

# 교통카드 태깅정보를 활용한 교통혼잡도 예측 그래프 모델

## Graph Model for Predicting Traffic Chaos Using Transportation Card Tagging Information

박경현, 서석희

(인하대학교, 학부생), (인하대학교 학부생)

Key Words : Subway, GAT, Wasserstein Distance

### 목 차

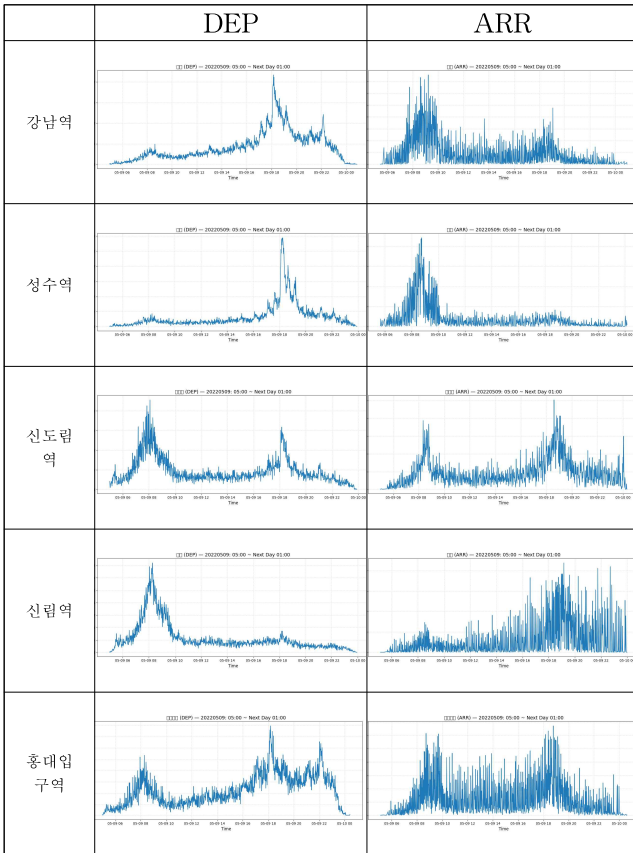
- I. 연구 개요
- II. 실험
  1. 데이터셋
  2. Earth Mover's Distance (Wasserstein metric)
- III. 실험 구성
- III. 결론
  1. 결론 및 향후 과제
  2. 모델 코드

## I. 연구 개요

최근 도시 교통망의 복잡성이 증가함에 따라, 대중교통의 효율적인 운영과 혼잡도 관리가 도시 계획의 핵심 과제로 부상하고 있다. 특히 지하철과 같은 대중교통 수단은 특정 시간대나 노선에서 승객 수요의 급격한 변동이 발생하기 때문에, 이를 사전에 예측하는 것은 교통 자원 배분 및 운영 효율 향상에 매우 중요하다. 이러한 예측 문제를 해결하기 위해 인공지능(AI) 기반의 시계열 예측 및 그래프 신경망(Graph Neural Network, GNN) 모델이 활발히 연구되고 있다.

본 연구에서는 그래프 어텐션 네트워크(Graph Attention Network, GAT)의 구조적 장점을 활용하여, 시간적·공간적 상관관계를 동시에 학습할 수 있는 혼잡도 예측 모델을 제안한다. 특히 기존의 ASTGNN[1] 모델을 기반으로, 교통 네트워크의 복잡한 공간적 의존성을 더욱 정교하게 반영하기 위해 지구 거리 기반의 EMD(Earth Mover's Distance)를 결합한 ASTGNN-EMD 모델을 개발하였다.

제안된 모델은 교통카드 태깅 데이터를 입력으로 받아 지하철 노선 간 상호 연관성과 시간적 흐름을 동시에 고려함으로써, 기존 모델 대비 향상된 예측 성능을 달성하였다. 또한 실제 데이터 환경에서의 실험을 통해 모델의 일반화 성능과 안정성을 검증하였으며, 이를 통해 보다 신뢰성 있는 지하철 혼잡도 예측 시스템을 구축하고자 한다.



