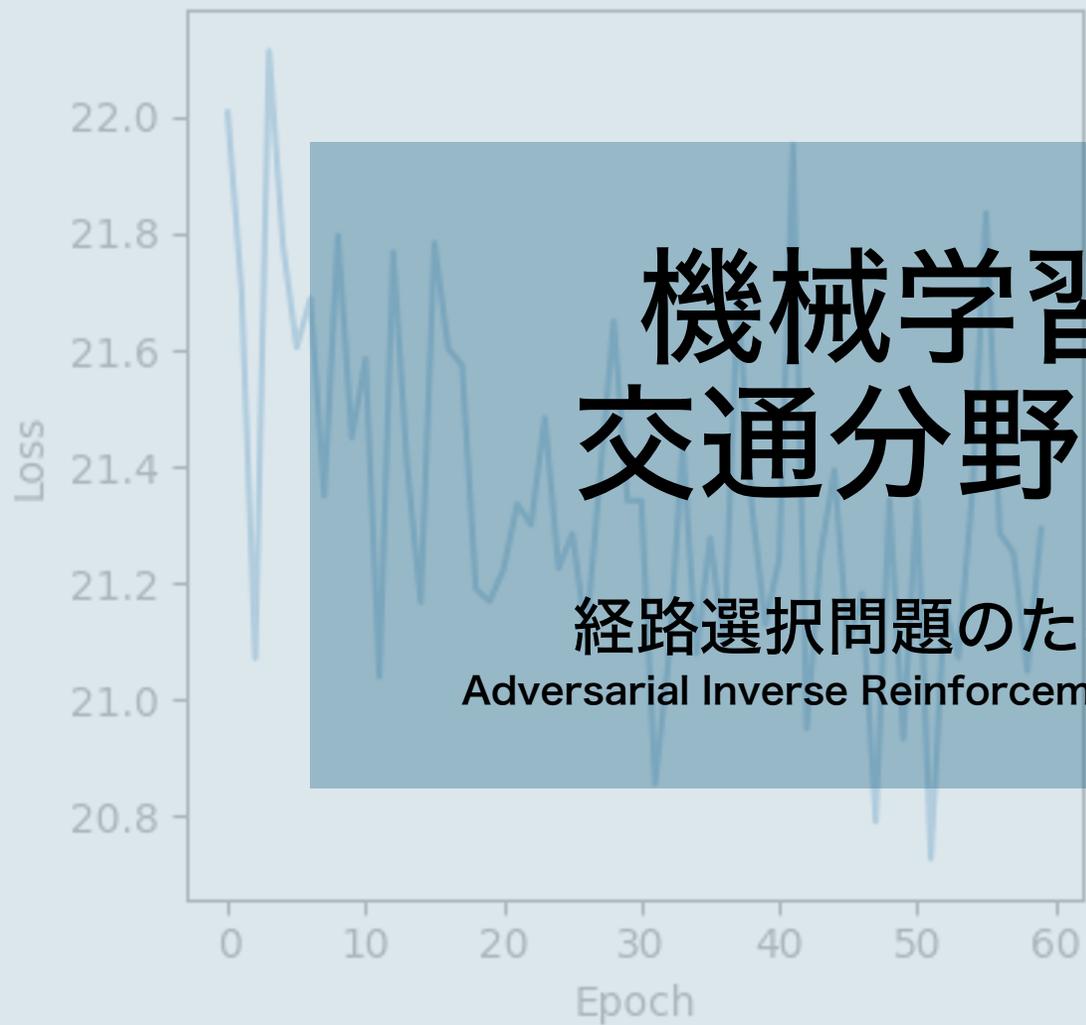
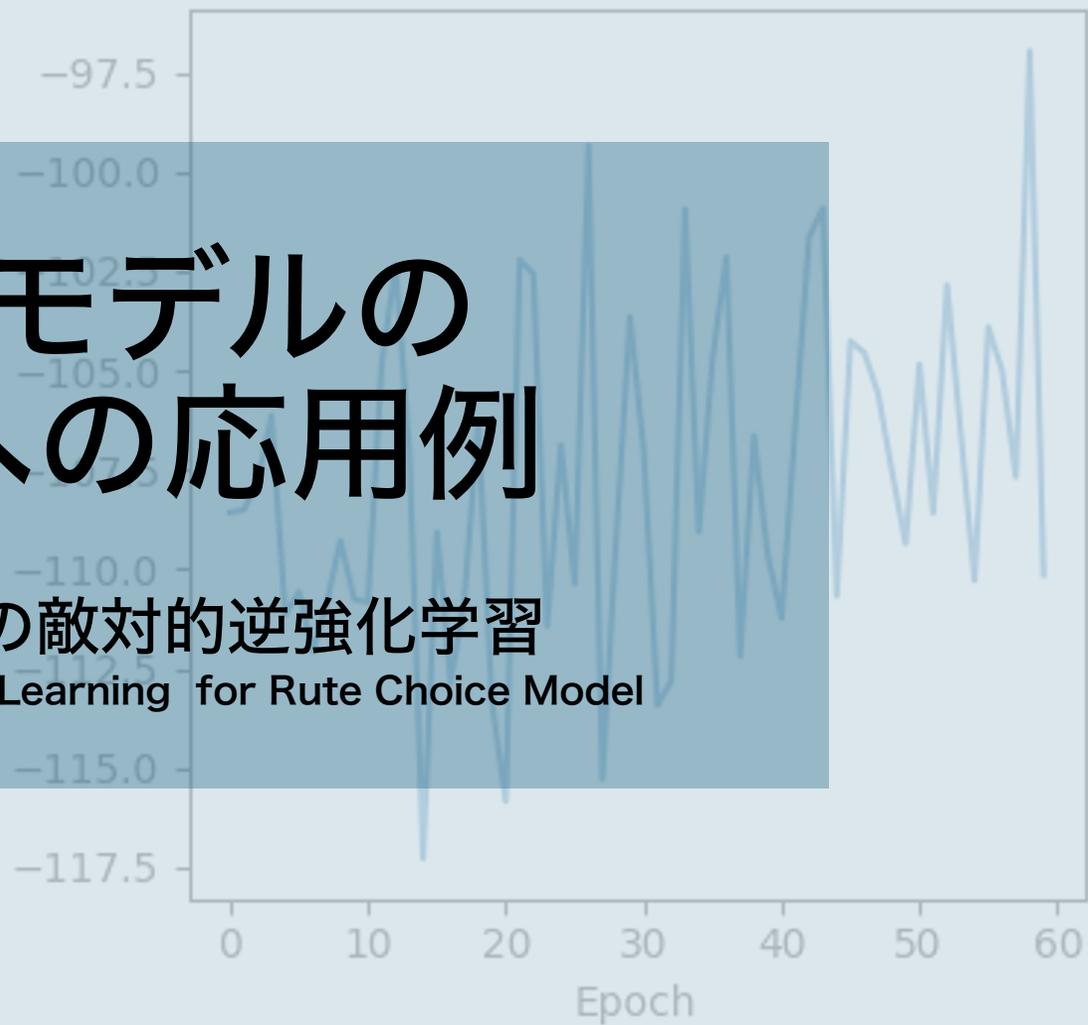


Ped Discriminator Loss



Ped Generator Loss



# 機械学習モデルの 交通分野への応用例

経路選択問題のための敵対的逆強化学習  
Adversarial Inverse Reinforcement Learning for Route Choice Model

# 目次

---

- AIRL for RCMの概要
- 定式化
- 逆強化学習とは
- 逆強化学習としての定式化
- AIRLとは
- Generator
- Discriminator
- 実装（機械学習の実装）
- AIRLの実装
- 実証実験の結果
- 機械学習の説明可能性(SHAP)



- **経路選択問題で機械学習をどう使うのか？**

How is the AIRL used in route choice problem?

