

# 表現学習の確率熱力学的解釈 (1/6)

機械学習と物理学の接点（理論の形式的類似性）を紹介します. その応用例として, 物理学側のツールである「準静的過程に沿った状態操作」を機械学習に導入します.

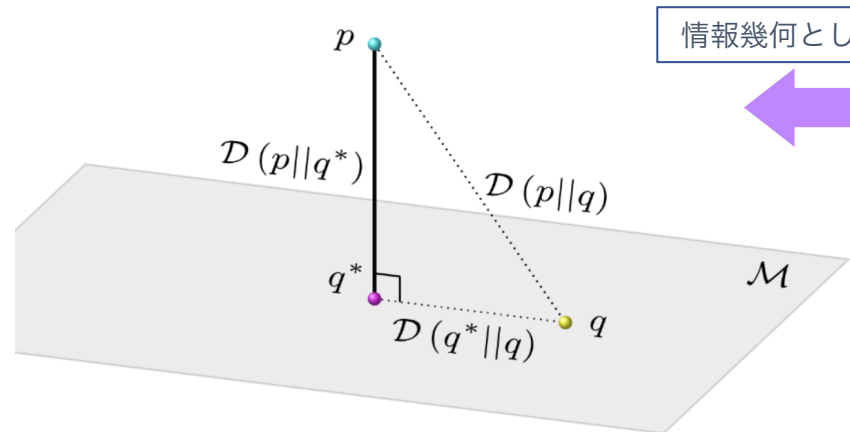
- 確率熱力学：ゆらぐ系の熱力学. 微小系が状態  $X$  から状態  $Z$  に遷移する過程（順過程と呼ぶ）の確率を  $p$  とし, その逆過程の確率を  $q$  とするとき,  $p$  と  $q$  との距離が, 微小系の状態変化に伴うエントロピー生成量  $\sigma_{\text{tot}}$  を決定する.

$$\mathcal{D}(p(s)||q(s)) = \mathcal{D}(p(s)||q^*(s)) + \mathcal{D}(q^*(s)||q(s)).$$

$$\sigma_{\text{tot}}^S := \sigma_{\text{sys}}^S + \sigma_{\text{bath}}^S.$$

特徴量の生成プロセスの学習

ゆらぐ系のエントロピー生成  
entropy production (2nd law)



情報幾何として共通の幾何構造

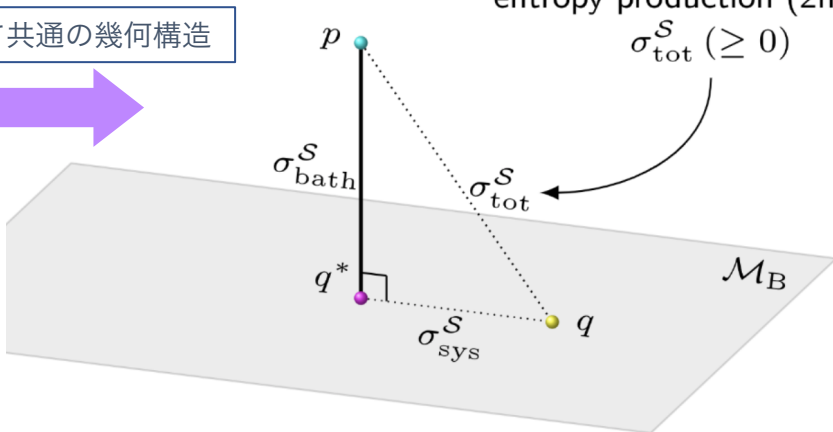
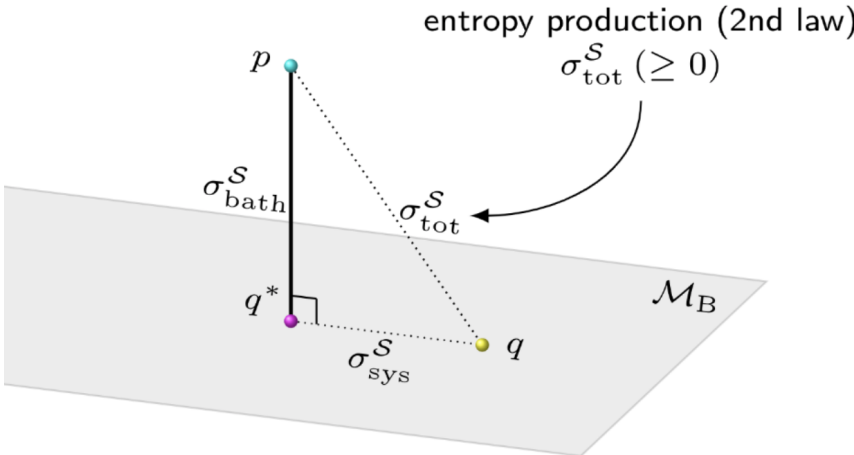
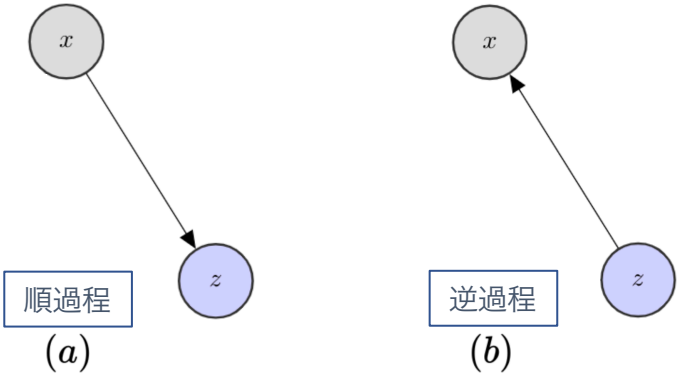


