

First Meeting Presentation

Project A2

Giulia Mezzadri Ettore Modina Oswaldo Jesus Morales Lopez
Federico Angelo Mor Abylaikhan Orynassar Federica Rena

Politecnico of Milano
Bayesian Statistics course

Project revision of
November 23, 2023

Presentation Flow

- ① Project Overview
 - Goal and Definition
 - Data Exploration
- ② Models
 - General model construction
 - Models from literature
- ③ Expected workflow
- ④ References
- ⑤ Text Examples
- ⑥ Table and Figure Examples
- ⑦ Mathematics
- ⑧ Referencing

The goal of the project

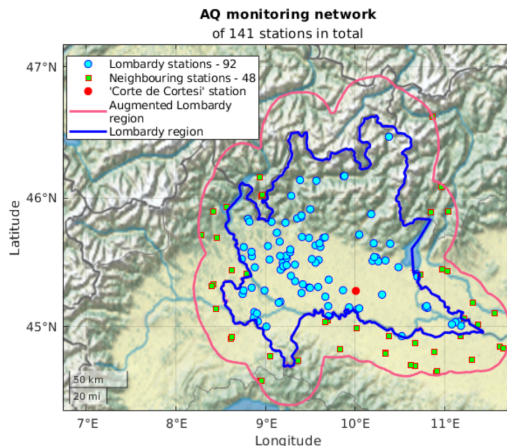
Goal: Clustering weekly data of one year of PM10 (plus covariates)

Dataset: AGRIMONIA project, at

<https://zenodo.org/records/7563265>

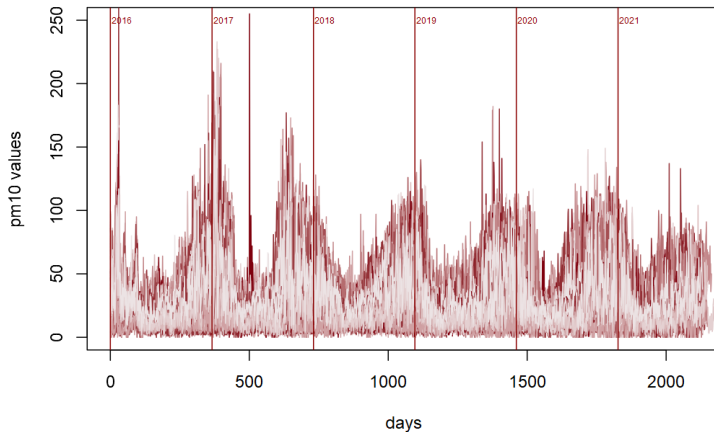
Spatial Exploration

We have 141 stations, which recorded data for 6 years (from 01/01/2016 to 31/12/2021).



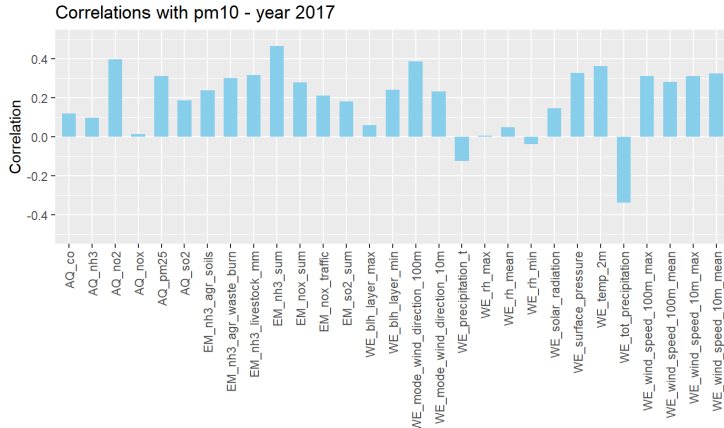
Temporal Exploration

We have 141 stations, which recorded data
for 6 years (from 01/01/2016 to 31/12/2021).

