

イントロ (1/n)

今日は、表現学習と熱力学の形式的類似性を紹介した後、その拘束条件について議論します。
そのお話をする前に、生成モデルの基礎についてお話します。

生成モデルと認識

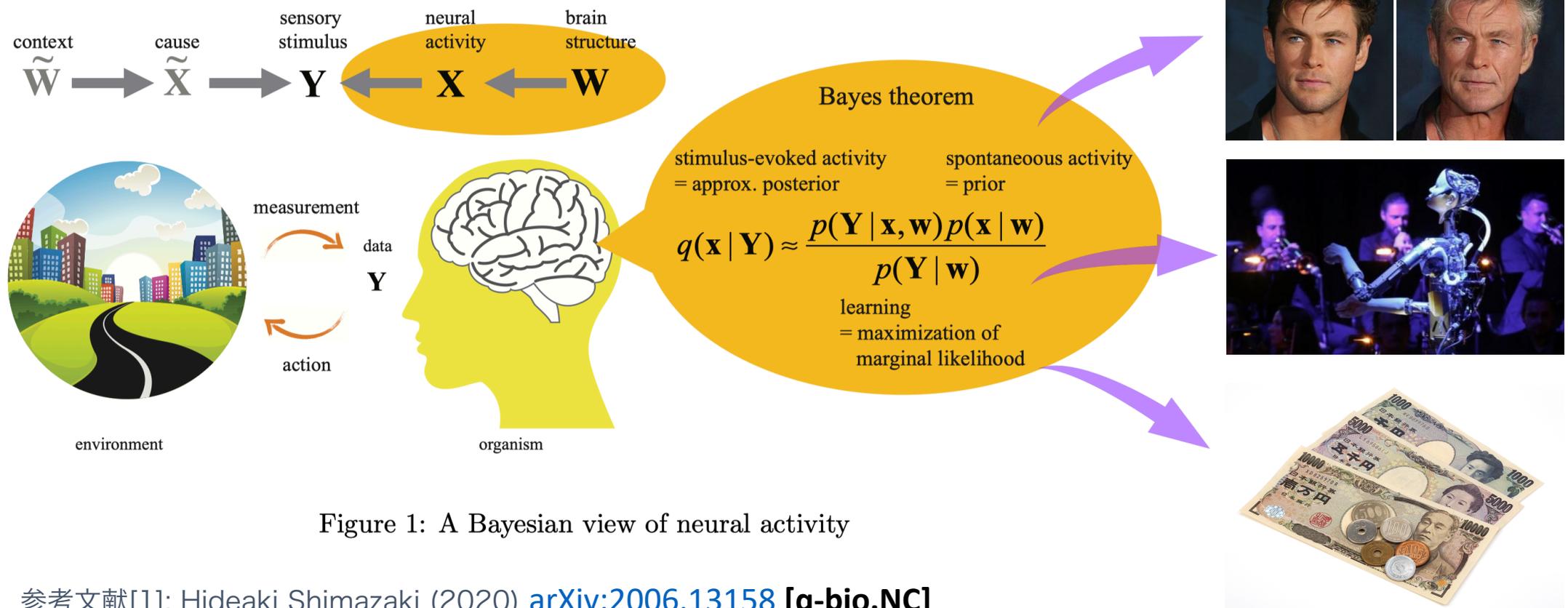


Figure 1: A Bayesian view of neural activity

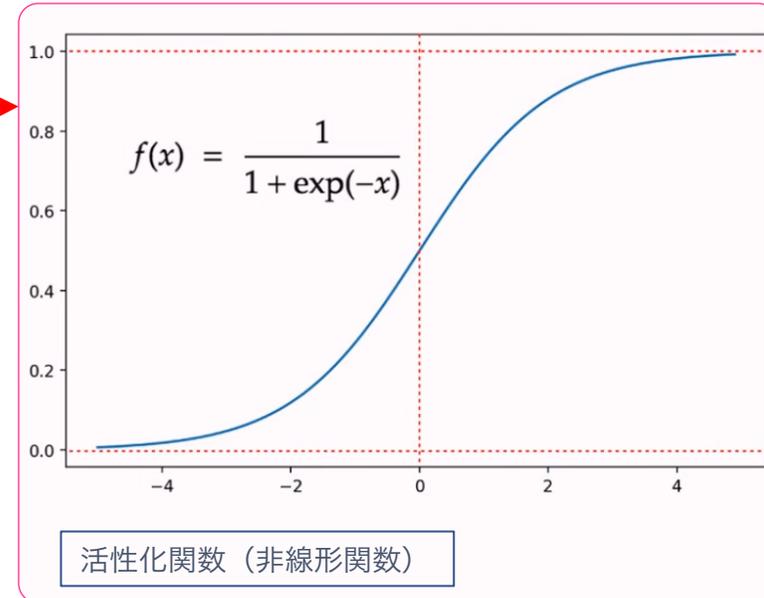
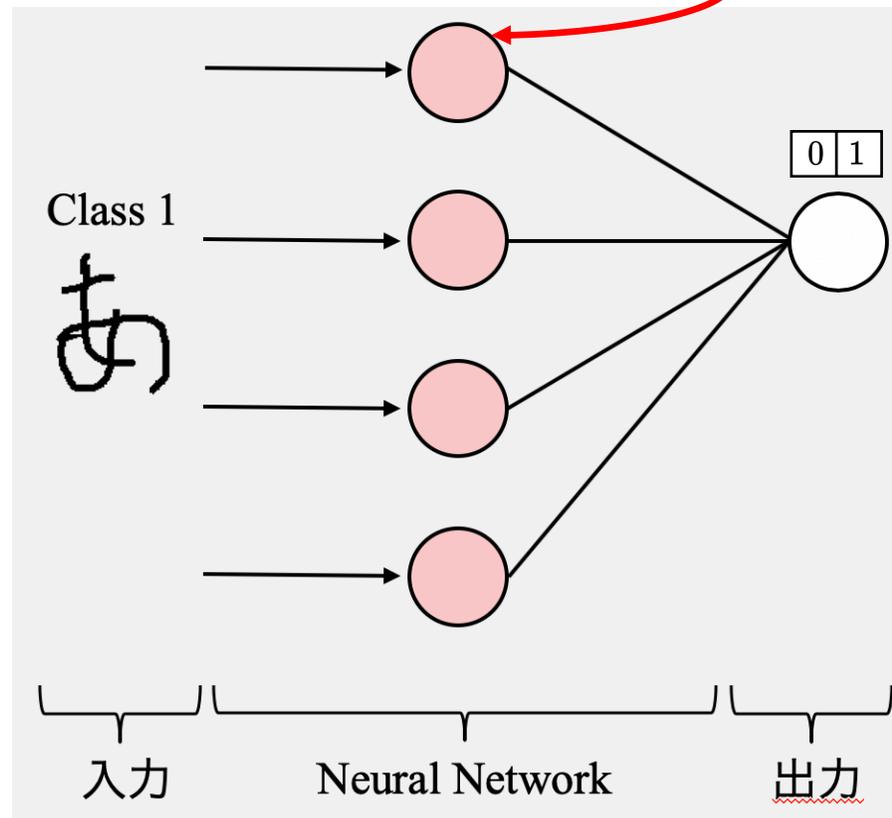
参考文献[1]: Hideaki Shimazaki (2020) [arXiv:2006.13158](#) [q-bio.NC]

参考文献[2]: Sumio Watanabe, Mathematical Theory of Bayesian Statistics, CRC Press (2018)

イントロ (2/n)

- 機械学習でやっていること

- 入力 x と出力 y を結びつける関数 $f(x)$ を推定する。



$$\theta_{k+1} = \theta_k - \eta \frac{\partial \mathcal{L}(y'_k, y_k)}{\partial \theta_k}$$

$$\mathcal{L}(y', y)$$

$$y' \begin{bmatrix} 0 \\ 1 \end{bmatrix} \longleftrightarrow \begin{bmatrix} 0 \\ 1 \end{bmatrix} y = f(\text{あ})$$

目的関数 (誤差関数)

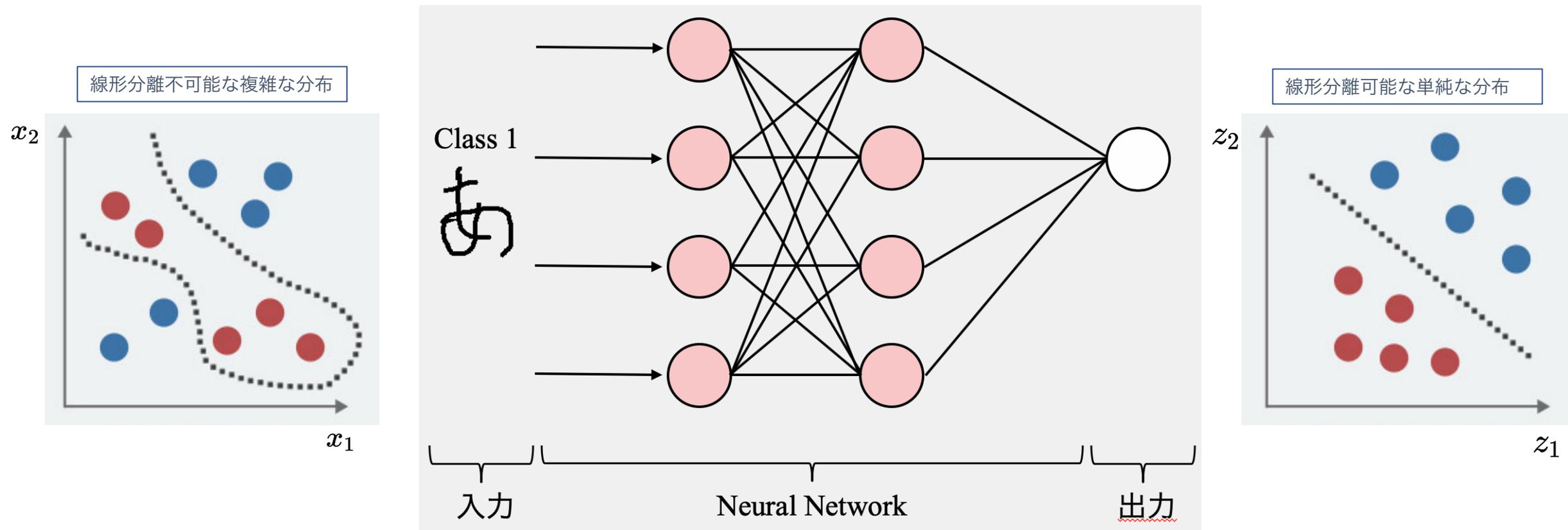
真の教師ラベル

モデルの予測

イントロ (3/n)

- ニューラルネットの層を一つ増やしたときに何が起こるか.

$$z = f_1(x) \quad y = f_2(z) = f_2(f_1(x))$$



- 統計モデルの内部では, 入力 x を, より単純な分布をもつ特徴量 z に変換している.
中間層無しのニューラルネットは, ロジスティック回帰モデルと同値.

