

Standard course – Modelling

Lesson SM1.3 – How to build a data-driven model

Marta Magnani – Institute of Geoscience and Earth Resources (CNR, Italy)



CERTH
CENTRE FOR
RESEARCH & TECHNOLOGY
HELLAS



CREAF



Consiglio Nazionale
delle Ricerche

Data-driven models

Empirical models are a useful tool in the pioneering exploration of processes. Data-driven models are based on empirical relationships between measured (or estimated) variables and aim at representing a variable (Y) as a function of other measured variables ($X_i, i=1,2,3,\dots$), called predictors or explanatory variables. The steps to build a data-driven model are:

1. Identify different model formulations, possibly inspired by literature or by the characteristics of data.
1. Test different candidate predictors
1. Use statistical criteria to select the best model:
 - ✓ the predictors should not be cross-correlated (excluded by a partial correlation analysis)
 - ✓ the regression should be significant (i.e. P-value of predictors lower than a chosen threshold of significance)
 - ✓ the modelled variable (\hat{Y}) should be as close as possible to the measured one (Y_t).
 - ✓ the model should be efficient (i.e. it maximizes representativeness while minimizing the number of parameters), for instance minimizing the Akaike Information Criterion (Akaike, 1974; Burnham and Anderson, 2002)

QUANTITATIVE MODELS FOR TERRESTRIAL CO₂ FLUXES

Measurements at the Alpine site

Arctic measurement site at Bayelva, Spitsbergen, Norway. In situ samplings of:

1. CO₂ fluxes with portable accumulation chambers: net flux (NEE) & emissions (ER) → uptake ($GPP = NEE - ER$)
1. Meteoroclimatic variables (Air pressure, Air temperature, air moisture, solar irradiance, soil temperature, soil moisture)
1. Vegetation green fractional cover
1. Vegetation classification
 - Vascular vs non-vascular
 - Different vascular species

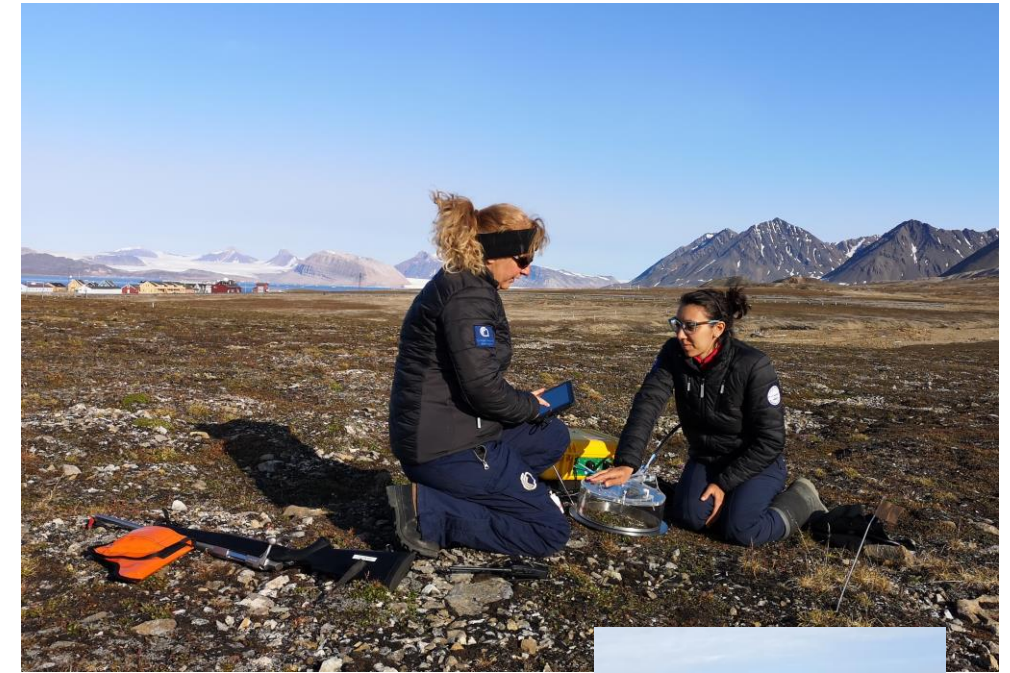
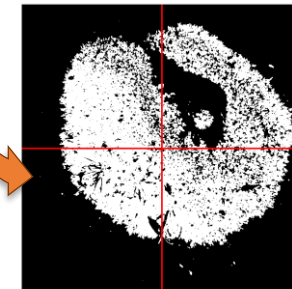
4. Vegetation classification