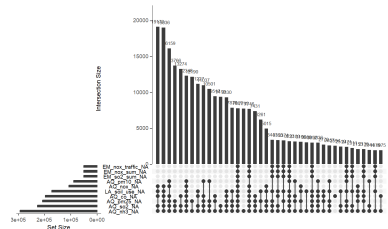


Correlation: pm10 and others

We compute the spearman's correlation index between the pm10 and the other covariates, which in the functional framework quantifies with a value in $[-1,1]$ the tendency of 2 r.v. X_t and Y_t to be perfect monotone functions one of each other



Plots to identify a general pattern on missing values



Combinations of variables with the most missing values

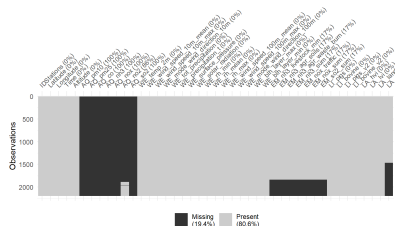
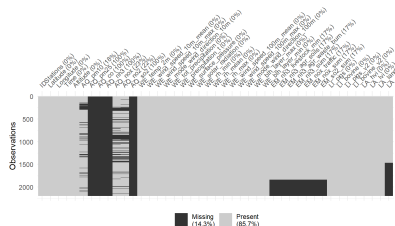
Missing value exploration

Plot of missing data divided by station to see particular patterns to help select variables

Possible to remove some stations that are not measuring PM10 values

Some columns present missing values in most columns, so we can remove corresponding covariates

Some missing observations concentrated in specific periods such as during last years



Missing value exploration

```
[1] "1274"      "1888"      "1881"      "581"       "584"
[6] "514"       "529"       "539"       "544"       "545"
[11] "551"       "552"       "573"       "588"       "682"
[16] "683"       "686"       "687"       "626"       "652"
[21] "656"       "657"       "665"       "672"       "682"
[26] "694"       "696"       "698"       "699"       "788"
[31] "782"       "787"       "STA.IT1518A" "STA.IT1751A" "STA.IT1924A"
[36] "STA.IT2282A"
```

First selection to remove non informative stations, then count of missing values for column, removing those above a chosen threshold

```
IDstations      Latitude
0               0
Longitude        Time
0               0
Altitude         AQ_pm10
0               12719
AQ_pm25         AQ_co
136782         136568
AQ_rh3          AQ_nov
227540         72988
AQ_rh2          AQ_soi
21266         156285
WE_temp_2m      WE_wind_speed_10m_mean
0               0
WE_wind_speed_10m_max WE_wind_direction_10m
0               0
WE_tot_precipitation WE_precipitation_t
0               0
WE_surface_pressure WE_solar_radiation
0               0
WE_rh_min       WE_rh_mean
0               0
WE_rh_max       WE_wind_speed_100m_mean
0               0
WE_wind_speed_100m_max WE_wind_direction_100m
0               0
WE_bih_layer_max WE_bih_layer_min
0               0
EH_rh3_livestock_gm EH_rh3_agr_sol16
38325            38325
EH_rh3_agr_waste_burn EH_rh3_sum
38325            38325
EH_nov_traffic      EH_nov_sum
38325            38325
EH_soi2_sum         LT_plgs
38325            6576
LT_bovine_v1       LT_plgs_v2
6576              0
LT_bovine_v2       LA_hvi
0                 0
LA_lv1             LA_land_use
0                 0
LA_soil_use        136656
```

Models: a complex task

Considering the nature of the data, our models should account for different levels of information:

- spatial context
- temporal context
- covariates

which is a not-so-trivial task.

Now we see the general incremental idea to build such models.

Purely spatial model

- We have n distinct locations $\mathbf{s}_1, \dots, \mathbf{s}_n$, where $\mathbf{s}_i = (\text{lat}, \text{long})$.
- There we record data y_i and (possibly) covariates \mathbf{x}_i , for $i = 1, \dots, n$. The goal is to define a model for partitioning them into k groups.
- So we define $\rho = \{S_1, \dots, S_k\}$ the cluster set variable (with $S_h \subseteq \{1, \dots, n\}$ for $h = 1, \dots, k$).
An equivalent formulation is possible through some cluster indicator variables c_1, \dots, c_n ; where $c_i = h \iff i \in S_h$ for $i = 1, \dots, n$.
- In general, the law for ρ follows a spatial Product Partition Model (sPPM):

$$p_\rho(\tilde{\rho}) \propto \prod_{h=1}^{k_n} C(\tilde{S}_h, \mathbf{s}_h^*)$$

where $\tilde{\rho} = \{\tilde{S}_1, \dots, \tilde{S}_k\}$, $\mathbf{s}_h^* = \{\mathbf{s}_i : i \in \tilde{S}_h\}$ and $C(\tilde{S}_h, \mathbf{s}_h^*)$ is a cohesion function.

