

# Cross-validation methods for quality control, cloud screening, etc.

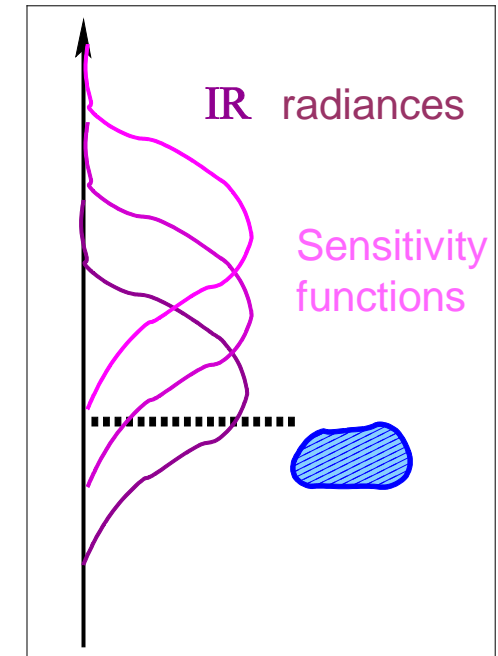
Olaf Stiller, Deutscher Wetterdienst

*Are observations consistent*

➤ *with the other observations ?*

*given the*

- *background*
- *assumed error covariances*
- *observation operator*



*Which observations are affected by the cloud???*

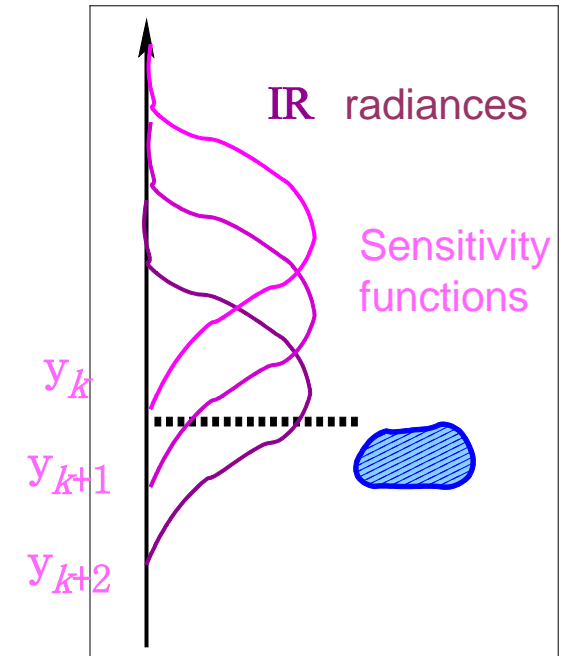
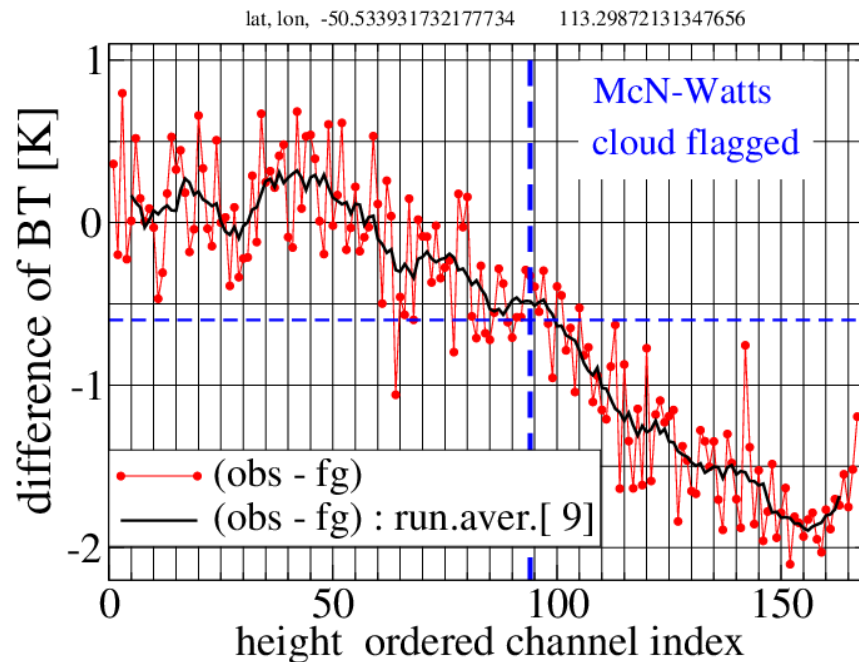
# Inspiration:

## McNally&Watts scheme



Diagnose clouds from the observations [ i.e., from (obs-fg) ]

1. Look whether a FoV is cloudy: (obs-fg) threshold
2. Find upper edge of cloud : gradient criterion



Diagnose clouds from the observations [ i.e., from (obs-fg) ]

Question: Can we do this more systematically?