Figure 1: Information-Geometric Pythagorean theorem.

Figure 2: Total entropy production.

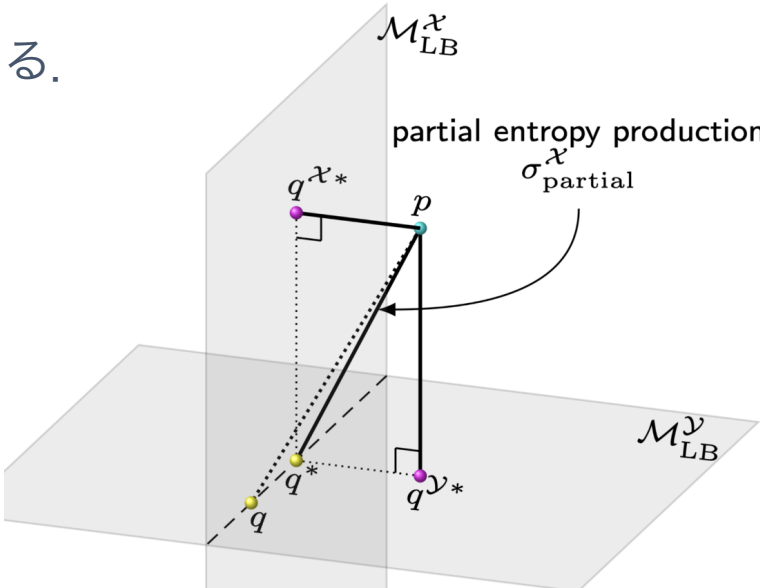
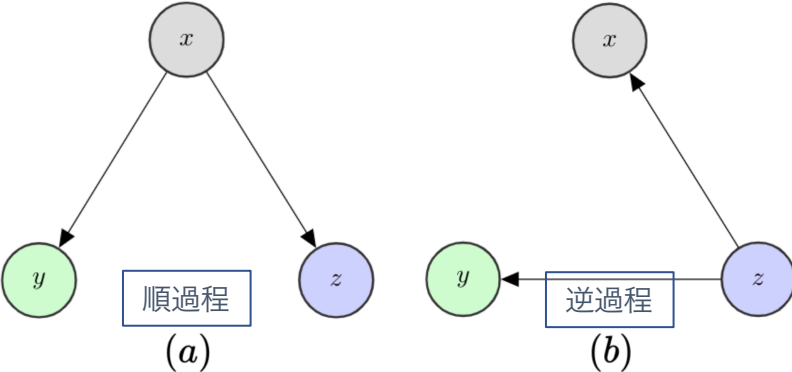
# 表現学習の確率熱力学的解釈 (2/6)

- 全系のエントロピー生成：順過程と逆過程とで，グラフィカルモデルの全要素について矢印の向きが逆になる.



$$p(s) := T(z, x) p(x) = p(z|x) p(x), q(s) := T(x, z) q(z) = q(x|z) q(z) \quad \mathcal{M}_B = \{q(s) \mid q(s) := T(x, z) q(z) = q(x|z) q(z)\}.$$