- **データ処理から実装までの大体的流れ** The process of AIRL for RCM

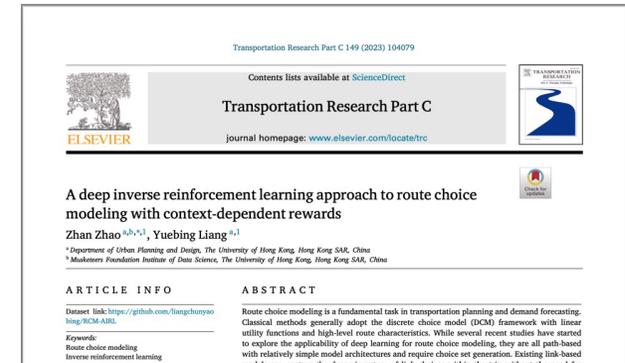
# AIRL for RCM

## Adversarial Inverse Reinforcement Learning for Route Choice Model

- AIRLという逆強化学習アルゴリズムを経路選択問題に適用した論文
- Zhan Zhao, Yuebing Liang, 2023 (←新しい！)

### ■ 従来モデルと比較したメリット

- 効用関数が線形に限定されず非線形も表せるため人間の複雑な嗜好を表現可能 → 精度が高い結果に shows better result  
実証分析（中国の広範囲ネットワークで、RLモデルや他の機械学習モデルと比較しても精度が最も高いことが示された）
- 深層学習：膨大なデータ量，データ種類の増加に強い  
→ 天気や個人属性などを考慮することが可能



中国のネットワーク図

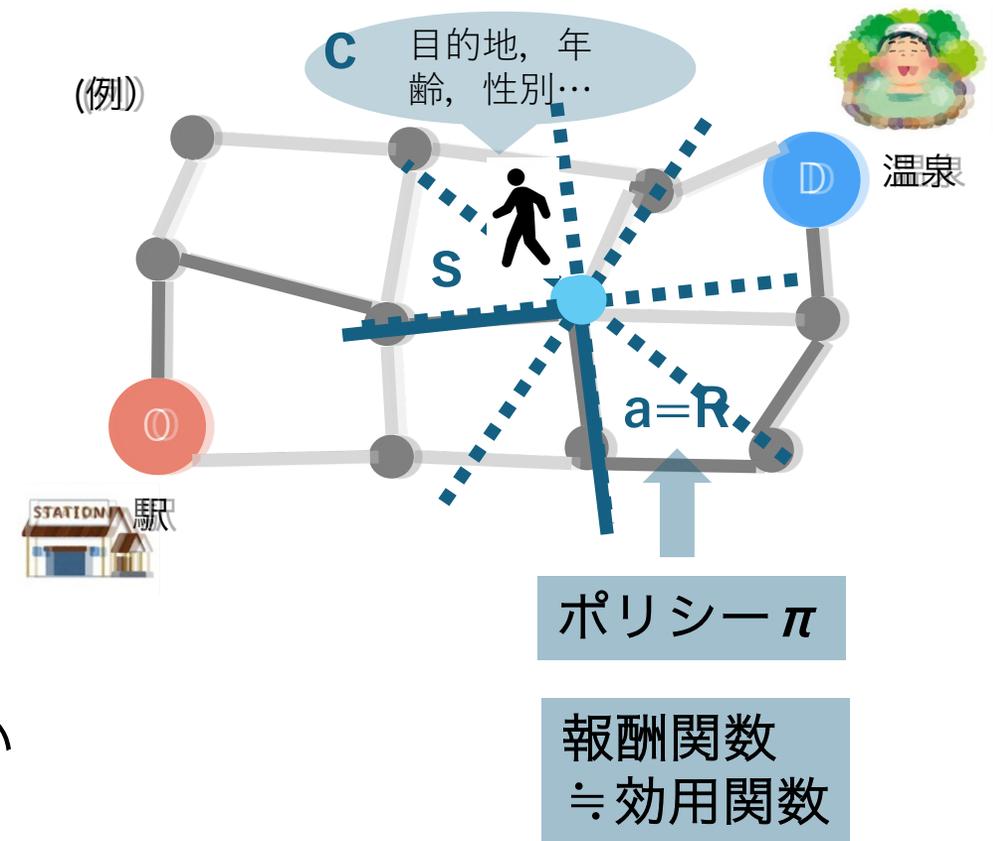


# 経路選択問題の定式化

経路選択問題をマルコフ決定過程(MDP)として考える

- **State(状態s)**: 道路ネットワークにおいて旅行者の現在の位置を示すリンク(current link)
- **Action(行動a)**: リンク間の移動選択 (8方向のいずれか)  
R,Fなど
- **Context(コンテキストc)**: トリップ中変化せず, 個人の経路選択に影響を及ぼす要素群ベクトル(consistent factor)  
(目的地や目的、個人属性、日時、天気など ex; distination)
- **Policy(ポリシー  $\pi$ )**: 状態sとコンテキストcが与えられた時に行動aを決定する関数( $\pi(a|s,c)$ と表記) ( $\pi^*$ は最適な選択行動を返す)
- **Reward function(報酬関数R)**: cが与えられた時、sの状態であの行動を選択した時の効用 ( $R(s,a|c)$ と表記)

←リンクベースの逐次的な経路選択という点ではRLの定式化と同じ!!



ポリシーより報酬関数の方がより根本的で汎用性が高い  
→報酬関数を求めたい (=逆強化学習)

# 逆強化学習とは？

補足

## 強化学習 Reinforcement Learning

予め定義された報酬関数を最大化する行動を生み出す決定プロセス(戦略, ポリシー)を学習することが目的

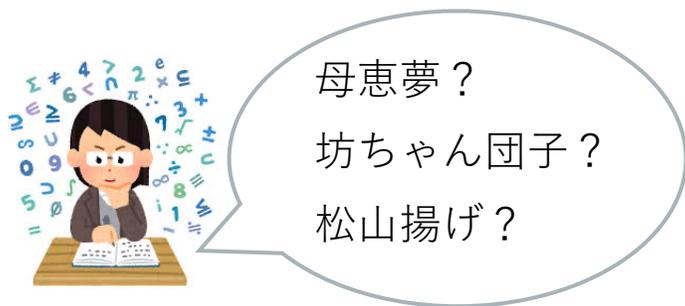


## 逆強化学習 Inverse Reinforcement Learning

データから、観察された人間の行動を説明する報酬関数を復元することが目的

(例)研究室にお土産買う時

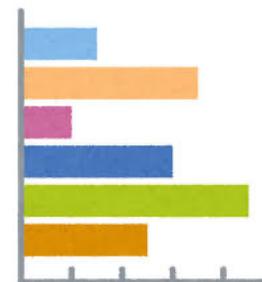
報酬関数の方が汎用性が高い



どのお土産が  
いいかな



見た目が豪華  
大きさがでかい  
数が少ない



# 逆強化学習として定式化

■ 目標：行動データXから、報酬関数Rを推定する Goal ; To recover Reward function from observed behavior

• 行動データXの定義式: 行動データXをs,a,cペアに分解

$$\underline{x}_i = \{(s_1^{(i)}, a_1^{(i)}), (s_2^{(i)}, a_2^{(i)}), \dots, (s_{T_i}^{(i)}, a_{T_i}^{(i)}) \mid c^{(i)}\},$$

←エージェントiの0からDまでをリンク遷移(s,a)ペアに分割  
 ↑i人目の軌跡 ↑tステップ目



• 報酬関数Rの定義式 パラメータ  $\theta$  関数として定義し、 $\theta$  を学習させることで、Rの関数を形にする  
 RLの効用関数と同じ！

$$R_\theta(x_i) = \sum_{t=1}^{T_i} \gamma^t R_\theta(s_t^{(i)}, a_t^{(i)} \mid c^{(i)}).$$

エージェントiの報酬関数 0からDまでの各移動ステップの報酬関数を足す



↓ 全エージェントの行動データから確率の立式  
 (観測されたx(軌跡)の起こる確率はRのexpに比例するという法則)

$$P_\theta(x) = \frac{1}{Z} \exp(R_\theta(x))$$

行動xの起こる確率 Zは全てのx (軌跡) の  $R_\theta(x)$  の積分和



↓ 全エージェントは最大の確率を選択して行動するものとする

$$\max_{\theta} \sum_{i=1}^N \log(P_\theta(x_i)).$$

→IRL問題は観測された軌跡に基づいて尤度を最大化する問題として組み立てられる

maximum likelihood problem

# 逆強化学習として定式化

## ■ 目標：行動データから、報酬関数Rを推定する

- 行動データXの定義式: Xにはs,a,cの情報が含まれる

$$\underline{x}_i = \{(s_{\underline{1}}^{(i)}, a_{\underline{1}}^{(i)}), (s_{\underline{2}}^{(i)}, a_{\underline{2}}^{(i)}), \dots, (s_{T_i}^{(i)}, a_{T_i}^{(i)}) \mid c^{(i)}\},$$

←エージェントiのOからDまでをリンク遷移(s,a)ペアに分割  
↑i人目の軌跡 ↑tステップ目



- 報酬関数Rの定義式 パラメータ  $\theta$  関数として定義し、 $\theta$  を学習させることで、Rの関数を形にする

$$R_{\theta}(x_i) = \sum_{t=1}^{T_i} \gamma^t R_{\theta}(s_t^{(i)}, a_t^{(i)} \mid c^{(i)}).$$

