Inverse Modeling of Surface Carbon Fluxes

Please read Peters et al (2007) and Explore the CarbonTracker website



- Given N measurements y_i for different values of the independent variable x_i , find a slope (m) and intercept (b) that describe the "best" line through the observations
 - Why a line?
 - What do we mean by "best"
 - How do we find m and b?
- Compare predicted values to observations, and find m and b that fit best
- Define a total error (difference between model and observations) and minimize it!







Minimizing the Error (cont'd)
(4)
$$m = \frac{N\sum x_i y_i - \sum x_i \sum y_i}{N\sum x_i^2 - (\sum x_i)^2}$$
 Plug (4) into (2) and simplify:
 $b = \frac{\sum y_i}{N} - \frac{\sum x_i}{N} \left\{ \frac{N\sum x_i y_i - \sum x_i \sum y_i}{N\sum x_i^2 - (\sum x_i)^2} \right\}$
(5) $b = \frac{\sum y_i \sum x_i^2 - \sum x_i \sum x_i y_i}{N\sum x_i^2 - (\sum x_i)^2}$
Now have simple "Least Squares" formulae for "best" slope and intercept given a set of observations

Geometric View of Linear Regression

- Any vector $\vec{a} = (x, y, z)$ can be written as a linear combination of the *orthonormal basis* set $(\hat{i}, \hat{j}, \hat{k})$
- set $(\hat{i}, \hat{j}, \hat{k})$ • This is accomplished by taking the dot product (or inner product) of the vector with each basis vector to determine the components in each basis direction
- Linear regression involves a 2D mapping of an observation vector into a different vector space
- More generally, this can involve an arbitrary number of basis vectors (dimensions)

Linear Regression Revisited

This notation can be rewritten in subscript notation: $d_i = \sum_{j=1}^{M} G_{ij} m_j$

and applied to a familiar problem. Imagine that there are 2 data points $(d_1,\,d_2)$ and 2 model parameters $(m_1,\,m_2).$

Then the system of equations could be explicitly written as:

 $\mathbf{d}_1 = \mathbf{G}_{11}\mathbf{m}_1 + \mathbf{G}_{12}\mathbf{m}_2$

 $\mathbf{d}_2 = \mathbf{G}_{21}\mathbf{m}_1 + \mathbf{G}_{22}\mathbf{m}_2$

Or in matrix form $\vec{d} = G\vec{m}$

With two points, this is just two slope-intercept form equations:

$$y_1 = m x_1 + b$$

 $y_2 = m x_2 + b$

This is an "even-determined" problem - there is exactly enough information to determine the model parameters precisely, there is only one solution, and there is zero prediction error.















Tracer	Description	Seasonal Terrestrial Biosphere		
		Fung	CASA	SiB2
T ₃	Subarctic Atlantic	-0.40	-0.30	+0.05
T ₄	North Atlantic Gyre	+0.69	+0.37	+0.57
T ₅	North Pacific	+1.26	+0.64	+0.80
T ₆	Equatorial Oceans	+1.62	+1.62	+1.62
T7	Southern Gyres	-0.11	+0.15	+1.56
T ₈	Antarctic Oceans	-0.50	-0.43	-0.30
Total Ocean Flux		+2.56	+2.06	+4.31
T ₂	Tropical Deforestation	4.75	+6.04	12.98
T ₁₀	NPP-based CO2 Fertilization	-4.22	-4.41	-9.56
T ₁₁	Water stress CO2 Fertilization	-2.69	-5.62	-14.41
T ₁₂	Temperate Forest Sink	+0.39	+0.14	+3.83
T ₁₃	Boreal Forest Sink	-8.29	-3.68	-4.45
T ₁₄	Tundra Sink	+4.59	+2.55	+4.19
Total Terrestrial Flux		-5.47	-4.97	-7.22
r.m.s. error (ppm)		0.39	0.35	0.41











Uncertainty in Flux Estimates $\vec{m}_{est} = \vec{m}_p + \left(\hat{G}^T \hat{C}_d^{-1} \hat{G} + \hat{C}_m^{-1}\right)^{-1} \hat{G}^T \hat{C}_d^{-1} \left(\vec{d}_{obs} - \hat{G} \vec{m}_p\right)$ $C_m^* = \left(\hat{G}^T C_d^{-1} \hat{G} + C_m^{-1}\right)^{-1}$ • A posteriori estimate of uncertainty in the estimated fluxes

- Depends on transport (G) and a priori uncertainties in fluxes (C_m) and data (C_d)
- Does not depend on the observations per se!











Sensitivity to Priors

- Flux estimates and a posteriori uncertainties for data-constrained regions (N and S) are very insensitive to priors
- Uncertainties in poorly constrained regions (tropical land) very sensitive to prior uncertainties
- As priors are loosened, dipoles develop between poorly constrained regions