Purely spatial model

- We have n distinct locations $\mathbf{s}_1, \dots, \mathbf{s}_n$, where $\mathbf{s}_i = (\text{lat}, \text{long})$.
- There we record data y_i and (possibly) covariates \mathbf{x}_i , for $i = 1, \dots, n$. The goal is to define a model for partitioning them into k groups.
- So we define $\rho = \{S_1, \dots, S_k\}$ the cluster set variable (with $S_h \subseteq \{1, \dots, n\}$ for $h = 1, \dots, k$).
An equivalent formulation is possible through some cluster indicator variables c_1, \dots, c_n ; where $c_i = h \iff i \in S_h$ for $i = 1, \dots, n$.
- In general, the law for ρ follows a spatial Product Partition Model (sPPM):

$$p_\rho(\tilde{\rho}) \propto \prod_{h=1}^{k_n} C(\tilde{S}_h, \mathbf{s}_h^*)$$

where $\tilde{\rho} = \{\tilde{S}_1, \dots, \tilde{S}_k\}$, $\mathbf{s}_h^* = \{\mathbf{s}_i : i \in \tilde{S}_h\}$ and $C(\tilde{S}_h, \mathbf{s}_h^*)$ is a cohesion function.

Purely spatial model

- We have n distinct locations $\mathbf{s}_1, \dots, \mathbf{s}_n$, where $\mathbf{s}_i = (\text{lat}, \text{long})$.
- There we record data y_i and (possibly) covariates \mathbf{x}_i , for $i = 1, \dots, n$. The goal is to define a model for partitioning them into k groups.
- So we define $\rho = \{S_1, \dots, S_k\}$ the cluster set variable (with $S_h \subseteq \{1, \dots, n\}$ for $h = 1, \dots, k$).

An equivalent formulation is possible through some cluster indicator variables c_1, \dots, c_n ; where $c_i = h \iff i \in S_h$ for $i = 1, \dots, n$.

- In general, the law for ρ follows a spatial Product Partition Model (sPPM):

$$p_\rho(\tilde{\rho}) \propto \prod_{h=1}^{k_n} C(\tilde{S}_h, \mathbf{s}_h^*)$$

where $\tilde{\rho} = \{\tilde{S}_1, \dots, \tilde{S}_k\}$, $\mathbf{s}_h^* = \{\mathbf{s}_i : i \in \tilde{S}_h\}$ and $C(\tilde{S}_h, \mathbf{s}_h^*)$ is a cohesion function.

Purely spatial model

- We have n distinct locations $\mathbf{s}_1, \dots, \mathbf{s}_n$, where $\mathbf{s}_i = (\text{lat}, \text{long})$.
- There we record data y_i and (possibly) covariates \mathbf{x}_i , for $i = 1, \dots, n$. The goal is to define a model for partitioning them into k groups.
- So we define $\rho = \{S_1, \dots, S_k\}$ the cluster set variable (with $S_h \subseteq \{1, \dots, n\}$ for $h = 1, \dots, k$).
An equivalent formulation is possible through some cluster indicator variables c_1, \dots, c_n ; where $c_i = h \iff i \in S_h$ for $i = 1, \dots, n$.
- In general, the law for ρ follows a spatial Product Partition Model (sPPM):

$$p_\rho(\tilde{\rho}) \propto \prod_{h=1}^{k_n} C(\tilde{S}_h, \mathbf{s}_h^*)$$

where $\tilde{\rho} = \{\tilde{S}_1, \dots, \tilde{S}_k\}$, $\mathbf{s}_h^* = \{\mathbf{s}_i : i \in \tilde{S}_h\}$ and $C(\tilde{S}_h, \mathbf{s}_h^*)$ is a cohesion function.

Spatial and temporal model

We have n distinct locations $\mathbf{s}_1, \dots, \mathbf{s}_n$, where $\mathbf{s}_i = (\text{lat}, \text{long})$.

There we record data y_i and (possibly) covariates \mathbf{x}_i , for $i = 1, \dots, n$.

The goal is to define a model for partitioning them into k_t groups, with t spanning over $1, \dots, T$.

So we define $\rho_t = \{S_{1,t}, \dots, S_{k_t,t}\}$ the cluster set variable (with $S_{h,t} \subseteq \{1, \dots, n\}$ for $h = 1, \dots, k_t$).

In general, the law for ρ_t follows a spatial Product Partition Model (sPPM) updated to account for the time relation (stPPM); meaning that we need a formulation of a joint probability model for ρ_1, \dots, ρ_T .