イントロ (4/n)

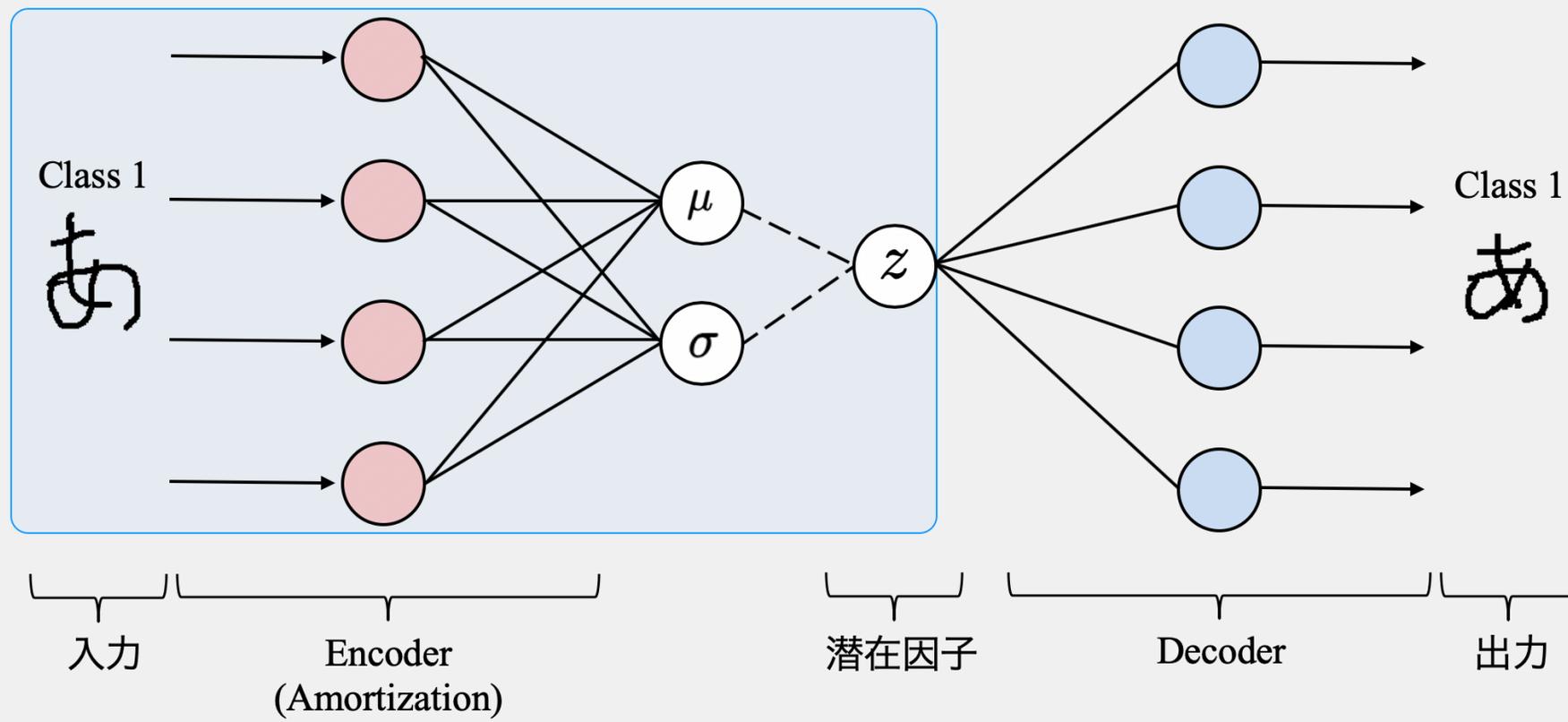
- 生成モデル (VAE)

$$z \sim \mathcal{N}(\mu(x), \sigma^2(x))$$

ガウス分布

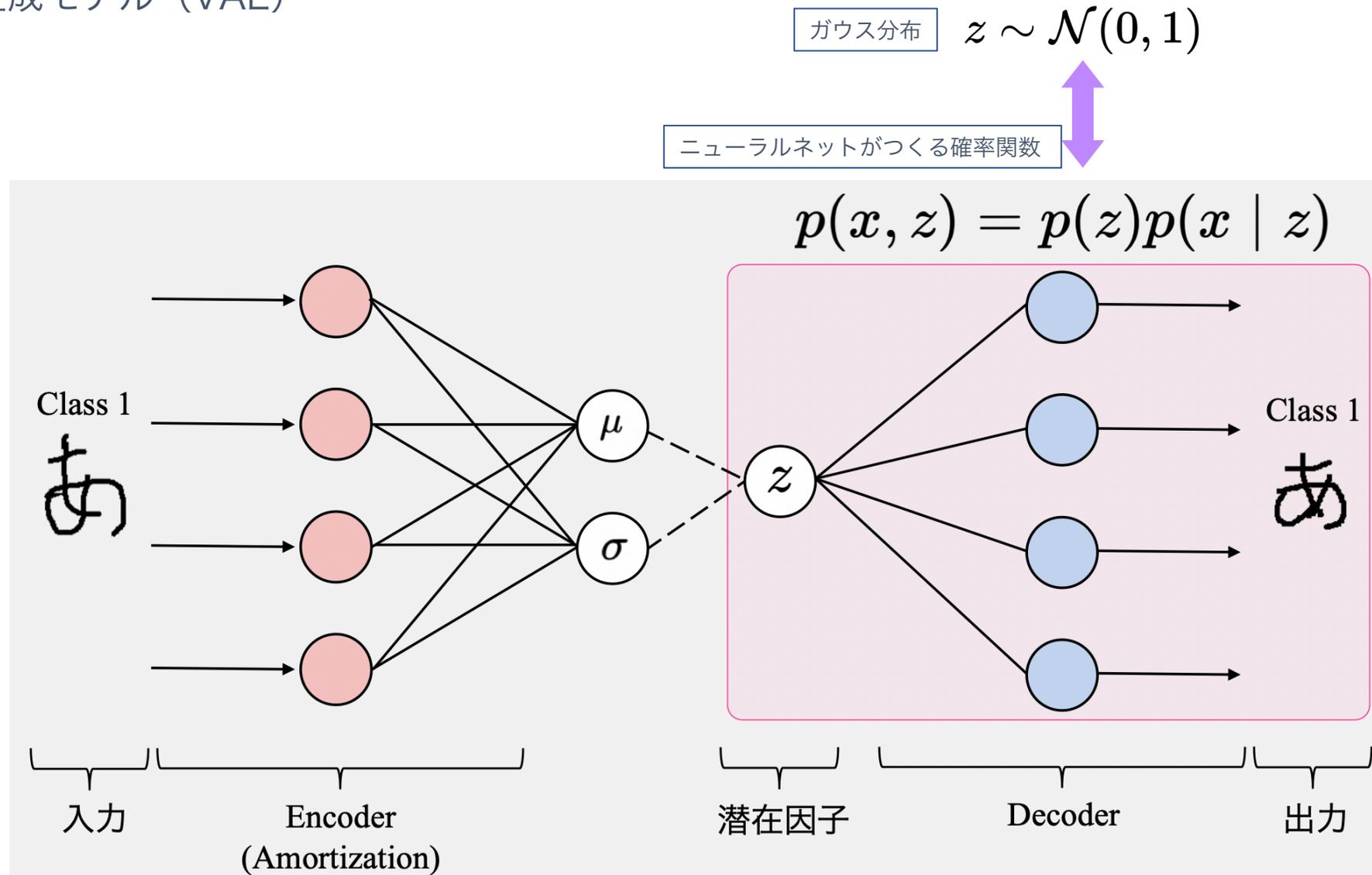
ニューラルネットがつくる確率関数

$$q(x, z) = q(x)q(z | x)$$

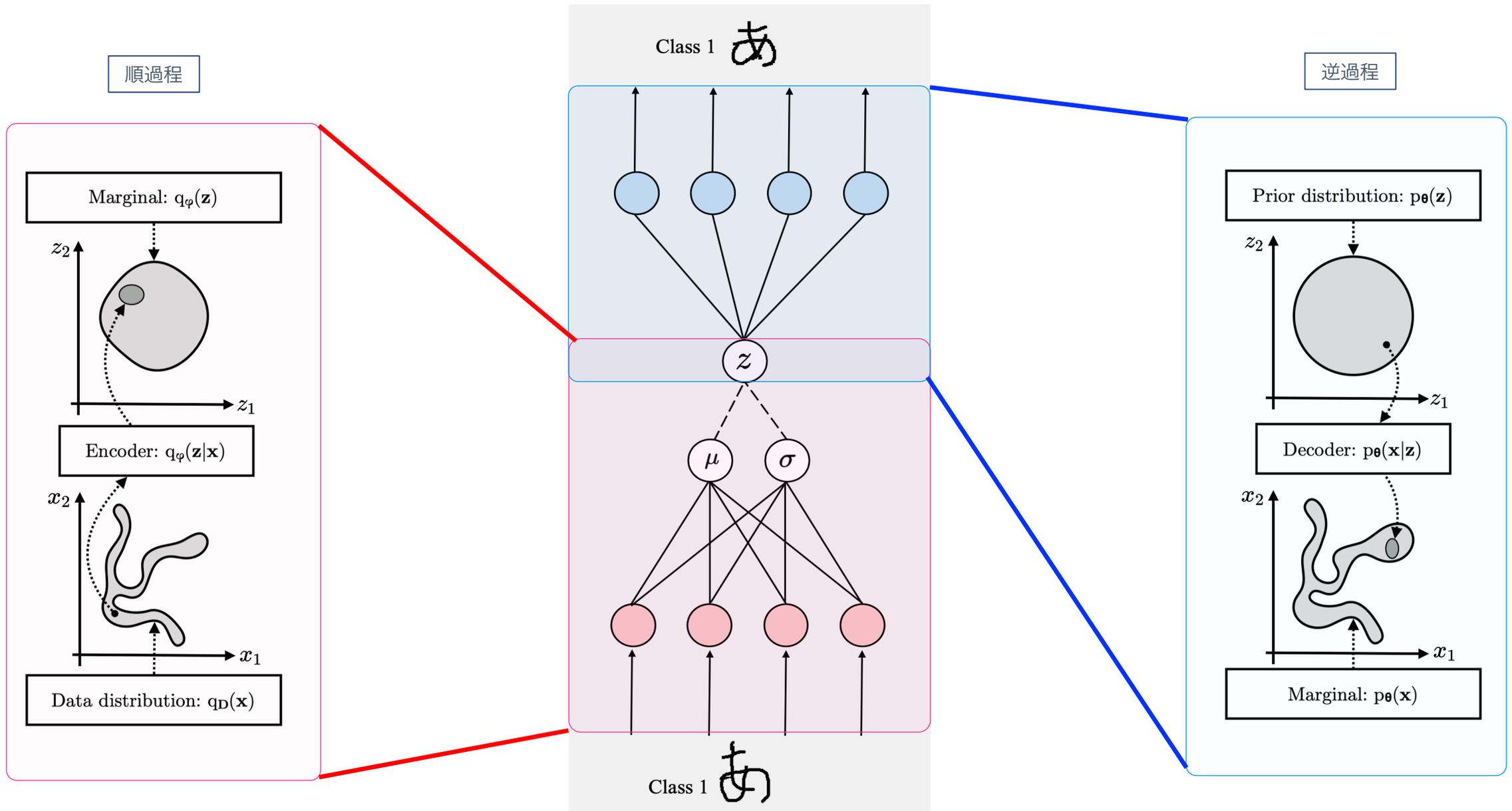


イントロ (5/n)

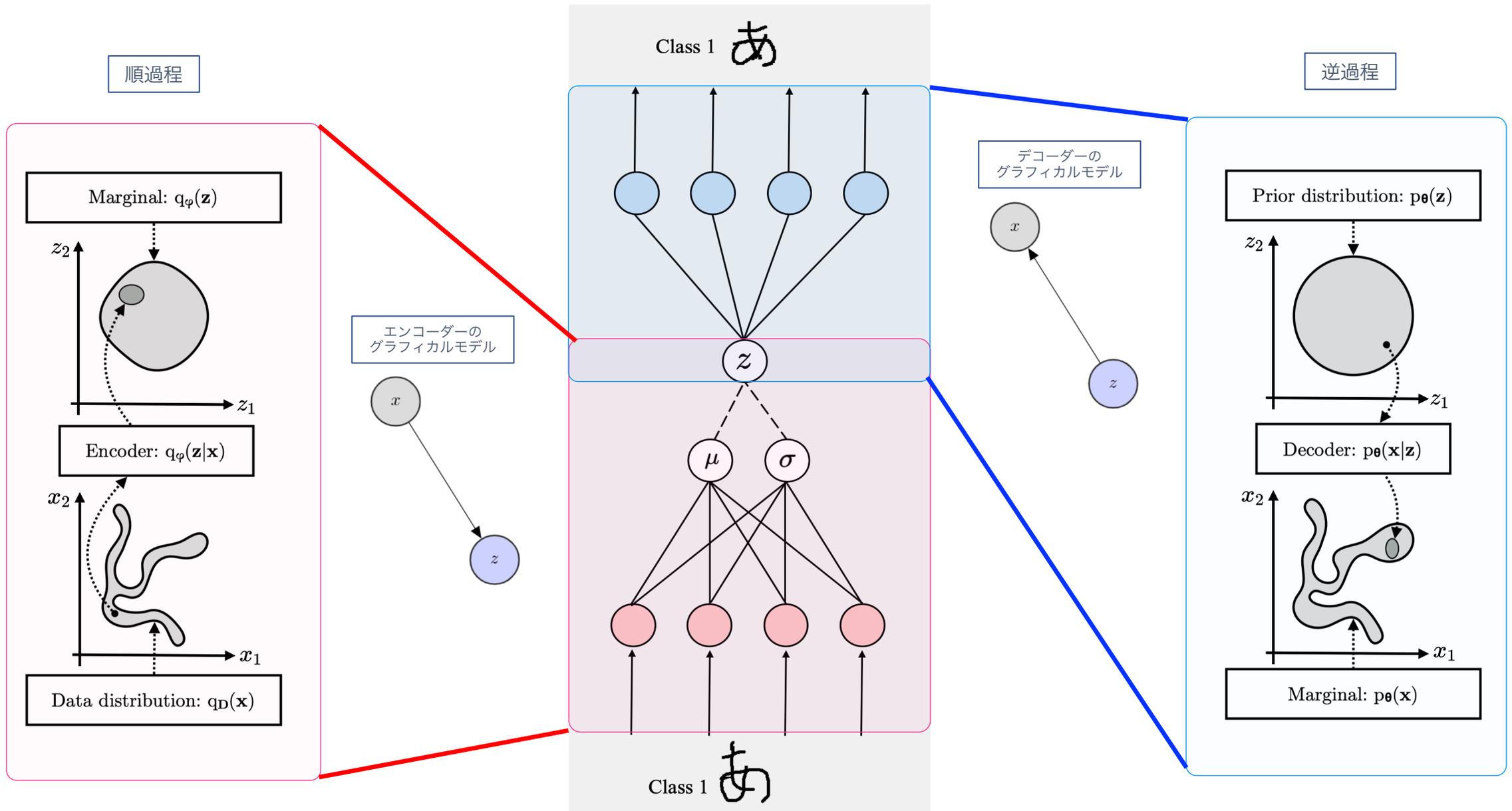
- 生成モデル (VAE)



