# Data Analysis and Fitting:
# Errors and Goodness of Fit

## Ashton S. Reimer

[1]Center for Geospace Studies
SRI International
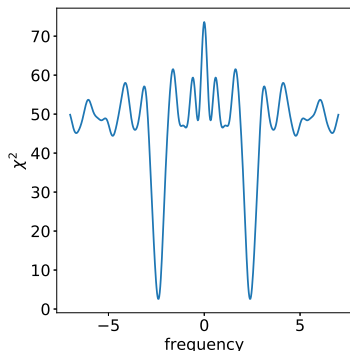
July 2021

## Chi-Squared

We can use least-squares to solve inverse problems:

$$\chi^2(\mathsf{p}) = [\mathsf{y} - f(\mathsf{p})]^T \, \Sigma_e^{-1} \, [\mathsf{y} - f(\mathsf{p})]$$

where $\hat{\mathsf{p}}_{LS}$ are the "best-fit" model parameters, those that minimizes $\chi^2(\mathsf{p})$

Great! But:

- What are the errors in the fitted parameters $\hat{\mathsf{p}}_{LS}$?
- Is the fit meaningful? Does the model accurately reproduce the measurements?

## Error Propagation (e.g. Linear Least-Squares)

For a linear forward model:

$$y = f(p) + e \qquad f(p) = Hp$$

The Least-Squares solution is:

$$\hat{p}_{LS} = \left[ H^T \Sigma_e^{-1} H \right]^{-1} H^T \Sigma_e^{-1} y$$

Given that jointly Gaussian random variables have the following property:

$$Y = AX \quad \Rightarrow \quad \Sigma_Y = A \Sigma_X A^T$$

it can be shown that:

$$\Sigma_{\hat{p}_{LS}} = \left[ H^T \Sigma_e^{-1} H \right]^{-1}$$

## Error Propagation (e.g. Nonlinear Least Squares)

For a non-linear forward model, guess a $p_i$, linearize, and step towards minimum:

$$y = f(p) + e \qquad f(p_i + \Delta p) \approx f(p_i) + J_i \Delta p \qquad J_i = \frac{\partial f}{\partial p_i}$$

J is known as the Jacobian:

$$J = \begin{pmatrix} \frac{\partial f_0}{\partial p_0} & \frac{\partial f_0}{\partial p_1} & \cdots & \frac{\partial f_0}{\partial p_{N-1}} \\ \frac{\partial f_1}{\partial p_0} & \frac{\partial f_1}{\partial p_1} & \cdots & \frac{\partial f_1}{\partial p_{N-1}} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial f_{M-1}}{\partial p_0} & \frac{\partial f_{M-1}}{\partial p_1} & \cdots & \frac{\partial f_{M-1}}{\partial p_{N-1}} \end{pmatrix}$$

J is $M \times N$ (tall and skinny)

Non-linear fitting process:

- iterate until $p_{i+1} = \hat{p}_{LS}$: that which minimizes $\chi^2$
- The covariance of $\hat{p}_{LS}$ is:

$$\Sigma_{\hat{p}_{LS}} = \left[ J^T \Sigma_e^{-1} J \right]^{-1}$$

Note the similarity to the linear case!

# Error Propagation

The covariance of the fitted parameters is the covariance of the input data propagated through the least-squares operation:

$$\Sigma_{\hat{p}_{LS}} = \left[ J^T \Sigma_e^{-1} J \right]^{-1}$$

"Error bars" for fitted parameters:

- Assumption: measurement errors are **Gaussian** distributed with covariance $\Sigma_e$, denoted $\mathcal{N}(0, \Sigma_e)$
- The "errors" in the fitted parameters are related to confidence intervals
- Confidence intervals are constructed from $\Sigma_{\hat{p}_{LS}}$
- $\Sigma_{\hat{p}_{LS}}$ may look reasonable, even if the fit is meaningless

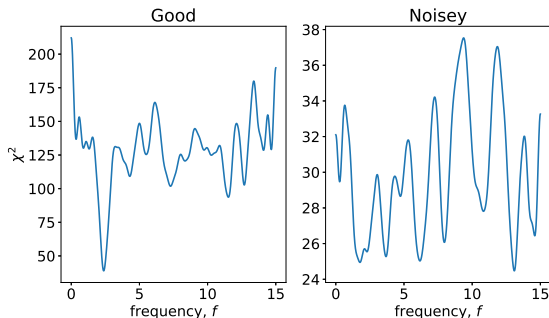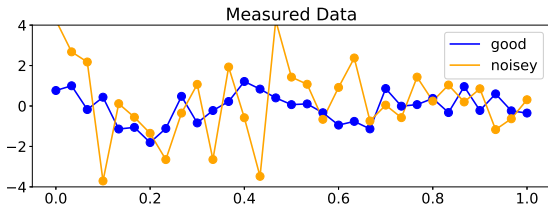## Constructing Confidence Intervals: From Fitted Covariance

Error bars, $\delta p_m$, for a fitted parameter can be constructed from the covariance $\Sigma_{\hat{p}_{LS}}$ and a $\Delta\chi^2$:

$$\delta p_m = \pm\sqrt{\Delta\chi^2}\sqrt{\Sigma_{mm}}$$

The value of $\Delta\chi^2$ selects the "significance level":

- $\Delta\chi^2$ is found in lookup tables calculated from the CDF of the $\chi^2$ distribution
- Single parameter fit, $N = 1$:
    - a 68% significance: $\Delta\chi^2 = 1$
    - a 95.4% significance: $\Delta\chi^2 = 4$
- Two parameter fit, $N = 2$:
    - a 68% significance: $\Delta\chi^2 = 2.3$

# Challenges With Constructing Confidence Intervals

## Validity of Confidence Intervals

Only quantitatively valid when:

- measurement errors are Gaussian, and
  - the model $f(p)$ is linear in for all p, or
  - measurement errors are small enough that $f(p)$ can be accurate approximated by a linear model in the region around p

Otherwise, alternative fitting methods are required: Monte Carlo, Bayesian, etc.

## Goodness of Fit

How do we know if the fit is even meaningful? The standard goodness of fit test involves computing the "reduced chi-squared":

$$\chi_\nu^2 = \chi^2/(m - n + 1)$$

Then, typically:

- $\chi_\nu^2 \approx 1$: a good fit
- $\chi_\nu^2 << 1$: an "over fit"
- $\chi_\nu^2 >> 1$: a poor fit

The $\chi_\nu^2$ could also be slightly larger or smaller than 1 depending on how accurately one is able to estimate the input measurement errors.

## Summary

Now we can answer the question: Are the fitted parameters meaningful?

- What is the uncertainty in the fitted parameters?
  - Error bars correspond to confidence intervals (CI)
  - CIs are constructed from covariance of the fitted parameters
  - For a 68% CI, interpretation is: "If we could hypothetically make and infinite set of new measurements and fit each of those, 68% of the time the 'true' value of the parameter would lie within the CI."
- Is the fit good?
  - Compute the reduced chi-squared
  - $\chi_\nu^2 \approx 1$: usually means the model accurately represents the data
- All of this error analysis depends on the assumption that measurement errors are **Gaussian** distributed with covariance $\Sigma_e$ such that $(y_m - f_m)/\sigma_m$ are $\mathcal{N}(0, 1)$