- 部分系のエントロピー生成：順過程と逆過程とで，グラフィカルモデルの一部の要素について矢印の向きが逆になる.



$$p(s) := p(z|x)p(y|x)p(x) \quad q(s) := q(x|z)q(y|z)q(z)$$

# 表現学習の確率熱力学的解釈 (3/6)

- 表現学習の学習プロセス：データ  $\mathbf{x}$  から潜在変数（特徴量） $\mathbf{z}$  の生成と教師ラベル  $\mathbf{y}$  の予測を行う順過程の確率分布  $p$  と，潜在変数（特徴量） $\mathbf{z}$  からデータ  $\mathbf{x}$  と教師ラベル  $\mathbf{y}$  をデコードする逆過程の確率分布  $q$  の距離を可能な限り近づけようとするプロセス．部分系のエントロピー生成と解釈可能．

$$H := \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p(\mathbf{x}, \mathbf{y})} [-\log p(\mathbf{x}, \mathbf{y})],$$

データのエントロピー

$$D := \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \left[ - \int e(\mathbf{z}|\mathbf{x}) \log d(\mathbf{x}|\mathbf{z}) d\mathbf{z} \right],$$

再構成誤差

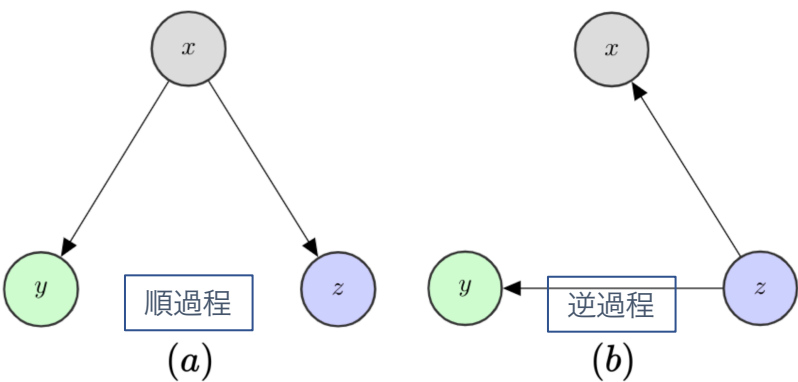
$$R := \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \left[ \int e(\mathbf{z}|\mathbf{x}) \log \frac{e(\mathbf{z}|\mathbf{x})}{q(\mathbf{z})} d\mathbf{z} \right],$$

正則化項

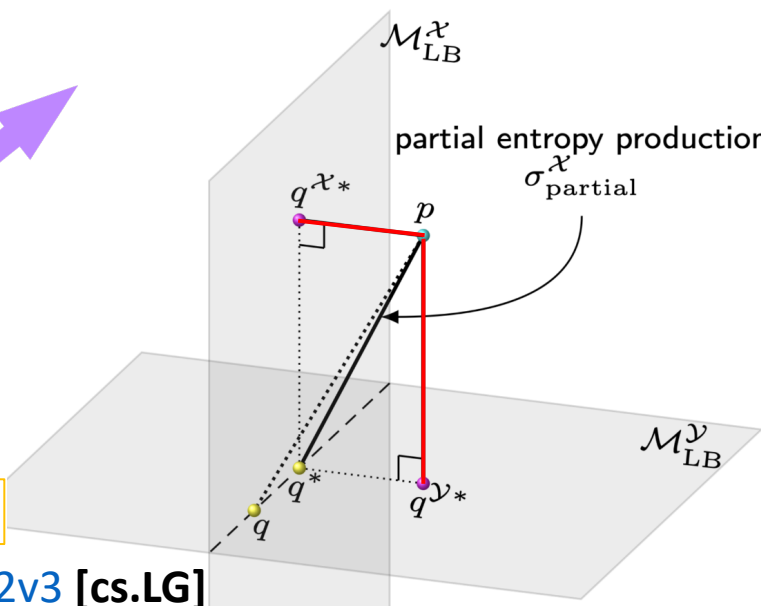
$$C := \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \left[ - \int e(\mathbf{z}|\mathbf{x}) \log c(\mathbf{y}|\mathbf{z}) d\mathbf{z} \right],$$