エージェントiの報酬関数



↓ 全エージェントの行動データから確率の立式  
(観測されたx(軌跡)の起こる確率はRのexpに比例するという法則)

$$P_{\theta}(x) = \frac{1}{Z} \exp(R_{\theta}(x))$$

行動xの起こる確率 Zは全てのx(軌跡)の $R_{\theta}(x)$ の積分和



## ■ 課題: Zを計算すること

Challenge: computation of Z

従来はRLモデルのように線形と決めることで、パラメータ推定  
→本研究ではZを直接求めず別のアプローチをとる(AIRLを利用)

# 目次

## 小結①

---

- AIRL for RCMの概要
- 定式化
- 逆強化学習とは
- 逆強化学習としての定式化
- AIRLとは
- Generator
- Discriminator
- 実装（機械学習の実装）
- AIRLの実装
- 実証実験の結果
- 機械学習の説明可能性(SHAP)

## まとめ

- **経路選択問題は逆強化学習の問題として定式化できる**

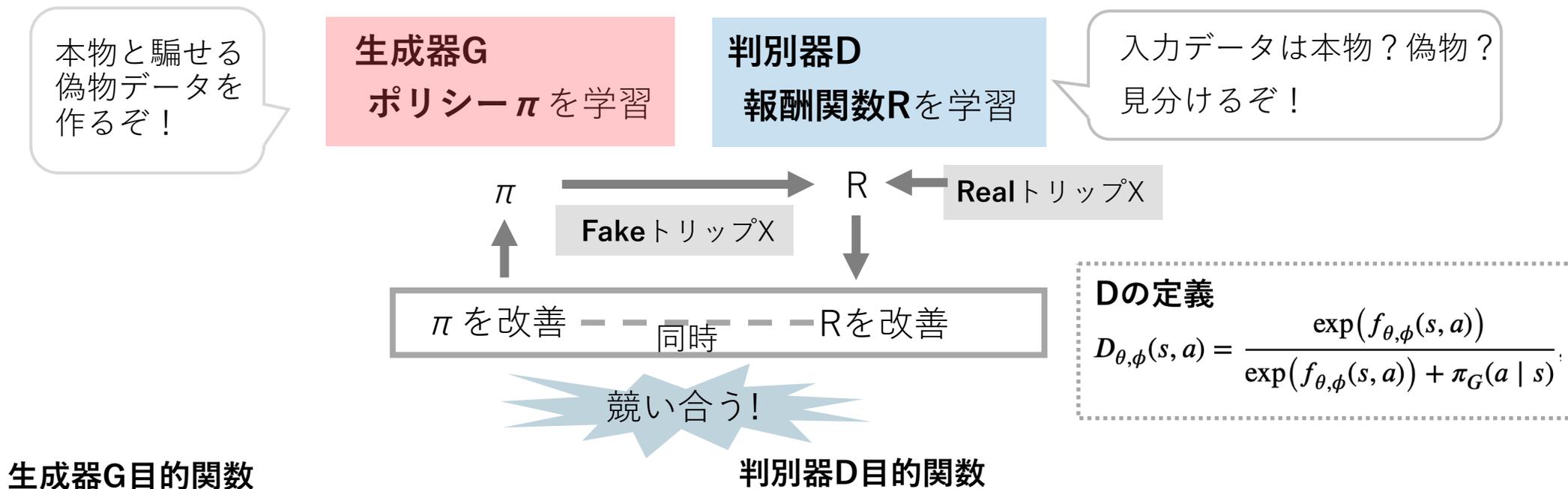
Formulation as IRL

- **逆強化学習を解くのにAIRLアルゴリズムを用いる**

use AIRL to solve IRL problem

# AIRLって？ Adversarial IRL：敵対的逆強化学習

- 逆強化学習のアルゴリズムの一種
- 生成ネットワークを援用し、方策関数を決定する生成器Gと報酬関数を決定する判別器Dという2つのモデルが互いに競い合うアルゴリズムで同時に学習  
Competition improves both models, discriminator and generator



$$\max_{\pi_G} [E_{\pi_G} [\log(D_{\theta, \phi})] - \log(1 - D_{\theta, \phi})]$$

$$\min_{\theta, \phi} -E_D [\log(D_{\theta, \phi}(s, a))] - E_{\pi_G} [\log(1 - D_{\theta, \phi}(s, a))]$$

実データ 偽データ

# Generator

**Gの目標**：OD（出発地と目的地）のペアから経路選択を生成するポリシー $\pi_G$ を学習すること

Goal; To learn policy

● **Gの入力** input

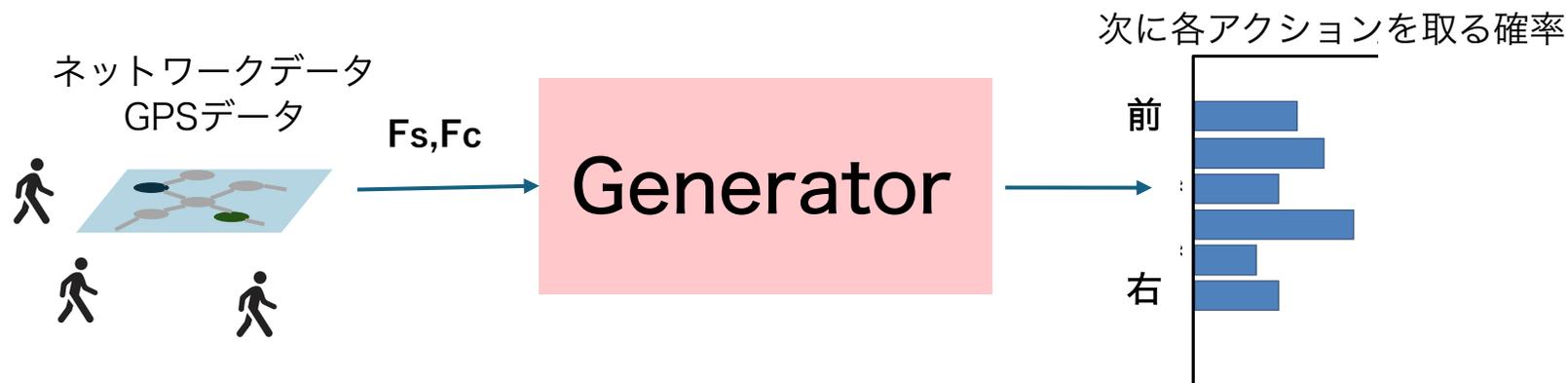
① **状態特徴量** $F_s$ ：現在のリンクの特徴（例：リンクの長さ）

② **コンテキスト特徴量** $F_c$ ：目的地に関する特徴（例：目的地までの距離）、その他（例：旅行者属性）

● **Gの出力** output

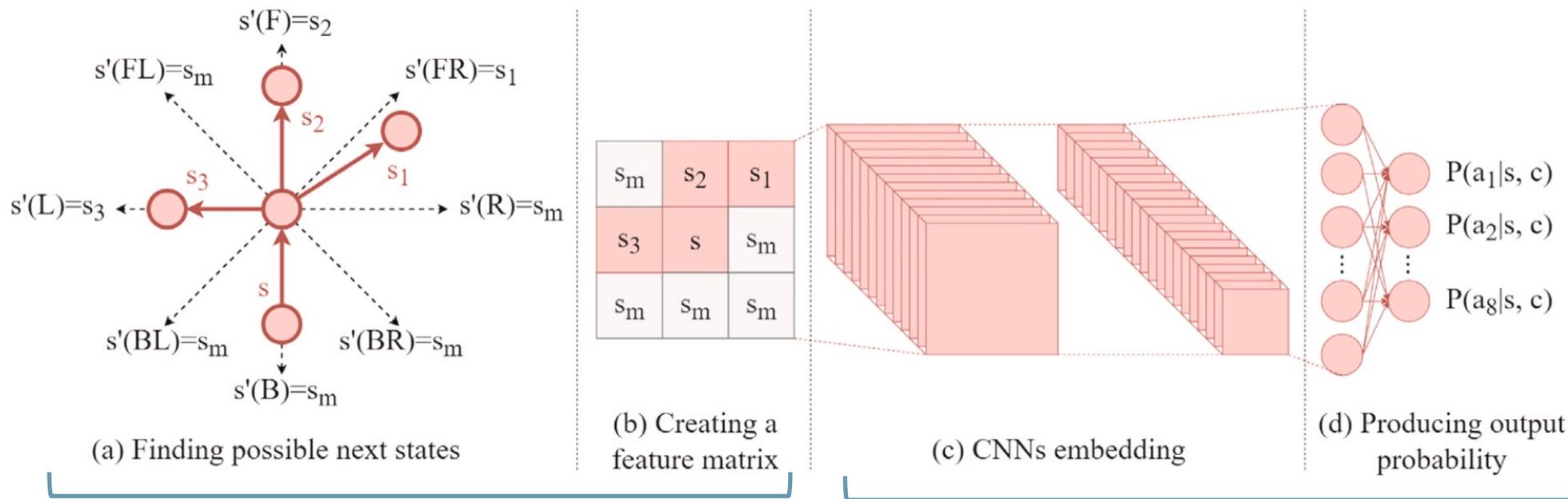
現在の状態 $s$ とコンテキスト $c$ が与えられた場合に旅行者が各アクションを選択する確率分布