<표 1 : 역별 승하차 인원>

## II. 실험

### 1. 데이터셋

본 연구에서 사용한 데이터셋은 지하철 교통카드 태깅 정보를 기반으로 구축하였다. 해당 데이터는 승객이 교통카드를 이용해 지하철 개찰구에서 태그한 승·하차 기록을 포함하며, 이를 통해 시간대별 승객 흐름과 역 간 이동 패턴을 분석할 수 있다. 연구 대상 노선은 서울 지하철 2호선으로 한정하였다. 2호선은 순환선 구조를 가지며, 서울의 주요 상업·업무 중심지를 연결하는 핵심 노선이다. 이러한 구조적 특성으로 인해 국내에서 가장 많은 유동 인구를 보유한 노선으로 알려져 있으며, 출퇴근 시간대에 승·하차 인구의 급격한 변동이 나타난다.

2호선의 데이터셋은 (timestamp, node, feature) 형태로 구성하였다.

#### 1) timestamp

2022년 5월 한 달간의 1분 단위 데이터로, 시간 흐름에 따른 세밀한 혼잡도 변화를 반영한다.

#### 2) node

2호선의 대표 순환역 50개를 기준으로 설정하였으며, 각 노드는 하나의 지하철역을 의미한다.

#### 3) feature

각 시간대별 승차 인원 수와 하차 인원 수로 구성되어 있으며, 이를 통해 시간·공간적 승객 분포 변화를 학습할 수 있다.

이러한 특성으로 인해 기존 도로 기반 모델에서 사용되던 단일 feature 융합 학습(attention 기반 혼합)은 지하철의 특성을 반영하기 어렵다. 따라서 본 연구에서는 승차와 하차 feature를 분리하여 학습하는 구조를 제안한다.

### 2. Earth Mover's Distance (Wasserstein metric)

EMD는 두 확률 분포 간의 차이를 측정하기 위한 거리로, Wasserstein Distance라고도 불린다. 직관적으로는 한 분포를 다른 분포로 변환하기 위해 이동해야 하는 “질량의 최소 이동 비용”으로 해석할 수 있다. 즉, 분포의 형태가 다를수록 더 많은 “이동량”이 필요하므로 EMD 값이 커진다.

EMD는 단순한 유클리드 거리나 코사인 유사도와 달리, 분포의 전체적인 구조적 차이와 공간적 거리를 함께 고려한다는 점에서 장점이 있다. 이로 인해 영상 처리, 자연어 처리, 그리고 교통 네트워크 분석 등 공간적 상관성이 중요한 분야에서 널리 활용된다.

교통 예측 문제에서는 각 노드(예: 역 또는 교차로) 간의 거리나 연결 강도를 단순한 고정 값으로 두는 대

신, EMD를 통해 실제 승하차 분포나 이동 패턴의 유사도를 계산할 수 있다. 본 연구에서는 이러한 EMD의 특성을 활용하여, 지하철역 간의 혼잡도 분포 차이와 공간적 연관성을 정량적으로 반영함으로써 그래프 모델의 학습 효율성을 향상하였다.

### 3. 실험 구성

#### 1) Feature 분리 학습 (Separated Feature Learning)

본 연구에서는 교통 흐름 예측에서 일반적으로 사용되는 PEMS 도로 교통 데이터셋과 달리, 지하철 교통 데이터셋의 구조적 특성을 고려하였다. 지하철 데이터는 각 노드(역)에 대해 출발(Departure, DEP)과 도착(Arrival, ARR) 정보가 명확히 구분되어 있으며, 이는 하나의 노드에서 발생하는 승차 인원이 다른 노드의 하차 인원으로 이어지는 정규화된 이동 구조(normalized mobility structure)를 가진다.

모델은 한 노드의 승차 정보를 입력으로 받아 다른 노드들의 하차 정보를 예측하도록 학습한다. 이는 실제 교통 흐름에서 승차 인원이 다른 역의 하차 인원 분포를 결정 짓는다는 점을 반영한다.