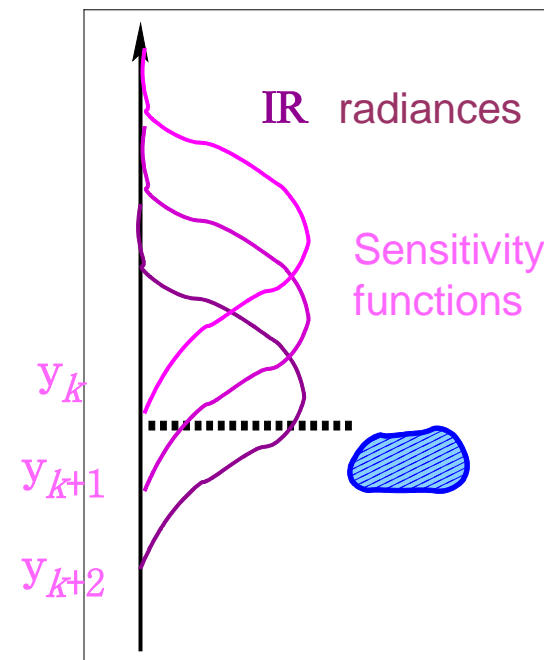
Aim : Identify observations which are not consistent with .....

*Are observations consistent*

➤ *with the other observations ?*

*given the*

- *background*
- *assumed error covariances*
- *observation operator*



*Which observations are affected by the cloud???*

# Assumed uncertainties in data assimilation



## Assumptions about Obs and FG errors:

$$J(\mathbf{x}) = \frac{1}{2} \left[ \mathbf{x}^T \mathbf{B}^{-1} \mathbf{x} + (\mathbf{y}^o - \mathbf{H}\mathbf{x})^T \mathbf{R}^{-1} (\mathbf{y}^o - \mathbf{H}\mathbf{x}) \right]$$

## Are FG departures consistent with these assumptions?

- $\mathbf{y}^o = \mathbf{Y}^o - \mathbf{Y}^b$   
obs - first guess

$$\langle (\mathbf{y}^o)^T \mathbf{y}^o \rangle = \mathbf{H}^T \mathbf{B} \mathbf{H} + \mathbf{R}$$

## Checking diagonal:

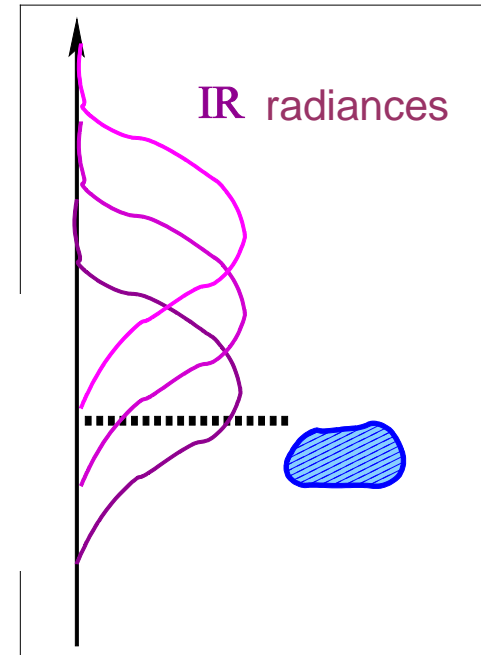
$$\langle (\mathbf{y}_k^o)^2 \rangle = [\mathbf{H}^T \mathbf{B} \mathbf{H} + \mathbf{R}]_{kk} = \sigma_k^2$$

Conditional probability of the observations  $\mathbf{y}_k^o$  (given the background):

$$P(\mathbf{y}_k^o | \mathbf{X}^b) \propto \exp -\frac{1}{2} \left( \frac{\mathbf{y}_k^o}{\sigma_k} \right)^2$$

Cross-Validation with background (standard Quality Control check):

$$n \text{ sigma check: } \left| \frac{\mathbf{y}_k^o}{\sigma_k} \right| < n$$



# Assumed uncertainties in data assimilation



## Assumptions about Obs and FG errors:

$$J(\mathbf{x}) = \frac{1}{2} \left[ \mathbf{x}^T \mathbf{B}^{-1} \mathbf{x} + (\mathbf{y}^o - \mathbf{H}\mathbf{x})^T \mathbf{R}^{-1} (\mathbf{y}^o - \mathbf{H}\mathbf{x}) \right]$$