イントロ (6/n)

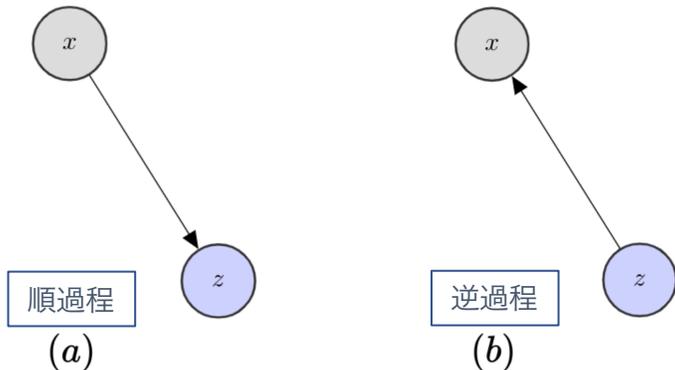


イントロ (7/n)

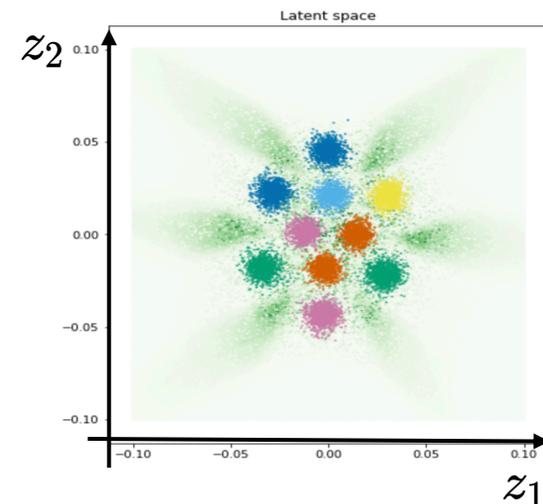


イントロ (8/n)

- 生成モデルの学習は、順過程と逆過程、それぞれの関数を推定するプロセスと理解できる。
その目的関数 (損失関数) は、順過程と逆過程、それぞれの同時分布を近づける形になっている。



$$p(\mathbf{s}) := T(\mathbf{z}, \mathbf{x}) p(\mathbf{x}) = p(\mathbf{z}|\mathbf{x}) p(\mathbf{x}) \quad q(\mathbf{s}) := T(\mathbf{x}, \mathbf{z}) q(\mathbf{z}) = q(\mathbf{x}|\mathbf{z}) q(\mathbf{z})$$



$$\begin{aligned} \mathcal{D}(p(\mathbf{s})||q(\mathbf{s})) &= \int p(\mathbf{s}) \log \frac{p(\mathbf{s})}{q(\mathbf{s})} d\mathbf{s} \\ &= \int p(\mathbf{s}) \log p(\mathbf{x}) d\mathbf{s} + \int p(\mathbf{s}) \log \frac{p(\mathbf{z}|\mathbf{x})}{q(\mathbf{z})} d\mathbf{s} - \int p(\mathbf{s}) \log q(\mathbf{x}|\mathbf{z}) d\mathbf{s} \\ &= -H(\mathbf{x}) + R + D \geq 0 \end{aligned}$$

データのエントロピー

$$H(\mathbf{x}) := - \int p(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x}$$

$\mathcal{L}_{\text{ELBO}}$

VAEの目的関数 変分下限 (ELBO)

$$\begin{aligned} \mathcal{D}(p(\mathbf{s})||q(\mathbf{s})) &= \int p(\mathbf{s}) \log \frac{p(\mathbf{s})}{q(\mathbf{s})} d\mathbf{s} \\ &= \int p(\mathbf{s}) \log \frac{p(\mathbf{z}|\mathbf{x})p(\mathbf{x})}{q(\mathbf{x}|\mathbf{z})p(\mathbf{z})} d\mathbf{s} + \int p(\mathbf{s}) \log \frac{p(\mathbf{z})}{q(\mathbf{z})} d\mathbf{s} \\ &:= \sigma_{\text{bath}}^S + \sigma_{\text{sys}}^S = \sigma_{\text{tot}}^S \geq 0 \end{aligned}$$

全系のエントロピー生成率

この変分下限を<自由エネルギー>とみなす
エントロピー生成と自由エネルギーには

$$\sigma_{\text{tot}}^S = \beta[W - \Delta F]$$

の関係式が成り立つので

$$\beta W = -H(\mathbf{x})$$

$$-\beta \Delta F = R + D$$

を同一視すると約束すれば

$$\sigma_{\text{tot}}^S = -H(\mathbf{x}) + R + D$$

表現学習の熱力学 (9/n)

ここから表現学習と熱力学の形式的類似性について議論していきます。後半で、転移学習のプロセスに準静的過程のアナロジーを導入して、サイクルとその熱効率について議論します。

- 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 σ_{tot} を決定する。

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

特徴量の生成プロセスの学習

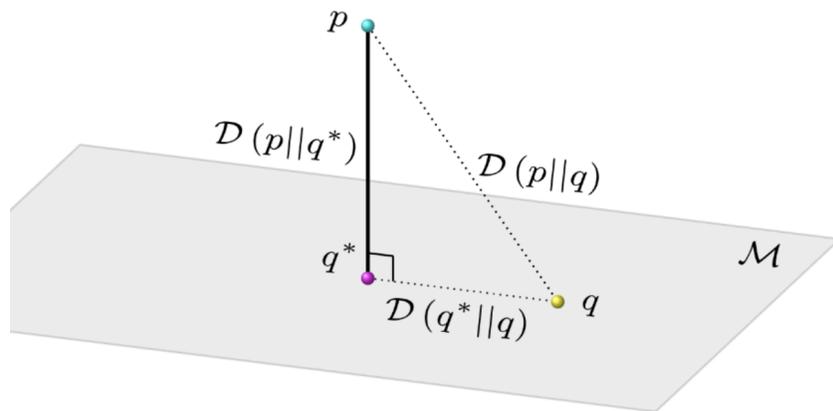


Figure 1: Information-Geometric Pythagorean theorem.

情報幾何として共通の幾何構造



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

ゆらぐ系のエントロピー生成

entropy production (2nd law)

$$\sigma_{\text{tot}}^S (\geq 0)$$

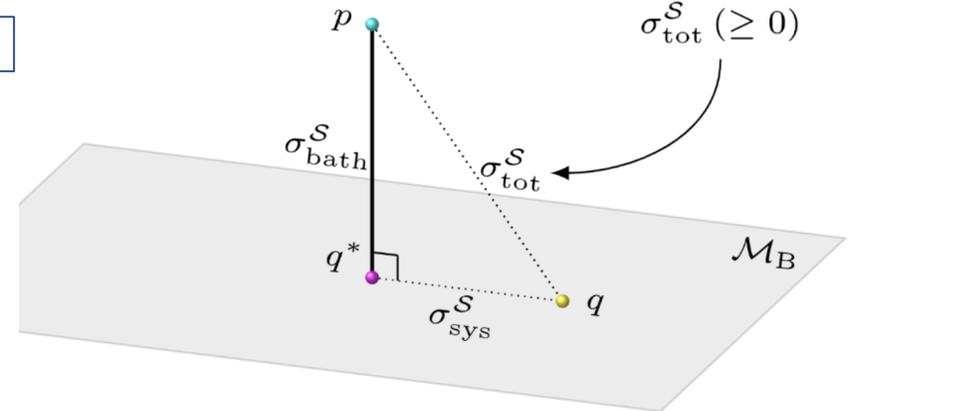
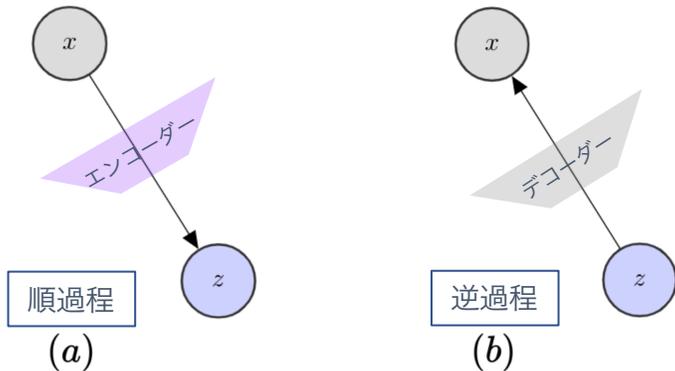


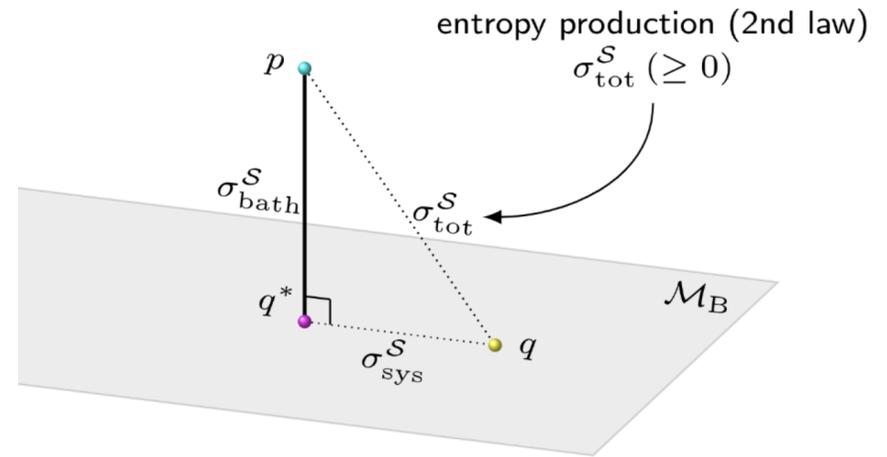
Figure 2: Total entropy production.