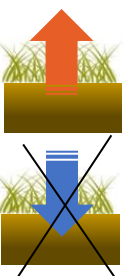
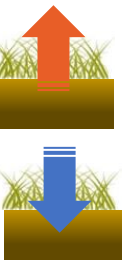
3. Green Fractional Cover



method: Liu & Pattey, 2010



1. Flux measurement NEE & ER



Experimental design

Flux temporal variability

Measurements performed at a fixed location (i.e. point) for 24 hours every 2 hours.



Point scale
(0cm)

Flux spatiotemporal variability

fluxes measured at different points randomly distributed over the site. Measurements repeated in consecutive days



Site
scale
(100m)

Models of the Arctic fluxes

Well-Known
drivers

Temporal variability

$$\begin{aligned} \text{Ecosystem Respiration} &= a e^{b \text{Temperature}} \\ \text{Gross Primary Production} &= \frac{F \alpha \text{Radiation}}{F + \alpha \text{Radiation}} \end{aligned}$$



Point scale
(0cm)

Spatiotemporal variability

+ ???
+ ???



Site scale
(100m)

At a fixed point along 24h: the temporal variability of fluxes was well reproduced by the classical drivers, temperature for ER and solar irradiance for GPP.

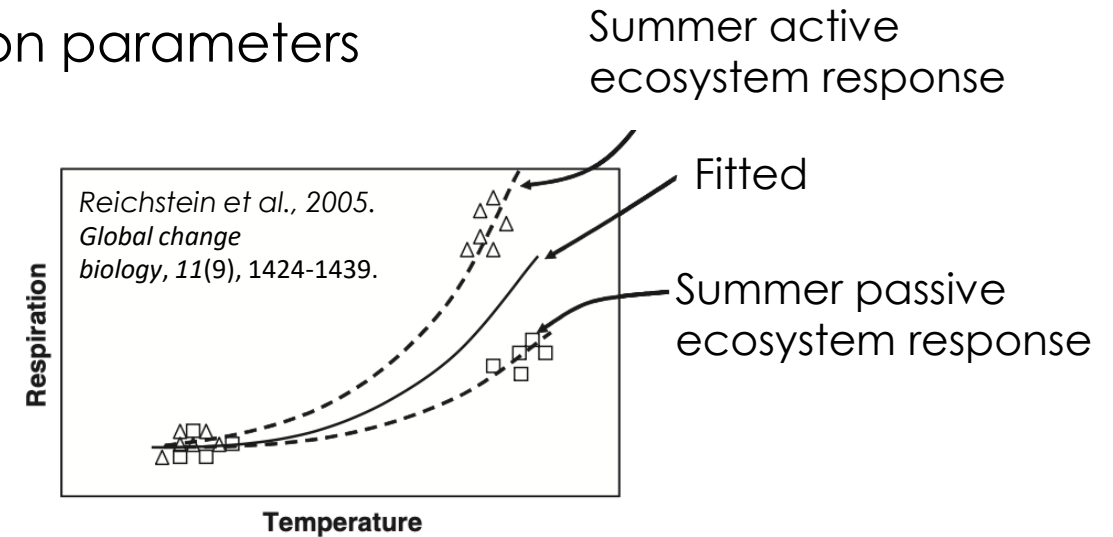
zooming out to the site scale: only a small part of flux variability was explained by the classical drivers and multi regression models were needed.