This update can be explicated for example by

- supposing a Markov Chain structure, letting ρ_t depend just on ρ_{t-1} ;
- introducing some cluster reallocation variable $\gamma_{i,t} \in \{0, 1\}$.

Model 1



Garritt L. Page, Fernando A. Quintana, David B. Dahl (2022)

Dependent Modeling of Temporal Sequences of Random Partitions. [Journal of Computational and Graphical Statistics](#), 31:2, 614-627.

$$\begin{aligned}
 Y_{it} | \boldsymbol{\mu}_t^*, \boldsymbol{\sigma}_t^{2*}, \mathbf{c}_t &\stackrel{\text{ind}}{\sim} \mathcal{N}(\mu_{c_{it}t}^*, \sigma_{c_{it}t}^{2*}) \quad i = 1, \dots, n \quad \text{and} \quad t = 1, \dots, T \\
 (\mu_{jt}, \sigma_{jt}) | \theta_t, \tau_t^2 &\stackrel{\text{ind}}{\sim} \mathcal{N}(\theta_t, \tau_t^2) \times \mathcal{U}(0, A_\sigma) \quad j = 1, \dots, k_t \\
 (\theta_t, \tau_t) &\stackrel{\text{iid}}{\sim} \mathcal{N}(\phi_0, \lambda^2) \times \mathcal{U}(0, A_\tau) \quad t = 1, \dots, T \\
 (\phi_0, \lambda) &\sim \mathcal{N}(m_0, s_0^2) \times \mathcal{U}(0, A_\lambda) \\
 \{\mathbf{c}_t, \dots, \mathbf{c}_T\} &\sim \text{tRPM}(\boldsymbol{\alpha}, M) \quad \text{with} \quad \alpha_t \stackrel{\text{iid}}{\sim} \text{Beta}(a_\alpha, b_\alpha)
 \end{aligned}$$

Adding covariates

Given a partition, now we can easily design models which also account for covariates. For example we can update the previous model into

$$\begin{aligned}
 Y_{it} | \beta_t^*, \sigma_t^{2*}, \mathbf{c}_t &\stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{x}_{it}^T \beta_{c_{it}t}^*, \sigma_{c_{it}t}^{2*}) \\
 (\beta_{jt}, \sigma_{jt}) | \theta_t, \tau_t^2 &\stackrel{\text{ind}}{\sim} \mathcal{N}(\theta_t, \tau_t^2) \times \mathcal{U}(0, A_\sigma) \\
 &\vdots
 \end{aligned}$$

or we can further characterize the time dependance with some AR(.) model

$$\begin{aligned}
 Y_{it} | Y_{it-1}, \beta_t^*, \sigma_t^{2*}, \mathbf{c}_t &\stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{x}_{it}^T \beta_{c_{it}t}^* + Y_{it-1} \eta_i, \sigma_{c_{it}t}^{2*} (1 - \eta_i)^2) \\
 Y_{i1} | \beta_1^*, \sigma_1^{2*}, \mathbf{c}_1 &\sim \mathcal{N}(\mathbf{x}_{i1}^T \beta_{c_{i1}1}^{2*}, \sigma_{c_{i1}1}^{2*}) \\
 &\vdots
 \end{aligned}$$

Model 2



Mozdzen A., Cremaschi A., Cadonna A., Guglielmi A., Kastner G. (2022)

Bayesian modeling and clustering for spatio-temporal areal data: An application to Italian unemployment. *Spatial Statistics* 52, 100715.

$$Y_{it} | \mathbf{x}_{it}, \beta_{s_i}^*, w_{it}, \sigma^2, s_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{x}_{it}^T \beta_{s_i}^* + w_{it}, \sigma^2)$$

$$\mathbf{w}_t | \mathbf{w}_{t-1}, \xi_s^*, \mathbf{s}, \tau^2, \rho, W \sim \mathcal{N}_I(\text{diag}(\xi_s^*) \mathbf{w}_{t-1}, \tau^2 Q(\rho, W)^{-1})$$

$$\mathbf{w}_1 | \tau^2, \rho, W \sim \mathcal{N}_I(\mathbf{0}, \tau^2 Q(\rho, W)^{-1})$$