分類誤差

$$\mathcal{D}(p(\mathbf{s})||q(\mathbf{s})) = -H + R + \lambda D + \gamma C \geq 0$$
$$:= -H + \mathcal{J}(\lambda, \gamma) \geq 0$$



$$p(\mathbf{s}) := e(\mathbf{z}|\mathbf{x})p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$$
$$q(\mathbf{s}) := d_{\theta}(\mathbf{x}|\mathbf{z})^{\lambda} c_{\theta}(\mathbf{y}|\mathbf{z})^{\gamma} q_{\theta}(\mathbf{z})$$



# 表現学習の確率熱力学的解釈 (4/6)

- 深層生成モデルを使った表現学習：エンコーダー，デコーダーとして深層学習モデル（例えばCNN）を使い，そのパラメータを  $\theta$  とする.

$$H := \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p(\mathbf{x}, \mathbf{y})} [-\log p(\mathbf{x}, \mathbf{y})],$$

データのエントロピー

$$D := \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \left[ - \int e(\mathbf{z}|\mathbf{x}) \log d(\mathbf{x}|\mathbf{z}) d\mathbf{z} \right],$$

再構成誤差

$$R := \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \left[ \int e(\mathbf{z}|\mathbf{x}) \log \frac{e(\mathbf{z}|\mathbf{x})}{q(\mathbf{z})} d\mathbf{z} \right],$$

正則化項

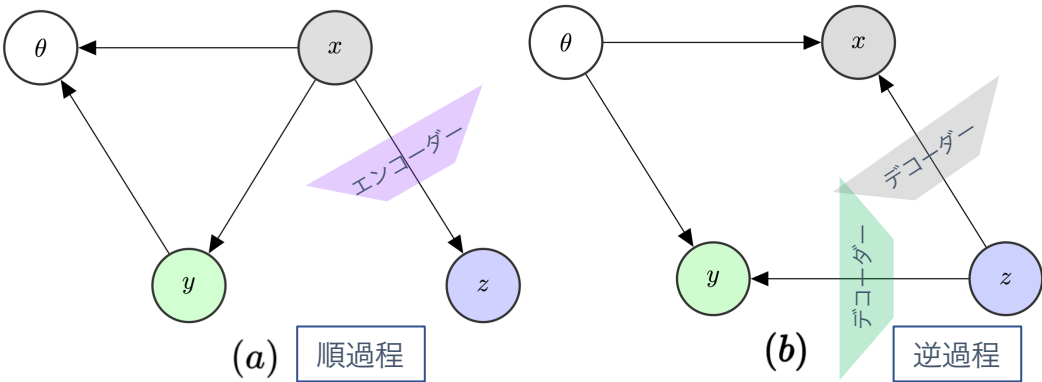
$$C := \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \left[ - \int e(\mathbf{z}|\mathbf{x}) \log c(\mathbf{y}|\mathbf{z}) d\mathbf{z} \right],$$

分類誤差

$$S := \mathbb{E}_{\theta, \mathbf{x}, \mathbf{y} \sim p(\theta, \mathbf{x}, \mathbf{y})} \left[ \log \frac{p(\theta|\{\mathbf{x}, \mathbf{y}\})}{q(\theta)} \right]$$