表現学習の熱力学 (10/n)

- 生成モデルの単純なモデル：変分オートエンコーダー (VAE) を考える。



$$p(\mathbf{s}) := T(\mathbf{z}, \mathbf{x}) p(\mathbf{x}) = p(\mathbf{z}|\mathbf{x}) p(\mathbf{x}), \quad q(\mathbf{s}) := T(\mathbf{x}, \mathbf{z}) q(\mathbf{z}) = q(\mathbf{x}|\mathbf{z}) q(\mathbf{z})$$



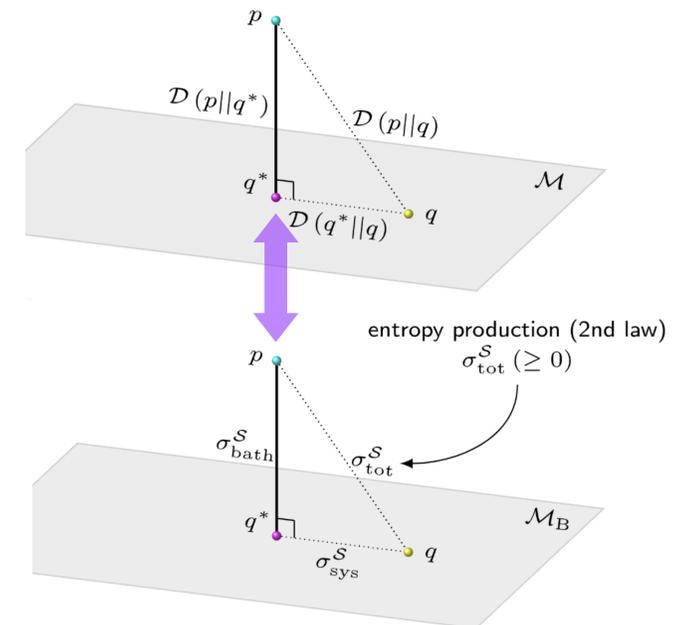
$$\mathcal{M}_B = \{q(\mathbf{s}) \mid q(\mathbf{s}) := T(\mathbf{x}, \mathbf{z}) q(\mathbf{z}) = q(\mathbf{x}|\mathbf{z}) q(\mathbf{z})\}$$

$$\begin{aligned} \mathcal{D}(p(\mathbf{s})||q(\mathbf{s})) &= \int p(\mathbf{s}) \log \frac{p(\mathbf{s})}{q(\mathbf{s})} d\mathbf{s} \\ &= \int p(\mathbf{s}) \log p(\mathbf{x}) d\mathbf{s} + \int p(\mathbf{s}) \log \frac{p(\mathbf{z}|\mathbf{x})}{q(\mathbf{z})} d\mathbf{s} - \int p(\mathbf{s}) \log q(\mathbf{x}|\mathbf{z}) d\mathbf{s} \\ &= \underbrace{-H(\mathbf{x}) + R + D}_{\mathcal{L}_{\text{ELBO}}} \geq 0 \end{aligned}$$

VAEの目的関数 変分下限 (ELBO)

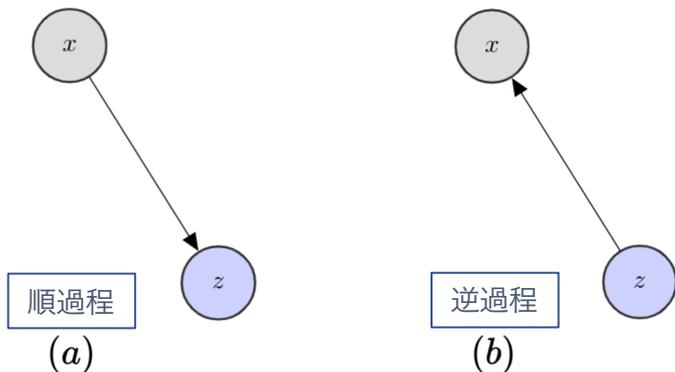
$$\begin{aligned} \mathcal{D}(p(\mathbf{s})||q(\mathbf{s})) &= \int p(\mathbf{s}) \log \frac{p(\mathbf{s})}{q(\mathbf{s})} d\mathbf{s} \\ &= \int p(\mathbf{s}) \log \frac{p(\mathbf{z}|\mathbf{x})p(\mathbf{x})}{q(\mathbf{x}|\mathbf{z})p(\mathbf{z})} d\mathbf{s} + \int p(\mathbf{s}) \log \frac{p(\mathbf{z})}{q(\mathbf{z})} d\mathbf{s} \\ &:= \sigma_{\text{bath}}^S + \sigma_{\text{sys}}^S = \sigma_{\text{tot}}^S \geq 0 \end{aligned}$$

全系のエントロピー生成率

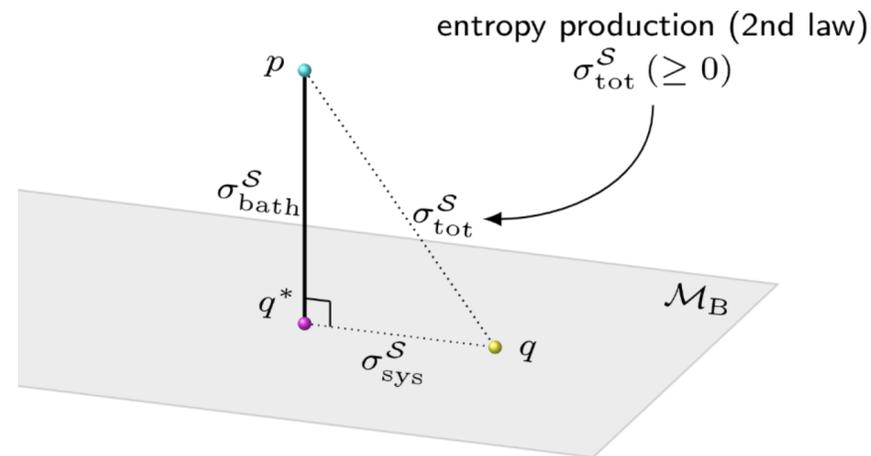


表現学習の熱力学 (11/n)

- 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。

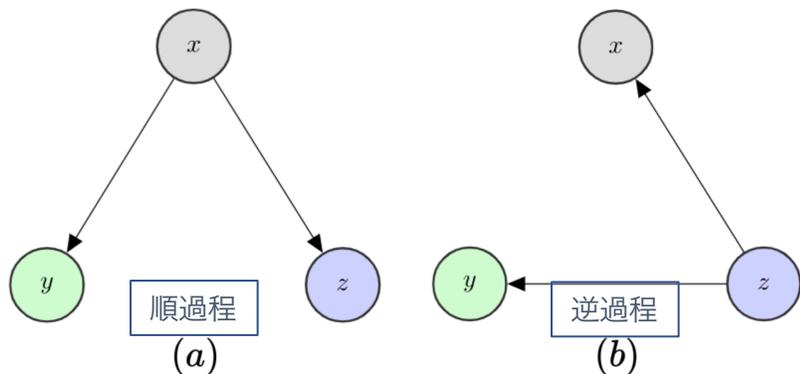


$$p(\mathbf{s}) := T(\mathbf{z}, \mathbf{x}) p(\mathbf{x}) = p(\mathbf{z}|\mathbf{x}) p(\mathbf{x}), \quad q(\mathbf{s}) := T(\mathbf{x}, \mathbf{z}) q(\mathbf{z}) = q(\mathbf{x}|\mathbf{z}) q(\mathbf{z})$$

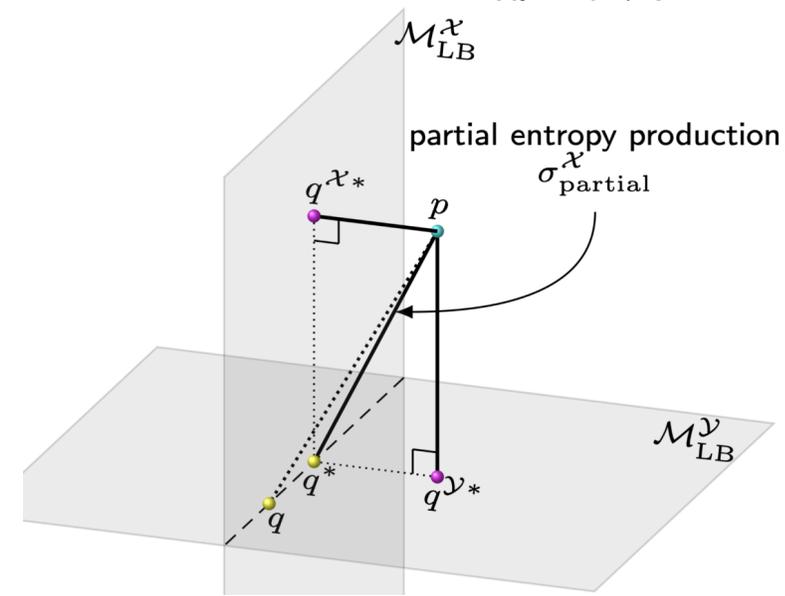


$$\mathcal{M}_B = \{q(\mathbf{s}) \mid q(\mathbf{s}) := T(\mathbf{x}, \mathbf{z}) q(\mathbf{z}) = q(\mathbf{x}|\mathbf{z}) q(\mathbf{z})\}$$

- 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。



$$p(\mathbf{s}) := p(\mathbf{z}|\mathbf{x}) p(\mathbf{y}|\mathbf{x}) p(\mathbf{x}), \quad q(\mathbf{s}) := q(\mathbf{x}|\mathbf{z}) q(\mathbf{y}|\mathbf{z}) q(\mathbf{z})$$



表現学習の熱力学 (12/n)

- 生成モデルの学習プロセス：データ \mathbf{x} から潜在因子（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，潜在因子（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。

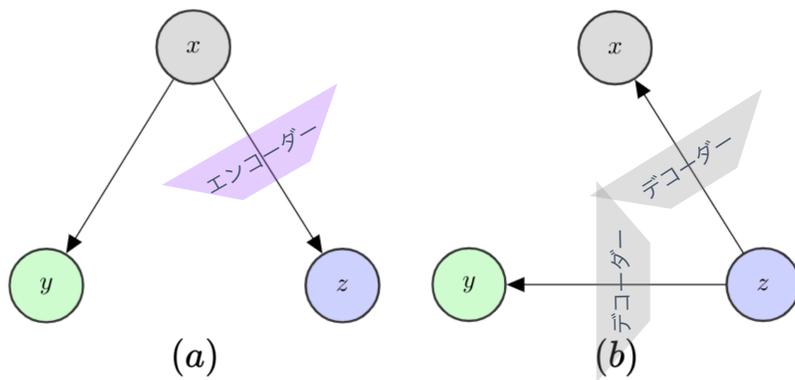
⇒ 部分系のエントロピー生成と解釈可能。

$$H := \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p(\mathbf{x}, \mathbf{y})} [-\log p(\mathbf{x}, \mathbf{y})], \quad \boxed{\text{データのエントロピー}}$$

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

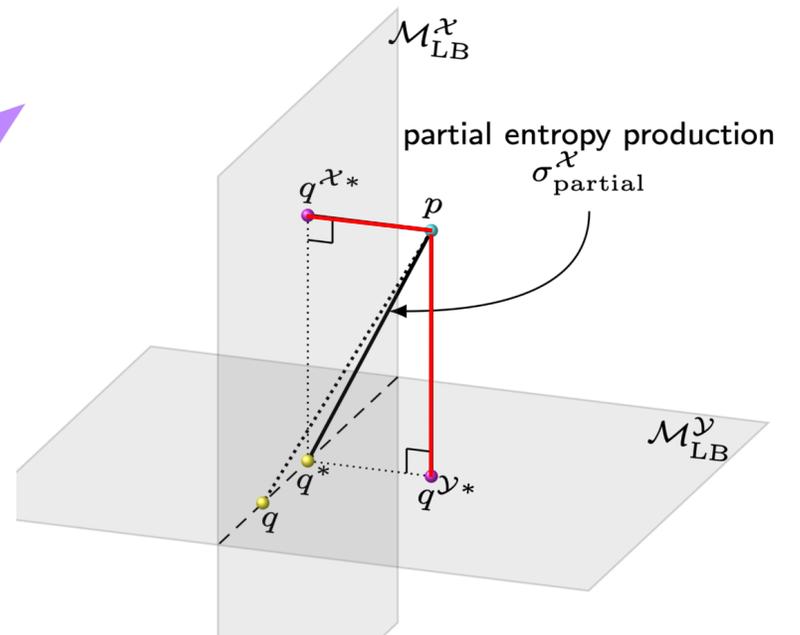
$$R := \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \left[\int e(\mathbf{z}|\mathbf{x}) \log \frac{e(\mathbf{z}|\mathbf{x})}{q(\mathbf{z})} d\mathbf{z} \right], \quad \boxed{\text{正則化項}}$$

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



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

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



表現学習の熱力学 (13/n)

- 表現学習における熱力学諸法則の導出

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

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

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

データのエントロピー

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

再構成誤差

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

正則化項

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

分類誤差

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

$$:= -H + \mathcal{J}(\theta, \lambda, \gamma) \geq 0$$

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

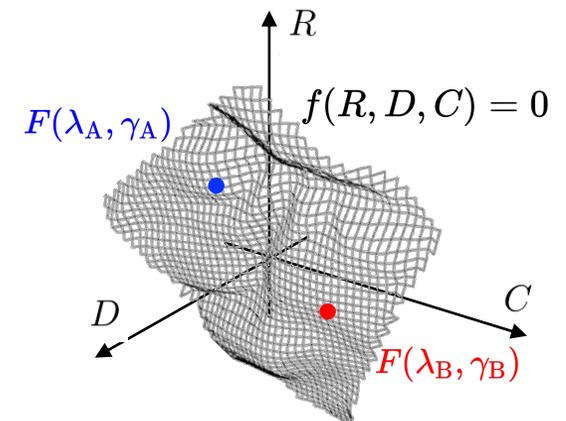
$$\nabla_\theta \mathcal{J}(\theta^*, \lambda, \gamma) = 0$$

$$f(R, D, C) = 0$$

$$F(\lambda, \gamma) = \min_\theta \mathcal{J}(\theta, \lambda, \gamma)$$

$$= \min_\theta R + \lambda D + \gamma C$$

$$F(\lambda, \gamma) = \mathcal{J}(\theta^*, \lambda, \gamma), \quad \theta^* \in \Theta$$



- 自由エネルギーが定義される最適曲面の上では, 以下が成り立つ.