$$\sigma^2 \sim \text{Inv-Gamma}(a_{\sigma^2}, b_{\sigma^2})$$

$$\tau^2 \sim \text{Inv-Gamma}(a_{\tau^2}, b_{\tau^2})$$

$$\rho \sim \text{Beta}(\alpha_\rho, \beta_\rho)$$

$$\mathbf{s} | \alpha \sim \text{PólyaUrn}(\mathbf{s} | \alpha)$$

$$\alpha \sim \text{Gamma}(a_\alpha, b_\alpha)$$

$$\phi_1^*, \dots, \phi_{K_I}^* | \mu_\beta, \Sigma_\beta, \alpha_\xi, \beta_\xi \stackrel{\text{iid}}{\sim} P_0, \quad \phi_j^* = (\beta_j^*, \xi_j^*) \quad j = 1, \dots, K_I$$

$$P_0(d\phi^*) = \mathcal{N}_{p+1}(d\beta^* | \mu_0, \Sigma_0) \text{Beta}_{(-1,1)}(d\xi^* | \alpha_\xi, \beta_\xi)$$

Model 3

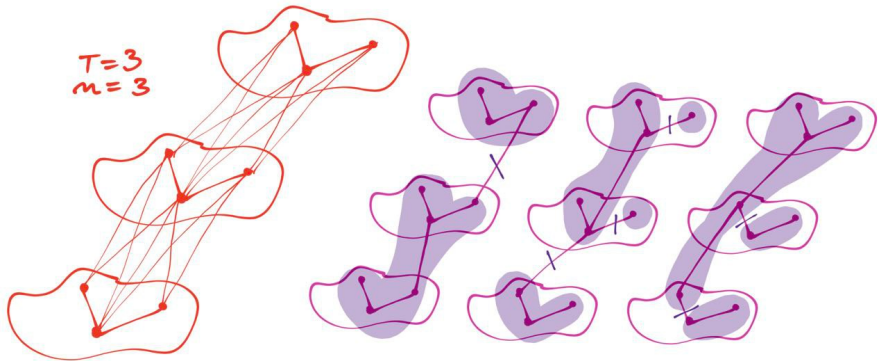


Leonardo V. Teixeira, Renato M. Assunção, Rosangela H. Loschi (2019)
Bayesian Space-Time Partitioning by Sampling and Pruning Spanning Trees.
[Journal of Machine Learning Research 20, 85, 1–35.](#)

This model works on a graph structure, which incorporates together space and time. That is, from data Y_{jt} for $j = 1, \dots, n$ and $t = 1, \dots, T$, we now move to Y_i for $i \in I = \{1, \dots, nT\}$ by stacking T times the spatial map. So we have a graph $\mathcal{G} = (V, E)$ of nT nodes and edges built according to time and space connections.

The idea is to search a partition $\pi = \{\mathcal{G}_1, \dots, \mathcal{G}_c\}$ of I (with \mathcal{G}_k subgraphs for \mathcal{G}), on randomly selected spanning trees \mathcal{T} of \mathcal{G} , on which we set cluster-specific parameters $\beta_{\mathcal{G}_1}, \dots, \beta_{\mathcal{G}_c}$.

Model 3



Model 3



Leonardo V. Teixeira, Renato M. Assunção, Rosangela H. Loschi (2019)
 Bayesian Space-Time Partitioning by Sampling and Pruning Spanning Trees.
[Journal of Machine Learning Research](#) 20, 85, 1–35.

$$Y_i | \mathcal{T}, \pi, \beta_{\mathcal{G}_k} \stackrel{\text{iid}}{\sim} f(Y_i | \beta_{\mathcal{G}_k}; \mathbf{x}_i) \quad i \in \mathcal{G}_k$$

$$\beta_{\mathcal{G}_1}, \dots, \beta_{\mathcal{G}_c} | \mathcal{T}, \pi \sim \prod_{k=1}^c f(\beta_{\mathcal{G}_k})$$

$$p(\pi = \{\tilde{\mathcal{G}}_1, \dots, \tilde{\mathcal{G}}_c\} | \mathcal{T}) \sim \prod_{k=1}^c \kappa(\tilde{\mathcal{G}}_k)$$

$$\mathcal{T} \sim \mathcal{U}(\text{St}(\mathcal{G}))$$

Expected workflow

- PM10 covariates analysis (physical/chemical relation)
- First models implementation and comparison
- Implementation of variations of those simple models, or more complex models from literature
- Gif/Video interactive plots for displaying results

References



Garritt L. Page, Fernando A. Quintana, David B. Dahl (2022)
Dependent Modeling of Temporal Sequences of Random Partitions. *Journal of Computational and Graphical Statistics*, 31:2, 614-627.