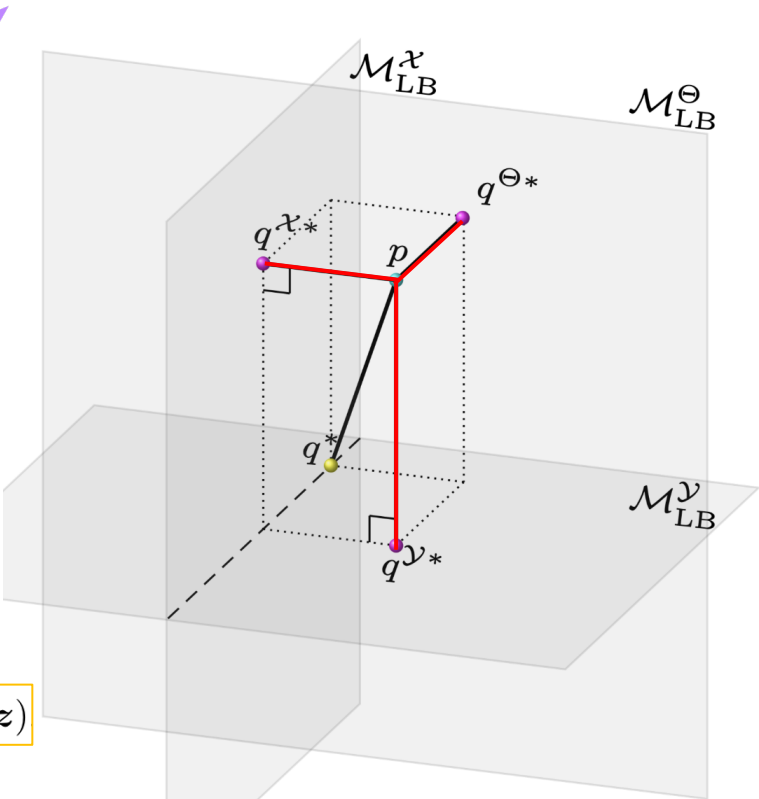
パラメータの  
相対エントロピー

$$D(p(\mathbf{s})||q(\mathbf{s})) = -H + R + \lambda D + \gamma C + \sigma S \geq 0$$
$$:= -H + \mathcal{J}(\lambda, \gamma, \sigma) \geq 0$$



$$p(\mathbf{s}) := p(\theta|\{\mathbf{x}, \mathbf{y}\})^\sigma e(\mathbf{z}|\mathbf{x}) p(\mathbf{y}|\mathbf{x}) p(\mathbf{x})$$

$$q(\mathbf{s}) := q(\theta)^\sigma d_\theta(\mathbf{x}|\mathbf{z})^\lambda c_\theta(\mathbf{y}_x|\mathbf{z})^\gamma q_\theta(\mathbf{z})$$



# 表現学習の確率熱力学的解釈 (5/6)

- 表現学習と熱力学諸法則の形式的対応

$$p(s) := p(\theta|\{\mathbf{x}, \mathbf{y}\})^\sigma e(\mathbf{z}|\mathbf{x})p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$$

$$q(s) := q(\theta)^\sigma d_\theta(\mathbf{x}|\mathbf{z})^\lambda c_\theta(\mathbf{y}_x|\mathbf{z})^\gamma q_\theta(\mathbf{z})$$

$$H := \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p(\mathbf{x}, \mathbf{y})} [-\log p(\mathbf{x}, \mathbf{y})],$$

データのエントロピー

$$D := \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \left[ -\int e(\mathbf{z}|\mathbf{x}) \log d(\mathbf{x}|\mathbf{z}) d\mathbf{z} \right],$$

再構成誤差

$$R := \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \left[ \int e(\mathbf{z}|\mathbf{x}) \log \frac{e(\mathbf{z}|\mathbf{x})}{q(\mathbf{z})} d\mathbf{z} \right],$$

正則化項

$$C := \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \left[ -\int e(\mathbf{z}|\mathbf{x}) \log c(\mathbf{y}|\mathbf{z}) d\mathbf{z} \right],$$

分類誤差

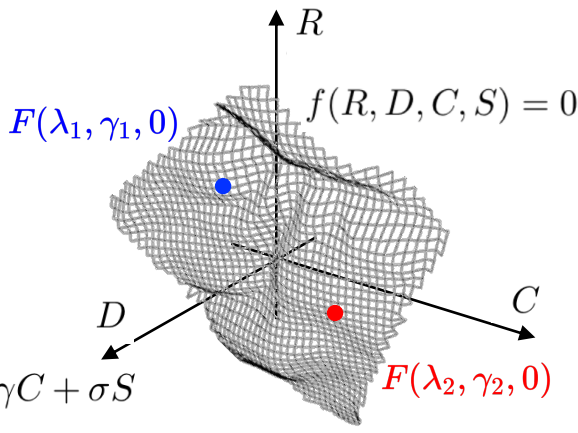
$$S := \mathbb{E}_{\theta, \mathbf{x}, \mathbf{y} \sim p(\theta, \mathbf{x}, \mathbf{y})} \left[ \log \frac{p(\theta|\{\mathbf{x}, \mathbf{y}\})}{q(\theta)} \right]$$

パラメータの  
相対エントロピー

$$\mathcal{D}(p(s)||q(s)) = -H + R + \lambda D + \gamma C + \sigma S \geq 0$$
$$:= -H + \mathcal{J}(\lambda, \gamma, \sigma) \geq 0$$

$$\mathcal{J}(\lambda, \gamma, \sigma) \geq H$$

$$F(\lambda, \gamma, \sigma) := \min_{e(\mathbf{z}|\mathbf{x}), q(\mathbf{z}), d(\mathbf{x}|\mathbf{z}), c(\mathbf{y}|\mathbf{z})} \mathcal{J}(\lambda, \gamma, \sigma)$$
$$= \min_{e(\mathbf{z}|\mathbf{x}), q(\mathbf{z}), d(\mathbf{x}|\mathbf{z}), c(\mathbf{y}|\mathbf{z})} R + \lambda D + \gamma C + \sigma S$$