## Are FG departures consistent with these assumptions?

- $\mathbf{y}^o = \mathbf{Y}^o - \mathbf{Y}^b$   
obs - first guess

$$\langle (\mathbf{y}^o)^T \mathbf{y}^o \rangle = \mathbf{H}^T \mathbf{B} \mathbf{H} + \mathbf{R}$$

## Checking diagonal:

$$\langle (\mathbf{y}_k^o)^2 \rangle = [\mathbf{H}^T \mathbf{B} \mathbf{H} + \mathbf{R}]_{kk} = \sigma_k^2$$

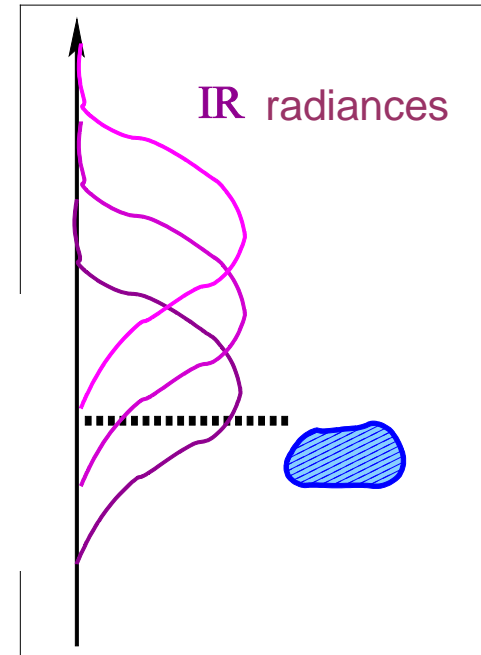
Conditional probability of the observations  $\mathbf{y}_k^o$  (given the background):

$$P(\mathbf{y}_k^o | \mathbf{X}^b) \propto \exp -\frac{1}{2} \left( \frac{\mathbf{y}_k^o}{\sigma_k} \right)^2$$

## Decompose observations: $\{\mathbf{y}_{\tau C}^o, \mathbf{y}_{\tau}^o\}$

Conditional probability of observations  $\mathbf{y}_{\tau}^o$   
(given the background and observations  $\mathbf{y}_{\tau C}^o$ ):

$$P(\mathbf{y}_{\tau}^o | \mathbf{y}_{\tau C}^o, \mathbf{X}^b) \propto \exp -\frac{1}{2} \left\{ (\mathbf{y}_{\tau}^o - \bar{\mathbf{y}}_{\tau})^T \mathbf{D}_{\tau} (\mathbf{y}_{\tau}^o - \bar{\mathbf{y}}_{\tau}) \right\}$$



# Special case: Observations can be ordered



$$P(\mathbf{y}_k | \mathbf{y}_{\{l < k\}}, \mathbf{x}^b) = \mathcal{N}^{-1} \exp -\frac{1}{2} \{ \mathbf{Y}_k^2 \}$$

$$\mathbf{Y} = (\mathbf{T}_l)^{-1} \mathbf{y}$$

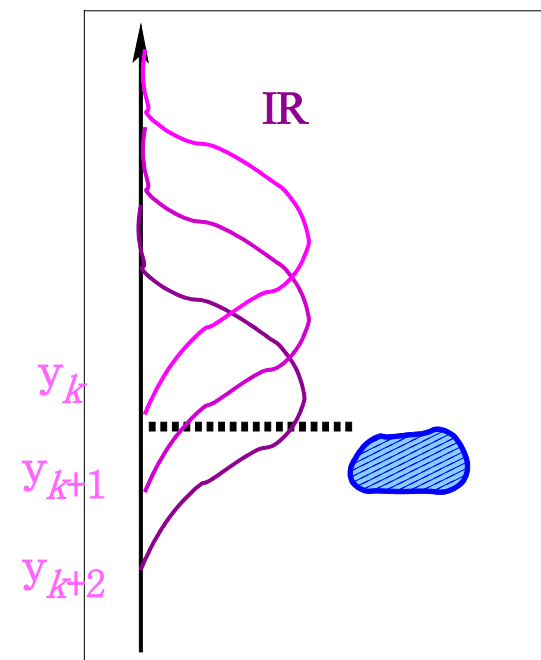
$$Y_k = \frac{y_k - y_k^{a*}}{\sqrt{\epsilon_k^{obs} + \epsilon_k^{a*}}}$$