이를 수식으로 표현하면 식 (1)과 같다.

$$\hat{Y}_t^{arr} = f_{\theta}(X_{t-H:t-1}^{dep}) \quad (1)$$

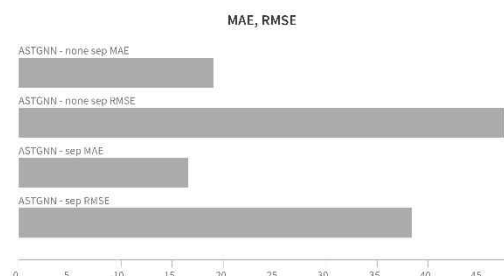
$X_{t-H:t-1}^{dep}$  : 과거 H 타임스텝 동안의 승차(dep) 입력

$\hat{Y}_t^{arr}$  : 시점 t의 하차(arr)인원 예측

$f_{\theta}(\cdot)$  : 그래프 신경망 기반의 예측 함수

실험 결과, feature를 분리하지 않은 경우 MAE =  $19.2 \pm 0.5$ , RMSE =  $48 \pm 1.2$  였던 반면, feature 분리 학습을 적용한 경우 MAE =  $16.4 \pm 0.4$ , RMSE =  $38.2 \pm 1.6$  으로 MAE는 약 14.6% 감소, RMSE는 약 20.4% 감소하여 성능이 향상되었다.

하차 정보의 경우 1분 단위의 데이터 수집 과정에서 열차 정차 시 순간적으로 인원이 몰리는 현상 때문에 이상치(0 또는 0에 가까운 값)가 다수 발생하였다. 이에 따라 MAPE 측정 시 편차가 과도하게 커지므로, 해당 지표는 제외하였다.



<그림 1 : 분리 학습 결과>

2) EMD 기반 보조 손실 (Auxiliary Loss with Earth Mover's Distance)

지하철역 간의 이동 관계를 단순한 거리나 정적인 인접 행렬로 정의하기 어렵다는 점에서, 본 연구는 Earth Mover's Distance (EMD)를 활용하여 역 간의 승하차 분포 유사도를 정량적으로 모델링하였다.

EMD의 수식은 아래 식 (2)와 같다.

$$W[u, v] = \inf_{\gamma \in \Pi[u, v]} \int_x \int_y \gamma(x, y) d(x, y) dx dy \quad (2)$$

$\Pi[u, v]$  :  $u(x), v(y)$ 를 각각 주변 분포로 갖는 모든 결합 분포의 집합

$\gamma(x, y)d(x, y)$  : 질량의 비율과 이동비용

이 연구에서는 각 역의 승차 시계열 분포와 하차 시계열 분포 사이의 EMD를 계산하여, 역 간 이동 패턴의 유사도를 나타내는 EMD 참조 행렬을 사전 구축하였다. 이 행렬은 모델 학습 시 보조 손실로 사용된다.

모델이 직접 EMD를 계산하거나 파라미터로 학습하는 것이 아니라, 사전에 계산된 EMD 참조값을 목표로 두고, 모델의 예측 결과가 해당 분포 패턴에 가까워지도록 학습한다.

총 손실 함수는 (3)과 같이 정의된다.

$$L_{total} = L_{pred} + \lambda_{emd} L_{emd} \quad (3)$$

$L_{pred} = \|Y^{arr} - \hat{Y}^{arr}\|_{L1}$  : 기본 예측 손실

$L_{EMD}$  : EMD 기반 보조 손실

$\lambda_{EMD}$  : 가중 계수 본 논문에서는 0.5 사용

이어서 EMD 손실은 샘플링된 노드쌍 (i, j)에 대해 계산되며, 이산정보로 이루어진 승하차 분포를 각 쌍에 대해 soft Gaussian 기반 미분가능 히스토그램을 사용하여 DEP 입력과 ARR 예측의 분포를 비교한다. 식 (4)는 이를 나타낸 함수다.

$$L_{EMD} = \frac{1}{K} \sum_{(i,j) \in \rho} \left( \|h(X_i^{dep}) - h(\hat{Y}_j^{arr})\|_{L2} - EMD_{target(i,j)} \right)^2 \quad (4)$$

$EMD_{target(i,j)}$  : 사전 계산된 EMD 참조값(.csv파일)

$\rho$  : 무작위로 샘플링된 노드 쌍의 집합

$h(\cdot)$  : gaussian soft-bin으로 구성된 미분가능 히스토그램 함수 식(5)로 전개된다.

$$h_m(x) = \frac{\sum_{t=1}^T \exp\left(-\frac{(x_t - c_m)^2}{2\sigma^2}\right)}{\sum_{m=1}^M \sum_{t=1}^T \exp\left(-\frac{(x_t - c_m)^2}{2\sigma^2}\right)} + \epsilon \quad (5)$$