- 隣接する状態間の空間的相関を考慮するために、畳み込みニューラルネットワーク（CNN←DNNの一種）を使用して、現在の状態と潜在的な次の状態のために $F = [F_s; F_c]$ を集約する



# Generator

使用しているネットワーク構造 (データ処理の流れ) Generator network



データ処理(Fs,Fcの整形)

NNでの学習

Step1

Step2

Step3

Step4

可能な次の状態を見つける。状態 $s$ が与えられた場合、各アクション（すなわち方向） $a \in A$ に対して次の状態 $s'(s, a)$ を見つける。アクションが有効な状態につながらない場合や、 $a \notin A(s)$ の場合、 $s'(s, a)$ はマスク状態として $s_m$ で示される。

特徴行列の作成。3x3の特徴行列に、状態特徴量 $Fs$ と可能な次の状態 $S'(s)$ のコンテキスト特徴量 $Fc$ を集約

畳み込みニューラルネットワーク（CNN）の埋め込みを学習する。カーネルサイズ2の2層のCNNを使用して、特徴行列から潜在空間ベクトルを学習

アクションの確率を生成する。CNNから学習した潜在空間ベクトルを入力として、2層の順伝播ニューラルネットワークとsoftmax関数を使用して出力を生成する。

# Discriminator

**Dの目標**： 実際の人間の軌跡とGeneratorが作った軌跡を見分けること

Goal; To discriminate real behavior from fake behavior generated from G.

## ● Dの入力

①状態特徴量 $F_s$ ： 現在のリンクの特徴（例：リンクの長さ）

②コンテキスト特徴量 $F_c$ ： 目的地に関する特徴（例：目的地までの距離）、その他（例：旅行者属性）

③行動ベクトル $F_a$ ：

## ● Dの出力

本物であると判断する確率D

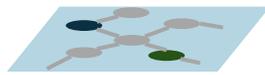
$$D_{\theta,\phi}(s, a|c) = \frac{\exp(f_{\theta,\phi}(s, a|c))}{\exp(f_{\theta,\phi}(s, a|c)) + \pi_G(a | s, c)}$$

$$f_{\theta,\phi}(s, a|c) = g_{\theta}(s, a|c) + \gamma h_{\phi}(s'|c) - h_{\phi}(s|c).$$

$g_{\theta}(s, a|c)$  Reward approximator

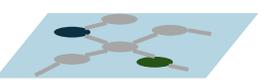
$h_{\phi}(s|c)$ . 不必要な更新を防ぐ役割

GPSデータ



$F_s, F_c, F_a$

Generator



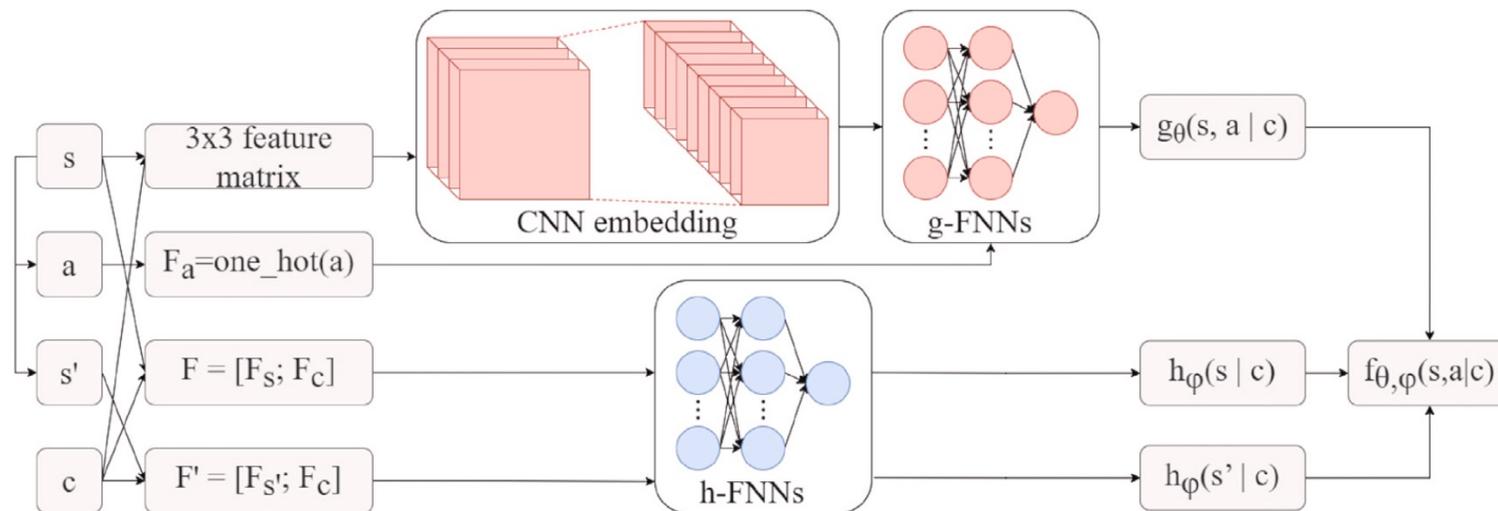
偽 $F_s, F_c, F_a$

Discriminator

本物である確率  
(0~1の間の値)

# Discriminator

## モデルイメージ

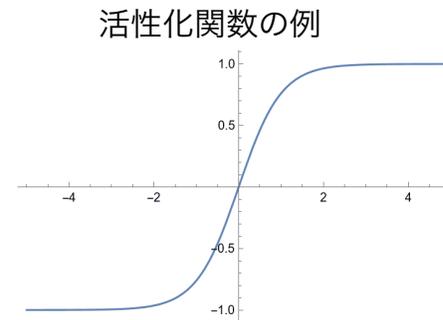
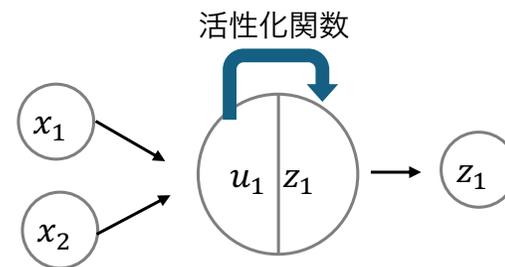
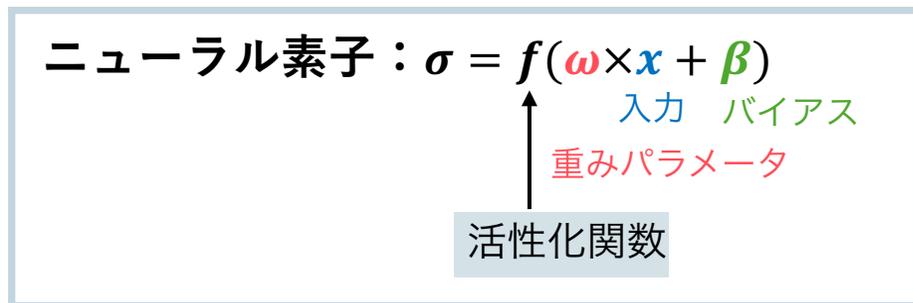


- $f$ を知るには $g$ ( $R$ の近似関数)と $h$ が必要
- $R$ は次の状態に依存するので今の状態と、隣接する状態の特徴を畳み込んで入力するためCNNネットワークを使う(CNNに $F_s, F_c$ がいる) ( $G$ の時と同じ感じ)
- $R$ にはactionも影響を及ぼす (人はまっすぐが好きとか、右とか左とかが好きとか)
- Actionの特徴はone-hotベクトルとして表し、CNNに $F = [F_s; F_c]$  がインプットされたレイヤーと繋ぐ
- $g$ はこのCNNレイヤーによって学習される
- $h$ もCNNレイヤーを使って学習できるが計算の都合上、 $g, h$ から $f$ を出し、 $D$ を求めたら、 $R$ を以下の式で求める

