# MaL-Air [Machine Learning for Air]

**Data fusion con metodi Machine Learning**

**Esempio di Applicazione sull'area urbana torinese**

Umberto Giuriato, Alessandro D'Ausilio, Camillo Silibello

# Data Fusion and Downscaling of Air Quality Deterministic models in the Turin area with Random Forest

U. Giuriato[1], A. D'Ausilio[1], C. Silibello[1], R. De Maria[2], S. Bande[2], C. Cascone[2], M. Maringo[2]

[1]ARIANET srl, 20159 Milano, Via Benigno Crespi 57, Italy.
[2]ARPA Piemonte, Regional Environmental Protection Agency of Piemonte, 10135 Torino, Italy.

**14th International Conference on Air Quality - Science and Application**
*Helsinki. 13 May - 17 May 2024*

# Introduction and motivation

**Supervised learning** can be successfully employed to spatialize concentrations measured by an **observation network**

**CTMs** concentration fields play the role of **predictors**, allowing the **data fusion** between deterministic models and actual observations

The combination with other High-Resolution predictors leads to the **downscaling** of concentration maps

A large and **representative sensor network** is fundamental to improve the learning process and enhance **generalization capabilities**

Since in real-case scenarios this is not always the case, we need to engineer a way to
- increase the representativeness of training set
- quantify the uncertanties

# Supervised training with Random Forest

- Training is performed on the monitoring station time series.
- The trained **Random Forest** model is then applied to each cell of the domain for each day to *infer concentration maps*

**PREDICTORS**



*Spatio-temporal*
CTM concentration fields
Leaf Area Index (LAI)

`.netcdf`

*Temporal (homogeneous)*
Periodic functions of
Julian day, day of week

*Spatial (stationary)*
Distance from roads
Population density
Land use
..........

`.netcdf`

**PRE-PROCESSING**

- **Extraction** of predictors fields at stations locations

- **Resampling** *in time*

- *Combination of predictors and target in the* **training set**

**TARGET**

Concentrations measured by the sensors

`.csv`

# Supervised training with Random Forest

- Training is performed on the monitoring station time series.
- The trained **Random Forest** model is then applied to each cell of the domain for each day to *infer concentration maps*



**PREDICTORS**

*Spatio-temporal*
CTM concentration fields
Leaf Area Index (LAI)

*Temporal (homogeneous)*
Periodic functions of
Julian day, day of week

*Spatial (stationary)*
Distance from roads
Population density
Land use
..........

**PRE-PROCESSING**
*training set*

MAL-Aria

**RANDOM FOREST TRAINING**

**TARGET**

Concentrations measured by
the sensors

**TRAINED MODEL**

(Hyperparameter tuning
and cross-validation also
in this phase)

# Supervised training with Random Forest

- Training is performed on the monitoring station time series.
- The trained **Random Forest** model is then applied to each cell of the domain for each day to *infer concentration maps*

**PREDICTORS**

**Spatio-temporal**
CTM concentration fields
Leaf Area Index (LAI)

**Temporal (homogeneous)**
Periodic functions of
Julian day, day of week

**Spatial (stationary)**
Distance from roads
Population density
Land use
..........

**TRAINED MODEL**

*RANDOM FOREST INFERENCE*

Data fusion maps between models and observation data

# Use case: SPoTT project

*Surveillance on POpulation health around the Turin waste-of-energy plant*

- Pollutants under study: **PM10** and **PM2.5**
- Year of study*: **2019**
- Target resolution: **200m**
- **14** monitoring stations
    - **8** measure PM25, **13** measure PM10
    - Majority of **urban** stations, just **2 rural** stations

| Nome Stazione | PM$_{10}$ | PM$_{2.5}$ |
|---|:---:|:---:|
| Baldissero T. (ACEA) | ● | |
| Beinasco TRM – Aldo Mei | ● | ● |
| Borgaro T. - Caduti | ● | ● |
| Carmagnola | ● | |
| Chieri - Bersezio | | ● |
| Collegno - Francia | ● | |
| Druento – Parco la Mandria | ● | |
| Leini (ACEA) - Grande Torino | ● | ● |
| Settimo T. - Vivaldi | ● | ● |
| Torino Consolata | ● | |
| Torino Grassi | ● | |
| Torino Lingotto | ● | ● |
| Torino Rebaudengo | ● | ● |
| Torino Rubino | ● | ● |