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

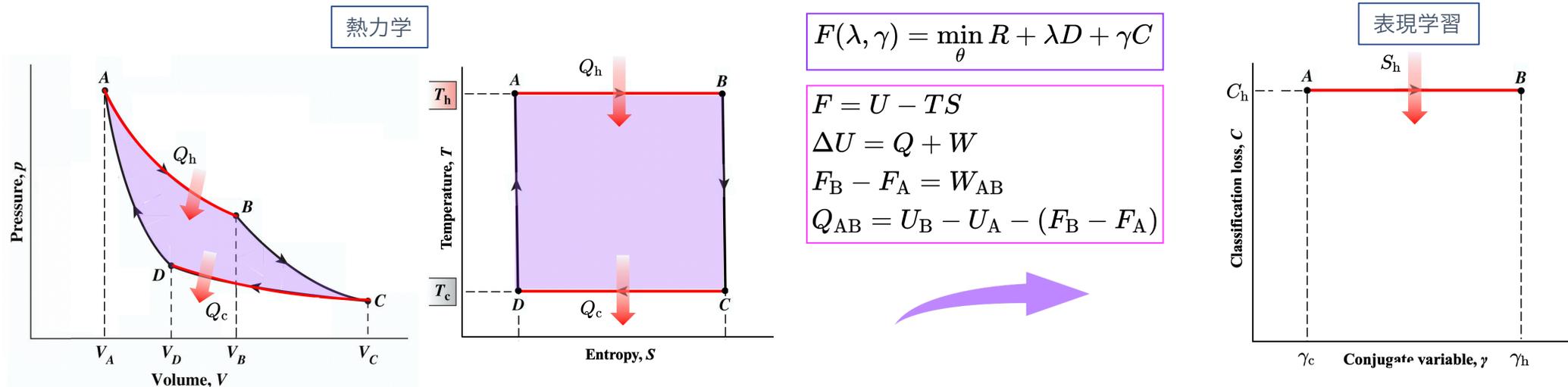
$$dR = -\lambda dD - \gamma dC \quad \lambda = -\left(\frac{\partial R}{\partial D}\right)_C \quad \gamma = -\left(\frac{\partial R}{\partial C}\right)_D \quad D = \left(\frac{\partial F}{\partial \lambda}\right)_\gamma \quad C = \left(\frac{\partial F}{\partial \gamma}\right)_\lambda$$

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

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

表現学習の熱力学 (14/n)

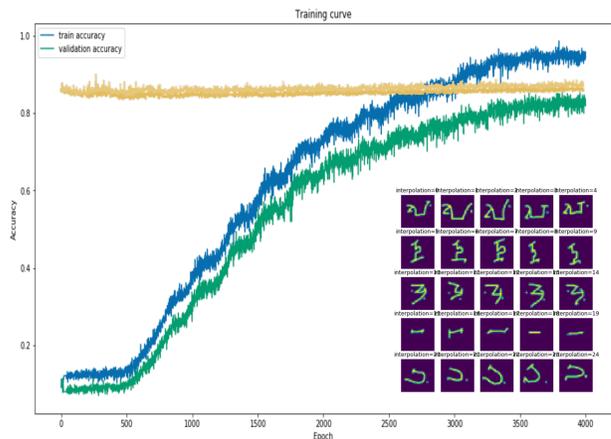
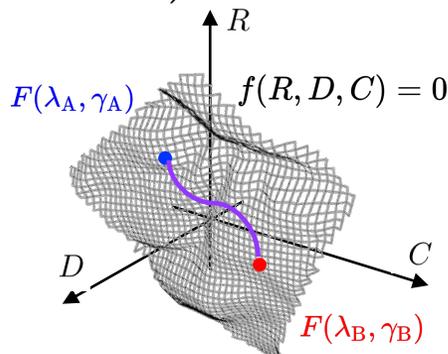
- 等温準静過程のアナロジーを考えると、統計モデルの分類誤差を低い値に抑えたままデータのドメインを変化させること (= 転移学習) はできないか？



- データのドメインを変化させる間、分類誤差の値が一定となるような“熱力学的操作”を考える。

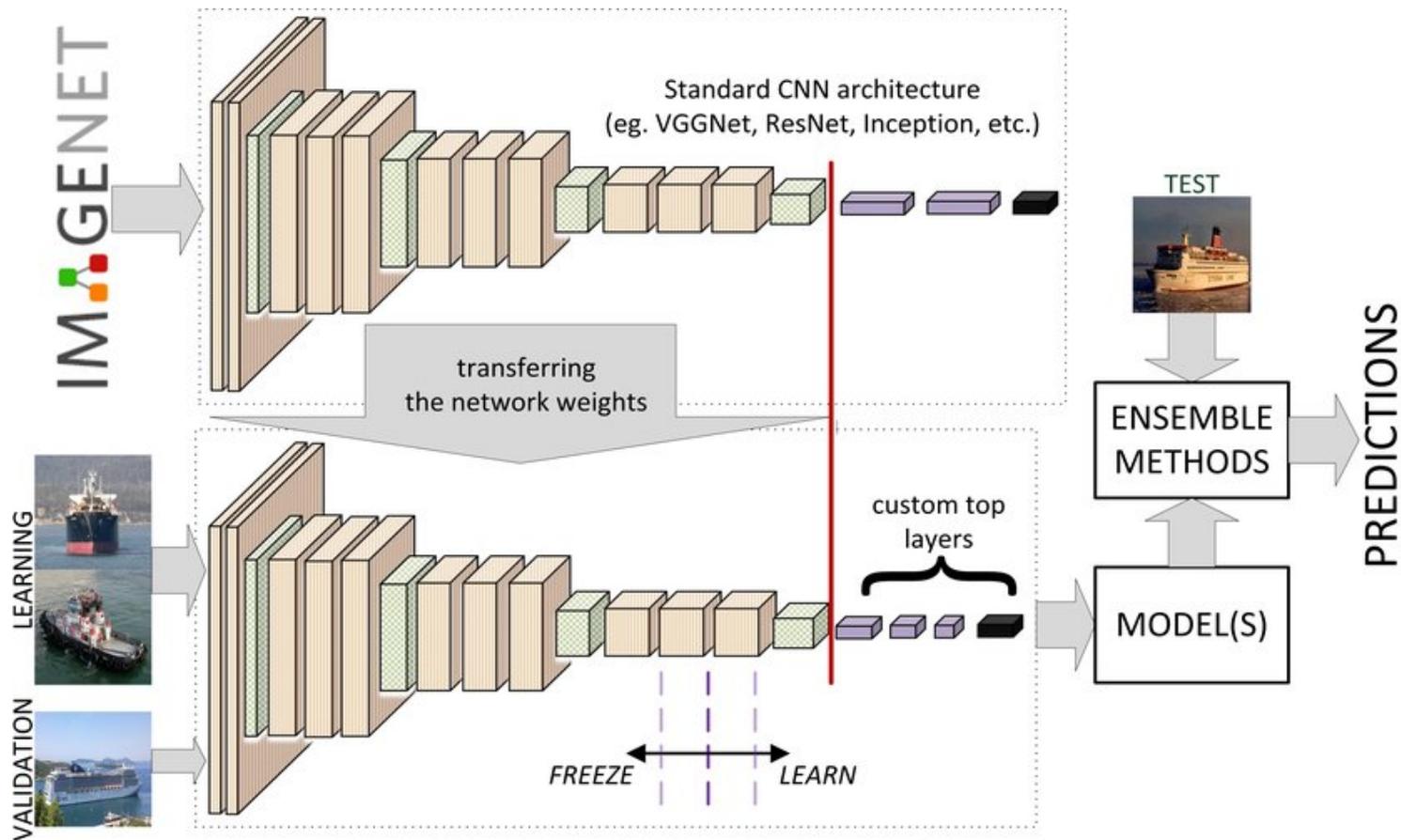
$$p(x(t)) = \operatorname{argmin}_{p(x)} (1-t)W_2^2(p(x^A), p(x)) + tW_2^2(p(x), p(x^B))$$

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



転移学習とはなにか (15/n)

特定のドメインで訓練が終了した学習済みモデルを、別のドメインの学習に再利用すること。少ない訓練データで高い分類性能を引き出すアルゴリズムを考えたりすることができる。

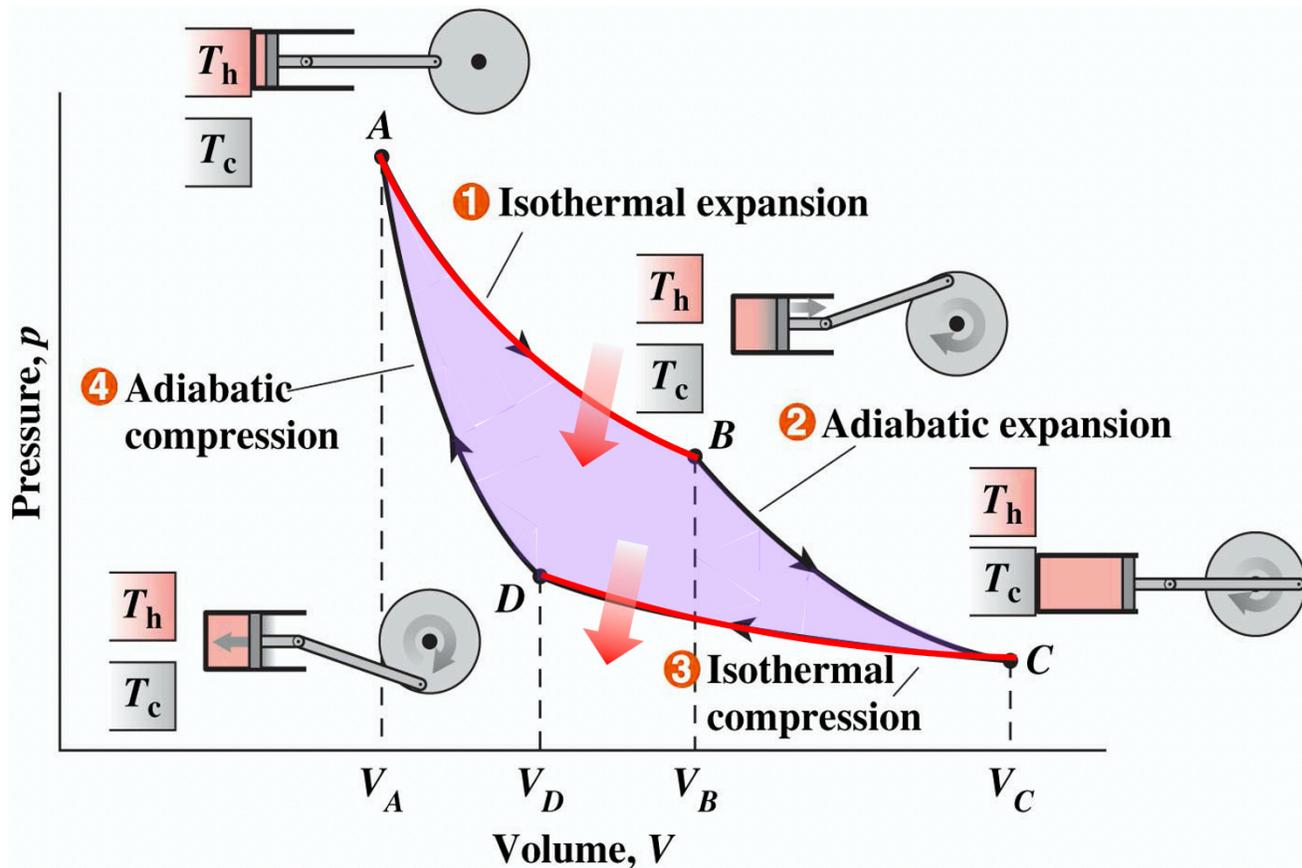


<https://medium.com/@subodh.malgonde/transfer-learning-using-tensorflow-52a4f6bcde3e>

学習済みモデルのパラメータを、別のドメインのデータの学習に再利用

熱力学エンジンの拘束条件 (16/n)