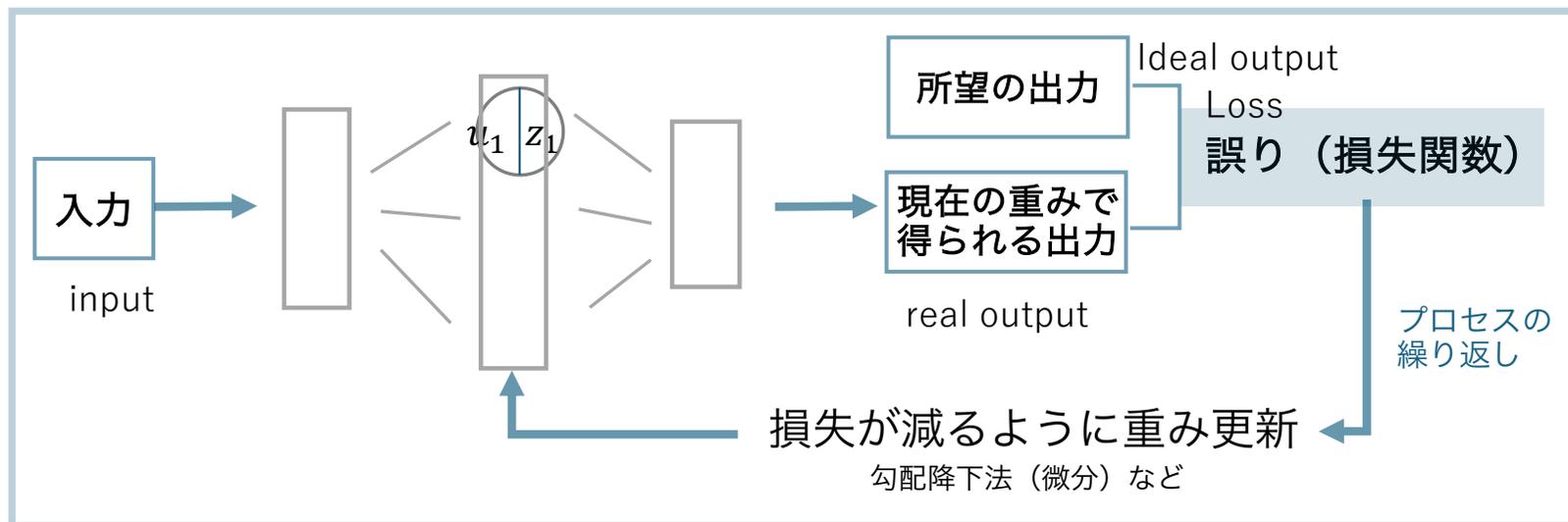
$$R_{\theta, \phi}(s, a) = \log(D_{\theta, \phi}(s, a)) - \log(1 - D_{\theta, \phi}(s, a)) = f_{\theta, \phi}(s, a) - \log \pi_G(a | s).$$

# 機械学習の実装 復習

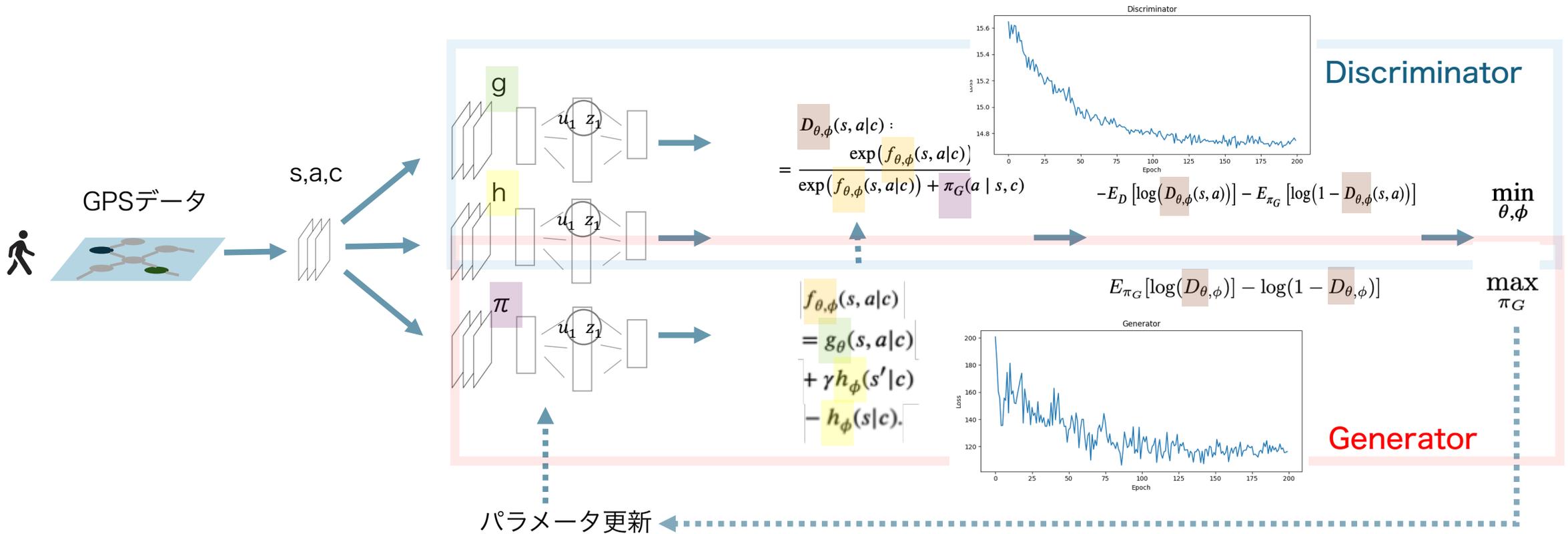
- ニューラルネットワーク：深層学習の基本的な仕組み Neural Networks



- モデル構築の基本（パラメータの学習プロセス）



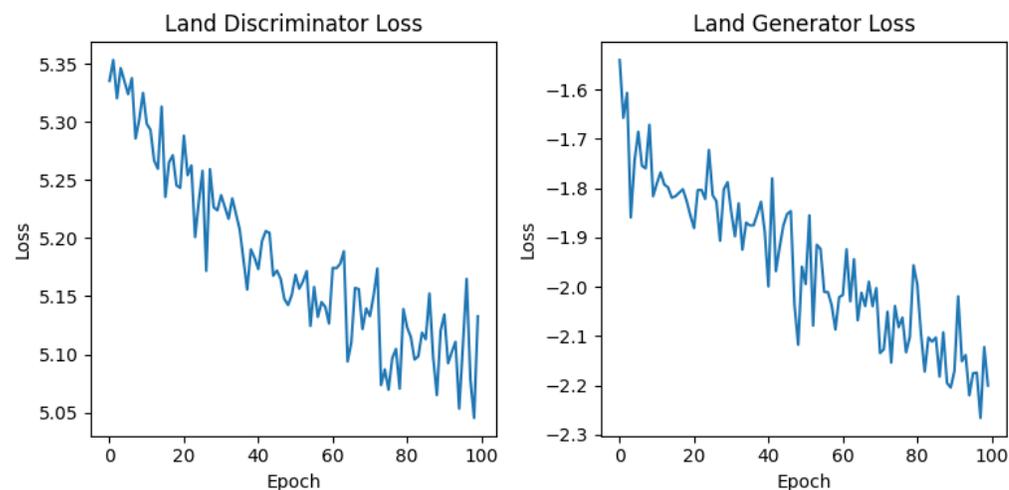
# AIRL実装の流れ



➡  $R_{\theta,\phi}(s, a | c) = \log(D_{\theta,\phi}(s, a | c)) - \log(1 - D_{\theta,\phi}(s, a | c))$  : 目標の報酬関数の復元達成