Additional variables

Idea: additional drivers modify the regression parameters

→ Taylor expansions

$$\frac{ER}{GPP} = \frac{a(x_1, x_2) e^{b(x_1, x_2)T}}{F(y_1, y_2) + \alpha(y_1, y_2)R}$$



Where x_1 and x_2 are chosen among the measured variables. Statistical selection of both model and variables. The model thus obtained is:

$$\left. \begin{array}{l} \text{emissions} \\ \text{uptake} \end{array} \right\} \left\{ \begin{array}{l} ER = (a_0 + a_1 VWC + a_2 GFC) e^{b_0 T} + \varepsilon \\ GPP = \frac{F_0 \alpha_0 R}{F_0 + \alpha_0 R} (A_0 + A_1 VWC + A_2 GFC) + \varepsilon \end{array} \right|$$

T =temperature
 VWC =soil moisture
 GFC =Green Fractional Cover
 R =solar irradiance



Multi regression models

Well-Known
drivers

$$\text{EcosystemRespiration} = ae^{b\text{Temperature}}$$

Lloyd & Tylor, 1994. Funct. Ecol. 315–323.

$$\text{GrossPrimaryProduction} = \frac{F\alpha\text{Radiation}}{F + \alpha\text{Radiation}}$$

Ruimy et al., 1995. Adv. Ecol. Res. 26, 1–68

Inserting additional descriptors in multi regression models

$$ER = (a| |0 + a_1\text{GreenFractionalCover} + a_2\text{soilmoisture})e^{bTemp}$$

$$GPP = (A| |0 + A_1\text{GreenFractionalCover} + A_2\text{soilmoisture})\frac{F\alpha Rad}{F + \alpha Rad}$$

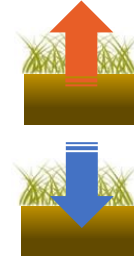
Ref to: Magnani et al., 2022. Scientific reports, 12(1), 1-14.

In both the Alpine and the Arctic case study the soil moisture (VWC) and a vegetation descriptor, namely the *GFC* here and *DOY* - that was interpreted as a proxy of phenology - in the Alps, were identified as additional predictors of the fluxes beyond the classical drivers.

Role of different plant types

$$ER = (a_0 + a_1 \text{GreenFractionalCover} + a_2 \text{soilmoisture}) e^{b \text{Temp}}$$

$$GPP = (A_0 + A_1 \text{GreenFractionalCover} + A_2 \text{soilmoisture}) \frac{F \alpha \text{Rad}}{F + \alpha \text{Rad}}$$



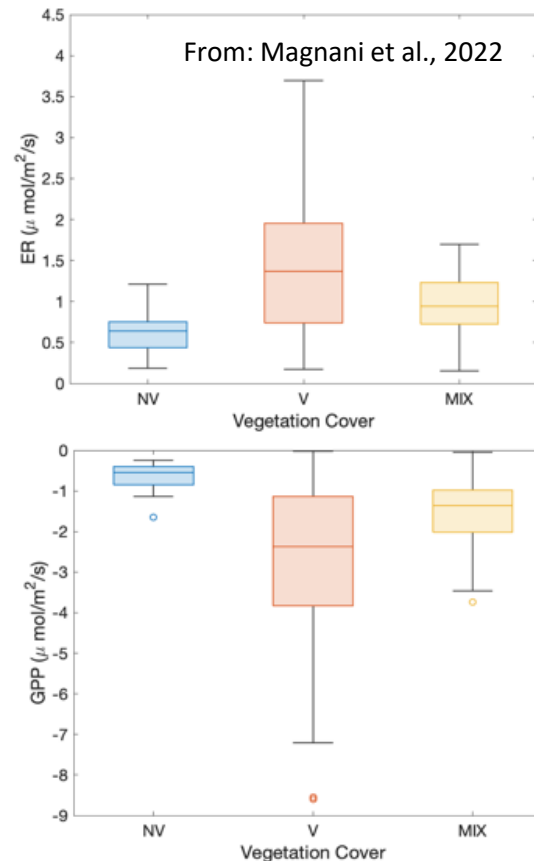
Ph: MS Giamberini

Non vascular



Ph: MS Giamberini

vascular



Exploring more in depth the role of different plant types, the dataset was divided into three subsets according to the prevailing vegetation cover type observed at each sampling point. The three classes:

V: vascular vegetation (herbs and prostrated plants),

NV: non-vascular vegetation (lichens, mosses and bacterial soil crust) MIX: mix of vascular and non-vascular vegetation (no clear prevalence).

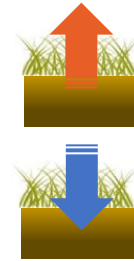
Reparametrizing the models for these three classes, significant differences were observed in the flux dependence on the *GFC* (i.e. significantly different a_1 and A_1 parameters) between classes.