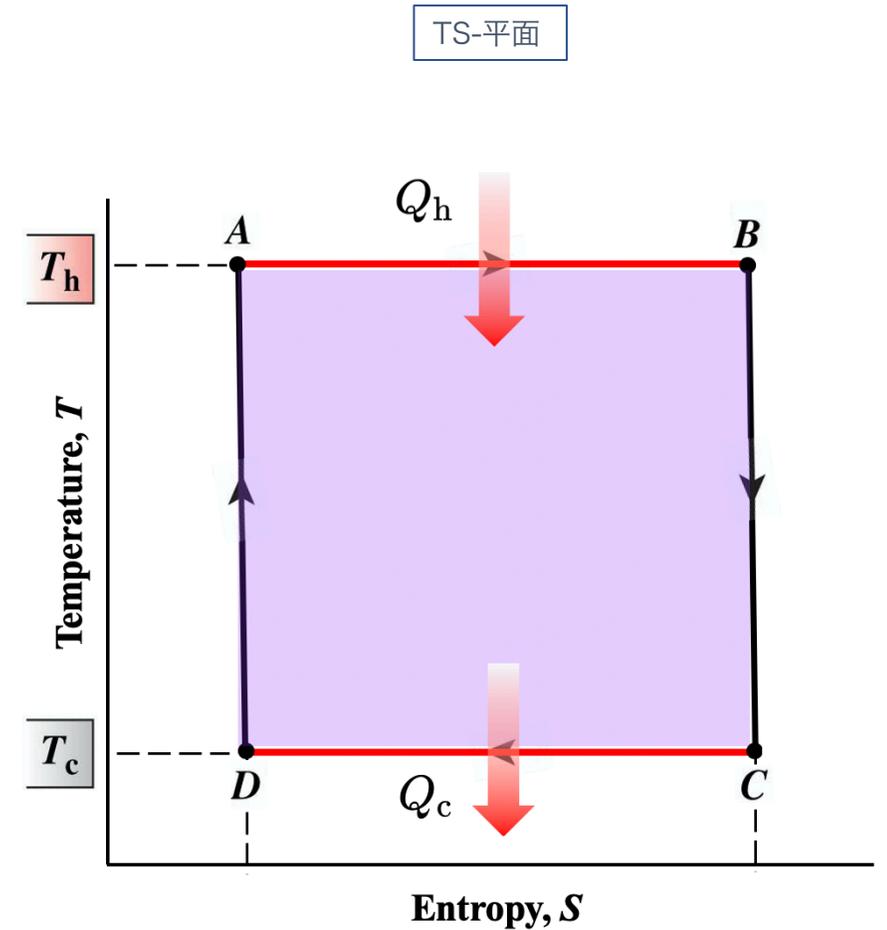
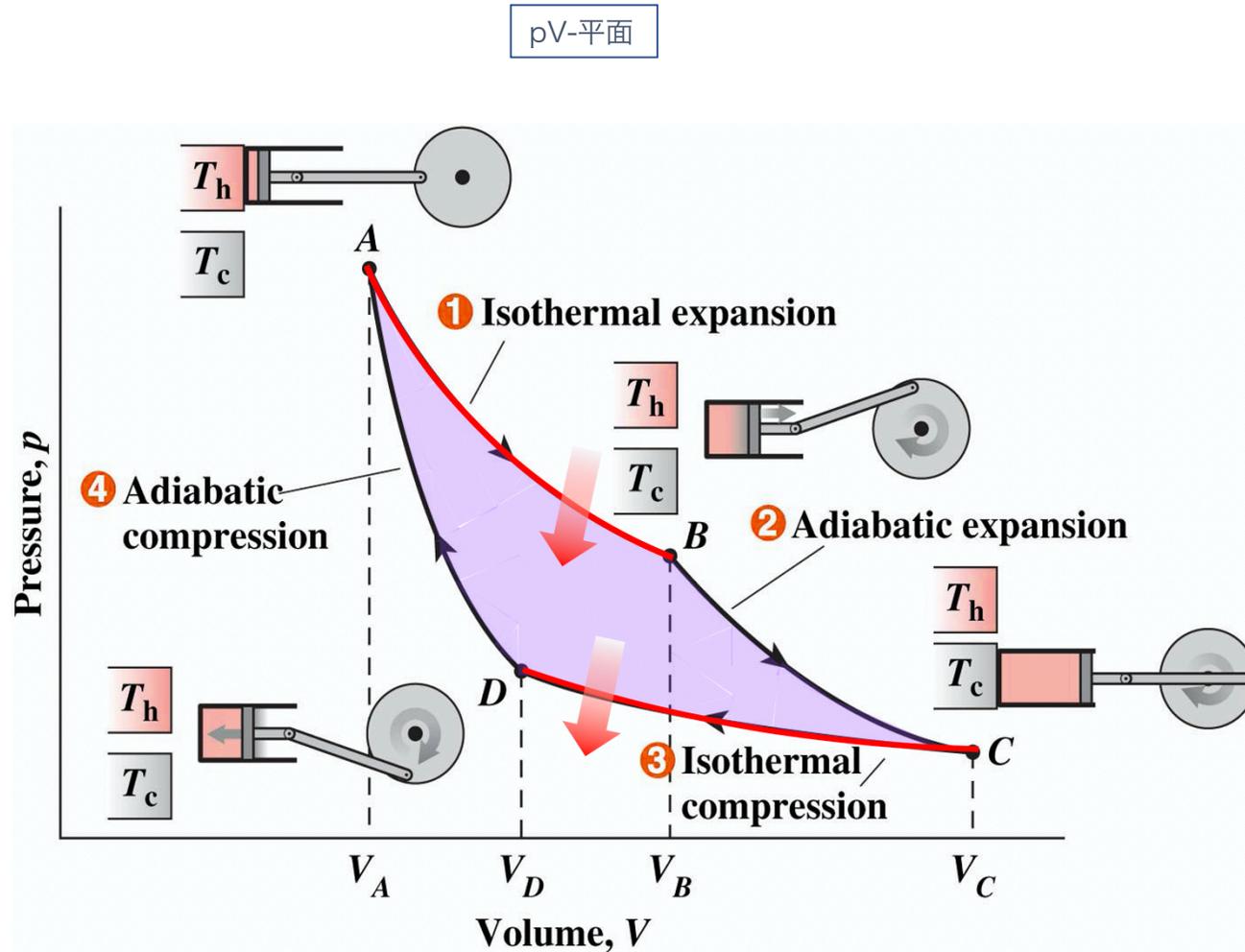
自由エネルギーや準静過程の意味は、熱力学エンジンの最大熱効率を求めるカルノーの定理を考えるとはっきりしてくる。カルノーサイクルに関する議論を振り返る。



- ① 等温準静過程 (膨張)
高温熱浴に接したまま膨張
- ② 断熱準静過程 (膨張)
エンジンを熱浴から離して断熱膨張
- ③ 等温準静過程 (圧縮)
低温熱浴に接したまま圧縮
- ④ 断熱準静過程 (圧縮)
エンジンを熱浴から離して断熱圧縮

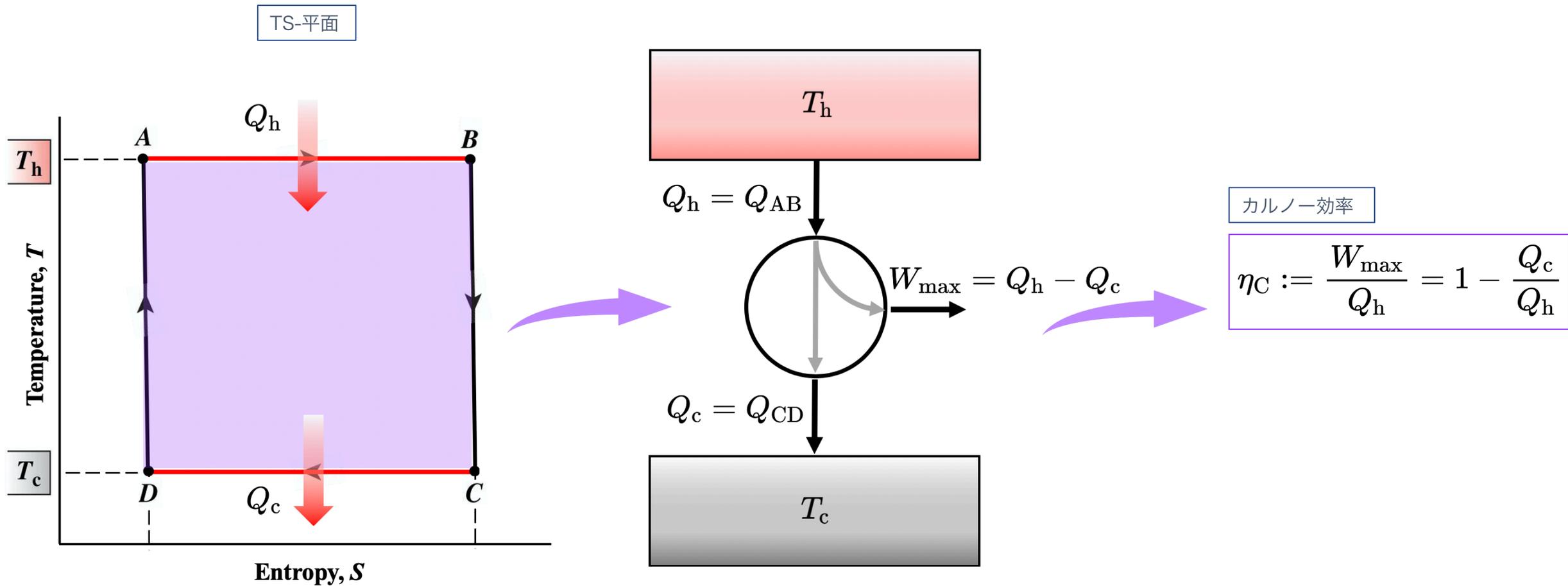
熱力学エンジンの拘束条件 (17/n)

- pV-平面上に描かれた準静サイクルを, TS-平面上に描き直すと平行四辺形があらわれる.



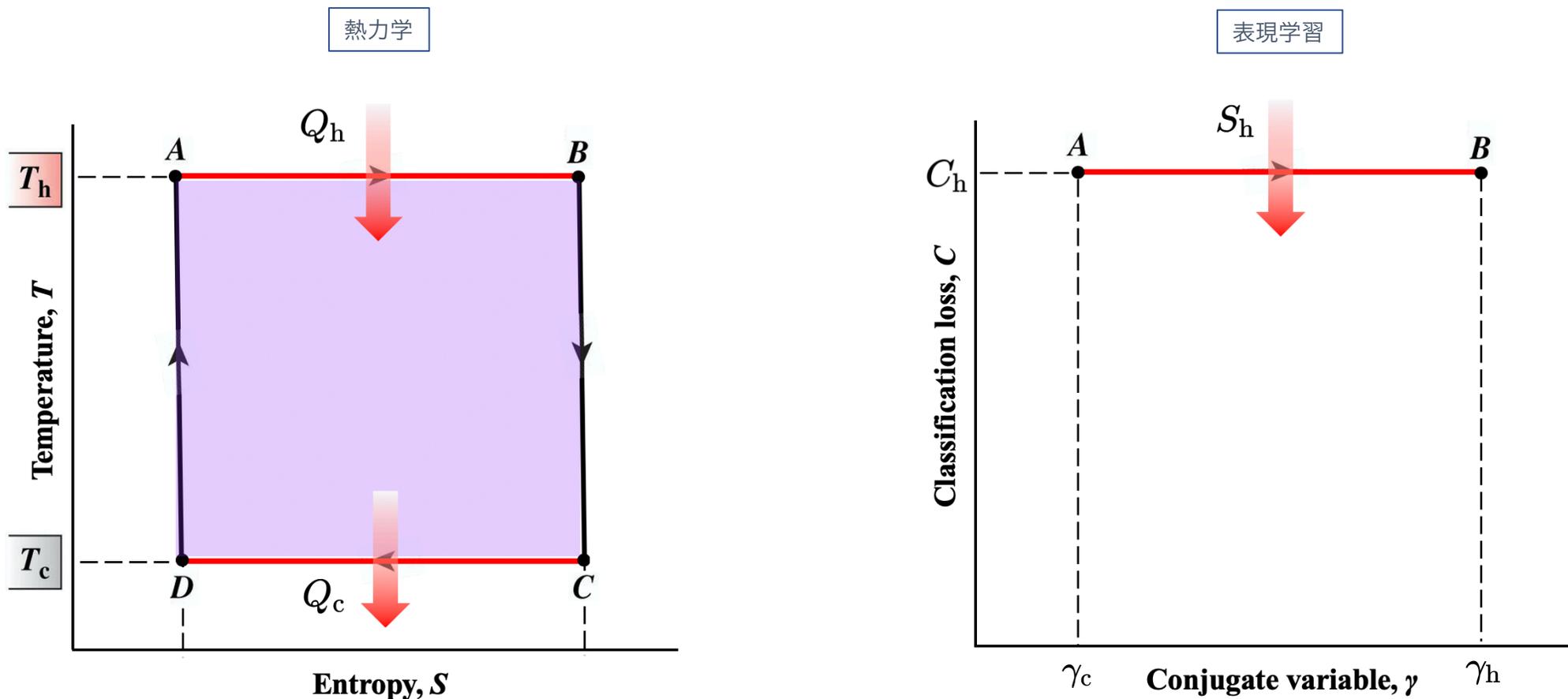
熱力学エンジンの拘束条件 (18/n)

- カルノーの定理より, 準静過程を組み合わせた準静サイクルが最大熱効率 (カルノー効率) を与えることが示される.