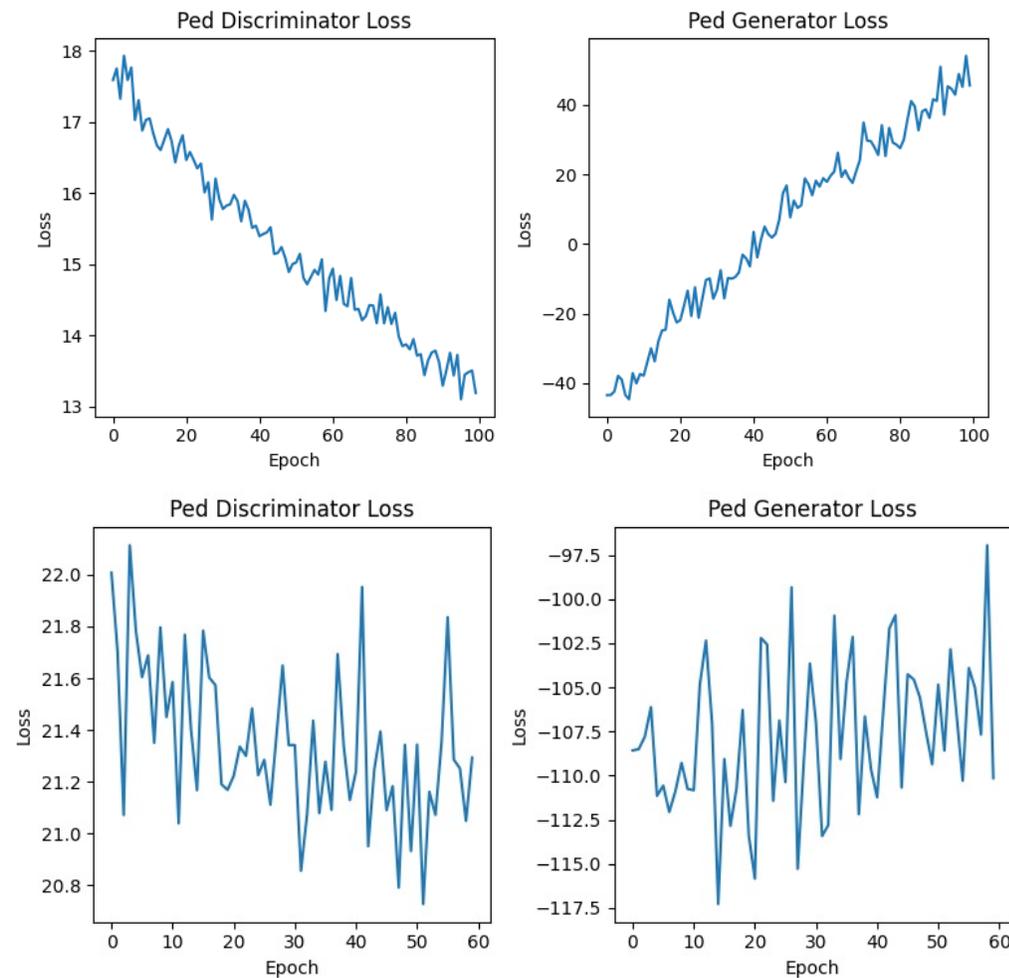
# AIRL学習の可視化

- 成功例



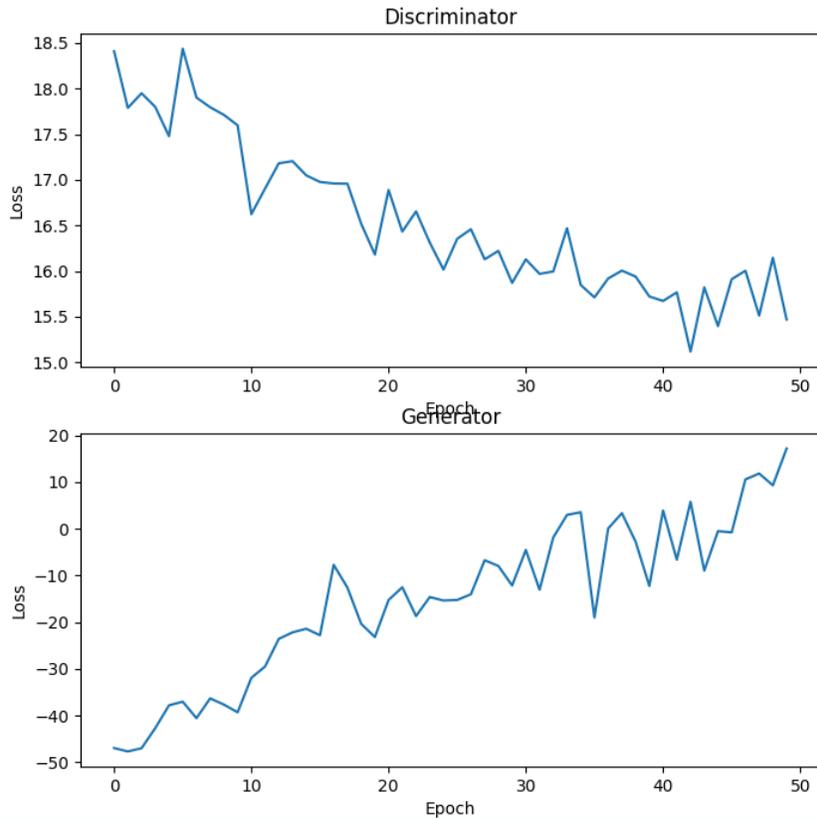
初期値, 学習率の変更で調整

- 失敗例

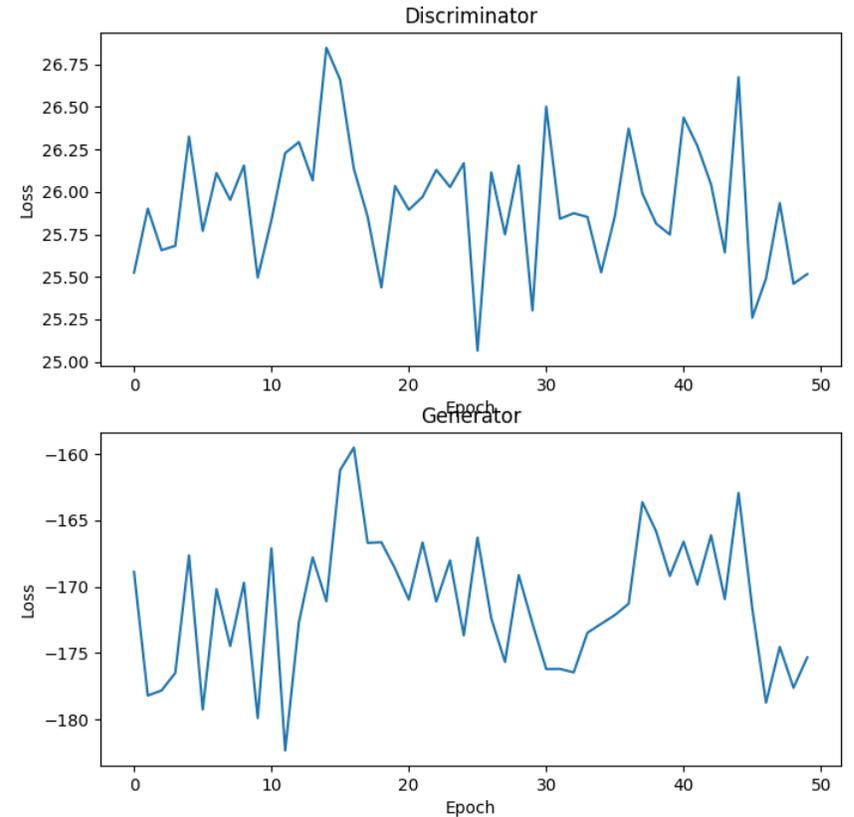


# AIRL学習の可視化

## 学習率変更



学習率0.1倍に



```
self.optimizer_util = optim.Adam(self.util.parameters(), lr=0.000001)#lr=0.000003
self.optimizer_value = optim.Adam(self.value.parameters(), lr=0.000001)#lr=0.00005
self.optimizer_policy = optim.Adam(self.policy.parameters(), lr=0.000001)#lr=0.00005
self.optimizer_ext = optim.Adam(self.ext.parameters(), lr=0.000001)#lr=0.00005
```

# 目次

## 小結②

---

- AIRL for RCMの概要
- 定式化
- 逆強化学習とは
- 逆強化学習としての定式化
- AIRLとは
- Generator
- Discriminator
- 実装（機械学習の実装）
- AIRLの実装
- 実証実験の結果
- 機械学習の説明可能性(SHAP)

## まとめ②

- AIRLはGeneratorとDiscriminatorの2つのモデルが競い合いながら精度を高める
- DやGを構成する3つの関数がニューラルネットワークで学習され（機械学習）、報酬関数が求まる

# 実証実験 データ

## 実証実験で用いたデータ

- 上海の大手タクシー会社の1つから提供されたデータセット
- 10609台のタクシーのGPS追跡, 24470のトリップデータ

## RCM-AIRLモデル比較する他のモデル

- Path Size Logit(PSL) : パスベースモデルの一つ。
- DNN-PSL : DNNを使用するPSLの拡張版。非線形な関係や柔軟なコンテキストを含めることができるのが特徴
- Recursive Logit : リンクベースモデル。価値関数は線形関数。パスの量は膨大だが、リンクベースなので全てのパスを抽出する必要がない
- RCM-BCモデル
- RCM-GAILモデル



Fig. 5. Selected road network in Shanghai.

