

吉澤尚輝（東北大学） 井料隆雅（東北大学）

はじめに

- 強化学習を使用した適応型信号制御が近年注目されている。
- 適応型信号制御と利用者の経路選択の相互作用を考慮した適応型信号制御は少ない。
- 強化学習を使用しない、信号制御と経路選択の相互作用を考慮した制御にPolicy P_0 がある。
 - Policy P_0 には、飽和交通流率の小さな枝に交通量とスプリットが集中し、ネットワークの交通容量がネットワークへ流入する総需要より少なくなることを防ぐ **capacity maximization**という性質がある。
 - Policy P_0 以外の制御にこの性質があるか不明。
- 強化学習を使用した、信号制御と経路選択の相互作用を考慮した制御は知る限り存在しない。
- 適応型信号制御に使用する強化学習の学習時、および学習済みの適応型信号制御の性能評価時の両方で、利用者の経路選択を考慮する。
- この時、強化学習を使用した適応型信号制御が **capacity maximization**を達成するか調べる。

シミュレーション設定

- 飽和交通流率 S_{we} は小さいが旅行時間が短い経路と、飽和交通流率 S_{ns} は大きい旅行時間が長い経路、の2経路が交差点に流入するネットワーク。
- スプリットの決定には、**Deep Q network (DQN)**を使用し、このネットワークを DQN_s と呼ぶ。
- DQN_s の学習時に流入する需要は経路選択を行う。
- 学習時の経路選択方法は、**ロジットモデル**と、経路選択用の DQN (DQN_r) の2種類。
 - DQN_r の報酬は DQN_s の報酬に -1 を掛けたものになっており、 DQN_r と DQN_s は互いに敵対的な行動を選択する。
 - これは **Generative Adversarial Networks (GAN)** の考え方を応用している。
- DQN_s の性能評価に使用する需要はロジットモデルに従って経路変更する。

シミュレーション結果 ロジットモデルに従う需要で学習

- 表1に、学習時に流入する需要の分散パラメータ θ が各値の DQN_s に対して、流入流率の期待値 l と分散パラメータ θ が各値の需要が流入した時の R を示す。
 - R は適応型信号制御が実現する交通容量に対するネットワークへの流入流率の割合で、**capacity maximization** が達成されている場合、 $R \leq 1.0$ となる。
 - 表中の色は、 $R=0.1$ に近いほど緑が、 $R=2.0$ に近いほど赤が濃くなる。 $R=1.0$ の時は白である。
 - 学習時の $\theta=0.25$ の行内にある nan の表記は、 DQN_s の性能を評価するシミュレーションが終了しなかったため、各枝の平均スプリットが算出できなかったことを示す。この現象は、車両が存在する枝に対してスプリットが割り振られず、車両が交差点から流出できなかったため発生した。
- 流入する需要が各 θ の時に実現する制御の R は、学習時に使用した需要の θ によって、大きく異なった。
 - 多くの条件で $R > 1.0$ となり、 l が大きい条件では $R=1.8$ 程度まで達した。

GANを応用した手法で学習

- 表2は学習時の流入流率の期待値が各値の DQN_s に対して、 l と θ が各値の需要が流入した時の R で、表中の色は表1と同じ意味である。
- 多くの条件で、 $R > 1.0$ となっていた。
 - 学習時の需要の流入流率の期待値が1440台/時、2160台/時、3600台/時の場合、 $R=2.0$ に近い値となっている場合も存在した。
 - 一方、2880台/時では、 $l=2880$ 、3060の時、どの θ が流入しても R は比較的1.0に近い値で安定している。

まとめ

- **ロジットモデル**に従う需要で学習した場合、各 θ の需要が流入した時に実現するスプリットは、学習時の需要の θ や評価時の流入流率の期待値によって、大きく異なった。
 - 学習時の θ の値によって R が1.2程度に収まるケースから1.8程度に達するものまで存在した。
- **GAN**を応用した学習の場合、学習時の需要の流入流率の期待値によっては、どのような θ や流入流率の期待値の需要が流入しても、1.2程度の相対的に小さな R が実現した。