analysis considering only obs  $y_l$  with  $l < k$

error of this analysis

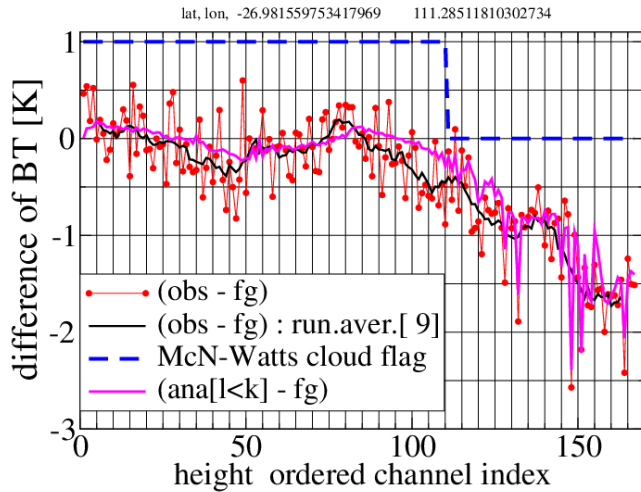
Cholesky decomposition:

$$\mathbf{T}_L \mathbf{T}_U = \begin{pmatrix} t_{11} & 0 & & & & \\ t_{21} & t_{22} & 0 & & & \\ & \cdot & \cdot & 0 & & \\ t_{k1} & t_{k2} & & t_{kk} & 0 & \\ & & \cdot & \cdot & \cdot & \\ & & & & \cdot & 0 \\ t_{p1} & & & & & t_{pp} \end{pmatrix} \begin{pmatrix} t_{11} & t_{21} & t_{k1} & t_{p1} \\ 0 & t_{22} & t_{k2} & \\ & 0 & \cdot & \cdot \\ & & 0 & t_{kk} \\ & & & 0 & \cdot \\ & & & & \cdot & \cdot \\ & & & & & 0 & t_{pp} \end{pmatrix} = [\mathbf{R} + \mathbf{H}\mathbf{B}\mathbf{H}^T]$$



# Application: IASI cloud screening

(see also 8P.05)



$$Y_k = \frac{y_k - y_k^{a*}}{\sqrt{\epsilon_k^{obs} + \epsilon_k^{a*}}}$$

analysis considering only obs  $y_l$  with  $l < k$

**Problem:** Standard deviation dominated by obs error  
Single observation not sensitive enough

**Need to detect systematic perturbations**

Consider joint probability:

$$P(\mathbf{y}_k, \mathbf{y}_{k+1}, \dots, \mathbf{y}_{k+s} | \mathbf{y}_{\{l < k\}}, \mathbf{x}^b) \propto \exp -\frac{1}{2} \{ \mathbf{Y}_k^2 + \mathbf{Y}_{k+1}^2 + \dots + \mathbf{Y}_{k+s}^2 \}$$

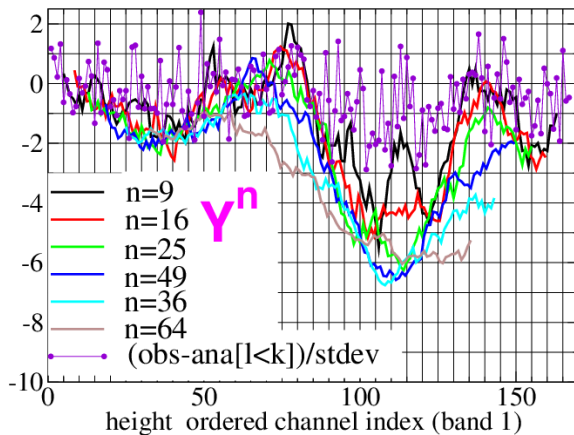
$\mathbf{Y}_k^n = \sum_{j=k}^{k+n} \mathbf{Y}_j / \sqrt{n}$  is also: stochastic variable with variance 1

Generalization (for any vector  $\vec{h}_l$ ):

$$\vec{Y} \rightarrow \tilde{Y}_l \equiv \frac{\vec{h}_l * \vec{Y}}{\|\vec{h}_l\|}$$

stochastic variable with variance 1

Targeted approach: project on most relevant directions  $\vec{h}_l$



Generalization (for any vector  $\vec{h}_l$ ):

$$\vec{Y} \rightarrow \tilde{Y}_l \equiv \frac{\vec{h}_l * \vec{Y}}{\|\vec{h}_l\|} \quad \text{stochastic variable with variance 1}$$