従来モデル

AIRL類似  
モデル

# 実証実験

## 評価指標 Evaluation

モデルの精度の良さを測る指標 (実際のデータとモデルが予測するデータを比較する)

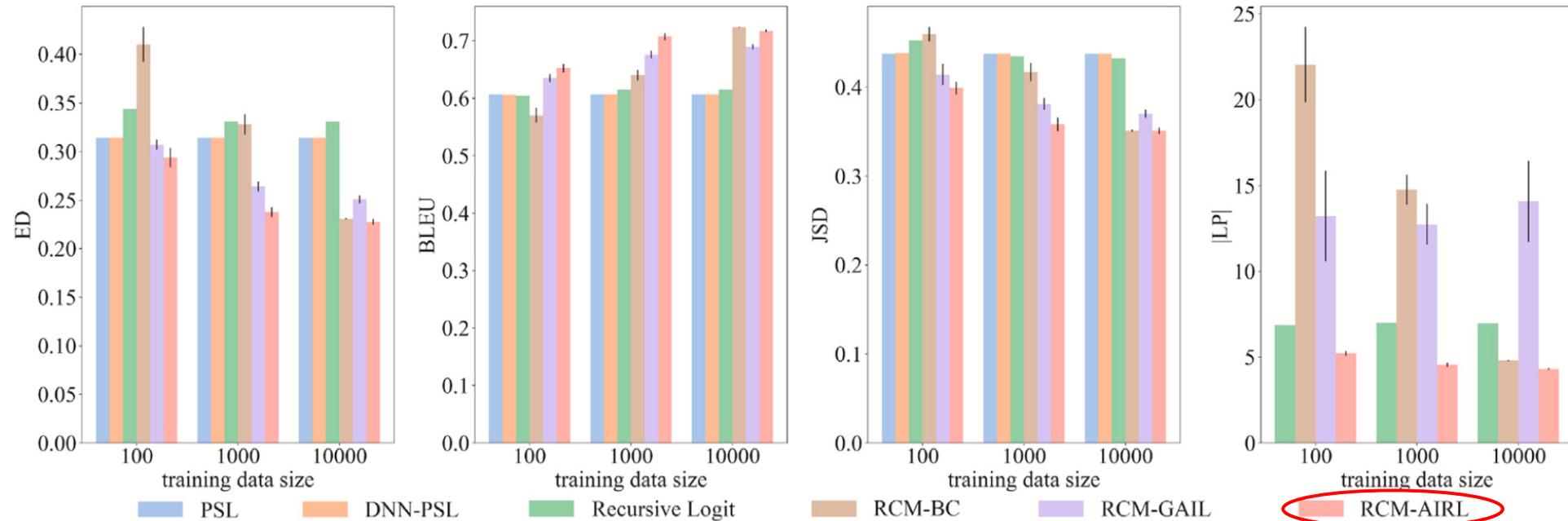
- Edit Distance(ED): 2つの関連するものがどれくらい**異なっているか**を測る指標。  
一方を他方に変換する最小の操作回数の計算によって求める。
- BiLingual Evaluation Understudy score (BLEU): スコアは2つがどの程度似ているか、**近い**かを示す
- Jensen-Shannon Distance (JSD): 2つのものの確率分布の**近さ**を測ることができる
- Log Probability (LP): 与えられたモデルにおいて実データの経路の確率のlogをとったものの平均値

## データ成形

- データをランダムに5つに分け、1つをテスト用、残り4つトレーニング用
- 学習用のデータサイズを100,1000,10000の3種類で実験
- 実験成功のために重要だとわかったこと：コンテキストに依存する報酬関数の設計

←ゴールに関連する情報をきちんと与えることで旅行者が合理的な行動として経路選択すると意味付けられるので、何度もループしたりゴールに到達しないという課題を解消する

# 結果 モデル比較



- RCM-BC、RCM-GAIL、RCM-AIRLは、既存のベースラインモデルよりも優れたパフォーマンスを示した  
→深層IRL（逆強化学習）/IL（模倣学習）手法がルート選択モデリングにおいて効果的であることを示唆  
(特にAIRLは全ての評価指標に関して良い結果)
- モデルの潜在的な不果実性を考慮するため、3回実験を行なったものの標準偏差をエラーバーとして示したところ、AIRLは割と安定

# 結果

## 交通流のシミュレーション結果

$R^2$  : 回帰モデルの予測の適合度を評価する統計的指標

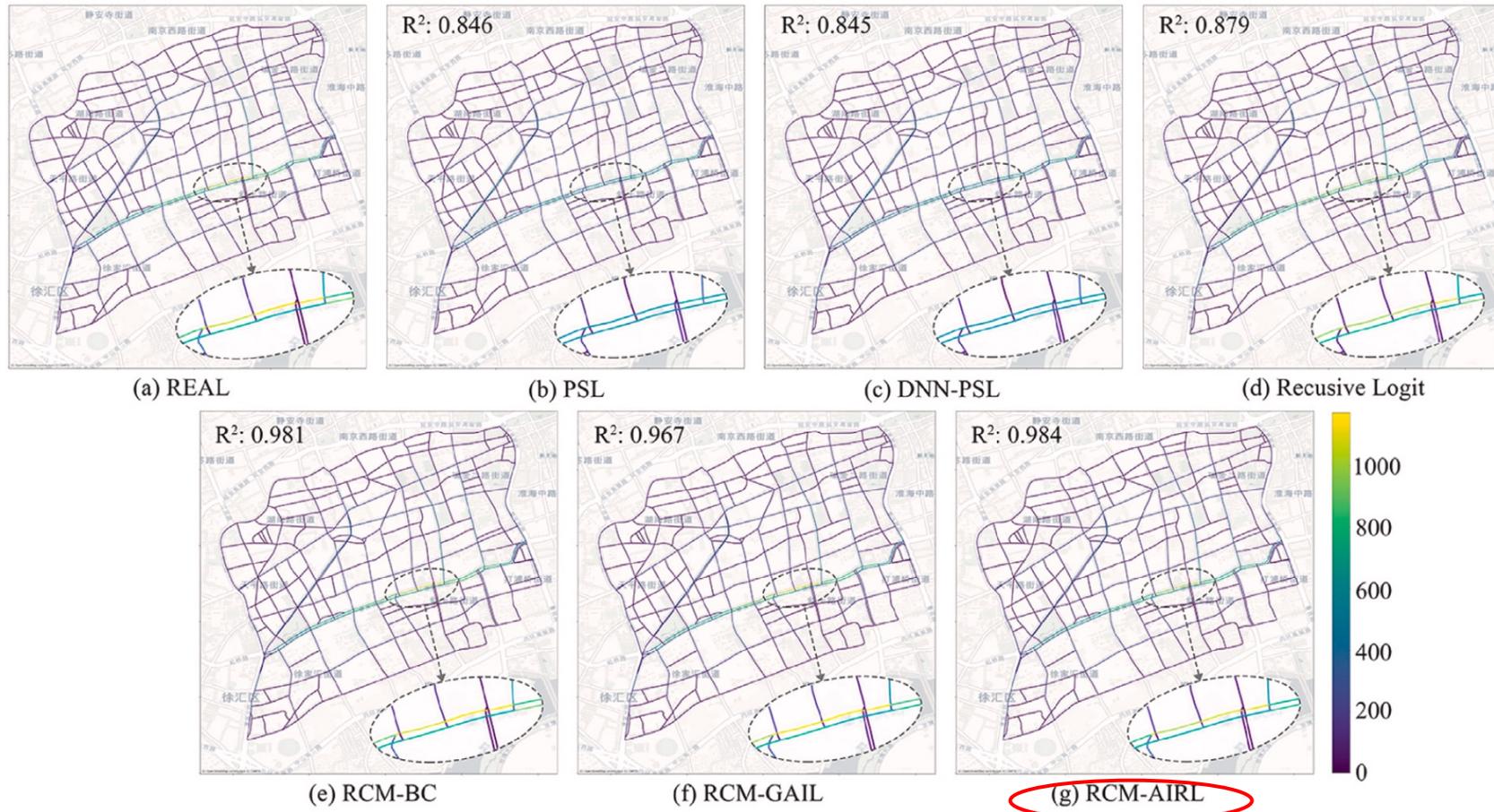


Fig. 7. Predicted link flow distribution using different models.

# 結果

## あるODペアに対する実際のルート選択と予測



# 結果

## 計算時間

- FE: 特徴量エンジニアリングに必要な時間
  - MO: モデルの最適化に必要な時間
  - PG: 経路生成に必要な時間
- モデルの訓練 (学習)
- モデルのテスト

Table 3

Computation time for feature engineering (FE), model optimization (MO), and path generation (PG) of different models. The results are averaged over 5-fold cross validation.