# Use case: SPoTT project

*Surveillance on POpulation health around the Turin waste-of-energy plant*

- Pollutants under study: **PM10** and **PM2.5**
- Year of study: **2019**
- Target resolution: **200m**
- **14** monitoring stations
  - **8** measure PM25, **13** measure PM10
  - Majority of **urban** stations, just **2 rural** stations

| Nome Stazione | PM$_{10}$ | PM$_{2.5}$ |
|---|:---:|:---:|
| Baldissero T. (ACEA) | ● | |
| Beinasco TRM – Aldo Mei | ● | ● |
| Borgaro T. - Caduti | ● | ● |
| Carmagnola | ● | |
| Chieri - Bersezio | | ● |
| Collegno - Francia | ● | |
| Druento – Parco la Mandria | ● | |
| Leini (ACEA) - Grande Torino | ● | ● |
| Settimo T. - Vivaldi | ● | ● |
| Torino Consolata | ● | |
| Torino Grassi | ● | |
| Torino Lingotto | ● | ● |
| Torino Rebaudengo | ● | ● |
| Torino Rubino | ● | ● |

# Spatio-temporal Predictors: FARM and Leaf Area Index

### FARM: PM$_{10}$ – daily mean



### FARM: PM$_{2.5}$ – daily mean



**FARM** simulations have been performed on the domain of interest at the resolution of **1km**

**Leaf Area Index** satellite images are monthly sampled at the resolution of **200m**

### FARM: NO$_2$ – daily mean



### FARM: O$_3$ – daily mean



### LAI - monthly

# Stationary (Spatial) Predictors

All the stationary predictors are at 200m resolution

# Stationary (Spatial) Predictors: Corine Land Cover

# PM$_{10}$ annual mean concentration



## Feature importances

| Feature | Importance (%) | Cumulative importance (%) |
|---|---|---|
| FARM PM10 | 0.38 | 0.38 |
| cos julian day | 0.16 | 0.54 |
| sin julian day | 0.16 | 0.7 |
| FARM NO2 | 0.08 | 0.78 |
| FARM O3 | 0.05 | 0.83 |
| Leaf Area Index | 0.04 | 0.87 |
| Dist primary roads | 0.02 | 0.89 |
| cos day of week | 0.02 | 0.91 |
| Dist second roads | 0.02 | 0.93 |
| Population | 0.01 | 0.94 |
| Elevation | 0.01 | 0.95 |
| Light at Night | 0.01 | 0.96 |
| sin day of week | 0.01 | 0.97 |

## Hyperparameters

| Hyperparameter | Value |
|---|---|
| N_trees | 800 |
| Max depth | 400 |
| Min samples split | 2 |
| Min samples leaf | 2 |
| Max features | 0.85 |
| Max leaf nodes | 1200 |

## Cross-validation RMSE

| Model | Value |
|---|---|
| FARM | 20.5 ± 0.74 |
| Random Forest | 7.2 ± 0.21 |

# PM$_{10}$ annual mean concentration



**FARM**

**Random Forest**

**Cross-validation RMSE**

| Model | Value |
|---|---|
| FARM | **20.5 ± 0.74** |
| Random Forest | **7.2 ± 0.21** |

- Globally, Random Forest inference *fixes the bias* of FARM predictions with respect to observations

- *Downscaling* due to stationary predictors is an *increment of concentration on roads* at the local scale

# Data augmentation: imputing PM$_{2.5}$ time series

## Lack of representativeness of observation network

- Only **8** stations measuring PM$_{2.5}$

- Lack of rural stations **Druento**, **Baldissero** and some **traffic** stations in Turin's urban area

It is possible to increase the training set size for PM$_{2.5}$ by **learning information from PM$_{10}$ stations**

- Preliminary random forest model (RF$_{st}$) is trained at stations where both species are sampled

- Then it is applied to infer timeseries of the ratio **PM$_{2.5}$ / PM$_{10}$** at stations where only PM$_{10}$ data are available



| Nome Stazione | PM$_{10}$ | PM$_{2.5}$ |
|---|---|---|
| Baldissero T. (ACEA) | ● | |
| Beinasco TRM – Aldo Mei | ● | ● |
| Borgaro T. - Caduti | ● | ● |
| Carmagnola | ● | |
| Chieri - Bersezio | | ● |
| Collegno - Francia | ● | |
| Druento – Parco la Mandria | ● | |
| Leini (ACEA) - Grande Torino | ● | ● |
| Settimo T. - Vivaldi | ● | ● |
| Torino Consolata | ● | |
| Torino Grassi | ● | |
| Torino Lingotto | ● | ● |
| Torino Rebaudengo | ● | ● |
| Torino Rubino | ● | ● |

$$\frac{PM_{2.5,obs}}{PM_{10,obs}} \sim RF_{st}\left(PM10_{obs}, \frac{PM_{2.5,\text{FARM}}}{PM_{10,\text{FARM}}}\right)$$