$c_m$  : m 번째 bin 중심값

$\sigma$  : 가우시안 너비

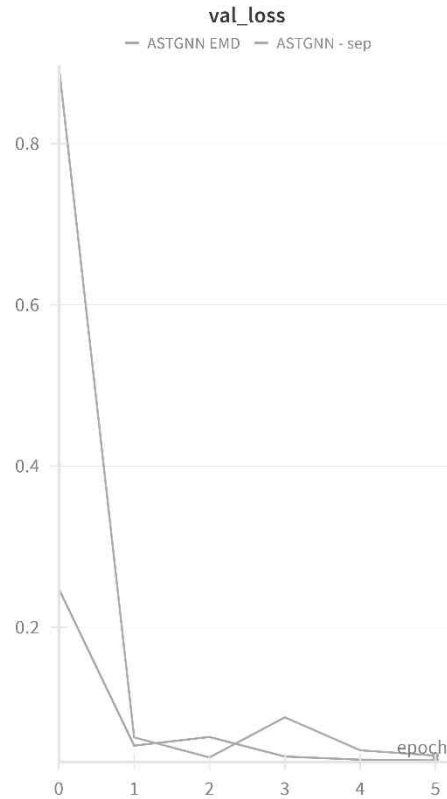
$T$  : 시계열 길이

$M$  : bin 개수

이 접근은 전체  $O(N^2)$  노드쌍 계산의 비용을 줄이기 위해 확률적 샘플링(stochastic approximation)을 적용하며, 매 스텝마다 K개의 다른 쌍을 샘플링함으로써 전체 분포 특성을 평균적으로 학습하도록 한다. 이를 통해  $O(K^2)$ 의 계산 비용을 가진다. ( $K \ll N$ )

Soft Gaussian 기반의 히스토그램은 hard binning 방식에 비해 미분 가능하며 경계 불연속이 없어, 역전파를 통한 gradient 전파가 가능하다. 이를 통해 모델 파라미터는 EMD 분포 패턴을 간접적으로 반영하게 된다.

결과적으로 본 연구의 모델은 단순한 시계열 예측을 넘어, 공간적 유사성(EMD)을 내재적으로 학습하여 <그림2>와 같이 평균적으로 보다 빠른 학습 수렴을 달성하였다.



<그림 2 : EMD 학습 그래프>

### III. 결론

#### 1. 연구 결론 및 향후 과제

본 연구에서는 서울 지하철 2호선을 대상으로 교통카드 태깅 기반의 승하차 데이터를 활용하여, 시공간 그래프 신경망에 기반한 혼잡도 예측 모델을 설계하였다. 특히, 모델의 계산 복잡도를 효율적으로 줄이면서도 예측 성능을 유지하기 위한 구조적 개선을 통해, 제한된 계산 자원에서도 안정적인 학습이 가능함을 보였다. 현재 연구는 2호선 내부의 구간별 혼잡도 예측에 국한되어 있으나, 제안한 경량화 구조는 전체 노선망 확장 시에도 효율적인 학습과 일반화 성능 향상을 기

대할 수 있다. 향후 연구에서는 서울시 전역의 다중 노선 데이터를 통합하여, 보다 정교하고 실시간성 있는 도시 교통 혼잡도 예측 모델로 발전시키는 것을 목표로 한다.

## 2. 모델 코드

<https://github.com/Arcana7642/ASTGNN-EMD>

## 참고문헌

1. S. Guo, Y. Lin, H. Wan, X. Li and G. Cong, "Learning Dynamics and Heterogeneity of Spatial-Temporal Graph Data for Traffic Forecasting," in IEEE Transactions on Knowledge and Data Engineering, vol. 34, no. 11, pp. 5415-5428, 1 Nov. 2022, doi: 10.1109/TKDE.2021.3056502
2. Lan, S., Ma, Y., Huang, W., Wang, W., Yang, H. & Li, P.. (2022). DSTAGNN: Dynamic Spatial-Temporal Aware Graph Neural Network for Traffic Flow Forecasting. Proceedings of the 39th International Conference on Machine Learning, in Proceedings of Machine Learning Research 162:11906-11917