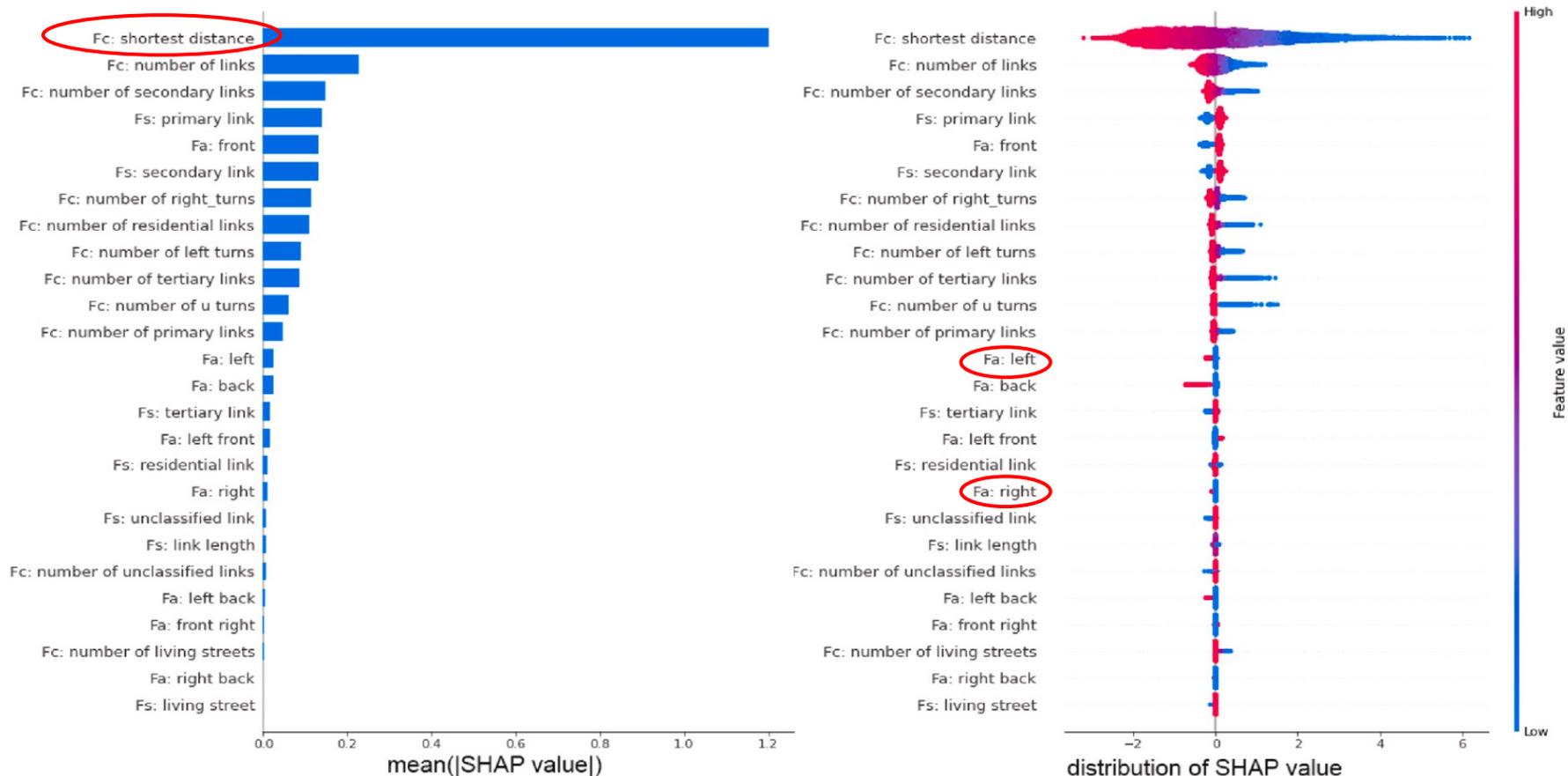
Models		100 training trips			1000 training trips			10000 training trips		
		FE(s)	MO(s)	PG(s)	FE(s)	MO(s)	PG(s)	FE(s)	MO(s)	PG(s)
Path-based	PSL	28.8	0.3	1370.5	277.8	1.8	1370.5	3166.7	12.6	1370.4
	DNN-PSL	28.8	0.8	1370.4	277.8	5.7	1370.4	3166.7	42.3	1370.4
Link-based	Recursive logit	0.0	22.2	4.0	0.0	44.4	4.0	0.0	238.8	4.1
	RCM-BC	3840.0	124.8	1.7	3840.0	1174.1	1.6	3840.0	11 487.4	1.6
	RCM-GAIL	3840.0	6112.8	1.7	3840.0	9305.3	1.6	3840.0	18 270.9	1.6
	RCM-AIRL	3840.0	5990.4	1.6	3840.0	8962.5	1.5	3840.0	18 307.0	1.5

モデル訓練は1度でよく、交通流の予測やシミュレーションにはモデルのテストにかかる時間が重要

# SHAP

## 経路選択行動の”解釈可能性”について

- SHapley Additive exPlanations (SHAP) を使用して、入力特徴量の影響を理解
- SHAPは、機械学習モデルの説明を行うためのゲーム理論的なアプローチを用いる説明可能なAI技術



# 参考 定式化

## 3.2. Route choice as an inverse reinforcement learning problem

Recall that IRL for route choice modeling seeks to infer the reward function  $R(s, a | c)$  given a set of observed trajectories  $X = \{x_1, x_2, \dots, x_N\}$  that are assumed to be drawn from an optimal policy  $\pi^*(a | s, c)$ . Note that each observed trajectory is a sequence of station–action pairs with a certain context. Let us denote the  $i$ th trajectory as  $x_i = \{(s_1^{(i)}, a_1^{(i)}), (s_2^{(i)}, a_2^{(i)}), \dots, (s_{T_i}^{(i)}, a_{T_i}^{(i)}) | c^{(i)}\}$ , where  $(s_t^{(i)}, a_t^{(i)})$  is the  $t$ th state–action pair of  $x_i$ ,  $c^{(i)}$  the context, and  $T_i$  the length of the trajectory.

In large road networks, there are many possible combinations of  $(s, a, c)$ , making it difficult to obtain a robust estimate of  $R(s, a | c)$ . A common solution is to approximate it with a parameterized reward function  $R_\theta(s, a | c)$ , where  $\theta$  is the function parameters to be learned. The specific function approximation can take many different forms, among which DNNs are generally more flexible and adept in dealing with high-dimensional features. As the total reward of a whole trajectory  $x_i$  is simply the sum of the discounted rewards for each state–action pair in the trajectory, the parameterized reward function can be expressed as

$$R_\theta(x_i) = \sum_{t=1}^{T_i} \gamma^t R_\theta(s_t^{(i)}, a_t^{(i)} | c^{(i)}). \quad (1)$$

While the underlying policy that generates the observed trajectories is assumed to be “optimal”, actual routing behavior may be suboptimal and vary across individuals. Even for the same OD pair, different travelers may choose different routes. To address this issue, the principle of maximum entropy is adopted to handle behavioral suboptimality as well as stochasticity by operating on the distribution over possible trajectories. Following the MaxEntIRL formulation (Ziebart et al., 2008), the probability of observing any given trajectory  $x$  is proportional to the exponential of its cumulative reward:

$$P_\theta(x) = \frac{1}{Z} \exp(R_\theta(x)), \quad (2)$$

where the partition function  $Z$  is the integral of  $R_\theta(x)$  over all possible trajectories. Therefore, we can frame the IRL problem as solving the maximum likelihood problem based on the observed trajectories:

$$\max_{\theta} \sum_{i=1}^N \log(P_\theta(x_i)). \quad (3)$$

# 参考

---

## 本

- 岡谷貴之著 深層学習 改訂第2版 (機械学習プロフェッショナルシリーズ)
- 松尾豊著 人工知能は人間を超えるか
- 斎藤 康毅 ゼロから作るDeep Learning —Pythonで学ぶディープラーニングの理論と実装

## 論文

- A deep inverse reinforcement learning approach to route choice modeling with context-dependent rewards, Zhan Zhao, Yuebing Liang, 2023  
<https://www.sciencedirect.com/science/article/abs/pii/S0968090X23000682>
- Learning Robust Rewards with Adversarial Inverse Reinforcement Learning, Justin Fu, Katie Luo, Sergey Levine, 2017  
<https://arxiv.org/abs/1710.11248>