Mozdzen A., Cremaschi A., Cadonna A., Guglielmi A., Kastner G. (2022)
Bayesian modeling and clustering for spatio-temporal areal data: An application to Italian unemployment. *Spatial Statistics* 52, 100715.



Leonardo V. Teixeira, Renato M. Assunção, Rosangela H. Loschi (2019)
Bayesian Space-Time Partitioning by Sampling and Pruning Spanning Trees. *Journal of Machine Learning Research* 20, 85, 1-35.

Paragraphs of Text

Sed iaculis **dapibus gravis**. Morbi sed tortor erat, nec interdum arcu. Sed id lorem lectus. Quisque viverra augue id sem ornare non aliquam nibh tristique. Aenean in ligula nisl. Nulla sed tellus ipsum. Donec vestibulum ligula non lorem vulputate fermentum accumsan neque mollis.

Sed diam enim, sagittis nec condimentum sit amet, ullamcorper sit amet libero. Aliquam vel dui orci, a porta odio.
— *Someone, somewhere. . .*

Nullam id suscipit ipsum. Aenean lobortis commodo sem, ut commodo leo gravis vitae. Pellentesque vehicula ante iaculis arcu pretium rutrum eget sit amet purus. Integer ornare nulla quis neque ultrices lobortis.

Lists

Bullet Points and Numbered Lists

- Lorem ipsum dolor sit amet, consectetur adipiscing elit
 - Aliquam blandit faucibus nisi, sit amet dapibus enim tempus
 - Lorem ipsum dolor sit amet, consectetur adipiscing elit
 - Nam cursus est eget velit posuere pellentesque
 - Nulla commodo, erat quis gravida posuere, elit lacus lobortis est, quis porttitor odio mauris at libero
-
- ① Nam cursus est eget velit posuere pellentesque
 - ② Vestibulum faucibus velit a augue condimentum quis convallis nulla gravida

Blocks of Highlighted Text

Block Title

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Integer lectus nisl, ultricies in feugiat rutrum, porttitor sit amet augue.

Example Block Title

Aliquam ut tortor mauris. Sed volutpat ante purus, quis accumsan.

Alert Block Title

Pellentesque sed tellus purus. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos himenaeos.

Suspendisse tincidunt sagittis gravida. Curabitur condimentum, enim sed venenatis rutrum, ipsum neque consectetur orci.

Multiple Columns

Subtitle

Heading

- ① Statement
- ② Explanation
- ③ Example

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Integer lectus nisl, ultricies in feugiat rutrum, porttitor sit amet augue. Aliquam ut tortor mauris. Sed volutpat ante purus, quis accumsan dolor.

Table

Subtitle

Treatments	Response 1	Response 2
Treatment 1	0.0003262	0.562
Treatment 2	0.0015681	0.910
Treatment 3	0.0009271	0.296

Table: Table caption

Figure

Figure: Space for a possible image.

Definitions & Examples

Definition

A **prime number** is a number that has exactly two divisors.

Example

- 2 is prime (two divisors: 1 and 2).
- 3 is prime (two divisors: 1 and 3).
- 4 is not prime (**three** divisors: 1, 2, and 4).

You can also use the theorem, lemma, proof and corollary environments.

Theorem, Corollary & Proof

Theorem (Mass–energy equivalence)

$$E = mc^2$$

Corollary

$$x + y = y + x$$

Proof.

$$\omega + \phi = \epsilon$$



Equation

$$\cos^3 \theta = \frac{1}{4} \cos \theta + \frac{3}{4} \cos 3\theta \quad (1)$$

Verbatim

Example (Theorem Slide Code)

```
\begin{frame}  
\frametitle{Theorem}  
\begin{theorem}[Mass--energy equivalence]  
$E = mc^2$  
\end{theorem}  
\end{frame}
```

Slide without title.

Citing References

An example of the `\cite` command to cite within the presentation:

This statement requires citation [Smith, 2022, Kennedy, 2023].

References



John Smith (2022)

Publication title

Journal Name 12(3), 45 – 678.



Annabelle Kennedy (2023)

Publication title

Journal Name 12(3), 45 – 678.

Acknowledgements

Smith Lab

- Alice Smith
- Devon Brown

Cook Lab

- Margaret
- Jennifer
- Yuan

Funding

- British Royal Navy
- Norwegian Government

The End

Questions? Comments?