表1 各 θ の流入需要で学習した DQN_s に、各 θ の需要が流入した時に実現する R

		評価時に流入する 需要の分散パラメータ $\theta=0.125$	評価時に流入する 需要の分散パラメータ $\theta=0.25$	評価時に流入する 需要の分散パラメータ $\theta=0.5$	評価時に流入する 需要の分散パラメータ $\theta=0.75$	評価時に流入する 需要の分散パラメータ $\theta=1.0$	評価時に流入する 需要の分散パラメータ $\theta=\infty$
学習時に 流入する需要の 分散パラメータ $\theta=0.0$	$l=720$ 台/時	0.35	0.35	0.35	0.35	0.35	0.35
	$l=2160$ 台/時	0.92	0.93	0.96	0.98	0.91	0.92
	$l=2880$ 台/時	1.08	1.05	1.05	1.08	1.02	1.07
	$l=3060$ 台/時	1.14	1.11	1.13	1.09	1.09	1.06
	$l=3600$ 台/時	1.23	1.21	1.20	1.21	1.19	1.21
学習時に 流入する需要の 分散パラメータ $\theta=0.125$	$l=720$ 台/時	0.35	0.35	0.35	0.35	0.35	0.35
	$l=2160$ 台/時	1.08	1.11	1.10	0.92	1.09	1.23
	$l=2880$ 台/時	1.04	1.46	1.62	1.43	1.65	1.13
	$l=3060$ 台/時	1.79	1.16	1.79	1.17	1.15	1.13
	$l=3600$ 台/時	1.23	1.28	1.30	1.33	1.44	1.25
学習時に 流入する需要の 分散パラメータ $\theta=0.25$	$l=720$ 台/時	0.35	0.35	0.35	0.35	0.35	0.35
	$l=2160$ 台/時	nan	nan	nan	nan	nan	nan
	$l=2880$ 台/時	nan	nan	nan	nan	nan	nan
	$l=3060$ 台/時	nan	nan	nan	nan	nan	nan
	$l=3600$ 台/時	nan	nan	nan	nan	nan	nan
学習時に 流入する需要の 分散パラメータ $\theta=0.5$	$l=720$ 台/時	0.35	0.35	0.35	0.35	0.35	0.35
	$l=2160$ 台/時	0.92	0.95	0.94	0.96	0.95	0.95
	$l=2880$ 台/時	1.09	1.05	1.08	1.06	1.07	1.06
	$l=3060$ 台/時	1.11	1.19	1.12	1.11	1.11	1.07
	$l=3600$ 台/時	1.19	1.19	1.19	1.27	1.20	1.24
学習時に 流入する需要の 分散パラメータ $\theta=0.75$	$l=720$ 台/時	0.35	0.35	0.35	0.35	0.35	0.35
	$l=2160$ 台/時	0.93	0.93	0.94	0.87	0.95	0.96
	$l=2880$ 台/時	1.07	1.02	1.06	1.06	1.11	1.07
	$l=3060$ 台/時	1.04	1.07	1.12	1.05	1.06	1.07
	$l=3600$ 台/時	1.20	1.21	1.22	1.18	1.20	1.16
学習時に 流入する需要の 分散パラメータ $\theta=1.0$	$l=720$ 台/時	0.35	0.35	0.35	0.35	0.35	0.35
	$l=2160$ 台/時	1.17	1.16	1.13	1.21	1.28	1.16
	$l=2880$ 台/時	1.56	1.59	1.59	1.56	1.57	1.56
	$l=3060$ 台/時	1.68	1.66	1.69	1.62	1.71	1.65
	$l=3600$ 台/時	1.88	1.91	1.82	1.92	1.94	1.90
学習時に 流入する需要の 分散パラメータ $\theta=\infty$	$l=720$ 台/時	0.35	0.35	0.35	0.35	0.35	0.35
	$l=2160$ 台/時	0.98	1.01	0.95	0.97	0.96	0.94
	$l=2880$ 台/時	1.10	1.13	1.10	1.05	1.07	1.12
	$l=3060$ 台/時	1.13	1.18	1.15	1.13	1.11	1.17
	$l=3600$ 台/時	1.32	1.28	1.30	1.30	1.29	1.28

表2 GANを応用した手法において、各流入流率の期待値を持つ流入需要で学習した DQN_s に、各 θ の需要が流入した時に実現する R

		評価時に流入する 需要の分散パラメータ $\theta=0.125$	評価時に流入する 需要の分散パラメータ $\theta=0.25$	評価時に流入する 需要の分散パラメータ $\theta=0.5$	評価時に流入する 需要の分散パラメータ $\theta=0.75$	評価時に流入する 需要の分散パラメータ $\theta=1.0$	評価時に流入する 需要の分散パラメータ $\theta=\infty$
学習時に 流入する需要の 流入流率期待値 1440台/時	$l=720$ 台/時	0.35	0.35	0.35	0.35	0.35	0.35
	$l=2160$ 台/時	0.99	0.96	0.97	0.96	0.98	0.95
	$l=2880$ 台/時	1.28	1.62	1.42	1.17	1.43	1.28
	$l=3060$ 台/時	1.69	1.57	1.40	1.45	1.24	1.22
	$l=3600$ 台/時	1.44	1.83	1.84	1.39	1.80	1.82
学習時に 流入する需要の 流入流率期待値 2160台/時	$l=720$ 台/時	0.35	0.35	0.35	0.35	0.35	0.35
	$l=2160$ 台/時	0.98	0.95	1.00	0.98	0.96	0.95
	$l=2880$ 台/時	1.21	1.64	1.16	1.68	1.57	1.59
	$l=3060$ 台/時	1.22	1.27	1.26	1.72	1.21	1.74
	$l=3600$ 台/時	1.31	1.47	1.29	1.30	1.31	1.30
学習時に 流入する需要の 流入流率期待値 2880台/時	$l=720$ 台/時	0.35	0.35	0.35	0.35	0.35	0.35
	$l=2160$ 台/時	0.94	0.92	0.94	0.91	0.94	0.96
	$l=2880$ 台/時	1.04	1.04	1.02	1.04	1.06	1.02
	$l=3060$ 台/時	1.09	1.09	1.12	1.08	1.05	1.07
	$l=3600$ 台/時	1.21	1.21	1.19	1.18	1.19	1.19
学習時に 流入する需要の 流入流率期待値 3600台/時	$l=720$ 台/時	0.35	0.35	0.35	0.35	0.35	0.35
	$l=2160$ 台/時	0.96	0.97	0.97	0.97	0.98	0.96
	$l=2880$ 台/時	1.15	1.10	1.12	1.05	1.10	1.09
	$l=3060$ 台/時	1.25	1.12	1.15	1.12	1.15	1.10
	$l=3600$ 台/時	1.41	1.25	1.45	1.29	1.31	1.38

今後の課題

- 本研究で DQN_r の選択する行動は使用する経路のみだが、流入間隔も DQN_r に選択させることで、学習時の需要の流入流率の期待値を設定する必要がなくなるかもしれない。