熱力学エンジンの拘束条件 (19/n)

- 等温準静過程のアナロジーを考えることで、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 (= 転移学習) はできないか?



表現学習の熱力学 (20/n)

- 分類誤差を“温度”にみたてて, 等温準静過程に沿ってモデルのパラメータ等を動かしていく.

$$p(x(t)) = \underset{p(x)}{\operatorname{argmin}} (1-t)W_2^2(p(x^A), p(x)) + tW_2^2(p(x), p(x^B))$$

要請①

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

$$\frac{dC(t)}{dt} = \dot{\theta}^\top \nabla_{\theta} C + \dot{\lambda} \frac{\partial C}{\partial \lambda} + \dot{\gamma} \frac{\partial C}{\partial \gamma} = 0$$

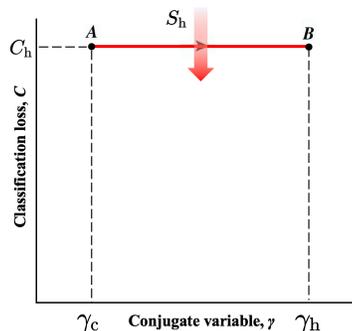
要請②

$$\nabla_{\theta} \mathcal{J}(\theta^*(t), \lambda(t), \gamma(t)) = 0$$

$$\frac{d}{dt} \nabla_{\theta} \mathcal{J}(\theta^*, \lambda, \gamma) = \dot{\theta}^\top \nabla_{\theta}^2 \mathcal{J}_{\theta^*} + \dot{\lambda} \frac{\partial}{\partial \lambda} \nabla_{\theta} \mathcal{J}_{\theta^*} + \dot{\gamma} \frac{\partial}{\partial \gamma} \nabla_{\theta} \mathcal{J}_{\theta^*} = 0$$

要請③

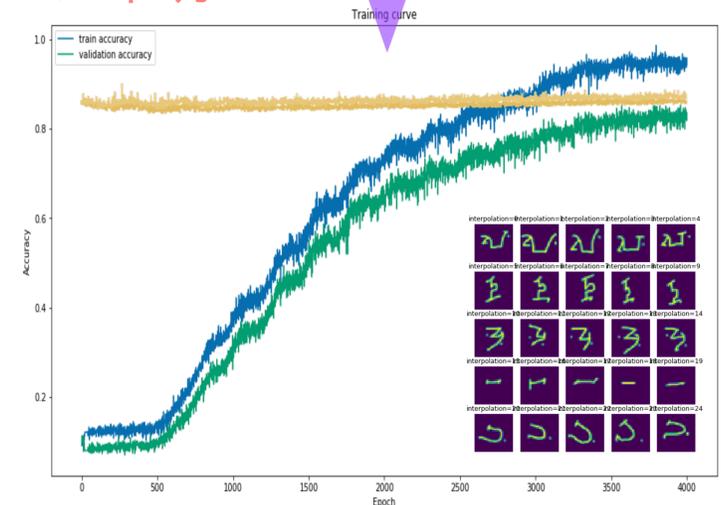
$$\dot{\lambda} = k$$



$$\begin{aligned} \dot{\theta}^\top \nabla_{\theta} C + \dot{\lambda} \frac{\partial C}{\partial \lambda} + \dot{\gamma} \frac{\partial C}{\partial \gamma} &= 0 \\ \dot{\theta}^\top \nabla_{\theta}^2 \mathcal{J} + \dot{\lambda} \frac{\partial}{\partial \lambda} \nabla_{\theta} \mathcal{J} + \dot{\gamma} \frac{\partial}{\partial \gamma} \nabla_{\theta} \mathcal{J} &= 0 \\ \dot{\lambda} &= k \end{aligned}$$

$(\theta, \lambda, \gamma)$ についての連立一次微分方程式

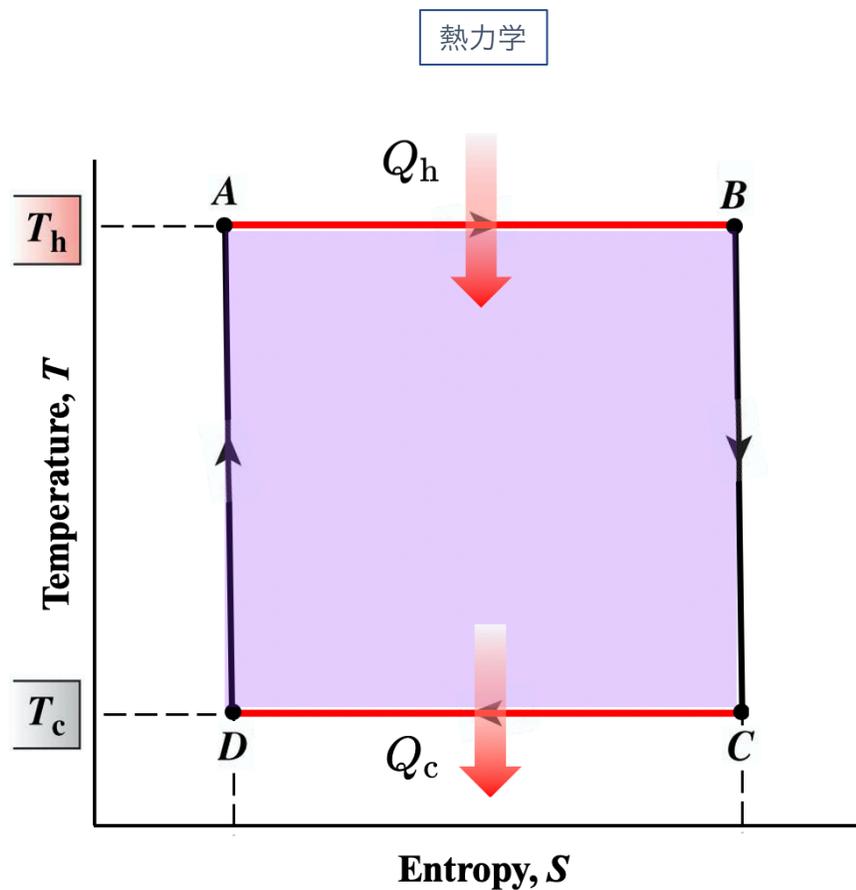
どういふときに転移学習が成功して
どういふときには失敗するのかわ不明



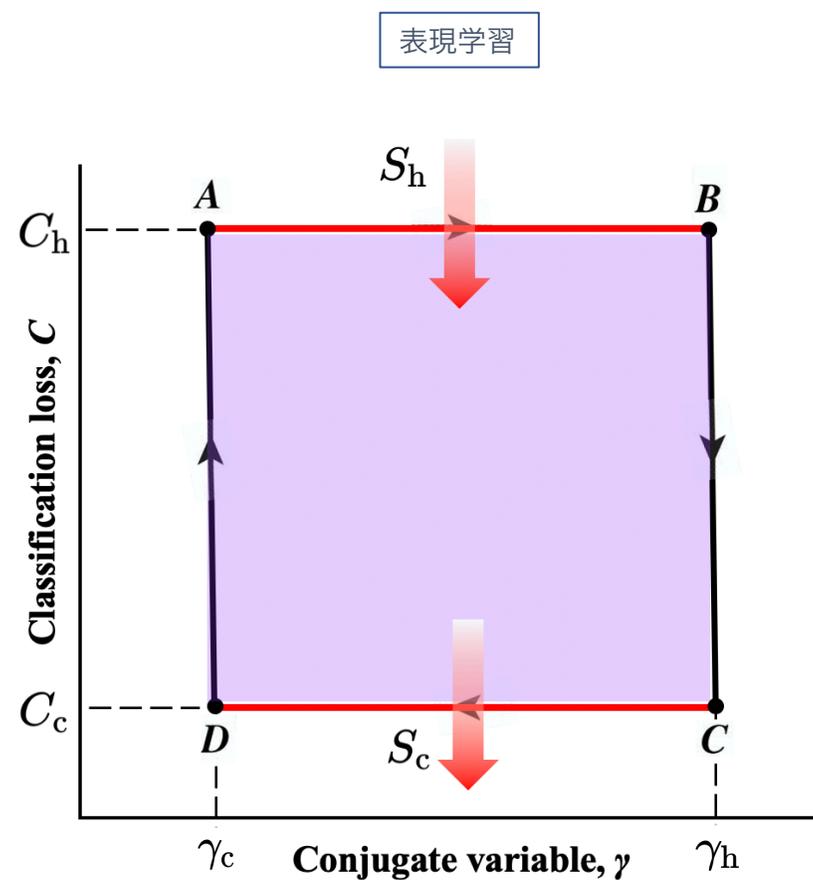
MNIST の 0~4 までの数字 (ドメイン) について訓練
⇒ 残った 5~9 までの数字 (ドメイン) について転移学習

表現学習の熱力学 (21/n)

- 等温準静過程に加え, 断熱準静過程に対応する熱力学的操作も導入することができれば表現学習版のカルノーサイクルを考えることができる.



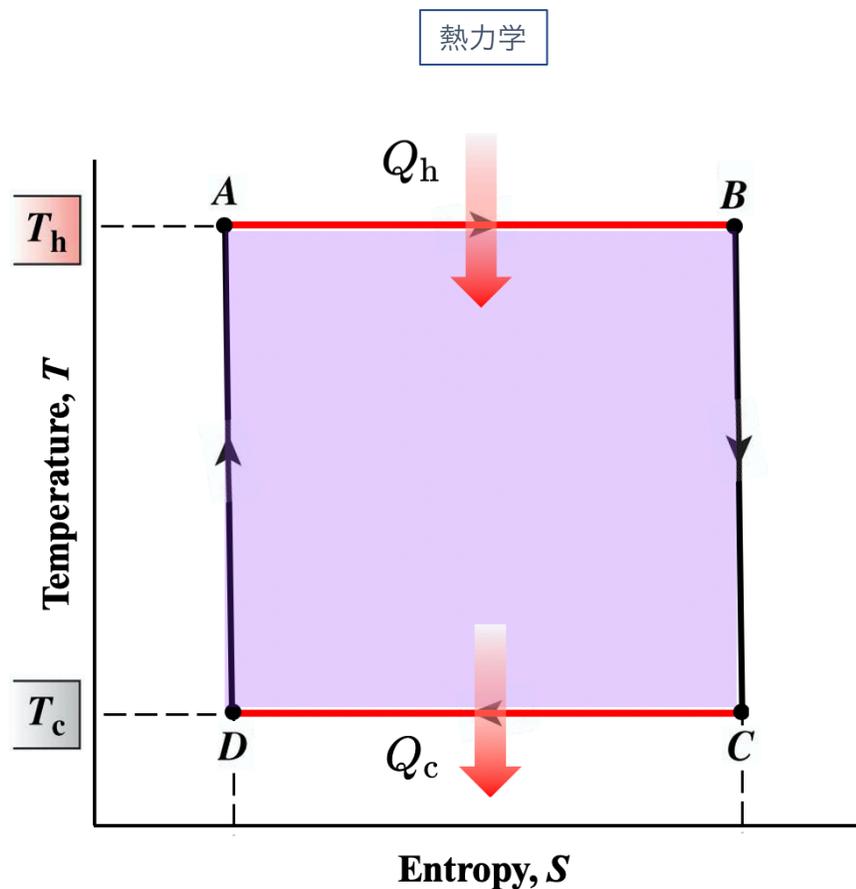
$$\eta_C := \frac{W_{\max}}{Q_h} = 1 - \frac{Q_c}{Q_h}$$