→ **different physiological functioning of vegetation groups**

Role of different species

$$ER = (a_0 + a_1 \text{GreenFractionalCover} + a_2 \text{soilmoisture}) e^{b \text{Temp}}$$

$$GPP = (A| \quad |0 + A_1 \text{GreenFractionalCover} + A_2 \text{soilmoisture}) \frac{F \alpha \text{Rad}}{F + \alpha \text{Rad}}$$



Vascular class was the most variable in terms of green fractional cover → species-specific samplings.

The models were parameterized using the species-specific samplings and significant differences were observed in the flux dependence on the *GFC* between two groups of species: *SI-DR* vs *SL-SX*, while *CX* showed hybrid characteristics between these two groups of species.

→ **different physiological functioning of vegetation species groups**

5 classes of most abundant vascular plants identified



Silene acaulis



Carex Spp.



Salix polaris



Dryas octopetala

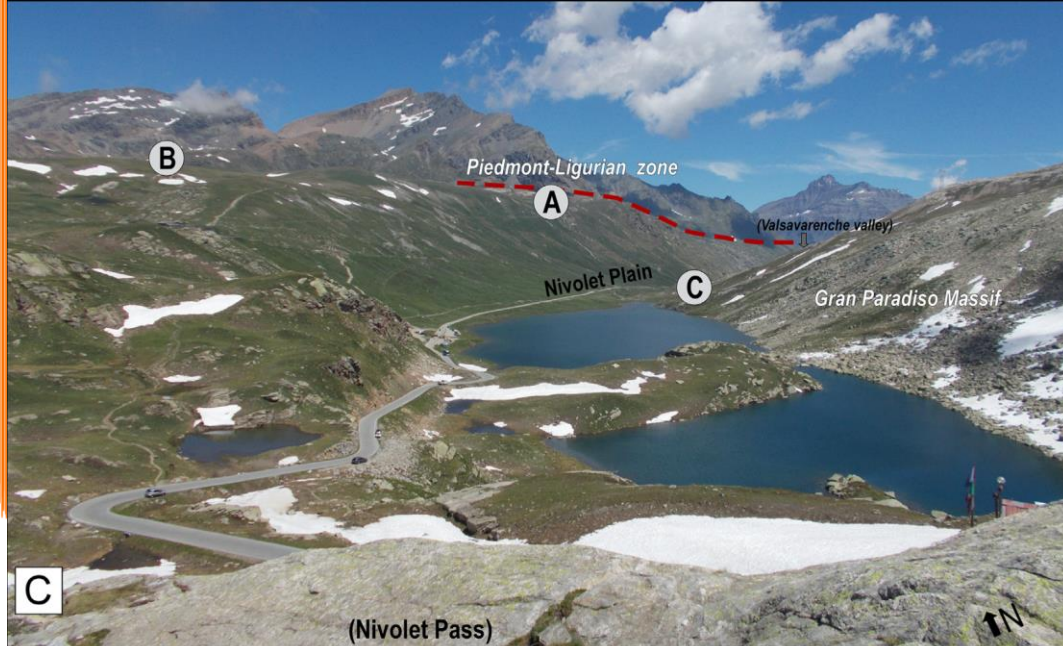


Saxifraga oppositifolia

Measurements at the Alpine sites



Alpine measurements site at the Nivolet plain, in the Gran Paradiso National Park, western Italian Alps (45° 28' 43" N 7° 08' 32" E).



Measurements in summers 2017->2021

- ❖ Carbon dioxide fluxes
- ❖ Soil temperature and moisture
- ❖ Air temperature, moisture and pressure
- ❖ Solar irradiance
- ❖ Vegetation information from pictures

On **3 sampling plots** characterized by different parental material:

Site A) carbonate rocks

Site B) glacial till

Site C) gneiss rock

Multi regression models

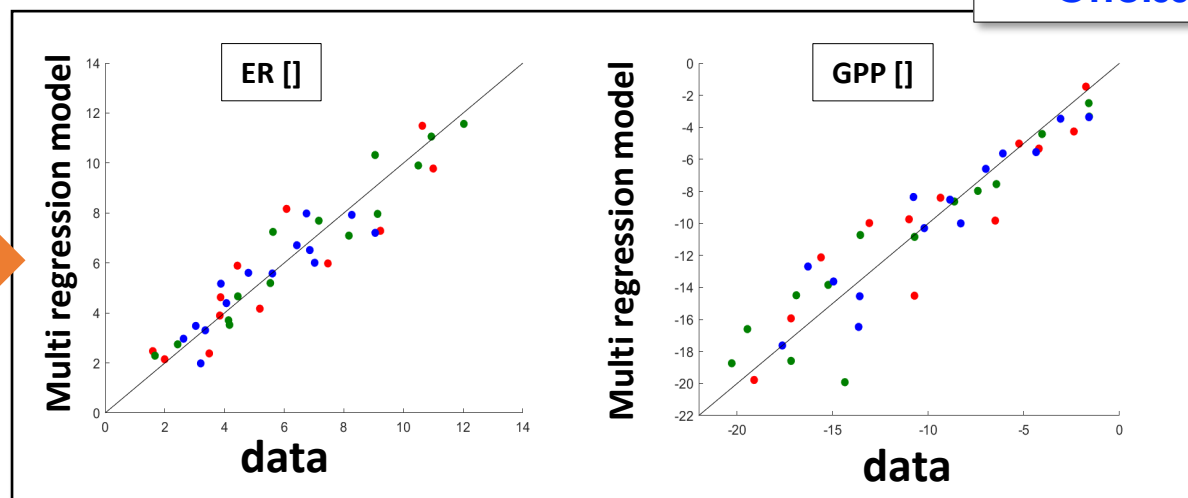
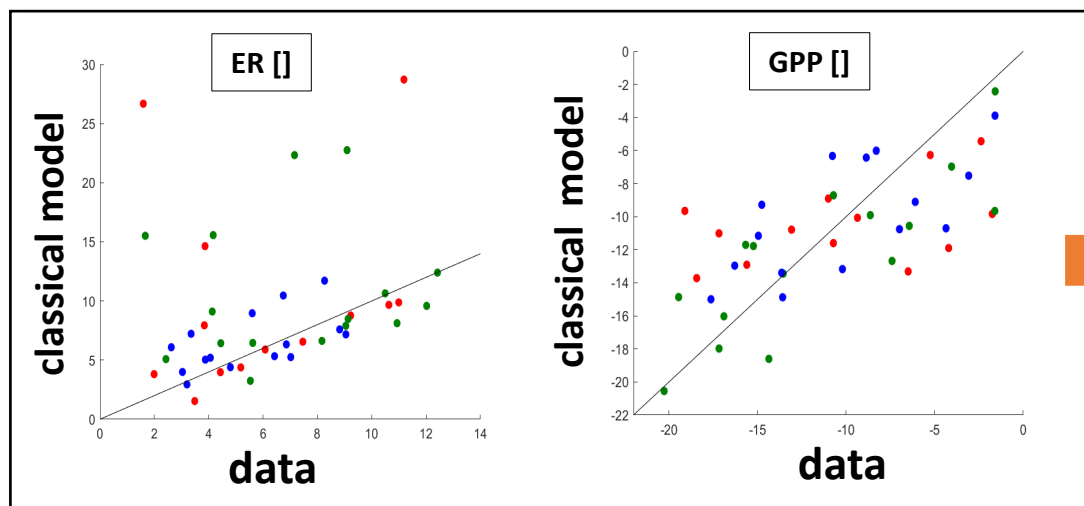
emissions
uptake

$$\begin{cases} ER = ae^{bT} \\ GPP = \frac{F\alpha R}{F + \alpha R} \end{cases}$$



$$\left\{ ER = (a_0 + a_1VWC + a_2Pr + a_3DOY e^{b_0T} \right\} \left| GPP = \frac{F_0\alpha_0}{F_0 + \alpha_0} \right.$$

• Carb.
• Glac.
• Gneiss



T=temperature; R=solar irradiance; VWC=soil moisture; Pr=air pressure; DOY = Day Of the Year

- **Soil moisture** limits both carbon emissions and uptake
- **Air pressure** may have a leaking effect on emissions
- Significant differences among parameters from different sites → **Soil characteristics** are additional constraints
- The **Day of the Year** explained a large part of fluxes variability → What it is a proxy for?

Ref to Magnani et al., STotEn, 2020