Using **PM$_{2.5}$ / PM$_{10}$** as target forbids unphysical values with **PM$_{2.5}$ > PM$_{10}$**

Stafoggia M, Johansson C, Glantz P, Renzi M, Shtein A, de Hoogh K, Kloog I, Davoli M, Michelozzi P, Bellander T. A Random Forest Approach to Estimate Daily Particulate Matter, Nitrogen Dioxide, and Ozone at Fine Spatial Resolution in Sweden. *Atmosphere*. 2020; 11(3):239. https://doi.org/10.3390/atmos11030239

# Data augmentation: imputing PM$_{2.5}$ time series



The "*surrogate*" PM$_{2.5}$ observation time series are highly correlated with PM$_{10}$ ones, with small variability in their ratio

# PM$_{2.5}$ annual mean concentration

**FARM**



**Random Forest WITH augmented data**



**Cross-validation RMSE**

| Model | Value |
| --- | --- |
| FARM | 14.0 ± 0.64 |
| Random Forest | 5.9 ± 0.52 |

**Random Forest WITHOUT augmented data**



**Cross-validation RMSE**

| Model | Value |
| --- | --- |
| FARM | 14.0 ± 0.64 |
| Random Forest | 6.02 ± 0.48 |

Augmenting data based on FARM PM ratio and PM$_{10}$ observations:

- Thins out the difference with FARM field, still healing the bias
- Keeps similarity with PM$_{10}$ map, without concentration increase outside urban area

# BONUS: Coverage Score

A **Coverage Score** can be defined to quantify how much the predictors' distribution in the target domain matches the one *"seen"* by the observation network.

*In every cell ij:* $\quad C_{ij} = \sum_{k=1}^{N_{predictors}} \text{importance}_k \cdot Q_{ij}$

*Where $Q_{ij}$ is:*
- *1 if the values belong to the distribution sampled by the observation network*
- *0 otherwise*

### *Distribution of FARM PM10 values (stations vs all cells)*

### *Mean Coverage map for PM$_{10}$ - Turin*

# BONUS: Coverage Score for PM₂.₅



Mean Coverage map for PM₂.₅ WITHOUT augmented data

Mean Coverage map for PM₂.₅ WITH augmented data

Adding the *"surrogate"* stations improves the coverage in regions outside Turin, increasing the reliability of the model.

Thank you for your attention !

Umberto Giuriato   Alessandro D'Ausilio   Camillo Silibello

# Nested K-fold cross validation



Validation
(outer loop)

Hyperparameter
tuning
(inner loop)

**Validation of the Random Forest training is performed with a K-fold nested cross-validation**

The evaluation metric used is *RMSE*

- Split dataset into K equal folds
- One fold is used as **test** set and K-1 remaining as **training**
  - In each inner training fold, perform an extra K-fold splitting and use it to perform **hyperparameter tuning**
- Train the model on the training set for each iteration independently
- Validate the model on the test set for each iteration
- The final score is the average obtained from all K iterations to get the final score

# PM$_{2.5}$ annual mean concentration



## Feature importances

| Feature | Importance (%) | Cumulative importance (%) |
|---|---|---|
| sin julian day | 0.25 | 0.25 |
| cos julian day | 0.23 | 0.48 |
| FARM O3 | 0.22 | 0.70 |
| FARM PM25 | 0.12 | 0.82 |
| FARM NO2 | 0.03 | 0.85 |
| cos day of week | 0.02 | 0.87 |
| Leaf Area Index | 0.02 | 0.89 |
| Dist primary roads | 0.01 | 0.90 |
| Elevation | 0.01 | 0.91 |
| Population | 0.01 | 0.92 |
| sin day of week | 0.01 | 0.93 |
| Dist second roads | 0.01 | 0.94 |

## Hyperparameters

| Hyperparameter | Value |
|---|---|
| N_trees | 400 |
| Max depth | 30 |
| Min samples split | 2 |
| Min samples leaf | 2 |
| Max features | 0.85 |
| Max leaf nodes | 700 |

## Cross-validation RMSE

| Model | Value |
|---|---|
| FARM | 14.0 ± 0.64 |
| Random Forest | 5.9 ± 0.52 |