化的香农信道 $q(z|x)$ 的表现

- 1. 混合模型，最小化F2时： $q(z|x) \approx P(z|x)$

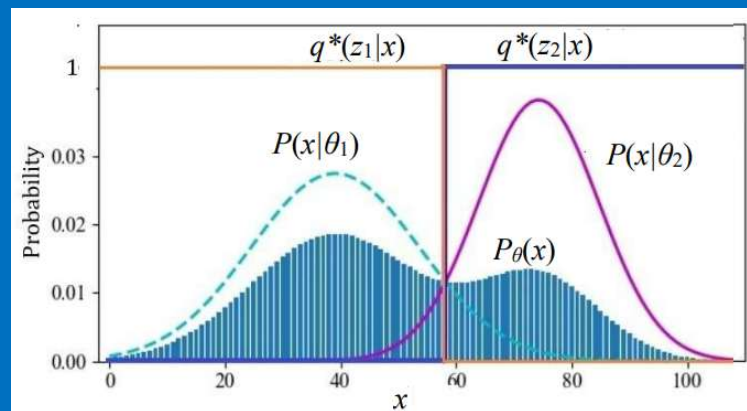
最小化 F2

这时, $F_2=0$, $G/R=1$, $s=1$,
 $q(x) \approx P_\theta(x) \approx P(x)$



- 2. 最大互信息分类

最小化 F1

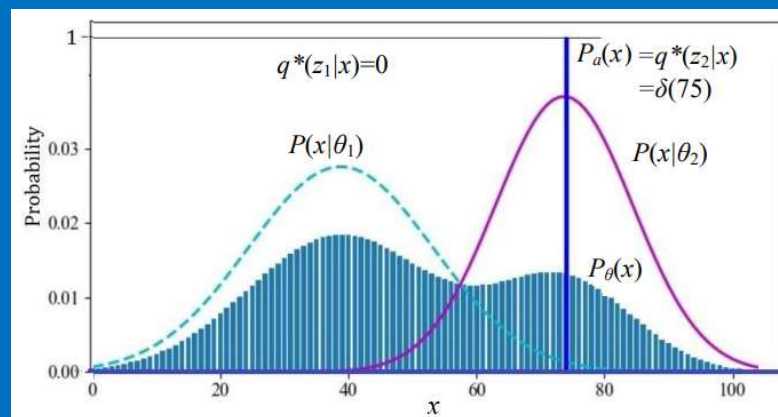


G/R 下降

- 3. 主动推断或约束控制时，最大F1
或最大化语义信息

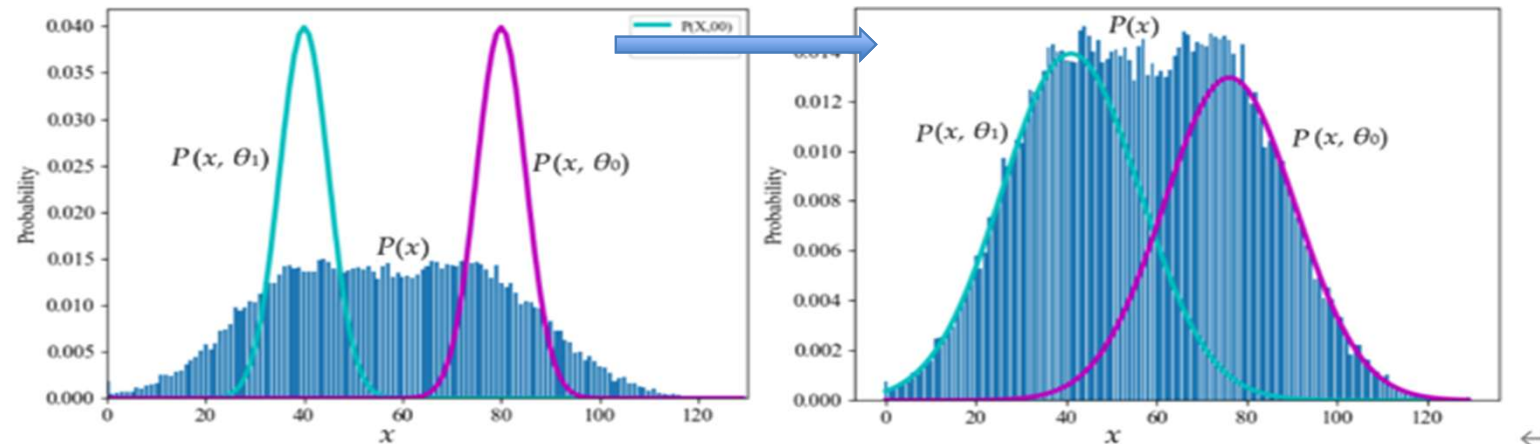
$P_a^*(x) = g^*(z_2|x) = \delta(75)$, 不需要 $q(x) \approx P(x)$.