**Remark 1 (“the first law” of learning).**

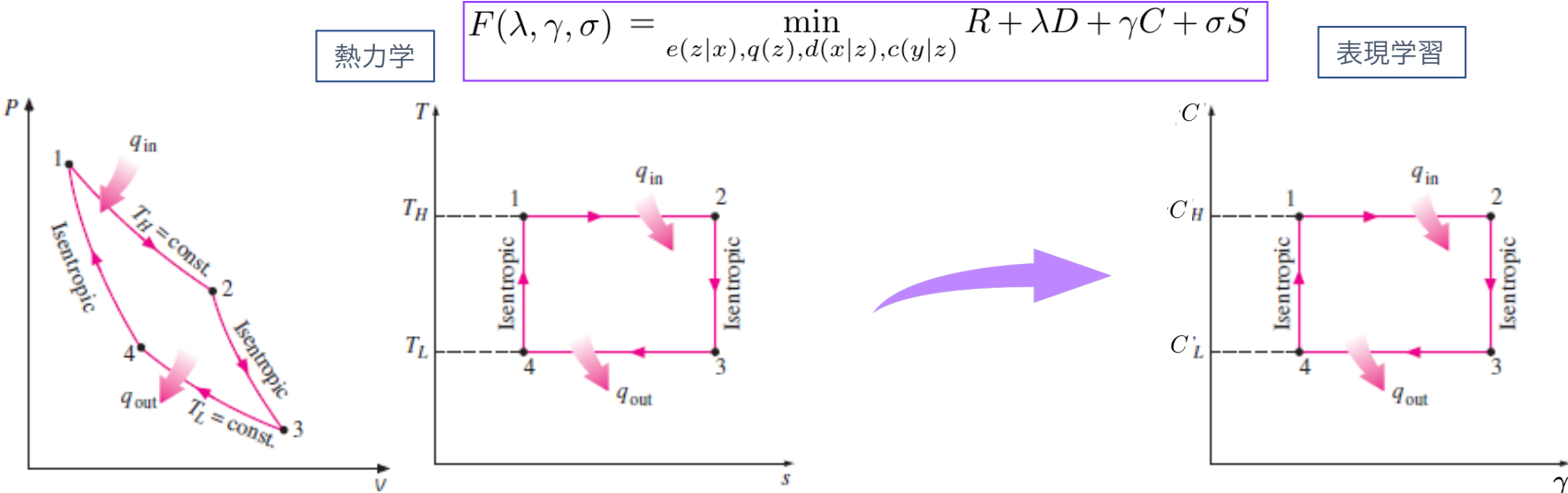
$$dR = -\lambda dD - \gamma dC - \sigma dS$$

**Remark 2 (“the second law” of learning).**

$$\lambda \equiv -\left(\frac{\partial R}{\partial C}\right)_{D,S} \quad \gamma \equiv -\left(\frac{\partial R}{\partial D}\right)_{C,S} \quad \sigma \equiv -\left(\frac{\partial R}{\partial S}\right)_{C,D}$$
$$\mathcal{D}(p(s)||q(s)) = -H + R + \lambda D + \gamma C + \sigma S := -H + \mathcal{J}(\lambda, \gamma, \sigma) \geq 0$$

# 表現学習の確率熱力学的解釈 (6/6)

- 「準静的過程に沿った状態操作」を考慮することで、統計モデルの分類誤差を低い値に抑えたまま、データのドメインを変化させることはできないか？



- データのドメインを変化させる間、分類誤差の値が不変となることを要請する

$$p(x, t) = \operatorname{argmin}_{p(x)} (1 - t)W_2^2(p(x^{ini}), p(x)) + tW_2^2(p(x), p(x^{fin}))$$

$$p(x, t) = \sum_{i=1}^{n_{ini}} \sum_{j=1}^{n_{fin}} \Gamma_{ij} \delta_{x - (1-t)x_i^{ini} - tx_j^{fin}}$$

$$y(x, t) = (1 - t)\delta_{y - y_{x_i^{ini}}} + t\delta_{y - y_{x_j^{fin}}}$$

$\frac{dC(t)}{dt} = 0$