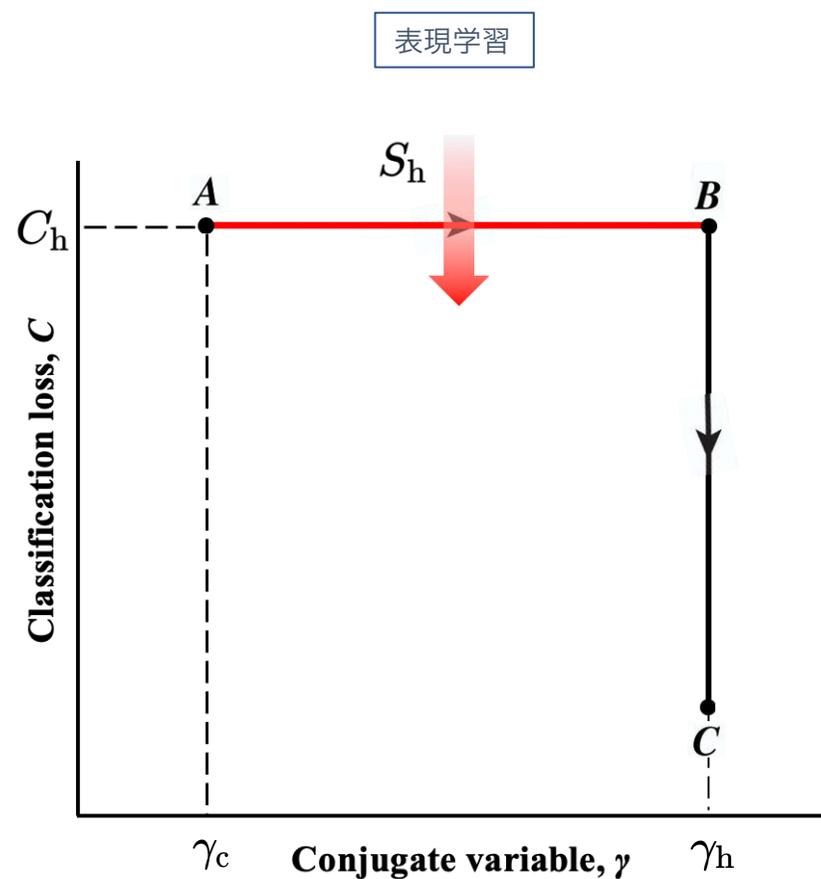
$$\eta_C := \frac{W_{\max}}{S_h} = 1 - \frac{S_c}{S_h} \quad ?$$

表現学習の熱力学 (22/n)

- 表現学習版のカルノーサイクルを考えるために, 断熱準静過程に対応する熱力学的操作を導入したい.



$$\eta_C := \frac{W_{\max}}{Q_h} = 1 - \frac{Q_c}{Q_h}$$



$$\eta_C := \frac{W_{\max}}{S_h} = 1 - \frac{S_c}{S_h} \quad ?$$

表現学習の熱力学 (23/n)

- 断熱準静過程のアナロジーを表現学習に持ち込む.

$$p(x(t)) = \underset{p(x)}{\operatorname{argmin}} (1-t)W_2^2(p(x^B), p(x)) + tW_2^2(p(x), p(x^C))$$

要請①

$$\gamma = - \left(\frac{\partial R}{\partial C} \right)_D$$

$$\dot{\gamma} = \alpha \frac{\partial C}{\partial \lambda} = \alpha \frac{\partial^2 F}{\partial \lambda \partial \gamma} = \alpha \frac{\partial^2 \mathcal{J}_{\theta^*}}{\partial \lambda \partial \gamma} = 0$$

要請②

$$\nabla_{\theta} \mathcal{J}(\theta^*(t), \lambda(t), \gamma(t)) = 0$$

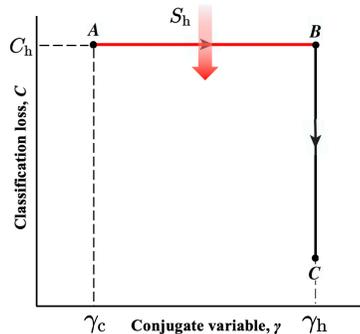
$$\frac{d}{dt} \nabla_{\theta} \mathcal{J}(\theta^*, \lambda, \gamma) = \dot{\theta}^{\top} \nabla_{\theta}^2 \mathcal{J}_{\theta^*} + \dot{\lambda} \frac{\partial}{\partial \lambda} \nabla_{\theta} \mathcal{J}_{\theta^*} + \dot{\gamma} \frac{\partial}{\partial \gamma} \nabla_{\theta} \mathcal{J}_{\theta^*} = 0$$

要請③

$$\dot{\lambda} = k$$

要請④ (断熱条件)

$$S_{BC} = R_C - R_B - (F_C - F_B) = 0$$



$$\dot{\gamma} = \alpha \frac{\partial^2 \mathcal{J}}{\partial \lambda \partial \gamma} = 0$$

$$\dot{\theta}^{\top} \nabla_{\theta}^2 \mathcal{J} + \dot{\lambda} \frac{\partial}{\partial \lambda} \nabla_{\theta} \mathcal{J} + \dot{\gamma} \frac{\partial}{\partial \gamma} \nabla_{\theta} \mathcal{J} = 0$$

$$\dot{\lambda} = k$$

$$S_{BC} = R_C - R_B - (F_C - F_B) = 0$$

$(\theta, \lambda, \gamma)$ についての連立一次確率微分方程式

$$S_{AB} = R_B - R_A - (F_B - F_A)$$

$$= - \int_{t_A}^{t_B} \lambda(t) D(t) + \gamma(t) C(t) dt$$

$$\int_{t_A}^{t_B} \gamma(t) C(t) dt = C \int_{t_A}^{t_B} \gamma(t) dt$$

$$p(x^A) \longrightarrow p(x^B)$$

表現学習の熱力学 (24/n)

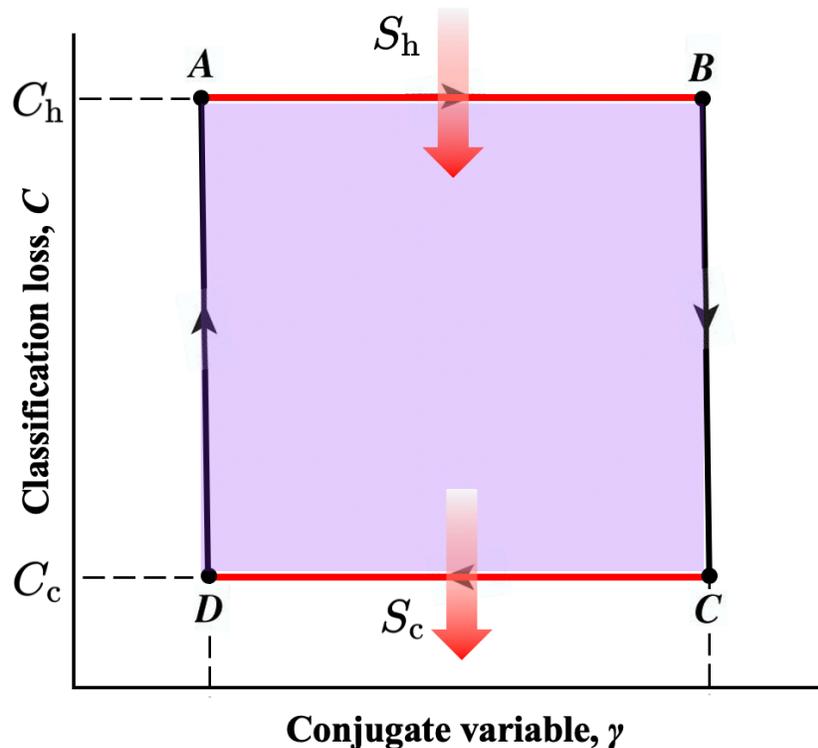
- 表現学習版の熱力学第1法則・第2法則から、カルノーサイクルが最大熱効率を与えることを確認。

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

$$dR = -\lambda dD - \gamma dC \quad \lambda = -\left(\frac{\partial R}{\partial D}\right)_C \quad \gamma = -\left(\frac{\partial R}{\partial C}\right)_D \quad D = \left(\frac{\partial F}{\partial \lambda}\right)_\gamma \quad C = \left(\frac{\partial F}{\partial \gamma}\right)_\lambda$$

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

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



最大熱効率 (カルノー効率)

$$\eta_C := \frac{W_{\max}}{S_h} = 1 - \frac{S_c}{S_h}$$

表現学習の熱力学 (25/n)

- まとめ

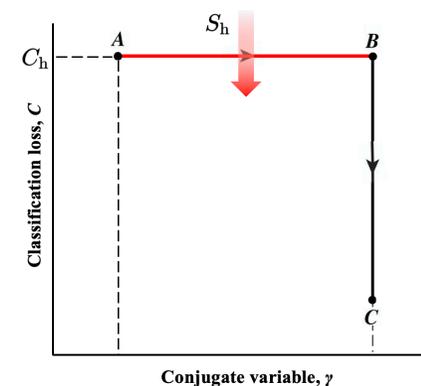
- エントロピー生成と生成モデルの損失関数を同一視することで、変分下限と自由エネルギーを対応させる関係式を発見.

$$\sigma_{\text{tot}}^S = \beta[W - \Delta F] \quad \sigma_{\text{tot}}^S = -H + \underbrace{R + \lambda D + \gamma C}_{\mathcal{L}_{\text{ELBO}}}$$

- “断熱準静過程”を表現学習に導入.

ドメインBからCへの転移

$$p(x(t)) = \underset{p(x)}{\operatorname{argmin}} (1-t)W_2^2(p(x^B), p(x)) + tW_2^2(p(x), p(x^C))$$



- 表現学習版のカルノーサイクルを考え、これが最大熱効率を与えることを確認.

$$\eta_C := \frac{W_{\text{max}}}{S_h} = 1 - \frac{S_c}{S_h}$$

Appendix (26/n)

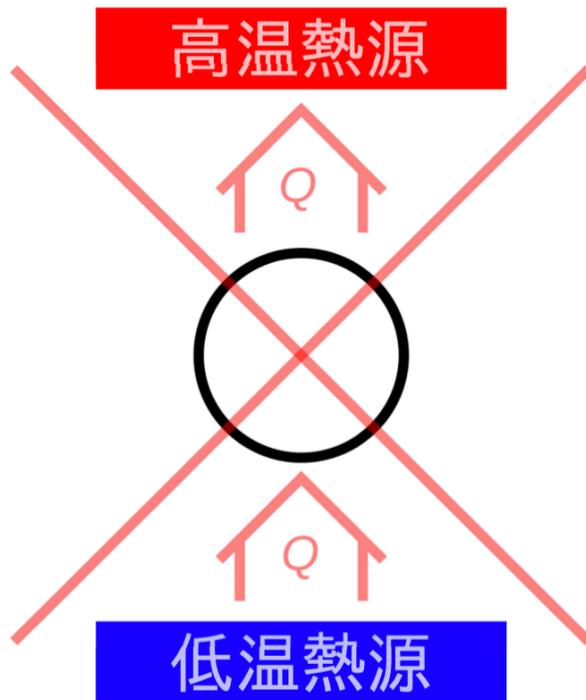


図 1.3 クラウジウスの原理

<https://camellia.net/study.html>

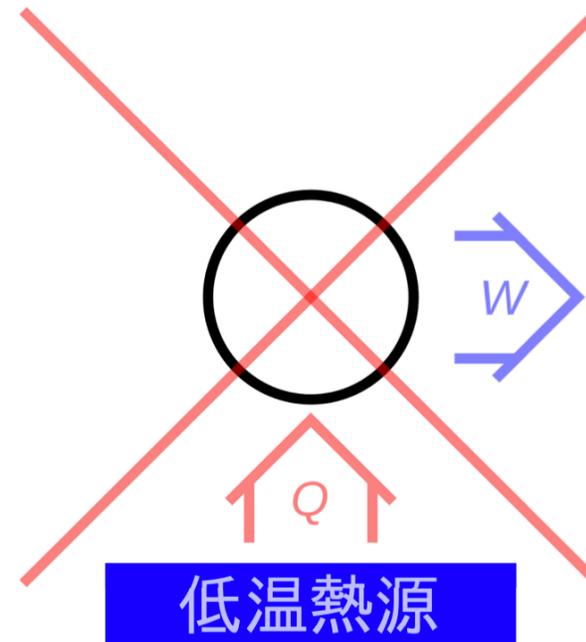


図 1.4 トムソンの原理