最小化 F1



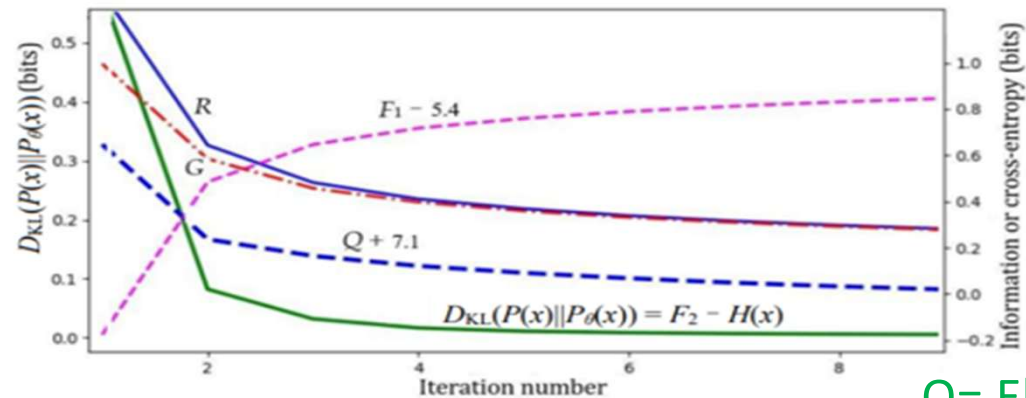
实验2: F2下降而F1上升的例子

- 存在反例, F2下降时F1可能增加, 说明两种自由能并不等价。
F2减小, F1增大



(a) The iteration starts

(b) The iteration converges



$$Q = E \log P_{\theta}(x, z)$$

(c) F_2 , $R - G$, F_1 , and Q change in the iterative process.

Fig. 5. During the convergence of the mixture model, $D_{KL}(P(s) \parallel P_{\mu}(s))$, $R - G$, and F_2 decreased, while F_1 increased, and Q decreased.

主要参考文献

1. Hinton, G.E.; van Camp, D. Keeping the neural networks simple by minimizing the description length of the weights. In *Proc. 6th Annu. Conf. Comput. Learn. Theory*; 1993; pp. 5–13.
2. Hinton, G.E.; Zemel, R.S. Autoencoders, minimum description length and Helmholtz free energy. In *Proc. 6th Int. Conf. Neural Inf. Process. Syst.*; 1993; pp. 3–10.
3. Neal, R.; Hinton, G. A view of the EM algorithm that justifies incremental, sparse, and other variants. In *Learning in Graphical Models*; Jordan, M.I., Ed.; MIT Press: Cambridge, MA, USA, 1999; pp. 355–368.
4. Wikipedia. Variational Bayesian methods. Available online: https://en.wikipedia.org/wiki/Variational_Bayesian_methods
5. Friston, K. The free-energy principle: A unified brain theory? *Nat. Rev. Neurosci.* **2010**, *11*, 127–138. doi:10.1038/nrn2787.
6. Friston, K.J.; Parr, T.; de Vries, B. The graphical brain: Belief propagation and active inference. *Netw. Neurosci.* **2017**, *1*, 381–414. doi:10.1162/netn_a_00018.
7. Parr, T.; Pezzulo, G.; Friston, K.J. *Active Inference: The Free Energy Principle in Mind, Brain, and Behavior*; MIT Press: Cambridge, MA, USA, 2022.
8. Lu, C. Improving the minimum free energy principle to the maximum information efficiency principle. *Entropy* **2025**, *27*, 684. doi:10.3390/e27070684.
9. Lu, C. A semantic generalization of Shannon's information theory and applications. *Entropy* **2025**, *27*, 461. doi:10.3390/e27050461.
10. Higgins, I.; et al. beta-VAE: Learning basic visual concepts with a constrained variational framework. In *ICLR*; 2017.

总结1：核心逻辑链

- 两种变分 $q(z)=q(z|x)$, $q(z)\leftrightarrow q(z|x)$



- 两种期望自由能 $F2$, $F1$



- 两种优化目标：最小化 $F2$ 即最小化 $R-G$ ，最小化 $F1 \approx H(x|z, \theta)$ 即最大化语义信息 G



混合模型，无监督学习需要 $F2$ ；主动推断实际用 $F1$



- 改进 $VB \rightarrow SVB$, $F2 \rightarrow R-sG$, 重视权衡语义信息和信息效率

总结2：结论

- VB和FEP中变分不同，变分自由能也不同。
- 用 $q(z)$ 作为 $q(z|x)$ 的简写，得到期望的自由能 F_2 ，用 $q(z) \llcorner q(z|x)$ 作为变分，得到 F_1 。
- 最小化 F_2 包含最小化 F_1 ，反之不然。
- 混合模型需要最小化 F_2 ，减小不确定性或优化预测模型需要最小化 F_1 。
- Friston等人结论，变分自由能表示既预测不准确性 ($H(x|z, \theta)$)，也表示平均惊奇 $H_\theta(x)$ ，这是错的，错在使用 $q(z) \llcorner q(z|x)$ 作为变分得到 $F_1 \approx H(x|z, \theta)$ ，客观上漏掉了VB追求信息效率的努力。
- VB用 $q(z)$ 作为后验分布 $q(z|x)$ 的简写不便于表达香农互信息，也不便于计算信息效率。
- SVB用信息差 $R - sG$ 作为目标函数代替 F_2 ，可以继承VB重视信息效率的思想，并能在最大语义信息和最大信息效率之间权衡。

感谢熊楚瑜博士，为他五年前提醒我关注Friston的最小自由能原理。

欢迎批评交流，更多讨论见我主页：<http://www.survivor99.com>