Let:  $c_{cfr}$  be a **model** state variable for **cloud fraction** in a layer  
 $\mathbf{H}_{cfr}$  corresponding part of **observation operator** matrix

$$\mathbf{H}_T = \left( \begin{array}{c|c} \mathbf{H} & \begin{matrix} \cdot \\ \mathbf{H}_{cfr} \\ \cdot \end{matrix} \end{array} \right) \quad \mathbf{B}_T = \left( \begin{array}{c|c} \mathbf{B} & \begin{matrix} \cdot \\ 0 \\ \cdot \end{matrix} \\ \hline \text{---} & \sigma_{cfr} \end{array} \right)$$

Then, in the limit of large  $\sigma_{cfr}$ , one finds:

$$c_{cfr}^a \rightarrow \left[ \mathbf{H}_{cfr}^T [\mathbf{R} + \mathbf{H}\mathbf{B}\mathbf{H}^T]^{-1} \mathbf{H}_{cfr} \right]^{-1} \mathbf{H}_{cfr}^T [\mathbf{R} + \mathbf{H}\mathbf{B}\mathbf{H}^T]^{-1} [\mathbf{y}^o - H(\mathbf{x}_c^b)] = \frac{\mathbf{h}^T \mathbf{Y}}{[\mathbf{h}^T \mathbf{h}]}$$

cloud fraction in layer  $k$

$$\frac{c_{cfr}^a[k]}{\sqrt{\langle (c_{cfr}^a[k])^2 \rangle}} = \frac{\vec{h}_k * \vec{Y}}{\|\vec{h}_k\|}$$

stochastic variable with variance 1

$$\mathbf{Y} = \mathbf{T}_L^{-1} [\mathbf{y}^o - H(\mathbf{x}_c^b)]$$

$$\mathbf{h} = \mathbf{T}_L^{-1} \mathbf{H}_{cfr}$$



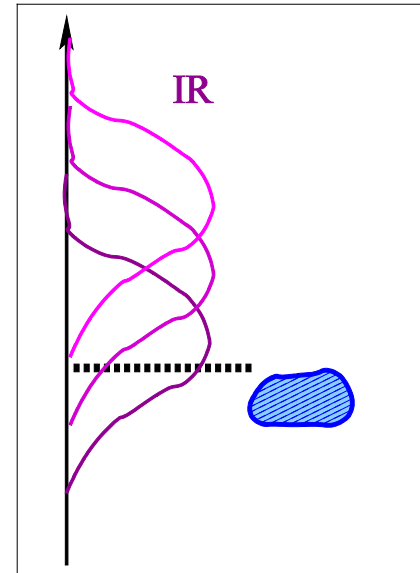
- A **cross validation method** for observations has been developed which

- works within the probabilistic framework of the DA system:  $H^T B H + R$ 
  - **disadvantage**: employed error matrixes are far from perfect
  - **advantage** : method will develop and improve systematically with improved DA systems

- is **cheap enough** to be run in **preprocessing step**
- requires that observation operators sufficiently overlap
  - good for IASI

- Diagnostics have to be tailored for systematic perturbations

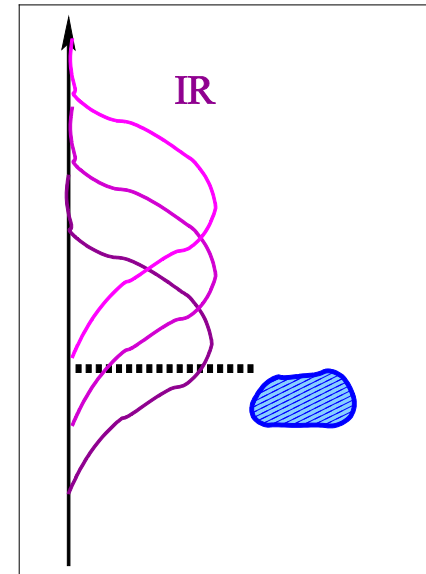
- project on relevant directions  $\vec{h}_i$ 
  - employed error matrixes are (*probably*) not good enough to flag more generally perturbed observations



- Which influences can be diagnosed from obs-fg increments?

- impact has to be generally strong (scale separation – weak signal must be rare)
- $\|\vec{h}_i\|$  must be large for typical signal
  - very low clouds can not be detected from IASI radiances

- The cross validation method
  - is planned to be run as a preprocessing system
    - flagging of bad observation **before** they enter into the analysis
  - possibly within a **1D Var** preprocessing step (important for strongly nonlinear observation as, e.g., the water vapor channels of IASI)
  - will profit from improved **B** matrix from Ensemble Kalman Filter
- The cross validation method may be useful for testing also other influences
  - which the observation operator does not represent properly
  - like, e.g., surface emissivity
- CV diagnostics good for comparing compatability of different observation types
  - collecting statistics of targeted diagnostics





Thank you for listening

